

Dynamic Mathematics for Automated Machine Learning Techniques

Nicholas I-Hsien Kuo

A thesis submitted for the degree of
Doctor of Philosophy of
The Australian National University

August 2021

© Nicholas I-Hsien Kuo 2021

Preliminary version – 31 August 2021

Except where otherwise indicated, this thesis is my own original work.

Nicholas I-Hsien Kuo
31 August 2021

to myself, to those who supported me, and to those who inspired me

Acknowledgments

I would like to express my sincere gratitude towards everyone who made this thesis possible. I was extremely lucky and privileged to have five brilliant and kind supervisors – Prof. Hanna Suominen, Dr. Mehrtash Harandi, Dr. Nicolas Fourrier, Dr. Gabriela Ferraro, and Dr. Christian Walder. I would also like to thank my family, to always be supportive and proud of me whether my manuscripts were accepted or rejected. In addition, I cannot possibly describe how much I am grateful to my friends; besides laughters, you really helped me to stay focus when life is harsh.

This research was supported by the Australian Government Research Training Program (AGRTP) Scholarship. Thank you for providing the financial support to help me pursue knowledge.

Nic¹, you are still young ... go do a PhD, it will probably be the only time in your life when you are allowed to think freely without worrying about other stuff...

This was suggested to me from Nicolas when I was half way through my Masters at the University of Auckland. Besides my parents, Nicolas has supported me for the longest; in both my Masters and later in my Doctoral studies at the Australian National University. He is a great data scientist and is someone who I really admire.

I was also fortunate to meet my supervision chair Hanna. Hanna always seemed to have endless energy for research. She is always punctual, disciplined, and yet is also a very amiable person. During the beginning of my doctoral studies, she patiently helped me to construct my own style of critical thinking. Furthermore, Hanna always give me good advises for advancing my career.

Through Hanna, I met a lot of talented researchers including Mehrtash, Gabriela, and Christian. Mertash always gives me the motivation to think big. He is highly aware of the lastest trend of machine learning research and has always pointed me towards interesting research directions. Gabriela provides me with a lot of good conversations. She helps me to navigate the landscape of natural language processing research and informs me of important yet incompleted research questions. Christian also spends a lot of time with me discussing and brainstorming ways to characterise neural network behaviours in a mathematical fashion. My research was only made possible because I had a great team of caring mentors.

I would also like to express my gratitude towards the Australian National University (ANU) and Canberra. Canberra is quite the small place – actually, much smaller than I could have first imagined. Yet this was the place where I met some of

¹Note, the student-author's name is Nic; not to be confused with the supervisor Nicolas.

the kindest, smartest, but most humble people in this world. During my stay, I had great opportunities to interact with professionals of all different background. I made friends with mathematicians, computer scientists, logicians, physicists, musicians, linguists, geopolitical analysts, economists, and medical experts. People thought I was a workaholic; but the truth was, I was simply happy to learn more from all these experts that I am fortunate enough to call friends. More importantly, they provided me with mental stability when my loved ones passed away, during the horrible 2019 Australian bushfire, and during this still-ongoing COVID19 pandemic.

I would also like to thank the Centre of Big Data in Health at the University of New South Wales (UNSW). Near the end of my PhD with the ANU, I was fortunate to be accepted as a post-Doc research fellow at the UNSW. I am extremely happy that I can continue to pursue knowledge, and apply my machine learning skills for greater good.

Many thanks for this wonderful experience,
- Nic

Abstract

Data is everywhere in the modern world. From a heartwarming grin to a bulk international commercial order, our lives can be easily reflected and captured by our activities on social media platforms, street security cameras, and internet footprints. As data becomes more accessible, our society inevitably turns to data-mining out of curiosity for self-exploration and value extraction for financial competition. However and as our data becomes both quantitatively and qualitatively complicated, methods that attempt to understand data are also changing. As of the time that this thesis is written, we are currently living in a world where *Neural Networks* are slowly becoming the new de-facto standard for understanding industry big data.

Machine learning consists classes of modelling techniques that are capable of finding patterns directly from data itself. Of its many classes, neural networks have been gaining a large amount of following as they can learn in a near-automatic fashion with minimal amount of prior knowledge on the data. This property makes neural networks highly versatile and thus they are heavily favoured in becoming the new industry standard models. Despite their popularity, several neural network methods remain opaque and the research community still struggles to fully explain their behaviours. The lack of interpretability not only halts innovation in designing better models, but it also complicates the difficulty to *automate machine learning* (AutoML) for neural networks to work in a dynamic environment.

To this end, this thesis is devoted towards increasing the understanding of neural network properties. More specifically, we will analyse neural networks through the mathematical lens of *dynamical systems* in order to clarify how they transform their input data. Then, the mathematical insights will assist us in addressing three interconnected topics – *recurrent neural networks*, *meta-learning*, and *continual learning* – to increase the flexibility and utility of neural network models.

Our first topic of discussion is recurrent neural networks (RNNs) – neural network models that are capable of handling sequential data. Advance RNNs possess an internal state, also known as the *cell*, to represent events that have occurred in the past. RNNs are intricate models, and the cell is recursively updated with several internal variables as the RNN receives more data. However, the variables that are employed to update the cell *do not have immediately clear meanings*. Thus in order to increase RNN interpretability, we compare RNNs to well-studied dynamical systems. By doing so, we demonstrate that we can establish connections between literature of machine learning and theoretical physics. This contribution also allows us to highlight the core inferential mechanism in RNN processing. It in turn enables us to propose a systematic simplification to create small-size RNNs that removes up to 85.4% secondary parameters while attains competitive performances.

Our second topic concerns meta-learning, the AutoML technique of rapid knowl-

edge acquisition. We focus on the meta-learning technique of *learning to learn* (L2L) which employs an RNN to update another base learner network. The RNN of L2L explores optimisation strategies superior than handcrafted alternatives to configure base learners in less time and with fewer data. However, current L2L algorithms are limited in practicality. While they are efficient, a meta-trained RNN can *only update one unique base learner for one specific dataset*. Our contribution on meta-learning is two-fold and we aim to increase meta-learning utility. Leveraging on our mathematical understanding, we propose an alternative RNN design to address input-domain heterogeneity. This enables a meta-trained RNN to update a base learner on datasets unseen during the meta-training phase. In addition, we design an L2L algorithm that prunes up to 83.7% base learner parameters while configuring them. Hence showing that learnt optimisation algorithms can be also be used to explore the internal limits of knowledge capacity of a base learner for learning a given task.

Our third topic focuses on continual learning, the AutoML technique for the sequential acquisition of skills for a single neural network. Under the conventional experimental setup, a trained neural network can achieve near-human level inference on a single task. Though impressive, this ability alone is not enough for neural networks to be employed in a dynamic environment. In such an environment, the context that surrounds a learning agent is subjected to continuous changes and it hence presents unforeseen challenges to a learning agent while requiring it to learn in a lifelong fashion. However, most existing neural networks are known to suffer from *catastrophic forgetting* – the abrupt loss of knowledge that solves previous tasks while assimilating information for solving new tasks. The catastrophic forgetting phenomenon hence limits neural networks to stably gain additional functions over time. Our contribution to continual learning is that, we formulate catastrophic forgetting as a network configuration problem. Leveraging our mathematical interpretation for meta-learning, we devise a learnt algorithm that enables a single network to robustly sequentially learn multiple tasks under a noisy dynamic environment. In addition, we also design a neural network architecture of which itself makes forgetting difficult to occur, hence attaining most of its previously learnt knowledge.

The three topics of concerned in this thesis are all intertwined as illustrated in Figure 1. Furthermore, they delve into two important aspects of neural network modelling. First, our RNN chapter aims to provide a mathematic understanding on neural network inferences. Second, our meta-learning chapter is devoted towards clarifying neural network configuration. Then, our continual learning chapter combines the insights of the two prior chapters to cohesively create novel machine learning algorithms that automate tedious hyper-parametric tuning that have been traditionally configured manually.

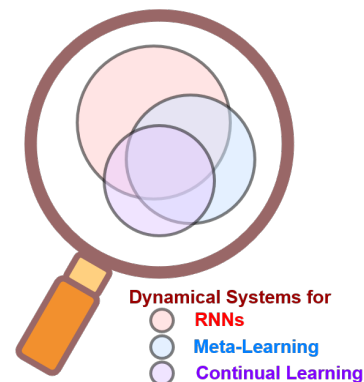


Figure 1: Our Studies

Keywords: Automated Machine Learning,

Dynamical Systems, Recurrent Neural Networks, Meta-Learning, Continual Learning

Publications

The majority of the thesis has been published in conference proceedings and in peer-reviewed workshops. The codes for replicating the results of this thesis are made publicly available and our readers can find them in the student-author's GitHub repository <https://github.com/Nic5472K>.

The papers and codes will form as bases for our future work.

On RNNs

Kuo, N. I.; Harandi, M. T.; Fourrier, N.; Walder, C.; and Ferraro, G.; Suominen, H., 2018. **DecayNet: A Study on the Cell States of Long Short Term Memories**. *An archived manuscript*.

Kuo, N. I.; Harandi, M. T.; Fourrier, N.; Walder, C.; and Ferraro, G.; Suominen, H., 2020. **An Input Residual Connection for Simplifying Gated Recurrent Neural Networks**. In the *International Joint Conference on Neural Networks*.

On Meta-Learning

Kuo, N. I.; Harandi, M. T.; Fourrier, N.; Walder, C.; and Ferraro, G.; Suominen, H., 2020. **MTL2L: A Context Aware Neural Optimiser**. In the *Workshop on Automated Machine Learning of the International Conference on Machine Learning*.

Kuo, N. I.; Harandi, M. T.; Fourrier, N.; Walder, C.; and Ferraro, G.; Suominen, H., 2020. **M2SGD: Learning to Learn Important Weights**. In the *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop on Continual Learning in Computer Vision*.

On Continual Learning

Kuo, N. I.; Harandi, M. T.; Fourrier, N.; Walder, C.; and Ferraro, G.; Suominen, H., 2021. **Learning to Continually Learn Rapidly from Few and Noisy Data**. In the *Meta-Learning and Co-Hosted Competition of the AAAI Conference on Artificial Intelligence*.

Kuo, N. I.; Harandi, M. T.; Fourrier, N.; Walder, C.; and Ferraro, G.; Suominen, H., 2021. **Plastic and Stable Gated Classifiers for Continual Learning**. In the *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop on Continual Learning in Computer Vision*.

Basic Notation Conventions

Below is a list for the language of mathematics used in this thesis. Most of our notations agree either completely or partially to the traditional notations. There are, however, cases such as the use of subscript of which we abuse the notation.

Of regular usage, unless specified otherwise

For network architectural descriptions

Bold lower case	a	a vector
Bold font upper case	W	a matrix
Symbol then brackets	$\mathcal{H}(\cdot)$	a function
Vertical line in function	$\mathcal{H}(\cdot)$	a function given θ
Upper case double r	\mathbb{R}	a real value
Real value with superscript	$\mathbb{R}^{ \theta }$	a real value of size $ \theta $
The cross symbol	$d_1 \times d_2$	for dimensionality, d_1 by d_2
The tanh symbol	tanh	the hyperbolic tangent function
The lower case Greek sigma	σ	the sigmoid function

For optimisation descriptions

Lower case Greek theta	θ	parameters of a neural network
Arrow pointing left	$A \leftarrow B$	A is updated into B with B
Nabla symbol	∇	gradient, specifically for loss
Italic upper case J	J	the Jacobian
Carl Jacobi's curly d	∂	the partial derivative symbol
Function with an apostrophe	$\mathcal{H}'(\cdot)$	the derivative of a function

Other common functions

Upper case double i	\mathbb{I}	an identity matrix
Upper case double e	\mathbb{E}	the mean function
The odot symbol	\odot	the element-wise product function
The log symbol	log	the log function
The sgn symbol	sgn	the sign function
The diag symbol	diag(a)	a diagonal matrix with a as the diagonal entries

For set or vector descriptions

Italic upper case b	B	a batch (set) of data
Some common notations	\in	set membership
	\forall	for all
	\cup	union
	\sim	sample from

Other common notations

Vector with apostrophe	\mathbf{a}'	the time derivative of a vector
Italic lower case e	e	the exponential constant
Lower case Greek pi	π	Archimedes' constant
Upper case Greek sigma	Σ	a summation
Upper case Greek pi	Π	a product

Of irregular usage, context always specified in text

Subscript	\mathbf{a}_{γ_L}	a vector of the γ_L th layer
	\mathbf{a}_t	a vector at time step t
	\mathbf{a}_j	the j element of a vector
	etc.	
Vertical bars	$ B $	the cardinality of a batch
	$ \theta $	the cardinality of the parameters
	$ \mathbf{a} $	the element-wise absolute values of a vector
	etc.	
Superscript with brackets	$\mathcal{X}_t^{(\mathcal{Q})}$	See Equation (3.13), \mathcal{X} specific to \mathcal{Q}
	$\mathcal{L}^{(n)}$	See Equation (3.15), the n th element of \mathcal{L}
	etc.	

Contents

Acknowledgments	vii
Abstract	ix
Publications	xi
Basic Notation Conventions	xiii
1 Introduction	1
1.1 An Overview on Automated Machine Learning	2
1.2 The Conventional Setup of Neural Network Experiments	3
1.2.1 Before a Neural Network Experiment	3
1.2.2 The Forward Pass	3
1.2.3 The Backward Pass	3
1.3 The Alternative Setups to Be Found therein this Thesis	4
1.4 Our Perspective towards Understanding AutoML	5
1.4.1 Why Study Shortcut Connections?	5
1.4.2 Studying Shortcut Connections via RNNs via Dynamical Systems	6
1.5 Outline, Contributions, Publications, and Codes	7
1.5.1 Chapter 2: Preliminaries	7
1.5.2 Chapter 3: Related Work	7
1.5.3 Chapter 4: Understanding RNNs with Dynamical Systems	8
1.5.4 Chapter 5: Learning to Learn beyond Batch-Efficient Updates	9
1.5.5 Chapter 6: Topics on Simple and Efficient Continual Learning	10
1.5.6 Chapter 7: Thesis Conclusion	11
2 Preliminaries: Past and Present on Deep Models	13
2.1 Residual Networks and Highway Networks	15
2.1.1 Training Deep Networks is Difficult	15
2.1.2 Related Ideas: Why Do We Need Deep Learning?	16
2.1.2.1 Kernel-Based Methods	16
2.1.2.2 Decision Trees	17
2.1.2.3 Are Neural Networks Over-Fitted Models?	17
2.1.2.4 So Why do Neural Networks Need so Many Parameters?	18
2.1.2.5 Complementary Studies to Architectural Modifications	18
2.1.3 Section Conclusion	21

2.2	Loss Landscape & Loss Minima	23
2.2.1	Loss Landscape	23
2.2.2	Visualising Parametric Updates	23
2.2.3	Good Loss, Bad Loss?	24
2.2.3.1	On the Depth of Minima	24
2.2.3.2	On the Width of Minima	25
2.2.3.3	Smoothness within a Minimum	25
2.2.3.4	What Constitutes the Outline of the Loss Landscape?	25
2.2.3.5	Limitations of the Loss Landscape Representation	26
2.2.4	Related Ideas: How were the Loss Minima Studied?	26
2.2.4.1	Consequences of Large Batch Sizes	27
2.2.4.2	Alternative Opinions on Minima	28
2.2.5	Section Conclusion	28
2.3	Some Fundamentals on Dynamical Systems	31
2.3.1	Differential System and Difference System	31
2.3.2	Visualising Difference Systems with Cobweb Diagrams	31
2.3.3	Visualising Differential Systems with Vector Fields	32
2.3.4	Equilibrium	33
2.3.4.1	Equilibria in Difference Systems	33
2.3.4.2	Equilibria in Differential Systems	34
2.3.5	Related Ideas: Existing Applications of Dynamical Mathematics in Understanding Neural Networks	34
2.3.5.1	Dynamical Studies in the Backward Pass	35
2.3.5.2	Non-Dynamical Studies in Understanding the Forward Pass of Shortcut Connections	36
2.3.5.3	Dynamical Studies in Understanding the Forward Pass of Shortcut Connections	36
2.3.6	Section Conclusion	38
2.4	Chapter Conclusion	39
3	Related Work	41
3.1	Recurrent Neural Networks and Their Interpretations	43
3.1.1	Early RNN Designs and the Gradient Vanishing Problem (Again)	43
3.1.2	The Long Short-Term Memory Design	44
3.1.3	A Brief History on LSTM's Development	45
3.1.3.1	Error Propagation through a Carousel	45
3.1.3.2	Learning to Predict while Learning to Forget	46
3.1.3.3	Learning to Predict with Direct Access to the Cell	47
3.1.4	The Gated Recurrent Unit Design and the Non-Immediately Clear Meanings of Gated Units	47
3.1.4.1	GRU: The Other De-Facto Standard RNN	47
3.1.4.2	A Comparison between GRU and LSTM	48
3.1.5	Section Recap	49
3.1.6	Strategies: Remove, Grow, and Observe	50

3.1.6.1	Removing Network Components: What to Keep? . . .	50
3.1.6.2	An Extensive Search on Sibling Models: What Else to Include?	50
3.1.6.3	An Analysis on the Network Input to Output: What Has been Learnt?	52
3.1.7	Examples of Existing Work on Remodelling LSTM and GRU . .	53
3.1.8	Section Conclusion	55
3.2	Meta-Learning & On Learning to Optimise	59
3.2.1	Offline Meta-Learning: Model Agnostic Meta-Learning	59
3.2.2	Online Meta-Learning: Metric-Based Meta-Learning	60
3.2.3	Online Meta-Learning: Model-Based Meta-Learning	61
3.2.4	Online Meta-Learning: Optimisation-Based Meta-Learning . . .	63
3.2.5	Details on Learning to Optimise	64
3.2.5.1	Method Introduction	64
3.2.5.2	The LSTM-Optimiser – Replacing Gradient Descent . .	64
3.2.5.3	The LSTM-Optimiser – Training the Neural Optimiser .	65
3.2.5.4	The LSTM-Optimiser – The Meta-Learner Input	65
3.2.5.5	The LSTM-Optimiser – The Meta-Learner Inference and Output	66
3.2.5.6	The LSTM-Optimiser – The Default Hyper-parameteric Settings	66
3.2.5.7	Related Work to the LSTM-Optimiser	67
3.2.6	Section Conclusion	68
3.3	Continual Learning & On Episodic Memories	71
3.3.1	The Framework of Continual Learning	71
3.3.2	Comparing Continual Learning to Other AutoML Directions . .	71
3.3.3	Continual Learning and Potentially Related Concepts in RNNs .	72
3.3.4	Scenarios for Continual Learning	73
3.3.5	Archetypes in Continual Learning	74
3.3.5.1	Continual Learning via Regularisation	74
3.3.5.2	Continual Learning via a Dynamic Architecture	75
3.3.5.3	Continual Learning via Memory Replay	76
3.3.6	Our Method of Choice	78
3.3.7	Section Conclusion	78
4	Understanding RNNs with Dynamical Systems	81
4.1	A Study on the Cell States of Long Short-Term Memories	83
4.1.1	Section Introduction	83
4.1.2	LSTMs as Difference Equations	83
4.1.3	The Cell State Revisited	84
4.1.3.1	The Importance of the LSTM Cell State	84
4.1.3.2	A Piece-wise Analysis on the Cell State	84
4.1.3.3	The Forget Gate and the Catch-and-Release Dynamics .	85

4.1.3.4	The Near-Random Dynamics of a Vanilla LSTM Cell State	86
4.1.3.5	Interpretable Cell State Dynamics under Controlled Forget Gates	86
4.1.4	DecayNets	87
4.1.4.1	From 1 to 0: A Monotonic Decay	87
4.1.4.2	Using Decay to Increase Interpretability	88
4.1.5	Section-wise Additional Related Work	88
4.1.6	Experiments and Results	89
4.1.6.1	Permuted Sequential MNIST	89
4.1.6.2	Learning to Optimise	90
4.1.7	Section Conclusion	92
4.1.8	An Evaluation of this Study	93
4.2	An Input Residual Connection for Simplifying Gated Recurrent Neural Networks	96
4.2.1	Section Introduction	96
4.2.2	Connecting GRUs to Hopfield Networks	96
4.2.2.1	Recursive Feedbacks as Systems of Differential Equations	97
4.2.2.2	Deriving the GRU from the Dynamic Hopfield Network	97
4.2.3	An Input Residual Connection	98
4.2.3.1	An Analysis on the Non-linearity of \mathbf{a}_t	99
4.2.3.2	An Analysis on the Non-linearity of \mathbf{i}_t	99
4.2.3.3	An Analysis on the Recurrent Feedback	99
4.2.3.4	An Input Residual Connection	100
4.2.4	IRC-/IHC-GRNNs	101
4.2.4.1	An Input Highway Connection	103
4.2.5	Experiments and Results	103
4.2.5.1	Language Modelling	104
4.2.5.2	Image Generation	106
4.2.5.3	Learning to Optimise	106
4.2.6	Section-wise Additional Related Work	107
4.2.7	Section Conclusion	109
4.3	Chapter Summary	110
4.4	Chapter Conclusion	111
5	Learning to Optimise beyond Rapid Knowledge Acquisition	113
5.1	MTL2L: A Context Aware Neural Optimiser	117
5.1.1	Section Introduction	117
5.1.2	Reasons of the Limitations of the Vanilla Neural Optimiser	117
5.1.3	Problems of Employing RNNs as Generally Applicable Optimisers	119
5.1.4	Context Awareness through HyperNetwork	120
5.1.4.1	The MTL2L Formulation	121

5.1.4.2	MTL2L Implementation Details: SVD	122
5.1.4.3	MTL2L Implementation Details: Naïve Hypernetwork Eigenvalues	122
5.1.4.4	MTL2L Implementation Details: Stable Hypernetwork Eigenvalues	122
5.1.4.5	MTL2L Implementation Details: On Multi-Task Learning	123
5.1.5	Experiments and Results	123
5.1.6	Section-wise Additional Related Work	125
5.1.7	Section Conclusion	126
5.1.8	An Evaluation on Our Solved Problem	127
5.2	M ² SGD: Learning to Learn Important Weights	129
5.2.1	Section Introduction	129
5.2.2	On MetaSGD	130
5.2.2.1	The Main Differences to the LSTM-Optimiser	130
5.2.2.2	Why Changing from LSTM-Optimiser to MetaSGD?	130
5.2.3	Extensions to MetaSGD	132
5.2.3.1	The Dynamic MetaSGD (DMSGD)	133
5.2.3.2	The Masked MetaSGD (M ² SGD)	133
5.2.4	Experiments and Results	134
5.2.5	Section-wise Additional Related Work	138
5.2.6	Section Conclusion	139
5.3	Chapter Conclusion	140
6	Topics on Simple and Efficient Continual Learning	141
6.1	Learning to Continually Learn Rapidly from Few and Noisy Data	143
6.1.1	Section Introduction	143
6.1.2	Replay with Episodic Memory	144
6.1.2.1	Episodic Memory Sampled from a Tiny Memory Storage	145
6.1.3	MetaSGD for Continual Learning	145
6.1.3.1	An Analysis on the Update of Experience Replay for Continual Learning	146
6.1.4	Implementation Details on MetaSGD-CL	147
6.1.5	Experiments and Results	148
6.1.6	Section-wise Additional Related Work	155
6.1.7	Section Conclusion	155
6.1.8	Re-imagining Continual Learning	156
6.2	Highway Classifier Networks for Plastic yet Stable Continual Learning	158
6.2.1	Section Introduction	158
6.2.2	Thesis Recap: Highway Networks vs Gated RNNs	159
6.2.3	A Dynamic Theory in the Forward Pass: A Common Vector Field with Many Stable Equilibria	159

6.2.4	Evidence in the Backward Pass: HCNs Prevent Parametric Over-Modification	160
6.2.5	Experiments and Results	162
6.2.6	Section-wise Additional Related Work	167
6.2.7	Section Conclusion	168
6.3	Chapter Conclusion	169
7	Thesis Conclusion	171
7.1	Thesis Summary	171
7.2	Thesis Contributions	173
7.2.1	On Increasing the Interpretability for RNNs	173
7.2.2	On Extending Meta-Learning for Better Optimisation Algorithms	176
7.2.3	On Improving Continual Learning Performances	179
7.3	Limitations	182
7.4	Future Work	184
7.5	Final Remarks	185
A:	Datasets	213
B:	Additional Notes for the Preliminaries	217

List of Figures

1	Our Studies	x
1.1	MLP Forward Pass	3
1.2	MLP Backward Pass	4
1.3	The Structure of this Thesis	7
1.4	Modification on LSTM	8
1.5	Modification to L2L	9
1.6	L2L for Continual Learning	10
2.1	Three Connections	16
2.2	Sample Illustrations	23
2.3	Parametric Update	23
2.4	The Consequences of Different Learning Rates	24
2.5	Three Minima	24
2.6	Minima Scenarios	25
2.7	Minima are not Static	26
2.8	Simple Examples of Cobwebs	32
2.9	32
2.10	Progress Update	39
3.1	An Early RNN Design	43
3.2	The LSTM Design	44
3.3	The GRU Design	47
3.4	Progress Update	56
3.5	The L2L Experimental Setup	64
3.6	Progress Update	68
3.7	Progress Update	79
4.1	Examples of Difference Equations Variable Propagation	84
4.2	The Cell State Propagation of Controlled and Uncontrolled Forget Gates	86
4.3	L2L vs Gradient Descent	91
4.4	GRU vs IRC-GRU	101
4.5	DRAW for Generating SVHN Images	106
4.6	L2L for Updating ResNet-34s	107
4.7	Progress Update	111
5.1	Selecting Minimum	114

5.2	Neural Optimiser Applicability	117
5.3	Input Variance	119
5.4	The Well-behaving Meta-Testing Results of Scenario 1	123
5.5	Meta-tested on an	124
5.6	MetaSGD vs M ² SGD	132
5.7	Loss curve (first 5000 steps) for FashionMNIST	135
5.8	Loss curve (first 1000 steps) for FashionMNIST	135
5.9	Loss Curve for Updating	137
5.10	Loss Curve for Updating	137
5.11	Progress Update	140
6.1	Exiting a Local Minimum	143
6.2	Task-wise Directional	146
6.3	Performances for Permuted MNIST	150
6.4	Performances for	150
6.5	Ring Buffer of Different Sizes for Permuted MNIST	151
6.6	Permuted MNIST but with Reduced Iteration	152
6.7	Permuted MNIST but with Noise Injection	152
6.8	Transition in Task-wise β Magnitudes	153
6.9	Utilising Different Basins of Attractors	159
6.10	Different Designs of Continual Learning Classifier Networks	161
6.11	Loss vs Iteration (MNIST)	162
6.12	Gradient vs Iteration (MNIST)	162
6.13	ACC of Naïve	163
6.14	ACC of Naïve	163
6.15	ACC with Different Techniques on Incremental Cifar100	164
6.16	HCNs Provided Robustness and Decreased Variability	166
6.17	Sequential Learning while Only Fine-Tuning Layer 1 of HCNs	167
6.18	Progress Update	169

List of Tables

2.1	Forward and Backward Passes of Residual and Highway Connections .	15
3.1	A Comparison between GRU and LSTM	48
3.2	A Comparison between GRU and LSTM	49
4.1	Results on Perm-SeqMNIST	89
4.2	A comparison between baseline GRNNs against IRC-GRNNs	102
4.3	Perplexity on the PTB dataset.	104
5.1	Performance Details of MLP learners	136
5.2	Performance Details of ResNet20 learners	137
6.1	Metrics for Permuted MNIST	149
6.2	Metrics for	150
6.3	Ratios of Extreme Task-wise β s with $\kappa = 0.1$	153
6.4	Metrics for the Ablated Version of Permuted MNIST	154
6.5	An Overview of the Results	165
6.6	An Overview of the Results	167

Introduction

Machine Learning and *Neural Networks* have been gaining popularity and are widely considered as the driving force of the Fourth Industrial Revolution [BMBF, 2014]¹. However, modern machine learning techniques such as backpropagation training was firmly established in 1986 (Rumelhart et al.) while computer vision was revolutionised in 2012 (Krizhevsky et al.) with AlexNet. Why are neural networks still not an integral part of our society? Because they are difficult to implement in practice.

I'd like to use machine learning, but I can't invest much time.

The term *Automated Machine Learning* (AutoML) [Hutter et al., 2019] was first coined by Professor Frank Hutter of the University of Freiburg. Machine learning is not simple; it requires a practitioner to have a thorough understanding on their data and the components which their model entails. AutoML aims to automate all tedious aspects of machine learning to form a clean data analysis pipeline.

This PhD thesis in computer science develops and understands ways to automate machine learning. Specifically, we focus on *Recurrent Neural Networks* (RNNs) [Hochreiter and Schmidhuber, 1997b], *Meta-Learning* [Andrychowicz et al., 2016], and *Continual Learning* [Lopez-Paz and Ranzato, 2017]. We study continual learning to enable a network to sequentially acquire skills in a dynamic environment; we study meta-learning to understand how a network can be configured efficiently; and we study RNNs to understand the consequences of consecutive actions in inference.

Our RNN-study focuses on the mathematical interpretability of network components (like that in [Sneyd et al., 2017]). We will describe a large variety of RNNs as one mathematical class to understand their core mechanism. This will enable us to extend meta-learning beyond configuration for network pruning and continual learning. This will also provide insights to understand how a single network should be consecutively configured and will lead us to create a simple generic patch compatible to several existing continual learning archetypes. This patch will enhance the robustness of continual learning techniques to achieve better generalisation.

By and large, this thesis presents a series of extensions to enable AutoML to be made simple, efficient, and robust. More importantly, all of our methods are motivated with mathematical understandings through the lens of dynamical systems. Thus, we also increase the interpretability of AutoML concepts.

¹By the German Federal Ministry of Education and Research (BMBF).

1.1 An Overview on Automated Machine Learning

Automated machine learning (AutoML) [Hutter et al., 2019] is an umbrella term for several machine learning disciplines. It includes, but is not limited to, *Transfer Learning* [Sharif Razavian et al., 2014], *Bayesian Optimisation* [Feurer and Hutter, 2019], *Few Shot Learning* [Vinyals et al., 2016b], *Zero-Shot Learning* [Norouzi et al., 2013], *Neural Architecture Search* [Elsken et al., 2019], *Reinforcement Learning* [Sutton and Barto, 2018], *Meta-Learning* [Vanschoren, 2019], and *Continual Learning* [Parisi et al., 2019]. From this list of items, we can see that AutoML has a simple yet ambitious goal: *to automate everything*.

Modelling data with neural networks can be a difficult and time-consuming task. Practitioners first need a thorough understanding on the properties of the data. Then, practitioners will proceed to select modules of their choice. For deep learning, it is also necessary to consider whether to combine different modules, such as attaching an RNN to a *Convolutional Neural Network* (CNN). This is further complicated when practitioners need to decide the order of the placements of the modules and their respective depths. Practitioners also have to consider whether they will need to employ specific initialisation and normalisation² to enable the training of large models. Finally and once an initial model is proposed, practitioners will then go through tedious rounds of hyper-parameter search for the optimal network performances.

AutoML provides a toolbox of machine learning disciplines that aims to ease the difficulty in training neural networks. Transfer learning borrows and reuses components of a trained working model as the starting point of a new task. Bayesian optimisation and meta-learning are designed to achieve batch-efficient machine learning. Few shot learning aims to conduct fast inferences from learning only from a handful of data; whereas zero-shot learning intends to recognise objects whose instances may not have been seen during training. Neural architecture search seeks to automate the process of architectural engineering, while reinforcement learning is a paradigm which enables a learning agent to find non-intuitive working policies via maximising rewards to positive actions taken by the agent. Moreover, there is continual learning, the study on the sequential acquisition of new functionalities in a single model thereby making the said single model adaptive to new ideas.

On the surface, it may seem as if AutoML is overly generalised and that it lacks a certain purpose. However and at their core, each of the aforementioned AutoML disciplines shares the common interest in exploiting domain specific knowledge to effectively guide the learning process to help the design and simplification of the next generation of learning systems. Furthermore, it is common in machine learning literature to see several of these disciplines being put into practice in conjunction with another. For this thesis, we will put a focus on proposing novel methods for meta-learning and continual learning; and explore the possibilities in leveraging meta-learning for continual learning to yield new methods that inherit the advantages of both fields.

²Interested readers should consult Section 2.1.2.5 for more information.

1.2 The Conventional Setup of Neural Network Experiments

This thesis is written with the assumption such that our readers are already familiar with neural networks. Nonetheless, several AutoML disciplines, including meta-learning and continual learning, modify the conventional setup³ of a neural network experiment consisting the network *forward pass* and network *backward pass*. Hence we will briefly describe the conventional setup in this section.

1.2.1 Before a Neural Network Experiment

Neural networks are powerful algorithms that heuristically learn to approximate functions. To make use of these algorithms, a practitioner is required to specify a number of hyper-parameters. This thesis will pay an especially close attention to the choice of network architecture, the loss function, and the optimisation function. These specified arguments dictate how a neural network model is iteratively updated through alternatively executing the forward and the backward passes.

1.2.2 The Forward Pass

The forward pass refers to the data driven inferential process, which takes place on a defined network architecture. Consider the *Multiple Layer Perceptron* (MLP) neural network illustrated in Figure 1.1. The MLP receives input x and outputs prediction \hat{y} . Internally, the MLP consists a set of transformations given as

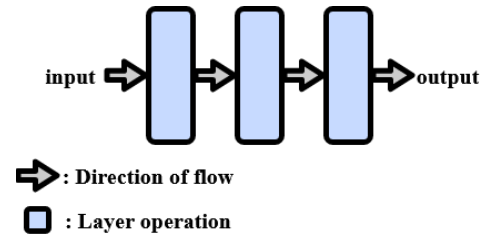


Figure 1.1: MLP Forward Pass

$$\mathbf{a}_{L_\gamma} = \mathcal{H}(\mathbf{a}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma}) \quad (1.1)$$

where layer γ has input $\mathbf{a}_{L_{(\gamma-1)}}$, output \mathbf{a}_{L_γ} , and a non-linear transformation \mathcal{H} parameterised by weights \mathbf{W}_{L_γ} . Traditionally, \mathbf{a}_{L_0} is denoted by x and that \mathbf{a}_{L_Γ} is denoted by \hat{y} ; we use Γ to represent the total amount of layer operations.

A deep model is consisted of multiple layers; and deep learning refers to the practice of combining various machine learning modules for abstracting different levels of details from the input data [Szegedy et al., 2015]. This thesis mainly focuses on 4 different types of modules – the feedforward connections of Equation 1.1, the *Residual Connections* of [He et al., 2016a], the *Highway Connections* of [Srivastava et al., 2015], and various recurrent connection designs in different RNNs.

1.2.3 The Backward Pass

The backward pass refers to the procedure of updating neural network parameters. For instance, the term \mathbf{W}_{L_γ} of Equation 1.1 needs to be configured correctly to make the MLP function properly. As depicted in Figure 1.2, the backward pass is executed

³See Tieleman and Hinton [2012] as an example for the conventional neural network setup.

after that of the forward pass. Our work will focus on the error back-propagation [Rumelhart et al., 1986] approach for parametric fine-tuning.

There are 3 steps to batch error backpropagation – first, we sample a batch of data B from the training data and send it through the forward pass; then, we measure the error between the network output \hat{y} to the desired output y with some cost function \mathcal{C} ; and last we update the network based on the error using gradient descent methods [Robbins and Monro, 1951]. Equivalently, our model of choice F has parameters θ_t at time t and generates output

$$\hat{y} = F(x|\theta_t) \quad (1.2)$$

of which has a batch loss of

$$\mathcal{L} = \frac{1}{|B|} \sum_{(x_i, y_i) \in B} \mathcal{C}(\hat{y}_i, y_i) \quad (1.3)$$

for the update

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L} \quad (1.4)$$

with learning rate α and derivatives of the loss with respect to the parameters $\nabla_{\theta} \mathcal{L}$.

1.3 The Alternative Setups to Be Found therein this Thesis

As mentioned in the foreword of this chapter, this thesis visits the AutoML disciplines of meta-learning and continual learning. The chosen techniques of our studies deliver modifications to the backward pass to influence the update of a network.

Learning to Optimise

On one hand, our meta-learning approach is heavily influenced by the work of Andrychowicz et al. [2016]. The core concept of their paper is to replace $-\nabla_{\theta} \mathcal{L}$ of Equation (1.4) with an RNN. They demonstrated that their *RNN-optimiser* updated a base learner more efficiently than handcrafted gradient descent algorithms. Our RNN study in Chapter 4 will provide insights to extend meta-learning in Chapter 5.

Sequentially Learning via Replay

On the other hand, our continual learning study is based on the work of Lopez-Paz and Ranzato [2017]. In their work, the authors preconditioned $-\nabla_{\theta} \mathcal{L}$ of Equation (1.4) to create new packages of updates that enabled the sequential acquisition of knowledge in a single network. In Chapter 6, we will introduce how meta-learning can greatly enhance the performance of sequential learning.

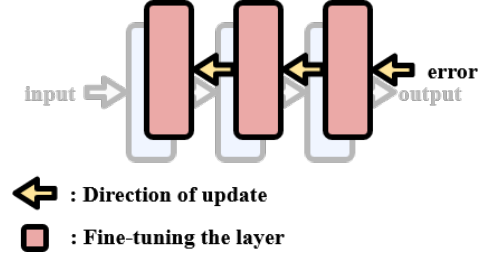


Figure 1.2: MLP Backward Pass

1.4 Our Perspective towards Understanding AutoML

One permeating theme of this thesis surrounds *Shortcut Connections* [Schraudolph, 1998]. Shortcut connection is a popular component shared among several state of the art neural network architectures. Hence in order to effectively automate machine learning, one would need to have a better understanding on the properties of this common machine learning feature.

1.4.1 Why Study Shortcut Connections?

Many advancements in machine learning can be found in *Computer Vision* (CV); and this is because that CV is one of the most well-studied topics in machine learning. In the past decade and in a near-simultaneous fashion, the introduction of *Residual Networks* (ResNets) [He et al., 2016a] and *Highway Networks* [Srivastava et al., 2015] have dramatically changed the practice of architectural designs. Both architectures employ shortcut connections to enable the training of deep networks over a thousand layers. This is a major milestone in machine learning as there is a plethora amount of empirical evidence showing that neural networks perform better when there are more parameters [Szegedy et al., 2015]. The significance of this milestone can be furthermore elaborated considering that the previous 2012 breakthrough in AlexNet [Krizhevsky et al., 2012] only consisted of eight layers.

Furthermore, CV tasks have long been studying in parallel with other important machine learning tasks such as *Natural Language Processing* (NLP). Hence, significant findings in CV tasks usually benefit NLP tasks, and vice versa; and this was also the case for shortcut connections. For one instance, the notable *Transformer* architecture [Vaswani et al., 2017]⁴ in NLP employs shortcut connections to ease the training for both of its encoders and decoders. For another instance, *Temporal Convolutional Network* (TCN) [Bai et al., 2018]⁵, which has demonstrated impressive temporal modelling power without recurrent units employs shortcut connections within each of its block-wise layers. This shows that by gaining a better understanding on the shortcut connection, we will gain insights that are applicable to both CV and non-CV tasks.

An extended comparison between ResNets and Highway Networks can be found in Section 2.1 of our preliminaries chapter. The same section will also discuss previous difficulties in building deep networks and motivate for the need to adapt shortcut connections into neural network models.

⁴Transformer [Vaswani et al., 2017] is a popular neural network algorithm that heavily relies on the *attention mechanism* [Bahdanau et al., 2015]. Many important works in NLP, including *Bidirectional Encoder Representations from Transformers* (BERT) [Devlin et al., 2019] and *Generative Pre-trained Transformer 3* (GPT-3) [Brown et al., 2020], are based on the Transformer model. Shortcut connections can be found in Transformers separately after the executions of multi-head attention and position-wise feedforward transformation.

⁵Inspired by Google's *WaveNet* architecture [Van den Oord et al., 2016], Bai et al. [2018] empirically evaluated the performance of a CNN-only network on an array of sequential tasks. TCN demonstrated impressive results (especially regarding speedup) for DNA/RNA nanopore basecalling [Zeng et al., 2020]. The network is based on multiple layers of dilational 1D-CNN block-modules with ResNet-like architectures which employ shortcut connections.

1.4.2 Studying Shortcut Connections via RNNs via Dynamical Systems

Of the two mentioned networks, the shortcut connection of Highway Networks was inspired by the *Long Short-Term Memory* (LSTM) [Hochreiter and Schmidhuber, 1997b] RNN for its ability to withstand the *Gradient Vanishing Problem* [Pascanu et al., 2013]. It is hence necessary to study LSTMs for understanding shortcut connections. The gradient vanishing problem will be elaborated in Section 2.1 of the preliminaries chapter, and more about the LSTM RNN will be discussed in detail in Section 3.1 of the background chapter.

LSTM RNN is one of the most versatile and extensively applied neural network modules. To name a few applications, it can be found in language modelling [Merity et al., 2018], image generation [Gregor et al., 2015], and meta-learning [Andrychowicz et al., 2016]. However, LSTMs are also intricate modules that consist of a lot of internal components that do not bear immediately clear meanings. Most existing research in LSTMs regard their applications instead of understanding their internal properties. Among those studies that aim to increase the architectural transparency of LSTMs, most of them analyse *what* components contribute to the success of LSTMs instead of *why* these components are important.

Existing literature in RNN understanding can mainly be classified as one of the following three categories: ablation study, search study, and observation study. Some examples include the work of Greff et al. [2016] in which they compared eight ablated versions of the LSTM; Jozefowicz et al. [2015] randomly mutated operators and evolved the network in attempt of searching for a better model; and that Karpathy et al. [2016] presented visualisations to observe the dependencies of LSTM outputs to their respective inputs. All of these studies have their own merits (see elaborations in Section 3.1.6); however, these studies do not answer the fundamental question

What purpose does each LSTM component serve?

Since we study LSTMs to forward our understanding on shortcut connections, it is then especially important for us to find ways which allow us to attach meaning to individual LSTM components in a first principle's manner. To achieve this, we propose to study LSTM components in a theoretical mathematical point of view.

Dynamical systems [Glendinning, 1994] study the geometric relationship between space and time. This field of mathematics is most commonly applied in physics, and can be found in the movements of celestial bodies [Neishtadt, 1984] and fluid dynamics [Rusak et al., 2012]. It is also found in biology for describing binding sites of enzymes [Keener and Sneyd, 1998] and in chemistry for narrating particle collisions [Sneyd et al., 2017]. Dynamical systems provide a rich set of vocabularies to describe qualitative changes in variables as the system unfolds in time. Using dynamical systems to describe the LSTM RNN will help us to:

- (i) gain an understanding on how shortcut connections propagate neural values;
- (ii) modify LSTM RNNs to improve meta-learning; and
- (iii) in turn leverage meta-learning for continual learning.

We will discuss some visualisation techniques for dynamical systems in Section 2.3 of the preliminaries chapter; and later introduce our novel contribution for understanding and simplifying RNNs in Chapter 4.

1.5 Outline, Contributions, Publications, and Codes

The rest of this thesis is organised as Figure 1.3 below.

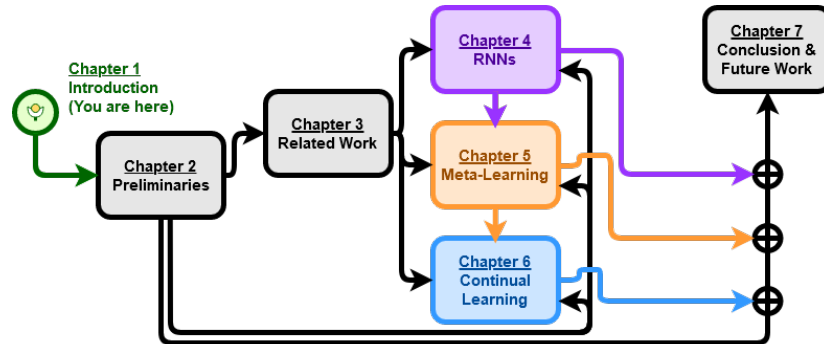


Figure 1.3: The Structure of this Thesis

1.5.1 Chapter 2: Preliminaries

This chapter introduces three fundamental machine learning aspects – residual networks & highway networks, loss landscapes & loss minima, and some fundamentals on dynamical systems. We would like to emphasise that these are not the main topics of our research. Instead, the covered materials therein that chapter only serve as the preliminaries for understanding our novel work. For instance, our work on RNNs will frequently refer to residual networks and highway networks. For another instance, both our work on meta-learning and continual learning concern the procedure of updating a target network; and hence the concepts of loss landscape and loss minima becomes useful for visualising the difficulties in training a network.

A second purpose of this chapter is to share some important adjacent ideas of which are related to our preliminaries. We also tried our best to provide our readers with a broader view on how some machine learning perspectives have changed, or have been utilised differently, over the past decade.

1.5.2 Chapter 3: Related Work

This is the chapter where we cover a broad discussion on RNNs, meta-learning, and continual learning. A portion of this chapter on meta-learning and continual learning will discuss and compare the main archetypes in their respective fields. In this chapter, we will also provide a detail description on the architectural designs of the LSTM RNN [Hochreiter and Schmidhuber, 1997b], the experimental setup of meta-learning for learning to optimise [Andrychowicz et al., 2016], and on the modifications therein the backward pass of continual learning [Lopez-Paz and Ranzato, 2017]. Our novel work, however, will be described over the next three chapters.

1.5.3 Chapter 4: Understanding RNNs with Dynamical Systems

This chapter is based on

Paper 1 [Kuo et al., 2018]

Kuo, N.I., Harandi, M., Suominen, H., Fourrier, N., Walder, C., & Ferraro, G. (2018). **DecayNet: A Study on the Cell States of Long Short Term Memories**. *An archived paper*.

Paper 2 [Kuo et al., 2020b]

Kuo, N.I., Harandi, M., Fourrier, N., Walder, C., Ferraro, G., & Suominen, H. (2020). **An Input Residual Connection for Simplifying Gated Recurrent Neural Networks**. In the *International Joint Conference on Neural Networks*.

and our complementary codes can be found in

https://github.com/Nic5472K/IJCNN2020_IRC.

We will start with the content in Paper 1 where we analyse LSTM RNNs as the discretised dynamical systems of difference equations. From a difference equation perspective, we aim to understand the *atomic* structure of the LSTM cell. That is, for describing the individual behaviours of the internal LSTM components to further our understanding on how the RNN stores and updates its representations of occurred events. We will show that the LSTM memory is created via a coupled dynamical behaviour with two competing objectives. We refer to this as the *Catch & Release* phenomenon and discuss the recursive properties of memory deletions versus memory amendments.

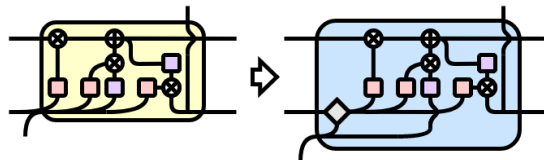


Figure 1.4: Modification on LSTM

In contrast to the *atomic* approach in Paper 1, Paper 2 will take a *schematic* approach and describe LSTM as the continuous dynamical system of differential equations. The schematic approach no longer considers the isolated behaviour of a single neuron in the network; instead, it looks at the systematic behaviour of the network from a macroscopic point of view. We will establish links between a class of state of the art RNNs (including LSTMs) with dynamic Hopfield networks [Hopfield, 1982; Sussillo and Abbott, 2009]. Using this connection, we will identify a section of LSTM components that are secondary in importance and empirically demonstrate that by removing these redundant gadgets, not only can we simplify the network but also increase LSTM's ability to generalise on unseen test data. We will show that in the most extreme case, we can remove up to 85.4% parameters from LSTM for image generation [Gregor et al., 2015] to yield a small and properly functioning model. Furthermore, this fashion of network simplification is generally applicable to several de-facto standard RNNs including *Gated Recurrent Units* (GRUs) [Cho et al., 2014] and *Simple Recurrent Units* (SRUs) [Lei et al., 2018a]. A concept art of our network modification is depicted in Figure 1.4.

1.5.4 Chapter 5: Learning to Learn beyond Batch-Efficient Updates

This chapter is based on

Paper 3 [Kuo et al., 2020c]

Kuo, N.I., Harandi, M., Fourrier, N., Walder, C., Ferraro, G., & Suominen, H. (2020) **MTL2L: A Context Aware Neural Optimiser**. In the *Workshop on Automated Machine Learning of the International Conference on Machine Learning*.

Paper 4 [Kuo et al., 2020a]

Kuo, N.I., Harandi, M., Fourrier, N., Walder, C., Ferraro, G., & Suominen, H. (2020). **M2SGD: Learning to Learn Important Weights**. In the *Workshop on Continual Learning in Computer Vision of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

and our complementary codes can be found in

https://github.com/Nic5472K/AutoML2020_ICML_MTL2L and

https://github.com/Nic5472K/CLVISION2020_CVPR_M2SGD.

This chapter documents our research findings on meta-learning which are based on our insights gained from studying RNNs in Chapter 4. We propose extensions to the optimisation-based approach of *Learning to Learn* (L2L) [Andrychowicz et al., 2016] and its concept art is depicted in Figure 1.5.

The chapter begins with the research outcome of Paper 3. Previously, Andrychowicz et al. [2016] have demonstrated that by replacing handcrafted gradient descent algorithms with LSTM, the LSTM-optimiser could learn a set of update rules that configured another

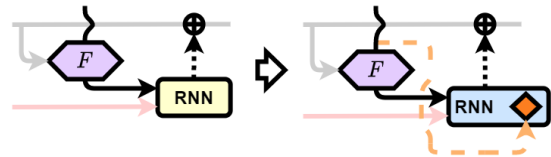


Figure 1.5: Modification to L2L

neural network to lower losses in less iterations. However, such learnt update rules were limited to updating *one* model architecture on *one* dataset. We will address this input-domain heterogeneity and extend a learnt LSTM-optimiser’s capability to update on unseen datasets thereby making LSTM-optimisers more generally applicable.

Later in this chapter, we will discuss our work in Paper 4 based on *Meta-SGD* [Li et al., 2017], of which is itself an extended method of L2L. Our idea is to modify Meta-SGD beyond the purpose of batch-efficient machine learning for network pruning. One basic principle of network pruning is to remove unimportant weights to achieve better generalisation and to improve processing speed. However, traditional network pruning is usually done in an iterative fashion [Han et al., 2015], and hence it is a time consuming task. Our contribution lies in connecting L2L to network pruning to create a learnt optimisation algorithm which not only learns to configure the weights of a network but also become aware of the underlying redundancies within the target network. Our novel algorithm is capable of updating ResNet20s while removing 83.71% parameters in a fast and non-recurrent fashion.

1.5.5 Chapter 6: Topics on Simple and Efficient Continual Learning

This chapter is based on

Paper 5 [Kuo et al., 2021b]

Kuo, N.I., Harandi, M., Fourrier, N., Walder, C., Ferraro, G., & Suominen, H. (2021) **Learning to Continually Learn Rapidly from Few and Noisy Data**. In the *Meta-Learning and Co-Hosted Competition of the AAAI Conference on Artificial Intelligence*.

Paper 6 [Kuo et al., 2021a]

Kuo, N.I., Harandi, M., Fourrier, N., Walder, C., Ferraro, G., & Suominen, H. (2021). **Plastic and Stable Gated Classifiers for Continual Learning**. In the *Workshop on Continual Learning in Computer Vision of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

and our complementary codes can be found in

https://github.com/Nic5472K/MetaLearnCC2021_AAAI_MetaSGD-CL and

https://github.com/Nic5472K/CLVISION2021_CVPR_HCN.

Continual learning concerns the sequential acquisition of skills. Or more specifically, it aims to prevent a trained network from forgetting how to solve an old task when the trained network learns a new skill. This remains a difficult problem for neural networks as parametric updates (see Equation (1.4)) modify stored knowledge in the network weights [McCloskey and Cohen, 1989]. This is thus a major hurdle towards achieving general intelligence [Legg and Hutter, 2007].

The first half of this chapter introduces an L2L meta-learner that sequentially configure a network for the sequential acquisition of skills. By leveraging meta-learning for continual learning, we surpass previous state of the art results on continual learning and also utilise the meta-learning advantage of rapid knowledge acquisition to robustly train a network under a highly noisy environment. The concept is presented in Figure 1.6, where L2L updates a backbone from one class of weights for solving a single task to an adequate region in the parametric space that intersects two of classes of family of weights for solving (both the old and the new) tasks⁶.

The second half of this chapter explores continual learning from the perspective of shortcut connections and RNNs. We present a neural network architecture that makes it inherently difficult to forget learnt knowledge. Our novel design can be easily employed in conjunction with the three main archetypes of continual learning approaches (and more); and it also enhances the robustness of all these techniques.

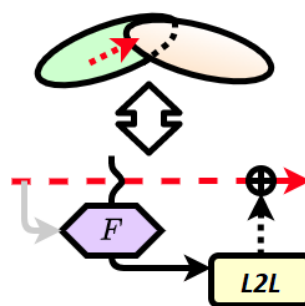


Figure 1.6:
L2L for Continual Learning

⁶The inspiration of our concept art is influenced by Figure 1 in page 3 of Kirkpatrick et al. [2017].

1.5.6 Chapter 7: Thesis Conclusion

We end this thesis with a conclusion chapter. Therein that chapter, we will summarise the contribution of this thesis, and we will also mention some related popular research that are different to our own study. We emphasise that we will cover those topics specifically for their orthogonality to our work. We will highlight the importance of those work; but at the same, we will state why our approach and the analytical framework of this thesis is equivalently important. Then we will show in our proposal for future directions how these different perspectives could possibly supplement each other.

Preliminaries: Past and Present on Deep Models

This is the preliminary chapter of the thesis. Unlike the related work (Chapter 3) where we provide an in-depth analysis on the background on the focus of our studies, this chapter aims to cover the following aspects:

- (i) necessary prerequisites for understanding our novel work and
- (ii) some interesting but loosely related ideas explored by the community.

The 3 prerequisites of this chapter are

Section 2.1: Residual Networks & Highway Networks,

Section 2.2: Loss Landscape & Loss Minima, and

Section 2.3: Some Fundamentals on Dynamical Systems.

For those readers who are already familiar with the preliminaries, we suggest them to go directly to the related ideas section of this chapter¹. Note, since there is already an abundance of work related to these three areas of study, we will mainly narrate the related ideas rather than delving into them in depth².

Section 2.1 will address two types of shortcut connections that are typically found in CNNs. In the related ideas, we will provide some brief discussions on why deep neural networks are needed. Furthermore, we will discuss how general trends in neural network practices have changed and mention some (older) techniques that enable deep learning.

Section 2.2 will address useful visualisations that describe the procedure of network updates to conceptualise some difficulties in training neural networks. In the related ideas, we will discuss how these visualisation techniques were first applied.

Section 2.3 will provide a brief introduction to the mathematical studies of dynamical systems (also see Section 1.4.2). They will allow us to define variables in a first principle's fashion and associate meaning towards individual (or sections of) neurons in a large network.

To increase clarification, we will provide additional *in-chapter clarifications* (ICCs) for this chapter in Appendix B. Places with additional clarification will be noted with the superscript ^{ICC}.

¹The related ideas can be found in Sections 2.1.2, 2.2.4, and 2.3.5, respectively.

²Also not all related ideas directly feeds into the novel work conducted by the student-author.

*The following section discusses shortcut connections.
This is an old technique from 1998 and
was adopted by the community in 2015 and 2016a for deep learning.
It allowed real deep learning over a thousand layers and
has ever since changed the way how we think of network modelling.*

Section Content:

- Section 2.1.1: Training Deep Network is Difficult,
- Section 2.1.2: Related Ideas: Why Do We Need Deep Learning?
- Section 2.1.3: Prerequisite 1 Conclusion.

2.1 Residual Networks and Highway Networks

This section will attend to the two canonical feedforward designs of *Residual Connections* [He et al., 2016a], and *Highway Connections* [Srivastava et al., 2015].

2.1.1 Training Deep Networks is Difficult

Table 2.1: Forward and Backward Passes of Residual and Highway Connections

Residual Connection (see He et al. [2016a] for the forward and Ling and Qiu [2019] for the backward pass)	
Forward pass	$\mathbf{n}_{L_\gamma} = \mathbf{n}_{L_{(\gamma-1)}} + \mathcal{H}(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})$ (2.1)
Backward pass	$J = \frac{\partial \mathbf{n}_{L_\gamma}}{\partial \mathbf{n}_{L_{(\gamma-1)}}} = (\mathbb{I} + \mathbf{D}^\gamma \mathbf{W}_{L_\gamma})$ (2.2) with diagonal \mathbf{D}^γ such that $\mathbf{D}_{(r,r)}^\gamma = \mathcal{H}'(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})_{(r)}$
Highway Connection (see Srivastava et al. [2015] for both the forward and backward passes)	
Forward pass	$\mathbf{k}_{L_\gamma} = (1 - \mathbf{T}_{L_\gamma}) \odot \mathbf{k}_{L_{(\gamma-1)}} + \mathbf{T}_{L_\gamma} \odot \mathcal{H}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})$ (2.3) with gated unit: $\mathbf{T}_{L_\gamma} = \mathbf{T}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{Q}_{L_\gamma})$
Backward pass	$J = \frac{\partial \mathbf{k}_{L_\gamma}}{\partial \mathbf{k}_{L_{(\gamma-1)}}} = \begin{cases} \mathbb{I}, & \text{if } \mathbf{T}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{Q}_{L_\gamma}) = \mathbf{0} \\ \mathcal{H}'(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma}), & \text{if } \mathbf{T}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{Q}_{L_\gamma}) = \mathbf{1} \end{cases}$ (2.4)

Conventional wisdom tells us that the depth of a neural network is crucial to its performance [Szegedy et al., 2015]. However, deep neural networks are infamously difficult to train due to the *Gradient Vanishing* [Pascanu et al., 2013] problem.

Consider an MLP F with the layer transformations of Equation (1.1) where

$$\mathbf{a}_{L_\gamma} = \mathcal{H}(\mathbf{a}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma}).$$

Adapting from Equation (1.4), the weight update of network F at layer γ is

$$\mathbf{W}_{L_\gamma}(t+1) = \mathbf{W}_{L_\gamma}(t) - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{L_\gamma}}$$

at time step t and by chain rule we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{L_\gamma}} = \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}} \prod_{\zeta=\gamma}^{\Gamma-1} \frac{\partial \mathbf{a}_{L_{(\zeta+1)}}}{\partial \mathbf{a}_{L_\zeta}} \frac{\partial \mathbf{a}_{L_\gamma}}{\partial \mathbf{W}_{L_\gamma}} \quad (2.5)$$

$$= \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{y}}} \prod_{\zeta=\gamma}^{\Gamma-1} \mathcal{H}'(\mathbf{a}_{L_\zeta}, \mathbf{W}_{L_{(\zeta+1)}}) \frac{\partial \mathcal{H}(\mathbf{a}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})}{\partial \mathbf{W}_{L_\gamma}}. \quad (2.6)$$

Hence the overall update of $\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{L_\gamma}}$ will quickly diminish to 0 as the total amount of layers³ increase and that individual $|\mathcal{H}'(\mathbf{a}_{L_\zeta}, \mathbf{W}_{L_{(\zeta+1)}})| < 1$.

In order to withstand the gradient vanishing problem, the community introduced a shortcut connection (or information highway) towards Equation (1.1). The modified

³Either with the direct increase in Γ or with the increasing difference between $(\Gamma - 1) - \gamma$.

forward passes can be found in Equations (2.1) and (2.3) in Table 2.1, and they are illustrated in Figure 2.1 on the right.

Residual connections can be seen as a special (or simplified) version of highway connections and both are theoretically robust against the vanishing gradient problem. This is because as stated in Table 2.1, both forward passes directly introduce identity matrix \mathbf{I} to Jacobian J in their backward passes and thus prevent the undesirable condition of $\|\mathcal{H}'(\mathbf{a}_{L_{\zeta}}, \mathbf{W}_{L_{(\zeta+1)}})\| < 1$.

Though residual connections and highway connections had only been introduced in 2015, they were not completely novel. The so-called gated unit in Equation (2.3) was taken from the LSTM RNN [Hochreiter and Schmidhuber, 1997b]^{ICCO1} and that the practice of shortcut connections can be dated back to the work of Schraudolph [1998]^{ICCO2}. Overall, the success of residual connections and highway connections can be attributed to the advances in computational power and also to a dedicated community that has continued to carefully verify the applications of existing ideas.

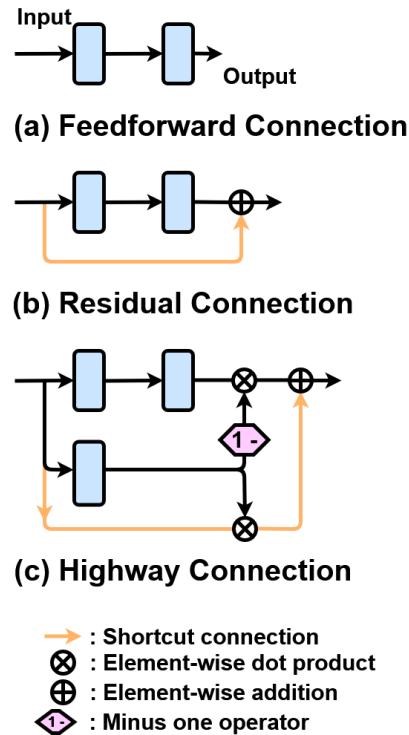


Figure 2.1: Three Connections

2.1.2 Related Ideas: Why Do We Need Deep Learning?

Residual networks and highway networks enable the training of deep learning models with more than 1000 layers. However and prior to their introduction, the 8-layer AlexNet of Krizhevsky et al. [2012] was once considered deep and that deep models were hard to train. So why did the community continue to develop deeper models instead of just using shallow ones? In order to answer this question, we will need to discuss some shortcomings of shallow models. Overleaf, we will follow the work of Bengio et al. [2013]⁴ and discuss kernel-based methods and decision trees.

2.1.2.1 Kernel-Based Methods

Kernel-based methods such as *Locally Linear Embedding* [Roweis and Saul, 2000], *Isomap* [Tenenbaum et al., 2000], and *Kernel Principle Component Analysis (PCA)* [Bengio et al., 2004] are bounded to the *locality of learning* [Bengio and Monperrus, 2005]. This fundamental weakness is a direct consequence on how kernel methods apportion the high dimensional feature space based on the leading eigenvalues and eigenvectors.

⁴Mainly adapted from page 3 under Section 3.2 “Smoothness and the Curse of Dimensionality” of Bengio et al. [2013].

These methods construct a tangent plane centred about a point x spanned by the leading eigenvectors of which can be closely approximated by a *linear combination* of the difference vectors $(x - x_j)$ for x_j near x . It is thus difficult for such techniques to generalise well to new data points that are *far* away from the centre point – with *far* used in the context of the amount of noises, the degrees of curvature of the manifold, and the manifold’s dimensionality. Noises displace data off the manifold thus causing inferential difficulties; whereas curvature and dimensionality make it difficult to create a manifold that captures all characters in the features space.

2.1.2.2 Decision Trees

Decision trees [Breiman et al., 1984] are models that explores sequential consequences of different decisions. A tree describes features of the entirety of the input space by internal decision nodes which lead to constant output leaves. The nodes (and their subsequent nodes) hence partition the input space and each partitioned region can be seen as if they were dictated by a piece-wise constant function.

Decision trees have shown to degrade in generalisation when the amount of nodes (descriptive functions) increase [Vilalta et al., 1997]. Bengio et al. [2010] stated that, like kernel-based methods, decision trees do not generalise well to unseen variations because of how they were built. They argued that this fundamental weakness stems from the node-branching protocol that governs the topology of the tree. To elaborate, the nodes partition the input space and in order for the partition to work for the testing dataset, the testing dataset must have at least one example in all partitioned regions of which were originally spawned from the features therein the training dataset.

2.1.2.3 Are Neural Networks Over-Fitted Models?

In the last two sections, we discussed why kernel-based methods and decision trees do not generalise well to unseen combinations of data variations. But how do neural networks overcome this difficulty? Aren’t they over-parameterised? Shouldn’t they hence overfit on the training data? As a simple response: neural networks are only considered over-parameterised if you also considered them as statistical models.

But are neural networks really statistical models? Undeniably, many advances in neural networks were based on *traditional* statistical principles. For instance, *Variational Autoencoders* [Kingma and Welling, 2014] were based on the idea of auto-encoding variational Bayes [Fox and Roberts, 2012]. However, it is becoming more apparent that neural networks exhibit several behaviours that deviate from *classic* statistical models. For example, deep neural networks can achieve super-human level of accuracy on image classification not only on its training dataset but also on its testing dataset [Szegedy et al., 2015]. Deep models are also known to exhibit the interesting *Double Descent* phenomenon [Nakkiran et al., 2020] – such that during training, the testing loss decreases then increases but eventually decreases again^{ICC03}. This phenomenon hence challenges the conventional statistical wisdom of *larger model hurts*

(formally known as the *Bias-Variance Trade Off* [Hastie et al., 2009]).

Deep networks have also shown to be more invariant towards translations and rotations in input data [Goodfellow et al., 2009] and possess exponentially greater expressiveness as depth increases [Delalleau and Bengio, 2011; Bengio and Delalleau, 2011]. Hence though it may be true that neural networks perform poorly as according to complexity measurements such as Akaike’s Information Criterion [Akaike, 1974] and that they do not encode information with the shortest description length [Shannon, 2001], but it should also be noted that neural networks outperform *traditional* statistical inferential tools on many important real life applications such as mortality prediction [Che et al., 2018].

2.1.2.4 So Why do Neural Networks Need so Many Parameters?

One commonly accepted explanation to why deep neural networks offer high generalisation is based on their inferential process. A conventional layer of a neural network consists a number of input nodes and a number of output nodes. Each input node is *fully connected* to and concurrently processed by all output nodes. This hence allows a neural network to encode *distributed features* among its many synapses.

Distributed features *reuse* a large set of common operations in the shallower layers. In the work of Montufar et al. [2014], the authors noted that deep *linear* neural networks⁵ have two special properties. First, the layer-wise composition of functions enable subsequent layers (deeper layers) to reuse low-level computations in a near-exponential fashion. Second, deep linear neural networks are able to sequentially map *portions* of each layer’s input-space to the same output. Combining these two factors, the subsequent layers thus significantly narrow down the vast amount of features generated in shallower layers. This many-to-1 relationship between intermediate layer functions, also known as *space folding*, is believed to create an implicit bottleneck which increases network generalisation while preventing over-fitting.

Besides examining the amount of individual weights, it is also important to analyse the naturally occurring activation sparsity when studying neural networks. The recent work of Kurtz et al. [2020] found that layer-wise activation sparsity decreases as depth increases. The findings of this work fits well with space folding, and the decrease in the effective amount of activated neurons serves as an alternative explanation to why neural networks do not over-fit on their training data and thus they generalise well to unseen data.

2.1.2.5 Complementary Studies to Architectural Modifications

Examples of Early Studies

As discussed in Section 2.1.1, residual connections and highway connections are shortcut connections that withstand the gradient vanishing problem. Before redis-

⁵By deep *linear* neural networks, we meant those neural networks with ReLU activation functions [Nair and Hinton, 2010] instead of sigmoid nor the hyperbolic tangent function.

covering the effectiveness of shortcut connections in modern machine learning technologies, multiple sub-disciplines were concurrently studied to enable deep learning.

One of the earliest breakthroughs in training a deep network was conducted by [Hinton et al. \[2006\]](#) to efficiently train the older-generation neural network of *Deep Belief Networks* (DBNs) [[Hinton, 2009](#)]. Since weight modifications in lower layers impact those in deeper layers, [Hinton et al. \[2006\]](#) proposed to train DBN parameters via greedy layer-wise optimisation with unsupervised learning. A supervised version of this work can be found in [Bengio et al. \[2007\]](#). With the advances in neural network research, these strategies are now less applied in modern machine learning.

Examples of More Recent Studies

In addition to greedy layer-wise training, more recent work looked at planned training in multiple stages of network updates. The work of [Weston et al. \[2012\]](#), [Cho et al. \[2013\]](#), and [Lee et al. \[2015\]](#) added supervision to the intermediate layers of deep networks to assist their training. Some popular alternatives are *Curriculum Learning* strategies [[Bengio, 2009](#)] and *Knowledge Distillation* [[Hinton et al., 2015](#)]. Curriculum learning provides a learning agent with tasks of a changing training distribution by first solving easy problems and then on more difficult problems. Whereas knowledge distillation transfers the knowledge in an ensemble of deep networks to a single shallow network via matching the logits of the softmax layer.

Examples of Timely Studies: Weight Initialisation

Network initialisation is one important study that is frequently practised. It is also common to see network initialisation employed in conjunction with shortcut connections. This study was first proposed to understand why standard gradient descent from random initialisation^{ICC04} was doing so poorly with deep neural networks. In the seminal paper of [Glorot and Bengio \[2010\]](#), the authors showed that randomly initialised deep neural networks caused logistic sigmoid activation to saturate quickly. This was partially due to the topological outline of the sigmoid function^{ICC05} and also because that the Jacobian associated with each layer could sway far from 1. Regarding the former issue, the authors suggested the use of hyperbolic tangent functions and softsign functions for learning. As for the latter issue, the author derived their normalised initialisation (or more commonly known as the Xavier initialisation^{ICC06}) to keep equal magnitudes for all input variances.

[He et al. \[2015\]](#) extended the study of [Glorot and Bengio](#) and proposed an alternative initialisation (better known as Kaiming initialisation^{ICC07}). Prior to the re-introduction of shortcut connections, Kaiming initialisation enabled the successful training of 30-layer deep linear neural networks of which was previously unachieved with Xavier initialisation.

Examples of Timely Studies: Normalisation

Another important pillar that continues to support deep learning is the study of normalisations. *Batch Normalisation* (BN) [[Ioffe and Szegedy, 2015a](#)] is arguably the most popular normalisation technique; and it was also employed in default in the

original residual connection paper [He et al., 2016a]. BN provides several benefits – it allows practitioners to use larger learning weights and be less careful about network initialisation, and it also acts as an explicit regularisation which in some cases eliminates the need for implicit regularisation such as *Dropout* [Srivastava et al., 2014]. These benefits were achieved by normalising each individual hidden feature by their statistics in the mini-batch during training^{ICC08}.

The effectiveness of BN was initially attributed to its ability to remedy *Internal Covariate Shift* (ICS) [Shimodaira, 2000]. ICS refers to the difficulty in training a multi-layer network when subtle changes in shallow-layer activation values cause serious changes to the distributions of those activation values in deeper layers. This occurs because that the changes in shallow layers could regurgitate and be amplified over the many intermediate-layer transformations.

It was later shown in Santurkar et al. [2018] that BN is indeed beneficial for deep learning but BN may not be reducing the shift at all. Instead, the authors found that BN invoked implicit modifications in the backwards pass and reduced the effective learning rate of higher order terms in the derivative of the loss functions^{ICC09}, and that it potentially encouraged the training process to converge to flatter minima. The concept of minima will later be explored in Section 2.2.

BN is however not without its own problems. Its dependency on mini-batches meant that it is non-obvious for how BN could be applied on recurrent connections; this dependency also meant that BN performances vary as according to the choices of the batch size. Large batch sizes will increase the amount of computational resources needed for calculating the batch statistics in each iterative update, whereas small batch sizes can cause the statistics to be highly influenced by outlier features. The issue of large batch sizes will later be elaborated in Section 2.2.4.1.

Two alternative normalisation designs were proposed to overcome this issue. On one hand, *Weight Normalisation* (WN) [Salimans and Kingma, 2016] proposed to normalise the weights^{ICC10} instead of the feature values. This direct modification towards the weight values do not require feature information from the batch. When compared to BN, WN can be applied to RNNs and that it provides a slight speed-up. It should be noted that by normalising the weights, the neural features in deeper layers are still theoretically susceptible to small changes in shallow layers^{ICC11}. However and as we previously mentioned, Santurkar et al. found that ICS is not a real problem for deep learning and it is instead the implicit modifications in the backward pass which causes normalisation techniques to be successful. In their own paper, Salimans and Kingma also validated Santurkar et al.'s claim and found that WN rescaled the gradient and projected the weight such that the gradient covariance matrix became closer to identity of which was beneficial for optimisation.

The second improved variant of BN is *Layer Normalisation* (LN) [Ba et al., 2016b]. LN was proposed to stabilise the hidden state dynamics in recurrent neural networks and it could substantially reduce the RNN training time. Like WN, LN normalises the features without batch-wise statistics. Instead of normalising every individual feature as according to their batch statistics, all features in one layer are normalised by the aggregated features of the same layer^{ICC12}. LN was later analysed in detail

by Xu et al. [2019]; and again, the effectiveness of this normalisation could also be attributed to the implicit modifications made in the backward pass^{ICC13}.

Two more noteworthy normalisations that extended BN are *Instance Normalisation* (IN) [Ulyanov et al., 2016] and *Group Normalisation* (GN) [Wu and He, 2018]. On one hand, IN was proposed to improve the quality of transferring the style of an image to another [Gatys et al., 2016]. The main difference between BN and IN is that it removed the batch-wise aggregation of feature statistics^{ICC14} (like that in LN but for CNN instead of RNN), and this resulted in better stylisation by ignoring contrast of the content image. On the other hand and as its name implies, GN groups several features of the same layer as a set and normalise the set with their corresponding aggregated set-wise features. Hence similar to LN, GN normalises parts of its features with other features therein the same layer. However at the same time, GN does not normalise each of its underlying features with all features therein the same layer, this allows each individual feature to keep a certain amount of its individual characteristics⁶.

2.1.3 Section Conclusion

In this prerequisite, we mentioned that residual connections [He et al., 2016a] and highway connections [Srivastava et al., 2015] are the main architectures that allow for successful deep learning. In Section 2.1.1, we also showed that these shortcut connections are mainly successful because of how they change the network updates in the backwards pass. Similarly, we also mentioned how the success of normalisation techniques were also attributed to their implicit modifications to the backward pass. Such modifications in the backwards pass will later be discussed in Chapter 6 when we study continual learning. Furthermore, while the purpose of shortcut connections are clear from the perspective of the backwards pass, their corresponding modifications in the forwards pass are less clear. For instance, what does the gated unit in Equation (2.3) do during the network inferential process? This will be discussed in detail in Chapter 4 when we study state-of-the-art RNNs.

⁶Interested readers may also want to visit Figure 2 in page 3 of Wu and He [2018]. In that figure, the authors of GN presented very nice illustrations to highlight the differences among BN, LN, IN, and GN.

*The neural network inferential process and their update procedure are dynamic in nature.
The following section discusses loss landscape & loss minima,
and we explore visualisation techniques that describes
morphological modifications that occur in the network backwards pass.*

Section Content:

- Section [2.2.1](#): Loss Landscape,
- Section [2.2.2](#): Visualising Parametric Updates,
- Section [2.2.3](#): Good Loss, Bad Loss?
- Section [2.2.4](#): Related Ideas: How was the Loss Minima Studied?
- Section [2.2.5](#): Prerequisite 2 Conclusion.

2.2 Loss Landscape & Loss Minima

Neural network parameters are updated in a high dimensional space. Loss landscapes visualise the parametric space in three dimensions; and though simple, they provide intuitive depictions of important concepts such as *Loss Minima* [Hochreiter and Schmidhuber, 1997a] to capture the substantive difficulties in network training.

2.2.1 Loss Landscape

An example loss landscape is presented in Figure 2.2. The (x, y) -dimensions represent two arbitrary weights $\theta_{\{1\}}$ and $\theta_{\{2\}}$ of the parameters of a network \mathcal{F} ; whereas the z -direction is the loss \mathcal{L} invoked by network \mathcal{F} when parameterised with $\theta_{\{1\}}$ and $\theta_{\{2\}}$. The landscape is hence a terrain which describes network performances with candidate weights.

The terrain can be scattered with hills and valleys. The former refers to contiguous regions of high loss values whereas the latter refers to contiguous regions of low loss values. Since each position on the landscape refers to one specific combination of parameters, the contiguity of nearby terrain represents a *class* of parameters. To elaborate, the aforementioned valleys, better known as loss minima, are classes of parameters that result in low loss values.

2.2.2 Visualising Parametric Updates

One easy way to conceptualise parametric *modification* on the loss landscape is to think of it as a marble placed on a table filled with pit holes. The marble represents the status of network configuration and the table corresponds to the loss landscape.

Figure 2.3 shows two scenarios of marble moments. The scenario in subplot (a) is caused when an external force is applied; like when someone purposely lifted the table. As a consequence of the applied force, the marble rushes up the hill but it eventually falls down into a valley. For this particular path, our target of optimisation (the marble) crosses a high loss region (the hill) before it finally settles in a low loss region (the valley). Since every update occurs at a discrete time, this particular path visits an unnecessary hill and is hence inefficient. We can do better with the gradient descent shown in subplot (b).

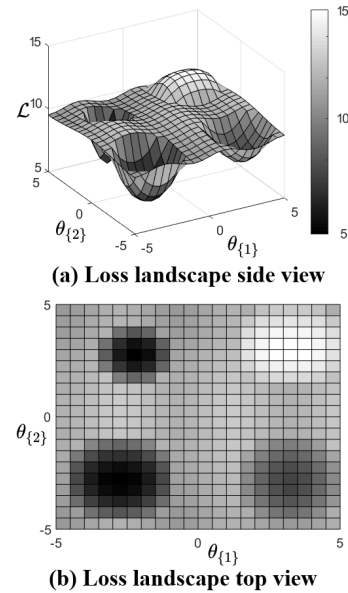


Figure 2.2: Sample Illustrations

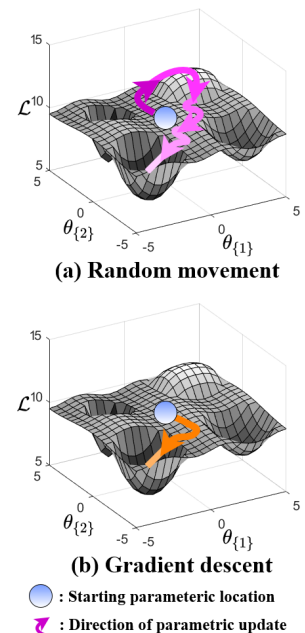


Figure 2.3: Parametric Update

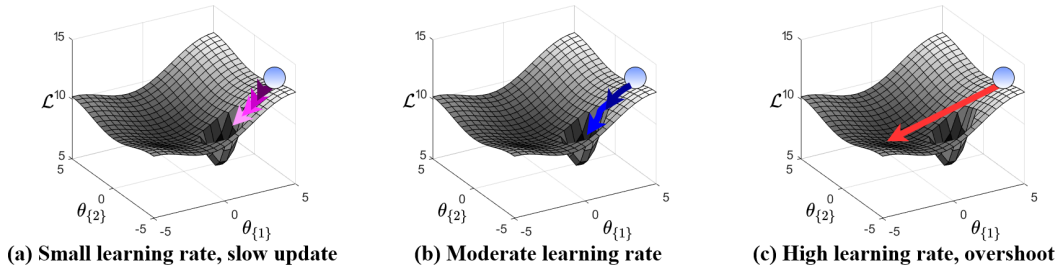


Figure 2.4: The Consequences of Different Learning Rates

Consider the gradient descent of Equation (1.4)

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}.$$

It consists the two terms of θ_t and $-\alpha \nabla_{\theta} \mathcal{L}$. The former is represented by the position of the marble whereas the latter (including the negative sign) is represented by the coloured arrows. In addition, the *trajectories* of the coloured arrows represent the iterative updates of the network subjected under the rules of gradient descent.

There are two elements in $-\alpha \nabla_{\theta} \mathcal{L}$. Element $-\nabla_{\theta} \mathcal{L}$ is the *direction* of the arrow and it points towards the lowest loss within the vicinity of the marble. Furthermore, α describes the increment of the update. Hence along with the magnitude $|\nabla_{\theta} \mathcal{L}|$, $\alpha |\nabla_{\theta} \mathcal{L}|$ would be the *total size of the stride*. A small loss enables the marble to carefully detect nearby terrain and slowly but surely drops into a minimum. Large α can parameterise a network quickly, but it can also overshoot and miss desirable minima. An elaboration of the effects of the sizes of learning rates is shown in Figure 2.4.

2.2.3 Good Loss, Bad Loss?

When multiple minima concurrently exist, *which minima is the most desirable?* We highlight the three minima of Figure 2.2 in blue in Figure 2.5. A minimum can be judged by several factors. The most obvious quality is its *lowest point* – the lower the loss, the better the generalisation⁷. In addition, the quality of a minimum can also be decided by its *accessibility* – on how easy and efficient it is to reach the minimum given the terrain around its neighbourhood.

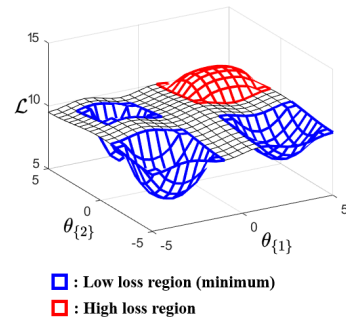


Figure 2.5: Three Minima

2.2.3.1 On the Depth of Minima

Figure 2.6 (a) elaborates the depth of the minima. We added two black lines both of loss $\mathcal{L} = 5$ below the loss landscape. The lines touch the lowest points of the two minima positioned left⁸ of $\theta_{\{1\}} = 0$ while they miss the minimum positioned right of $\theta_{\{1\}} = 0$. Updating towards the latter minimum would hence cause worse generalisation on the training data.

⁷At least on the training dataset.

⁸Centred at some $\theta_{\{1\}} = \rho$ where $\rho < 0$.

2.2.3.2 On the Width of Minima

Figure 2.6 (b) elaborates the width of the minima. We highlighted the two minima left of $\theta_{\{1\}} = 0$; though these two minima have the same depth, they differ in two types of widths. First, the widths of the mouths are different. Minima with larger mouths, *Wide Minima*, makes it easier for a marble to fall in. Second, the widths at the bottom of the valley are different. Minima with wider bottoms, *Flat Minima*, are harder for a marble to be forced out⁹. Flat minima are desirable, as leaving a minimum severely degrades network performances as we exit a low loss region.

2.2.3.3 Smoothness within a Minimum

It is also important to consider *the ease to continually update* the network once we enter the periphery of a minimum. In Figure 2.6 (c), we present an example minimum surrounded with a treacherous terrain. In the vicinity of the example minimum, we highlight those region with loss $\mathcal{L} < 10$ in cyan, and those those region with $\mathcal{L} > 10$ in magenta. Due to the inconvenient terrain, the marble can initially roll towards the vicinity, but it then requires more iterations before finally entering a low loss region.

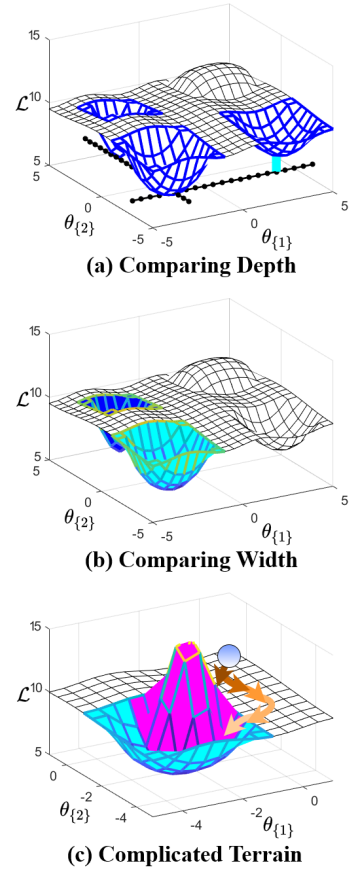


Figure 2.6: Minima Scenarios

2.2.3.4 What Constitutes the Outline of the Loss Landscape?

The loss landscape is the visualisation of loss \mathcal{L} about two arbitrary parameters $\theta_{\{1\}}$ and $\theta_{\{2\}}$. As stated in Equation 1.3, the loss of a batch is given as

$$\mathcal{L} = \frac{1}{|B|} \sum_{(x_i, y_i) \in B} \mathcal{C}(\hat{y}_i, y_i)$$

and hence the loss is dependent on several factors including, the choice of network architecture \mathcal{F} , the choice of cost function \mathcal{C} , the feature distribution therein the randomly sampled inputs of current batch $(x, y) \in B$, and also on batch size $|B|$. This same set of factors thus also contributes to the shape of the landscape.

The loss landscape is not static for batched supervise learning. Though the overall construct is determined by the choice of \mathcal{F} , \mathcal{C} , and $|B|$, random changes on the loss landscape can occur due to the stochasticity of the randomly sampled input x of batch B . This inconsistent nature makes optimisation difficult.

⁹For instance from those batches containing outliers data points.

Figure 2.7 elaborates this problem. Each of the coloured surfaces represents an *approximation*¹⁰ of the position where the loss minimum exists based on their respective mini-batches. If the inputs were sampled from a single distribution with guaranteed stationarity, all coloured surfaces would appear *roughly at the same spot*. That is, they should behave similarly to that in subplot (a), where the 3 surfaces mostly overlay each other. The solid dots on the subplot correspond to the centres of the minima; and the curved lines are the interceptions of different candidate minima.

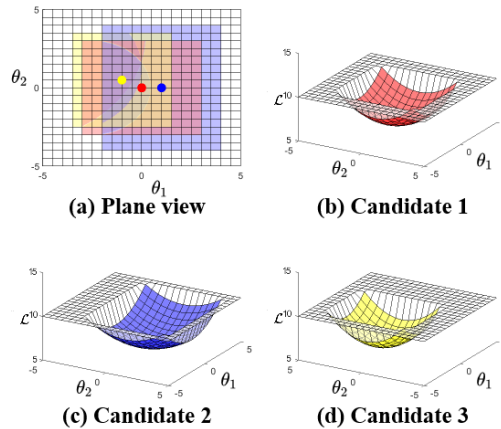


Figure 2.7: Minima are not Static

2.2.3.5 Limitations of the Loss Landscape Representation

In addition to its inherently random nature described in Section 2.2.3.4, loss landscape analyses are also made difficult by the network dimensionality. Neural networks are composed of thousands if not billions of weights, thus loss landscape alone is an ineffective way of understanding the true nature of updates because we can only gaze whether our system is behaving correctly in the high dimensional parametric space via pondering through a much lower three dimensional representation. Nonetheless, these simple and elegant conceptualisations are capable of reflecting much of the difficulties in network optimisation.

2.2.4 Related Ideas: How were the Loss Minima Studied?

One of the most important early studies in loss minima is the influential work by Hochreiter and Schmidhuber [1995]. This particular paper provided a rough (and imprecise) definition of the concept of flat minima; and it argued that flat minima were more advantageous than their sharp counterparts. However, they did not study loss minima for the purpose of understanding the topology on the loss landscape. Instead, their algorithmic search for the loss minima was formulated as a mean to implicitly regularise the learning of a neural network.

At the time that Hochreiter and Schmidhuber presented their paper, machine learning research was largely compared with statistical studies. Like statistical models, machine learning practitioners of the time were interested in creating simple models that generalise well to unseen data (arguably due to the lack of computational resources at that time). This largely followed the school of thoughts of information theory [Shannon, 2001] such as the AIC [Akaike, 1998, 1974] and TIC [Takeuchi, 1976]. We can see this trend of thoughts because Hochreiter and Schmidhuber fre-

¹⁰In contrast, the exact geometry of the location can be unveiled if we choose to provide just one update after seeing all the data.

quently compared their method to *Weight Decay* [Weigend et al., 1991] and *Optimal Brain Surgeon* [Hassibi and Stork, 1993].

Furthermore, Hochreiter and Schmidhuber derived their method from KL divergence and the authors claimed that their method was hence equivalent to minimising the minimal description length. Though this particular study was later shown to be insufficient for fully understanding the utility of flat minima, it was nonetheless one of the most important and pioneering work in machine learning that has influenced how we mostly think of flat minima as the de-facto desirable minima.

2.2.4.1 Consequences of Large Batch Sizes

As we previously discussed in Section 2.2.3.4, the topology of loss minima and that of the loss landscape are largely influenced by a number of factors including the batch size. Batch size is a particularly interesting hyper-parameter for network optimisation. This is because that gradient descent is an inherently sequential process and thus cannot be parallelised well for fast computation. Updating a network with large batch sizes hence becomes one of the few options to accelerate learning. However, a plethora amount of work (such as LeCun et al. [2012]) have noticed that network generalisation can be negatively impacted from employing large batch sizes. Literature that focuses on batch sizes are hence important for understanding the trade-off for computational efficiency in machine learning versus network generalisation.

In order to understand the effects of batch sizes on optimisation, recent studies theorise that the choice of batch size directly modifies the geometric shape of loss minima thus influencing the quality of network generalisation. Furthermore and heavily influenced by the style of thinking of Hochreiter and Schmidhuber [1995], a large portion of this literature claimed that large batch sizes attributed to the phenomenon of minima-sharpening which consequentially led to sub-optimal network parameterisation. Most of these work studied the minima-sharpening phenomenon through the lens of the Hessian of the loss function. This is because that given the analogy of the marble on the loss landscape of Figure 2.3 and Equation (1.4), while $-\nabla\mathcal{L}$ represents the velocity and direction of the marble, $\nabla^2\mathcal{L}$ represents the corresponding acceleration and hence the curvature of the terrain; $(-\nabla\mathcal{L}, \nabla^2\mathcal{L})$ thus provides a more detailed description of the terrain around the minima.

Since neural networks can easily have more than thousands of parameters, most literature in the Hessian study of minima topology are restricted to focus on the leading eigenvalues. Keskar et al. [2017] approximated the eigenvalues by taking the average gradients of random projections^{IC15} on the loss. Similarly, Yao et al. [2018] used white box adversarial attacks [Ge et al., 2015] based on the top 20 eigenvalues of the Hessian spectrum to converge their network to a position on the landscape with positive curvature hence increasing robustness^{IC16}. Both of these two studies (and their related work) found that employing large batch sizes could lead to the increase of the leading eigenvalues by more than an order of magnitude and thus resulting in the formation of extremely sharp minima. Furthermore, these papers also found that improper choices of batch sizes could easily lead to the drop of more than 5% of

testing accuracy on the standard benchmark datasets of Cifar100 [Krizhevsky, 2009] and TIMIT [Garofolo et al., 1993].

Currently the best practice regarding batch sizes follows the influential work of Jastrzebski et al. [2017]. In their experiments, they demonstrated that SGD performances were mainly determined by the learning rate α (see Equation (1.4)) and the batch size $|B|$ (see Equation (1.3)); and that the larger the ratio $\frac{\alpha}{|B|}$ the better the generalisation. In a recent paper, Gurbuzbalaban et al. [2021] studied the $\frac{\alpha}{|B|}$ ratio through the lens of implicit renewal theory. They related this ratio to SGD’s tendency in converging base learner weights to a heavy tail distribution. In another paper [Simsekli et al., 2020], the same panel of authors further showed that the heavy-tailed behaviour improved generalisation in network.

2.2.4.2 Alternative Opinions on Minima

However, not every literature agree that flat minima are necessary. In a highly cited work by Dinh et al. [2017], the authors claimed that the notion of sharp and flat minima are especially problematic for understanding network optimisation and one should hence apply these concepts cautiously without further context^{ICC17}. They argued that the conjecture of Hochreiter and Schmidhuber [1995]’s work had a rather imprecise definition on minimum itself^{ICC18}. Alternatively, Dinh et al. presented an argument based on the non-identifiability^{ICC19} of simple ReLU networks, and they derived a proof showing that each minimum was actually of infinite volume^{ICC20} hence negating the descriptive power of *flatness*.

There also exist some dispute over whether large batch sizes are responsible for the creation of sharp minima (and hence responsible for poor network generalisations). In the recent work of Jastrzebski et al. [2019], the authors of that paper presented a new experimental setup for measuring the Hessian spectrum of the loss¹¹. Instead of measuring the Hessian spectrum of a fully optimised network (like that in Yao et al. [2018] and Keskar et al. [2017]), Jastrzebski et al. measured the Hessian spectrum from the *very beginning* of the network optimisation procedure to understand how the top eigenvalues continued to change throughout the entire scope of training. They observed two interesting remarks. Remark 1: that the eigenvalues of the spectrum started low and became larger throughout the course of training regardless of the batch size and learning rate^{ICC21}. Remark 2: that smaller batch sizes and larger learning rates caused a network to have considerably smaller Hessian spectra at the beginning of training. This paper hence showed that it was natural for a neural network to parameterise well in a sharp minima.

2.2.5 Section Conclusion

Empirically, neural networks can train well even without a practitioner having full understanding of the loss landscape (and the minima of which it entails). However,

¹¹The paper by Jastrzebski et al. [2019] looked at the top 50 eigenvalues and they approximated the eigenvalues through the Lanczos algorithm [Lanczos, 1950].

there are a couple of interesting and important work that suggest that we should first consider before ruling out the importance of studying minima.

Regardless of the lack of precise notation from the seminal paper of [Hochreiter and Schmidhuber \[1995\]](#), flat minima continues to be one of the desiderata in the machine learning community. The reason is simply because that flat minima search usually results in empirically superior networks. For instance, the work of [Chaudhari et al. \[2019\]](#) proposed *Entropy-SGD* where their optimisation algorithm had access to a fixed amount of look ahead Langevin dynamics to account for the drift in weights during parameterisation. Such a minimisation in drift indirectly seeks for a flat terrain. In addition, though the Entropy-SGD paper was not formulated as a meta-learning paper, its look ahead mechanism shared a high degree of similarity to the work of MAML [[Finn et al., 2017](#)] and it also yielded iteration-efficient updates.

Large batch sizes should always be used cautiously. In the work by [Golmant et al. \[2018\]](#), the authors of that paper found that though increasing the batch size led to faster training, their experiments also demonstrated that the speed-up was not always linear. That is, there existed important threshold(s) that slowed the speed-up from linear to sub-linear, and from sub-linear to no virtual improvements. Aside from the non-guaranteed linear speed-up, their experiments found that there was almost a guaranteed decrease in network generalisation with large batch sizes.

To conclude this chapter, let us consider the work of [Zhang et al. \[2017a\]](#). This work presented an interesting finding, such that neural networks could still be trained well on images that had fixed but randomised labels¹². This unexpected finding showed that instead of *inferring* the labels like a statistical model, neural networks were actually *remembering* the mapping between inputs and outputs. The experimental results of this paper suggested that neural networks (or at least deep models) should not be considered as extensions of statistical models and hence analysis in the backward pass (thus including studies on the loss minima) should only be considered in parallel with the network architectural designs of the forward path. The work of [Li et al. \[2018\]](#) supported this idea and the authors showed that the geometric shape of a loss minima could change dramatically with small architectural modifications. For instance, residual connections [[He et al., 2016a](#)] yielded smoother loss minima than those without the connections^{ICC22}.

¹²To an astonishing 97.5%, see Table 2 in Appendix B of [Zhang et al. \[2017a\]](#).

Neural networks have long been studied under the light of Bayesian inferential processes. However like their optimisation in the backwards pass, neural network inferences pass through multiple layers to add more details to the extracted features. The following section discusses dynamical systems; and we introduce an alternative mathematical perspective to conceptualise the incremental changes that occur within a network.

Section Content:

- Section 2.3.1: Differential System & Difference System,
- Section 2.3.2: Visualising Difference Systems with Cobweb Diagrams,
- Section 2.3.3: Visualising Differential Systems with Vector Fields,
- Section 2.3.4: Equilibrium,
- Section 2.3.5: Related Ideas: Existing Applications of Dynamical Mathematics in Understanding Neural Networks,
- Section 2.3.6: Prerequisite 3 Conclusion.

2.3 Some Fundamentals on Dynamical Systems

Dynamical systems study the geometric relationship between space and time. Their applications can be found in the movements of celestial bodies [Neishtadt, 1984], in fluid dynamics [Rusak et al., 2012], in the binding sites of enzymes [Keener and Sneyd, 1998], and in narrating particle collisions [Sneyd et al., 2017]. Dynamical systems are applied in a wide range of disciplines because they provide a rich set of vocabularies to describe qualitative changes in variables as systems unfold in time. This thesis will use dynamical mathematical concepts to study how neural values propagate between layers.

Note, this section should only serve as a brief introduction to dynamical systems. Interested readers should consult Glendinning [1994] for a broad and complete view on dynamical techniques and implications.

2.3.1 Differential System and Difference System

Dynamical systems can be either continuous or discrete. *Differential systems*

$$u' = \mathcal{Q}(u) \quad (2.7)$$

describe *continuous* changes on variable u as \mathcal{Q} . While *difference systems*

$$u_{t+1} = \mathcal{H}(u_t) \quad (2.8)$$

describe the non-continuous *discrete* update of \mathcal{H} occurring to variable u_t at time t . Differential Systems can be approximated with difference equations and they have an equivalent discretised version of

$$u_{t+1} = u_t + \zeta \mathcal{Q}(u_t) \quad (2.9)$$

with a small positive constant $0 < \zeta \ll 1$ that denotes instantaneous changes.

Let us now compare the two dynamical systems to the transformations that take place in a neural network. The non-linear transformation of an MLP in Equation (1.1) can be described as the difference system of Equation (2.8); whereas the iterative parametric updates from gradient descent of Equation (1.4) is a form of differential system described in Equation (2.9). Also note that the residual connections [He et al., 2016a] of Equation (2.1) also follows the same formulation as that of Equation (2.9). These similarities indicate that dynamical mathematics serve as an alternative to Bayesian inferential processes for understanding properties of neural networks.

2.3.2 Visualising Difference Systems with Cobweb Diagrams

The difference equation of $u_{t+1} = \mathcal{H}(u_t)$ is a repetitive process with some starting value $u_0 = \kappa$. At all time, variable u_t serves both as the output of instance t and as the input of instance $t + 1$. A sequence of propagation from a low dimensional difference system can be visualised with cobweb diagrams.

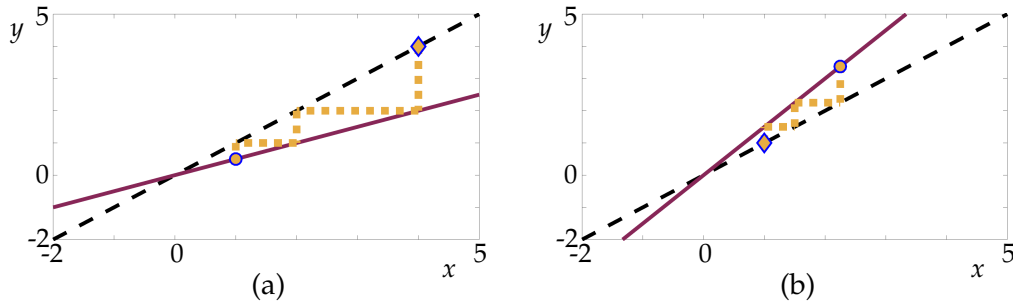


Figure 2.8: Simple Examples of Cobwebs

We present 2 example cobweb diagrams in Figure 2.8, with subplot (a) representing the difference equation $u_{t+1} = 0.5u_t$ and (b) for $u_{t+1} = 1.5u_t$. We use y to represent u_{t+1} and x to represent u_t . A cobweb diagram consists 4 fundamental elements. The 45-degree $y = x$ line, the characteristic shape of the operator $y = \mathcal{H}(x)$, the initial point $u_0 = \kappa$, and a trail of propagation trajectories. The dashed line is $y = x$, the solid line is $y = \mathcal{H}(x)$. The initial positions are in rhombuses, the solution trajectories are dotted, and the final positions are in circles.

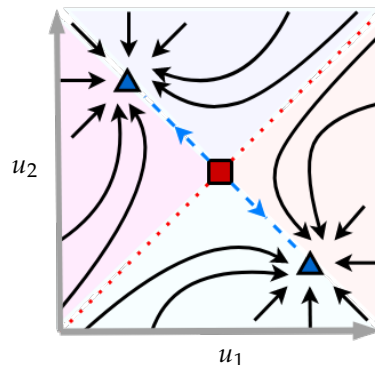
Cobwebbing visualises the repeated insertion of functional outputs. Its trajectory is composed of two different types of mappings: a vertical mapping that initiates from $y = x$ and terminates on $y = \mathcal{H}(x)$; and a horizontal mapping which starts from $y = \mathcal{H}(x)$ and ends on $y = x$. Vertical mappings represent for functional computations, i.e., given the input of $u_{t-1} \rightarrow$ compute for the output of $\mathcal{H}(u_{t-1})$; whereas horizontal mappings represent for system updates, i.e., given the finished computation of $\mathcal{H}(u_{t-1}) \rightarrow$ update the network variable of $u_t = \mathcal{H}(u_{t-1})$. This is the reason why the trajectories in Figure 2.8 are zig-zagged.

Not all trajectories are zig-zagged and their shapes are highly influenced by the topological outline of the transformation function \mathcal{H} . Interested readers can find descriptions on the tent map propagation in Boeing [2016]. The trajectories and the ending locations of the propagation will be discussed later in Section 2.3.4.

2.3.3 Visualising Differential Systems with Vector Fields

The differential equation of $u' = \mathcal{Q}(u)$ is a different kind of repetitive process to difference equations. Instead of providing a sequence of *transformations*, differential equations provide a sequence of *incremental updates*. By comparing Equation (2.8) to Equation (2.9), the “ u_t+ ” component indicates that differential equations provide an inherently smoother change than those from difference equations. A sequence of propagation from a low dimensional differential system can be visualised with vector fields.

Vector fields are visual representations of dynamical flows that describe the rules of propagation. Figure 2.9 is an example vector field of the two arbitrary dimensions of (u_1, u_2) of vari-

Figure 2.9:
An Example Vector Field

able u . In this figure, we see a vector field with four separate coloured regions with arrow trajectories. The arrow trajectories describe the routes for information propagation; furthermore, those arrows therein the same coloured regions behave qualitatively similar. That is, depending on the initial values of $(u_1, u_2) = (\kappa_1, \kappa_2)$, variable u would start from one of the 4 coloured regions and propagate as according to the rules dictated by the corresponding coloured region. Hence depending on the starting location, the trajectory will converge (or diverge) to different regions in the solution space.

Like the last section, we will defer our discussion on the ending location of a propagating variable in Section 2.3.4. Furthermore, interested readers should consult [Tennyson \[1989\]](#) for some common scenarios in differential systems.

2.3.4 Equilibrium

One important concept in dynamical mathematics is equilibrium. Equilibria are singularities on the vector field which stop variables from changing either if solution trajectories start on, or enter such singularities. That is, an equilibrium achieves

$$x_{t+1} = x_t \quad (2.10)$$

in a difference system; and they achieve

$$x' = 0 \quad (2.11)$$

in a differential system. These singularities do not always exist in an arbitrary dynamical system, but when they do, they have huge implications on both the long-term and short-term behaviours of variable propagation. Solution trajectories move towards attractive / stable equilibria known as sinks, and moves away from repulsive / unstable equilibria known as sources. In this thesis, we will only consider those scenarios where equilibria exist; and we will now revisit [Figure 2.8](#) and [Figure 2.9](#) to describe how equilibria influence these example scenarios.

2.3.4.1 Equilibria in Difference Systems

The equilibria of a difference system satisfy the condition of Equation (2.10). Due to this reason, the equilibria for the difference systems illustrated in [Figure 2.8](#) can be found figuratively by intersecting the 45-degree line $y = x$ and the characteristic shape of $y = \mathcal{H}(x)$. Recall that the rhombuses in the figure represent for the starting positions and that the circles represent for the ending positions. Hence, the origin of subplot (a) is a stable equilibrium because the trajectory converges towards the origin; and that of subplot (b) is an unstable equilibrium because the trajectory diverges away from it.

For difference systems, *the gradient around the equilibria* is important. The gradient around the origin of subplot (a) is < 1 and the equilibrium is attractive; whereas for subplot (b), the gradient around the origin is > 1 and the equilibrium is repulsive.

More importantly, the attractiveness and repulsiveness of these singularities dictate whether the network variables grow or shrink. While travelling towards the sink, the variable of subplot (a) decreases in magnitude, i.e., the magnitudes of the y -coordinates of the trajectory decrease and hence so do the network outputs u_t . As for the scenario in subplot (b), the trajectory travels away the source and hence the variable of subplot (b) increases in magnitude.

More interesting and complicated dynamics may emerge when the topological outline of $y = \mathcal{H}(x)$ becomes more intricate. Furthermore, these visual insights also enable a dynamical mathematical practitioner to craft suitable transformation function \mathcal{H} for achieving their desirable variable propagation.

2.3.4.2 Equilibria in Differential Systems

The equilibria of a differential system satisfy the condition of Equation (2.11). Hence they would need to be derived from setting Equation (2.7) with $u' = \mathcal{Q}(u) = 0$. In the example given in Figure 2.9, we presented a scenario where the blue triangles denoted sinks and that the red square represented a source.

Locations on the vector field corresponded to the magnitudes (and positivities) of the variables. As previously mentioned, this sample vector field consisted four coloured regions each describing qualitatively similar variable propagation. The four regions were separated by the red dotted line and the blue dashed line of which represented null-clines – boundaries with zero movement in specific directions. Hence this figure presented a scenario where a boundary separated the two basins of attractors. Using vector fields, a dynamical mathematical practitioner could understand the long term behaviour of an unperturbed system. For instance, the scenario presented in Figure 2.9 showed that depending on the starting position of the variable (on either side of the red null-cline), the variable would eventually propagate towards one of the two values entailed by either stable equilibria.

2.3.5 Related Ideas: Existing Applications of Dynamical Mathematics in Understanding Neural Networks

As we previously mentioned in Section 2.3.1, dynamical mathematics have close connections to non-linear MLP transformations, ResNet [He et al., 2016a] shortcut connections, and also gradient descent. However, much of the early endeavour in connecting dynamical mathematics to neural networks were conducted in the network backward pass instead of the forward pass. This is because that the usefulness and practicality of shortcut connections were only re-demonstrated recently in Srivastava et al. [2015] and in He et al. [2016a]¹³. In this related ideas section, we will first narrate how dynamical system has been used to forward our understanding in hastening gradient descent; and then we will discuss the links between dynamical systems and a trend of study in search for more advanced neural architectures.

¹³See our narration on the research for enabling deep learning in Section 2.1.

2.3.5.1 Dynamical Studies in the Backward Pass

There has been a long history of connecting optimisation to dynamical systems. See [Schropp and Singer \[2000\]](#) for instance, where the authors of that paper analysed the convergence condition for an augmented family of systems for constrained minimisation. During training, neural networks are optimised for some loss function through the iterative process of gradient descent (see Equation (1.4))

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}$$

of which is an inherent differential process. Some popular choices of optimisation schemes are *Stochastic Gradient Descent* (SGD), SGD with *Momentum* [[Rumelhart et al., 1986](#)], the *Nesterov Accelerated Gradient* (NAG) [[Nesterov, 1983](#)], and the *Adam* optimiser [[Kingma and Ba, 2015](#)]^{ICC23}¹⁴. All of these optimisers have previously been analysed using dynamical mathematics.

We have previously discussed gradient descent and visualised network updates in Section 2.2. If SGD were comparable to that of the motion of a marble on a table, then the machine learning community has commonly seen momentum as the technique which addresses the notion for an accelerating gradient descent that accumulates a velocity in directions of “persistent reduction in the object across iterations” [[Sutskever et al., 2013](#)]¹⁵. A more mathematical description of the effectiveness of SGD and momentum can be found in the work of [Qian \[1999\]](#).

[Qian](#) connected SGD to the Newtonian motion of a particle moving through a viscous medium in a conservative force field^{ICC24}. Under this analogy, momentum corresponds to the mass of the particle represented by SGD. The most important part of this work was that by using dynamical mathematics, the authors found that there existed *Lyapunov Functions* [[La Salle et al., 1962](#)] for SGD and its momentum variant^{ICC24}. In brief, a Lyapunov function maps a dynamical system to a transformed space with two properties: a set of strictly positive first derivatives with locally negatively definite second derivatives. Thus, if one could find the existences of Lyapunov functions for a dynamical system, one has also geometrically shown the existence of a stable equilibrium point in the original dynamical system. Hence [Qian](#) showed that SGD and its momentum variant could always find a minimum on the loss landscape, albeit a non-optimal local minimum. Interested readers should consult [Glendinning \[1994\]](#) for precise descriptions on Lyapunov functions.

A dynamical mathematics analysis on NAG can be found in [Su et al. \[2014\]](#). This study found that the first-order optimiser of NAG can be expressed as a time-dependent *non-autonomous second-order*^{ICC25}¹⁶ ordinary differential equation. Their formulation described NAG as an over-damped system during the early phases of optimisation of which transitioned into an under-damped oscillatory system during the later phases of optimisation.

Interestingly, the work of [Barakat and Bianchi \[2018\]](#) also described Adam as a

¹⁴The respective formulations are shown in the ICCs.

¹⁵This quote is taken from paragraph 1 on Section 2 of page 2 of [Sutskever et al. \[2013\]](#).

¹⁶Note that [Qian \[1999\]](#)’s Newtonian description is also a second-order description for SGD and momentum.

non-autonomous ordinary differential equation. In addition, Adam was expressed as a multivariate differential system^{ICC26} to factor in the changes of which approximated the low order moments. Their analysis provided interesting insights, such as the adaptive moments could be interpreted as convex combinations of past gradients. More importantly, they found that Adam could almost surely not converge; and that one would need to continually decrease the learning rate to assist convergence.

2.3.5.2 Non-Dynamical Studies in Understanding the Forward Pass of Shortcut Connections

Since the effectiveness of shortcut connections [Schraudolph, 1998] was demonstrated for deep learning [He et al., 2016a], the community has dedicated a lot of efforts in understanding them. For instance, the paper by Veit et al. [2016] introduced the unraveled view of ResNet to describe the network as a combination of many paths instead of a single deep network. Furthermore and by conducting an ablation study on a trained ResNet where connections of different depths were randomly deleted, Veit et al. also showed that ResNets could be seen as an ensemble of shallow networks. In a similar study conducted by Huang et al. [2016], the authors introduced the concept to train ResNets by stochastic depths – to *train* short networks and employ deep networks at *test* time. This was achieved by randomly replacing a subset of layers with the identity function during training. The authors found that both the training time and network generalisability could be improved.

In addition to these two work which studied the flow of neural information, He et al. [2016b] experimented with different ResNet designs and reordered deep learning modules. They found that significant testing improvements could be achieved when the residual connection was applied post activation¹⁷. Another noteworthy study is the work conducted by Zoph and Le [2017]. The authors took a neural architectural search approach where they used reinforcement learning to train an RNN to heuristically and sequentially construct a deep learning network one layer after another. For image recognition, their RL-RNN yielded a deep CNN with shortcut connections of various lengths. This is particularly fascinating because the qualitative analysis of Veit et al. [2016] and Huang et al. [2016] provided clear implications such that that the ResNet architecture only received minor benefits when it consisted long range connections that span through more than 2 layers of modules; but when constructed by a learning agent, the RL-RNN somehow deemed it necessary to include various levels of such *redundancies*¹⁸.

2.3.5.3 Dynamical Studies in Understanding the Forward Pass of Shortcut Connections

While the community have accumulated a substantive amount of insights on residual networks via ablation studies (see the last section), there have also been a large

¹⁷He et al. [2016b] also named this ordering the “pre-activated” ResNet of which can be found in Figure 1 on page 2 of their paper.

¹⁸Interested readers should also refer to Figure 7 in page 15 of Zoph and Le [2017].

amount of dedicated effort aimed in designing better neural network architectures through exploiting shortcut connections. *PolyNet* [Zhang et al., 2017b], *FractalNet* [Larsen et al., 2017], and *Shake-Shake* [Gastaldi, 2017], to name a few, were all extend from the shortcut connection designs from ResNets.

PolyNet took inspiration from the Inception paper [Szegedy et al., 2015] and focused on exploring structural diversity in designing deep network. More specifically, they experimented with component $\mathcal{H}(n_{L(\gamma-1)}, W_{L_\gamma})$ of Equation (2.1) and proposed several different compositions^{CC27}. The relationship between PolyNet and ResNet is hence similar to that of the *FitzHugh–Nagumo* model [FitzHugh, 1955; Nagumo et al., 1962] and the *Hodgkin–Huxley* model [Hodgkin and Huxley, 1952]¹⁹. However, instead of simplifying ResNet, PolyNet looked at more complicated non-linear transformations in hope of increasing network computational power.

FractalNet was proposed as an alternative network of which *achieved* deep learning without shortcut connection. The paper found that instead of explicitly summing features via a shortcut connection, deep learning could also be achieved via *concatenating* multiple branches of self-similar networks. This fashion of combining parallel features was similar to the relaxation oscillation of the *Van der Pol* oscillator [Van der Pol, 1926], where a slow dynamic complements a fast dynamic to describe the periodicity of an oscillatory motion.

Shake-Shake combined different aspects from stochastic depth [Huang et al., 2016] and PolyNet. First and like stochastic depth, it was a regularisation technique that aimed to increase the overall generalisation of ResNet. Second and similar to PolyNet, its regularisation was provided through network modification. More specifically, Shake-Shake experimented with component $\mathcal{H}(n_{L(\gamma-1)}, W_{L_\gamma})$ of Equation (2.1) by adding a second non-linear transformation. This changed a ResNet into a parallel multi-branch network and the features of the 2 branches were summed with a stochastic affine combination^{CC28}.

Recognising the increasing popularity of network modification, Lu et al. [2018] conducted a thorough analysis on this class of ResNet-inspired networks. They started by stating the similarity between ResNets and ordinary differential equations; then using this as the basis of their argument, they found that interestingly, stochastic depth, FractalNet, and Shake-Shake could all be referred to well-established dynamical mathematical studies. They connected stochastic depth and Shake-Shake to *Ito calculus* [Kunita, 2010], and FractalNet to the 2nd-order *Runge-Kutta* [Butcher, 1963].

Besides extending the existing ResNet architecture, Chen et al. [2018] crafted the novel *Neural Ordinary Differential Equations* (Neural ODEs) architecture utilising differential equation solvers. The aim of Neural ODEs was to learn the underlying vector field transformation (see Figure 2.9) to encapsulate the rules that govern variable propagation. The importance of Neural ODEs was that, arguably, it bridged RNNs to CNNs. However, this interesting work is still in its infancy.

¹⁹Hodgkin and Huxley [1952]’s celebrated Nobel-prize winning work described how action potentials in neurons were initiated and propagated in a large squid axon. Their intricate mathematical model was later simplified by FitzHugh [1955] & Nagumo et al. [1962] into a 2-variate differential system.

2.3.6 Section Conclusion

The prerequisite of this section looked at the differences between the discrete-time difference system and the continuous-time differential system. Then, we introduced cobweb diagrams and vector fields as visual cues for the propagation of information. We also compared the similarities between MLP non-linear transformations to difference systems, and likewise the similarities between ResNets and differential systems. Furthermore, we discussed how gradient descent is inherently a differential dynamical process.

While in the related ideas, we discussed how popular optimisers such as SGD with momentum, NAG, and Adam have been analysed with dynamical mathematics. Following the existing analyses in the backwards pass, we covered some literature that used dynamical mathematics to understand and extend ResNets. Specifically, we addressed PolyNet, FractalNet, and Shake-Shake. Though these architectural studies have showed slight improvements on standard machine learning benchmarks, the student-author would suggest our readers to not pursue this direction of research. Our justifications will be provided below.

First, the performance gain of PolyNet, FractalNet, and Shake-Shake came at the cost of making ResNets more complicated. This is perplexed when we recall that, second, ResNet was originally proposed to resolve the gradient vanishing problem in the backward pass. Hence its shortcut connection was actually only introduced as a byproduct in the forward pass for addressing the complications in the backward pass. Shortcut connections thus challenged conventional interpretations on the forward pass of feedforward neural networks. The student-author hence thinks that it is harder to make meaningful architectural changes when the understandings on the fundamental properties of the conventional ResNet are still lacking. Third and as the Inception paper [Szegedy et al., 2015] has shown, though network modification can yield substantive performance differences, they can only be made meaningful when the modifications themselves are backed by educated guesses. For instance, the Inception module consisted convolutional kernels of various sizes because the previous findings of Lin et al. [2014] showed that lower layers create clusters of correlated units of which concentrated in local regions.

2.4 Chapter Conclusion

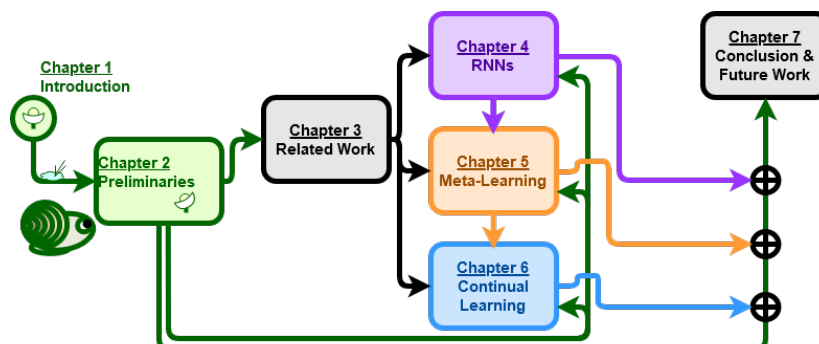


Figure 2.10: Progress Update

In this preliminary chapter, the student-author covered a variety of interconnected topics. We discussed how residual connections [He et al., 2016a] and highway connections [Srivastava et al., 2015] became the dominant architectures for deep learning because their shortcut connections enabled gradient information to flow through many layers. We also delved into the concepts of loss landscape and loss minima to visualise the backward pass of network parameterisation where gradient was required. Then, we briefly covered some fundamentals on dynamical mathematical techniques and elaborated on how dynamical understanding could be utilised to gain insights on both the forward pass and backward pass.

To conclude this chapter, the student-author would suggest that more research is required in RNNs for understand ResNets. This is because that ResNets can be seen as special cases of Highway Networks, and that highway connections were adapted from the de-facto RNN design of LSTM [Hochreiter and Schmidhuber, 1997b]. By explaining the expressiveness of LSTMs with dynamical systems, we will be able to gain insights on how deep learning models process their information, thereby understanding how to design better models and how to interpret the entwined relationship between the forward pass and the backward pass. The next chapter begins with our background section on RNNs; and the thesis will start exploring dynamical mathematical insights for creating better AutoML techniques.

Related Work

This is the background chapter of the thesis. Therein this chapter, we will present a discussion each for the three following topics:

- Section 3.1: Recurrent Neural Networks & Their Interpretations,
- Section 3.2: Meta-Learning & On Learning to Optimise, and
- Section 3.3: Continual Learning & On Episodic Memories.

The RNN section will begin with an introduction on comparing simple RNNs with *Long Short-Term Memories* (LSTMs) [Hochreiter and Schmidhuber, 1997b] and *Gated Recurrent Units* (GRUs) [Cho et al., 2014]. Then we will provide a deep analysis on a selective few papers on *understanding* gated RNN components. The section will conclude with a wide range of papers that have previously re-designed their RNN connections. Our research on RNN can be found in Chapter 4.

The meta-learning section will start with a clarification on the differences between a task specific base learner and a meta-learner. It will cover three of the most common approaches in existing practice. Then, we will go in depth to discuss the archetype of our choice – *Learning to Optimise* [Andrychowicz et al., 2016], and provide a detail documentation of its experimental setup. Our research on meta-learning can be found in Chapter 5.

The continue learning section will follow a similar format to the meta-learning section. We will first discuss the continual learning experimental setup and relate it to neighbouring learning paradigms in AutoML. The section will cover the three main methods that taxonomises the discipline. We will compare the advantages and disadvantages of the three approaches and cover our archetype of choice – *Episodic Memory-Based Continual Learning* [Lopez-Paz and Ranzato, 2017] and its most recent research developments. Our continual learning research can be found in Chapter 6.

All of our content are heavily related to the preliminaries of Chapter 2 and we will frequently refer to previously discussed concepts. The content therein this chapter will also show how extensively RNNs have been applied in meta-learning and continual learning; and why understanding the inner workings of recurrent information processing can greatly benefit the design for better AutoML techniques.

*In biological neural networks,
Lamme [2006] argues that
recurrent processing is necessary and sufficient for consciousness.*

*In artificial neural networks,
RNNs are structurally different to feedward networks and
they trade iteration over depth for neural processing.*

*Understanding the procedure behind RNN's recursive processing is difficult, and
the following section reviews ways to understand the inner workings of RNNs.*

Section Content:

- Section 3.1.1: Early RNN Designs and the Gradient Vanishing Problem (Again),
- Section 3.1.2: The Long Short-Term Memory Design,
- Section 3.1.3: A Brief History on LSTM's Development,
- Section 3.1.4: The Gated Recurrent Unit Design & The Non-Immediately Clear Meanings of Gated Units,
- Section 3.1.5: Section Recap,
- Section 3.1.6: Strategies: Remove, Grow, & Observe,
- Section 3.1.7: Examples of Existing Work on Remodelling LSTM and GRU,
- Section 3.1.8: Section Conclusion.

3.1 Recurrent Neural Networks and Their Interpretations

As we mentioned in the introduction and in the preliminaries, our main purpose for understanding RNNs is for gaining deeper insights on the shortcut connections of highway networks [Srivastava et al., 2015] in turn to forward our knowledge in AutoML. This section will start with the early RNN design and its problems; then on LSTM, the RNN design which inspired highway networks. Afterwards, we will introduce GRU and emphasise the difficulties which make RNNs hard to interpret.

3.1.1 Early RNN Designs and the Gradient Vanishing Problem (Again)

RNNs are one of the most widely applied neural network modules. However, there were difficulties in applying early RNN designs (such as [Elman, 1990]) in practice. An early RNN design can be found in Figure 3.1, of which is very similar to the MLP architecture that we have previous seen in Figure 1.1. The only discrepancy is that RNNs appended the MLP design with recurrent connection(s) such that time delayed representation of occurred events can be reused for data inference. This yields the network formulation¹ of

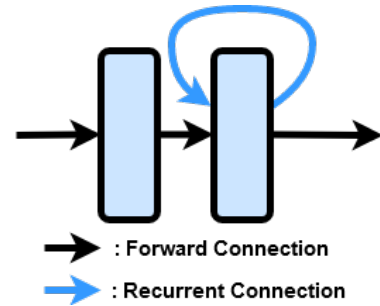


Figure 3.1: An Early RNN Design

$$\mathbf{a}_t = \mathcal{H}(\mathbf{W}\mathbf{a}_{t-1} + \mathbf{U}\mathbf{h}_t) \quad (3.1)$$

at time step t with hidden state \mathbf{h}_t , the forward synaptic connections of \mathbf{W} , and the recurrent synaptic connections of \mathbf{U} .

In principle, the recurrent connection of an RNN makes it possible to conduct sequential inferences over several time steps. Thus, the loss of an RNN is similar to the batch loss of Equation (1.3) but aggregated over T fixed steps such that

$$\mathcal{L} = \sum_{1 \leq t \leq T} \mathcal{L}_t = \sum_{1 \leq t \leq T} \frac{1}{|\mathcal{B}_t|} \sum_{(x_{i,t}, y_{i,t}) \in \mathcal{B}_t} \mathcal{C}(\hat{y}_{i,t}, y_{i,t}). \quad (3.2)$$

Like the chain rule found in Equation (2.5), the gradient of Equation (3.2) follows

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}} = \sum_{1 \leq t \leq T} \frac{\partial \mathcal{L}_t}{\partial \mathbf{W}}, \text{ such that} \quad (3.3)$$

$$\frac{\partial \mathcal{L}_t}{\partial \mathbf{W}} = \sum_{1 \leq \tau \leq t} \left(\frac{\partial \mathcal{L}_t}{\partial \mathbf{a}_t} \frac{\partial \mathbf{a}_t}{\partial \mathbf{a}_\tau} \frac{\partial \mathbf{a}_\tau}{\partial \mathbf{W}} \right). \quad (3.4)$$

From here, the term $\frac{\partial \mathbf{a}_t}{\partial \mathbf{a}_\tau}$ echos with $\frac{\partial \mathbf{a}_{L(\zeta+1)}}{\partial \mathbf{a}_{L_\zeta}}$ of Equation (2.5). Thus as described in

¹The formulation is adapted from Equation 2 in page 1 of Pascanu et al. [2013].

Section 2.1.1, the existence of this Jacobian-like term causes difficulties in training a neural network due to gradient vanishing². For the MLPs of Equation (1.1), gradient vanishing emerges from training through ζ layers; whereas for the RNNs of Equation (3.1), gradient vanishing is from retracing past copies in τ time steps.

3.1.2 The Long Short-Term Memory Design

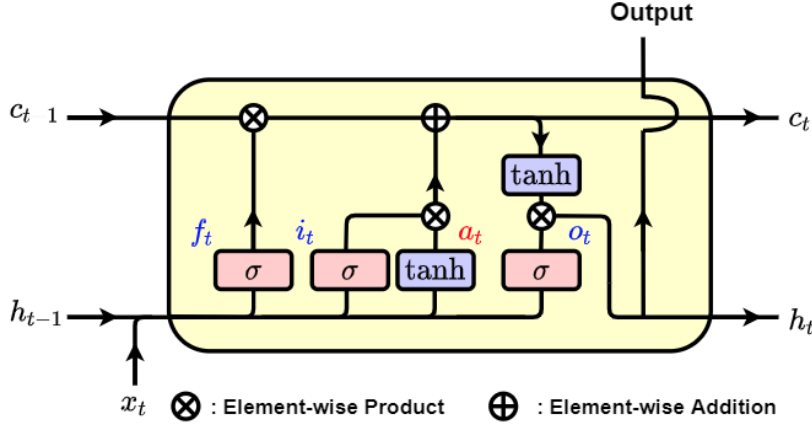


Figure 3.2: The LSTM Design

Long Short-Term Memory (LSTM) [Hochreiter and Schmidhuber, 1997b] has established itself as the de-facto stand RNN as it has continually shown empirical effectiveness in withstanding gradient vanishing. That is, it is preferred among practitioners because it is easier to train and that it demonstrates superior performances.

An LSTM unit is shown in Figure 3.2; and it is formally written as

$$\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,I,O\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{F,I,O\}}), \quad (3.5)$$

$$\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{h}_{t-1} + \mathbf{b}_A), \quad (3.6)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and} \quad (3.7)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \quad (3.8)$$

There are two hidden variables to conduct inferences – the hidden state \mathbf{h}_t and the cell state \mathbf{c}_t . The hidden state and the network input \mathbf{x}_t are reused to create three *Gated Units* – the forget gate \mathbf{f}_t , the input gate \mathbf{i}_t , and the output gate \mathbf{o}_t , of which are activated by sigmoid function σ with values $\in [0, 1]$. They also create the internal update \mathbf{a}_t activated via hyperbolic tangent \tanh with values $\in [-1, 1]$. Matrices \mathbf{W} s are the forward synapses, \mathbf{U} s are the recurrent synapses; while vectors \mathbf{b} s are the biases. Furthermore, \odot denotes the element-wise product.

Thus overall, an LSTM employs gated units to fine-tune cell state \mathbf{c}_t . Forget gate \mathbf{f}_t releases a section of old cell memory whereas input gate \mathbf{i}_t adjusts the magnitude of \mathbf{a}_t to replenish the cell. The network output is also the hidden state, of which is a transformed version of the cell state with control exposure from output gate \mathbf{o}_t .

²Another difficulty for training RNNs is gradient explosion when individual Jacobians > 1 . This can be resolved by *Gradient Clipping* [Pascanu et al., 2013] – by setting an upper bound for gradients.

3.1.3 A Brief History on LSTM's Development

Highway Network's shortcut connection in Equation (2.3)

$$\mathbf{k}_{L_\gamma} = (1 - \mathbf{T}_{L_\gamma}) \odot \mathbf{k}_{L_{(\gamma-1)}} + \mathbf{T}_{L_\gamma} \odot \mathcal{H}(\mathbf{k}_{L_{(\gamma-1)}} \mathbf{W}_{L_\gamma})$$

is inspired by LSTM's incremental update of the cell. However, the vanilla LSTM formulation of which we just introduced in the previous page is actually very different to its original formulation provided in Hochreiter and Schmidhuber [1997b].

3.1.3.1 Error Propagation through a Carousel

The original LSTM therein Hochreiter and Schmidhuber [1997b] consisted only 2 gated units: the input gate \mathbf{i}_t and output gate \mathbf{o}_t ; and it was given as

$$\begin{aligned} \mathbf{i}_t, \mathbf{o}_t &= \sigma(\mathbf{W}_{\{I,O\}} \mathbf{x}_t + \mathbf{U}_{\{I,O\}} \mathbf{h}_{t-1} + \mathbf{b}_{\{I,O\}}), \\ \mathbf{a}_t &= \tanh(\mathbf{W}_A \mathbf{x}_t + \mathbf{U}_A \mathbf{h}_{t-1} + \mathbf{b}_A), \\ \mathbf{c}_t &= \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and} \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \end{aligned}$$

The main novelty of this original LSTM design was that it introduced the concept of *multiple* hidden variables. By adopting $\mathbf{c}_t = \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$ to learn an increment instead of learning a transformation, this helped to introduce the identity matrix to the Jacobian (see Section 2.1.1). That is, and we emphasise this because it is not well recognised in the community, that the modifications made towards early RNN designs of which we now come to know as the LSTM forward pass were originally based on mathematical properties which were nice to satisfy therein the backward pass. Furthermore, this particular act of *summation* (the + sign in $\mathbf{c}_t = \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$) was referred to as the *Constant Error Carousel* by Hochreiter and Schmidhuber for trapping the gradient of the error through the cell state as a medium.

However, if the simple act of error-trapping was enough, then what was the original purpose of the output gate? This was actually not well elaborated therein Hochreiter and Schmidhuber [1997b]; but we may be able to get an intuition of their decisions in some of Schmidhuber's older work.

Meta-Learning for Achieving General Intelligence

In his diploma thesis [Schmidhuber, 1987], Schmidhuber conducted an in-depth discussion with his readers on the concept of "Giving a system the ability to learn the methods how to learn, too." and referred to learning as *First Occurrence* with *Confirmation* which adds unexpected events to the information processing system³. This

³Below are quotes from Von Weizsäcker [1985] quoted in Schmidhuber [1987] (continues overleaf).
Keyword: first occurrence (*Erstmaligkeit*); confirmation (*Bestätigung*)

Quote 1: On the effectiveness of coding unexpected events:

*Nahe dem Grenzfall hundertprozentiger Bestätigung kann jede Neuigkeit registriert werden
[..]die Verfasser schlagen vor, in diesem Grenzfall die Erstmaligkeit direkt durch die*

was because that Schmidhuber’s work was heavily influenced by [Von Weizsäcker \[1985\]](#)’s concept of evolution³, of which conceptualised learning as a process that enforces modifications of a (already working) subsystem to discover structures in tasks. From this narrative alone, we can see that Schmidhuber’s diploma thesis saw a huge overlap on the different AutoML disciplines of which we now call reinforcement learning, meta-learning, and continual learning.

We can continue to see traces of this narrative in other early work of Schmidhuber. For instance in [Schmidhuber \[1993\]](#), he experimented with *Self-Inferential Weight Modifications* at two places of a neural network. First, there were some reading weights of which served the purpose to address the *currently connected weights* of the network. Second, there were another set of weights which addressed the *outputting connections* of the network. It is hence arguable that, the weights that were designed to attend the currently connected weights gave rise to the input gate \mathbf{i}_t of LSTM, whereas those weights that attended to the outputting connections manifested themselves as the output gate \mathbf{o}_t in LSTM.

In another line of study, Schmidhuber’s concurrent interests in meta-learning and RNNs also produced a prototypical class of neural networks which use *Fast Weights* to modify their embedded knowledge (or policies) [[Schmidhuber, 1992](#)]. In a relatively recent and highly cited paper, *Highway Networks* [[Ha et al., 2017](#)] demonstrated that deep neural networks with fast weights can achieve competitive generalisation ability while removing a huge section of network parameters.

3.1.3.2 Learning to Predict while Learning to Forget

The vanilla LSTM that we would come to know today was introduced in [Gers et al. \[1999\]](#). The rationale behind that paper for adding the third and final gate, the forget gate \mathbf{f}_t , was based on the concern such that $\mathbf{c}_t = \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$ would cause the cell to grow indefinitely. An arbitrarily large cell state would severely limit the expression of the hidden state because of $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$. That is, an arbitrarily large cell state could cause component $\tanh(\mathbf{c}_t)$ of the hidden state to saturate at \tanh ’s boundaries and in turn limit the ability to represent for different inputs fed into the network.

Besides narrating the ill consequences in the forward pass, the authors of that paper also provided a deeper mathematical conjecture on the corresponding problems in the backwards pass. They stated that, a squashed hidden state would not only make it difficult to represent different inputs, but it would also limit the amount of modifications to be provided in the gradient. This is because that the unique topological outline of the derivatives of σ and \tanh are near-zero on both extrema such that $\lim_{|\epsilon| \rightarrow \infty} \sigma'(\epsilon) = 0$ and $\lim_{|\epsilon| \rightarrow \infty} \tanh'(\epsilon) = 0$. Without sufficient values to be provided through the backward pass, gradient descent would not make meaningful modifications towards the RNN weights and thus training will be made difficult.

*Information im Sinne Shannons zumessen.[...] Nimmt aber der Bruchteil der **Bestätigung** ab, so kann nicht mehr jede Neuigkeit pragmatisch effektiv registriert werden*

Quote 2: On the necessity on having pre-existing bias to conceptualise more
*Blosse **Bestätigung** entspricht der Karikatur des Spezialisten: er weiss alles über nichts;
 blosse **Erstmaligkeit** entspricht der Kankatur des Generalisten: er weiss nichts über alles.*

3.1.3.3 Learning to Predict with Direct Access to the Cell

After proposing the standard vanilla LSTM that we currently know, the same panel of authors proposed the LSTM-variant with *Peephole Connections* [Gers and Schmidhuber, 2000]. The modification proposed in that paper was rather straight forward but non-mathematical – the idea was to simply supply gated units of the LSTM with cell state information connected by an extra recurrent synapse \mathbf{V} such that $\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,I,O\}}\mathbf{h}_{t-1} + \mathbf{V}_{\{F,I,O\}}\mathbf{c}_{t-1} + \mathbf{b}_{\{F,I,O\}})$. The authors of that paper showed that the main advantage of peephole connection was that the network became more precise to finer time lags and can even differentiate between 49 and 50 lag steps. Since its introduction, this technique has gained some interest in the robust speech recognition community (like Kang et al. [2018]), but it is not commonly adopted in other major machine learning disciplines. Thus in our thesis, we will refer to the LSTM with forget gate of Gers et al. [1999] as *the* LSTM or the *vanilla* LSTM.

3.1.4 The Gated Recurrent Unit Design and the Non-Immediately Clear Meanings of Gated Units

Thus far in this section, we have compared simple RNNs to LSTMs and covered the history of the incremental changes on the early LSTM designs. Most importantly, we clarified that the forget gate and the input gate were added to the forward pass to satisfy certain mathematical properties in the backward pass; whereas the output gate was built in the network based on the original author’s heuristics without elaborated explanations. In this subsection, we will introduce another widely adopted RNN called *Gated Recurrent Unit* (GRU) [Cho et al., 2014] and further emphasise on the difficulties for understanding RNNs.

3.1.4.1 GRU: The Other De-Facto Standard RNN

Originally introduced for neural machine translation and called the *RNN-Encoder* [Cho et al., 2014], GRU is simpler than LSTM with 1 instead of 2 hidden variables, and 2 instead of 3 gated units. Since its introduction, GRU has become the second de-facto standard RNN for its similar computational power to LSTM despite being simpler [Chung et al., 2014]⁴. A GRU is depicted in Figure 3.3.

Formally, GRU is written as

$$\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{I,R\}}\mathbf{x}_t + \mathbf{U}_{\{I,R\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{I,R\}}), \quad (3.9)$$

$$\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A), \text{ and} \quad (3.10)$$

$$\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t. \quad (3.11)$$

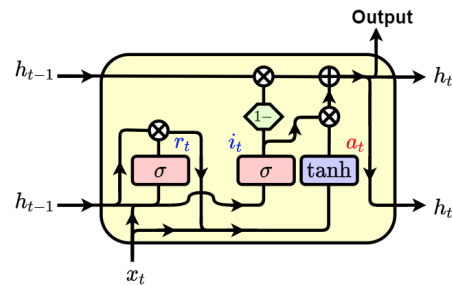


Figure 3.3: The GRU Design

⁴Chung et al. [2014] was the first to refer the network as GRU, and to finalised the GRU design.

Table 3.1: A Comparison between GRU and LSTM

GRU Cho et al. [2014]	LSTM Hochreiter and Schmidhuber [1997b]
$\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{L,R\}}\mathbf{x}_t + \mathbf{U}_{\{L,R\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{L,R\}})$ $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A), \text{ and}$ $\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t.$	$\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,I,O\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{F,I,O\}}),$ $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{h}_{t-1} + \mathbf{b}_A),$ $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and}$ $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t).$

3.1.4.2 A Comparison between GRU and LSTM

GRU (of 2014) was introduced many years after LSTM (of 1997b). However, it was not originally proposed as a modification towards LSTM, and was instead presented as neural machine translation paper to recurrently encode and decode context⁵ for sequence-to-sequence learning. Hence, the internal mechanics of GRU were not mathematically analysed and were only textually described.

Difference 1: Hidden Variables

The GRU RNN employs one hidden variable, the hidden state \mathbf{h}_t , unlike the LSTM RNN which uses both \mathbf{h}_t and the cell state \mathbf{c}_t . However, the GRU hidden state is updated similarly to the cell state of LSTM. That is, GRU also incrementally updated its \mathbf{h}_t with gated units to remove old values and add new information.

Difference 2: Style of Variable Synchronisation

The sole gated unit responsible for the fine-tuning process is the input gate \mathbf{i}_t . GRU replaces LSTM's forget gate \mathbf{f}_t with $1 - \mathbf{i}_t$; or in other words, that the GRU design binds its forget gate with its input gate.

Difference 3: Exposure Control on the Network Output

The network outputs of both gated RNNs are hidden state \mathbf{h}_t . However, one major difference is that, the LSTM output $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$ is much more carefully adjusted. First, LSTM uses tanh to set upper and lower bounds on the extreme values of the network output; and second, that LSTMs utilise an additional gated unit \mathbf{o}_t to limit excessive outgoing information.

Difference 4: Creating an Internal Update

Similar to how the LSTM RNN provided careful control over its network output \mathbf{h}_t , the GRU RNN instead directed more resources towards implementing a more rigorous control over its internal update \mathbf{a}_t . Therein the GRU \mathbf{a}_t , Cho et al. employed a *Reset Gate* \mathbf{r}_t to fine tune the strength of the recursive usage of the representation of

⁵Year 2014 also saw another very successful neural machine translation paper of Sutskever et al. [2014]. However, differ to the GRU paper of Cho et al. [2014], Sutskever et al. [2014] used LSTMs as their choice of recurrent encoder and decoder.

old occurred events \mathbf{h}_t . [Cho et al. \[2014\]](#) described that this was for the purpose of selectively encoding important sequential content⁶

In this formulation, when the reset gate is close to 0, the hidden state is forced to ignore the previous hidden state and reset with the current input only. This effectively allows the hidden state to drop any information that is found to be irrelevant later in the future, thus, allowing a more compact representation.

As we see, the explanation provided in [Cho et al. \[2014\]](#) was not as mathematically elaborated than the justification of the inclusion of the forget gate \mathbf{f}_t in [Gers et al. \[1999\]](#) for [Hochreiter and Schmidhuber \[1997b\]](#). However, the student-author of this thesis found that the GRU reset gate plays a similar role to the LSTM output gate. That is, if we were to return to the formulations provided in [Table 3.1](#) and substitute every LSTM instance of \mathbf{h}_{t-1} with $\mathbf{o}_{t-1} \odot \tanh(\mathbf{c}_{t-1})$, then we would yield the following comparison in [Table 3.2](#).

Table 3.2: A Comparison between GRU and LSTM

GRU Cho et al. [2014]	LSTM (hidden state substitution)
$\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{I,R\}}\mathbf{x}_t + \mathbf{U}_{\{I,R\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{I,R\}})$ $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A), \text{ and}$ $\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t.$	$\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t$ $+ \mathbf{U}_{\{F,I,O\}}(\mathbf{o}_{t-1} \odot \tanh(\mathbf{c}_{t-1}))$ $+ \mathbf{b}_{\{F,I,O\}}),$ $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t$ $+ \mathbf{U}_A(\mathbf{o}_{t-1} \odot \tanh(\mathbf{c}_{t-1}))$ $+ \mathbf{b}_A),$ $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and}$ $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t).$

From the magenta coloured-content in the table, we can see an old LSTM output gate \mathbf{o}_{t-1} serving the same role as the GRU reset gate \mathbf{r}_t for the internal update \mathbf{a}_t . However, from the cyan coloured-content in the table, we can understand that the old LSTM output gate \mathbf{o}_{t-1} not only fine-tunes the recurrent strength for \mathbf{a}_t but also for \mathbf{f}_t , \mathbf{i}_t , and even the new output gate of \mathbf{o}_t . This is a novel perspective of the student-author's, and as much as his knowledge goes, there isn't another paper of which has made this connection. Still, this is a relatively non-mathematical analysis.

3.1.5 Section Recap

Thus far in this section, we mentioned the similarities and differences between simple RNN units, LSTMs, and GRUs. We highlighted on the history of their incremental changes; but more importantly, we emphasised on the non-mathematically justified modifications made towards the de-facto standard RNNs. The next subsection will start reviewing some few but existing work in *understanding* RNN properties.

⁶Quoted from the last paragraph on page 3 (left column) of [Cho et al. \[2014\]](#).

3.1.6 Strategies: Remove, Grow, and Observe

The three main strategies in understanding GRUs and LSTMs (and RNNs in general) are ablation study, search study, and observation study. Ablation study trains a sub-network which consists only a subset of components of the vanilla gated RNN; then after training, measures the performance difference between the sub-network to its vanilla counterpart baseline. The more the performance lesion, the more important the removed component could be for the well-being of the network. Search studies look for systematic ways to evolve the architectural design of the network in search of important yet non-obvious methods to boost existing network performances. Observation studies are qualitative studies that analyses the mapping relationships of the network inputs to their outputs to deduce the underlying mechanics of a network. Note, there are not many studies related to *understanding* RNN properties; instead, most RNN papers are related to their applications.

3.1.6.1 Removing Network Components: What to Keep?

The work by Greff et al. [2016] is one of the only large scale ablation studies on the vanilla LSTM design. The authors found that peephole connections and the binding of input gate to forget gate (like that in GRU) almost always functioned comparatively to the vanilla LSTM design across all tasks. This provided potential reasons to why peephole was not picked up as a standard practice outside of the robust speech recognition community⁷. In particular, peephole introduced a lot of redundant network parameters, and this paper showed that an LSTM functioned well as long as there is a *minimal threshold* of shared information among all components.

Greff et al. also demonstrated that the removal of any single gated unit from the vanilla LSTM resulted in significant performance lesion. However, the most severe damages were caused by the removal of the forget gate and the removal of the output gate. However, this paper did not provide an in-depth mathematical analysis between the interactions of the forget gate and input gate, nor did they comment much on Cho et al. [2014]’s GRU⁸. Hence it remained unclear on the implicit compounded effects between any 2 arbitrary components of the LSTM network.

3.1.6.2 An Extensive Search on Sibling Models: What Else to Include?

This subsection will discuss the three search studies of Pascanu et al. [2014], Bayer et al. [2009], and Jozefowicz et al. [2015]. In brief, the study of Pascanu et al. [2014] investigated RNN connections and questioned whether these connections could be manually redesigned to enhance network expressiveness. While in the other two work of Bayer et al. [2009] and Jozefowicz et al. [2015], the authors of both papers

⁷And in the adjacent NLP community of language modelling, it is common practice to use the vanilla LSTM design as the default LSTM design. See more details in Merity et al. [2018] and their code repository in <https://github.com/salesforce/awd-lstm-lm>.

⁸The only comments from Greff et al. [2016] to Cho et al. [2014] was that on the last paragraph of Section III on page 3, of which they commented as “...a simplified variant of the LSTM...”.

chose to randomly mutate network components in search of a best performing model.

Pascanu et al. [2014] investigated three types of connections in a deep RNN. First, there was the *input-to-hidden* functions; for instance the $\mathbf{W}_A \mathbf{x}_t$ term in the LSTM internal input of $\mathbf{a}_t = \tanh(\mathbf{W}_A \mathbf{x}_t + \mathbf{U}_A \mathbf{h}_{t-1} + \mathbf{b}_A)$. Second, there was the *hidden-to-hidden* functions, the $\mathbf{U}_A \mathbf{h}_{t-1}$ of \mathbf{a}_t . Then also the *hidden-to-output* functions which discussed inter-LSTM layer relationships where the series of outputs of a preceding LSTM layer was treated as input to a subsequent LSTM layer. This paper was partially motivated by the need to increase the representation capacity of LSTM. It was not motivated from a mathematical perspective; and instead, it was for finding interpretable intermediate representations which can serve as easily understandable cues of the heuristic knowledge learnt by the network. The authors attempted to increase an RNN's representation power by adding more parameters. However, this did not make their RNNs more interpretable nor did it boost the RNN performance.

Inspired by the incremental development of the LSTM unit (see Section 3.1.3), Bayer et al. [2009] made the assumption that existing RNNs could still be non-optimal. Motivated by this thought, the authors decided to start from a small network and then applied random mutations on the starting network to add new connections and new components. The authors reported some surprising findings. One particularly important remark⁹ was that, one mutation started with a simple RNN and developed into a 3-node fully recurrent unit without gated units. This candidate mutation showed that good RNN designs did not necessarily require gated units.

Jozefowicz et al. [2015] took a similar evolutionary route to Bayer et al. but at a much larger scale. The author of this paper also raised the importance note that RNNs were only reparameterised as LSTMs to address the vanishing gradient problem and proposed that LSTMs actually *do not have any representational advantage*¹⁰. This served as a reason to why Pascanu et al. was not able to significantly increase LSTM expressiveness via redesigning connections. Influenced by this line of thoughts, their mutation scheme started strictly with either LSTM or GRU to have starting models that were already less prone to gradient vanishing.

None of Jozefowicz et al.'s final candidate networks diverged far from the original LSTM and GRU designs¹¹. The authors of this paper also concluded their paper with one interesting remark – that none of their best performing candidate networks always outperformed LSTM nor GRU. This finding could be interpreted in two different ways. First and an optimistic interpretation, that LSTMs and GRUs are indeed “optimised” RNNs. Second and a pessimistic interpretation, that LSTMs and GRUs are perhaps incorrect starting points for mutating RNNs and that any RNNs based on these two designs only yield less optimised networks.

The combining results of Bayer et al. [2009], Jozefowicz et al. [2015], and Pascanu et al. [2014] showed that it is difficult to improve LSTM by naively adding more

⁹Here we are referring to Section 3.1 on page 9 and Figure 5 of Bayer et al. [2009].

¹⁰Refer to 2nd paragraph on page 2 (right column) in Jozefowicz et al. [2015].

¹¹Refer to page 7 (top left corner) in Jozefowicz et al. [2015] for their top 3 mutated models. All of their 3 yielded networks remained largely the same as GRUs and only with minor modifications appearing in either the input-to-hidden or hidden-to-hidden connections.

parameters. These studies motivated this thesis to take a mathematical approach to understand *how* to create a better network in a first principle's manner.

3.1.6.3 An Analysis on the Network Input to Output: What Has been Learnt?

This subsection will discuss the two qualitative analysis studies of [Wu and King \[2016\]](#) and [Karpathy et al. \[2016\]](#). These two papers demonstrated that LSTMs could be made interpretable *on the network outputs* of natural language processing tasks. The former tested LSTMs for statistical parametric speech synthesis and reported interesting purposes of the forget gate and cell; while the latter reported LSTM behaviours on character-based language modelling and focused on the hidden state.

Motivated by the need to understand the learnt heuristics of LSTM for speech processing, [Wu and King \[2016\]](#) proposed a 2-part study to analyse LSTM components. First and inspired by [Greff et al. \[2016\]](#) (see Section 3.1.6.1), the authors trained different ablated variations of LSTM for speech synthesis and compared the performances. Specifically, they found that an LSTM with no output gate and no input gate¹² yielded acceptable performances for speech synthesis. Second, the authors analysed this simplest ablated variant with the least lesion and provided a qualitative analysis on the remaining variables in the said model.

Without the output gate and the input gate, their network was left only with the forget gate and the cell state¹³. [Wu and King](#) then showed that the average activation value of all forget gate neurons served nicely as the phoneme boundaries of speech utterances. They also showed that *one specific* dimension of the cell state exhibited strong correlation to the Mel-Cepstral coefficients (MCC)¹⁴.

One of the most common applications of RNNs is the NLP of text for language modelling. RNNs are convenient tools to manifest the idea of Bayesian inferential processes because they can be continually updated to learn the probability of a token given the conjoint probability of a series of preceding tokens. [Karpathy et al. \[2016\]](#) aimed to understand the learnt heuristics of LSTMs by analysing and visualising their outputs and gated units when trained as a character-level language model.

[Karpathy et al.](#) compared LSTMs to the two baseline models of basic *n-gram* language models [[Heafield et al., 2013](#)] and *n-NN*, where a fully connected neural network was fed a concatenation of the features of *n* consecutive characters. LSTMs significantly outperformed *n-gram* models; but more interestingly, it was shown that *n-gram* models outperformed *n-NNs*. The authors thus conjectured that normal neural networks would quickly overfit on text data, and that it required a special architecture design like the LSTM in order to prevent this from happening. LSTMs were also shown to have a higher probability to correctly predict special characters and

¹²[Wu and King \[2016\]](#) bound the s-LSTM input gate to forget gate $i_t = 1 - f_t$ like how a GRU bound its forget gate to its input gate. Refer to page 3 (top left corner) for more details.

¹³Without the output gate, the hidden state becomes $h_t = \tanh(c_t)$ of which is a mere transformation and near identical of the cell.

¹⁴For those readers who are unfamiliar with speech engineering, MCC refers to a short-term power spectrum of a sound and thus possesses some memory qualities. This was hence the reason why [Wu and King \[2016\]](#) compared MCC to the cell state, the network memory, of their simplified LSTM.

were implied to have a longer memory than n-gram models in the sense that¹⁵ they always possessed significantly higher probability over n-gram models in creating “}” after creating “{”.

Besides visualising the performances of LSTMs to n-gram models, [Karpathy et al.](#) also provided images for observing some heuristically learnt LSTM mechanisms. On the LSTM hidden state, the authors found that the magnitude of \mathbf{h}_t decreased monotonically over a fixed length before resetting if it never encountered a special character such as “{”. The authors hence conjectured that the hidden state acted as a line-length counter. However, upon encountering a “{”, \mathbf{h}_t values remained close to -1 until it encountered a “}”. This indicated that character-level LSTMs learned either to highlight or to ignore all non-special characters in between special characters.

The authors also provided insights to the mean activation values of every neuron in the gated units across all time steps. This served as an interesting analytical tool because it could visualise the variability of every single neuron. From there, they found two interesting remarks. First, they found that for a 3-layer LSTM model, the neurons of all gated units were more saturated to either 0 or to 1 as the layer deepened. Second and they found that regardless of the layer, all gated units had some neurons that were either consistently 1s or 0s at all time. The authors found that these results to be interesting; but they also admitted that they were not able to understand what had actually been learnt¹⁶.

3.1.7 Examples of Existing Work on Remodelling LSTM and GRU

The present subsection provides more literature review on RNNs. This bundle of papers were not originally written for understanding RNN designs; instead, they were originally proposed for modifying and simplifying the gated RNN architectures. There will be a mix of positive and negative results, but these papers will provide us with valuable information to avoid previously found modelling pit-holes.

There exist many papers that modify LSTMs and GRUs to enhance computational speed. In the case of LSTMs, training is made time-inefficient because that each dimension of the cell state \mathbf{c}_t of the current time t depends on all entries of its previous instance of \mathbf{c}_{t-1} . That is, the computation is required to wait until \mathbf{c}_{t-1} is fully computed. This internal dependency makes LSTMs, and likewise for other gated RNN designs including GRU, difficult to be computed in parallel. In text processing, [Bradbury et al. \[2017\]](#) proposed to interleave RNN-layers with convolutional transformation layers and treated RNNs only as pooling layers with trainable parameters. In another paper, [Lei et al. \[2018b\]](#) proposed to reconnect the pre-activated

¹⁵The question on whether LSTMs (and RNNs) really possess long-term memories remain a popular and contentious topic. Toy examples in [Hochreiter and Schmidhuber \[1997b\]](#) showed that LSTMs were capable of possessing long term memories and that [Karpathy et al. \[2016\]](#) (in the student-author’s opinion) seem to agree to this idea. However, some recent papers such as [Zhao et al. \[2020\]](#) and [Daniluk et al. \[2017\]](#) seem to indicate otherwise. Such a stigma has also been exacerbated with the rise of popular alternative models such as TCN [[Bai et al., 2018](#)] and Transformer [[Vaswani et al. \[2017\]](#)].

¹⁶Refer to page 5 of [Karpathy et al. \[2016\]](#), ... *We struggle to explain this finding but note that it is present across all of our models.* The student-author also struggled to understand the fascinating phenomenon of which they remarked on Figure 3 of page 5 of their paper.

gated neural values as $\mathbf{W}\mathbf{x}_t + \omega \odot \mathbf{c}_{t-1} + \mathbf{b}$ with a trainable vector of ω to replace the ordinary $\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}$. The introduction of ω made each dimension independent of each other and hence parallelisable. Otherwise, one could also speed up RNNs like in [Kuchaiev and Ginsburg \[2017\]](#) where the authors factorised matrices \mathbf{W} and \mathbf{U} into smaller components to enable parallel updates of smaller blocks. Though we study RNNs and LSTMs to gain a better understanding on the utility of shortcut connections, there also exist papers with an opposite approach to ours. That is, there exist papers of which re-introduce shortcut connections to the LSTM for increasing network performances. [Yao et al. \[2015\]](#) introduced highway connections (see Section 2.1) which introduced cells in layer l with those of the preceding layer as $\mathbf{c}_t^l = \mathbf{f}_t^l \odot \mathbf{c}_{t-1}^l + \mathbf{i}_t^l \odot \mathbf{a}_t^l + \mathbf{d}_t^l \odot \mathbf{c}_t^{l-1}$ with the depth gate \mathbf{d}_t^l . Using this method, the authors yielded some performance increase; but since their paper lacked a mathematical analysis, the said performance increase could be simply due to the increase in trainable parameters from \mathbf{d}_t^l . [Yao et al. \[2015\]](#)'s paper was then followed up and extended by the work of [Kim et al. \[2017\]](#). [Kim et al. \[2017\]](#) also connected parts of a preceding LSTM layer to its immediately next LSTM layer. However, their connection was introduced in the hidden state instead of the cell state and with a residual connection instead of a highway connection. Though also lacking a mathematical explanation, the authors acknowledged that it was much better to append the extra information outside of the cell state because that the cell was originally introduced to solely address gradient vanishing.

There are several papers which removed the tanh non-linearity from the internal update of \mathbf{a}_t . In the work of [Balduzzi and Ghifary \[2016\]](#), the authors claimed that it was unnatural to for the LSTM cell to be formulated as $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$. Their explanation was that this formulation violated the physics principle of dimensional homogeneity such that 2 objects of different units are being added together (like of kg to volts). To address this problem, the authors proposed to reformulated any connection of the form $\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}$ to $\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{x}_{t-1} + \mathbf{b}$ to prevent the mixture of temporal information to non-temporal information. Furthermore, they removed the activation function of \mathbf{a}_t for a strictly linear internal input. Though their paper lacked a mathematical analysis on the backward pass, their reformulated LSTM outperformed the vanilla LSTM; however, when they applied the same modification to GRUs, the performance of their reformulated GRU went down. Thus showing that this method required further investigation and that it wasn't yet generally applicable.

Many papers also made the similar conjecture to GRU to bind the LSTM input gate to the forget gate. [Miao et al. \[2016\]](#) based their LSTM modifications on a qualitative analysis similar to that of [Karpathy et al. \[2016\]](#). In their paper, they showed that for acoustic modelling, the saturation (see Section 3.1.6.3) of neurons of the LSTM input gate exhibited strong negative correlation to that of neurons of the LSTM forget gate. However, while the authors bound the input gate to the forget gate, they also provided an extra trainable vector \mathbf{g} to increase the input gate expression as $\mathbf{i}_t = \mathbf{g} \odot (1 - \mathbf{f}_t)$. In another paper which also bound the two gated units, [Kusupati et al. \[2018\]](#) reformulated the GRU hidden state as $\mathbf{h}_t = \mathbf{f}_t \odot \mathbf{h}_{t-1} + (v(1 - \mathbf{f}_t) + \kappa) \odot \mathbf{a}_t$. This paper provided a much more in-depth

mathematical explanation on their modification. They analysed how neural information diverges and appended the original GRU hidden state with the trainable vectors of ν and κ to provide expressive yet stable condition numbers.

Thus far in this subsection, we introduced papers that based their RNN extension on either LSTM or GRU. There are, however, other ways to improve network computational power but without utilising gated units – both [Le et al. \[2015\]](#) and [Koutnik et al. \[2014\]](#) introduced modifications targeting simple RNN units. The former paper noted that, as long as we initialise the recurrent synapse (\mathbf{U} in $\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1}$) as (scaled versions) of the identity matrix, then the RNN will start with backpropagation with identity-like gradient highways. Their results demonstrated that such RNNs initialised with identity recurrent matrices were indeed easy to train and that they were good at modelling long-range dependencies. Whereas in the work of [Koutnik et al. \[2014\]](#), the authors took inspiration from *Echo State Networks* [[Jaeger, 2001](#)] and partitioned neurons of the simple RNN hidden states into multiple mutually exclusive clusters. All clusters were fully connected to the input neurons and output neurons; however and while the neurons that shared the same cluster were internally connected, those that were of 2 arbitrarily different clusters were purposely designed to not be connected. Furthermore, each cluster was updated at a different *rate* (some every 2 steps, some every 4, 8, or 16 steps). The network was designed as such to force the network to learn temporal hierarchical features from the data.

3.1.8 Section Conclusion

As we motivated in the introduction chapter, we study RNN and especially LSTMs to pursue a deep understanding on the properties of shortcut connections. We first covered the differences between simple RNNs, LSTMs, and GRUs. Then, we emphasised that the LSTM as well as the GRU reparameterisation came from the need to satisfy mathematical properties in the backward pass to mitigate gradient vanishing. However, by introducing new components to the forward pass, the gated RNN designs became less interpretable and hence we started to look for means to understand RNN mechanics.

Unfortunately, most papers on RNNs concern their applicability instead of their interpretation. Thus in Section 3.1.6, we analysed those few and existing important papers that aimed to understand RNNs via ablation [[Greff et al., 2016](#)], search [[Bayer et al., 2009](#); [Pascanu et al., 2014](#); [Jozefowicz et al., 2015](#)], and qualitative observations [[Karpathy et al., 2016](#); [Wu and King, 2016](#)]. However, all of these three mentioned approaches in understanding RNNs share one commonality – that they focused on interpreting the components of a *trained* network instead of understanding the *purpose* of why each component was included prior of training.

Thus in order to complement existing work and to forward our own knowledge on how to AutoML, we have decided that it is necessary for ourselves (and hence this thesis) to conduct a mathematical analysis to describe how RNNs and LSTMs

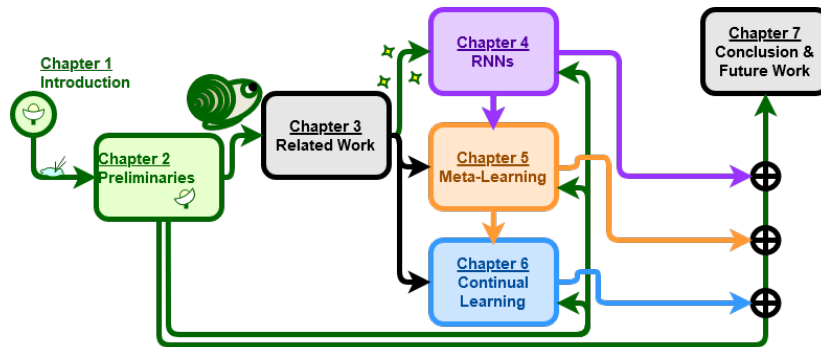


Figure 3.4: Progress Update

dynamically encode occurred events as memories for future inferences. After conducting our analysis, we will then aim to deliver a simplification to the LSTM design based on our analytical insights. Please refer to Chapter 4 for our research on RNNs.

Machine learning algorithms are traditionally hardwired by practitioners. However, recent developments have moved from learning classifiers and policies towards learning representations that enhance the operations of classifiers and policies. The following section discuss meta-learning; and the aim is to create self-referential and self-modifying code that enables the network itself to explore for better inferential strategies.

Section Content:

- Section 3.2.1: Offline Meta-Learning: Model Agnostic Meta-Learning,
- Section 3.2.2: Online Meta-Learning: Metric-Based Meta-Learning,
- Section 3.2.3: Online Meta-Learning: Model-Based Meta-Learning,
- Section 3.2.4: Online Meta-Learning:
Optimisation-Based Meta-Learning,
- Section 3.2.5: Details on Learning to Optimise,
- Section 3.2.6: Section Conclusion.

3.2 Meta-Learning & On Learning to Optimise

Meta-learning is the AutoML technique that achieves rapid knowledge acquisition. Its utility can be interpreted in two ways. First, it can refer to the efficiency in updating a task specific base learner. Second, it could also be formulated as how correctly does one network adapt existing knowledge to perceive new and low resource data. To elaborate, the former regards base learner optimisation; and that the latter is generally tested on few-shot, one-shot, and zero-shot learning [Vinyals et al., 2016b].

Hospedales et al. [2020] grouped meta-learning algorithms into *offline* approaches and *online* approaches. Both approaches dichotomise their experimental setup to compare local information under the influence of meta-knowledge. While local information refers to how to fit a data, meta-knowledge is about how to solve a task. Offline meta-learning provides preconditioned optimisation. Prior to a formal update, the algorithm gathers information in the parametric space for optimising an optimisation step. Whereas online meta-learning conjointly trains a second auxiliary meta-learner along with the base learner. The purpose of the meta-learner is to accumulate knowledge in assisting the base learner to solve a task.

We will discuss offline meta-learning and three main approaches in online meta-learning. As each approach takes a very different perspective in achieving meta-learning, readers may directly refer to Section 3.2.4 to see our method of choice.

3.2.1 Offline Meta-Learning: Model Agnostic Meta-Learning

Finn et al. [2017]’s work on *Model Agnostic Meta-Learning* (MAML) is the pioneering work for offline meta-learning. MAML employs a double-loop setup to optimise a base learner. For each iteration of the inner loop, the base learner with its unmodified weights is fed a small batch of data and an internal-optimisation is executed to find a candidate update. The outer loop then aggregates on all candidates to perform a formal update to the current parameters of the base learner. Due to this reason, MAML is also often described as a meta-learning algorithm that *optimises an initialisation*. Via the parametric search of the inner loop, the outer loop is capable of delivering more iteration-efficient optimisation to a base learner.

Since its introduction, MAML has played an influential role in meta-learning and it continues to be extensively studied. For one instance, Vuorio et al. [2019] studied task-dependent representations to increase MAML’s modulation for few-shot learning. In another paper, Sun et al. [2019] studied scaling and shifting parameters under the MAML paradigm for transferring learnt heuristics across large-scale datasets. MAML’s double loop setup also inspired work such as Grant et al. [2018] to see meta-learning from the perspective of Bayesian approach.

However and from a more dynamical mathematical point of view, MAML’s inner loop can be seen as a “look-ahead” mechanism commonly found in the implicit and explicit iterative methods for approximating the numerical solutions of ordinary differential equations. For instance, both the *Simpson’s Rule* and *Runge-Kutta’s Method* (see either Süli and Mayers [2003] or Atkinson [2008]) purposely overshoot

in a direction of update (the inner loop) to precondition the direction of update (the outer loop) for propagating the variables of a dynamical system. A similar idea was proposed in [Ayadi and Turinici \[2021\]](#); however, the authors did not formulate their paper as a meta-learning paper and instead as extensions to traditional handcrafted optimisations (see Section 2.3.5.1).

In a recent work, [Raghu et al. \[2019\]](#) found that MAML’s effectiveness lied in *feature reuse* instead of *rapid learning*. Their paper demonstrated that network parameters before and after an MAML update were highly similar. Hence, the paper argued that MAML’s capability in few-shot learning was due to the meta-initialised parameters’ effectiveness in finding high quality low level presentations that could be easily transferred across domains.

3.2.2 Online Meta-Learning: Metric-Based Meta-Learning

One of the main approaches in online meta-learning is through *Metric-Learning* [[Schultz and Joachims, 2003](#); [Weinberger et al., 2006](#)]. This archetype is mostly inspired by *Kernel Density Estimations* [[Davis et al., 2011](#)] and by *Nearest Neighbour* algorithms [[Fix, 1951](#)]. They can either take the identical setup as supervised learning to estimate probability density functions, or be formulated as a procedure for finding maximal separability between transformed inputs of dissimilar (pairs or query-key of) data.

Many influential work in metric-based meta-learning are conducted under the famous *Siamese Network* [[Bromley et al., 1993](#)] architectural setup. In brief, a Siamese network employs a pair of *identical* neural networks to encode two separate inputs and then uses a distance measure (such as L1 loss or cosine distance) to measure similarity/dissimilarity. The two encoder sub-networks are required to be identical; and this particular design ensures that inputs of the exact same content cause the least amount of dissimilarity score. Though not directly related, a similar idea can be found in the transfer learning concept of *Learning to Borrow* [[Lim et al., 2011](#)] of which measures contextual similarities between instances and target classes.

In the original Siamese network paper, [Bromley et al.](#) employed the *Time Delayed Neural Networks* [[Lang et al., 1990](#)] as their twin encoder sub-networks. More recently, many work have adopted the Siamese network experimental setup but substituted deep neural networks as their twin encoders; this direction of research has yielded many successes in meta-learning for few-shot learning. For instance, [Koch et al. \[2015\]](#) took a 3-step process for deep metric learning – first, they used the weighted L1 distance between two sets of features encoded by a pair of deep networks; second, that they mapped the differences to $[0, 1]$; and third, trained their network with a cross-entropy objective loss. In another paper, [Chopra et al. \[2005\]](#) took a 2-step approach and neglected the mapping. That is, their network directly compared the similarity between two inputs with the differences of their deep network encoded features. In other words, their approach was to directly learn the metric itself. This however, made the experimental setup slightly more tedious and required practitioners to use the contrastive energy function to train the network.

Some influential metric-based meta-learning techniques also developed means

to provide semantic meaning towards the extracted features before the differences were taken. This generally meant to conjointly train multiple qualitatively different networks to form a complicated encoding protocol. For one instance, *Matching Networks* [Vinyals et al., 2016b] extended traditional metric learning with a *Full Context Embedding* (FCE) LSTM RNN (see LSTM in Section 3.1.2). Their FCE mechanism employed a bi-directional LSTM to encode the features extracted by another deep (typically forward) neural network. RNNs were involved in FCE to encode concurrent features extracted by the backbone network. The RNNs of FCE were required to be bi-directional because in Vinyals et al. [2016a], the same panel of authors found that RNNs were susceptible to sequence invariances (permutations) and that order mattered for the extracted features. In another important work, *Relation Networks* [Sung et al., 2018] took a similar approach to matching networks but replaced the RNN embedding with a pooling layer. Matching network was applied to the few-shot learning scenario; and that a similar approach in the context of unsupervised clustering can be found in Salakhutdinov and Hinton [2007].

Another significant work in this category of meta-learning is the *Prototypical Network* [Snell et al., 2017]. Similar to Chopra et al. [2005]’s work, prototypical networks aim to directly learn metrics. However, while Chopra et al. employed neural networks as an instrument to replace traditional metrics, Snell et al.’s network focused on learning a metric *space*. On one hand, Chopra et al. took the differences between encoded features as a direct measurement of the similarity between two inputs. On the other hand, Snell et al. first encoded their inputs and defined *prototypical features* – the mean of all extracted features of those data within the same class; then, their objective was to minimise the distance between an element and its corresponding prototypical features in the metric space. Hence, prototypical networks can be interpreted as a bridge between deep learning and clustering techniques.

3.2.3 Online Meta-Learning: Model-Based Meta-Learning

Another major approach in meta-learning is through learning to compartmentalise recent (short term) memory from long term memory. There are two philosophical gains from such a setup. First, long term memory can be thought as the internalisation and stabilisation of knowledge while short term memory can be thought as small conceptual modifications to fine-tune a grander ideology (see Von Weizsäcker [1985]’s quotes in Section 3.1.3). Second, long term memory can be thought as higher level concepts which are generally applicable across tasks while short term memory serve as specialisations in distinct domains. While both analogies possess a short-term and a long-term memorial feature, these two schools of thoughts are actually very different from each other. Below, we will describe how these ideas manifest themselves in very different forms of experimental setups.

A good example of the first idea is Ha et al. [2017]’s work on *HyperNetworks*. Hypernetworks are auxiliary networks that dynamically modify the weights of a much larger network. See Section 3.1.3, this research can be regarded as a natural extension to Schmidhuber [1992] and the LSTM RNN [Hochreiter and Schmidhuber, 1997b]

in pursuit of a network which can intelligently embed its own context-dependent weights. This school of thought equates the ability to solve a task as a network's long term memory. That is, they treat the weights of a neural network as its long term memory. Hence in contrast, the layer-wise neural activation values which temporarily occur in the forward pass during instantaneous inferences are treated as the short term memory. While hypernetworks were not able to achieve large performance gains for deep neural networks, [Ha et al.](#) showed that this idea is capable of achieving a large parametric reduction for both deep RNNs and deep feedforward neural networks.

An adjacent idea to hypernetworks can be found in [Ba et al. \[2016a\]](#). The authors proposed that beyond neural activities and weights, synapses also played a significant role in the neural network inferential process. Weights are the modifiers to inputs and neural activities are the aggregated results of a layer function, while synapses are an intermediate substance that lie between the two – synapses refer to the aggregated strength from all weights to all inputs for computing one specific neural activity. Hence synapses have dynamics at many different time-scales and the authors conjectured that neural networks may benefit from regularising this apparatus which controls the rate of information assimilation. The authors hence provided fast weights and a decay mechanism to slowly diminish the significance of feedback of distinct network memory.

As mentioned, the first school of thought sees the weights of a network as its long term memory. In contrast, the second school of thought sees the entirety of a neural network as a short term memory. They put forward the idea that, neural networks are incapable of memorising long term memory because they are bounded by their pre-selected architectures and their pre-defined amount of parameters prior to training. Thus in order to create an intelligent agent capable of utilising *real* long term memory, the said agent would need to first, be able to store memory outside of itself (the network), and second, be able to sensibly access the stored external memory and mix it with its internal short term memory [[Hochreiter et al., 2001](#)]. Models capable of achieving so are also commonly referred to as *Memory Augmented Neural Networks* (MANN).

One important MANN is the *Memory Network* study of [Weston et al. \[2015\]](#) for the NLP task of question answering. Memory networks consist of an external memory storage and four learning components responsible for feature extraction, generalisation, output mapping, and inference. Generalisation is a learnt update protocol which modifies the externally stored memory with new extracted features. Output mapping reads from the externally stored memory and preconditions it before it is transformed by the inference step for creating a response. [Weston et al.](#)'s memory network employs an LSTM for output mapping and inference. Hence, we see that aside from an external memory, the model itself also possesses an internal memory. That is, while the internal memory is used for constructing a Bayesian inferential language model, the purpose of the external memory is to perceive less tangible meta-knowledge regarding the overall context of the question.

Another important MANN work was presented by [Santoro et al. \[2016\]](#), and the

novelty of that paper was on adapting *Neural Turing Machines* (NTMs) [Graves et al., 2014] for few-shot learning on image classification. Like Weston et al.’s memory network, NTM learns two levels of knowledge each associated with different time scales and is hence capable of shifting its bias by rapid cache representations. NTMs access their externally cached representations via a controller neural network. An interesting feature of the NTM controller is that their weights are updated at each time step by decaying their previous usage weights and by adding new weights. This is similar to the LSTM cell state update seen in Equation (3.7) but occurring at the weight level instead of the neural activity level. An architecture similar to NTM was proposed in Munkhdalai and Yu [2017]’s *Meta Networks*. Like NTM, meta networks also use weights to represent knowledge acquired at different time scales. However, the meta-learner and base learner are mutually communicative – the meta-learner and the base learner creates weights for each other.

3.2.4 Online Meta-Learning: Optimisation-Based Meta-Learning

Thus far in this section, we introduced offline meta-learning and two online meta-learning archetypes. The approach which we focus in this thesis is *Learning to Optimise* and it is also one of the major approaches in online meta-learning. Below, we will provide a higher level conceptualisation of the technique; and the reader can find a detailed experimental setup documented overleaf in Section 3.2.5.

Learning to optimise, also known as *Learning to Learn* (L2L), was first proposed by Andrychowicz et al. [2016] in the name-sake paper of *Learning to Learn by Gradient Descent by Gradient Descent*. That paper explored the possibility in updating one neural network with a second neural network. Under this dual-network experimental setup, the first network is the task-specific base learner while the second network is the meta-learner. Instead of updating the base learner with handcrafted gradient descent algorithms such as SGD [Robbins and Monro, 1951], the authors of that paper cast optimisation as the learning problem for the meta-learner to learn the decision procedures to **update the parameters of the base learner**. Since their meta-learner neural optimiser was itself **trained by gradient descent**, they called their method learning to learn **by gradient descent by gradient descent**. As a reminder, the concept of gradient descent can be found in Section 1.2.3 and a discussion on different variants of handcrafted optimisers can be found in Section 2.3.5.1.

We dedicated our efforts to this archetype of meta-learning for three reasons. First Andrychowicz et al.’s paper selected LSTM RNNs as their meta-learner neural optimiser. Thus, we get to reuse our insights gained through understanding LSTMs in Section 3.1 and Chapter 4. Second and as shown in Chapter 2, we are studying AutoML through the backward pass. Hence we are naturally attracted to a method which allows us to understand how to craft an learned optimisation scheme. Third and as shown in our comments on Qian [1999] in Section 2.3.5.1, gradient descent are one special type of dynamical systems. Since we are already studying LSTMs as dynamical systems and that LSTMs are the neural optimisers of this optimisation-based approach, this technique is a natural fit to all our other studies.

3.2.5 Details on Learning to Optimise

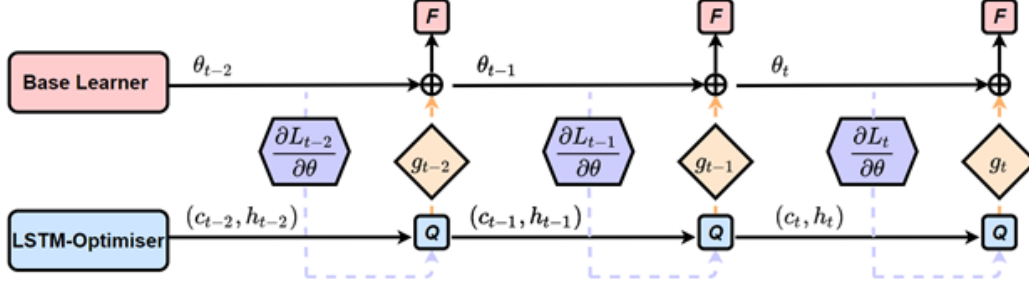


Figure 3.5: The L2L Experimental Setup

Keywords of this subsection will be colour-coded to assist readability.

3.2.5.1 Method Introduction

In their paper, [Andrychowicz et al.](#) explored meta-learning via optimisation. Unlike the handcrafted optimisations of SGD, momentum [[Nesterov, 1983](#)], Adam [[Kingma and Ba, 2015](#)], and RMSprop [[Tieleman and Hinton, 2012](#)], they propose to cast optimisation as a learning problem for the meta-learner. The aim for the neural optimiser meta-learner is to directly learn the rules for updating learner parameters θ_t of time t . After training, the neural optimiser is then used to substitute gradient descent of Equation (1.4) where

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}.$$

[Andrychowicz et al.](#) demonstrated in their work that the learnt neural optimisers could update learner parameters more rapidly than their handcrafted rivals. More importantly, learners updated with neural optimisers also yielded lower losses.

3.2.5.2 The LSTM-Optimiser – Replacing Gradient Descent

[Andrychowicz et al.](#) selected the LSTM RNN [[Hochreiter and Schmidhuber, 1997b](#)] as their meta-learner. Recall from Section 3.1.2, LSTM is formally written as

$$\begin{aligned} \mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t &= \sigma(\mathbf{W}_{\{F,I,O\}} \mathbf{x}_t + \mathbf{U}_{\{F,I,O\}} \mathbf{h}_{t-1} + \mathbf{b}_{\{F,I,O\}}), \\ \mathbf{a}_t &= \tanh(\mathbf{W}_A \mathbf{x}_t + \mathbf{U}_A \mathbf{h}_{t-1} + \mathbf{b}_A), \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and} \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \end{aligned}$$

with hidden variables \mathbf{c}_t and \mathbf{h}_t and a set of trainable weights (the \mathbf{W} s, \mathbf{U} s, and \mathbf{b} s).

Given a base learner \mathcal{F} with parameters θ and an LSTM \mathcal{Q} with parameters η , the LSTM-optimiser \mathcal{Q} is trained to replace gradient descent of Equation (1.4) to

prepare updates g_t for the learner parameters via

$$\theta_{t+1} = \theta_t + g_{t+1} \quad \text{and} \quad (3.12)$$

$$[g_{t+1}, \mathbf{c}_{t+1}, \mathbf{h}_{t+1}] = \mathcal{Q}(\mathcal{X}_t^{(\mathcal{Q})}, \mathbf{c}_t, \mathbf{h}_t | \eta). \quad (3.13)$$

That is, g_t is used in place of $-\alpha \nabla_{\theta} \mathcal{L}$ to update θ_t with some input $\mathcal{X}_t^{(\mathcal{Q})}$. The input will be discussed later in this subsection. From Equation (3.12), we see that the rules to be learnt not only consist the directions to update the parameters $-\nabla_{\theta} \mathcal{L}$ but also the scales for the said directions $-\alpha \nabla_{\theta} \mathcal{L}$. The parameters of the neural optimiser are updated along learners parameters but only once every “unrolling” instances; these concepts will be introduced in the next subsection. Refer to Figure 3.5 for a depiction on L2L.

The readers should note that usually, the LSTM hidden state (of the last layer) \mathbf{h}_t serves as the output of the network. However, we specifically separated g_t from \mathbf{h}_t in Equation (3.13). This was because that Andrychowicz et al. provided specific schemes to pre-process and post-process the inputs and outputs of the LSTM-optimiser. Along with $\mathcal{X}_t^{(\mathcal{Q})}$, these details will be discussed later in this subsection.

3.2.5.3 The LSTM-Optimiser – Training the Neural Optimiser

The neural optimiser learns through updating \mathcal{V} learners for \mathcal{S} steps. The objective loss \mathcal{L} of \mathcal{Q} is dependent on the objective loss \mathcal{L} of \mathcal{F} . (Note the differences between the two \mathcal{L} s.) For instance, an image classifier \mathcal{F} will have cross-entropy loss as \mathcal{L} . For every Ξ steps, the neural optimiser *unrolls* and parameters η of \mathcal{Q} are updated based on their respective gradients to \mathcal{L} , which is defined as

$$\mathcal{L}(\eta) = \mathbb{E} \left[\sum_{\tau=1}^{\Xi} \mathcal{L}_{\tau+1}(\theta_{\tau}) \right]. \quad (3.14)$$

While the learner loss $\mathcal{L}_{t+1}(\theta_t)$ is dependent on the parameters, Equation (3.14) shows that the objective of \mathcal{Q} depends on the *trajectory* of the optimisation in the base learner with Ξ as the maximum memory bandwidth for the LSTM-optimiser.

3.2.5.4 The LSTM-Optimiser – The Meta-Learner Input

Inputs $\mathcal{X}_t^{(\mathcal{Q})}$ of Equation (3.13) are related to the learner gradients. The raw gradients $\nabla_{\theta} \mathcal{L} = \nabla_{\theta} \mathcal{L}_{t+1}(\theta_t)$ are subjected to the following pre-processing scheme¹⁷ of

$$\mathcal{X}_{n,t}^{(\mathcal{Q})} \leftarrow \begin{cases} \left(\frac{\log(|\nabla_{\theta} \mathcal{L}^{(n)}|)}{p}, \text{sgn}(\nabla_{\theta} \mathcal{L}^{(n)}) \right) & \text{if } |\nabla_{\theta} \mathcal{L}^{(n)}| \geq e^{-p}, \\ (-1, e^p \nabla_{\theta} \mathcal{L}^{(n)}) & \text{otherwise} \end{cases} \quad (3.15)$$

where $\mathbb{R}^{|\theta_t|} \rightarrow \mathbb{R}^{|\theta_t| \times 2}$ for each element of the gradient for $n = 1, \dots, |\theta_t|$.

¹⁷Quoted from page 11 of Andrychowicz et al. [2016] under “A. Gradient preprocessing”.

Andrychowicz et al. devised this scheme to ensure that the magnitudes of every dimension is on the same order. They mentioned that this is necessary because neural networks, and including neural optimisers

*naturally disregard small variations in input signals and
concentrate on bigger input values.*

In addition, hyper-parameter $p > 0$ controls how small gradients are disregarded. Our thesis follows their default setting with $p = 10$ in all of our experiments.

3.2.5.5 The LSTM-Optimiser – The Meta-Learner Inference and Output

As noted in the last subsection, the LSTM-optimiser input is the pre-processed gradient of dimensionality $\mathcal{X}_t^{(\mathcal{Q})} \in \mathbb{R}^{|\theta_t| \times 2}$. We emphasise the dimensions because the meta-learner neural optimiser treats the amount of base learner weight of $|\theta_t|$ as the batch size of $\mathcal{X}_t^{(\mathcal{Q})}$ and treat the feature size as 2. Hence naturally, the input dimension of the LSTM-optimiser is also 2.

This is an important detail, and the reason is to enforce the LSTM-optimiser to provide *Coordinate-wise Updates*¹⁸ to the learner parameters of θ_t . The core idea for coordinate-wise updates is to ensure that each entry is updated only with information of its own gradient, i.e., $\theta_{n=1,t}$ is only updated with $\mathcal{X}_{n=1,t}^{(\mathcal{Q})}$ without needing the additional information of $\mathcal{X}_{n=n^*,t}^{(\mathcal{Q})}$ with any $n^* \neq 1$. Note that gradient descent of Equation (1.4) also do not update one variable with the gradient of another.

Caution, this is a particularly important detail and that failing to set the experiment correctly would cause an unnecessarily large network. Say that we set the hidden dimension of the LSTM-optimiser as M and also say that we implemented the network correctly, then the forward synapses of the first layer of the LSTM (the \mathbf{W} s of the formulation in Section 3.2.5.2) will be of dimensionality $\mathbb{R}^{M \times 2}$. One common mistake is that a practitioner incorrectly reshaped $\mathcal{X}_t^{(\mathcal{Q})}$ as dimensionality $\mathbb{R}^{2|\theta_t|}$ instead of $\mathbb{R}^{|\theta_t| \times 2}$ and hence mistakenly yielding $\mathbf{W} \in \mathbb{R}^{M \times 2|\theta_t|}$. This will make the LSTM-optimiser large and near-impossible to execute; what we desire is a small and efficient meta-learner that can correctly train a much larger base learner network.

Finally, the neural optimiser output of g_t in Equation (3.12) is re-scaled before adding to the learner parameters θ_t . It is timed by 0.1.

3.2.5.6 The LSTM-Optimiser – The Default Hyper-parameteric Settings

This thesis follows the default setting provided in Andrychowicz et al. [2016] and we use Algorithm 1 overleaf to clarify the setup. The LSTM-optimisers (and any other RNN-optimisers of this thesis) are 2 layers deep with input dimension 2 (see Section 3.2.5.5) and hidden dimension 20. Unless specified otherwise, the meta-training phase updates $\mathcal{V} = 20$ randomly initialised base learners (line 4 of Algorithm 1) for $\mathcal{S} = 100$ steps (line 6 of Algorithm 1); and that the neural optimisers themselves

¹⁸Refer to Section 2.1 on page 4 of Andrychowicz et al. [2016].

Algorithm 1 Training the LSTM-optimiser of [Andrychowicz et al. \[2016\]](#)**Input:**

Learner \mathcal{F} with parameter θ , LSTM-optimiser \mathcal{Q} with parameter η ,
and training dataset $\mathcal{D}_{\text{Train}}$.

Hyper-parameters:

Training-trial \mathcal{V} , training-step \mathcal{S} , unroll-instance Ξ , and learning rate ω for \mathcal{Q} .

```

1:  $\tau = 0$ 
2:  $\mathcal{D} \leftarrow$  set up dataset
3:  $\eta \leftarrow$  random initialisation ▷ For the neural optimiser  $\mathcal{Q}$ 
4: for  $q = 1, \mathcal{V}$  do ▷ Updating  $\mathcal{V}$  learners
5:    $\theta_0 \leftarrow$  random initialisation ▷ For the  $q$ th  $\mathcal{F}$ 
6:   for  $t = 1, \mathcal{S}$  do ▷ Update for  $\mathcal{S}$  steps
7:      $\mathcal{X}_t, \mathcal{Y}_t \leftarrow$  sample batch from  $\mathcal{D}_{\text{Train}}$ 
8:      $\mathcal{L}_{t+1}(\theta_t) \leftarrow \mathcal{L}(\mathcal{F}(\mathcal{X}_t|\theta_t), \mathcal{Y}_t)$  ▷ Acquire  $\mathcal{L}$ 
9:      $\tau \leftarrow \tau + 1$ 
10:     $\mathcal{L}_{\tau+1}(\theta_\tau) \leftarrow \mathcal{L}_{t+1}(\theta_t)$  ▷ Record  $\mathcal{L}$ 
11:    if  $\tau = \Xi$  then ▷ Unroll  $\mathcal{Q}$ 
12:       $\mathcal{L}(\eta) = \mathbb{E} \left[ \sum_{\tau=1}^{\Xi} \mathcal{L}_{\tau+1}(\theta_\tau) \right]$  ▷ Acquire  $\mathcal{L}$ 
13:       $\eta \leftarrow \eta - \omega \nabla_{\eta} \mathcal{L}(\eta)$  ▷ Update  $\eta$ 
14:       $\tau = 0$  ▷ Reset records
15:    end if
16:     $\mathcal{X}_t^{(\mathcal{Q})} \leftarrow \nabla_{\theta} \mathcal{L}_{t+1}(\theta_t)$  ▷ Prepare the input of  $\mathcal{Q}$ 
17:     $g_{t+1} \leftarrow \mathcal{Q}(\mathcal{X}_t^{(\mathcal{Q})}|\eta)$  ▷ Compute the update to  $\theta_t$ 
18:     $\theta_{t+1} \leftarrow \theta_t + g_t$  ▷ Update  $\theta$ 
19:  end for
20: end for

```

are updated (unrolled) once every $\Xi = 20$ steps (line 11 of Algorithm 1). During meta-testing, the neural optimisers are tested to update initialised base learners for the lengthier 1000 steps to test their robustness against potential unknown sequential errors that can occur during the consecutive parameterisation.

3.2.5.7 Related Work to the LSTM-Optimiser

Thus far, we have provided a detail description on [Andrychowicz et al.](#)'s L2L task. However, the student-author have decided to not present the related work to the LSTM-optimiser in this section. The reasons are justified as follows. LSTM-optimisers possess a lot of undesirable features and most follow-up work (including our own work) were targeted to address those issues. Thus a discussion of those problematic characteristics will instead be covered in Section 5.1 of which we delay the related studies to Section 5.1.6.

3.2.6 Section Conclusion

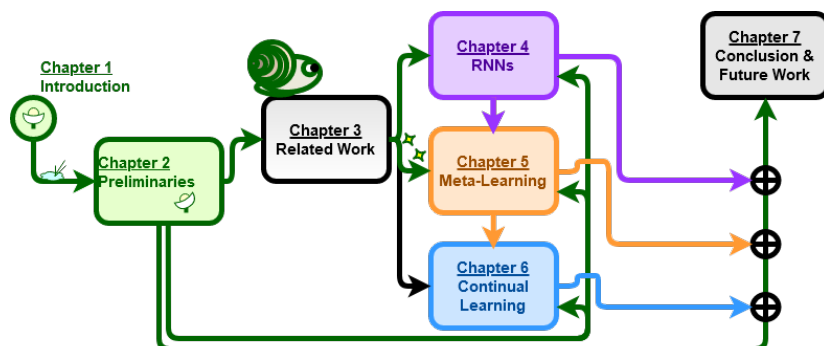


Figure 3.6: Progress Update

This section provided a general overview on meta-learning. We discussed that meta-learning techniques can be dichotomise as offline techniques and online techniques. We also covered the three main approaches in online techniques including the metric-based, model-based, and optimisation-based archetypes.

Most importantly, we defined the experimental setup of [Andrychowicz et al. \[2016\]](#)'s L2L. This is our meta-learning method of choice and this particular technique was chosen because that it fit well with our research interest in RNNs (see Chapter 4) and in understanding the backwards pass as a dynamical system (see Section 2.3.5.1 and Chapter 2).

Please refer to Chapter 5 on our novel contribution to meta-learning. Our research topics are on making neural optimisers more generally applicable (see Section 5.1), and on expanding the utility of neural optimisers beyond the sole purpose of rapid knowledge acquisition (see Section 5.2).

*“The distinction between thought and reality is possible
only when we know how to distinguish between
what the ‘I’ can change and what the ‘I’ cannot change. ...
What we then attain is knowledge of the lawful order in the realm of the real,
but only in so far as it is represented in the tokens within the system of sensory impressions”*
[Von Helmholtz \[1978/1903\]](#)¹⁹

*Humans have the ability to adapt new knowledge while not forgetting the old.
However, such is not the case for connectionalist models like neural networks.
The following section will discuss the phenomenon of catastrophic forgetting and existing
methods to prevent it from occurring.*

Section Content:

- Section 3.3.1: The Framework of Continual Learning,
- Section 3.3.2: Comparing Continual Learning to
Other AutoML Directions,
- Section 3.3.3: Continual Learning vs the Brain, RNNs, and the Self,
- Section 3.3.4: Senarios for Continual Learning,
- Section 3.3.5: Archetypes of Continual Learning,
- Section 3.3.6: Our Method of Choice,
- Section 3.3.7: Section Conclusion.

¹⁹Also see an analysis of [Von Helmholtz](#)’s ideas in [Westheimer \[2008\]](#).

3.3 Continual Learning & On Episodic Memories

Continual learning studies the sequential acquisition of functions for a single neural network. This remains a difficult problem because *Catastrophic Forgetting* [McCloskey and Cohen, 1989; French, 1999] occurs when a trained network learns a new skill. Catastrophic forgetting refers to the abrupt loss of knowledge for solving an old task when information for solving a new task is incorporated; and this is due to parametric modifications for satisfying the learning objectives of a new task.

Though the research community has had interest in continual learning since 1989 (and earlier), continual learning has been gathering a serious amount of popularity post-2015 due to the significant breakthroughs of AlexNet [Krizhevsky et al., 2012] and ResNet [He et al., 2016a]. After neural networks were shown to achieve human-level of image classification ability, the research community has now turned to pursue general artificial intelligence. Continual learning is seen as one (if not the only) route towards achieving this goal [Lomonaco et al., 2020].

3.3.1 The Framework of Continual Learning

In continual learning, we employ a task-specific backbone network \mathcal{F} to sequentially learn the n tasks of $\omega_1, \dots, \omega_n \in \Omega$ in a *continuum* of data

$$\mathcal{K} = \{(x_1^1, \omega_1, y_1^1), \dots, (x_\gamma^i, \omega_i, y_\gamma^i), \dots, (x_\Gamma^n, \omega_n, y_\Gamma^n)\} \quad (3.16)$$

to map inputs x_γ^i to labels y_γ^i from datasets $(x_\gamma^i, y_\gamma^i) \in \mathcal{D}_{\omega_i}$ for all Γ data where $\gamma = 1, \dots, \Gamma$. The aim is to yield a learner that generalises well on subsequent tasks while it maintains high accuracy for learnt tasks.

Continual learning is challenging [Kemker et al., 2018]. When a learner *only* observe data of a new task, its optimised parameters for an old task will be modified to suit the learning objective of the new task.

3.3.2 Comparing Continual Learning to Other AutoML Directions

Continual learning is widely inspired by immediately adjacent learning paradigms in AutoML research. Nonetheless, the continual learning experimental setup has some unique and defining characteristics. For instance, continual learning shares similarities to multi-task learning [Caruana, 1997] which aims to concurrently improve all learning tasks by aggregating the common extracted knowledge from all tasks. The extracted knowledge of different tasks serve as a form of regularisation for other tasks to prevent overfitting. While multi-task learning needs all data of all tasks, the (ideal) continual learning setup only requires data for the new tasks. Thus continual learning is also naturally linked to transfer learning [Pan and Yang, 2009] where some tasks are purposely reserved and not simultaneously optimised. However, catastrophic forgetting spontaneously occurs as a consequence of domain adaptation for fine-tuning feature extraction. In stark contrast, continual learning is

interested in keeping the skills for solving old tasks when (nearly) all training data of the original tasks of the source domains are not available²⁰.

There is also strong affinities between continual learning and meta-learning (see Section 3.2 and Chapter 5). On one hand, meta-learning uses past performances on different tasks to recommend the best algorithm for new tasks in order to achieve rapid knowledge acquisition. On the other hand, continual learning keeps the best single model and aims to continuously expand the proficiency of a single network. This is hence related to topic relevance methods, like knowledge distillation [Hinton et al., 2015], which aims to transfer knowledge between multiple networks.

Continual learning also has overlapping themes with research outside of AutoML. Below, we will briefly discuss some commonalities it shares with biology and philosophy, and similar ideas in early RNN studies (see Section 3.1 and Chapter 4).

3.3.3 Continual Learning and Potentially Related Concepts in RNNs

Sitting at the heart of continual learning is the *Stability-Plasticity Dilemma* [Mermilod et al., 2013]. An intelligent agent needs plasticity to learn new knowledge but it also requires stability to prevent forgetting. For humans, neural circuits exhibit profound plasticity during early life and of which are later stabilised [Takesian and Hensch, 2013]. In our windows of early life, there are critical periods when the brain aims to find an optimal ratio of excitatory and inhibitory circuit activity [Hensch, 2004]. Beyond such critical periods, there exist factors that limit excessive circuit rewiring [Bavelier et al., 2010]. Interested readers should consult the recent work of Hayes et al. [2021]; the said paper addressed the similarities and shortcoming of common continual learning experimental setups to bio-physiological facts.

The biological framework that has inspired many continual learning studies also aligns nicely to the bedrock philosophy of Schmidhuber’s RNN research. Recall from Section 3.1.3, some of Schmidhuber’s early work entertained the concept of creating networks that could implement their own parameters. In order to achieve that goal, the author required a network to access domain concept²¹ to temporarily modify the existing network configuration. Furthermore, such temporary modifications were also required to be (quasi-)recoverable so that a network would not forget what it had already learnt. We argue that Schmidhuber’s early work on memory-recoverable self-inferential networks paved the way for the creation of the LSTM RNN [Hochreiter and Schmidhuber, 1997b] (see Section 3.1.2). Consider Equation (3.7), where the cell state is incrementally updated via $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$. The packages of updates \mathbf{a}_t are as if new information that are assimilated into the old knowledge \mathbf{c}_{t-1} carefully fine-tuned by the gated units \mathbf{f}_t and \mathbf{i}_t ²². That is, the gated units function as if they

²⁰Interested readers can find more descriptions regarding the specificity of the differences on pages 2 to 4 in Li and Hoiem [2017].

²¹Such as the MANNs in meta-learning, see Section 3.2.3.

²²The LSTM cell \mathbf{c}_t is transformed as the hidden state \mathbf{h}_t to create the output gate \mathbf{o}_t to control output exposure like that of a biological system. As noted in Section 3.1.6.1, Greff et al. [2016]’s ablation study found that the output gate served as one of the most important gated units to the well-being of LSTMs.

were artificial organs that find the optimal rate of graceful forgetting to allow new inferences could be built on top of existing network extracted concepts.

As mentioned earlier, Schmidhuber’s work was also heavily influenced by philosopher Von Weizsäcker [1985]²³ of which they treated *the process of learning* as an inherent Bayesian inferential process. This style of conceptualisation is also congruous with continual learning – by theorising learning as a non-memory-less procedure, it shows that catastrophic forgetting is not inevitable. As we will later discuss in Section 3.3.5, many important work in continual learning had been heavily inspired by the Bayesian learning framework to address feature drift in observing a series of consecutive tasks.

Empirically however, catastrophic forgetting remains a major hurdle in machine learning. It is also a sign such that conventional models exhibit high plasticity for the integration of new knowledge but suffer from low stability to maintain their acquired experiences on solving previously learnt tasks. The converse is also shown to be true in the empirical studies of Mirzadeh et al. [2020b].

Continual learning is a fast evolving field with new nomenclatures, new benchmarks, and new evaluation metrics being constantly proposed. We will first discuss some common continual learning setups and then cover the 3 main archetypes in existing literature. Again, we will reserve the approach of our choice at the end (see Section 3.3.5.3).

3.3.4 Scenarios for Continual Learning

As mentioned, continual learning has become a topic of intense research post-2015. As a result, new jargon, evaluation metrics, and experimental setups have been constantly proposed. The continual learning research community has recently been working towards experimental standardisation. For one instance, the *Continual AI* research groups has been working closely with multiple top machine learning conferences to define research challenges and prepare datasets²⁴. For another instance, the work of Van de Ven and Tolias [2018] identified the 3 most common continual learning experimental setups and highlighted their differences. The discussion below will mostly follow the narrative of Van de Ven and Tolias.

Recent literature in continual learning sees both the usage of the “single-head” and the “multi-head” classifier setups. Under the former setup, one network is being continuously trained over a series of tasks. That is, both the feature extractor (everything from layer 1 to the penultimate layer) and the classifier (the last/softmax layer or decoder) are concurrently trained throughout the course of sequentially observing multiple tasks. While for the latter case, a unique classifier (a head) is defined for

²³See Von Weizsäcker’s quote in footnote 3 in Section 3.1.3.

²⁴Details of Continual AI can be found in <https://www.continualai.org>. The group has collaborated closely with the *Conference on Computer Vision and Pattern Recognition* (CVPR) since 2020 [Lomonaco et al., 2020] (also see <https://sites.google.com/view/clvision2020>), and collected a challenging continual learning dataset [Lomonaco and Maltoni, 2017]. Recently in 2021, the same group has further established connections with the *International Joint Conference on Artificial Intelligence* (IJCAI) (also see <https://sites.google.com/view/sscl-workshop-ijcai-2021/>).

every dataset. To elaborate, the multi-head setting only continually trains the feature extractor and this is hence much easier than the single-head setup because that catastrophic forgetting is partly prevented without the need for constant re-modification on a single unified classifier.

Besides the classifier, the community also lacks a consensus on how to present data points of different datasets to a continually learning backbone network. [Van de Ven and Tolias](#) noted that there are typically 3 different ways in feeding the input; and they are the *Task-Incremental Learning* (IL), *Domain-IL*, and *Class-IL*.

Task-IL aims to solve all observed tasks while knowing the exact *Task Boundary*. That is, the backbone \mathcal{F} of Equation (3.16) receives task labels ω_i along with inputs x_γ^i and targets y_γ^i . This provides the backbone network with information regarding its available choice of output classes. Of the two classifier setups, the multi-head classifier layout naturally adopts the Task-IL experimental setup as the backbone network requires the task ID to select the correct task-wise classifier during test time.

Domain-IL is slightly harder and is typically applied along the single-head classifier setup. The classifier output dimension is fixed and that each dataset consists of an equivalent amount of target classes. To elaborate, every single entry of the classifier output is reused for identifying multiple targets dependent on the context of the observed data. Task boundary is also not given.

Class-IL is the most difficult setup and its classifier subnetwork adopts an expanding architecture. Say that the classifier weights while observing V tasks is $\mathcal{W} \in \mathbb{R}^{N_I \times N_O}$ where N_I and N_O are the input and output dimension respectively, then when proceeding to observe Task $V + 1$ with 2 new classes, Class-IL defines a new weight $\overline{\mathcal{W}} \in \mathbb{R}^{N_I \times N_{O+2}}$ and distils knowledge from \mathcal{W} to $\overline{\mathcal{W}}$. Thus unlike Domain-IL, the amount of classes in each dataset do not necessarily need to be fixed. Furthermore, not only is the task boundary not given, the base learner also needs to infer the correct task ID from the input data in order to correctly identify the corresponding output dimensions used during test time²⁵.

In this thesis, we will target our attention on the simpler Task-IL and Domain-IL scenarios. Later in Section 6.1, we will show that many methods, including those that claim state of the art status, can easily perform abysmally with minor changes to their standard experimental setups.

3.3.5 Archetypes in Continual Learning

There are 3 main archetypes in continual learning literature. They are the *Regularisation* approach, the *Dynamic Architectural* approach, and the *Memory* approach.

3.3.5.1 Continual Learning via Regularisation

Regularisation-based continual learning assists a backbone network to sequentially learn via reformulating the loss function. That is, the original loss (say cross entropy

²⁵Interested readers can refer to a clean implementation of the Class-IL scenario in Arthur Douillard's repository https://github.com/arthurdouillard/incremental_learning.pytorch/blob/master/inclearn/lib/network/classifiers.py line 103 from his work [Douillard et al. \[2020\]](#).

loss for image classification) is appended by an auxiliary loss to ensure that knowledge encoded within a network is not lost through observing a stream of information with a changing distribution.

There are many ways to design the auxiliary loss. For one instance, [Jung et al. \[2016\]](#) defined this loss as the L2 norm of the backbone network features before and after a step of gradient update. In a similar spirit, [Li and Hoiem \[2017\]](#) adopted [Hinton et al. \[2015\]](#)'s knowledge distillation (KD) as the auxiliary loss. Knowledge distillation treated the data inferential process as a black box and its goal, regardless of how the weights had changed during sequential learning, was to force the backbone network to attain similar network output distributions. In another paper, [Zenke et al. \[2017\]](#) took a much more direct approach and defined the auxiliary loss as a function that measures the direct differences in the changes in backbone weights.

Another line of regularisation-based approaches provided more theoretical aspects towards continual learning with multiple target objectives. For instance, [Kirkpatrick et al. \[2017\]](#) considered continual learning as a Bayesian inferential process (also see Section 3.3.2). They first compared gradient descent as a probabilistic problem, and then extended this idea to compare sequential learning as a procedure of acquiring new posterior probabilities for the parameters. Though their auxiliary loss was similar to that in [Zenke et al.](#)'s work, they provided a more structured formulation of the continual learning problem. [Kirkpatrick et al.](#)'s framework was expanded by [Nguyen et al. \[2018\]](#), which the authors included a small representation dataset of old tasks to sequentially learn a variational autoencoder (VAE) generative model.

A commonality shared between [Nguyen et al.](#)'s VAE and [Li and Hoiem](#)'s KD was KL divergence. The geometric meaning of KL divergence for continual learning was recently studied in [Cha et al. \[2020\]](#). The authors linked the usage of KL divergence to single-task wide local minima studies such as [Pereyra et al. \[2017\]](#) (also see Section 2.2) and found that KL divergence could project the conditional probabilities of backbone outputs to uniform distributions. More importantly, they demonstrated that KL divergence could in fact be added as a general patch to any regularisation-based continual learning technique to enhance their performances. To elaborate, they suggested that a 3-loss combination including the original loss (cross entropy), the surrogate loss (such as L2 norm), and a KL divergence loss would always function better than just the original loss appended with the surrogate loss.

3.3.5.2 Continual Learning via a Dynamic Architecture

Another way to address the plasticity-stability dilemma is to increase the amount network parameters to host new knowledge whenever a new task is presented. A naïve network expansion was studied in [Rusu et al. \[2016\]](#). In that paper, the authors started sequential learning with a simple backbone network; then upon the introduction of a new task, they instantiated a new neural network while freezing the weights of the old backbone network. In addition, they laterally transferred features extracted from the previously learnt subnetwork.

Though [Rusu et al.](#)'s method prevented catastrophic forgetting, the model size

scaled linearly to the amount of new tasks and was hence unsuitable in practice. In [Lomonaco and Maltoni \[2017\]](#), the authors showed that network expansion could be limited to only instantiating a new classifier network per task, which substantially decreased the computational cost of sequential learning. This framework was further studied by [Yoon et al. \[2018\]](#) and the authors presented a protocol to dynamically extend the backbone network architecture. [Yoon et al.](#)'s approach reiteratively re-trained a sequential learning backbone network. And when the loss of the network stagnated (or worsened) for a pre-defined amount of time steps, the authors increased the network by a small set of parameters. Notably, the objective loss of [Yoon et al.](#)'s backbone appended the original loss (cross entropy for image recognition) with L1-regularisation to promote sparsity in the weights. This coerced the heuristic learning of the backbone to favour sparse representations to promote a selective retraining of learnt parameters. This in turn allowed less parameters to be modified and thus lowered the impact of catastrophic forgetting.

[Yoon et al.](#)'s sparsity promotion could also be interpreted as a way to reduce representational overlap. This is another popular strategy to overcome catastrophic forgetting and could be dated back to [French \[1991\]](#) before the recent revived interest in continual learning. In their paper, [French](#) proposed to post-process hidden units with activation sharpening to force the majority of network hidden units as zero while keeping only a few as non-zero. In another recent study, [Serra et al. \[2018\]](#) showed that neural values could also be made sparse via learning a nullification mask per task. Their masks were learnt from input embedding (like word embedding in language modelling [[Merity et al., 2018](#)]) to achieve task-dependent attention [[Bahdanau et al., 2015](#)]. [Parisi et al. \[2018\]](#) showed that similar nullification masks could also be acquired through RNNs. The authors designed RNNs that embedded Hebbian-like activations [[Parisi et al., 2017](#)] to recursively update the spatio-temporal sparsity of their backbone hidden states.

[Parisi et al.](#)'s method utilised RNNs in a similar fashion to memory networks in meta-learning (see Section 3.2.3) where a specifically crafted subnetwork was conjointly trained with the backbone to enhance the latter's decision making procedure. Their network design hence echoed with the dual-memory model of mammalian memory [[McClelland et al., 1995](#)]. This idea was further explored in the work of [Kemker and Kanan \[2018\]](#) which the authors introduced a framework for the backbone network that separated the feature extractors from the classifier and explicitly compartmentalised a subnetwork for storing long-term memory.

3.3.5.3 Continual Learning via Memory Replay

The third archetype in preventing catastrophic forgetting is through the concurrent training of data of past tasks with those of the new task. Metaphorically, memory is *replayed* to the backbone network. However, naïvely applying replay will incur high memory demand and severe computational cost. The former undesired feature emerges because an intelligent agent needs to externally store a set of data for every past task and memory is thus scaled linearly. The latter displeasing characteristic

appears because the concurrent learning from both the old and new datasets requires a backbone network to simultaneously process more data (from a larger batch) for every time step. Fortunately, many memory-based continual learning techniques have already addressed these problems.

In the pioneering work of [Robins \[1995\]](#), the author proposed to use *pseudo replay* to lower the memory cost. Pseudo replay shared many similarities to knowledge distillation (see Section 3.3.5.1), and the main idea was to preserve a backbone network’s learnt mapping via uniformly sampling random inputs and their corresponding outputs and replaying them along with the new task samples. This approach was adopted by other early researchers of continual learning; for instance, [French \[1997\]](#) extended pseudo replay with recurrent memory in a similar fashion to the architectural-based approaches mentioned in Section 3.3.5.2. However, pseudo replay was demonstrated to only be effective with inputs from small binary input spaces and that it failed to preserve the learnt mappings of backbone networks in high-dimensional spaces [[Kamra et al., 2017](#)].

One line of recent research abandoned the idea of using pseudo replay in favour of replaying the backbone network with exact past data points. In one influential study, [Rebuffi et al. \[2017\]](#) proposed to prevent the memory issue in exact replay by selecting a small *exemplar set* of data points with prototypical features. Their method consisted a representation learning protocol that updated the representative and externally stored exemplar set to remove and update content within the storage. However in another important study, [Lopez-Paz and Ranzato \[2017\]](#) showed that replaying data of all past tasks from all task-wise buffers could result in significant performance gain against [Rebuffi et al.’s](#) method. [Lopez-Paz and Ranzato’s](#) method introduced the idea of *gradient alignment and interference* – when the native gradients yielded from replaying past data misaligned with the native gradient from learning data of the current task, catastrophic forgetting would occur as the learnt parameters of the backbone network would have attempted to leave low loss regions that solved all previous tasks. In order to prevent that, [Lopez-Paz and Ranzato](#) pre-conditioned the native gradient of the new task with those gradients of the old tasks to project parametric update into a compromised zone with less parametric over-modification. Their method has since become one of the most popular methods in continual learning with several follow up studies. Nonetheless, [Lopez-Paz and Ranzato’s](#) approach remained computationally inefficient as all memory of all past tasks were required to be replayed per gradient update.

Since its introduction, [Lopez-Paz and Ranzato’s](#) framework has been continually streamlined and that two recent contributions from [Chaudhry et al.](#) have made memory-based continual learning computationally efficient. First, the authors showed in [Chaudhry et al. \[2018\]](#) that [Lopez-Paz and Ranzato’s](#) gradient pre-conditioning scheme did not necessarily require all past data. One could simply sampling a mini-batch from the union of all old memory and that this alone contained sufficient information to project native gradients of the new task to prevent catastrophic forgetting. In a second paper [[Chaudhry et al., 2019](#)], the same panel of authors further showed that continual learning could be prevented by co-training

as little as one past data per update. Their experiments revealed that by directly co-training the exact replay with gradient descent, without employing the gradient projection scheme in [Lopez-Paz and Ranzato](#), backbone networks could achieve near-0 catastrophic forgetting. Though the computational issue in memory-based techniques was successfully addressed by [Chaudhry et al. \[2019\]](#), their mini-batches were still required to be sampled from a large externally stored memory buffer. Their method hence suffered from being memory demanding.

Reducing memory consumption is one of the main research topics in *generative replay*. In this school of thoughts (including [Shin et al. \[2017\]](#), [Draeos et al. \[2017\]](#), and [Kamra et al. \[2017\]](#)), researchers build generative models into their backbone networks to synthesise pseudo data points as replacements for exact replays. This particular framework was mainly inspired by the human hippocampus and the *Complementary Learning Systems* (CLS) theory [[McClelland et al., 1995](#)]; and that the built-in generative model consolidated memory by providing the rest of the backbone with constructed data points that captured general semantics of observed tasks. Though memory consumption could be relieved via this approach, this technique would fail if the generative model itself suffered catastrophic forgetting. Furthermore, the inclusion of a generative model would also directly increase the computational cost for sequential learning.

3.3.6 Our Method of Choice

In this thesis, we will mainly focus on the memory-based techniques discussed in Section 3.3.5.3. This is mainly because of two reasons. First and as noted in [Lopez-Paz and Ranzato \[2017\]](#) and [Chaudhry et al. \[2019\]](#), regularisation-based and dynamic architectural-based continual learning are less effective in preventing catastrophic forgetting. More notably, [Van de Ven and Tolias \[2018\]](#) showed that regularisation-based techniques failed to function under the Class-IL experimental setup. Our second reason is that as according to a recent report on the continual learning challenge of the prestigious CVPR venue [[Lomonaco et al., 2020](#)], all top performing methods belonged to the memory-based archetype. More importantly, all top contenders had adopted an experimental setup similar to that in [Chaudhry et al.](#)'s work. The setup of [Chaudhry et al.](#)'s work will be addressed in Section 6.1.2.

3.3.7 Section Conclusion

This section discussed continual learning. We mentioned how catastrophic forgetting remains a major hurdle towards realising general artificial intelligence, and that this is mainly due to the plasticity-stability dilemma of connectionist models including machine learning algorithms. Furthermore, we covered the three main approaches of regularisation-based, dynamic architectural-based, and memory-based techniques.

In addition, we stated that our effort in investigating continual learning will mainly be targeted towards memory-based approaches. This is because that memory-based approaches have continuously outperformed other rival techniques and have

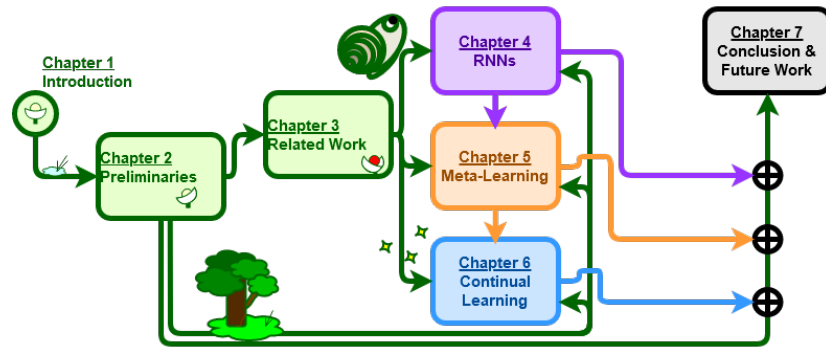


Figure 3.7: Progress Update

been constantly streamlined.

Refer to Chapter 6 on our novel contribution to continual learning. Our research topics are dedicated to reducing the memory consumption of replay methods, and on proposing a simple and general patch adaptable by most (if not all) existing continual learning techniques.

Understanding RNNs with Dynamical Systems

This chapter introduces our contribution on RNN research. The content therein this section serves as a continuation to Section 1.5.3 of the introduction, Section 2.1 of the preliminaries, and Section 3.1 of the related work.

In the introduction, we stated that we chose to understand shortcut connections [Schraudolph, 1998] to design better AutoML techniques. To continue, the preliminary explained that the state-of-the-art CNN of highway network [Srivastava et al., 2015] was inspired by the LSTM RNN [Hochreiter and Schmidhuber, 1997b] to include shortcut connections to facilitate deep learning. Thus we need to understand LSTMs in order to perceive the functionality of shortcut connections.

However in the related work, we mentioned that understanding RNNs was made difficult by the fact that most papers concerned RNNs for their utility over network interpretability. And for those work that analysed RNN behaviours, their analyses were centred around interpreting trained networks to conjecture *what has been heuristically learnt* (see Section 3.1.6). These papers contained a rich amount of information; but they did not directly provide answers to two of our questions on LSTMs to help us to perceive shortcut connections. Namely, we want to know

- (i) *What is the mathematical purpose of each component in the forward pass? &*
- (ii) *Can we mathematically explain how to design a simple high functional RNN?*

Question (i) shows that we are less interested in what can be heuristically learnt; but rather, how to ensure the importance of each LSTM component prior to an RNN being trained. Question (ii) indicates that after we find the core LSTM mechanism, we desire to simplify the network while maintaining competitive performances.

The two questions above are answered in our papers of Kuo et al. [2018] and Kuo et al. [2020b] (also see Section 1.5.3). The former paper takes an atomic-approach to analyse the individual importance of every neuron in the LSTM memory. However, we will discuss the disadvantages of the atomic-approach and eventually come to favour the system-approach of Kuo et al. [2020b] to describe LSTM mechanics from a macroscopic level. The commonality of both papers is that they are both based on dynamical systems (see Section 2.3) for describing the incremental changes that were made, learnt, and retained in the LSTM memory.

*“non sunt multiplicanda entia sine necessitate” Duns Scotus (1639)¹
– entities must not be multiplied without necessity.*

*The following section aims to follow the Occam’s Razor parsimony principle and
look at neural network from a atomic point of view.*

*By doing so,
we hope to understand the most significant interaction within a neural network
with neurons being the simplest unit.*

Section Content:

- Section 4.1.1: Section Introduction,
- Section 4.1.2: LSTMs as Difference Equations,
- Section 4.1.3: The Cell State Revisited,
- Section 4.1.4: DecayNets,
- Section 4.1.5: Section-wise Additional Related Work,
- Section 4.1.6: Experiments and Results,
- Section 4.1.7: Section Conclusion:
DecayNet and an Analysis on the LSTM Cell,
- Section 4.1.8: An Evaluation of this Study.

¹See Williams [2019] for an introduction of the Scotus’s character; and also see the interesting blog of “Who Sharpened Occam’s Razor” at <https://www.irishphilosophy.com/2014/05/27/who-sharpened-occams-razor/>.

4.1 A Study on the Cell States of Long Short-Term Memories

This section is based on our work in [Kuo et al. \[2018\]](#) to understand how an LSTM updates its recursively used cell state.

4.1.1 Section Introduction

The complicated LSTM design creates difficulties in network interpretability and thus its utility on real-world data is non-immediately clear. We will take a theoretical mathematical approach to study the atomic behaviours on consecutive transitions in an LSTM. The aim is to become able to build better designs with higher transparency.

The nature of this work is a numerical study which examines an LSTM as *a system of difference equations*. This is because that an LSTM repeatedly integrates new observations into its implicit hidden variables to model temporal representations of input sequences. Like dynamical systems, the network undergoes iterative computations on a same set of operators. The transitional dynamics can be visualised through cobweb diagrams (see Section 2.3.2); and numerical analyses allow us to predict *shrinkage* and *growth* in the variables therein the network.

The main result of this section is to show that the updating scheme of the LSTM cell state volatily alternates between a *catch* and a *release* phase. The former controls the shrinkage while the latter controls the growth of cell state magnitudes. Our study reveals that the forget gate predominantly dictates the dynamical alterations of the cell state. We hence introduce a *decay* mechanism to monotonically decrease forget gate values to stabilise cell state propagation for making LSTMs more interpretable. Our experiments showed that LSTMs with decaying forget gates yielded state of the art results on the fixed-length *Permuted Sequential MNIST* task [[Le et al., 2015](#)].

4.1.2 LSTMs as Difference Equations

Recall that LSTMs are recursive systems driven by inputs \mathbf{x}_t and propagate for a data-length of D units of time, i.e., $t = 1, \dots, D$. Each instance computes for

$$\begin{aligned} \text{the gated variables:} \quad & \mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,I,O\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{F,I,O\}}), \\ \text{the internal input:} \quad & \mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{h}_{t-1} + \mathbf{b}_A), \\ \text{the cell state:} \quad & \mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and} \\ \text{the hidden state:} \quad & \mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \end{aligned}$$

As we mentioned in Section 3.1.2, the main difference between an LSTM and a simple RNN unit is that it employs multiple gated units to synchronise the deletion and replenishment of new information for the update of cell state \mathbf{c}_t . In order to study how an LSTM repeatedly writes and rewrites the cell values, we propose to compare LSTMs as *difference equations* [[Glendinning, 1994](#)] (also see Section 2.3.1).

Difference equations are dynamical systems with *discrete* recurrent propagation. Their behaviours can be studied to determine whether network variables will grow or shrink as time unfold. By modelling the LSTM cells as a system of difference equations, this allows us to explain how the gated units increase and decrease the cell state values and thus helps us to clarify network mechanics.

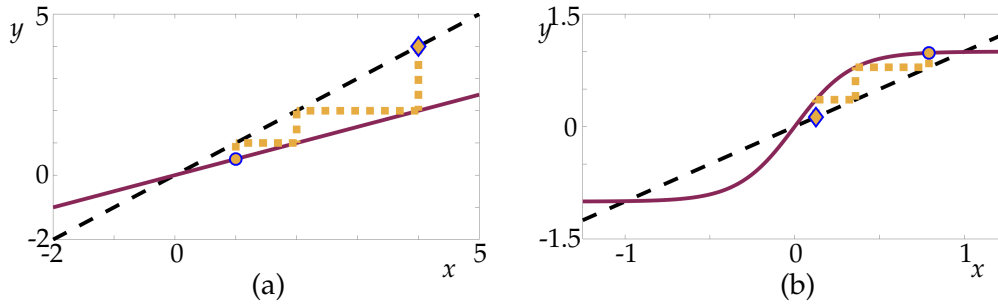


Figure 4.1: Examples of Difference Equations Variable Propagation

4.1.3 The Cell State Revisited

This section discusses three key concepts – the importance for studying the LSTM cell state; the two dynamic regimes for cell state values to grow and to shrink; and that the forget gate can be controlled to yield interpretable and predictable cell state motions. We will use cobweb diagrams to support our concept visualisations, and refer to Section 2.3.2 for a brief introduction of the technique.

4.1.3.1 The Importance of the LSTM Cell State

The work therein this section postulates that cell state \mathbf{c}_t acts as the source of LSTMs' modelling capability. Furthermore, we regard gated variables only as secondary variables. Our three justifications are as follows.

First consider $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$ and $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$, we see that gated variables only exist to read and to set the hidden variables of \mathbf{h}_t and \mathbf{c}_t . That is, they exist as auxiliary components to assist an LSTM to memorise better, but they are not reused iteratively as recurrent inputs and hence are not network memories. Next consider $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$, since \tanh is a one-to-one function then the dynamic properties of \mathbf{h}_t are inherited from \mathbf{c}_t . Thus by analysing \mathbf{c}_t , we will gain sufficient information for understanding \mathbf{h}_t . Last, we consider $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$ where we find \mathbf{a}_t included in the update of \mathbf{c}_t . Thus by studying changes in \mathbf{c}_t , we will automatically concern changes in \mathbf{a}_t .

For these reasons, we postulate that the cell state serves as the most important variable of LSTM. Since it also happens to be designed with discrete incremental changes, we thus choose to study its recurrent behaviours with the dynamical system of difference equations.

4.1.3.2 A Piece-wise Analysis on the Cell State

We study the cell state, $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$, separately as $\mathbf{c}_p = \mathbf{f}_t \odot \mathbf{c}_{t-1}$ and $\mathbf{c}_q = \mathbf{i}_t \odot \mathbf{a}_t$. Since we have asserted gated variables as secondary variables, we can think of \mathbf{f}_t and \mathbf{i}_t as constant vectors with elements between 0 and 1. Furthermore, the individual dimensions of \mathbf{c}_p and \mathbf{c}_q have the cobweb outlines analogous to those in Figures 4.1(a) and 4.1(b) respectively. In addition, we will refer to the marginal

dimensions with the subscript of m . For instance, each dimension of \mathbf{c}_p is $i_{mt} a_{mt}$ and thus inherits the characteristic outline of Figure 4.1(b) through a_{mt} 's tanh function.

Figure 4.1(a) and 4.1(b) are not the exact cobweb outlines of \mathbf{c}_p and \mathbf{c}_q ; nonetheless, they serve as easy and direct visual cues. The dynamic regimes of \mathbf{c}_p and \mathbf{c}_q are similar to (a): $x_t = T(x_{t-1}) = 0.5x_{t-1}$ and (b): $x_t = T(x_{t-1}) = \tanh(3x_{t-1})$. The variable propagation are shown in dots and they start from the rhombuses and end on the circles. (See Section 2.3.2 for a more elaborated explanation on the symbols.) The diagrams reflect that the origin in subplot (a) is a sink which attracts propagation and diminishes the magnitudes of variables; whereas its counterpart in subplot (b) is a source which repels propagation and increases variable magnitudes. Due to these reasons, we refer to the phenomenon in Figures 4.1(a) as the dynamic *catch regime* and refer to that in Figure 4.1(b) as the dynamic *release regime*.

Also note that the dynamic regimes are highly related to the gradients around the origin of their respective plots. This is because that the origin is an *equilibrium*, and that the gradient determines whether the said equilibrium is attractive or repulsive. More details regarding these concepts can be found in Section 2.3.4.

The update of the LSTM cell is thus composed of two opposing modifiers. With one that dedicates to its growth, and the other specialised in its shrinkage. It is thus crucial to understand which one between \mathbf{c}_p and \mathbf{c}_q serves as the more dominant evolution component. We will then reformulate the dominant component with deterministic properties to provide the network with affirmative qualities.

4.1.3.3 The Forget Gate and the Catch-and-Release Dynamics

A key difference between the opposing regimes of Figure 4.1 is the boundedness of the characteristic outlines. The sub-dynamic of \mathbf{c}_q is limited between ± 1 , while that of \mathbf{c}_p is unbounded. This implies that \mathbf{c}_q contributes less to the additive relationship of \mathbf{c}_t and is thus less impactful than \mathbf{c}_p . With the dominant component identified, our next target is to understand how \mathbf{c}_p affects vanilla LSTMs. The insight will allow us to propose modifications to make LSTM mechanics more interpretable.

Since $\mathbf{c}_p = \mathbf{f}_t \odot \mathbf{c}_{t-1}$ serves as the dominant modifier to \mathbf{c}_t and that gradient f_{mt} modifies the characteristic outline of $c_{m(t-1)}$, we conjecture that the forget gate controls the *overall* gradient of the outline of the cell state. That is, as \mathbf{f}_t predominantly controls the stability of the origin, it dictates the alterations between the catch and the release phases. As a consequence and regardless of \mathbf{c}_q , values of the cell state are more likely to grow with large \mathbf{f}_t , and more likely to shrink with small \mathbf{f}_t .

These conjectures are supported and highlighted with Figure 4.2 overleaf. Three qualitatively different scenarios for cell state propagation are investigated. The images largely follow the format of Figure 4.1 but we introduced additional dot-dash lines to show the characteristic shape of the $t - 1$ th instance. The characteristic shapes depict $y = T(x) = f_{mt}x + i_{mt}\tanh(x)$ to represent the cell state. From top to bottom, the first row presents the cell state dynamics of an unregulated LSTM; the second row illustrates that but under a consistently large forget gate; and the third row is for the scenario where we have a consistently small forget gate.

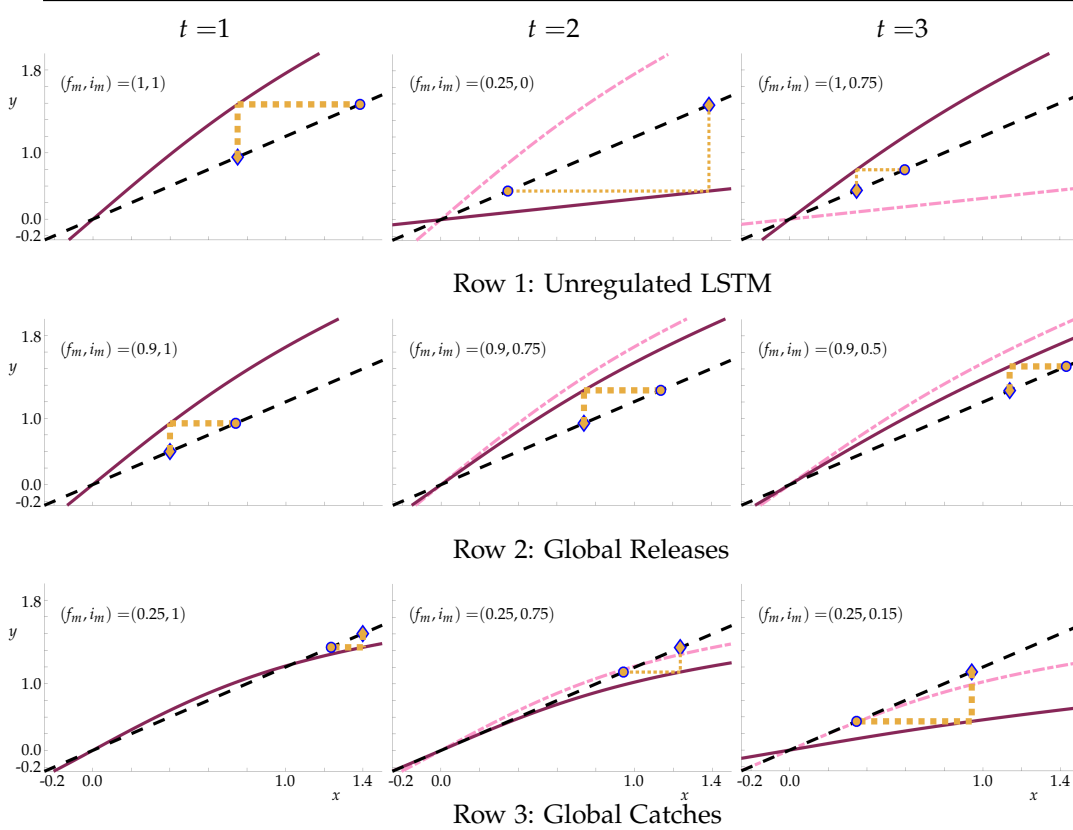


Figure 4.2: The Cell State Propagation of Controlled and Uncontrolled Forget Gates

4.1.3.4 The Near-Random Dynamics of a Vanilla LSTM Cell State

Gated units of LSTMs are not hardwired with deterministic propagation. As a result, they evolve in a near-chaotic fashion. The first row of Figure 4.2 simulates such behaviours of c_t under an unregulated f_t .

The three subfigures of the row are the consecutive instances of $t = 1, 2, 3$. We set parameters $f_{m(1...3)} = [1, 0.25, 1]$ and $i_{m(1...3)} = [1, 0, 0.75]$ they are carefully chosen to exaggerate impacts brought to the vanilla LSTM cell state by the near-chaotic evolution of gates. These settings lead to sudden catches and abrupt releases. The disordered dynamics make it difficult to predict growth and shrinkage in cell state magnitudes. It is also unclear how such drastic behaviours can assist us to understand the learned features of a trained LSTM.

4.1.3.5 Interpretable Cell State Dynamics under Controlled Forget Gates

Controllable forget gate values yield predictable cell propagation. The second row has $f_{mt} = 0.9$ with $i_{m(1...3)} = [1, 0.75, 0.5]$; and the third row has $f_{mt} = 0.25$ and with $i_{m(1...3)} = [1, 0.75, 0.15]$. Consistent growth and consistent shrinkage are observed in rows 2 and 3 respectively.

The magnitudes evolve monotonically because the controlled forget gates expose cell states to homogeneous dynamic regimes. Large constant f_{mt} s expose the cell state

to consistent (*global*) releases; whereas small constant f_{mt} s expose c_{mt} to *global* catches. Dynamic reversals may still occur in \mathbf{c}_t through \mathbf{c}_q , but only when \mathbf{i}_t and \mathbf{a}_t are simultaneously extreme. In addition, the dynamic reversals are merely temporary (*local*) due to the boundedness of \mathbf{c}_q .

4.1.4 DecayNets

From Figure 4.2, it is evident that forget gates deprived of random evolution yield predictable and stable cell state propagation. Thus based on the mathematical insights in Section 4.1.3, we present *DecayNets* as an extension of LSTM to calibrate cell state dynamics.

The forget gate of DecayNet continuously and monotonically decreases. It entails the forget-polar input:

$$\begin{aligned} \mathbf{p}_t &= \mathbf{p}_{t-1} - \frac{\pi}{2D} \frac{1}{6} (\mathbf{W}_F \mathbf{x}_t + \mathbf{U}_F \mathbf{h}_{t-1} + 3) \\ \text{with } \mathbf{p}_0 &= \frac{\pi}{2} \cdot \mathbf{1}, \text{ and} \\ \text{forget gate elements: } \mathbf{f}_{t_n} &= \Phi(\mathbf{p}_{t_n}) = \begin{cases} 1 & \text{for } \mathbf{p}_{t_n} > \frac{\pi}{2}, \\ 0 & \text{for } \mathbf{p}_{t_n} < 0, \\ \sin(\mathbf{p}_{t_n}) & \text{otherwise} \end{cases} \\ &\text{where } n = 1, \dots, N. \end{aligned}$$

We use constant N to denote the network dimension and use D for the total length of the input sequence. DecayNet inherits all remaining LSTMs variables and their corresponding dimensionalities. Note, DecayNet does not introduce any additional learnable parameters.

4.1.4.1 From 1 to 0: A Monotonic Decay

The DecayNet design initialises neural values of the forget gate as a vector of ones

$$\mathbf{f}_{0_n} = \sin(\mathbf{p}_{0_n}) = \sin\left(\frac{\pi}{2}\right) = 1$$

and decreases their values over time. By the end of a fixed sequence of D inputs, we remove a maximised cumulative amount of 1 and a minimised cumulative amount of 0 to constrain the final value $\mathbf{f}_{t_D} \in [0, 1]$. This is made possible by a combination of 2 factors. First, we select the sinusoidal function as the “activation” function of the DecayNet forget gate; and second, by engineering the pre-activated neurons of the DecayNet forget gate to mimic radians with values within the bounds of $[0, \frac{\pi}{2}]$. More on the latter is elaborated below.

DecayNet values all instantaneous decays with equivalent importance. The elements of the forget-polar input are bounded in $[0, \frac{\pi}{2}]$. And for each instance, a maximum amount of $\frac{\pi}{2D}$ is lost from \mathbf{p}_{t_n} 's initialised $\frac{\pi}{2}$. In other words, over an input-length of D units of time, an accumulated decay of $\in [0, \frac{\pi}{2}]$ will be taken away from \mathbf{p}_{0_n} . This hence induces the final values between $\mathbf{p}_{D_n} = \frac{\pi}{2}$ (equivalently $\mathbf{f}_{D_n} = 1$) and $\mathbf{p}_{D_n} = 0$ (equivalently $\mathbf{f}_{D_n} = 0$).

The amount of appropriate decay is learnt through \mathbf{W}_F and \mathbf{U}_F ; and we designed the forget-polar input $\mathbf{p}_{t-1} - \frac{\pi}{2D} \frac{1}{6} (\mathbf{W}_F \mathbf{x}_t + \mathbf{U}_F \mathbf{h}_{t-1} + 3)$ with offset “+3” based on experiments. Empirically, we found that when \mathbf{W}_F and \mathbf{U}_F had their entries initialised with uniform sampling within the interval of ± 1 , $\mathbf{W}_F \mathbf{x}_t + \mathbf{U}_F \mathbf{h}_{t-1}$ was nearly always

bounded within the range of $[-3, 3]$. Hence, “+3” was added as a scaling feature to merely shift the range of the distribution. In addition, “+3” is simply a constant term and thus it does not alter gradient computation in the backward pass.

4.1.4.2 Using Decay to Increase Interpretability

With the redesigned forget gate, DecayNets can learn to dichotomise cell state values conditioned on the input shown. By the end of observing a fixed-length sequence, DecayNet will separate a subset of the cell state values with large magnitudes to another subset of cell state values with small magnitudes. The former subset corresponds to those dimensions of the forget gate that maintain large magnitudes over time; whereas the latter subset corresponds to those that diminish as time unfold.

The monotonic decrease in the forget gate also serves as an interpretation to LSTM’s inferential mechanism. Since the cell $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$ is composed of memory \mathbf{c}_{t-1} and internal update \mathbf{a}_t , the irreversible decay clarifies which dimensions of the cell state gradually lose their recurrent nature over time and effectively become feedforward neurons.

4.1.5 Section-wise Additional Related Work

Similar Dynamical Analyses/Behaviours

[Tallec and Ollivier \[2018\]](#) offered another theoretical study which also treated the incremental updates of LSTMs as dynamical systems. They analysed the effects of multiple scales and showed that recurrent networks that incorporates invariance to time transformation naturally results in a gating mechanism used in LSTMs.

DecayNet also automatically adopts a practice advocated in [Gers et al. \[1999\]](#). [Gers et al. \[1999\]](#) found that LSTM performances could be elevated simply by initialising LSTM forget gate biases to 1s as a way to start the network inferential process with large forget values to “remember more by default”. Similarly, the forget gates of DecayNets are designed to start with the large values of a vector of 1s.

Decay as an Implicit Regularisation

The monotonic decrease of the DecayNet forget gate can be seen as a *hard* regularisation. The network learns to *irreversibly* suppress the recurrent nature of a selection of cell state values conditioned on the input shown. This differs to *Dropout* [[Srivastava et al., 2014](#)], *DropConnect* [[Wan et al., 2013](#)], and *Zoneout* [[Krueger et al., 2017](#)].

Of the three mentioned methods, Dropout randomly nullifies activated neural values. Similarly, DropConnect randomly nullifies some randomly selected synapses. Zoneout behaves like Dropout but freezes instead of nullifies its selected nodes. A commonality among these methods is that their impact on the network are all applied for a single time step and that the nullified and frozen components are recoverable (unless randomly re-selected) for the next time step. Since the modified effects are *temporary*, these three methods are not as harsh as our deterministically decaying forget gate; and hence we regard them as *soft* regularisation methods.

In a previous paper, [Burger and Neubauer \[2003\]](#) also used regularisation to boost neural network performances. However unlike Dropout, DropConnect, and Zoneout, they took a more mathematical approach and invested in inverse problem and integrated Tikhonov regularisation to enhance neural network stability.

4.1.6 Experiments and Results

This subsection will discuss our evaluations on DecayNets and LSTMs for *Permuted Sequential MNIST* [[Le et al., 2015](#)] and L2L [[Andrychowicz et al., 2016](#)] to verify the effectiveness of a monotonically decreasing forget gate.

4.1.6.1 Permuted Sequential MNIST

Table 4.1: Results on Perm-SeqMNIST

Model	Accuracy (%)	Source
RBN DecayNet	Best: 97.8 ; Mean: 97.4	Ours
Temporal Convolutional Network	97.2	Bai et al. [2018]
RBN Zoneout LSTM	95.9	Krueger et al. [2017]
3D tLSTM+CN	95.7	He et al. [2017]
Dilated GRU and Dilated CNN	94.6 and 96.7	Chang et al. [2017]
iRNN	82.0	Le et al. [2015]

Datasets and Experimental Setups

Our first experiment followed the experimental setup of [Le et al. \[2015\]](#) using the MNIST [[LeCun et al., 1998](#)] data. Refer to Appendix A for details on MNIST; but in brief, this dataset consisted handwritten number digits between [0-9] on 28×28 -pixel formats. The permuted sequential MNIST (Perm-SeqMNIST) presented each pixel, one at a time, as a $784 (= 28 \times 28)$ -unit pixel-by-pixel sequence to RNNs. For this image recognition task, the RNNs outputs were sent to a softmax layer (see Equation (1)) and we measured the accuracy.

Our DecayNets and LSTMs had a single layer of 100 hidden dimensions and were trained on the RMSProp optimiser [[Tieleman and Hinton, 2012](#)] over 150 epoches. The learning rate was 0.001, with 0.9 momentum and no weight decay. Gradient clipping was applied at 1 to avoid gradient exploding [[Pascanu et al., 2012](#)].

Results

The best prediction accuracy for our vanilla RNNs was 92.1% for DecayNet and 89.5% for LSTM (which matched the baseline model of 89.8% in [Krueger et al. \[2017\]](#)). When we integrated *Recurrent Batch Normalisation* (RBN) [[Cooijmans et al., 2017](#)] to DecayNet, we found that both the best performing RBN DecayNet of 97.8% and the average performance of 97.4% yielded better results than the state of the art performance reported in [Bai et al. \[2018\]](#)² at 97.2%. A list of competitive results of methods

²Also note that [Bai et al. \[2018\]](#)'s network employed CNNs instead of RNNs, and that CNNs are the most natural ML module of choice for modelling images.

in recent literature can be found in Table 4.1.

Remarks and Analyses

There were two remarks regarding this experiment. First, the superior performances exhibited by vanilla DecayNets (92.1%) over LSTMs (89.5%) could be contributed to the implicit regularisation of the monotonically decaying forget gate. Recall that in Section 4.1.4.2, we stated that the decay mechanism functioned as a form of regularisation because that it removed the recursive nature of a section of cell state neurons. This hence allowed the network to filter out unnecessary reiterative noises at the activation level in the cell state values.

Second, our experiments showed that DecayNet achieved impressive results (97.8%) when combined with RBN. As mentioned in Section 4.1.4.1, we found that the forget polar input was consistently bounded within the range of ± 3 and hence it was likely that neural values of other variables within the DecayNet also required normalisation at the activation level. Readers can find the inner workings of batch normalisation [Ioffe and Szegedy, 2015b] in Section 2.1.2.5 and defined in ICC-08.

Overall, the positive results of the first experiment indicated that we had correctly identified the cell state as the most important component within the LSTM design. More details regarding this experiment can be found in our paper Kuo et al. [2018]. Therein that paper, we also evaluated a similar experiment on the permuted sequence task but on the more difficult dataset of Fashion-MNIST [Xiao et al., 2017].

4.1.6.2 Learning to Optimise

In addition to the permuted MNIST test, we conducted the following L2L experiment to test RNNs of their robustness in making consecutive decisions with the potentiality of regurgitating their past erroneous outputs.

Datasets and Experimental Setups

Here we considered L2L and employed DecayNets and LSTMs as meta-learners to update two-layer MLP base learners on MNIST classification. The MLPs had 784 input dimensions and 20 hidden dimensions, and processed images in batch sizes of 32. The meta-learners were themselves trained with Adam with learning rate 0.01.

Recall that the original L2L [Andrychowicz et al., 2016] used LSTMs to replace gradient descent $\theta_t = \theta_{t-1} - \alpha \nabla_{\theta} \mathcal{L}$ to update base learner parameters θ_t . See Section 3.2.5 for an in-depth description. We modified the L2L experiment for DecayNets and only used DecayNets to infer new learning rates

$$\alpha_t, \mathbf{h}_t, \mathbf{c}_t = \text{DecayNet}(\nabla_{\theta_t} \mathcal{L}, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}). \quad (4.1)$$

Our DecayNet-optimisers updated their hidden variables for multiple steps before creating the final learning rates. For each time step, we re-initialised the forget-polar

input as $\frac{\pi}{2} \cdot \tilde{\mathbf{1}}$ and operated 2 internal iterations. Thus for each time step, we have

$$\mathbf{h}_{t-\frac{1}{2}}, \mathbf{c}_{t-\frac{1}{2}} = \text{DecayNet}(\nabla_{\theta_t} \mathcal{L}, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}) \text{ and} \quad (4.2)$$

$$\alpha_t, \mathbf{h}_t, \mathbf{c}_t = \text{DecayNet}(\nabla_{\theta_t} \mathcal{L}, \mathbf{h}_{t-\frac{1}{2}}, \mathbf{c}_{t-\frac{1}{2}}). \quad (4.3)$$

If required, practitioners could prolong the internal iteration to, for example, 20 steps.

Results

Figure 4.3 presents a comparison of optimisers on their progress of updating MLPs. The figure is a comparison of the loss curves – the cross entropy loss (see Equation (2)) of the baseline learners against time. The loss curve of a good optimiser is expected to have three attributes. First, a sharp decrease for its loss curve. The sharpness of slope the curve depicts how efficient an optimiser can update a baseline learner. Second, a low converged loss value. The converged value reflects the ability of the baseline learner to generalise on the observed data. Third, a low variability during the entire course of training updates. This implies how reliable the optimisers are at guiding the baseline learners towards minima on the loss landscape, and how easily they are at assisting the baseline learners in escaping sub-optimal minima.

We used colours green, grey, blue, yellow, and purple denote SGD with learning rate 0.05, Adam with learning rate 0.002, RMSProp with learning rate 0.001, a trained LSTM-optimiser, and a trained DecayNet-optimiser, respectively. The loss curve performances of all optimisers were aggregated over 20 seeds. The solid lines represented the mean loss performance at each iteration; while the coloured bands surrounding the solid lines represented the variability of the optimiser performances. Our results highly resembled those reported in [Andrychowicz et al. \[2016\]](#)³; and as we expected, the RNN-optimisers outperformed their hand-crafted rivals by a significant margin hence showing that they were able to discover better optimisation strategies than hand-crafted rules. In addition, our DecayNet-optimisers converged MLPs faster than the vanilla LSTM-optimisers.

Remarks and Analyses

It was not particularly surprising to see that DecayNets converged faster than LSTM RNNs for L2L. Figure 4.3 showed that the LSTM cell could benefit from the stabilisation in recurrent neural processing from the implicit regularisation of the decay mechanism. If there were noises in the neural optimiser outputs, the learnt algorithm would then parameterise the network incorrectly. Consider Figure 2.3(a), an incorrect parametric update would be equivalent to the action of pushing the marble (network weights) towards hills (high loss region) on the loss landscape. This meant

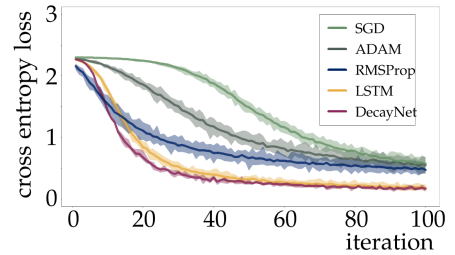


Figure 4.3: L2L vs Gradient Descent

³Refer to the middle subplot of Figure 4 on page 5 of [Andrychowicz et al. \[2016\]](#).

that it would require additional iterations to right the incorrect actions taken to reset the trajectory of parametric updates to the right path to guide the marble to fall into an appropriate loss minimum.

4.1.7 Section Conclusion

This section proposed DecayNets, LSTMs with monotonically decreasing forget gates. Our design was based on mathematical understandings on the cell state mechanics. We studied the cell state propagation as a system of difference equations and found an alternating catch-and-release regime which controlled its shrinkage and growth. The alteration was predominately dictated by the forget gate; and that interpretable cell state propagation could emerge by monotonically decreasing forget gate values. With RBN, DecayNets achieved state of the art performances on Perm-SeqMNIST. In addition, the decay mechanism also served as a promising addition to the meta-learning task of L2L for AutoML.

Our positive findings demonstrated that neural network performances could be improved at the level of neural activations and highlighted the importance in the need for understanding behavioural changes in individual neurons. However, we propose a counter-argument on this perspective overleaf.

4.1.8 An Evaluation of this Study

Before we move to our next RNN study, we would like to highlight some advantages and disadvantages of our approach taken in the DecayNet study.

The Advantages

The main advantage of this study is the clarification on the two dynamically different regime in the update of the cell state. By identifying these behaviours, we have paved the way for many more interesting future work. An example is given below.

Proposal: On Redesigning the Activation Function

As we mentioned in Section 2.3 and simulated in Figures 4.1 and 4.2, we found that the attractiveness of LSTM's equilibria is highly dependent on the topological outline of the activation function. Recall from Section 4.1.3.3, the only reason why the catch regime was more dominant than the release regime was because of the boundedness of $\mathbf{c}_q = \mathbf{i}_t \odot \mathbf{a}_t$ from \tanh . Perhaps cell state replenishment is actually more important than cell state forgetting; and in order to make the former the more dominant dynamic regime, one could engineer an appropriate topological outline by selecting a different activation function for \mathbf{a}_t .

The Disadvantages

LSTMs are still opaque. Thus far, we have only considered an atomic approach towards understanding the recurrent behaviour of the LSTM cell. Though this is definitely important, but by isolating the cell and making the assumptions that gated units are only secondary in importance (see Section 4.1.3.1), we have ignored potential data driven causative relationships between the different components of the LSTM RNN. For instance, this study cannot address the question “*How does the forget gate behave as according to the input gate?*” We believe that our study (with using dynamical systems) is still methodologically correct; however, we also acknowledge that we would need to steer our work in a more Bayesian inferential approach, to account for the entwined relationship of LSTM components in the forward path.

In order to achieve this, we have to approach to understand LSTMs from a much higher level of conceptualisation. Instead of investigating the individual behaviour of any one neuron, we should target to understand LSTM properties from a network level. Metaphorically, instead of examining just one specie of vegetation in a habitat, we should focus on understanding the properties of the entirety of the habitat as a unit of a well-defined and self-enclosed ecosystem.

Our Next Study

To fill the holes identified in **The Disadvantages**, our next study will be much more generalised. We will compare a list of successful RNN and highlight the *minimal ingredient* for making these networks successful. After finding the core recipe, our attention will be directed towards deleting unnecessary components in LSTMs. The aim is to create a general-use and simplified LSTM with equivalent performances to its vanilla counterpart while having a drastically lower amount of parameters.

*Understanding in what terms an RNN should really be formulated
is one of the basic needs and goals in machine learning.
Knowing some of the phenomena that can occur for strong coupling —
if one can know them without already knowing the good formulation!⁴
Many RNN designs look similar to LSTM design albeit
crafted specifically for a designated task.
Hence, we question whether it is possible to unify all designs.
The following section aims to find the underlying properties
that conform many gated RNNs as one mathematical class.*

*By doing so,
we hope to identify the core building blocks for achieving recurrent neural processing.*

Section Content:

- Section 4.2.1: Section Introduction,
- Section 4.2.2: Connecting GRUs to Hopfield Networks,
- Section 4.2.3: An Input Residual Connection,
- Section 4.2.4: IRC-/IHC-GRNNs,
- Section 4.2.5: Experiments and Results,
- Section 4.2.6: Section-wise Additional Related Work,
- Section 4.2.7: Section Conclusion:
IRC/IHC on Simplifying GRNNs,

⁴A playful twist on [Witten \[1995\]](#) opening remark for the S-duality between weak and strong coupling for the unification on string theory.

4.2 An Input Residual Connection for Simplifying Gated Recurrent Neural Networks

The research in Section 4.1 took an atomic approach to understand LSTM behaviours from the level of a single neuron. Though it helped us to identify the LSTM cell as the most important component of the network, that approach limited our ability in describing coordinated behaviours between any pair of LSTM gated units.

To complement our previous study, this section analyses RNNs from a macroscopic level to find generic behaviours among gated units transferable across different recurrent neural designs. The work herein is based on our paper [Kuo et al., 2020b] which describes LSTM, GRU, and many more gated RNNs as one class of models conforming to a common set of mathematical properties.

4.2.1 Section Introduction

This section aims to connect *Gated Recurrent Neural Networks* (GRNNs) to well-studied dynamical systems to understand the inner-workings of the internal components of LSTMs and GRUs [Cho et al., 2014]. The idea is to increase interpretability by returning to first principles and establish mathematical meanings for every GRNN variables. Specifically, our theoretical groundwork centres around linking the GRU RNN to Hopfield networks [Hopfield, 1982].

This mathematical framework will allow us to identify network redundancies of which we simplified with an *Input Residual Connection* (IRC) and an *Input Highway Connection* (IHC). Both IRC and IHC are generally applicable to most existing GRNN designs; besides the canonical LSTMs and GRUs, we also applied these methods to simplify *FastGRNNs* [Kusupati et al., 2018], *Simple Recurrent Units* (SRUs) [Lei et al., 2018a], and *Strongly-Typed Recurrent Neural Networks* (T-RNNs) [Balduzzi and Ghifary, 2016]. We tested all vanilla GRNNs against their IRC counterparts on language modelling [Merity et al., 2018]. Despite parameter reductions, all IRC-GRNNs showed either comparative or superior generalisation than their baseline models. In addition to NLP, we took the most widely applied GRNN of LSTM and tested it against IHC-LSTMs on the CV task of image generation [Gregor et al., 2015], and on L2L [Andrychowicz et al., 2016] to update another learner network. All IHC-LSTMs exhibited comparative performances to LSTMs; and notably, IHC-LSTMs also removed 85.4% of LSTM parameters on the image generation task.

4.2.2 Connecting GRUs to Hopfield Networks

This subsection will establish connections between GRUs and Hopfield Networks. As a reminder from Equation (3.9), GRU is formally written as

$$\begin{aligned} \mathbf{i}_t, \mathbf{r}_t &= \sigma(\mathbf{W}_{\{I,R\}} \mathbf{x}_t + \mathbf{U}_{\{I,R\}} \mathbf{h}_{t-1} + \mathbf{b}_{\{I,R\}}), \\ \mathbf{a}_t &= \tanh(\mathbf{W}_A \mathbf{x}_t + \mathbf{U}_A (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A), \text{ and} \\ \mathbf{h}_t &= (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t. \end{aligned}$$

Similar to our previous work in DecayNet (see Section 4.1), we will continue to model RNNs as dynamical systems. However this time, the update of GRU memories will be modelled with the continuous dynamical systems of *differential equations* [Glendinning, 1994]. See Section 2.3.1 for the differences between differential equations and difference equations. Through this analysis, we will describe some existing desirable properties that are already possessed by the vanilla GRU design.

4.2.2.1 Recursive Feedbacks as Systems of Differential Equations

The GRU hidden state \mathbf{h}_t is subjected to continual updates. The current value of \mathbf{h}_t is updated from \mathbf{h}_{t-1} with the additive package of \mathbf{a}_t . For this reason, we model \mathbf{h}_t as a system of differential equations with the autonomous function

$$\mathbf{h}'_t = \mathcal{H}(\mathbf{h}_{t-1}, \mathbf{a}_t). \quad (4.4)$$

Ideally, we would like function \mathcal{H} to provide the GRU memory with the two desirable properties of *robustness* and *plasticity*. We consider robustness as the ability for network memory to generalise data under the presence of unknown noises; whereas plasticity allows the network to shape its memory as according to the input data.

From the vanilla GRU formulation provided in the previous page, it is hard to ensure whether such desirable properties exist. However, these properties are known to be possessed by the well studied theoretical physics model of the dynamic spin-glass Hopfield network [Hopfield, 1982; Sompolinsky et al., 1988; Sussillo and Abbott, 2009]. We know that dynamic Hopfield network variables are both robust and plastic because the network possesses both *attractive equilibria* [Sompolinsky et al., 1988] and *phase transitions* [Dyre, 2006]. On one hand, attractive equilibria can provide stability; more descriptions were previously detailed in Section 2.3.4. On the other hand, phase transitions are parametric transitional boundaries for internal degrees of freedom to successively fall out of equilibrium, and thus they provide network variables with the plasticity to transmute among various configurations. If we could successfully establish a connection between GRUs and Hopfield networks, it would then be easier to analysis whether GRUs possess any of these desirable qualities.

4.2.2.2 Deriving the GRU from the Dynamic Hopfield Network

In order to establish a connection between the GRU memory and Hopfield networks, we formulate \mathcal{H} of Equation (4.4) as the dynamic Hopfield model [Sussillo and Abbott, 2009] below

$$\mathbf{h}'_t = \mathcal{H}(\mathbf{h}_{t-1}, \mathbf{a}_t) = -\mathbf{h}_{t-1} + \mathbf{a}_t. \quad (4.5)$$

Similar to that in Equation (2.9), Equation (4.5) has the equivalent discrete time difference equation of

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \zeta(-\mathbf{h}_{t-1} + \mathbf{a}_t), \quad (4.6)$$

with a small positive constant ζ serving as the discretised step size. However, it is unnecessary for the network to be limited to a single time scale. That is similar to the practice of scaling in Echo State Networks [Lukoševičius, 2012], we can enable a set of different timescales by rewriting Equation (4.6) as

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \zeta \odot (-\mathbf{h}_{t-1} + \mathbf{a}_t) \quad (4.7)$$

where ζ now becomes a vector of (trainable) positive values. From here, we can rearrange Equation (4.7) as below

$$\mathbf{h}_t = (1 - \zeta) \odot \mathbf{h}_{t-1} + \zeta \odot \mathbf{a}_t. \quad (4.8)$$

There is high similarity between Equation (4.8) and the GRU update of the hidden state in Equation (3.9). That is, the GRU update of \mathbf{h}_t conforms to the differential equation of (4.8) with function \mathcal{H} as the dynamic Hopfield network (see Equation (4.5)). Additionally from Equation (4.8), we observe that the input gate \mathbf{i}_t plays a similar role to ζ and serves as the dynamic step sizes for the network. This series of derivation thus increases the individual interpretability of GRU components.

However, our readers should note that Equation (4.8) is not equivalent to the GRU formulation. In fact, we can only derive the GRU from

$$\begin{aligned} \mathbf{h}'_t &= \mathcal{H}(\mathbf{h}_{t-1}, \mathbf{a}_t) = -\mathbf{h}_{t-1} + \mathbf{a}_t \quad \text{to} \\ \mathbf{h}_t &= (1 - \zeta) \odot \mathbf{h}_{t-1} + \zeta \odot \mathbf{a}_t \end{aligned}$$

while *also* setting

$$\zeta = \mathbf{i}_t = \sigma(\mathbf{W}_I \mathbf{x}_t + \mathbf{U}_I \mathbf{h}_{t-1} + \mathbf{b}_I), \quad (4.9)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_R \mathbf{x}_t + \mathbf{U}_R \mathbf{h}_{t-1} + \mathbf{b}_R), \quad \text{and} \quad (4.10)$$

$$\mathbf{a}_t = \tanh(\mathbf{W}_A \mathbf{x}_t + \mathbf{U}_A (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A). \quad (4.11)$$

To elaborate, Equation (4.8) is not a specific update for any RNN. Instead, it represents a *class* of RNN memory updates with Hopfield network-like robustness and plasticity. Furthermore, since the GRU falls under this class of differential systems, GRU memories possess both robustness and plasticity.

4.2.3 An Input Residual Connection

As we discussed in Section 4.2.2, the GRU RNN is simply *one design* that updates its memory in accordance to Equation (4.8). That is, there are infinitely many ways to formulate the update of an RNN memory to achieve the same, or potentially better, level of Hopfield network-like robustness and plasticity. In this subsection, we will look for the lightest possible simplification for the GRU.

Equation (4.8) highlights two types of GRU redundancies. One relates to the network non-linearities, and the second relates to the recursive network feedback.

4.2.3.1 An Analysis on the Non-linearity of \mathbf{a}_t

Equation (4.8) required neither step size ζ nor update \mathbf{a}_t to undergo non-linear transformation. The hyperbolic tangent function of \mathbf{a}_t can be limiting for the network because it bounds $\mathbf{a}_t \in [-1, 1]$. The upper and lower bound restrict \mathbf{a}_t from distinguishing pre-activated neurons of extreme values. Thus, the hyperbolic tangent function in \mathbf{a}_t restricts a diverse source of information to be added towards the hidden state \mathbf{h}_t . For this reason, our first simplification on the GRU is

$$\mathbf{a}_t = \mathbf{W}_A \mathbf{x}_t + \mathbf{U}_A (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A. \quad (4.12)$$

However, we choose to keep the sigmoid activation function in \mathbf{i}_t .

4.2.3.2 An Analysis on the Non-linearity of \mathbf{i}_t

We keep the sigmoid activation function in \mathbf{i}_t on the same grounds as having small learning rates for gradient descents. Recall from Equation (1.4) that gradient descent updates neural network parameters with

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L};$$

and it shares several similarities to Equation (4.7) of the formulation of

$$\mathbf{h}_t = \mathbf{h}_{t-1} + \zeta \odot (-\mathbf{h}_{t-1} + \mathbf{a}_t).$$

Both Equation (4.7) and gradient descent comply with the properties of a vector field (see Section 2.3.5.1). As we previously discussed in the preliminary chapter and illustrated in Figure 2.4, the vector field of gradient descent can be described with the loss landscape. Gradient descent first specifies the direction to update the network parameters via $-\nabla_{\theta} \mathcal{L}$, and then takes a step of the size of learning rate α towards this specific direction on the loss landscape. Larger learning rates can cover more ground in each step but risks overshooting a minimum. Similarly, the differential system formulation of Equation (4.7) means that it internally possesses attractive equilibria which risk overshooting with a large step size ζ .

As we detailed in Section 2.3.4, an equilibrium is said to be attractive if nearby solutions were to be updated towards it. For this reason, a small step size of ζ in Equation (4.7) is required in order for the network memory to traverse within the vicinity of an attractive equilibrium in the solution space to achieve robustness. Since derivation (4.8) have already shown that the input gate \mathbf{i}_t acts as the dynamic step size for the GRU, we decide to keep the sigmoid function for the input gate to restrict \mathbf{i}_t as small positive values within the range of $\in [0, 1]$.

4.2.3.3 An Analysis on the Recurrent Feedback

Through derivation (4.8), it is not immediately clear how the recursive pre-activated neural values of $\mathbf{U}\mathbf{h}_{t-1}$ should be designed in the GRU. Thus, it is necessary to

rethink whether \mathbf{a}_t and \mathbf{i}_t should take the forms of

$$\begin{cases} \mathbf{i}_t &= \mathbf{i}_t(\mathbf{x}_t, \mathbf{h}_{t-1}) \\ \mathbf{a}_t &= \mathbf{a}_t(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{r}_t) \end{cases} \quad (4.13)$$

or as

$$\begin{cases} \mathbf{i}_t &= \mathbf{i}_t(\mathbf{x}_t) \\ \mathbf{a}_t &= \mathbf{a}_t(\mathbf{x}_t) \end{cases}. \quad (4.14)$$

Since the formulations of Equation (4.13) and Equation (4.14) both suffice for derivation (4.8), we hypothesise that $\mathbf{U}\mathbf{h}_{t-1}$ is not directly required for formulating the GRU components of \mathbf{i}_t , \mathbf{r}_t , and \mathbf{a}_t . Hence, we simplify the gated units as

$$\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{I,R\}}\mathbf{x}_t) \quad (4.15)$$

and further simplify the internal input from Equation (4.12) to

$$\mathbf{a}_t = \mathbf{W}_A \mathbf{x}_t. \quad (4.16)$$

Note that we also removed the biases from all components; this was because like the recursive connections, derivation (4.8) did not explicitly require for them. However, after removing the recursive connections of $\mathbf{U}\mathbf{h}_{t-1}$, this simplified version of the GRU is no longer a real RNN. Below, we will describe where and how we restore the relinquished recurrent feedback.

4.2.3.4 An Input Residual Connection

Thus far, we have achieved two things in this section. First, we established connections between the dynamic Hopfield network and the GRU RNN; and second, we used this connection to simplify the GRU components of \mathbf{i}_t , \mathbf{r}_t , and \mathbf{a}_t . The removal of $\mathbf{U}\mathbf{h}_{t-1}$ and biases \mathbf{b} have resulted in a large amount of parametric reduction. However, the current simplification entirely renounced the recurrent feedback of \mathbf{h}_t and hence the network is no longer a real RNN. Here, we will introduce our novel input residual connection (IRC) design to restore \mathbf{h}_t while maintaining a light architecture.

Since Equation (4.14) suffices for derivation (4.8), we hypothesise that GRUs can function properly with $\mathbf{W}\mathbf{x}_t$ as the sole pre-activated neural values. For this reason, we see the recursive $\mathbf{U}\mathbf{h}_{t-1}$ as a form of regularisation to $\mathbf{W}\mathbf{x}_t$ that eases the training of the feedforward input in the activation functions. Under this interpretation, we propose to regularise $\mathbf{W}\mathbf{x}_t$ before they enter the activation functions. We replace \mathbf{x}_t in the GRU gated units with the regulated inputs of \mathbf{v}_t , $\mathbf{x}_t \leftarrow \mathbf{v}_t$, such that

$$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\mathbf{U}_V \mathbf{h}_{t-1}). \quad (4.17)$$

Variable \mathbf{v}_t connects the recursive $\mathbf{U}_V \mathbf{h}_{t-1}$ to the input \mathbf{x}_t with a residual connection [He et al., 2016a] to provide \mathbf{x}_t with a learnt discretised update based on the

network memory. Since the residual connection is applied to the input, we call this technique an IRC. The new variables have dimensionality $\phi \in \mathbb{R}^N$ and $\mathbf{U}_V \in \mathbb{R}^{N \times M}$ where N and M are the input dimension and hidden dimension respectively. In comparison, the recursive connections in the vanilla GRU design of Equation (3.9) have dimensionality $\mathbf{U} \in \mathbb{R}^{M \times M}$ and thus the original recursive connections are typically much larger than the size of that of \mathbf{U}_V . Vector ϕ is an idea of which we borrowed from FastGRNN [Kusupati et al., 2018] for controlling the network condition number (see Section 3.1.7). This vector contains trainable weights where $0 \leq \phi_n \leq 1$ for $n = 1, \dots, N$ and are parameterised by the sigmoid function. Vector ϕ is hence employed to limit the extent that the recurrent feedback modifies the features of \mathbf{x}_t .

To summarise, the IRC-GRU consists

i) the IRC regulated input of Equation (4.17)

$$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\mathbf{U}_V \mathbf{h}_{t-1}),$$

ii) which substitutes \mathbf{x}_t in the simplified gates of Equation (4.15)

$$\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{I,R\}} \mathbf{v}_t),$$

iii) the simplified internal input of Equation (4.16)

$$\mathbf{a}_t = \mathbf{W}_A \mathbf{x}_t, \text{ and}$$

iv) the hidden state update

$$\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot (\mathbf{r}_t \odot \mathbf{a}_t). \quad (4.19)$$

Readers should note that, unlike the gated units of Equation (4.18), the IRC-GRU \mathbf{a}_t does not substitute \mathbf{x}_t with \mathbf{v}_t . This is because that similar to the dynamic Hopfield network of Equation (4.5), we prefer memory \mathbf{h}_t to receive a new source of information \mathbf{a}_t that is dissimilar to the existing memory. See Figure 4.4 for a figurative comparison of the architectural differences between GRU and IRC-GRU.

4.2.4 IRC-/IHC-GRNNs

The previous subsection addressed IRC as a technique for simplifying GRU. In this subsection, we will show that IRC is generally applicable across multiple GRNN designs including, LSTMs, SRUs [Lei et al., 2018a], T-LSTMs [Balduzzi and Ghifary, 2016], and FastGRNNs [Kusupati et al., 2018]. A list of contrasting designs between the vanilla network and their IRC counterparts are provided in Table 4.2 overleaf.

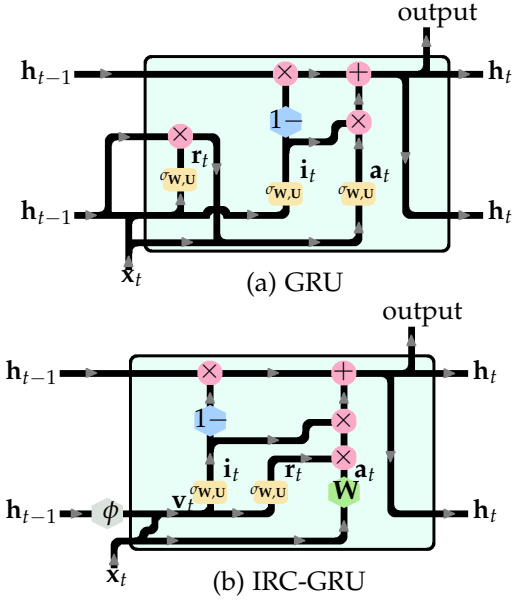


Figure 4.4: GRU vs IRC-GRU
(4.18)

Table 4.2: A comparison between baseline GRNNs against IRC-GRNNs

Traditional design	Simplified design
GRU [Cho et al., 2014]	IRC-GRU (Ours)
$\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{I,R\}}\mathbf{x}_t + \mathbf{U}_{\{I,R\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{I,R\}}),$ $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_A),$ and $\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t.$	$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\mathbf{U}_V \mathbf{h}_{t-1}),$ $\mathbf{i}_t, \mathbf{r}_t = \sigma(\mathbf{W}_{\{I,R\}}\mathbf{v}_t),$ $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t,$ and $\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot (\mathbf{r}_t \odot \mathbf{a}_t).$
LSTM [Hochreiter and Schmidhuber, 1997b]	IRC-LSTM (Ours)
$\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,I,O\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{F,I,O\}}),$ $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{h}_{t-1} + \mathbf{b}_A),$ $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t,$ and $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t).$	$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\mathbf{U}_V \mathbf{h}_{t-1}),$ (4.20) $\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{v}_t),$ (4.21) $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t,$ (4.22) $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t,$ and (4.23) $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t).$ (4.24)
SRU [Lei et al., 2018a]	IRC-SRU (Ours)
$\mathbf{f}_t, \mathbf{p}_t = \sigma(\mathbf{W}_{\{F,P\}}\mathbf{x}_t + \omega_{\{F,P\}} \odot \mathbf{c}_{t-1} + \mathbf{b}_{\{F,P\}}),$ (4.25) $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t$ (4.26) $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + (1 - \mathbf{f}_t) \odot \mathbf{a}_t,$ and (4.27) $\mathbf{h}_t = \mathbf{p}_t \odot \mathbf{c}_t + (1 - \mathbf{p}_t) \odot \mathbf{x}_t.$ (4.28)	$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\omega_V \odot \mathbf{h}_{t-1}),$ (4.29) $\mathbf{f}_t, \mathbf{p}_t = \sigma(\mathbf{W}_{\{F,P\}}\mathbf{v}_t),$ (4.30) $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t,$ (4.31) $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + (1 - \mathbf{f}_t) \odot \mathbf{a}_t,$ and (4.32) $\mathbf{h}_t = \mathbf{p}_t \odot \mathbf{c}_t + (1 - \mathbf{p}_t) \odot \mathbf{x}_t.$ (4.33)
T-LSTM [Balduzzi and Ghifary, 2016]	IRC-T-LSTM (Ours)
$\mathbf{f}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,O\}}\mathbf{x}_{t-1} + \mathbf{b}_{\{F,O\}}),$ (4.34) $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{x}_{t-1} + \mathbf{b}_A,$ (4.35) $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + (1 - \mathbf{f}_t) \odot \mathbf{a}_t,$ and (4.36) $\mathbf{h}_t = \mathbf{o}_t \odot \mathbf{c}_t.$ (4.37)	$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\omega_V \odot \mathbf{h}_{t-1}),$ (4.38) $\mathbf{f}_t, \mathbf{o}_t = \sigma(\mathbf{W}_{\{F,O\}}\mathbf{v}_t),$ (4.39) $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t,$ (4.40) $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + (1 - \mathbf{f}_t) \odot \mathbf{a}_t,$ and (4.41) $\mathbf{h}_t = \mathbf{o}_t \odot \mathbf{c}_t.$ (4.42)
FastGRNN [Kusupati et al., 2018]	IRC-FastGRNN (Ours)
$\mathbf{f}_t = \sigma(\mathbf{W}_F\mathbf{x}_t + \mathbf{U}_F\mathbf{h}_{t-1} + \mathbf{b}_F),$ (4.43) $\mathbf{a}_t = \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{h}_{t-1} + \mathbf{b}_A),$ (4.44) $\mathbf{h}_t = \mathbf{f}_t \odot \mathbf{h}_{t-1} + (\nu(1 - \mathbf{f}_t) + \kappa) \odot \mathbf{a}_t.$ (4.45)	$\mathbf{v}_t = \mathbf{x}_t + \phi \odot (\mathbf{U}_V \mathbf{h}_{t-1}),$ (4.46) $\mathbf{f}_t = \sigma(\mathbf{W}_F\mathbf{v}_t),$ (4.47) $\mathbf{a}_t = \mathbf{W}_A\mathbf{x}_t,$ and (4.48) $\mathbf{h}_t = \mathbf{f}_t \odot \mathbf{h}_{t-1} + \nu(1 - \mathbf{f}_t) \odot \mathbf{a}_t.$ (4.49)

In Table 4.2, the left column provides the vanilla design while the right column shows our IRC simplification. For those designs on the left, we colour redundant components in brown identified via derivation (4.8). As for those simplifications on the right, we highlight the standardised IRC variables in red.

A discussion on the vanilla networks in Table 4.2 was previously given in Section 3.1.7. For that reason, we will refrain ourselves from repeatedly comparing these networks. Alternatively, interested readers could also consult Section IV on page 4 of our work in Kuo et al. [2020b]. Therein that paper, we detailed the meaning of every component of the different architectures.

4.2.4.1 An Input Highway Connection

Here we introduce the more powerful IRC variant of input highway connection (IHC). As the name implies, IHC prepares a regulated input with a highway connection [Srivastava et al., 2015]. The newly introduced variable \mathbf{v}_t of Equation (4.17) is alternatively computed as

$$\mathbf{v}_t = (1 - \Lambda_t) \odot \mathbf{x}_t + \Lambda_t \odot (\mathbf{U}_V \mathbf{h}_{t-1}) \text{ where} \quad (4.50)$$

$$\Lambda_t = \sigma(\mathbf{W}_\Lambda \mathbf{x}_t + \Gamma \mathbf{h}_{t-1} + \mathbf{b}_\Lambda). \quad (4.51)$$

The components have dimensionality $\Lambda_t, \mathbf{b}_\Lambda \in \mathbb{R}^N$, and $\Gamma \in \mathbb{R}^{M \times N}$; hence similar to \mathbf{U}_V in Equation (4.17), the size of matrix Γ is typically much smaller than the ordinary recursive connection of $\mathbf{U} \in \mathbb{R}^{M \times M}$ in GRUs and LSTMs. The parametric vector of ϕ in Equation (4.17) is replaced with gate Λ_t . Gate Λ_t provides a two-fold advantage over ϕ . First, unlike the static control exerted by ϕ , Λ_t limits the extent of modification from $\mathbf{U}_V \mathbf{h}_{t-1}$ to \mathbf{x}_t in a dynamic manner based on the memory of the network. Second, \mathbf{x}_t is modified with $1 - \Lambda_t$; and in principle, this allows the network to learn a balance extrapolation of information based on the recurrent feedback $\mathbf{U}_V \mathbf{h}_{t-1}$ and on input \mathbf{x}_t .

4.2.5 Experiments and Results

This subsection compares all GRNNs included in Table 4.2. The original goal of our study was to show that IRC and IHC can enable GRNNs to achieve comparative performances to their baseline counterparts while removing network redundancies. It is a welcoming surprise however, that our results demonstrated that all IRC-GRNNs achieved superior generalisation performances when compared with their vanilla counterparts. We conjecture that the improvements could be due to the theoretical advantages of the low complexity of simpler models [Krizhevsky, 2009].

We evaluated the networks over three tasks. First, we tested GRNNs against their IRC counterparts on language modelling [Merity et al., 2018]. Then, we tested LSTM against IHC-LSTM on image generation [Gregor et al., 2015], and on L2L to update another learner network [Andrychowicz et al., 2016].

Image generation and L2L were reserved for LSTMs for fairness. GRNNs have been applied on a wide range of tasks and that most of the baseline GRNNs included in Table 4.2 have previously been tested on language modelling. However, such is not the case for image generation and for L2L. For this reason, we reserved image generation and L2L only for the most extensively applied GRNN of LSTM.

The three tasks were carefully chosen to analyse different aspects of GRNNs. The GRNNs were employed as generative language models to test for their abilities to represent and to manipulate high dimensional probability distributions. The GRNNs were employed as decoders in image generation to test for their reconstruction powers. Last, the GRNNs operated for variable lengths on L2L to test their robustness for regularising compounded prediction errors conditioned on unseen sequential context at test time.

Table 4.3: Perplexity on the PTB dataset.

Model	Test PPL	#RNN-Params	Reduction
Lower is better			
3-layers, (ED, HD) = (400, 1150)			
Quoted Score			
LSTM	57.3	24.9M	–
Our Evaluated Score			
LSTM	58.32	24.9M	–
IHC-LSTM	57.44	15.1M	39.2%
1-layers, (ED, HD) = (650, 650)			
Quoted Score			
FastGRNN	116.11	–	–
Our Evaluated Score			
FastGRNN	88.48	1.7M	–
IRC-FastGRNN	85.60	1.3M	25.0%
2-layers, (ED, HD) = (650, 650)			
Quoted Score			
GRU	93.44	5.1M	–
T-LSTM	81.52	5.1M	–
FastGRNN	106.23	–	–
Our Evaluated Score			
FastGRNN	89.30	3.4M	–
IRC-GRU	76.51	3.4M	33.4%
IRC-T-LSTM	79.56	2.5M	50.0%
IRC-FastGRNN	85.60	2.5M	25.0%
3-layers, (ED, HD) = (1150, 1150)			
Quoted Score			
SRU (Ours)	92.42	12.0M	–
Our Evaluated Score			
IRC-SRU (Ours)	86.70	12.0M	0.06%

4.2.5.1 Language Modelling

Dataset and Metrics

Language models (see Equation (3)) employ RNNs to construct probabilistic predictions of the next word given preceding ones. We tested the models on the English corpus dataset of *Penn Treebank* (PTB) [Marcus et al., 1993]; refer to Appendix A for more description on the dataset.

We evaluated the well-being of the models with the *Perplexity Loss* (PPL) [Brown et al., 1992] and the amount of parameters of each model. On one hand, PPL is a standard metric for language modelling and it is defined as the inverse probability of the test set normalised by the number of words (see Equation (5)); it hence checks for the quality of the probabilistic mapping learnt by language models. Since we want our model probabilities to be high, it is the lower the PPL the better the model. On

the other hand, the amount of parameters in each language model was considered as an auxiliary metric specifically for the study of this section. This is because that we wanted to see how much parametric reduction could we yield with our simplified IRC and IHC network architectures of Table 4.2.

Experimental Setups

The codes of our GRNNs were based on the popular AWD-LSTM [Merity et al., 2018] repository⁵, and all GRNNs and their IRC counterparts were trained for 50 epochs. The hyperparameters were not fixed across all models; instead, the dimensionality and the amount of layers were taken from their respective original paper (if provided⁶). The settings are listed in Table 4.3, and we use the tuple of (ED, HD) to denote the combination of word embedding dimension and hidden dimension.

There were three special remarks on the experimental setups. First, the FastGRNN paper of Kusupati et al. [2018] tested language models with significantly different training settings to Merity et al. [2018]⁷. Thus we implemented 1-layer and 2-layer FastGRNNs and IRC-FastGRNNs with (650, 650). Second, the SRU paper of Lei et al. [2018a] did not include word-level language modelling as a task. Thus we tested 3-layer SRUs and IRC-SRUs with (1150, 1150). Third, we thoroughly trained 3-layer LSTMs and IHC-LSTMs with (400, 1150) for 800 epochs.

Results

As shown in Table 4.3, all of our IRC-GRNNs achieved lower PPLs than their baseline counterparts. Thus, all IRC-GRNNs had better generalisation ability than their vanilla GRNN counterparts. In addition, this was achieved while IRC-GRNNs had lower amount of learnable parameters. Redundancy removal was especially remarkable when the hidden dimension was much larger than the word embedding dimension. For instance, the IRC-T-LSTM removed up to 50.0% of the T-LSTM redundancies; while the IRC-SRU only removed up to 0.06% SRU redundancies. Finally, our IHC-LSTM was capable of achieving a competitive sub-60 perplexity. This matched the performances in Merity et al. [2018] despite a near-40% parametric reduction.

Remarks and Analyses

Our simplified model not only delivered parametric reduction but also improved the overall performances of all GRNNs. This showed our study not only connected the less opaque GRNNs to the well-studied dynamical systems literature for Hopfield networks, but we were also able to leverage these insights to highlight the correct network redundancies and simplify existing GRNN designs with our novel IRC and IHC modifications.

⁵Refer to <https://github.com/salesforce/awd-lstm-lm>.

⁶This footnote provides the source information of the quoted results, amount of layers, and (ED, HD). LSTM – Merity et al. [2018], FastGRNN – Kusupati et al. [2018], T-LSTM – Balduzzi and Ghifary [2016], GRU – Balduzzi and Ghifary [2016].

⁷Refer to page 6 of Section 3.2.1 of Kusupati et al. [2018] for their training scheme, and Refer to page 22 of Kusupati et al. [2018] for additional hyperparametric changes.

4.2.5.2 Image Generation

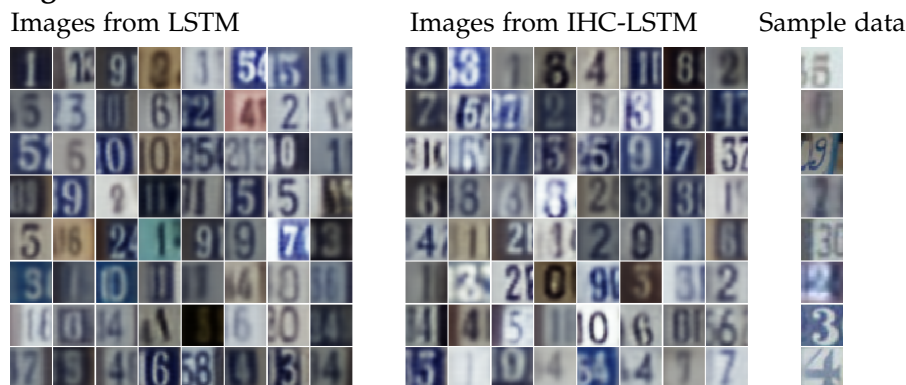


Figure 4.5: DRAW for Generating SVHN Images

Dataset and Experimental Setups

Gregor et al. [2015] proposed the *deep recurrent attention writer* (DRAW) architecture to generate highly realistic natural images. DRAW employed LSTM-decoders to sequentially adjust the generated pixels of an image in order to recursively provide more details to a canvas. This experiment used the multi-digit Street View House Numbers (SVHN) dataset [Goodfellow et al., 2014] which contains 231,053 training images, and 4,701 validation images.

Following Gregor et al. [2015], we employed single-layer LSTMs and IHC-LSTMs with 800 hidden dimensions to decode encoded pixel information. The encoded inputs were latent vectors of \mathbb{R}^{100} and the canvas was reassessed by the LSTM for 32 iterations as according to the paper. We used the pre-cropped Format 2 dataset provided by Goodfellow et al. [2014], and thus the canvas of our input images were of boxes of pixels of size 32×32 .

Results and Analyses

No metric was given in Gregor et al. [2015] for assessing the quality of the generated SVHN images. However, as we show in Figure 4.5, IHC-LSTMs were capable of generating images of similar qualities to those of the LSTM. The LSTM-decoders had 2.88 million parameters, while the IHC-LSTM-decoders had only 0.42 million, thereby removing 85.4% parameters.

The competitive performances exhibited by IHC-LSTM showed that, again, our identified network redundancies (the brown components of Table 4.2) were indeed secondary in nature. That is, our IHC simplification of which we derived from the dynamic Hopfield network had not deleted any important components that fine-tuned the recurrent neural processing within the original LSTM design.

4.2.5.3 Learning to Optimise

Besides our interest in AutoML, L2L was also selected as an experiment to test IHC-LSTM's ability to run online for a variable length of instances. See Section 3.2.5 for an in-depth task description for the L2L experimental setup.

Dataset and Experimental setup

We followed [Andrychowicz et al. \[2016\]](#) and trained base learners to classify images of the CIFAR-10 dataset [[Krizhevsky, 2009](#)] under the L2L framework. See Appendix A for more task description.

In the original experiment in [Andrychowicz et al. \[2016\]](#) employed meta-learners to update small learner networks of 3 convolutional layers with max pooling followed by a fully-connected layer. Their LSTM-optimisers was trained to update the learner for 100 steps; and during test time, was set to run freely to update the learner for 1000 steps. The amount of iterations in test time was set to be significantly longer than that in train time to test LSTM’s robustness against unseen error during training. For our own work, we selected the more difficult task to train the much larger ResNet-34 [[He et al., 2016a](#)] base learner networks.

Results

We tested LSTMs and IHC-LSTMs as neural optimisers. Following Section 3.2.5.5, we set both the input dimension and hidden dimensions as 20 for both architectural designs. The vanilla LSTM-optimisers had 3.28 thousands parameters, while our IHC-LSTM-optimiser had 2.82 thousands thereby removing 14.0%

of the total amount of trainable weights. The results shown in Figure 4.6 were the average performances of 50 ResNet-34s; and the figure was setup similarly to that in Section 4.1.6.2. Both neural optimisers significantly outperformed the handcrafted algorithms of SGD and Adam. Furthermore, the IHC-LSTM-based optimiser yielded ResNets with lower losses than those updated by the LSTM-based optimiser.

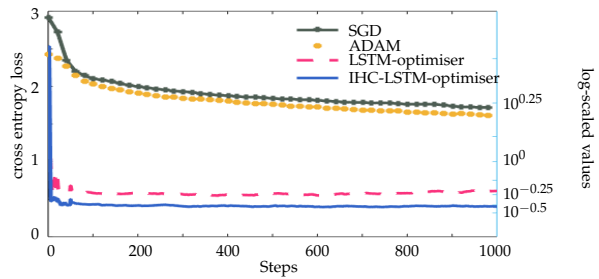


Figure 4.6: L2L for Updating ResNet-34s

Remarks and Analyses

The outcome of this experiment echoed to the language modelling findings. The superior performance exhibited by our IHC-LSTM design could be due to the simplification of LSTM of which in turn helped to increase model generalisation.

4.2.6 Section-wise Additional Related Work

This subsection will continue the grander discussion of the RNN related work in Section 3.1. Our main aim is to connect our findings via derivation (4.8) with other existing work. In particular, the mathematical insights that we have gained through our derivation echoed much with the observations remarked in the LSTM ablation study of [Greff et al. \[2016\]](#) (also see Section 3.1.6.1).

[Greff et al. \[2016\]](#) found that significant LSTM performance lesion occurred either when removing the output gate or by removing the forget gate. The authors considered the two lesions to be congenetic (caused by similar reasons). They hy-

pothesised that the output gate prevented unbounded cell states from destabilising learning; and that similarly, forget gate could suppress unbounded cell growths from relinquishing old cell values. They further conjectured that GRUs do not need output gates because their hidden state (the GRU equivalent to the LSTM cell) is bounded by the coupling of the input and the forget gate.

Our findings in this section can complement the observations made in [Greff et al. \[2016\]](#); but more importantly, we can further provide mathematical reasons to why that both GRUs and LSTMs function as successful GRNNs despite having distinct architectural designs. Recall that a GRU updates its hidden state via

$$\mathbf{h}_t = (1 - \mathbf{i}_t) \odot \mathbf{h}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t;$$

and that an LSTM updates its hidden state along with its cell state via

$$\begin{aligned} \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t). \end{aligned}$$

Our mathematical framework (see Section 4.2.2.2) from

$$\begin{aligned} \mathbf{h}'_t &= \mathcal{H}(\mathbf{h}_{t-1}, \mathbf{a}_t) = -\mathbf{h}_{t-1} + \mathbf{a}_t \text{ to} \\ \mathbf{h}_t &= (1 - \zeta) \odot \mathbf{h}_{t-1} + \zeta \odot \mathbf{a}_t \end{aligned}$$

found that the GRU input gate of \mathbf{i}_t served as the dynamic step size for updating the GRU \mathbf{h}_t . Hence, given the presence of \mathbf{i}_t , GRUs do not require an additional forget gate unit \mathbf{f}_t (if the two were bound). In fact, having a separate forget gate under the presence of an input gate can do more harm than good to the well-being of a network. This is because that when $\mathbf{f}_{mt} > 1 - \mathbf{i}_{mt}$ occurs for any neuron $m = 1, \dots, N$ in the forget and input gate, the relinquishment rate of $\mathbf{c}_{m(t-1)}$ is lower than the replenishment rate of \mathbf{a}_{mt} , and will cause \mathbf{c}_{mt} to increase uncontrollably and ultimately destabilising learning. This is hence the reason why LSTMs require the extra unit \mathbf{o}_t to restrict the exposure of the network memory from potential step size mismatches. This can also serve as a reason why [Greff et al. \[2016\]](#) found that the no-output-gate variant of LSTM caused significant performance lesion.

Our derivation (4.8) also complements the findings of [Jozefowicz et al. \[2015\]](#) (also see Section 3.1.6.2). Recall that [Jozefowicz et al. \[2015\]](#) conducted an extensive mutation scheme in hope to find better RNN designs. They were, however, unsuccessful in finding any one architecture that functioned consistently better than either GRUs and LSTMs. In contrast, our IHC-LSTMs outperformed LSTMs on the three challenging tasks of language modelling, image generation, and L2L.

We think that the main problem in [Jozefowicz et al. \[2015\]](#)'s approach was that they prohibited their candidate RNNs from developing new gated units and new hidden states. Their mutation scheme started strictly with the GRU and LSTM architectures and only permitted the modification of node-wise computation. As a result, all of their top three discovered architectures closely resembled that of the vanilla

GRU design. In contrast, our IRC and IHC designs in Section 4.2.4 introduced a new regulated variable \mathbf{v}_t and the new gated unit of Λ_t . Our IRC-GRNNs and IHC-GRNNs enjoyed better generalisation ability than their baseline GRNN counterparts and this suggests that the optimal RNN designs might be a network with multiple small units, instead of a network with few large units.

4.2.7 Section Conclusion

The work in this section modelled the GRU hidden states as discretised dynamic Hopfield networks. We interpreted the gated units as adaptive step sizes for network variables and highlighted redundant network components. Based on our insights, we presented the novel IRC as a technique to simplify GRNNs. We showed that theoretically and experimentally that the IRC is generally applicable to the GRNNs of GRU, LSTM, SRU, T-RNNs, and FastGRNN.

All IRC-GRNNs showed better generalisation abilities than their respective baselines for language modelling. These results were also achieved while the IRC-GRNNs all had fewer learnable parameters. Furthermore, IHC-LSTMs also showed competitive results for L2L and image generation. Most notably, the IHC-LSTM removed up to 85.4% LSTM parametric redundancies on the task of image generation.

Another important feature of this work was the successful establishment of connections between GRNNs to well-studied mathematical models. It is hence likely that more existing mathematical theories, such as *bifurcation analysis* [Glendinning, 1994], are also applicable for understanding and interpreting RNNs.

4.3 Chapter Summary

This Chapter consisted our two RNN studies of DecayNet [Kuo et al., 2018] (Section 4.1.8) and IRC [Kuo et al., 2020b] (Section 4.2). The aim for this chapter is to provide a mathematical study on the extensively applied GRU and LSTM RNNs. In contrast to ablation studies, search studies, and qualitative studies (see Section 3.1), our work returned to a first principles manner and attached mathematical meanings to the components of GRUs and LSTMs thereby increasing network interpretability.

While both studies were based on dynamical systems, DecayNet took an atomic approach whereas IRC took a architectural-based approach towards understanding the inner workings of RNNs. The former helped us to identify the cell state as the most important component towards LSTMs' ability in recurrent neural processing but it did not answer how the network components worked together collaboratively (see Section 4.1.8). Hence for our work on IRC, we aimed to establish connections between important gated RNNs, notably the GRU RNN, to the well-studied (dynamic) Hopfield network. We found that the LSTM input gate \mathbf{i}_t served a similar purpose to step sizes in dynamical systems (see Section 4.2.2.2); that GRNNs do not necessarily need a forget gate \mathbf{f}_t in the presence of an input gate (see Section 4.2.6); and that the output gate \mathbf{o}_t only existed to prevent the cell state from growing unboundedly when $\mathbf{f}_t > 1 - \mathbf{i}_t$. The mathematical connection also helped us to classify several GRNNs (including GRUs and LSTMs) as one class of mathematical models. This not only helped us to understand why various RNN architectures work despite their differences, but more importantly it highlighted several architectural redundancies. Based on our acquired insights, we proposed IRC as a generic simplification over 5 GRNN designs (see Table 4.2).

Our readers should note that though GRUs and LSTMs remain the most popular RNNs, RNN is an enormous neural network archetype that consists more than just GRNNs. For instance, Echo State Networks [Jaeger, 2001] and Clockwork RNNs [Koutnik et al., 2014] are both successful designs that function without gated units and are based on drastically different principles to GRUs and LSTMs. Hence we emphasise that our work presented in this chapter are only useful and valid for analysing those RNNs that utilise gated units. The reason why we have included GRNN studies in this thesis was for understanding AutoML via having a firmer grasp on the mathematical properties of shortcut connections (see Section 1.4).

4.4 Chapter Conclusion

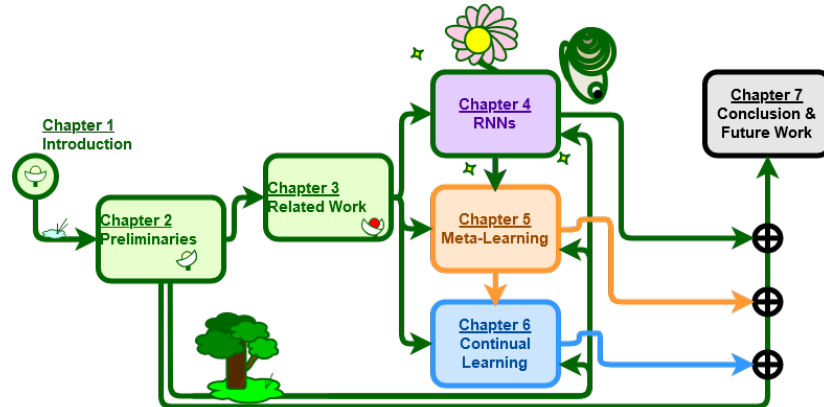


Figure 4.7: Progress Update

As we motivated in Section 1.4 and in Section 2.1, shortcut connections made deep learning possible and it continues to push the state of the art results for various machine learning disciplines. We studied RNNs because that the shortcut connection in the deep Highway Network [Srivastava et al., 2015] was inspired by the architectural design of the LSTM RNN [Hochreiter and Schmidhuber, 1997b]. Thus in order to develop better AutoML methods, it is necessary to understand shortcut connections for how to automate practical and deep learning models.

Our GRNN studies provide a two-fold benefit towards understanding the AutoML techniques considered in this thesis. First, it fits cohesively to our goals for understanding meta-learning in the next chapter (Chapter 5). In fact, we had already employed RNNs for the meta-learning task of L2L (see Section 3.2.5) in this chapter (see Sections 4.1.6.2 and 4.2.5.3). Second, we will later show in our continual learning chapter (Chapter 6) that deep networks with shortcut connections can be continually fine-tuned to acquired more skills by only modifying a small set of selective layers.

Our RNN research is yet to be extended below in the thesis. Though this chapter showed some of our work in understanding the theoretical principles behind recurrent neural processing, we have yet to discuss the more practical applications of RNNs in AutoML. These research will be covered in our next chapter.

Learning to Optimise beyond Rapid Knowledge Acquisition

This chapter introduces our contribution on meta-learning research. The content therein this chapter continues the discussion in Section 3.2 of the related work and heavily borrows concepts from Section 2.2 of the preliminaries.

This thesis focuses on the optimisation-based approach of [Andrychowicz et al. \[2016\]](#) as detailed in Section 3.2.5. Their paper proposed *Learning to Learn* (L2L) – to train an RNN to update the weights of another neural network in place of gradient descent. Their RNN of choice was the LSTM RNN [[Hochreiter and Schmidhuber, 1997b](#)] of which we studied carefully in Chapter 4.

Recall that gradient descent of Equation (1.4) provides the update of

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}$$

to a base learner \mathcal{F} with parameter θ ; and that an *RNN-optimiser* \mathcal{Q} substitutes this update by learning Equation (3.12) which prepares g_t to update the base learner via

$$\begin{aligned} \theta_{t+1} &= \theta_t + g_{t+1} \text{ and} \\ [g_{t+1}, \mathbf{c}_{t+1}, \mathbf{h}_{t+1}] &= \mathcal{Q}(\mathcal{X}_t^{(\mathcal{Q})}, \mathbf{c}_t, \mathbf{h}_t | \eta). \end{aligned}$$

L2L was already included as a part of our experiments in our RNN studies. And as shown in Figure 4.3, the RNN-optimisers updated base learner parameters faster than their handcrafted rivals of, say, SGD [[Robbins and Monro, 1951](#)].

Our work on meta-learning concerns on expanding L2L capability beyond rapid knowledge acquisition. Our contribution is two-fold. First, we showed that *RNN-optimisers* were extremely limited in their capability – they were only able to update one base learner architecture for one dataset. Our work of [Kuo et al. \[2020c\]](#) introduced a modifications that allowed the RNN-optimisers to update base learners on an unseen dataset thereby increasing their general applicability. Second, we propose a novel algorithm in [Kuo et al. \[2020a\]](#) to learn an optimisation scheme which removed redundant network connections on the fly. This allowed network pruning to be achieved much more efficiently without the need of repetitive re-training.

But Why Does an RNN Meta-learner Work Well?

Though the seminar work of [Andrychowicz et al. \[2016\]](#) showed that an RNN could learn to update another independent base learner, their paper lacked a mathematical explanation on why this was possible. Hence before we discuss our novel contribution, it is again beneficial to analyse and visualise the mathematics behind the L2L mechanism. In particular, we will use concepts on loss landscape and minima of which we previously covered in Section 2.2.

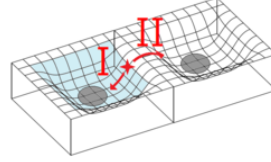


Figure 5.1: Selecting Minimum

Recall that a neural optimiser is updated as according to Equation (3.14) where

$$\mathcal{L}(\eta) = \mathbb{E} \left[\sum_{\tau=1}^{\Xi} \mathcal{L}_{\tau+1}(\theta_{\tau}) \right].$$

The neural optimiser \mathcal{Q} is hence dependent on the objective loss \mathcal{L} of \mathcal{F} over Ξ time steps. Say that our neural optimiser is meta-trained on two copies of randomly initialised base learners ($\mathcal{V} = 2$ for line 4 of Algorithm 1), we will ideally want

$$\mathcal{L}(\eta_{\mathcal{V}=2}) \leq \mathcal{L}(\eta_{\mathcal{V}=1}) \quad (5.1)$$

so that the neural optimiser can *further* fine-tune its performance on updating base learner 2 after previously providing erroneous updates to base learner 1. To elaborate even further, the fact that we desire the loss of base learner 2 to be smaller than the loss of base learner 1 in a fixed amount of updates¹ meant that conceptually, when the neural optimiser guides the base learner to traverse across the terrain on the loss landscape, we want it to aim for the shortest path towards the *closest minimum*.

To gain a deeper intuitions for L2L, consider the scenario presented in Figure 5.1 where multiple minima concurrently exist in the vicinity of θ in the parametric space of base learner \mathcal{F} . The figure illustrates the hypothetical scenario which a base learner \mathcal{F} has only been lightly trained and that now its parameters θ sits in the “entrance” of one minimum (shaded in blue) on the loss landscape. Given that there are now two distinct minima (shaded in grey), we question:

Which minimum will L2L update θ towards to?

L2L will update the parameters towards minimum I situated on the left of the figure, *the minimum closest to the current position*.

We based our inference on Equation (3.14) where the neural optimiser will favour a route towards any one non-specific minimum which happens to incur the least amount of cumulative loss of $\mathcal{L}(\eta)$ over the Ξ steps in the trajectory. To elaborate,

¹This refers to the fact that \mathcal{L} is designed as the mean of the loss \mathcal{L} of the base learner over the fixed bandwidth of Ξ iterative updates from \mathcal{Q} to \mathcal{F} . See Algorithm 1 for more details.

since the starting location marked with the **star** sign is already placed in a “downwards” slope towards minimum **I**, it will then avoid the need to “climb over a mountain” on the loss landscape of which will otherwise accumulate a larger total loss. That is, a trained neural optimiser will favour the route towards minimum **I** over that towards minimum **II**, even if minimum **II** is better than minimum **I**.

This inference tells us that since it is impossible to escape a minimum with L2L, the neural optimiser is hence capable of updating a base learner much faster than handcrafted optimisations of, say, SGD [Robbins and Monro, 1951]. However, this inference only supports that a neural optimiser will be able to find one minimum in the periphery of the starting location of θ_0 from a network random initialisation. That is, this argument does not guarantee the quality of the said minimum; the sought minimum is nearly always likely to be a local minimum instead of the global minimum. Also, this argument does not guarantee that a neural optimiser can successfully find a minimum if there were no minimum in the vicinity of an unlucky randomly initialised θ_0 . Nonetheless, there have been a lot of empirical evidences showing that neural optimisers perform much better than their handcrafted counterparts.

Learning to learn has brought us powerful strategies to update a base learner.

But what are their limitations?

The following section discusses input-domain heterogeneity from the perspective of neural recurrent processing.

Section Content:

- Section 5.1.1: Section Introduction,
- Section 5.1.2: Reasons of the Limitations of the Vanilla Neural Optimiser,
- Section 5.1.3: Problems of Employing RNNs as Generally Applicable Optimisers,
- Section 5.1.4: Context Awareness through HyperNetwork,
- Section 5.1.5: Experiments and Results,
- Section 5.1.6: Section-wise Additional Related Work,
- Section 5.1.7: Section Conclusion:
Increasing L2L's General Applicability through
Addressing Input-Domain Heterogeneity,
- Section 5.1.8: An Evaluation on Our Solved Problem.

5.1 MTL2L: A Context Aware Neural Optimiser

This section is based on our work in Kuo et al. [2020c] which investigates the general applicability of L2L neural optimisers. We will show that conventional L2L neural optimisers are limited to update one specific type of learner on one specifically designated dataset. To address this problem, we will provide modifications on the standard LSTM design covered in Chapter 4 to address input-domain heterogeneity.

5.1.1 Section Introduction

Learning to learn (L2L) trains a meta-learner to assist the learning of a task-specific base learner. This technique has shown to achieve rapid knowledge acquisition to converge a base learner with fewer amount of iterations than the handcrafted algorithms of SGD, Adam and *etc.* As a reference, see our previous experimental remarks in Figure 4.3. L2L achieves fast assimilation of knowledge by casting optimisation as the learning target for a RNN meta-learner to explore meta-knowledge, such as information on the loss landscape (see here), to describe the structure of the task in the task domain (see Section 3.2.5). However, we will show that existing neural optimisers are limited to update *one specific type of learner on one specifically designated dataset* and this largely restricts the practical utility of the neural optimisers.

The work in this section concerns *input-domain heterogeneity*. Our goal is to expand the general applicability of a trained neural optimiser across multiple different datasets. This will be achieved with *Multi-Task Learning to Learn* (MTL2L). MTL2L is an RNN-optimiser based on a modified LSTM design which achieves contextual awareness for input data by allowing self-modification on its optimisation rules dependent on the network input.

5.1.2 Reasons of the Limitations of the Vanilla Neural Optimiser

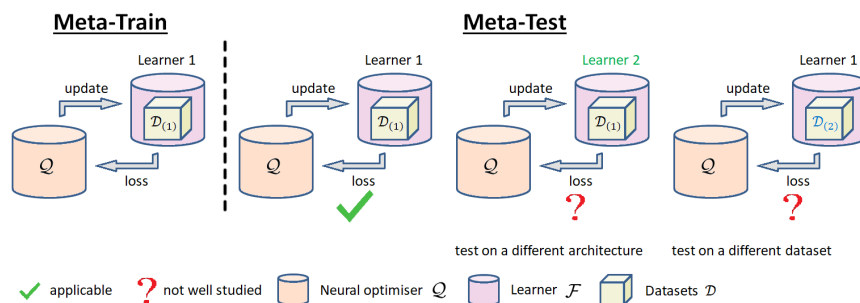


Figure 5.2: Neural Optimiser Applicability

Let us recall the meta-training scheme of the neural optimiser shown in Algorithm 1, consider its algorithmic inputs, and discuss 2 problematic premises. First, that the training dataset $\mathcal{D}_{\text{Train}}$ is specified for meta-training. And second, that the base learner architecture \mathcal{F} is also required for meta-training. This means that the vanilla neural optimiser in Andrychowicz et al. [2016] was thus designed to update

one specific learner architecture for one specifically designated dataset. The vanilla neural optimiser \mathcal{Q} will hence face 2 problems during meta-testing of which we illustrate in Figure 5.2.

- (i) If meta-trained on $\mathcal{F}^{(1)}$, it would then fail to update $\mathcal{F}^{(2)}$; &
- (ii) if meta-trained for $\mathcal{F}^{(1)}$ to solve $\mathcal{D}_{(1)\text{Train}}$, it would then fail to update $\mathcal{F}^{(1)}$ for solving $\mathcal{D}_{(2)\text{Train}}$.

Meta-learners were originally developed to encapsulate information of the structure of the task domain in order to assist the conjoint training of a task specific base learner [Rendell et al., 1987]; hence the student-author thinks that it is understandable that the vanilla neural optimisers struggle against problem (i). To elaborate more, let us consider a neural optimiser trained to update an MLP but later tested to update a ResNet [He et al., 2016a]. MLPs are known to use linear layers whereas ResNets are mostly defaulted with convolutional layers; hence it is questionable whether a trained neural optimiser could be used across different layer-functions. This concern was indirectly highlighted in the recent work of Sun et al. [2021] on parameter corruption; the authors of that paper found that convolutional layer weights were much easier to be corrupted than linear layer weights. To us, their findings hence indicated that different layer functions should be updated with different strategies. Further note, it also often requires some manual changes when we update two distinct networks with the same handcrafted gradient descent algorithm. For instance, it is commonly known that MLPs could be trained with the much larger learning rate of 0.01; while architectures with shortcut connections (including TCN [Bai et al., 2018]) often require the much smaller learning rate such as 2×10^{-3} .

Another noteworthy remark on problem (i) is that, the work of Andrychowicz et al. [2016] also found that if a neural optimiser was trained to update a MLP with ReLU activations, then that neural optimiser would also fail² to update the same MLP even only when the activation functions were changed to tanh. This is also understandable as architectural changes in networks not only explicitly change the forward pass but also implicitly impact the backward pass. See Section 2.1.1 on a comparison between the backward passes of MLPs and ResNets; and also see Section 2.1.2.5 (and their corresponding ICCs) on how including normalisation techniques could impact the backward passes. Recall that neural optimisers were originally employed to explore strategies to update a designate base learner, hence if the neural optimiser had successfully learned the 0 or 1-derivative of ReLUs, it is then natural that the same inferred logic would fail for the *bowling ball*-derivative of tanh.

Problem (i) has actually been addressed in a follow-up paper by the authors of Andrychowicz et al. [2016] of which we will discuss later. Therefore our first work on meta-learning is directed to focus on our identified problem (ii).

Compared to problem (i), the student-author sees problem (ii) as a conceptually more troubling complication. To the best of our knowledge, there is not any paper that has addressed this specific problem. This problem is troubling because that, meta-learners are supposed to gather information in the *task-domain*; for instance,

²Refer to page 6, and especially the right-most subplot of Figure 5, in Andrychowicz et al. [2016].

once learned to update a learner to classify images on MNIST [LeCun et al., 1998], the neural optimiser should in principle be able to update an identical learner \mathcal{F} but to classify images of Cifar10 [Krizhevsky, 2009]. This is because that both problems belong to the same task of image recognition. However, Figure 5.3 shows us that this is not the case. (The datasets will later be described in Section 5.1.5.)

Figure 5.3 addresses two meta-learning scenarios and has the same layout as Figure 4.3. Subplot (a) illustrates meta-learning under the common assumption: *to meta-test on the single identical input-domain used in meta-training*; and that subplot (b) depicts deteriorated performance in meta-testing on an unseen dataset. The failure shown in subplot (b) implies that the neural optimiser could not update the base learner due to the distinct feature distributions in the two datasets. To elaborate, the vanilla neural optimiser failed to learn an alternative strategy to gradient descent due to impacts from the forward pass of the base learner (the changes in the input features) rather than the changes therein the backward pass (no changes as the base learner architecture remains the same).

The work in this section is directed to remedy problem (ii) from an architectural perspective of the neural optimiser. As shown in Chapter 4 and Section 2.1, sensible changes in the network design can usually address difficult problems.

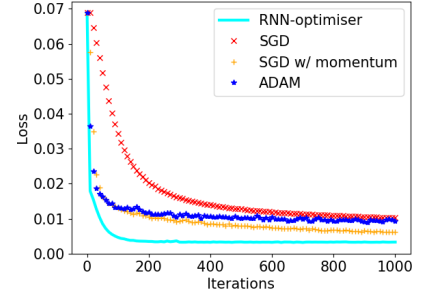
5.1.3 Problems of Employing RNNs as Generally Applicable Optimisers

As described in Section 3.2.5, Andrychowicz et al. chose the LSTM RNN as the meta-learner neural optimiser. Recall from Section 3.1.2 that an LSTM is formally written as

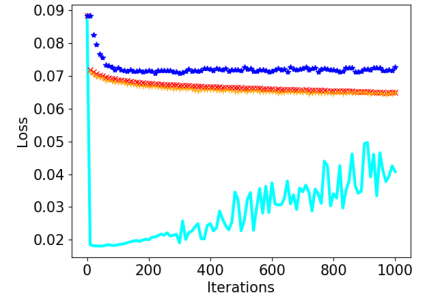
$$\begin{aligned} \mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t &= \sigma(\mathbf{W}_{\{F,I,O\}}\mathbf{x}_t + \mathbf{U}_{\{F,I,O\}}\mathbf{h}_{t-1} + \mathbf{b}_{\{F,I,O\}}), \\ \mathbf{a}_t &= \tanh(\mathbf{W}_A\mathbf{x}_t + \mathbf{U}_A\mathbf{h}_{t-1} + \mathbf{b}_A), \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \text{ and} \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t); \end{aligned}$$

and from Equation (3.12) that an LSTM-optimiser \mathcal{Q} updates a learner \mathcal{F} via

$$\begin{aligned} \theta_{t+1} &= \theta_t + g_{t+1} \text{ and} \\ [g_{t+1}, \mathbf{c}_{t+1}, \mathbf{h}_{t+1}] &= \mathcal{Q}(\mathcal{X}_t^{(\mathcal{Q})}, \mathbf{c}_t, \mathbf{h}_t | \eta). \end{aligned}$$



(a) Meta-trained on MNIST and Meta-tested on MNIST



(b) Meta-trained on MNIST and Meta-tested on Mod-Cifar10

Figure 5.3: Input Variance

LSTM was hence chosen for its network hidden variables. That is, the LSTM-optimiser was designed to encode sequential events in the task-domain into its hidden variables to infer parametric updates g_t for the learner parameters of θ_t .

Following this line of thoughts, we now know that the well-being of an trained LSTM-optimiser is highly dependent on its hidden state \mathbf{h}_t and cell state \mathbf{c}_t . Furthermore, the well-being of these hidden variables are determined by two factors: the LSTM-optimiser parameters of η , and the LSTM-optimiser input of $\mathcal{X}_t^{(\mathcal{Q})}$. Since the LSTM-optimiser is easily influenced by $\mathcal{X}_t^{(\mathcal{Q})}$, it is thus susceptible to problem (ii). This also reveals that one way to remedy problem (ii) is through re-designing η .

The LSTM-optimiser parameters of η include forward synapses \mathbf{W} s, recursive synapses \mathbf{U} s, and biases \mathbf{b} s; and they are our target of modifications. Compared to feedforward neural networks, RNNs trade *recursiveness* for *depth*. There are advantages and disadvantages to this trick of the trade.

On one hand, RNNs are much smaller than deep feedforward neural networks and require less parameters. Instead of utilising a large set of similar layer-functions (say 50 layers of [2 * (convolutional layers + ReLU) + shortcut connection]), RNNs can efficiently and elegantly encode information via incrementally updating their hidden variables (see Chapter 4) through recursions. On the other hand, the main complaint for RNNs regards their difficulty in training [Pascanu et al., 2013] and also their computational inefficiency [Lei et al., 2018a; Bai et al., 2018]. In this section however, we will look at a less obvious disadvantage of recursion – the limitation of RNN’s *expressiveness*.

Recall in Chapter 4 that we look at (gated) RNNs from a dynamical systems approach [Glendinning, 1994]. We described RNNs as systems of difference equations (see Section 4.1) and as systems of differential equations (see Section 4.2) because that RNNs recursively execute *a set of pre-defined operators* to update their hidden variables. The pre-defined operators of an LSTM consists of, say, the forget gate $\mathbf{f}_t = \sigma(\mathbf{W}_F \mathbf{x}_t + \mathbf{U}_F \mathbf{h}_t + \mathbf{b}_F)$. Note that for every instance, no matter for $t = 1$ or that for $t = 5$, the forget gate executes the same forward synapse of \mathbf{W}_F and the same recursive synapse of \mathbf{U}_F . Compared to a deep MLP with layer operator \mathcal{H}^L s that trades recursion with depth, the MLP’s first layer is executed with function $\mathcal{H}^{L=1}$ while its fifth layer is executed with function $\mathcal{H}^{L=5}$ where $\mathcal{H}^{L=1} \neq \mathcal{H}^{L=5}$. The repetitive RNN executions of its identical \mathbf{W}_F and \mathbf{U}_F thus limits the expressiveness of an RNN. This also means that an LSTM-optimiser reuses the same learnt logic to prepare for updates g_t for the learner parameters θ_t regardless of the source dataset for input $\mathcal{X}_t^{(\mathcal{Q})}$. This hence tells us that, to solve problem (ii) through re-designing the LSTM-parameters of η , we will need to re-structure LSTM such that it can implement its own weights dependent on the input. Or in other words, to become context aware of the input data.

5.1.4 Context Awareness through HyperNetwork

We took a *HyperNetwork* [Ha et al., 2017] approach towards dataset-agnostic meta-learning. HyperNetworks are small auxiliary networks that re-implement weights

for a larger network based on input data. Besides re-implementing its own weights, we desire our modified LSTM to satisfy a second property: to be able to compare and to contrast different datasets. That is, we want a modified LSTM that remembers the feature distributions of those datasets that it has observed during meta-training; and that during meta-testing, compare features of an unobserved dataset to those seen during meta-training and update a base learner with the meta-knowledge on the similarities and dissimilarities among different datasets. To achieve the second desiderata, we required our modifications to be geometrically interpretable.

5.1.4.1 The MTL2L Formulation

We propose to substitute the vanilla LSTM-optimiser with our novel *Multi-task Learning to Learn* (MTL2L) neural optimiser provided as below

$$\begin{bmatrix} \mathbf{f}_t \\ \mathbf{i}_t \\ \mathbf{o}_t \end{bmatrix} = \sigma \left(\begin{array}{l} \mathcal{W}_t^{(\mathcal{F})} \mathbf{x}_t + \mathcal{U}_t^{(\mathcal{F})} \mathbf{h}_{t-1} + \mathbf{b}_{\mathcal{F}} \\ \mathcal{W}_t^{(\mathcal{I})} \mathbf{x}_t + \mathcal{U}_t^{(\mathcal{I})} \mathbf{h}_{t-1} + \mathbf{b}_{\mathcal{I}} \\ \mathcal{W}_t^{(\mathcal{O})} \mathbf{x}_t + \mathcal{U}_t^{(\mathcal{O})} \mathbf{h}_{t-1} + \mathbf{b}_{\mathcal{O}} \end{array} \right), \quad (5.2)$$

$$\mathbf{a}_t = \tanh(\mathcal{W}_t^{(\mathcal{A})} \mathbf{x}_t + \mathcal{U}_t^{(\mathcal{A})} \mathbf{h}_{t-1} + \mathbf{b}_{\mathcal{A}}) \quad (5.3)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t, \quad \text{and} \quad (5.4)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t); \quad \text{where} \quad (5.5)$$

$$\mathcal{W}_t = \mathcal{Q} \mathcal{I}_t \mathcal{P}^T, \quad (5.6)$$

$$\mathcal{U}_t = \mathcal{S} \mathcal{K}_t \mathcal{J}^T, \quad \text{such that} \quad (5.7)$$

$$\mathcal{I}_t = \text{diag}(\mathcal{N}^{(1)}(\mathbf{x}_t)), \quad \text{and} \quad (5.8)$$

$$\mathcal{K}_t = \text{diag}(\mathcal{N}^{(2)}(\mathbf{x}_t)). \quad (5.9)$$

As shown above, the MTL2L formulation from Equation (5.2) to Equation (5.5) is largely the same as that of an LSTM. Whereas Equation (5.6) to Equation (5.9) entails the modifications made to the components coloured in brown. These modifications turn RNN weights into random variables to increase network expressiveness.

The brown components of the MTL2L are the synaptic connections of the RNN. We explicitly deconstruct these connections through the *singular value decomposition* (SVD) of Equations (5.6) and (5.7). We use SVD to split each synapse as fixed purple terms and orange components that require learning. The former are constant orthonormal matrices, and the latter are diagonal matrices with eigenvalues placed in the non-zero entries.

We infer the eigenvalues of the orange matrices with hypernetworks $\mathcal{N} = \mathcal{N}(\mathbf{x}_t)$. These hypernetworks are specifically designed in conjunction with SVD so that we can use different combinations of eigenvalues to represent for the different datasets that the MTL2L has observed during meta-training. The MTL2L RNN thus becomes context aware and self-modifies its optimisation rules based on input data.

It is natural for an MTL2L to support multi-task learning; and the purpose of the built-in SVDs is to explicitly describe the input features of all datasets. For all

datasets seen in meta-training and those to-be-observed during meta-testing, the SVD in the MTL2L describes all of their underlying features as elements belonging to one unified domain space spanned by the fixed purple eigenvectors. Through this design, the eigenvalues enable an MTL2L to compare the task structures of unseen inputs to those that it would have observed during meta-training. This will enable the MTL2L RNN to “think outside the box” and treat features of an unobserved dataset as a mixture of those features that it would have observed during meta-training. To elaborate and in face of unobserved input features, MTL2Ls will infer new sets of eigenvalues to create new parameterisation rules to prepare updates for solving unseen datasets.

5.1.4.2 MTL2L Implementation Details: SVD

Equations (5.6) and (5.7) describe the MTL2L synapses with SVD. Abiding to the SVD formulation, the purple components are orthonormal matrices, whereas the orange components are diagonal matrices with eigenvalues (or singular values, to be more precise). We emphasise again that the purple matrices are fixed prior to training, and only the orange components are trained with hypernetworks. More on the orange components will be addressed later.

The matrices are required to be specifically $\mathbb{R}^{20 \times 20}$. This is chosen as MTL2L inherits the default setup of the LSTM-optimisers of Section 3.2.5.6 and has the hidden dimension size of 20. Any means of orthonormal matrix generation is viable. The student-author did his majority of work in Python [Rossum, 1995], and chose to use the SciPy library [Virtanen et al., 2020]. Refer to https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.ortho_group.html for guidelines.

5.1.4.3 MTL2L Implementation Details: Naïve Hypernetwork Eigenvalues

The hypernetwork \mathcal{N} s for MTL2L can be simple MLPs. By default, we will use MLPs with two linear layers. One technical details regards the activation functions. The first layer is followed by a ReLU activation, and that there is no activation function for the second linear layer. This allows the inferred eigenvalues to be either positive or negative. A second technical detail is that the output dimension of the second linear layer is required to be specifically 20 to match the dimensionality of the orthonormal matrices of Section 5.1.4.2. The input dimension of the first layer is task-dependent, and we default the hidden dimensions of the hypernetwork-MLPs as 32.

5.1.4.4 MTL2L Implementation Details: Stable Hypernetwork Eigenvalues

Under this setting, we formulate the hypernetworks as

$$\mathcal{N}_t = 0.9\mathcal{N}_{t-1} + 0.1\text{MLP}(\mathbf{x}_t), \text{ with} \quad (5.10)$$

$\mathcal{N}_0 = \vec{\mathbf{0}} \in \mathbb{R}^{20}$. The hypernetwork-MLPs employed in Equation (5.10) were identical to those employed in Section 5.1.4.3. However and if needed, we can use this

particular setup to make the changes in the inferred eigenvalues act like momentum [Nesterov, 1983] to increase training stability.

5.1.4.5 MTL2L Implementation Details: On Multi-Task Learning

This subsection will borrow notations from Section 3.2.5.6. For the experiments of this section, we set $\Xi = 5$, $\mathcal{V} = 20$, and $\mathcal{S} = 100$ for all LSTM-optimisers. As for MTL2L, we set $\mathcal{V} = 30$ with prolonged learning trials. The learning trials were lengthened to support multi-task learning during the MTL2L meta-training phase to ensure enough exposure from each meta-training dataset. We meta-trained MTL2Ls with alternating datasets. For instance, if an MTL2L were to meta-train on \mathcal{D}_A and \mathcal{D}_B , we would train the MTL2L to update base learners to solve \mathcal{D}_A during the learner trials of $\mathcal{V} = 1, 3, \dots, 29$, and similarly to update base learners for solving \mathcal{D}_B during even-number \mathcal{V} s.

5.1.5 Experiments and Results

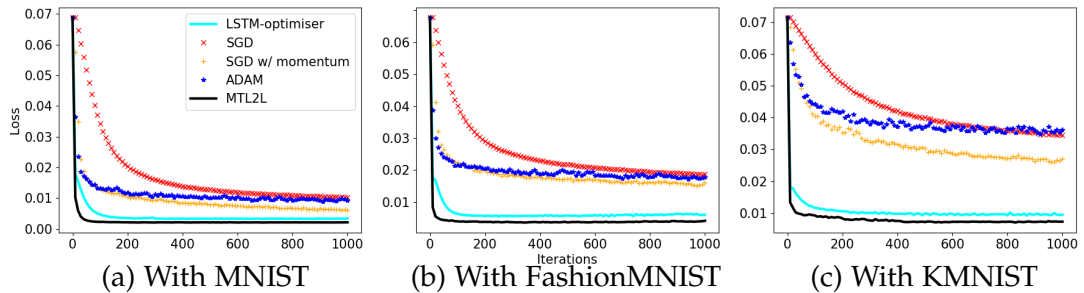


Figure 5.4: The Well-behaving Meta-Testing Results of Scenario 1

Experimental Setup

We repeated our experiments in Section 4.1.6.2 and employed neural optimisers to update MLP base learners. The neural optimisers were meta-trained to update the MLPs for $\mathcal{S} = 100$ steps (see Section 5.1.4.5), but meta-tested to update the MLPs for the much longer 1000 steps. The neural optimisers are judged based on two attributes of the loss curves of their base learners. The sharper the slope of the curve, the more efficient the optimiser was in converging a base learner; and the lower the converged loss, the better the final generalisability of a base learner.

In this section, we also extended previous experiments to consider the two experimental scenarios detailed as below.

Scenario 1: Test for the common meta-training assumption

First meta-train neural optimisers to update MLP learners to classify images on one single dataset; and then meta-test by employing the learnt neural optimisers to update MLP learners to classify *the identical dataset* used in meta-training.

Scenario 2: Test for input-domain heterogeneity

Unlike **Scenario 1**, meta-test to update learners to classify *an unseen dataset*.

We tested MTL2Ls against LSTM-optimisers, and against the hand-designed optimisation algorithms of vanilla SGD, SGD with momentum 0.9, and Adam; all with learning rates 0.01. We presented the learning curves of the learners and colour-coded them. Black solid lines are for MTL2Ls, cyan solid lines for LSTM-optimisers, red ‘×’ for vanilla SGD, yellow ‘+’ for SGD with momentum, and blue ‘*’ for Adam.

Datasets

We updated the MLPs on MNIST [LeCun et al., 1998], FashionMNIST [Xiao et al., 2017], KMNIST [Clanuwat et al., 2018], and Cifar10 [Krizhevsky, 2009]. Refer to Appendix A for descriptions on all datasets.

Scenario 1 considered MNIST, FashionMNIST, and KMNSIT whereas **Scenario 2** considered a synthesised Cifar10 dataset. Both FashionMNIST and KMNIST were specifically designed with the same format to MNIST. They have images with 10 classes of items on 28×28 pixels. However, Cifar10 contains 32×32 RGB-colour images. In order to fit these images in the MLP learner, we created a synthesised Cifar10 dataset consisting of greyscale 28×28 pixel images. The images were made greyscale by combining the coloured filters with ratio $0.30R + 0.59B + 0.11G$. The images were made 28×28 by only selecting the first 28 rows and the first 28 columns.

Results

The results for **Scenario 1** are shown in Figure 5.4. We implemented our hypernetworks naïvely as according to Section 5.1.4.3. Like what we have previously seen in Section 4.1.6.2, we found that the neural optimisers updated learners more rapidly than their handcrafted counterparts for all datasets. MTL2Ls also converged the learners more efficiently and yielded lower losses than LSTM-optimisers.

The results for **Scenario 2** are shown in Figure 5.5. For this scenario, we implemented the stable hypernetworks of Section 5.1.4.4. Unlike Figure 5.3, we see that our MTL2L optimiser in black could successfully update base learners on an unseen synthesised dataset.

Since LSTM-optimisers only supported single-task learning, they were meta-trained to update MLPs for KMNIST. In contrast and as MTL2Ls naturally supports multi-task learning, we meta-trained them to update MLPs on both Fashion-MNIST and KMNIST (see the setup in Section 5.1.4.5). During meta-testing, all optimisers updated MLPs on a modified re-sized greyscale Cifar10 dataset. As shown in the figure, MTL2Ls did not yield learners with deteriorating performances and displayed stability.

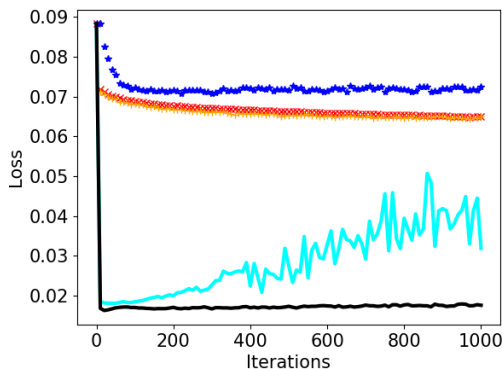


Figure 5.5: Meta-tested on an Unseen Mod-Cifar10 Dataset

Remarks and Analyses

For **Scenario 1**, it was interesting to see that Figure 5.4 reflected that the MTL2L-optimiser performed better than the LSTM-optimiser. This is because that we only introduced the modifications (see Section 5.1.4.1) to address input-domain heterogeneity. We attribute this welcoming surprise to the sequential multi-task setup addressed in Section 5.1.5. That is, by interleaving and exposing the MTL2L-optimiser to a variety of different datasets, the neural optimiser could have learnt a much more generalised update strategy from a diversified task domain.

For **Scenario 2**, it was also interesting to observe the *timing* of which the LSTM-optimiser exhibited instability. As shown in Figure 5.5, the LSTM-optimiser instability was less severe in the first 400 iterations of the base learner updates – there was a slight upraise in the loss curve but there were no volatile oscillations. The volatility was only amplified after the first 400 iterations; and this could hence be interpreted that there existed some underlying similarities between the task structures for updating a learner on KMNIST against that on the unseen synthesised Cifar10. Hence we conjectured that, though with imprecision, LSTM-optimisers were actually able to capture the “overall policy” in updating learners for an unseen dataset. Nonetheless, such imprecision introduced errors which reiterated in a vanilla LSTM for multiple instances and thus resulted in the eventual deteriorating performances.

5.1.6 Section-wise Additional Related Work

Remedies to Low Adaptability in Learnt Optimisers

Recall in Section 5.1.2, we discussed problem (i) such that it is not well studied whether neural optimisers are capable of updating base learners of a different architecture during the meta-testing phase. We mentioned how [Andrychowicz et al. \[2016\]](#) demonstrated in their own work that a learnt vanilla LSTM-optimiser was not able to handle the change of MLP learner activation functions from ReLU to tanh. This negative result motivated the subsequent work of [Wichrowska et al. \[2017\]](#) by the original panel of authors. Problem (i) was also independently studied in [Lv et al. \[2017\]](#). These 2 papers addressed different solutions to overcome this problem.

In [Lv et al.](#), the authors solved this issue by taking inspiration from the hand-crafted adaptive optimisation algorithms. They explicitly modified the input of their RNN-optimiser similarly to Adan to reduce current derivatives based on past derivatives. Whereas in [Wichrowska et al.](#), the authors presented a solution which took inspiration from the processes of relation classification in text [[Luo, 2017](#)]. They presented a hierarchical multi-layered RNN-optimiser capable of abstracting information of base learner gradients at different levels.

However, these 2 proposals also have some clear disadvantages. First, the proposed idea of [Lv et al.](#) required numerous manual engineering tricks and hence conflicted with the idea of “learning an optimiser from the scratch”. Second, the hierarchical RNN-optimiser of [Wichrowska et al.](#) significantly increases the meta-learner complexity. In contrast, our MTL2L optimiser was directly extended from [Andrychowicz et al.](#)’s LSTM-optimiser with minor modifications (see Equations (5.2)

to (5.9)) because our design was based on geometric assumptions of the task structures and on our understanding for RNN dynamics from Chapter 4.

We would like to emphasise that these mentioned remedies are all limited specifically for neural optimisers that are based on RNNs. For instance, MetaSGD [Li et al., 2017] is a neural optimiser without an RNN. It takes the amount of parameters of the base learner as a hyper-parameter and conducts meta-learning via training a unique learning rate per parameter for the base learner. Thus by design, a trained MetaSGD neural optimiser is inapplicable to meta-test on any alternative base learner architectures with mismatching total amount of parameters.

Other Types of Learnt Optimisers

Both the original work by Andrychowicz et al. [2016] and our novel study in this section have mainly focused on learning a neural optimiser for the supervised learning paradigm. However and since the seminal paper of Andrychowicz et al., there has also been a large amount of research devoted to learn neural optimisers for reinforcement learning (RL).

For instance, the original panel of authors of Andrychowicz et al. also extended their work and trained an RNN to learn RL in Chen et al. [2017]. Similarly and as shown in Wang et al. [2016], an RL algorithm could also be learnt by another independent RL algorithm. Their key concept was to use a standard RL to train an RNN which implements its own free-standing RL procedure. The latter idea has also been implemented in an adjacent field of meta-learning. In neural architecture search, Zoph and Le [2017] showed that an RNN could be trained with RL to maximised the expected accuracy of the generated architectures on a validation set.

5.1.7 Section Conclusion

In this section, we showed that neural optimisers are able to perform beyond the common meta-learning assumption: to meta-test on the single identical input-domain used during meta-training. We introduced our novel context aware neural optimiser MTL2L with synaptic connections described with SVD and hypernetworks. First, we used SVD to construct a geometric space within the neural optimiser to describe all input features with the same eigenvectors which spanned the solution space. Then, we employed hypernetworks to encode eigenvalues to describe feature directions of the input data. Together, this combination allowed our MTL2L optimiser to self-modify its own RNN synaptic connections thereby creating new optimisation rules to update base learners on unseen datasets. In our experiments, we meta-trained our MTL2L on Fashion-MNIST and KMNIST, then meta-tested it on the unobserved synthesised Cifar10 dataset. Unlike the vanilla LSTM-optimiser which exhibited instability, we demonstrated that MTL2Ls successfully converged learners on an unseen input-domain, thereby making neural optimisers more generally applicable.

5.1.8 An Evaluation on Our Solved Problem

This section provided a solution to input-domain heterogeneity. We derived our solution based on our understandings on RNNs (see Chapter 4) and on a mathematical perspective of encoding gradient features on the loss landscape (see Section 2.2). The studies that we have done thus far have mostly been theoretical. However, there are more questions that require to be answered for making L2L practical for solving real life AutoML problems. Two important problems are listed below.

Problem 1: On the Efficiency of the L2L Meta-Learning Setup

There is a plethora amount of evidence that AutoML techniques can yield better learner models. However techniques such as L2L, Bayesian optimisation [Hutter et al., 2019], and neural architectural search [Elsken et al., 2019] all require intensive computational resources. Take L2L as an example, though our experiments (for instance Figure 4.3) showed that neural optimisers could update base learners more efficiently than their handcrafted rivals, neural optimisers themselves also required training. Thus, the total amount of time for training both the neural optimiser plus that for using the neural optimiser to update a base learner might end up being more time consuming than just using SGD to update a base learner.

Problem 2: On Training Large Base Learners for Long Duration

Thus far in this thesis, we followed the experimental setup of Andrychowicz et al. of which we documented in Section 3.2.5 to use neural optimiser to update small learners (2-layer MLPs) for short duration (meta-tested for 1000 iterations). However, state of the art neural networks such as ResNets usually consist more than 20 layers and typically require more than 200 epochs to train on large size datasets with small batch sizes. Hence, we would need to further test for the well-being of our neural optimisers on similar settings to validate their applicability in real life scenarios.

In our next study for this chapter, we will come up with an efficient L2L algorithm that addresses both the listed research problems above. Furthermore, the developed technique in the next study is aimed to achieve more than just rapid knowledge acquisition. We will show that it is possible to do network pruning with L2L thereby hastening an otherwise algorithmic reiterative and time consuming task.

*In order to make neural networks computationally simple,
it is good for them to let go of unnecessary connections via sparse learning.
However, this process is made computationally heavy because that conventional sparse
learning required iterative training.*

This section will show that sparse learning can be achieved with L2L.

*Our technique will assist a dense network to let go
of burdensome redundant connections without the need of heavy computation.*

Section Content:

- Section 5.2.1: Section Introduction,
- Section 5.2.2: On MetaSGD,
- Section 5.2.3: Extensions to MetaSGD,
- Section 5.2.4: Experiments and Results,
- Section 5.2.5: Section-wise Additional Related Work,
- Section 5.2.6: Section Conclusion:
Leveraging L2L for Efficient Network Pruning.

5.2 M^2 SGD: Learning to Learn Important Weights

In Chapter 4, we analysed RNNs in order to understand properties on shortcut connections (see Section 2.1). Then in the last section, we modified the canonical LSTM RNN to address input-domain heterogeneity in the meta-learning task of L2L. Though we proposed a better neural optimiser to substitute gradient descent, we also noted in Section 5.1.8 that this thesis had largely remained theoretical and that many practical issues in real life AutoML are not yet addressed. In order to demonstrate the utility of AutoML, the work of this section continues to modify L2L algorithms and we show that AutoML can automate the tedious experimental setup of the machine learning task of network pruning. The work of this section is based on our study in Kuo et al. [2020a].

5.2.1 Section Introduction

Since its introduction in 2016, research that surrounds L2L were mostly targeted towards addressing the undesirable features in the early LSTM-optimiser designed of Andrychowicz et al. (see Section 5.1.2). However as meta-learning through optimisation (see Section 5.1.6) slowly matures, we believe that we should start looking at L2L applications beyond *just* rapid knowledge acquisition. That is, we would like to integrate L2L into a machine learning discipline of which can benefit from meta-learning’s ability to achieve rapid knowledge acquisition. We selected network pruning, and this section presents empirical evidences on a learnt optimisation rule that can sparsely update learner parameters while removing redundant weights.

Current state-of-the-art neural networks possess a large amount of parameters making them both memory intensive and computationally demanding. To address these limitations, *sparse learning* can be used to only learn the important connections. Han et al. [2015] introduced an influential paradigm which removed redundant weights via a cyclic training composed of the three steps of (1) trains, (2) prunes weights, and (3) retrains a dense network. They were able to reduce the parameters of VGG-16 [Simonyan and Zisserman, 2015] by 13 times while showing no loss in accuracy. However, cyclic training is itself a computationally expensive operation – it requires a dense network to be repeatedly pruned and retrained from the beginning of the training phase.

Our work in this section is based on *MetaSGD* [Li et al., 2017] which learns a unique learning rate for each parameter of the learner network. We present *Masked MetaSGD* (M^2 SGD) as an extension to *MetaSGD*. We introduce a mask learned with LSTM [Hochreiter and Schmidhuber, 1997b] to remove redundant weights from the base learner and ensure that only the important weights are being updated. The masking system is a part of M^2 SGD’s learnt optimisation algorithm, a small size sub-network within a dense learner architecture is thus sought *on the fly*. This means that a dense network does not need to be retrained for weight removal thereby hastening the pruning procedure. Our experiments showed that M^2 SGDs not only updated learners quickly, but also removed 83.71% weights for ResNet20s [He et al., 2016a].

5.2.2 On MetaSGD

Li et al. [2017] proposed *MetaSGD* to update learner parameters via

$$\theta_{t+1} = \theta_t - \beta \odot \nabla_{\theta} \mathcal{L}, \quad (5.11)$$

where $\beta \in \mathbb{R}^{|\theta_t|}$ and that symbol \odot represents the element-wise product. Compared to gradient descent of Equation (1.4), learning rates β are the meta-learning target of MetaSGD and they become a constant vector after training. MetaSGD is a considerably simple neural optimiser when compared to its RNN variants like Ravi and Larochelle [2017], but it has achieved better results for few-shot learning.

5.2.2.1 The Main Differences to the LSTM-Optimiser

Recall that gradient descent (Equation (1.4)) performs

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}$$

to update parameters θ for a base learner \mathcal{F} while an LSTM-optimiser \mathcal{Q} substitutes (Equation (3.12)) this with

$$\begin{aligned} \theta_{t+1} &= \theta_t + g_{t+1} \text{ and} \\ [g_{t+1}, \mathbf{c}_{t+1}, \mathbf{h}_{t+1}] &= \mathcal{Q}(\mathcal{X}_t^{(\mathcal{Q})}, \mathbf{c}_t, \mathbf{h}_t | \eta). \end{aligned}$$

Hence, the LSTM-optimiser learns to prepare updates g_t which encapsulates both the scale and the direction of the movement of the learner parameters in the parametric space (see Figure 2.2.2). In contrast, the MetaSGD of Equation (5.11) explicitly utilise the native gradient $\nabla_{\theta} \mathcal{L}$. That is, while an LSTM-optimiser learns to precondition the gradients, MetaSGD learns a set of auxiliary components that is compatible with the original gradient descent paradigm. Figuratively, employing MetaSGD is equivalent to learn to control a set of different step sizes on different directions on the loss landscape. See Figure 2.4 for more intuitions.

5.2.2.2 Why Changing from LSTM-Optimiser to MetaSGD?

As stated in the introduction, this study is extended from MetaSGD instead of the LSTM-optimiser. The reason for this was to satisfy the listed desiderata of Section 5.1.8 to create an efficient AutoML algorithm. Below, we will provide elaborations on why MetaSGD is our technique of choice for achieving this.

From an Algorithmic Perspective

There are 4 entwined reasons why we prefer MetaSGD over the LSTM-optimiser.

First, that an LSTM-optimiser requires repetitive meta-training.

Second, that the learning rate of gradient descent is typically a mere constant \mathbb{R} .

Third, that MetaSGD's learnt learning rates are $\beta \in \mathbb{R}^{|\theta_t|}$; and

fourth, that MetaSGD's learnt learning rates are typically initialised as $0.01 \cdot \tilde{\mathbf{1}}$.

Let us start with point 2 and point 3 – it is not necessarily the best practice to update all parameters with one unified learning rate. Consider *inductive* transfer learning³ for domain adaptation in images, where we first train a randomly initialised deep network (such as ResNet [He et al., 2016a]) on one dataset, and then reuse parts of the pre-trained deep network on another dataset. It is common to keep the low level features in the shallow layers of the pre-trained network and then replace and re-train the deeper layers of the original network. Upon training on the new dataset, the pre-trained layers are updated with a learning rate that is typically 5 to 10 times lower than a separately defined learning rate allocated for updating the replaced layers. This enables the replaced layers to be updated quickly while ensuring that the kept layers only receive minimal amount of modification. A similar study was proposed for language modelling in finding suitable learning rates with *discriminative fine-tuning* [Howard and Ruder, 2018]. This idea is also viable for L2L, as it updates important parameters more rapidly to converge base learners faster.

Let us now compare point 1 and point 4 – on the immediate availability of L2L neural optimisers. LSTM-optimisers require iterations of training and cannot be used without thorough meta-training. As mentioned in Section 3.2, the meta-training phase of an LSTM-optimiser updates multiple copies of randomly initialised base learners (see line 4 of Algorithm 1). The meta-training phase is made lengthy because RNNs weights require careful calibration – as any incorrectness in an RNN will be re-injected into the network and cause instability (see Figure 5.3(b), and also Ehret et al. [2021]). During our experiments, we found that an LSTM-optimiser can only be employed if it were to be meta-trained with at least $\mathcal{V} = 20$ copies of base learner networks. In contrast, MetaSGD can be naïvely employed without any training. This is because that MetaSGD explicitly utilises the native gradient and only fine-tunes the learning rates. Thus as long as we initialise the MetaSGD learning rates as a vector of small values, the algorithm is directly applicable on any given base learner without any training required. (MetaSGD will still perform better if meta-trained.)

From a Mathematical Perspective

A base learner \mathcal{F} is considered to be updated successfully if its parameters θ_t converge to some θ^* after a lengthy horizon \mathcal{T} . That is, $\theta_{\mathcal{T}}$ should fulfil

$$\|\theta_{\mathcal{T}+\mathcal{P}}^{(j)} - \theta^{*(j)}\| < \epsilon^{(j)}, \forall \mathcal{P} \in \mathbb{Z}^+, \exists \epsilon^{(j)} > 0$$

for all entries $n = 1, \dots, |\theta^*|$ of and remain close to θ^* for any additional arbitrary horizon of \mathcal{P} . This description adheres to stability in the Lyapunov sense, and can be proven given if a *Lyapunov function* [Glendinning, 1994] exists. A Lyapunov function $\mathcal{V}(\theta_t)$ satisfies the following conditions such that

$$\mathcal{V}(\theta_t) > 0, \wedge \mathcal{V}'(\theta_t) < 0, \forall \theta_t \in \Omega \setminus \{\theta^*\} \quad (5.12)$$

³Refer to Section 3 of pages 4 to 6 in Pan and Yang [2009].

on a domain Ω about a convergent point θ^* . Simply put, \mathcal{V} is a function which maps the original parametric space to another space where positive growth in the transformed parametric space is bounded. Through a physical analogy, Qian [1999] used the energy principle to show that gradient descents are Lyapunov functions thereby guaranteeing the successful convergence of $\theta_{\mathcal{T}}$ to some θ^* (see Section 2.3.5.1). In brevity, they found⁴ that if an arbitrary parametric update complied to the following formulation of

$$\theta_{t+1} = \theta_t + \alpha_{t+1} \odot (\mathcal{A}_t \odot \theta_t' + \mathcal{B}_t \odot \theta_t'') \quad (5.13)$$

constrained with $\mathcal{A}_t, \mathcal{B}_t \in \mathbb{R}^{|\theta_t|}$, both $\mathcal{A}_{n,t}, \mathcal{B}_{n,t} \geq 0$ and that $\mathcal{A}_{n,t} + \mathcal{B}_{n,t} > 0$ for all dimensions and all entries of $\alpha_t \geq 0$, then the update is guaranteed to converge parameters of a network to a minimum. Given that MetaSGD of Equation (5.11) has $\mathcal{A}_{n,t} = 1$ and $\mathcal{B}_{n,t} = 0$, it is hence guaranteed to converge a base learner as long as all β values are non-negative.

5.2.3 Extensions to MetaSGD

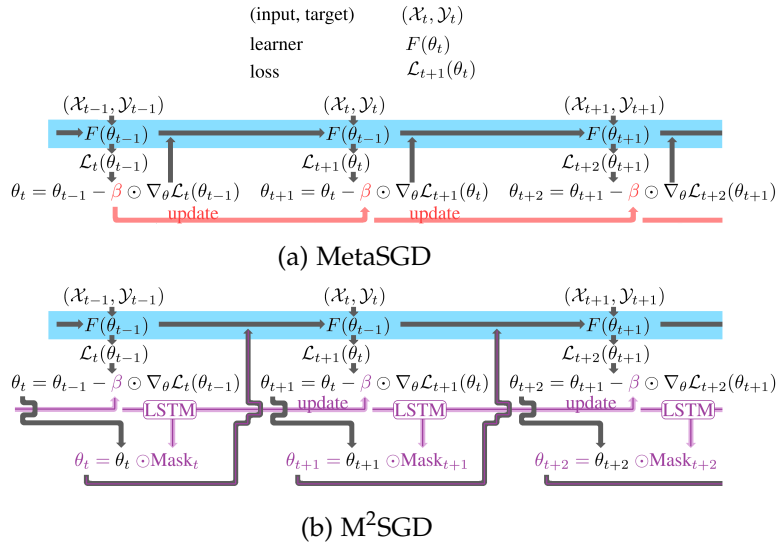


Figure 5.6: MetaSGD vs M²SGD

In this subsection, we introduce *Dynamic MetaSGD* (DMSGD) and *Masked MetaSGD* (M²SGD) as extensions to MetaSGD. We add two new features and depict the conceptual differences between MetaSGD and M²SGD in Figure 5.6. The modifications centre around the inferential process of the learning rates (see Equation (5.11)), and hence we highlight each algorithm’s β with different colours in the figure. First, we propose to learn the learning rates dynamically; and second, to learn a mask to remove redundant weights of the learner. Both of the features will be learnt via an LSTM RNN. Together, the two features will create an alternative to gradient descent to update a base learner that not only removes unimportant connections but also ensures that only the important weights are being updated.

⁴See their derivation in Section 3 on pages 2 and 3 of Qian [1999].

5.2.3.1 The Dynamic MetaSGD (DMSGD)

First, we present the DMSGD neural optimiser as follows

$$\theta_{t+1} = \theta_t + \beta_{t+1} \odot (-\nabla_{\theta} \mathcal{L}), \quad (5.14)$$

$$\beta_{t+1} = \text{ReLU}(\beta^{sta} + \gamma \odot \beta_{t+1}^{dyn}), \quad (5.15)$$

$$\beta_{t+1}^{dyn} = \tanh(U_{\beta} h_{t+1}), \text{ and} \quad (5.16)$$

$$[h_{t+1}, c_{t+1}] = \text{LSTM}(\nabla_{\theta} \mathcal{L}, h_t, c_t | \phi). \quad (5.17)$$

Similar to MetaSGD, the sole meta-learning target is the learning rates β_t . However unlike MetaSGD, the DMSGD learning rates are not constant values after meta-training. DMSGD is an optimisation algorithm which fine-tunes its own learning rates given the gradients of the learner parameters – in addition to the existing static learning rates of β^{sta} , DMSGD infers for an additional set of dynamic learning rates β_t^{dyn} with an LSTM. The balance between these 2 vectors of learning rates are controlled by a hyper-parametric vector $\gamma \in \mathbb{R}^{|\theta_t|}$ in Equation (5.15); and we default γ as a constant vector $\vec{1}$ unless stated otherwise.

There are two important notes for the DMSGD. First, like the LSTM-optimiser, DMSGD employs a coordinate-wise LSTM (Equation (5.17)) for inferring β_t^{dyn} . See Section 3.2.5.5 for the definition of coordinate-wise updates. Second, the DMSGD provides no update to the n th element of parameters θ_t when $|\beta_n^{sta}| < |\gamma_n \beta_{n,t}^{dyn}|$. Hence the DMSGD is trained to find and sparsely update the important weights of the learner while keeping the less important weights as constants.

5.2.3.2 The Masked MetaSGD (M^2 SGD)

The M^2 SGD neural optimiser is built on DMSGD and further extends MetaSGD with a masking system. *After* executing the DMSGD of Equation (5.14) to Equation (5.17), we execute the additional steps of

$$\theta_{t+1} = \theta_{t+1} \odot \text{Mask}_{t+1} \text{ where} \quad (5.18)$$

$$\text{Mask}_{n,t+1} = \max(0, \text{sgn}(\beta_{n,t+1} - \mathcal{U})) \quad (5.19)$$

for every entry n . We emphasise again that the mask is applied *after* parameters θ_t are updated; and its purpose is to nullify weights that are secondary in importance.

In this work, we consider less important weights as those that are updated with small learning rates. Our conjecture is based on two reasons. First, that base learner networks are encouraged to have their parameters randomly initialised with small magnitudes in order to prevent saturation in layer outputs [Glorot and Bengio, 2010] (see Section 2.1.2.5). Second, randomly initialised learners are non-optimal and have low accuracy. Hence following reasons 1 and 2, weights that are not updated with large learning rates are of low priority to the successful inferential process of a converged learner. Due to these reasons, we nullify the weights base on the magnitudes of their respective learning rates β_t as shown in Equation (5.19).

For M^2 SGD, we include hyper-parameter \mathcal{U} in Equation (5.19) to represent a non-negative lower bound for the learning rates. That is, it serves as a threshold for the minimal desired importance for those weights to be kept. We denote sgn as the sign function. When $\mathcal{U} = 0$, M^2 SGD nullifies weights that receive no update during training; and when $\mathcal{U} > 0$, we nullify both the non-updated weights and a portion of weights with less significant learning rates.

5.2.4 Experiments and Results

In this subsection, we repeat our experiments in Section 4.1.6.2 but employ MetaSGD, DMSGD, and M^2 SGD as our neural optimisers⁵. Recall from 5.1.5 that neural optimisers are judged based on the loss curves of their base learners. Loss curves with sharper slopes indicate that an optimiser is more efficient in converging a base learner; and loss curves with lower converged losses reflect that an optimisation has assisted a base learner to achieve better generalisation. In addition to these two attributes, we also aim to validate that M^2 SGD can update as efficiently as the vanilla MetaSGD while removing redundant weights for a base learner on the fly.

As shown in Sections 3.2.5 and 5.1.8, Andrychowicz et al. [2016]’s experimental setup was relatively simple. This was understandable because their sole purpose was to introduce the idea of learning to optimise. Nonetheless, they selected small sized 2-layer MLPs as base learners for the contextually simple MNIST dataset [LeCun et al., 1998], and that their meta-testing phase was considerably short with only 1000 iterations of updates. To ensure that L2L can be employed in more difficult situations, we proposed two scenarios. Our first scenario tested all neural optimisers as according to Andrychowicz et al.’s original setting but on the slightly more difficult FashionMNIST [Xiao et al., 2017] data. We then present a second scenario which tested L2L’s applicability in real world problems (see **Question 2** of Section 5.1.8) – which we updated deep ResNet20s [He et al., 2016a] as base learners on the low resource and contextually rich STL10 dataset [Coates et al., 2011] for multiple epochs until the large size base learners fully converged.

Base Learners and Datasets

The MLP base learners of scenario 1 were identical as those in Section 4.1.6.2 and had 25.4K parameters. In contrast, the ResNet20 base learners of scenario 2 followed the identical architecture as described in Shan et al. [2019] and had 270K weights. Readers can find descriptions on both dataset in Appendix A.

Experimental Setup

We meta-trained DMSGDs, M^2 SGDs, and MetaSGDs in a similar fashion to Algorithm 1 of Section 3.2.5.6. The readers should note that Lines 16 to 18 of the algorithm vary for each neural optimiser, and these lines should be substituted with the correct formulations: Equation (5.11) for MetaSGD, Equations (5.14) - (5.17) for DMSGD,

⁵Thus far in this thesis, all L2L experiments were conducted for RNN-based neural optimisers.

and the additional Equations (5.18) and (5.19) for M^2 SGDs. To minimise the neural optimiser loss \mathcal{L} in lines 12 and 13, we used Adam with learning rate 0.001.

For both scenarios, we set unroll $\Xi = 5$, learner-trials $\mathcal{V} = 1$, and learner-steps $\mathcal{S} = 100$. The combination of these hyper-parameters meant that, each neural optimiser was trained to update 1 learner for 100 steps; and that each neural optimiser memorised 5 continual steps of the depth of the optimisation trajectory.

After the neural optimisers were trained, we employed them to update the learners until they converged. Despite having the neural optimisers trained only to update the learners for $\mathcal{S} = 100$ steps, we updated the MLP base learners on MNIST for 5000 iterations and updated the ResNet20 base learners on STL10 for 4000 iterations. We would like to emphasise that for scenario 2, the specific number of 4000 steps was the least amount of iterations for MetaSGDs to converge the learners. That is, the setups of this section was purposely selected to favour the optimisation results for our rival MetaSGD neural optimiser.

We also compared our neural optimisers to the classic SGD. For both scenario, the classic SGD updated the base learners with learning rate 0.001 and with momentum 0.9. We also updated learners with the classic SGD until they converged. These results were included to compare base learner performances of those updated with L2L against those updated by the classic SGD.

Scenario 1: On Updating MLPs for FashionMNIST

After meta-training, the neural optimisers managed to converge the MLP learners in 10000 iterations; in contrast, it required the classic SGD 37500 iterations to converge the MLPs. The results are in Figures 5.7 and 5.8, and Table 5.1. Figure 5.7 and Figure 5.8 plot the first 5000 and first 1000 steps of the loss curves of the base learners respectively. All neural optimisers updated the MLP learners more rapidly than SGD; and our novel DMSGDs outperformed MetaSGDs in both efficiency and the base learners' final generalisation. In addition, M^2 SGDs outperformed MetaSGDs in terms of efficiency and removed a large amount of redundant network connections.

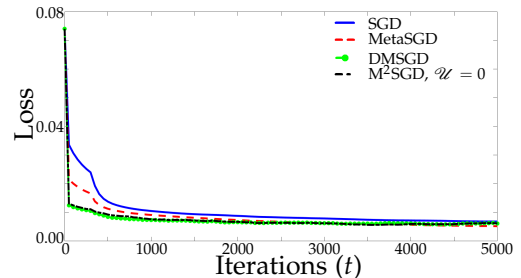


Figure 5.7: Loss curve (first 5000 steps) for FashionMNIST

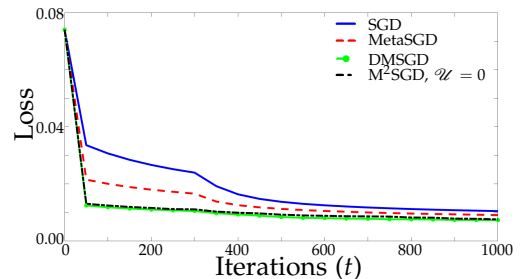


Figure 5.8: Loss curve (first 1000 steps) for FashionMNIST

Results

DMSGDs converged the learners with the highest accuracy (85.37%). MLPs that were updated by M^2 SGDs enjoyed a better level of accuracy than those that were updated

Table 5.1: Performance Details of MLP learners

Optimiser	% Accuracy	# Iterations	# Synapses (% ↓)
SGD	82.72 ± 0.55	37500	25.4K (-)
MetaSGD [Li et al., 2017]	85.06 ± 0.72	10000	25.4K (-)
DMSGD [Ours]	85.37 ± 0.55	10000	25.4K (-)
M ² SGD [Ours]			
$\mathcal{U} = 0$	84.80 ± 0.35	10000	13.0K (49.00%)
$\mathcal{U} = 0.1$	84.46 ± 0.43	10000	9.9K (60.68%)
$\mathcal{U} = 0.2$	83.81 ± 0.76	10000	7.8K (68.94%)

by the classic SGDs (82.72%) and their connections were removed dependent on the threshold of importance \mathcal{U} (see Section 5.2.3.2). When we set $\mathcal{U} = 0$, M²SGDs nullified nearly half (49.00%) of the unimportant MLP connections. While imposing the harsher constraint of $\mathcal{U} = 0.2$, we nullify up to 68.94% connections with a minimal amount of performance lesion (84.80% → 83.81%).

Remarks and Analyses

It was expected to see that both MetaSGD and our novel DMSGD algorithm outperformed the conventional SGD. Though DMSGD updated base learners to higher accuracy than MetaSGD, the differences were small. This showed that the heuristically learnt static learning rates of MetaSGD were sufficient for updating MLP base learners to their optimal accuracy for the FashionMNIST dataset. However and as shown in Figures 5.7 and 5.8, the additional dynamic learning rates (see Section 5.2.3.1) allowed DMSGD to update the MLPs more rapidly than MetaSGD. This meant that it was important, even for learnt optimisers, to fine-tune their own learning rates.

There were two remarks regarding M²SGD. First, the M²SGD loss curves were comparable to those of DMSGD. This showed that the difficulty in optimising a subnetwork within the MLP base learner was the same as that in optimising the entirety of the base learner. However and second, there was a slight performance lesion in only updating a subnetwork within the MLP base learners. Showing that though nearly 70% weights of the base learners were secondary in importance, there were still benefits in keeping network redundancies.

Scenario 2: On Updating ResNet20s for STL10

In comparison to scenario 1, scenario 2 tested the neural optimisers on a set of more difficult experimental settings to verify L2L’s applicability in real world scenarios as detailed in **Question 2** of Section 5.1.8. The objective was to replicate those successes seen in Figures 5.7 and 5.8; and we summarised the experimental results of this scenario in Figure 5.9, Figure 5.10, and Table 5.2 overleaf.

Additional Results

A similar set of observations can be drawn from scenarios 1 and 2. First and as shown in Figure 5.9, all neural optimisers updated learners more efficiently than SGDs. We found that M²SGDs performed comparably to MetaSGD, while DMSGDs outperformed all rival optimisers. DMSGDs updated the learners more rapidly, and converged learners to the least amount of loss.

Table 5.2: Performance Details of ResNet20 learners

Optimiser	% Accuracy	# Iterations	# Synapses (% ↓)
SGD	59.54 ± 0.23	4000	270K (-)
SGD	63.64 ± 0.73	7850	270K (-)
MetaSGD Li et al. [2017]	60.39 ± 2.27	4000	270K (-)
DMSGD [Ours]	65.52 ± 1.22	4000	270K (-)
M^2 SGD [Ours]			
$\mathcal{U} = 0$	65.10 ± 1.17	4000	140K (46.89%)
$\mathcal{U} = 0.1$	64.63 ± 1.50	4000	100K (61.31%)
$\mathcal{U} = 0.2825$	63.53 ± 0.38	4000	40K (83.71%)

Second and in Table 5.2, DMSGDs converged ResNets20s with the highest accuracy 65.52%. This was much higher than MetaSGDs’ converged accuracy of 60.39% and marginally higher than the classic SGDs’ 63.64%. In addition, DMSGDs were able to converge the learners much faster than SGDs and managed to converge ResNet20s in 4000 steps of update, while SGDs required 7850 steps. Note that with 4000 steps of update, SGDs were only able to yield ResNet20s with 59.54% accuracy. Thus we showed that the learning rates could be learnt to yield competitive learners via only updating the important weights.

Third and as recorded in Table 5.2, M^2 SGDs yielded competitive ResNet20s and also removed a significant portion of weights from the learners. M^2 SGDs with $\mathcal{U} = 0.2825$ removed 83.71% synaptic connections for ResNet20s and lowered the total amount of weights from 270K to 40K. In addition, M^2 SGDs converged these learners to 63.53% in 4000 steps, and matched those that were converged by SGDs to 63.64% in 7850 steps. This showed that for the more challenging setup in scenario 2, our proposed algorithm was still capable of detecting unnecessary learner connections via learnt learning rates. More importantly, these positive results showed that sparse sub-networks could be found without the need to retrain a dense network. Hence, M^2 SGD could serve as an efficient alternative to the work of Han et al. [2015].

Additional Remarks and Analyses

We found that lower losses in base learners did not necessarily mean higher accuracy. When Andrychowicz et al. [2016] introduced the influential concept of learning to optimise, their work focused on L2L’s ability to achieve rapid knowledge acquisition. That is, their main metric was the loss of the base learner networks and that

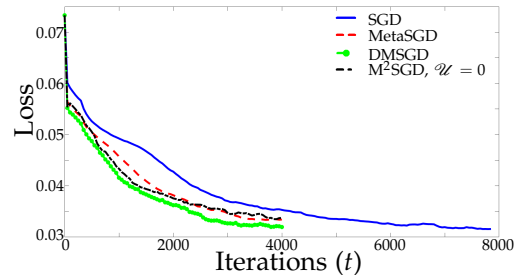


Figure 5.9: Loss Curve for Updating ResNet20s Until Convergence

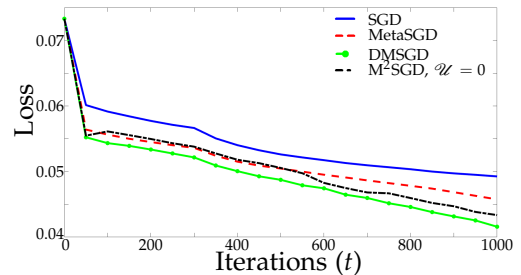


Figure 5.10: Loss Curve for Updating ResNet20s (first 1000 steps)

their second metric was the slope of the loss in the learning curve of the base learner networks. Their paper did not provide the accuracy of their base learners after being updated by neural optimisers. One potential reason for this could be related to their experimental setups (see Section 3.2.5) where they only updated base learners for 1000 steps. Our recorded performances in Table 5.2 reflected that though MetaSGDs updated ResNet20s more rapidly than the classic SGDs, the final accuracy of their learners were not as competitive (60.39%) as those that were updated with the classic SGDs (63.64%). Hence a neural optimiser should also be validated by their learners' accuracy when fully converged.

5.2.5 Section-wise Additional Related Work

Sparse learning

Sparse learning can be used to find small and sparse sub-networks within a dense architecture. Sparse networks enjoys multiple benefits from removing unimportant connections. For instance, the reduced size and computational cost makes it possible to build small but powerful networks on low-resource devices such as mobile phones.

In the early works, [LeCun et al. \[1990\]](#) proposed to prune weights via second-order derivatives, and [Ishikawa \[1996\]](#) modified the loss function to selectively prune weights with small magnitudes. A lot of these early work sees neural network as a natural extension of traditional statistical methods; and hence they took an information theoretic approach in simplifying complicated models. Their perspective and some conflicting experimental findings were previously reviewed in Section 2.2.4.

More recently, [Han et al. \[2015\]](#) showed that the number of parameters in a dense network could be reduce by an order of magnitude via the cyclic training of (1) train, (2) prune smallest weights, and (3) retrain. The cyclic training has since became the main paradigms in network pruning and was extended by the work of [Narang et al. \[2017\]](#) and [Carreira-Perpinán and Idelbayev \[2018\]](#). However, these approaches are computationally costly due to the need to repeatedly retrain dense networks.

The Lottery Ticket Hypothesis

A special variant of sparse learning is the idea of the *Lottery Ticket Hypothesis* proposed by [Frankle and Carbin \[2019\]](#). The lottery ticket refers to the smallest sub-network consisting all essential connections within a much larger dense network. [Frankle and Carbin](#) sought their candidate ticket via cyclic training. Their algorithm did not remove secondary weights but instead created fixed masks prior to the training phase and used the masks to hold a section of base learner parameters as constant values. The amount of weights to be masked were also decided prior of network training. They demonstrated that dense networks could be trained successfully even with up to 80% of weights unchanged during training. The authors also noted that cyclic training was computationally costly and required them to re-train a dense learner for multiple times over multiple trials.

Comparison

Our proposed DMSGD shares some similarities to the Lottery Ticket Hypothesis. Instead of removing the less important weights, DMSGDs learn to provide them with no updates (via the ReLU of Equation (5.15)) and thus indirectly keeping those less important weights as constants. Our proposed M^2 SGD is a trainable SGD-like optimiser which learns a mask (see Equation (5.19)) to nullify unnecessary connections in the base learner architecture. Hence after meta-training, M^2 SGD simultaneously updates the learner parameters and applies a mask to remove redundant weights for the learner on the fly. M^2 SGD therefore finds a lightweight sub-network within the dense architecture without the need of retraining thereby hastening the weight removal procedure.

We also found that M^2 SGD was actually more closely related to the early sparse learning works of LeCun et al. [1990] as both studies did not remove weights based on the weights' magnitudes. Instead, we influenced the weights through the magnitudes of their corresponding learning rates (via the sgn function of Equation (5.19)). This was because that we considered weights associated with small learning rates as parameters that were secondary in importance; as they were not required to be configured quickly for successful network inferences. The masking system of M^2 SGD was thus designed to nullify those weights with small learning rates.

5.2.6 Section Conclusion

This section introduced DMSGD and M^2 SGD as extensions to the optimisation-based MetaSGD meta-learner neural optimiser. Our proposed methods added two new features. First, we inferred dynamic learning rates with an LSTM RNN; and second, we added a masking system to nullify weights that were considered secondary in importance based on the magnitudes of the inferred learning rates. On the L2L task for classifying STL10 images, our novel algorithms updated ResNet20s base learners more rapidly than MetaSGDs, converged ResNet20s to lower losses, updated ResNet20s with better accuracy, and removed 83.71% weights from the dense ResNet20 architecture.

The positive results in this section provided us with more confidence such that meta-learning via optimisation could be employed in difficult real-life scenarios. As commented in Section 5.1.8, neural optimisers were previously limited to updating small size base learners to simple datasets without fully trained until parametric convergence. In this section, we showed that not only did our novel L2L techniques successfully update large size base learner networks on a contextually rich dataset, but they also pushed L2L capability beyond parameter configuration and accomplished non-iterative sparse learning.

5.3 Chapter Conclusion

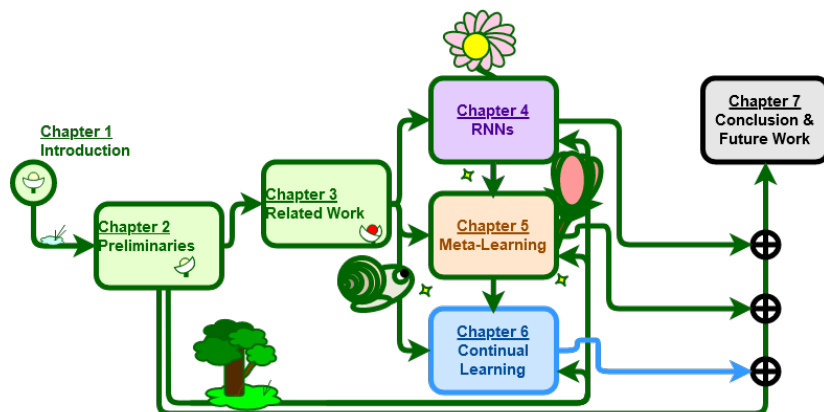


Figure 5.11: Progress Update

This chapter presented modifications to increase both the applicability and ability of the L2L meta-learning techniques. While Section 5.1 remedied a major undesirable feature in Andrychowicz et al. [2016]’s LSTM-optimiser, Section 5.2 proposed new algorithms based on Li et al. [2017]’s MetaSGD to a set of more difficult experimental settings that mimicked real world problems.

Specifically, our study in Section 5.1 looked at input-domain heterogeneity. We used LSTM insights from Chapter 4 to identify why RNN-based neural optimisers were unable to update a base learner over a different dataset in meta-testing than that seen during meta-training. Our solution was a simple hypernetwork [Ha et al., 2017] approach which addressed the geometry of the solution space for the RNN hidden variables for adjusting the dynamic properties that arise from RNN recursions. The inclusion of hypernetworks allowed our novel MTL2L-optimisers to re-implement their synaptic connections based on input data thereby alternating between different parametric configuration strategies.

Besides input-domain heterogeneity, we listed important issues that are required to be addressed for L2L to become applicable in real life in Section 5.1.8. As a remedy to the existing challenges in learning to optimise, we presented the M^2 SGD-optimiser based on MetaSGD [Li et al., 2017] in Section 5.2 which met all of our desiderata. M^2 SGD is an efficient algorithm capable of updating deep ResNet20s [He et al., 2016a] on the contextually rich STL10 [Coates et al., 2011] dataset until network convergence. It also leverages on L2L’s rapid knowledge acquisition to hasten network pruning thereby achieving non-iterative sparse learning.

Thus far in this thesis, we studied RNNs to devise more advanced L2L techniques. Our next chapter will discuss *Continual Learning* [Lopez-Paz and Ranzato, 2017], the AutoML technique for the sequential acquisition of skills. Sequential learning bears the potential to reach general artificial intelligence [Legg and Hutter, 2007] but is plagued by catastrophic forgetting (see Section 3.3). Our first study of the next chapter will combine continual learning techniques with neural optimisers; and we will show that this hybrid approach is more effective in preventing forgetting.

Topics on Simple and Efficient Continual Learning

This chapter introduces our contribution on continual learning research. The content therein this chapter continues the discussion in Section 3.3 in search for a better learning agent for achieving human-like sequential learning.

In standard practice (like the setup used in Chapters 4 and 5), it is assumed that the examples used to train a neural network are drawn independently and identically distributed (i.i.d.) from some fixed distribution. However in a changing environment, *continual learning* is required where an agent faces a continual stream of data, and must adapt to learn new predictions. This remains a difficult challenge because neural networks are known to suffer from *catastrophic forgetting* [McCloskey and Cohen, 1989], the phenomenon which they abruptly lose the capability to solve a previously learnt task when information for solving a new task is incorporated.

This thesis focuses on the memory-based techniques of Lopez-Paz and Ranzato [2017] and Chaudhry et al. [2019] to address sequential learning. Recall from Equation (3.16), a sequential learning backbone network \mathcal{F} is set to observe a continuum of data of n tasks $\omega_1, \dots, \omega_n \in \Omega$

$$\mathcal{K} = \{(x_1^1, \omega_1, y_1^1), \dots, (x_\gamma^i, \omega_i, y_\gamma^i), \dots, (x_\Gamma^n, \omega_n, y_\Gamma^n)\}$$

where each task contains Γ observations. Exact replay methods store a subset of memory \mathcal{M}_{ω_i} per past task and co-train the externally stored memory with data of a new task. Due to this reason, memory-based approaches suffer from high memory demand. One part of our contribution in Kuo et al. [2021b] leveraged on meta-learning (see Chapter 5) to create memory efficient algorithms.

This chapter also introduces our second contribution in Kuo et al. [2021a] which designs a simple and generally applicable patch compatible with every main continual learning archetypes. With our understanding on RNNs and insights on shortcut connections in Chapter 4, we present a simple architectural modification to sequential learning backbone classifiers that enhances the performance for all main archetypes including regularisation-based methods (see Section 3.3.5.1), dynamic architectural-based methods (see Section 3.3.5.2), and memory-based methods (see Section 3.3.5.3).

*The more that I research, the more that I realise the importance of redundancy.
Instead of declaring redundancy as secondary,
I have come to appreciate how redundancy in a model enables
representation, manipulation, and learnability (plasticity)
based on existing knowledge (see Section 3.3.3).
Over the next section,
we will explore the continual adaptation of modelling ability at the synaptic level
and inject meta-learning to fine-tune a set of continual learning network weights.*

Section Content:

- Section 6.1.1: Section Introduction,
- Section 6.1.2: Replay with Episodic Memory,
- Section 6.1.3: MetaSGD for Continual Learning,
- Section 6.1.4: Implementation Details on MetaSGD-CL,
- Section 6.1.5: Experiments and Results,
- Section 6.1.6: Section-wise Additional Related Work,
- Section 6.1.7: Section Conclusion:
Leveraging L2L for Efficient & Robust Continual Learning,
- Section 6.1.8: Re-imagining Continual Learning.

6.1 Learning to Continually Learn Rapidly from Few and Noisy Data

As discussed in Section 3.3, continual learning sequentially configures a network while aiming to maintain existing knowledge etched in the model parameters. *But, do all weights hold equivalent importance to all of the learnt tasks?* Perhaps Neuron 1 is indispensable for Task 3 but mayhaps it is irrelevant for Task 5. Since it is tedious for practitioners to fine-tune and judge the importance of individual neurons, we propose to learn an auxiliary network to facilitate the sequentially learning of a stream of tasks for a task-specific backbone base learner. This section is based on our work in Kuo et al. [2021b] which leveraged meta-learning for continual learning. The main contribution was to mitigate the high memory demand in replay-based continual learning techniques with L2L’s ability of rapid knowledge acquisition.

6.1.1 Section Introduction

Neural networks are powerful function approximators but suffer from catastrophic forgetting and are unable to sequentially learn new tasks without guaranteed stationarity in the distribution of a continuum of data (see Equation (3.16)). This is because the learnt parameters of a connectionist model are required to be continuously modified thus abruptly disrupting the model’s learnt mappings for solving past tasks.

Figure 6.1 illustrates the process of catastrophic forgetting using the concepts of loss landscape and loss minima (see Section 2.2). A backbone learner optimised for Task U is now sequentially updated for an independent Task V. Naïvely training the learner with gradient descent will cause the learner’s parametric configuration to leave the minimum (low loss region) on the loss landscape of Task U in pursuit for that of the new task [Kirkpatrick et al., 2017]¹.

Previously, it was shown that continual learning could be achieved via replay – by concurrently training externally stored old data while learning a new task (see Section 3.3.5.3). These *replay* [Robins, 1995] and *pseudo replay* [Shin et al., 2017] mechanisms have continuously been improved [Rebuffi et al., 2017; Lopez-Paz and Ranzato, 2017; Chaudhry et al., 2018]; and a recent paper from Chaudhry et al. [2019] have even shown that forget prevention could be made possible by co-training as little as one single sample of past data per parameteric update. Nonetheless, we found

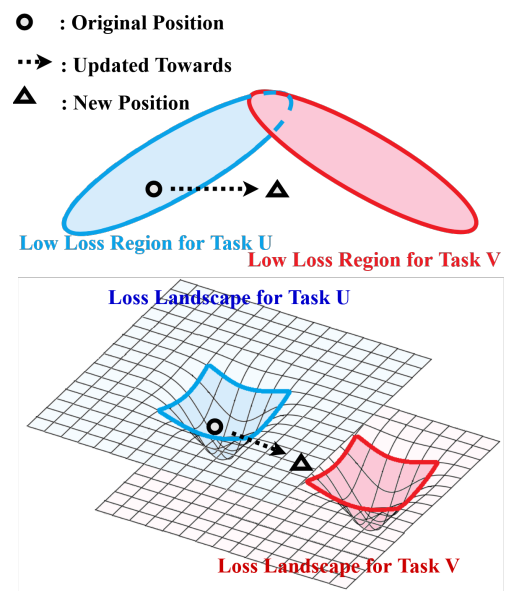


Figure 6.1: Exiting a Local Minimum while Continual Learning [Kirkpatrick et al., 2017]¹.

¹Inspiration taken from Figure 1 of Kirkpatrick et al. [2017].

that the replayed data was still required to be sampled from a large set of externally stored memory. Sampling from an insufficient amount of external memory would lead a network to overfit on early tasks thus generalise poorly on subsequent tasks.

To address memory-based continual learning under a small memory setting, we supplement replay mechanics with meta-learning (see Chapter 5) for rapid knowledge acquisition. Our research in this section combines the work of Chaudhry et al. [2019] with MetaSGD [Li et al., 2017] (also see Section 5.2) to maintain high performances under the constraint of learning with low memory consumption. We employed a meta-learner to assist the training procedure of a task-specific backbone network of which *learns a learning rate per parameter per past task*. We found that this hybrid approach not only produced strong results when less memory was available, but it also inherited several desirable meta-learning advantages for continual learning. Our approach demonstrated strong robustness against the realistic scenario which we continually learn under the presence of noises, and it could also optimise base learner parameters to achieve higher accuracy in less iterations.

6.1.2 Replay with Episodic Memory

Algorithm 2 ER Recipe

```

1:  $\mathcal{M}_{\omega_i} \leftarrow \{\}$   $\forall i$  such that  $\omega_i \in \Omega$  ▷ Define memory storage
2: for  $i = 1, \dots, |\Omega|$  do
3:   for  $(x_{\gamma}^i, \omega_i, y_{\gamma}^i) \in \mathcal{K}$  do
4:      $\mathcal{M}_{\omega_i} \leftarrow \mathcal{M}_{\omega_i} \cup (x_{\gamma}^i, y_{\gamma}^i)$  ▷ Update memory
5:     if  $i > 2$  then
6:        $B_{\mathcal{M}} \sim \cup_{u=1}^{i-1} \mathcal{M}_{\omega_u}$  ▷ Sample from all memory
7:        $\theta \leftarrow \text{SGD}((x_{\gamma}^i, y_{\gamma}^i) \cup B_{\mathcal{M}}, \theta)$  ▷ Update with sampled memory
8:     else
9:        $\theta \leftarrow \text{SGD}((x_{\gamma}^1, y_{\gamma}^1), \theta)$  ▷ Conventional update
10:    end if
11:  end for
12: end for

```

As discussed in Section 3.3.5.3, naïve replay causes high memory consumption and high computational cost. Lopez-Paz and Ranzato [2017] showed that replay could be made less memory demanding via *episodic memory* where not all data of past tasks are required. There were two steps to the episodic memory mechanism. Prior to the sequential training phase, Lopez-Paz and Ranzato first defined the size of the memory storage \mathcal{M}_{ω_i} allocated for each task. Then upon training a new task, all stored data from the external memory were used to precondition the gradients of the current task. Though this method yielded very strong results and only stored a section of all past data, all externally stored memory were required to be replayed at every time step and hence the algorithm was computationally costly.

In a recent and important work, Chaudhry et al. [2019] streamlined Lopez-Paz and Ranzato’s method with *Experience Replay* (ER) and showed that catastrophic for-

getting could be alleviated without utilising all data stored in the episodic memories. Algorithm 2 shows their ER approach; and the core mechanism of ER was to *sample* mini-batches $B_{\mathcal{M}}$ from all stored data of past tasks (see line 4) to effectively mitigate the computational burden from employing episodic memories. The sampled content was then concurrently trained with the existing batch of the current learning task and they used SGD to update parameters θ of backbone learner \mathcal{F} (see line 7). Impressively, they found that strong results for forget prevention can be achieved even when $B_{\mathcal{M}}$ only sampled 1 data from all memory $\cup \mathcal{M}_{\omega_i}$.

6.1.2.1 Episodic Memory Sampled from a Tiny Memory Storage

To employ episodic memory, the practitioner would need to explicitly define the size of the external memory storage;

But what are the disadvantages of selecting small sizes?

A small memory storage meant that the replay mini-batch $B_{\mathcal{M}}$ would have a higher chance of re-sampling the same past data for concurrent training (see line 7 of Algorithm 2). A base learner would hence overfit on earlier tasks through the consecutive re-exposure of a small set of past data and thus underfit on future tasks. This issue was lightly investigated in Chaudhry et al. [2019]², but they only compared medium-size to large-size memories³. In this current work, we limit the memory size to only a handful of data per past task to study its consequences for continual learning and propose meta-learning as a remedy.

6.1.3 MetaSGD for Continual Learning

In this study, we formulate continual learning with tiny episodic memory storage as a low resource problem. Though we do not lack data on subsequent tasks, the backbone network is still required to learn from a small memory unit for past tasks.

Our meta-learning technique of choice is the previously introduced optimisation-based approach of MetaSGD [Li et al., 2017]. Recall from Equation (5.11) that MetaSGD replaces gradient descent of Equation (1.4)

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}$$

with

$$\theta_{t+1} = \theta_t - \beta \odot \nabla_{\theta} \mathcal{L}_{t+1}(\theta_t)$$

to learn a learning rate per parameter to hasten the update procedure of a base learner network. And as we previously demonstrated in Section 5.2.4, our novel meta-learners extended from MetaSGD could rapidly update base learners in 10000

²See Section 5.5 and Appendix C of Chaudhry et al. [2019].

³See Table 7 in Chaudhry et al. [2019], where they compared 1000 to 20000 memory units. What happens when we use 100?

steps to similar accuracy of which otherwise required 37500 iterations of conventional SGD updates.

6.1.3.1 An Analysis on the Update of Experience Replay for Continual Learning

Under the conventional experimental setup, loss \mathcal{L} aggregates every individual cost \mathcal{C} therein a mini-batch B_n with

$$\mathcal{L} = \frac{1}{|B_n|} \sum_{i=1}^{|B_n|} \mathcal{C}(\hat{y}_i, y_i) \quad (6.1)$$

where \hat{y}_i is the prediction produced by learner \mathcal{F} and that y_i is the corresponding true label. Hence under the ER regime, every batch used in continual learning will be substituted with $B_n \leftarrow B^*$ where $B^* = B_n \cup B_{\mathcal{M}}$ (see line 7 of Algorithm 2). Then due to the linearity of summation, the new loss \mathcal{L}^* can be written as

$$\mathcal{L}^* = \frac{1}{|B^*|} \left[\sum_{i=1}^{|B_n|} \mathcal{C}(\hat{y}_i, y_i) + \sum_{j=1}^{|B_{\mathcal{M}}|} \mathcal{C}(\hat{y}_j, y_j) \right]. \quad (6.2)$$

If a base learner had observed V tasks, we can highlight task specificity and rewrite Equation (6.2) as

$$\mathcal{L}^* = \frac{1}{|B^*|} \left[\sum_{i=1}^{|B_n|} \mathcal{C}(\hat{y}_i, y_i) + \sum_{u=1}^{V-1} \sum_{j=1}^{|B_{\mathcal{M}_u}|} \mathcal{C}(\hat{y}_j, y_j) \right] \quad (6.3)$$

where $B_{\mathcal{M}} = \cup_{u=1}^{V-1} B_{\mathcal{M}_u}$. Then with \mathbb{L} as the task-wise loss we have

$$\mathcal{L}^* = \frac{1}{|B^*|} \left[|B_n| \mathbb{L}_V + \sum_{u=1}^{V-1} |B_{\mathcal{M}_u}| \mathbb{L}_u \right]; \quad (6.4)$$

and this will create the standard parametric update of

$$\theta_{t+1} = \theta_t - \left(\frac{\alpha}{|B^*|} |B_n| \right) \nabla \mathbb{L}_V - \sum_{u=1}^{V-1} \left(\frac{\alpha}{|B^*|} |B_{\mathcal{M}_u}| \right) \nabla \mathbb{L}_u \quad (6.5)$$

of which we could employ Meta-SGD and introduce β as

$$\theta_{t+1} = \theta_t - \beta_V \odot \nabla \mathbb{L}_V - \sum_{u=1}^{V-1} \beta_u \odot \nabla \mathbb{L}_u \quad (6.6)$$

to learn a learning rate per parameter per past task.

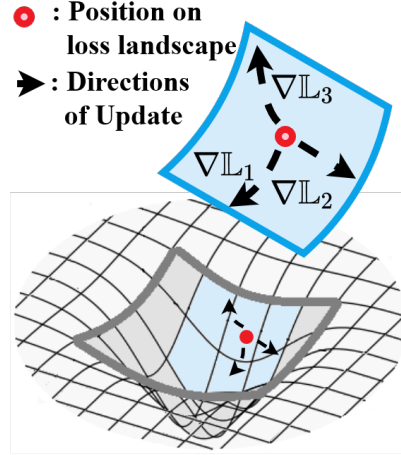


Figure 6.2: Task-wise Directional Updates

To elaborate, consider the scenario in Figure 6.2 where a learner observes its third task. The learner faces different directions of updates $\nabla\mathbb{L}$ from the competing learning objectives of different tasks. *Interference* [Riemer et al., 2018] hence occurs and that transfer in knowledge becomes difficult; this is further complicated as learner \mathcal{F} could easily surpass thousands of parameters in practice [Szegedy et al., 2015]. Thus instead of applying handcrafted rules such as SGD to update a base learner, we employ **MetaSGD for Continual Learning** (MetaSGD-CL) to explore a stream of changes in a series of task structures.

6.1.4 Implementation Details on MetaSGD-CL

Upon learning a new task V (including task 1), MetaSGD-CL initialises a trainable vector $b_V = 0.01 \cdot \tilde{\mathbf{1}} \in \mathbb{R}^{|\theta|}$. In the context of Equation (6.6), we set

$$\beta_{V_k} = \text{ReLU}(\min(b_{V_k}, \kappa)); \quad (6.7)$$

where κ is a positive constant hyper-parameter. Thus every individual learning rate β_{V_k} lies within the range of $[0, \kappa]$. We default κ as 0.02 unless specified otherwise. Though $\kappa = 0.02$ may appear to be an arbitrary number, it isn't. Instead, it is twice as much as the initialised value of 0.01. This will be emphasised later in the experimental setups; and we will highlight some properties of this hyper-parameter later in an ablation study.

When sequential learning a new task V (excluding task 1), MetaSGD-CL freezes its past learning rates β_u for all $u < V$ and updates the backbone network \mathcal{F} as according to

$$\theta_{t+1} = \theta_t - \beta_V \odot \nabla\mathbb{L}_V - \frac{1}{V-1} \sum_{u=1}^{V-1} \beta_u \odot \nabla\mathbb{L}_u. \quad (6.8)$$

Our readers may find that there are slight differences between Equation (6.6) and the formulation presented above. This is because that our design for Equation (6.8) also accounts for potential alignments in $\nabla\mathbb{L}$, of which can occur due to similarity among arbitrary task structures. The term $\frac{1}{V-1}$ is thus added to prevent the over-parametrisation in any weight dimension.

In order to train the MetaSGD-CL neural optimiser, we modified upon the original L2L meta-learner loss function from Andrychowicz et al. [2016] of which we described in Section 3.2.5.3. The modified meta-learner loss of MetaSGD-CL is

$$\mathcal{L}(\beta_V) = \mathcal{L}(\mathcal{F}_{\theta_{t+1}}(x_t), y_t | \beta_V, \beta_u \text{ for } u = 1, \dots, (V-1)). \quad (6.9)$$

Our readers should note that as described in Algorithm 1, the meta-learner loss is based on an updated learner $\mathcal{F}_{\theta_{t+1}}$ with the original data x_t . Thus, the objective for MetaSGD-CL is to further fine-tune the backbone base learner \mathcal{F} (see the foreword of Chapter 5). In addition, the meta-learning targets of β_s are updated with Adam with learning rate 0.01 (also see the setup in Section 5.2.4).

6.1.5 Experiments and Results

This section will include some benchmark sequential learning experiments. Our focus is to demonstrate the unique benefits of MetaSGD-CL over ER [Chaudhry et al., 2019]. Recall that ER is a replay-based technique (see Section 3.3.6) that has achieved state-of-the-art results; and that memory-based techniques are generally regarded as the most stable and most effective methods (also see Van de Ven and Tolias [2018]). In addition, this section will also include an ablation study to clarify some inner workings of our novel MetaSGD-CL.

Dataset

We conducted experiments on *permuted MNIST* [Kirkpatrick et al., 2017] and *incremental Cifar100* [Lopez-Paz and Ranzato, 2017]. Previously, Lopez-Paz and Ranzato [2017] and Chaudhry et al. [2019] showed that memory-based approaches outperformed continual learning techniques of other archetypes and achieved state-of-the-art status. However, we will show that when the amount of memory decreases, ER performances can drop sharply on the considerably simple permuted MNIST tasks. Refer to Appendix A for descriptions on the MNIST and the Cifar100 datasets.

Of the two experiments that we consider, permuted MNIST was derived from MNIST by first reshaping each image as a vector of 784(=28×28) pixels, then applied a unique permutation on all vectors of that task. On the other hand, the incremental Cifar100 dataset was generated via sequentially introducing 5 classes per task drawn without replacement from the dataset. For permuted MNIST, we sequentially trained backbone networks over 10 permuted MNIST tasks. While for incremental Cifar100, we reserved the first 50 classes from Cifar100 to pre-train a classifier, and then tested the said classifier on the remaining 50 classes⁴. There were hence 10 tasks in total for both experimental setups. More details will be discussed in the next subsection.

Backbone Learners

For permuted MNIST, we followed Lopez-Paz and Ranzato and trained MLPs as backbone learner. As for incremental Cifar100, we followed Mai et al. [2020]⁵ and used their best practice – continual learning with network compartmentalisation – which employed a fixed feature extractor and only sequentially trained the classifier networks. Our backbone network compartmentalisation included a ResNet18 [He et al., 2016a] feature extractor and a MLP classifier network. We first thoroughly pre-trained the ResNet18 on the first 50 classes of Cifar100. Then upon the sequential learning phase, we froze all weights of the ResNet18 and substituted the original softmax layer of the ResNet18 with an MLP and only continually updated the said MLP as a new classifier. The MLPs of both our experimental setups consisted 2 hidden layers with 100 dimensions followed by a softmax layer.

⁴By reserving the first 50 classes for thorough pre-training, we created a highly stable backbone learner with high accuracy on identifying objects of the said 50 classes. This allowed us to analyse for the full extent of the true severity of catastrophic forgetting on a well-trained network.

⁵The method addressed in Mai et al. [2020] won the CVPR CLVision 2020 challenge.

Following Lopez-Paz and Ranzato, we used the batch size of 10. Each task of permuted MNIST observed 100 batches of data; and each incremental Cifar100 task observed 1 epoch of data. Our MetaSGD-CL was described in Section 6.1.4; and all of our continual learning baselines were updated with SGD with learning rate 0.01.

Metrics

We used 2 metrics to evaluate the continual learning techniques. Both metrics were based on a matrix $\mathcal{R} \in \mathbb{R}^{n \times n}$ of which we employed to store all accuracy during the entire course of continual learning. Row i recorded the accuracy after training on the i th dataset; and column m recorded the accuracy of the m th task.

The severity of catastrophic forgetting was reflected via

Final Accuracy of Task 1 (FA1) $\mathcal{R}_{n,1}$; while

Average accuracy (ACC) $\frac{1}{n} \sum_{i=1}^n \mathcal{R}_{n,i}$

showed the overall performances and the severity of the overfitting from the backbone networks. For both metrics, the larger the better.

Continual Learning Techniques

Our benchmark experiments included six techniques. *Singular* was the scenario where we naïvely and sequentially trained a backbone network. Then, we tested our novel *MetaSGD-CL* against the baseline *ER* [Chaudhry et al., 2019]. In addition, we considered one method each from the three main continual learning archetypes (see Section 3.3.5). We chose the memory-based *Gradient Episodic Memory* (GEM) of Lopez-Paz and Ranzato [2017], the regularisation-based *Elastic Weight Consolidation* (EWC) of Kirkpatrick et al. [2017], and the dynamic architectural-based *Hard Attention to the Task* (HAT) of Serra et al. [2018].

Results on Preventing Catastrophic Forgetting

Our results on permuted MNIST are in Figure 6.3 and Table 6.1; and those for incremental Cifar100 are in Figure 6.4 and Table 6.2. The tables list the mean values of each metric for each method. From the FA1 metric, we validated Chaudhry et al. [2019] and Lopez-Paz and Ranzato [2017]’s finding that memory-based approaches

were more effective in preventing forgetting. The change of the metrics throughout the entire course of sequential learning are denoted as the solid lines in the figures; and they are accompanied by shades to reflect their variability.

Our memory-based approaches used *hard storage* [Lopez-Paz and Ranzato, 2017] and assigned 250 memory to each tasks. For each update, GEM inefficiently used all stored data to prevent forgetting; while both MetaSGD-CL and ER randomly sampled 10 data from each task-wise storage to append batch B_n with $B_{\mathcal{M}}$ (see Equation (6.2)). The old samples saved in the hard storage did not change throughout training.

Table 6.1: Metrics for Permuted MNIST

Methods	FA1 (%)	ACC (%)
MetaSGD-CL (Ours)	81.02	82.19
ER	80.60	69.42
GEM	80.71	79.43
EWC	64.80	71.84
HAT	74.03	76.17
Singular	62.18	71.34

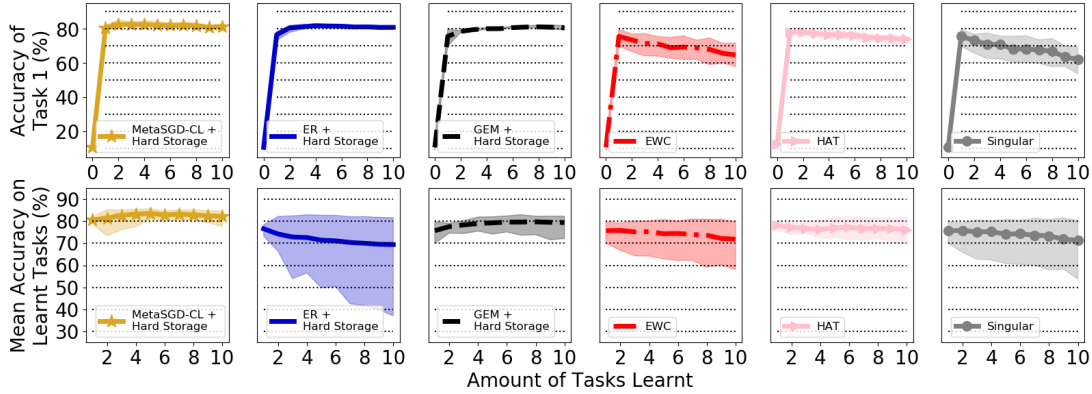


Figure 6.3: Performances for Permuted MNIST

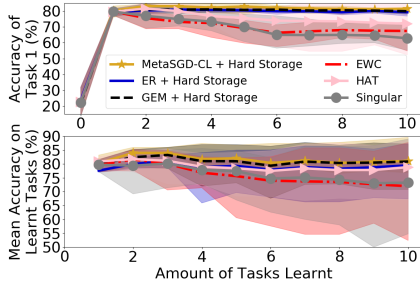


Figure 6.4: Performances for Incremental Cifar100

Table 6.2: Metrics for Incremental Cifar100

Methods	FA1 (%)	ACC (%)
MetaSGD-CL (Ours)	81.52	81.06
ER	79.28	79.50
GEM	79.92	80.88
EWC	67.46	71.98
HAT	71.98	78.86
Singular	62.86	73.27

Though all MLPs were meant to have 100 hidden dimensions, this was increased to 400 for HAT. By design, HAT could not prevent forgetting with low amount of synapses in the backbone. Similar findings were also reported in [Ahn et al. \[2019\]](#).

Remarks and Analyses

From Figure 6.3 and Table 6.1, we observed that ER was comparable to MetaSGD-CL in FA1 but much lower in ACC. The comparable performances in FA1 was expected as both techniques had direct access towards the exact replay memory. However, while ER maintained high accuracy on the first task (FA1), the overall accuracy of all tasks (ACC) continued to drop throughout the course of sequential learning. ER hence suffered from severe overfitting.

In contrast, our novel MetaSGD-CL maintained high FA1 and ACC scores at all times. The reason could be attributed to the design of the objective function of the meta-learner (see Equation (6.9)) where

$$\mathcal{L}(\beta_V) = \mathcal{L}(\mathcal{F}_{\theta_{t+1}}(x_t), y_t | \beta_V, \beta_u \text{ for } u = 1, \dots, (V-1)).$$

That is, the learning of the new learning rates β_V were based on past learning rates β_u and hence MetaSGD-CL provided a Bayesian inferential-like update to the backbone base learner via perceiving sequential changes in the task structure. This in turn achieved a good balance between solving both past tasks and the current task.

For incremental Cifar100 and Table 6.2, Mai et al. [2020]’s experimental setup limited some extent of overfitting by only continually training the classifier network. However, the bottom subplot of Figure 6.4 showed that ER’s variability (in blue) was still much larger than MetaSGD-CL’s (in yellow); and hence overfitting persisted.

3 Advantages of Meta-SGD over ER

Here we analyse the advantages of integrating meta-learning in continual learning.

More Remarks and Analyses: MetaSGD-CL Does Not Overfit on Old Tasks

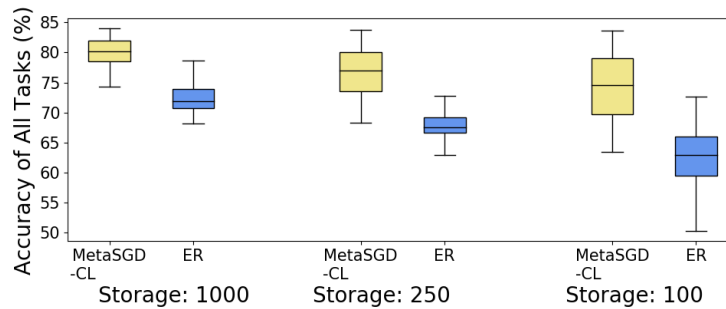


Figure 6.5: Ring Buffer of Different Sizes for Permuted MNIST

We further tested overfitting in replay-based continual learning with the alternative memory storage of *ring buffer* [Chaudhry et al., 2019]. Ring buffer shared Φ memory units *for all tasks*⁶. By design, the buffer was under-utilised in the early phases of sequential learning; and that its stored samples on observed tasks did not change. Compared to hard storage, the shared storage of ring buffer made sequential learning much more computationally efficient. However, it was also more likely to incur overfitting because data of all past tasks shared the same buffer thus increasing the chance of re-exposure to past data (see Section 6.1.2.1).

Figure 6.5 shows the results for MetaSGD-CL and ER on permuted MNIST with ring buffers of the sizes of 1000, 250, and 100 units. Note, these settings were much smaller than the scenarios tested in Chaudhry et al.³. Again, both techniques randomly sampled 10 past data as $B_{\mathcal{M}}$ for each iteration (see line 7 of Algorithm 2).

Unsurprisingly, the performances of either methods decreased as the sizes of their memories dropped. This was a natural consequence of overfitting by repeatedly re-sampling from a small set of past data. ER is one of the state-of-the-art techniques; but still, our experiments found that a considerably large amount of lesion could be incurred on the considerably simple permuted MNIST benchmark when the size of the ring buffers changed from 1000 to 100. In contrast, our MetaSGD-CL was more capable of maintaining higher accuracy. This showed that meta-learning’s ability to achieve rapid knowledge acquisition could be introduced to remedy the large memory cost of replay-based continual learning techniques. Our novel MetaSGD-CL was computationally efficient and only required a small external memory storage; MetaSGD-CL hence solved two major issues that plagued previous replay techniques.

⁶If $\Phi = 250$ and that $n = 10$ tasks, then it is 25 units per task.

More Remarks and Analyses: Rapid Knowledge Acquisition

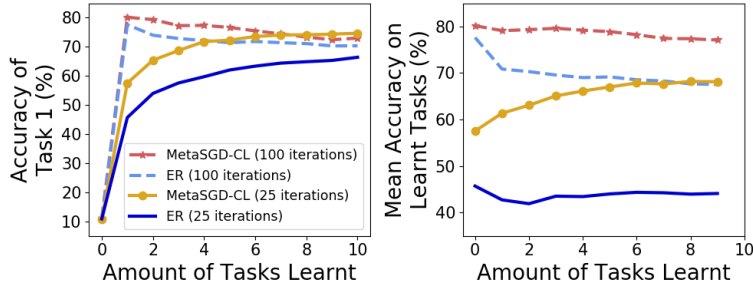


Figure 6.6: Permuted MNIST but with Reduced Iteration

We also tested 2 realistic continual learning scenarios. First, we considered when there was a limited amount of source data; and second, when data was collected from a noisy environment.

In Figure 6.6, we tested MetaSGD-CL and ER on permuted MNIST with ring buffer but this time with a reduced amount of training iterations. Unsurprisingly, we found that the test accuracy dropped when the available training data was decreased from 100 iterations to 25 iterations per task.

We continued to observe overfitting in ER. Furthermore, ER was not able to achieve high accuracy in Task 1 as early as MetaSGD-CL. This was expected, because that ER sequentially updated the backbone base learner with the conventional SGD (see line 7 of Algorithm 2). On the other hand, our novel MetaSGD-CL approach was able to leverage on meta-learning to yielded higher accuracy in few iterations (also see Chapter 5 and Section 5.2.4). More importantly, both the FA1 and ACC scores concurrently increased over all tasks for MetaSGD-CL – a sign showing no forgetting and no overfitting.

More Remarks and Analyses: Strong Robustness to Noise

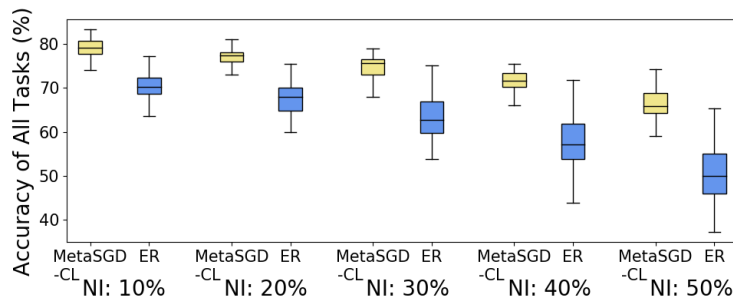


Figure 6.7: Permuted MNIST but with Noise Injection

In Figure 6.7, we tested MetaSGD-CL and ER on permuted MNIST with ring buffer with additional *Noise Injection* (NI) in both the training data and the externally stored memory. That is, after applying the task-wise fixed permutation mask on top of the original MNIST images, NI randomly shuffled a section of pixels of the reorganised MNIST images. Our NI setting ranged from injecting a low amount of 10% noise to

the high amount of injecting 50% noise.

We found that MLPs trained with MetaSGD-CL were less susceptible to NI than their counterparts trained with ER. The side-by-side box plots in Figure 6.7 showed that MetaSGD-CL yielded MLPs with higher accuracy and smaller variability. This was likely because that the meta-learned learning rates β of Equation (6.8) explored which backbone parameters were less prone to NI; and thus assigned higher learning rates to configure them quickly.

Behaviours on Learning Rate β

Here we study the progression of learning rate updates with the permuted MNIST task.

More Remarks and Analyses: Convergence in β values

Table 6.3: Ratios of Extreme Task-wise β s with $\kappa = 0.1$

Proportion of $\beta > 0.05$ (in %) in task-wise learning rate					
Tasks	2	4	6	8	10
Layer 1	1.39	0.63	0.47	0.46	0.36
Layer 2	2.26	1.41	1.33	1.32	1.30
Softmax Layer	3.64	2.28	1.88	1.76	1.74
Proportion of $\beta < 0.02$ (in %) in task-wise learning rate					
Tasks	2	4	6	8	10
Layer 1	1.43	3.90	6.56	8.31	8.79
Layer 2	4.29	6.59	7.71	10.51	11.27
Softmax Layer	6.02	7.50	8.38	8.92	11.58

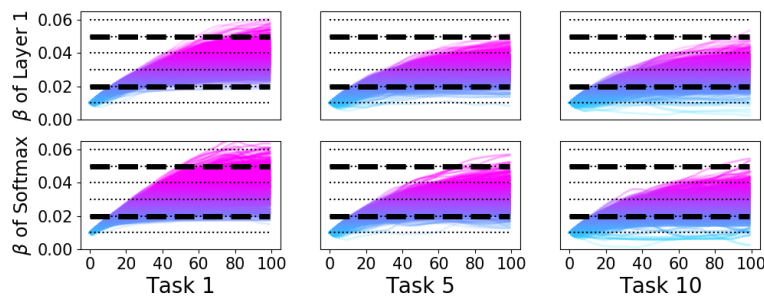


Figure 6.8: Transition in Task-wise β Magnitudes

As we previously mentioned, the maximal value of β was defaulted as $\kappa = 0.02$ (see Equation (6.7)). The value for κ was purposely set low because we hypothesised that there existed conflicted objectives between meta-learning and continual learning.

As shown in Figure 6.6, meta-learning is naturally capable of rapid knowledge acquisition. That is, we hypothesised that the meta-learned learning rates β would grow in an unbounded fashion. To elaborate, given large κ hyper-parametric values, we presumed that the MetaSGD-CL would learn unreasonably large β values and thus indirectly cause overfitting on old tasks. To our surprise, $\kappa = 0.1$ yielded 76.95 ACC and this was a minor lesion to $\kappa = 0.02$'s 82.19 ACC shown in Table 6.1.

We found that this was because that extreme β values converged as the base learner observed more tasks. As shown in Table 6.3 and Figure 6.8, only a minor proportion of β values were lower than 0.02 and larger than 0.05. Note, the figure tracks the change of each individual learning rate as a coloured line. Learning rates with large values are highlighted in colour magenta while those with small values are reflected with colour blue. This also means that locations on the figure with high saturations of the same colour indicate that many learning rates share similar values. The top and bottom rows of Figure 6.8 show those task-wise β values of the first MLP layer and the MLP softmax layer respectively. We found that as the backbone network observed more tasks, the proportion of large magnitude task-wise β decreased, while that of small magnitude task-wise β increased. Tasks that were observed later during the sequential training phase thus contributed less towards the parametric updates of a base learner.

Furthermore, a larger proportion of large values were found in the task-wise β values of the deeper base learner layers. This showed that our MetaSGD-CL algorithm prioritised the updates of those weights of the deeper layers rather than those of the shallower layers.

More Remarks and Analyses: Ablation Study

Table 6.4: Metrics for the Ablated Version of Permuted MNIST

Methods	FA1 (%)	ACC (%)
Un-ablated MetaSGD-CL	81.02	82.19
Singular	62.18	71.34
Replace all old β with 0s	69.81	74.41
Replace all old β with 0.01s	73.48	77.16
Replace all old β with 0.1s	72.57	75.80

To further understand MetaSGD-CL, we reported 3 ablation study scenarios in Table 6.4. Refer to see Figure 6.3 and Table 6.1 for the un-ablated performances. When learning a new task, we removed the learnt learning rates β of old tasks, and replaced them with 0, 0.01, and 0.1.

When β of the old tasks were replaced with 0, catastrophic forgetting occurred. This was expected as this particular setting meant that the exact replays had no effect on influencing the update of the backbone base learner. However, readers should note that the ablated performances of this scenario were still much better than the naïvely implemented sequential learning. This was likely because that by learning a learning rate per parameter per past task, MetaSGD-CL was able of reducing important overlapping features [French, 1991].

Interestingly, when those β of the observed tasks were replaced with non-zero values, the ablated performances were similar. We attributed this to the formulation of Equation (6.8); that is, as long as information of the past tasks were given, MetaSGD-CL was able to learn appropriate learning rates β for the new task to indirectly balance the importance of the new task in relation to all past observed tasks.

6.1.6 Section-wise Additional Related Work

Potential extensions

Though we designed MetaSGD-CL with access to externally stored memory, MetaSGD-CL actually shared some similarities to the regularisation-based techniques (see Section 3.3.5.1) of EWC [Kirkpatrick et al., 2017] and *Learning without Forgetting* (LwF) [Li and Hoiem, 2017]. To elaborate, all 3 techniques made direction modifications to the backward pass to alleviate catastrophic forgetting. Thus, it is likely that MetaSGD can also be integrated to continual learning in a purely regularisation-based fashion. In other words, one could think of this idea as *learning a regularisation term per parameter per past task*. This could completely remove the need to store external memory thereby making continual learning techniques even more efficient.

Meta-learning for continual learning

Our work in this section is not the only study that has entertained the idea of combining meta-learning with continual learning. Riemer et al. [2018] proposed to align GEM-like task-wise gradients with Reptile [Nichol et al., 2018], an efficient MAML algorithm [Finn et al., 2017]. While Javed and White [2019] also continually learned via representation learning with MAML-like techniques. Recall in Section 3.2.1 that though MAML is significant to the field of meta-learning, it belongs to the class of offline meta-learning and that there exist a plethora of techniques of the online meta-learning archetype. A large scale comparison for such hybrid approaches would thus be an interesting future study.

6.1.7 Section Conclusion

This section introduced **MetaSGD-CL** as an extension to the *ER* technique proposed by Chaudhry et al. [2019]. The aim was to remedy the overfitting on old tasks when memory-based continual learning techniques were allocated with small memory sizes. MetaSGD-CL leveraged on meta-learning by *learning a learning rate per parameter per past task*. Not only did MetaSGD-CL alleviate catastrophic forgetting, but it also prevented backbone networks from overfitting and achieved high performances for all tasks over the entire course of sequential learning. MetaSGD-CL also optimised backbone learners in fewer iterations, and demonstrated higher robustness against noises in the training data.

Along with Chaudhry et al.’s *ER* approach, we believe that the 2 biggest challenges in replay-based continual learning has now been addressed. As we mentioned in Section 3.3.5.3, naïve replay is both computationally expensive and memory demanding. The former problem is dealt by *ER*’s ability to prevent forgetting with exact replaying 1 past data point per network update; and that the latter problem is resolved in our novel MetaSGD-CL where sequential learning can be achieved under a low resource setting on previously externally stored memories.

6.1.8 Re-imagining Continual Learning

Thus far in this thesis, we described the three main-stream approaches in continual learning. Replay methods show externally stored memories to the network; regularisation methods impose extra constraints to preserve learnt mappings; and that dynamic architectures introduce more modules to host new knowledge. All of these approaches are inherently *external* to the sequentially learning base learner. That is, they are all adding elements to the base learner via changing the experimental setup.

In the next section, we will aim to introduce a new perspective to continual learning. Instead of adding foreign components to a sequentially learning base learner to prevent forgetting, we will aim to create a base learner that itself is robust against forgetting. We want a network that, when naïvely applied for continual learning, does not incur severe catastrophic forgetting.

It has become increasingly clearer that artificial neural networks function on a different set of principles to that of biological neural networks. Modern artificial neural networks lack dendrites, axons, and electro-chemical impulses; in replacement, they have weighted connections, layered functions, and activation functions. The in-organic components in artificial neural networks allow us to use mathematics as a tool to describe what we want and how to achieve them.

Section Content:

- Section 6.2.1: Section Introduction,
- Section 6.2.2: Thesis Recap: Highway Networks vs Gated RNNs,
- Section 6.2.3: A Dynamic Theory in the Forward Pass:
A Common Vector Field with Many Stable Equilibria,
- Section 6.2.4: Evidence in the Backward Pass:
HCNs Prevent Parametric Over-Modification,
- Section 6.2.5: Experiments and Results,
- Section 6.2.6: Section-wise Additional Related Work,
- Section 6.2.7: Section Conclusion:
A Plastic Yet Stable Backbone Learner Design.

6.2 Highway Classifier Networks for Plastic yet Stable Continual Learning

While researching continual learning literature for the last section, we realised one interesting phenomenon (see Section 3.3 and Section 6.1.8) – to our knowledge, most existing work aimed to address catastrophic forgetting through imposing *external* constraints on the backbone network. That is, the community currently lacks

mathematical interpretations on the ideal continual learning backbone architecture.

This section is based on our work in Kuo et al. [2021a] and we return to our understanding on shortcut connections gained through studying RNNs. Instead of proposing external means to prevent a learner from forgetting, this section *look inwards* and analyse neural network mechanisms to design new architectures for backbone networks to tread a fine line between plasticity and stability.

6.2.1 Section Introduction

Most if not all continual learning literature impose external means to prevent catastrophic forgetting. Regularisation (see Section 3.3.5.1) coerces the learner to simultaneously suffice multiple learning objectives and that replay (see Section 3.3.5.3) forces the backbone network to solve all tasks under a domain-imbalanced environment. While dynamic architectural approaches (see Section 3.3.5.2) reconstruct the backbone network to achieve sequential learning, most techniques require tedious compartmentalisation and hence increase network computational cost. We believe that this current discomfiture is due to a lack of mathematical description on the plasticity-stability dilemma. To elaborate, if we could precisely characterise the mathematical desiderata for sufficing both plasticity and stability, then we could design a backbone network of which itself is less prone to forgetting.

This section models neural networks as dynamical systems (Glendinning [1994], and see Section 2.3 and Chapter 4) to interpret synaptic plasticity. We will discuss desirable properties theoretically possessed by **Highway-Connection** [Srivastava et al., 2015] **Classifier Networks (HCNs)** and show that they are effective in preventing forgetting. HCNs provide several advantages for continual learning. We found that when implemented alone, HCNs exhibited strong robustness against forgetting; and when employed with continual learning techniques, the combined results achieved better performances than non-HCN baselines. We also observed that HCNs increased robustness and decreased variability in sequential learning; and helped memory-based continual learning techniques to achieve competitive performances when less memory units were allocated. More interestingly, we found that HCNs achieved competitive results when only the shallowest layer was subjected to continual learning techniques. Implying that some classifier parameters were more important than others. Leveraging on this observation, continual learning could also be made efficient by limiting forget-prevention techniques to specific layers.

6.2.2 Thesis Recap: Highway Networks vs Gated RNNs

Depth is important to the success of neural networks [Szegedy et al., 2015]. However, deep networks are hard to optimise and that highway connections [Srivastava et al., 2015] took inspirations from gated RNNs (see Section 3.1 and Chapter 4) to ease the training of very deep networks.

As stated in Section 2.1, highway connections are written as

$$\mathbf{k}_{L_\gamma} = (1 - \mathbf{T}_{L_\gamma}) \odot \mathbf{k}_{L_{(\gamma-1)}} + \mathbf{T}_{L_\gamma} \odot \mathcal{H}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})$$

with gated unit $\mathbf{T}_{L_\gamma} = \mathbf{T}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{Q}_{L_\gamma})$.

We also showed that if we were to trade depth in layers (γ) with recursion in time (t), then this formulation could be connected to multi-scale differential systems in Section 4.2 (from Equation (4.5) to Equation (4.8)) where

$$\begin{aligned} \mathbf{k}'_t &= -\mathbf{k}_{t-1} + \mathcal{H}(\mathbf{k}_{t-1}, \mathbf{W}_L) \\ \mathbf{k}_t &= \mathbf{k}_{t-1} + \zeta \odot (-\mathbf{k}_{t-1} + \mathcal{H}(\mathbf{k}_{t-1}, \mathbf{W}_L)), \text{ and} \\ \mathbf{k}_t &= (1 - \zeta) \odot \mathbf{k}_{t-1} + \zeta \odot \mathcal{H}(\mathbf{k}_{t-1}, \mathbf{W}_L). \end{aligned}$$

Highway networks can thus be recovered by substituting the discretised time step ζ with $\mathbf{T}_t = \mathbf{T}(\mathbf{k}_{t-1}, \mathbf{Q}_L)$. In the next section, we describe how differential systems can also be used to describe desirable properties for continual learning classifiers.

6.2.3 A Dynamic Theory in the Forward Pass: A Common Vector Field with Many Stable Equilibria

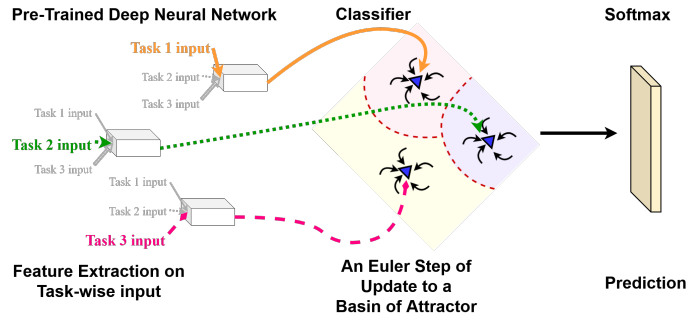


Figure 6.9: Utilising Different Basins of Attractors

Dynamical systems possess vector fields and (may have) equilibria. As stated in Section 2.3.3 and illustrated in Figure 2.9, those equilibria that are attractive are known as sinks and otherwise known as sources. We theorise that a classifier network with multiple sinks in its vector field can overcome the plasticity-stability dilemma. Network stability can be addressed via the vector field. *The existence of a common vector field* enables the classifier network to possess one *contiguous* solution space

for solving all sequential tasks; and this creates a firm structure that is robust to parametric modifications during continual network updates. In addition, plasticity can be resolved with multiple sinks. *The existence of multiple sinks* allows the solution space to be compartmentalised in order to *host* multiple skills.

We illustrate our theory as shown in Figure 6.9. If the tasks of the sequential learning continuum (ω_i of Equation (3.16)) were sampled from a non-i.i.d distribution, their features would have entered different regions of the vector field of the classifier network. In addition, these features will then take an Euler step (similar to that in the Neural Ordinary Differential Equation [Chen et al., 2018]) and sink towards their closest basins of attractors.

Two natural questions now emerge:

*What network design possess such a vector field? And,
is it possible for that vector field to possess multiple attractors?*

We hypothesised that highway connections possessed vector fields because our description in Section 6.2.2 showed that they could be linked to systems of differential equations. In addition, we hypothesised that there co-existed multiple attractors because that the high dimensional nature of deep learning models allowed them to be more expressive than traditional statistical models (see Section 2.1.2.3).

6.2.4 Evidence in the Backward Pass: HCNs Prevent Parametric Over-Modification

This subsection discusses HCNs in the perspective of the backward pass. In most literature, highway connections and residual connections [He et al., 2016a] are praised for their role in enabling deep model training (see Section 2.1). However, we will show that highway connections actually serves a second and less obvious purpose – to prevent the over-modification in parameters.

Let us first recall the well-known property of highway connections. Given a conventional feed forward neural network of Equation (1.1)

$$\mathbf{a}_{L_\gamma} = \mathcal{H}(\mathbf{a}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})$$

with Γ -layers, input $\mathbf{a}_{L_0} = \mathbf{x}$, and output $a_{L_\Gamma} = \hat{\mathbf{y}}$ for target \mathbf{y} , its derivative of a cost \mathcal{C} with respect to its parameters would be

$$\frac{\partial \mathcal{C}}{\partial \mathbf{W}_{L_\gamma}} = \frac{\partial \mathcal{C}}{\partial \hat{\mathbf{y}}} \left[\prod_{\zeta=\gamma}^{\Gamma-1} \frac{\partial \mathbf{a}_{L_{(\zeta+1)}}}{\partial \mathbf{a}_{L_\zeta}} \right] \left\{ \frac{\partial \mathbf{a}_{L_\gamma}}{\partial \mathbf{W}_{L_\gamma}} \right\} \quad (6.10)$$

$$= \frac{\partial \mathcal{C}}{\partial \hat{\mathbf{y}}} \left[\prod_{\zeta=\gamma}^{\Gamma-1} \mathcal{H}'(\mathbf{a}_{L_\zeta}, \mathbf{W}_{L_{(\zeta+1)}}) \right] \left\{ \frac{\partial \mathcal{H}(\mathbf{a}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})}{\partial \mathbf{W}_{L_\gamma}} \right\}. \quad (6.11)$$

That is, gradient vanishing would occur through $\prod_{\zeta=\gamma}^{\Gamma-1} \frac{\partial \mathbf{a}_{L_{(\zeta+1)}}}{\partial \mathbf{a}_{L_\zeta}}$ as the amount of layers increase and that individual $\left\| \mathcal{H}'(\mathbf{a}_{L_\zeta}, \mathbf{W}_{L_{(\zeta+1)}}) \right\| < 1$. Both the residual connection

and the highway connection mitigate this via introducing identity \mathbb{I} to individual $\frac{\partial \mathbf{a}_{L_\gamma}}{\partial \mathbf{a}_{L_{(\gamma-1)}}$ to allow gradients to flow through many layers (see Section 2.1.1).

The descriptions given above mainly targeted the $[\cdot]$ term of Equation (6.10). However, a less obvious change is also made in the $\{\cdot\}$ term as a consequence of the forward pass modification. Consider the residual connection of Equation (2.1) where

$$\mathbf{n}_{L_\gamma} = \mathbf{n}_{L_{(\gamma-1)}} + \mathcal{H}(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma}).$$

This leads to $\frac{\partial \mathbf{n}_{L_\gamma}}{\partial \mathbf{W}_{L_\gamma}} = \frac{\partial \mathcal{H}(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})}{\partial \mathbf{W}_{L_\gamma}}$ of which is identical to that in a simple feedforward network $\frac{\partial \mathbf{a}_{L_\gamma}}{\partial \mathbf{W}_{L_\gamma}} = \frac{\partial \mathcal{H}(\mathbf{a}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})}{\partial \mathbf{W}_{L_\gamma}}$ because that $\mathbf{n}_{L_{\gamma-1}}$ is independent of \mathbf{W}_{L_γ} for all arbitrary \mathbf{n}_{L_γ} .

On the other hand, highway connection networks⁷ have $\frac{\partial \mathbf{k}_{L_\gamma}}{\partial \mathbf{W}_{L_\gamma}} = \mathbf{T}_{L_\gamma} \frac{\partial \mathcal{H}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})}{\partial \mathbf{W}_{L_\gamma}}$. Since all units of gate \mathbf{T}_{L_γ} lie between⁸ $[0, 1]$, the gated units also limit the extent of update to parameters \mathbf{W}_{L_γ} . This is a desirable characteristic for sequential learning, because it makes catastrophic forgetting harder to occur as the knowledge stored within the classifier weights receives less modification.

Demonstration

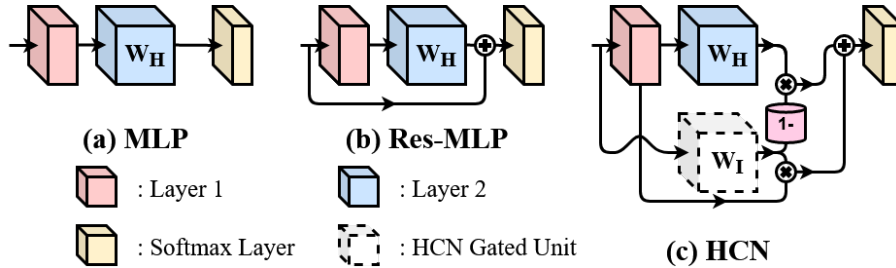


Figure 6.10: Different Designs of Continual Learning Classifier Networks

In order to validate our mathematical analysis on how architectural designs could impact sequentially learning, we tested three classifier network designs and compared their corresponding $\{\cdot\}$ term in Equation (6.10). The three different classifier networks were MLPs, MLPs with residual connections (Res-MLPs), and MLPs with highway connections (HCNs) illustrated above in Figure 6.10. The classifier networks were tested on the MNIST [LeCun et al., 1998] and the KMNIST [Clanuwat et al., 2018] datasets. Descriptions of both datasets can be found in Appendix A.

⁷Refer to Equation (14) of Hochreiter and Schmidhuber [1997b] for a similar backward pass of the LSTM cell. Our \mathbf{T}_{L_γ} corresponds to their $y^{ini}(t)$; and our $\mathcal{H}(\mathbf{k}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})$ corresponds to their $g'(net_{c_j^v}(t)) \frac{\partial net_{c_j^v}(t)}{\partial w_{lm}}$.

⁸It is standard practice to activate highway connection gated units with the sigmoid function; and this was because that highway connections were based on the LSTM RNN gated units.

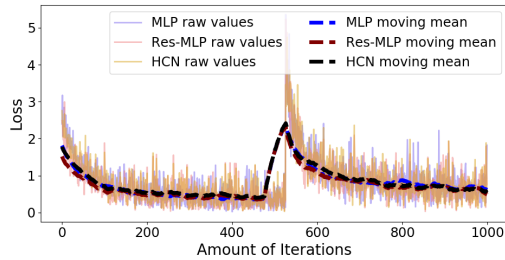


Figure 6.11: Loss vs Iteration (MNIST)

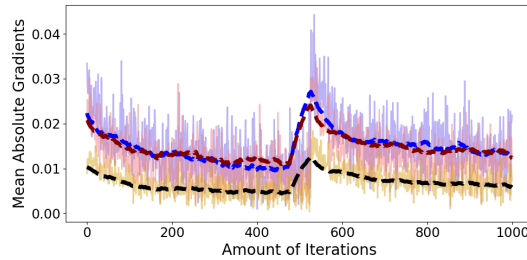


Figure 6.12: Gradient vs Iteration (MNIST)

All of our classifier networks had 100-hidden dimensions, and were updated with SGD with learning rate 0.01 for 1050 steps. The first 525 steps were on MNIST, and the remaining steps were on KMNIST. The purpose of this experiment was to inspect the behaviours of the classifier loss and gradient during training. The behaviours of the losses are in Figure 6.11, while those of their respective gradients are in Figure 6.12. The classifiers were updated for 1050 steps (instead of 1000) to calculate moving averages with window size 50 which we presented as the dashed lines in the figures.

The gradients of our interest were the mean absolute values of \mathbf{W}_{L_2} s of Equations (1.1), (2.1), and (2.3) for MLPs, Res-MLPs, and HCNs, respectively. That is, we recorded those parameters which were responsible for the transformation at layer 2. The gradients of layer 2 was specifically measured because that, as illustrated in Figure 6.10, layer 2 was the location where the output features were modified by the Res-MLP and HCN connections. We took the mean of the absolute values of the gradients because that first, individual gradients could be of either signs, and second, that there were more than just one weight in \mathbf{W}_{L_2} .

While we observed that there were virtually no differences in the losses, the gradients behaved qualitatively differently. The gradients of HCNs were much lower than those of their rivals; while those of Res-MLPs and MLPs behaved similarly. These findings verified our analysis such that the gates of highway connections lowered gradient values in the backward pass.

6.2.5 Experiments and Results

We tested three datasets on six techniques with three classifier designs and summarised our results in Table 6.5 and Table 6.6.

Datasets, Backbone Networks, and Sequential Tasks

We tested image permutation [Kirkpatrick et al., 2017] on MNIST [LeCun et al., 1998], and incremental classes [Rebuffi et al., 2017] on Cifar100 [Krizhevsky, 2009] and CUB200 [Wah et al., 2011]. Permuted MNIST and incremental Cifar100 were previously experimented in Section 6.1.5; the details could be found in that section. Our readers can find the dataset descriptions in Appendix A; and we followed Douillard et al. [2020]⁹ to process the incremental datasets.

⁹Refer to https://github.com/arthurdouillard/incremental_learning_pytorch; and the specificities can be found in script `incremental_learning_pytorch/inclearn/lib/data/datasets.py`.

We sequentially trained the classifier networks shown in Figure 6.10 over multiple permuted image tasks. As for incremental classes, we repeated our setup mentioned in Section 6.1.5 and adopted Mai et al. [2020]’s network compartmentalisation. That is, we employed fixed pre-trained feature extractors and only sequentially updated the classifier networks. The fixed feature extractor networks for incremental classes were ResNet18s [He et al., 2016a]; and we pre-trained them with the first 50 classes from Cifar100 and the first 100 classes from CUB200.

All of our classifiers were 2 layers deep followed by a softmax layer. The hidden dimensions (HDs) for the 3 classifiers were defaulted as 100. However, the HCN gated units introduced more parameters; thus in order to maintain a fair comparison, we also tested Large-MLPs – MLPs with 175-HD of which had a similar total amount of weights to our novel HCNs.

In this section, we tested 20 tasks for permuted MNIST. In addition, we tested our (semi-)pre-trained backbone networks with the remaining unseen classes in Cifar100 and CUB200. There were 10 tasks for incremental Cifar100 with 5 classes per task, and 10 tasks for incremental CUB200 with 10 classes per task.

Metrics

We adopted three metrics from Lopez-Paz and Ranzato [2017] with matrix \mathcal{R} defined as described Section 6.1.5:

$$\begin{aligned} \text{Average accuracy} \quad (\text{ACC}) & \quad \frac{1}{n} \sum_{i=1}^n \mathcal{R}_{(n+1),i} \\ \text{Backwards transfer} \quad (\text{BWT}) & \quad \frac{1}{n-1} \sum_{i=1}^{n-1} (\mathcal{R}_{(n+1),i} - \mathcal{R}_{(i+1),i}), \text{ and} \\ \text{Forward transfer} \quad (\text{FWT}) & \quad \frac{1}{n-1} \sum_{i=2}^n (\mathcal{R}_{i,i} - \mathcal{R}_{1,i}). \end{aligned}$$

ACC computes the mean accuracy after observing all tasks $\mathcal{R}_{(n+1),i}$, and BWT quantifies task-wise forget since $\mathcal{R}_{(i+1),i}$ immediately after training on a task’s own data. FWT measures task-wise zero-shot learning [Socher et al., 2013] from an initialised learner $\mathcal{R}_{1,i}$ to $\mathcal{R}_{i,i}$ immediately before observing the task’s own data. For all metrics, the larger and more positive the better.

Results on Naïve Sequential Learning

We naïvely trained HCNs, MLPs, Large-MLPs, and Res-MLPs on permuted MNIST with SGD with learning rate 0.01 without any continual learning techniques. The accuracy of Task 1 was plotted against the amount of tasks observed in Figure 6.13. *Oracle* denoted a single MLP thoroughly trained on all data. The least amount of forgetting was exhibited by HCNs.

Similar results could be found in Figure 6.14 when we naïvely trained the (semi-)pre-trained backbone networks for the incremental Cifar100 tasks. All classifier networks of this scenario were also naïvely updated with SGD with learning rate 0.01.

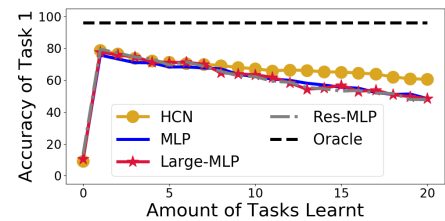


Figure 6.13: ACC of Naïve Sequential Learning (MNIST)

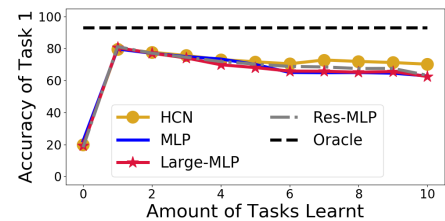


Figure 6.14: ACC of Naïve Sequential Learning (Cifar100)

Remarks and Analyses

The **metrics** for these performances are listed in Table 6.5. We found that Res-MLPs only performed comparably to MLPs. This was expected because that as demonstrated in Section 6.2.4, residual connections incurred the same amount of gradient as that of a conventional MLP without it. Hence the backward pass of a Res-MLP modified its parameters to the same extent as that in the backward pass of an MLP; and as we conjectured in Section 6.2.4, highway connections were crucial and irreplaceable with residual connections.

Interestingly, we also found that Large-MLPs did not outperform the normal-sized MLPs. Hence simply adding more parameters was not successful in preventing forgetting even though Large-MLPs were supposed to consist more synaptic connections. Since our novel HCNs achieved the best results, our following experiments will only proceed by comparing HCNs against MLPs.

Results on HCNs vs MLPs with Continual Learning Techniques Applied

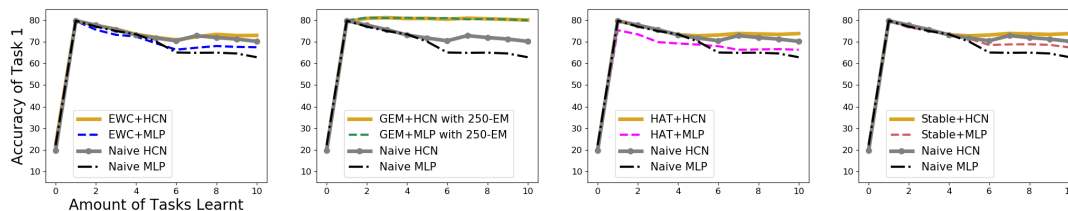


Figure 6.15: ACC with Different Techniques on Incremental Cifar100

After naïvely testing the classifier networks, we then tested HCNs against MLPs in conjunction with different archetypes of continual learning strategies including EWC [Kirkpatrick et al., 2017], GEM [Lopez-Paz and Ranzato, 2017], HAT [Serra et al., 2018], and *Stable continual learning* (Stable) [Mirzadeh et al., 2020b]. The first three techniques were previously employed in Section 6.1.5, and see Section 3.3.5 on different approaches for sequential learning. Mirzadeh et al.’s *Stable continual learning* slowly decreased the learning rate for the backbone networks across tasks. The combined efforts were tested on permuted MNIST, incremental Cifar100, and incremental CUB200. In total, there were 18 **comparisons** (see Table 6.5); and we found that not only was HCNs compatible with all techniques, but it performed favourably with nearly all techniques.

More Remarks and Analyses

We selected the typical sequential learning behaviours of this scenario and presented them in Figure 6.15. The specifically selected results were from training (semi-)pre-trained backbone networks on incremental Cifar100. From the images, we observed that HCNs outperformed MLPs on incremental Cifar100 regardless of the applied method. In addition, those with naïvely trained HCNs also achieved better performances than those with MLPs trained with EWC, HAT and Stable. This showed that though continual learning techniques were important for sequential learning, an appropriate backbone architecture was more impactful for learning many skills.

Table 6.5: An Overview of the Results

Larger & positive is better for ACC, BWT, and FWT.
EM – Episodic Memory; HD – Hidden Dimension.

Method	ACC	BWT	FWT
Naïve Implementation for <i>Permuted MNIST</i>			
HCN (Ours)	71.49 ± 11.57	-7.71	4.57
MLP	65.54 ± 19.76	-14.39	4.01
Large-MLP	68.06 ± 18.81	-13.29	3.31
Res-MLP	65.36 ± 19.85	-15.28	2.96
Naïve Implementation for <i>Incremental Cifar100</i>			
HCN (Ours)	74.82 ± 15.48	-6.10	0.26
MLP	73.08 ± 17.49	-8.40	0.25
Large-MLP	73.44 ± 17.84	-7.92	-0.09
Res-MLP	70.84 ± 18.78	-11.34	-0.85
In-Class Comparison for <i>Permuted MNIST</i>			
EWC + HCN (Ours)	71.74 ± 7.84	-7.02	4.65
EWC + MLP	66.72 ± 13.78	-12.67	4.10
GEM + HCN (Ours)	74.19 ± 8.24	6.22	4.74
GEM + MLP	74.52 ± 4.27	3.52	3.93
HAT + HCN (Ours)	70.72 ± 8.32	-3.39	3.69
HAT + MLP	62.50 ± 20.41	-5.70	2.12
HAT with 500-HD + HCN (Ours)	78.00 ± 5.44	-2.02	4.30
HAT with 500-HD + MLP	76.57 ± 7.99	-2.36	2.70
Stable + HCN (Ours)	67.38 ± 6.35	-4.43	5.49
Stable + MLP	64.71 ± 7.54	-7.71	5.48
In-Class Comparison for <i>Incremental Cifar100</i>			
EWC + HCN (Ours)	75.47 ± 13.47	-5.34	0.19
EWC + MLP	73.12 ± 18.16	-8.27	0.15
GEM with 25 EM + HCN (Ours)	75.19 ± 13.82	-5.53	0.16
GEM with 25 EM + MLP	74.64 ± 15.05	-6.14	0.30
GEM with 250 EM + HCN (Ours)	80.84 ± 9.18	0.70	-0.30
GEM with 250 EM + MLP	80.96 ± 9.71	0.61	0.10
HAT + HCN (Ours)	76.04 ± 14.15	-2.75	1.36
HAT + MLP	73.23 ± 16.54	-3.69	0.83
Stable + HCN (Ours)	69.66 ± 16.81	-3.23	-0.04
Stable + MLP	69.51 ± 14.25	-5.69	0.25
A-GEM + HCN (Ours)	81.07 ± 8.99	1.36	0.17
A-GEM + MLP	80.96 ± 9.65	1.30	0.19
ER + HCN (Ours)	80.35 ± 8.66	0.76	-0.65
ER + MLP	80.00 ± 10.66	0.31	0.01
In-Class Comparison for <i>Incremental CUB200</i>			
EWC + HCN (Ours)	44.16 ± 33.26	-11.42	0.44
EWC + MLP	38.91 ± 29.54	-15.83	0.19
GEM + HCN (Ours)	44.67 ± 32.61	-10.66	0.51
GEM + MLP	39.96 ± 30.52	-15.57	0.28
HAT with 400-HD + HCN (Ours)	48.16 ± 32.54	-3.73	0.15
HAT with 400-HD + MLP	42.63 ± 30.31	-2.81	-0.17
Stable + HCN (Ours)	40.71 ± 29.38	-7.47	-0.05
Stable + MLP	38.22 ± 27.02	-11.51	0.66

Results on Unique HCN Advantages for EWC, GEM, and HAT

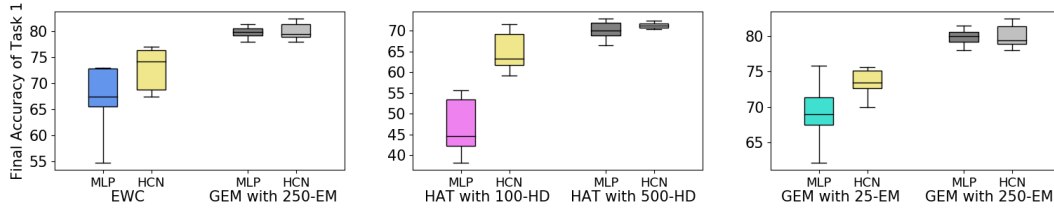


Figure 6.16: HCNs Provided Robustness and Decreased Variability

There were several implicit advantages in employing HCNs that were unreflected through the metrics recorded in Table 6.5 in the previous page. In order to fully investigate the performance gains brought by our HCNs, we plotted the final accuracy of Task 1 after training on all incremental Cifar100 tasks in Figure 6.16. We will mainly be focused on the coloured box-plots, of which we find HCNs outperformed MLPs on all scenarios. We kept the grey-scale box plots only as control groups.

More Remarks and Analyses

We found that EWC could lead to unstable MLP performances; and that MLPs with low hidden dimensions (HDs) could not yield high accuracy with HAT (also see Section 6.1.5). HCNs were able to increase both accuracy and performance stability (low variance) for both methods. We also found that GEM only achieved exceptionally good results when they stored large volumes of data in their episodic memories (EMs), and that MLP performances dropped sharply when 250 EMs were decreased to 25 EMs¹⁰. In contrast, HCNs were able to still maintain higher accuracy.

The successes could be attributed to highway connections ability to scale down gradients (see Section 6.2.4). Recall from Section 2.2 of the preliminaries chapter, trained neural networks perform well because their parametric configurations lie within a minimum on the loss landscape. And by reducing the gradient for a subsequent task, this means that HCNs are capable of moving out of a local minimum slower than their MLP rivals. This allows the HCN parametric configurations to readjust itself for every iteration to look for compromising *minimum within a minimum* to achieve good results for both the new and the old tasks.

Results on Continual Learning with Extreme Network Compartmentalisation

This experiment tested *split-HCN* for permuted MNIST. We trained the entire HCN for the first task but then froze the parameters of both the gated unit and layer 2 to only sequentially train layer 1 of the classifier network. Since the frozen parameters belonged to the highway connection, this meant that we tested whether sequential learning could be made successful when the vector field remained unmodified (see theory in Section 6.2.2). Table 6.6 revealed that *split-HCN* yielded competitive metrics against the fully sequential learning HCNs of Table 6.5. In addition, Figure 6.17 showed that *split-HCN*s were better than HCNs in preventing forgetting.

¹⁰This was unsurprising, and was previously investigated with results shown in Figure 6.5 on low memory continual learning.

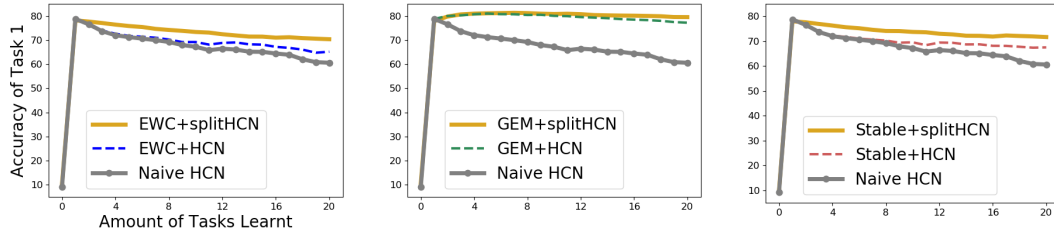


Figure 6.17: Sequential Learning while Only Fine-Tuning Layer 1 of HCNs

Table 6.6: An Overview of the Results

Method	ACC	BWT	FWT
Extreme Network Compartmentalisation on HCNs for <i>Permuted MNIST</i>			
EWC + split-HCN (Ours)	66.30 ± 6.21	-5.85	6.13
GEM + split-HCN (Ours)	76.12 ± 7.07	7.32	5.20
Stable + split-HCN (Ours)	59.94 ± 10.04	-3.04	6.10
Baselines results from in Table 6.5			
EWC + HCN (Ours)	71.74 ± 7.84	-7.02	4.65
GEM + HCN (Ours)	74.19 ± 8.24	6.22	4.74
Stable + HCN (Ours)	67.38 ± 6.35	-4.43	5.49

More Remarks and Analyses

By fixing the weights of the gated unit in the HCN, we prohibited the network to learn new ways to mask the activation of layer 1. Thus, it was not surprising to see that the results in Table 6.6 were inferior to those of the baselines. However, it should be noted that the lesions were not severe. Hence it is possible that not all components within the HCN required sequential training for acquiring new skills.

6.2.6 Section-wise Additional Related Work

HCN vs dynamic architectural approaches

Though our HCN introduced modifications to the traditional sequential learning classifier networks, we do not consider it as a part of the dynamic architectural continual learning archetype (see Section 3.3.5.2). This is because that HAT, progressive neural networks [Rusu et al., 2016], and incremental denoising autoencoders [Zhou et al., 2012] all introduced new parameters to host newly acquired knowledge. Our HCN architecture did not expand in size; and its successes were based on the dynamical mathematical formulation of the plasticity-stability dilemma.

HCN and related dynamical machine learning literature

The recent work of Yang et al. [2020] may provide alternative explanations to why residual connections [He et al., 2016a] were not effective in preventing forgetting. They showed that residual connections can be modelled as ordinary differential equations (ODE) but they were not as robust as damped ODE models of which followed a similar formulation to Equation (4.5)¹¹.

¹¹Refer to page 3 Equation (17) of Yang et al. [2020], where damped ODE: $\frac{dx}{dt} = -\lambda x(t) + \rho(\lambda)f(x(t))$.

HCN and representation learning

The effectiveness of HCNs could also be explained by [Mirzadeh et al. \[2020a\]](#). They showed that dropout [[Hinton et al., 2012](#)] prevented forgetting via learning sparse activation patterns as an implicit gating mechanism. Hence HCNs' effectiveness could also be attributed to its explicit gating mechanism. In addition, [Kurtz et al. \[2020\]](#) showed that shallow layers generally had lower sparsity than later layers; implying that neural networks are mainly driven by features generated in the shallower layers. This could explain why our extreme split-HCN network compartmentalisation did not incur large performance lesion.

Continual learning effects on the forward pass and backward pass

We emphasise that HCN took a different approach to prevent forgetting than EWC [[Kirkpatrick et al., 2017](#)] and GEM [[Lopez-Paz and Ranzato, 2017](#)]. On one hand, EWC is a regularisation continual learning technique that introduced a surrogate loss to minimise modification in important learnt weights (see more in Section 3.3.5.1). On the other hand, GEM influenced parametric updates in the backward pass by preconditioning the native gradients into a compromised direction that minimises gradient misalignment (also see Section 3.3.5.3). In stark contrast, HCN influenced variable propagation in the backward pass by explicitly adding gated units as extra variables to the forward pass of the backbone network. As shown in Table 6.5, HCN thus serves as a simple and generally applicable patch of which is compatible to all existing archetypes of continual learning techniques.

6.2.7 Section Conclusion

This paper introduced **Highway-Connection Classifier Networks** (HCNs) to prevent forgetting in sequential learning multiple tasks. We theorised that HCNs could overcome the stability-plasticity dilemma by possessing a vector field with multiple basins of attractors. We demonstrated that HCNs were more effective than MLPs in sequential learning; and that highway connections were specifically required and are irreplaceable with residual connections. HCNs were also versatile and could be employed in conjunction with the regularisation-based EWC, memory-based GEM, and network-based HAT continual learning techniques. More importantly, naïve HCNs also performed favourably against MLPs with EWC and HAT. Hence implying that the architectural design of a sequential learning backbone could be more important than the applied continual learning techniques.

6.3 Chapter Conclusion

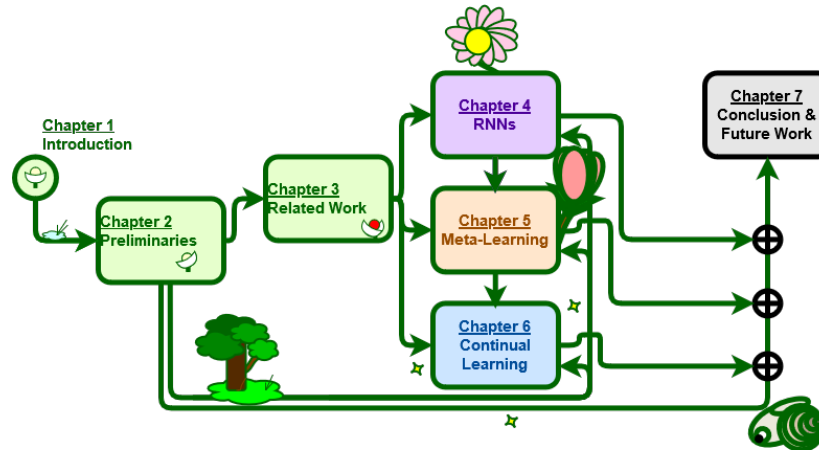


Figure 6.18: Progress Update

The continual learning chapter of this thesis was built directly on top of our results from RNN understandings (see Chapter 4) and meta-learning extensions (see Chapter 5). The former allowed us to mathematically highlight what network components could incur less catastrophic forgetting; and the latter allowed us to devise more powerful continual learning techniques.

The first section of this chapter showed that the replay continual learning archetype could benefit greatly when conjointly trained with a meta-learner. By learning a learning rate per parameter per past task, our novel MetaSGD-CL technique achieved rapid knowledge acquisition with a set of small externally stored data. Furthermore, MetaSGD-CL enabled a backbone to sequentially learn robustly with data collected from a noisy environment.

The second section of this chapter emphasised on the importance of the correct architecture of choice for the sequential acquisition of skills. Unlike conventional continual learning techniques which impose external constraints on a sequential learning backbone, we made internal modifications to the backbone. Our internal modifications were based on a mathematical analysis on the backward pass with insights on properties of the shortcut connection during parametric updates. We showed that by adopting a highway connection and using HCNs, we could achieve superior performances over conventional MLP classifier networks.

Last, this chapter showed that all machine learning disciplines considered in this thesis could work in conjunction with another. As illustrated in Figure 6.18, all concepts are entwined and we have demonstrated that they complement each other cohesively to create better AutoML algorithms.

Thesis Conclusion

This thesis used the mathematical language of *dynamical systems* [Glendinning, 1994] to examine machine learning algorithms. Our studies were devoted towards *automated machine learning* (AutoML) [Hutter et al., 2019] for the goal of automating tedious aspects in deep learning such as hyper-parameter tuning to create a clean engineering pipeline. This chapter summarises our approach, highlights our contributions, and connects them with existing literature. We will also state some of our limitations and propose ideas for future work.

7.1 Thesis Summary

This thesis aimed to devise better algorithms to automate complicated machine learning models (AutoML, see Chapter 1). Modern neural networks often include *shortcut connections* [Schraudolph, 1998] – a design in the forward pass that has important implications in the backward pass. Hence to make AutoML practical, it required us to gain a better understanding on neural network *inference* and *optimisation*. To that end, we studied *recurrent neural networks* (RNNs) to understand properties of shortcut connections; and on how RNNs could be applied to update another network. These knowledge then served as the foundation of our novel AutoML contributions. Below, we provide a short description for the elements of our research framework; and we will elaborate each element in Section 7.2.

Methodology Revised

A progression of the literature in inference and optimisation was presented in Chapter 2 as the preliminaries of this thesis. On one hand, Section 2.1 mentioned how the state-of-the-art deep learning models of *residual networks* (ResNets) [He et al., 2016a] and *highway networks* [Srivastava et al., 2015] owed their origins to the *long short term memory* (LSTM) RNN [Hochreiter and Schmidhuber, 1997b] and shortcut connections. On the other hand, Section 2.2 discussed how *loss landscape* and *loss minima* [Hochreiter and Schmidhuber, 1997a] highlighted the difficulties in network parameterisation. In Section 2.3, we emphasised that the recursive nature of inference and optimisation can both be described by the rich vocabularies of the mathematical study of dynamical systems.

Understanding and Simplifying RNNs Inferences

The recursive nature of inference and optimisation prompted us to study AutoML through RNNs, *meta-learning* [Andrychowicz et al., 2016], and *continual learning* [Lopez-Paz and Ranzato, 2017]. In Chapter 4, we studied the quasi-recursive inferential nature of deep feedforward models with RNNs. Our work had a heavy focus on LSTMs and *gated recurrent units* (GRUs) [Cho et al., 2014]. Through a dynamical system analysis, we returned to a first principle’s manner and attached meaning to each RNN element. This helped us to distinguish the roles of individual RNN variables and identify redundant parameters that exist in current de-facto standard RNN designs. Not only were we able to propose a systematic simplification applicable to multiple well-known gated RNN architectures; our study also provided us with deeper insights on the properties of shortcut connections which facilitated the deep learning of ResNets and highway networks.

Extending Optimisation-based Meta-Learning Approaches with RNN Insights

Though our RNN chapter focused solely on network inference, they produced insights towards our understanding on optimisation. In Chapter 5, we studied the AutoML technique of meta-learning which enabled a base learner network to be configured quickly with fewer data. The effort of that chapter was dedicated towards Andrychowicz et al. [2016]’s *learning to learn* (L2L) algorithm of *RNN-optimisers*. An RNN-optimiser substituted the conventional network configuration scheme of *gradient descent* [Rumelhart et al., 1986] with an RNN. Their setup leveraged the recurrent inferential processes of RNNs to explore and learn alternative optimisation strategies that were more powerful than the handcrafted *stochastic gradient descent* (SGD) [Robbins and Monro, 1951]. We investigated the disadvantages of current L2L algorithms with dynamical systems and proposed remedies based on our RNN insights. We also addressed a list of desiderata to make L2L algorithms practical in real world scenarios. Then based on our requirements, we further developed a novel neural optimiser that achieved efficient non-iterative network pruning [Han et al., 2015] thereby making the pruning procedure less computationally expensive. This demonstrated that L2L utility could be extended beyond just rapid knowledge acquisition.

Conceptualising and Addressing the Causes of Catastrophic Forgetting

In Chapter 6, we combined our findings on inference and optimisation to address the AutoML technique of continual learning. Continual learning concerns the setting when a single learning agent is placed in a dynamic environment and is required to adapt to learn new predictions. This is a difficult task as neural networks are known to suffer from *catastrophic forgetting* [McCloskey and Cohen, 1989] where they abruptly lose previously learnt skills while assimilating knowledge for solving a new task. Sitting at the heart of continual learning is the *stability-plasticity dilemma* [Mermillod et al., 2013]. Most conventional neural networks have plastic weights that can learn skills quickly, but they are unstable such that their learnt knowledge can be easily corrupted due to parametric modifications. One effective approach to mit-

igate forgetting is via storing and replaying data of an old task with those of the new problem at hand [Lopez-Paz and Ranzato, 2017]. With our insights on optimisation, we introduced meta-learning for continual learning to configure a backbone learner quickly with fewer data thereby greatly reducing the computational memory cost. Then with our understanding of shortcut connections, we identified network architectures that incurred less parametric modifications during optimisation while causing no lesions to feature representation during inference.

7.2 Thesis Contributions

In this section, we will highlight the contributions made in each chapter and connect them with existing literature.

7.2.1 On Increasing the Interpretability for RNNs

RNNs are important neural network modules that have pushed the state-of-the-art results across multiple machine learning disciplines. However, it was unclear how RNNs recursively process their input data. Moreover, de-facto standard RNNs have intricate designs. Both LSTMs and GRUs possess internal variables known as *gated units* to synchronise the *cell* – a hidden variable to represent occurred events. Though important, gated units do not have immediately clear meanings and this obscures how RNNs conduct feature representation.

Previous Studies

Since its inception, the research community has proposed several approaches to explain RNN behaviours. As we described in Section 3.1.6, RNN behaviours were mostly studied through *ablation studies*, *search studies*, and task-specific *qualitative analyses*. While these research have created some important and useful insights, most of their findings were non-conclusive. For one instance, the popular ablation study of Greff et al. [2016] compared the performances of the vanilla LSTM against its variants with only a subset of all LSTM gated units. Across multiple tasks, the authors of that paper found that none of the ablated LSTMs could perform on par to the full LSTM design. For another instance, the search study of Jozefowicz et al. [2015] conducted a large scale RNN mutation study. By randomly replacing node computations and network connections, they hoped to discover architectural designs that outperform LSTMs and GRUs consistently across multiple experiments. Though their work was important, they were not successful in finding a superior RNN. Let us further consider the qualitative RNN study of the NLP in text by Karpathy et al. [2016]. By visualising the long-term behaviours of the LSTM gated units and hidden variables, their paper found that character-based LSTM models reacted differently to a subset of special characters. However, the authors of that paper also admitted that they were uncertain of how to interpret the alterations observed in their LSTM gated units¹.

¹Refer to Footnote 16 of Section 3.1.6.3.

Ablation studies, search studies, and qualitative analyses all share a commonality – that they aimed to understand the behaviours of a *trained* RNN. Though these research have found valuable insights, they were heavily dependent on their experimental setup and their results hence could not be guaranteed to generalise across different tasks. In order to explore superior RNN designs, we argued that we must have a deeper understanding on the mathematical properties of RNNs. That is, we need to return to a first principle’s manner and clarify the meaning of each component therein the complicated RNN design. To this end, we modelled RNNs’ recurrent transformations with dynamical systems to understand how they process input data.

Our Novelty

Chapter 4 provided two different interpretations for LSTMs with dynamical systems. We first studied LSTMs as systems of difference equations in Section 4.1; and then analysed LSTMs as systems of differential equations in Section 4.2. Both of our studies focused on the LSTM cell of $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{a}_t$ (see Equation (3.7)) where the network represented the continual exposure of context from input data.

Section 4.1: On Local and Discretised LSTM Variable Propagations

Since the cell was updated for every timestep t , we studied the forget gate \mathbf{f}_t and input gate \mathbf{i}_t as the discrete dynamical systems of difference equations. We found that alterations in the cell state values were controlled by two opposing dynamic regimes. A more influential *catch regime* dictated by \mathbf{f}_t and a less contributive *release regime* influenced by \mathbf{i}_t ; the former diminished cell state magnitudes while the latter granted growth in cell states instead.

However, the forget gate and input gate could oscillate as if they were two independent variables. That is, the reason and timing why the cell state magnitudes contracted or expanded were made difficult to interpret. In order to increase network interpretability, we modified LSTMs with a monotonically decreasing forget gate which we called *DecayNet*. This architectural modification has three features. First, the monotonic decrease in the forget gate enabled DecayNet cells to slowly converge to different sets of values dependent on the input data. Second and expanding on feature one, DecayNets dichotomised its cell states as recurrent cells and *quasi-feedforward cells*. Those forget gates with values near-zero would effectively turn their respective cell state propagation to $\mathbf{c}_t = \mathbf{i}_t \odot \mathbf{a}_t$. Third, the monotonicity of the decreasing forget gate also acted as a *hard regularisation*. Unlike the soft regularisation of Dropout [Srivastava et al., 2014] which stochastically introduced noises to the network during training, the monotonic decay in the forget gate was deterministic by design and that it occurred both during training and in test time. Features one and two clarified RNN mechanisms, while feature three increased RNN generalisation. Empirically, we found that DecayNets achieved state-of-the-art results on the challenging permuted sequential MNIST task [Le et al., 2015].

The insights of this section could serve as nice supplements to search studies. The main focus of search studies were to mutate and explore for unknown superior architectures. However, instead of redrafting an RNN design from a *clean canvas*, it

would be beneficial to have prior knowledge on the desired features for each RNN component. For instance, we could conduct a search study that specifically targeted the influential catch dynamic regime of the LSTM cell (the $\mathbf{f}_t \odot \mathbf{c}_{t-1}$ term) to drastically limit search space. Search space limitation has already shown to be a practical technique in *neural architectural search* (NAS) [Zoph and Le, 2017]. In the work of Pham et al. [2018b], the authors found that mutations with parametric sharing could still lead to successful discovery of competitive models while drastically reducing search time.

Section 4.2: On the Desired Properties for LSTM Memory Updates

While we continued to examine LSTMs through the lens of dynamical systems, this section took a very different approach to the last section. The previous section was an *atomic* mathematical analysis which we studied the individual and *unpaired relationships* between the forget gate and input gate. However and to fully understand LSTMs, we aimed to gain a deeper understanding on how different components could work cooperatively. To this end, we established connections between LSTMs with the well-studied *Hopfield network* [Hopfield, 1982] for a higher level interpretation of the roles of the LSTM components.

Hopfield networks are well-studied in the realm of theoretical physics [Sompolinsky et al., 1988] as well as in the artificial neural network discipline of *reservoir computing* [Jaeger, 2001; Sussillo and Abbott, 2009; Kuo, 2017]. Dynamic Hopfield networks have two very nice properties – they possess both *attractive equilibria* and *phase transitions* [Glendinning, 1994]. On one hand, the existence of attractive equilibria provides the network solution space with stability. On the other hand, the existence of phase transitions signifies that there are parametric transitional boundaries for the internal degrees of freedom to successively fall out of an equilibrium. One of our main contributions in this section was a derivation that started with the dynamic Hopfield network which eventually arrived at the LSTM formulation (see Section 4.2.2.2). This connection revealed that the LSTM cell was both stable and plastic. The recursive transformations in an LSTM were thus robust to noises but also possess high flexibility and could be shaped differently according to the input data.

Beyond establishing the connection, our mathematical framework also found three important insights on the purpose of the LSTM gated units in the forward pass. First, we revealed that the importance of the input gate to an LSTM was equivalent to the roles of timescales for differential systems. That is, the input gate controlled the rate of propagation for cell state values. Second, we found that it was unnecessary for an LSTM to possess a forget gate given the presence of the input gate in the network. This echoed with one of the findings in the ablation study of Greff et al. [2016]. In that paper, the authors found that LSTMs did not suffer great performance lesions when the forget gate was bound to the input gate like that in the GRU design. Third, we found that the co-existence of a forget gate and an input gate was undesirable. The potentiality of $\mathbf{f}_t + \mathbf{i}_t \geq 1$ meant that the cell state could be introduced with excessive information thereby growing uncontrollably. This indirectly highlighted the importance of the LSTM hidden state of $\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$ (see Systems of equa-

tions 3.5). Note that the hidden state provided a two-fold regulation on the cell state values. The extreme values of the cell state were clamped by the tanh function and then scaled down by the output gate \mathbf{o}_t . This echoed with another finding reported in Greff et al.; in that paper, the authors found that an ablated LSTM without \mathbf{o}_t caused some of the worst lesions in their experiments. Overall, our mathematical framework highlighted the utilities of all LSTM gated units and hence increased the interpretability of the inferential process of the LSTM RNN.

Besides connecting LSTMs to theoretical physics and highlighting the purposes of gated units, our derivation unveiled a third and even more significant finding – that multiple important gated RNNs could be grouped as one unified class of mathematical models. Our novel classification included LSTMs, GRUs, SRUs [Lei et al., 2018a], T-LSTMs [Balduzzi and Ghifary, 2016], and FastGRNNs [Kusupati et al., 2018]; and it could be potentially expanded to include more designs. In addition, our mathematical framework revealed redundant parameters that were only secondary in importance. Based on our analysis, we proposed *input residual connection* (IRC) and *input highway connection* (IHC) to systematically simplify all aforementioned gated RNNs. Under the most extreme scenario, we empirically showed that our IHC-LSTM were able to remove up to 85.4% parameters from the vanilla LSTM for image generation [Gregor et al., 2015] while not affecting model performances. Similar results were also demonstrated across a variety of tasks including language modelling [Merity et al., 2018] and L2L [Andrychowicz et al., 2016].

7.2.2 On Extending Meta-Learning for Better Optimisation Algorithms

After acquiring insights on network inferences through our studies on RNNs, we redirected our focus to the meta-learning technique of L2L [Andrychowicz et al., 2016] to understand network optimisation. Andrychowicz et al.’s LSTM-optimiser replaced the gradient descent of $\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}$ (see Equation 1.4) with an LSTM. Their LSTM took the gradient of a base learner as input and returned packages of incremental updates to the base learner weights. They demonstrated that their meta-learner was more efficient than the traditional SGD in configuring a base learner. Our goal in this chapter was to mathematically conceptualise why LSTM-optimisers achieved such superior performances.

Previous Studies

Since Andrychowicz et al.’s seminal paper, most of the work in L2L aimed to address known problems that plagued the original meta-learner. More specifically, the vanilla LSTM-optimiser was notorious for its limited utility where it could only *update one specific learner for one designated dataset*. To elaborate, if a vanilla LSTM-optimiser was meta-trained to update an MLP, its learnt optimisation scheme would not be applicable for updating a ResNet. Likewise, if it was meta-trained to update an MLP on MNIST [LeCun et al., 1998], it could not be generalised to update an MLP for Cifar10 [Krizhevsky, 2009].

Several remedies were proposed by the research community. One class of study

changed the paradigm to train the meta-learner neural optimiser. They changed from learning to learn, to *reinforcement learning to learn* [Chen et al., 2017], and *reinforcement learning to reinforcement learn* [Wang et al., 2016]. Another line of research argued that the problems caused by the vanilla LSTM-optimiser were associated with the level of feature abstraction in the domain of meta-knowledge (i.e., the task structure). Hence the work of Wichrowska et al. [2017] proposed to overcome these issues via increasing the layers of the vanilla LSTM-optimiser in order to extract more detailed information to track the progress of consecutive updates in the base learner weights.

Our Novelty

Similar to our approach in the last chapter, our studies aimed to supplement existing literature through the lens of dynamical systems. Moreover, our goal was to clarify the inferential process of the L2L meta-learner to devise better optimisation algorithms. To that end, we presented a thought experiment in the [foreword](#) of the chapter to conceptualise the progress of updating base learner weights via a learnt optimiser. From the thought experiment, we found that a learnt optimiser would nearly never update the base learner weights to a global minimum. This was because that the meta-learner objective function in Andrychowicz et al. [2016] (see Equation 3.14) discouraged a meta-learner to seek for alternative loss minima upon encountering the very first minimum that it had found for the base learner. Since the base learner weights never had the chance to escape the first loss minimum, it would then never need to be subjected to the trials and errors of alternative minima. By omitting the further trials and errors, L2L could update a base learner more rapidly; but this was done at the cost of not finding a potentially better minimum of which could provide better generalisation for the base learner.

Section 5.1: A Trained RNN-optimiser for Multiple Tasks Across Multiple Datasets

A trained vanilla LSTM-optimiser suffered from *input-domain heterogeneity* and could not update a base learner on an unobserved dataset during the meta-training phase. We claimed that L2L’s root cause of input-domain heterogeneity was closely tied to the architectural design of the LSTM RNN. LSTMs were designed to transform its input data independently of the input context. Thus when a vanilla LSTM-optimiser was meta-trained to configure an MLP for MNIST, it would execute the same heuristics in test time to seek for weights to configure the MLP even if the said MLP were to be presented with the contextually different KMNIST dataset [Clanuwat et al., 2018]. Consequentially, it was not unexpected such that the incorrectly configured MLP would then fail to classify KMNIST images. To address this problem, we proposed that the LSTM neural optimiser needed to freely change its inferential process with awareness to the input data.

In order to grant an LSTM-optimiser the ability to achieve context awareness, we integrated *hypernetworks* [Ha et al., 2017] to our novel design called *multi-task learning to learn* (MTL2L). The hypernetworks were small size neural networks that changed the RNN weights of MTL2L dependently on input data. This meant that MTL2L could change its strategy to update a base learner for different datasets. An impor-

tant detail was that we compartmentalised the MTL2L RNN weights with *singular value decomposition* and only learned the singular values with the hypernetworks. This made it natural for the MTL2L optimiser to learn from multiple datasets and represent different optimisation strategies with different combinations of singular values. Furthermore, this also meant that MTL2L could synthesise new optimisation strategies to update a base learner on an unknown dataset. If the said unknown dataset were to possess shared characteristics with any one of the datasets observed during meta-training, the MTL2L optimiser could assemble different combinations of singular values to generalise updates for the base learner on the unseen dataset. Empirically, we demonstrated that not only was MTL2L capable of meta-testing under the common assumption: to meta-test on the single identical input-domain used in meta-training; but it could also update base learner networks on an unseen dataset in meta-test time.

More importantly, our work showed that [Andrychowicz et al.](#)'s LSTM-optimiser was not completely incapable of addressing input-domain heterogeneity. Though a vanilla LSTM-optimiser would eventually trigger unstable updates on an unseen dataset during meta-test time, it actually provided some correct updates to a base learner during the early stages of network configuration (see the cyan line of [Figure 5.3](#)). This meant that the vanilla LSTM-optimiser was still able to generalise its configuration strategy to unobserved datasets but with *significant impreciseness*. This observation thus explained why [Wichrowska et al. \[2017\]](#)'s approach was successful in addressing input-domain heterogeneity. In that paper, the authors extended [Andrychowicz et al.](#)'s original work with a multi-layer LSTM as the neural optimiser; the additional LSTM modules thus made it possible to refine details in the task domain thus decreasing the error rate of updating a base learner incorrectly. However, their method greatly increased the complexity of their meta-learner. In contrast, our MTL2L showed that simple architectural modifications made towards the RNN could also increase the preciseness in configuration generalisation to an unseen dataset. This prevented the meta-learner to become more complicated than the base learner thereby making L2L computationally less costly.

Section 5.2: A Light Weight Neural Optimiser that Configures and Prunes

At the end of the last section (see [Section 5.1.8](#)), we outlined two desirable properties to make L2L methods more practical in real world scenarios. A meta-learner should be easy to employ, and it should also be able to configure large size deep learning base learners. These reasons led us to switch from [Andrychowicz et al.](#)'s LSTM-optimiser to [Li et al. \[2017\]](#)'s *MetaSGD*. The LSTM-optimiser required multiple rounds of meta-training before it could learn a mapping that substituted gradient descent to update base learner weights. In contrast, the *MetaSGD* meta-learner augmented existing gradient descent schemes by learning a learning rate per parameter. That is, *MetaSGD* also used the native gradients of the base learner to configure the base learner weights. As explained in [Section 5.2.2.2](#), *MetaSGD* thus sufficed [Qian \[1999\]](#)'s conditions and behaved as a *Lyapunov function* [[Glendinning, 1994](#)] and could

properly update a base learner even when it was not meta-trained². MetaSGD was thus more readily available to be employed for meta-learning tasks.

In this section, we extended MetaSGD to expand its meta-learning capabilities beyond just rapid knowledge acquisition. In particular, we concerned the connections between meta-learning and network pruning. The reason was two-fold. On one hand, network pruning echoed with our RNN work as both studies aimed to eliminate redundant connections for network inference. On the other hand, network pruning techniques would help us to understand the minimal amount of network capacity required to solve tasks thereby increasing our knowledge for optimisation.

Most of the existing pruning techniques were based on the influential work of Han et al. [2015]. In that paper, the authors introduced a cyclic training regime which (1) trained, (2) pruned weights, and (3) retrained a dense network. Though effective, the need to retrain a dense network made this paradigm time costly. In this section, we extended MetaSGD and presented a learnt optimiser called *Masked MetaSGD* (M²SGD) to configure ResNet20s as well as nullify redundant connections on the fly. Not only did M²SGD leveraged meta-learning to rapidly configure the deep base learner, it also removed up to 83.71% synaptic connections in a non-cyclic manner. Our approach did not need to retrain deep networks and was consequently more efficient.

In addition to M²SGD, we also introduced a second neural optimiser called *Dynamic MetaSGD* (DMSGD). Unlike M²SGD which removed redundant connections, DMSGD simply kept the redundant connections unmodified. Empirically, we found that DMSGD assisted base learners to achieve slightly better generalisation than those updated by M²SGD. This was an important finding, because it showed that network redundancy wasn't completely harmful and could instead act as a natural source of regularisation for a dense network.

Our DMSGD could be linked with the *lottery ticket hypothesis* of Frankle and Carbin [2019]. The hypothesis posed a fundamental question in optimisation – *were all parameters created equally?* In that paper, a subset of subnetworks were referred as the winning tickets within a dense and randomly initialised backbone network. When trained alone, these special subnetworks could achieve competitive performances equivalent to a entirely and thoroughly fine-tuned dense backbone learner. Through our DMSGD algorithm, we demonstrated that these special subnetworks could be discovered on the fly through meta-learning. That is, we showed that optimisation could be achieved by only fine-tuning a subset of important weights in the dense base learner.

7.2.3 On Improving Continual Learning Performances

Thus far, this thesis had accomplished two important objectives. We had probed network inferential properties via our RNN studies and had gained a better understanding on network optimisations through our meta-learning work. With these

²But as we also mentioned in Section 5.2.2.2, a meta-trained MetaSGD will always outperform one that was not meta-trained.

insights, we then aimed to address the difficult AutoML problem of continual learning to help a single backbone learner to acquire skills for solving multiple tasks in a sequential fashion. Importantly, what separates our continual learning study with most existing literature is that we aim to mathematically describe the source cause of catastrophic forgetting.

Previous Studies

Catastrophic forgetting is a long-lasting problem [French, 1991]. As discussed in Section 3.3.5, most of the recent continual learning papers could be categorised as one of the three archetypes of *replay*, *regularisation*, and *dynamic architecture*. Replay methods such as Lopez-Paz and Ranzato [2017] memorise a section of data from past tasks and then show the stored data to the backbone network when it learns a new task. Regularisation approaches like Kirkpatrick et al. [2017] introduce additional learning objectives to the backbone in hope to indirectly influence the procedure of network configuration to preserve previously learnt mappings. Dynamic architectures, for instance Rusu et al. [2016], expand the backbone with more modules whenever a new task is observed for hosting new knowledge.

Each approach has its own disadvantages. The model size of the dynamic architectural backbone network scales linearly with the amount of tasks learnt. This approach thus requires high memory cost. Similarly, replay-based methods also require high memory cost for storing past data. In addition, they are also computationally costly because the external data are replayed with data of the new task for every iteration. While the regularisation archetype does not incur much memory cost and is computationally more efficient, Van de Ven and Tolias [2018] found that they could fail to mitigate forgetting under more difficult sequential learning setups. In addition, their paper showed that only replay-based techniques delivered robustly performances in their experiments.

Our Novelty

Our contribution in continual learning aimed to provide a mathematical interpretation on why forgetting occurs for sequential learning. More specifically, our goal was to develop algorithms that mitigate forgetting from both an optimisation and an inference perspective. Our understanding on optimisation enabled us to create methods to provide a backbone network with the correct updates; and that our insights on inference helped us to stop our models to forget old skills while maintaining high plasticity for new tasks.

Section 6.1: A Low Memory and Computationally Efficient Replay Method

Our first work in this chapter focused on the *episodic memory* [Lopez-Paz and Ranzato, 2017] experimental setup. While replay typically requires large amount of memory and is computationally costly, it has been greatly streamlined in the recent work of Chaudhry et al. [2019]. In that paper, the authors demonstrated that forgetting could be significantly decreased by replaying as low as one *exact* copy of past data. Their method thus made replay computationally efficient. However, they still required to

sample past data from a large size external memory. This was because that small size memories would frequently re-exposure a small set of past data to the backbone network thus causing it to overfit on old tasks while underfit on future tasks. That is, Chaudhry et al.’s method still incurred a large memory burden.

To that end, we extended Chaudhry et al.’s exact replay technique with meta-learning. We formulated continual learning as a data imbalanced problem where the backbone base learner was required to learn past tasks from a low resource environment with data stored in small size external memory units. Leveraging on meta-learning’s ability to achieve rapid knowledge acquisition, we successfully configured the backbone learner with only a handful of past data.

Our work in this section was heavily influenced by Frankle and Carbin [2019]’s lottery ticket hypothesis. It was shown in that paper that dense networks could achieve high accuracy for a task by only training a subset of their parameters. This meant that, if we were to only fine-tune a subset of all parameters for a sequential learning setup, we would then be able *reserve* network capacity to host new knowledge for later tasks. Inspired by this idea, we developed an L2L meta-learner that could continually explore for different subnetworks in a single dense backbone for sequential learning. We achieved this with our novel neural optimiser MetaSGD-CL – by *learning a learning rate per parameter per past task*, we demonstrated that exact replay with small size external memory could be achieved by hierarchically updating the backbone parameters at different timescales. We provided empirical evidence that MetaSGD-CL could help backbone base learners to achieve both higher stability and plasticity than Chaudhry et al.’s approach. MetaSGD-CL thus achieved new state-of-the-art status and addressed the two major complications in the replay archetype thereby making continual learning practical for tackling difficult real world scenarios.

Section 6.2: On Creating A Backbone Architecture with Inherently High Stability
 Replay, regularisation, and dynamic architecture all mitigate forgetting by adding external constraints on the sequential learning backbone learner. Replay methods show the backbone external data of the past task; regularisations subject the backbone to extra learning objectives; and dynamic architectures introduce external modules to host new knowledge. That is, the research community have not yet explored for solutions *within* the backbone learner to stop forgetting. To elaborate, there is currently a gap in the literature for developing neural networks with inherently high stability.

Designing a model with inherently high stability is difficult. This is because that every modification made in the forward pass results in a direct change in inference and creates an implicit alteration in optimisation. Consider *Batch Normalisation* (BN) [Ioffe and Szegedy, 2015b] which enabled neural networks to train faster by directly reshaping feature representations in inference with batch data statistics. However, it was shown in Santurkar et al. [2018] that BN also altered higher order terms in the derivative of the loss function for optimisation. This was the same for *Weight Normalisation* (WN) [Salimans and Kingma, 2016]. Though WN normalised neural network weights in the forward pass, the authors of that paper also showed that WN affected optimisation by influencing the gradient covariance matrix to be-

come closer to an identity matrix. Refer to Section 2.1.2.5 for more on normalisation techniques. Further consider residual connections [He et al., 2016a]. Residual connection changed the inference of feedforward neural networks from learning *complete transformations* to learning *incremental feature updates*. However and as discussed in Section 2.1, residual connections facilitated the training of deep models because it introduced an identity matrix to the first derivative term of the gradients between immediate layers. Designing a model with inherently high stability thus required a practitioner to have a good understanding on both inference and optimisation.

To that end, we inspected the chain rule of gradient descent to understand what type of architectures caused less parametric updates during optimisation. Leveraging on our knowledge for shortcut connections through studying RNNs, we found that highway connections [Srivastava et al., 2015] could naturally mitigate forgetting. This was because that highway connections indirectly introduced their gated units to gradient descent via the chain rule and thus scaled down parametric modifications that helped to preserve the learnt mappings for solving old tasks. Impressively, we showed that when naïvely applied for sequential learning, *MLPs with highway connections* (HCNs) incurred less forgetting than MLPs with most continual learning techniques applied. Notably, we demonstrated that highway connections were specifically required and that residual connections were not able to stop forgetting. More importantly, our HCNs was compatible to all major archetypes of continual learning techniques including replay, regularisation, and dynamic architecture.

7.3 Limitations

Our thesis used dynamical systems as the mathematical instrument to specifically analyse the process of neural network inference and optimisation. Hence, there were two layers of limitations to our work – on the *mathematical language* used, and on the *level of granularity* of our interpretation.

On the Choice of Mathematical Language

Consider continual learning, the difficult AutoML technique that requires practitioners to have a good understanding on both inference and optimisation. In this thesis, we studied inference via RNNs and optimisation with meta-learning. Furthermore, we accumulated insights on RNNs and meta-learning through dynamical systems. However, the area of dynamical systems is not the only school of mathematical ideas that could be used to probe the properties of inference and optimisation. This could also be achieved through a *Bayesian framework* and an *information theoretic* approach. In Kirkpatrick et al. [2017], the authors addressed catastrophic forgetting with a Bayesian framework and conceptualised the proper trajectories to update their backbone weights in the backward pass. Whereas in the work of Ahn et al. [2019], the authors took an information theoretic approach and gave a fresh interpretation on the Kullback–Leibler (KL) divergence [Kullback and Leibler, 1951] to interpret how parametric modifications destabilise the learnt mappings in the forward pass.

Both [Kirkpatrick et al.](#) and [Ahn et al.](#) were hence able to address new loss functions to mitigate forgetting through regularisation. It was unlikely that we could have discovered similar findings to their work with only dynamical systems.

However, there are also merits that are unique and specific to a dynamical systems approach. Consider the interpretation of the LSTM RNN mechanism, where there existed multiple gated units that collaboratively configure the cell state. A plethora of work treated RNNs as Bayesian predictive models where the cell state generated an output based on a history of observed events. Though the Bayesian narrative helped practitioners to connect LSTMs with directed graphical models in statistics, the narrative could not clarify the individual purposes of each gated unit. For one instance, [Greff et al. \[2016\]](#) observed that there were no lesions when the LSTM forget gate was bound to its input gate in their ablation study, but they were unable to provide a precise reason on why this was observed in their experiments. In contrast, our dynamic systems analysis in Section 4.2 established a link between LSTMs to dynamic Hopfield networks and our derivation revealed that a forget gate was unnecessary in the architecture given the presence of an input gate. Moreover, our analysis found that the co-existence of the forget gate and input gate was actually harmful and thus the LSTM RNN required an output gate to prevent the cell state from growing uncontrollably.

With the examples given above, we showed that the choice of mathematical language was important. It dictated the types of research that we conducted, and it influenced the ways how we interpreted our findings. While dynamical systems lacked Bayesian inference’s ability to describe non-trivial behaviours of an arbitrary black box operator, we demonstrated that they were helpful in clarifying the individual purposes of the components in a well-defined model.

On the Granularity of Feature Interpretation

Again, let us consider our dynamical systems analyses on RNNs. Our studies revealed the behaviours of each LSTM gated unit, and how to systematically remove redundant parameters from the network. However, our work only focused on the *mechanical side* of inference but did not concern the *utility side* of the propagation in neural values. For one instance, [Gould et al. \[2021\]](#)’s *deep declarative networks* studied the desiderata to define intermediate neural network layers to create more robust models that could generalise better from less training data. For another instance, [Higgins et al. \[2016\]](#)’s β -VAE addressed the learning of non-linear principal components encoded by neural networks with a bottleneck design in hope for these factors to be humanly-interpretable in the simulated data of a generative model. Since the granularity of our dynamical systems analysis was larger than the studies in intermediate layer features, our work was able to clarify the specific components that were necessary for inference but it was not capable of conceptualising the mapped relationship between two neural network layers.

In addition, our dynamical systems analyses had a granularity smaller than *neural architectural search* (NAS). Consider the influential work of [Zoph and Le \[2017\]](#), the authors of that paper first trained an RNN with reinforcement learning (RL) and then

employed the said RNN to construct a task-specific base learner in a layer-by-layer fashion. Their generated designs often included strange features that were out-of-the-box to human practitioners. For instance, their RL-RNN generated deep convolutional networks with different lengths of shortcut connections. Nonetheless, their generated base learners demonstrated competitive results on standard benchmarks. While a dynamical systems approach could help us to define necessary components in inference, it was unlikely that our work could have suggested such novel designs.

Our dynamical systems approach is still important. Though it did not comprehend the direct meanings of intermediate neural values nor did it automatically explore for novel architectures, our work highlighted the fundamental elements that must exist in a neural network. This meant that our study could serve as valuable supplements to either studies. For one instance, [Sepliarskaia et al. \[2021\]](#) noted that it was extremely difficult to measure the level of disentanglement for latent features in intermediate layers because that dense neural networks consisted an abundance of redundancies. By highlighting the necessary components with a dynamical systems approach, practitioners could potentially simplify the process of understanding feature disentanglement by solely studying the relationships among essential network components. For another instance, [Pham et al. \[2018a\]](#) pointed out that it was computationally costly to search for novel architectures with a RL-RNN. This was because that the RL-RNN would need to build several copies of candidate base learners from the scratch and separately verify for their performances. Instead of building a base learner from a clean canvas, our dynamical systems approach could highlight what must exist in a neural network architecture. That is, we could start from a simple template network that possessed all necessary components to limit the search space of the RL-RNN thereby hastening the construction procedure.

7.4 Future Work

As reflected in our limitations, our current work has a moderate level of granularity in conceptualising neural network mechanism. Thus, it will be beneficial to re-examine our current work from a higher-level perspective. Below, we will discuss how our current continual learning algorithms can be further extended.

Learning to Learn a Neural Optimiser that Scans the Loss Landscape

In Section 6.1, we introduced our novel MetaSGD-CL neural optimiser to configure a base learner for sequential learning. By learning a learning rate per parameter per past task, we showed that it was possible to mitigate forgetting in a base learner via carefully fine-tuning each weight with a different set of time scales. The current version of our MetaSGD-CL hence optimised one weight w_i independently of another weight w_j where $i \neq j$.

As a future work, we will like to craft an optimiser that updates w_i with open information on the gradients of all other w_j . While a similar research was previously conducted by [Martens and Grosse \[2015\]](#), their work only considered to use the open

information to hasten configuration. By freely sharing the status of all weights, this elevated level of optimisation perception can perhaps transform MetaSGD-CL from a tool for rapid knowledge acquisition into a device for navigating neighbouring terrains on the loss landscape. That is, we can use a history of continuous changes in base learner weights to understand whether the parametric configuration is stuck in an undesirable loss minimum. This will then help us to create an AutoML technique that investigates the early signs of sub-optimal generalisation and reset a base learner with a set of last saved weights.

A Stable Classifier with Robust Feature Representations to Parametric Corruption

In Section 6.2, we showed that neural networks classifiers with highway connections (HCNs) were highly stable and could incur less forgetting in sequential learning. The reason was because that highway connections introduced their gated unit to gradient descent via chain rule and scaled down the magnitudes of their parametric updates.

However, we argue that the gradient rescaling scheme from highway connections can be further improved. This is because that the forgetting of a base learner network can be conceptualised in many different ways. On one hand, forgetting can be measured through the amount of content changed in the base learner weights; but on another hand, forgetting can also be measured through the feature representation drift in the learnt mappings of a pre-trained network [Desai et al., 2019]. Currently, our HCN only addresses the first type of forgetting. As a future work, we will propose novel metrics to reflect different types of forgetful measures; and we will work towards creating neural networks that excel in lowering both types of forgetting.

7.5 Final Remarks

The work in this thesis used dynamical systems to analyse neural network behaviours. We demonstrated that this choice of mathematical language was helpful for us to understand inference through RNNs and to study optimisation with meta-learning. Together, our dynamical systems approach generated useful insights to assist us to develop better AutoML techniques for L2L and continual learning. Though the machine learning community has started to adopt dynamical systems into their analytical toolbox, the frontiers of machine learning applications are still mostly advanced by the Bayesian framework. Dynamical systems can offer more than first order and second order differential equation interpretations for ResNets. For one instance, bifurcation diagrams can help practitioners to understand how a high dimensional system functions when selected parameters are purposely modified [Tsaneva-Atanasova et al., 2010]. For another instance, centre manifold theorems [Krauskopf et al., 2006] can help practitioners to design smooth feature spaces with desirable topology. Hence through this thesis, we hope to encourage more people to invest in dynamical system approaches.

Bibliography

- AHN, H.; CHA, S.; LEE, D.; AND MOON, T., 2019. **Uncertainty-Based Continual Learning with Adaptive Regularisation**. In *Neural Information Processing Systems*. (cited on pages [150](#), [182](#), and [183](#))
- AKAIKE, H., 1974. **A New Look at the Statistical Model Identification**. *IEEE Transactions on Automatic Control*, (1974), 716–723. (cited on pages [18](#) and [26](#))
- AKAIKE, H., 1998. **Information Theory and an Extension of the Maximum Likelihood Principle**. In *Selected Papers of Hirotugu Akaike*, 199–213. Springer. (cited on page [26](#))
- ANDRYCHOWICZ, M.; DENIL, M.; GOMEZ, S.; HOFFMAN, M. W.; PFAU, D.; SCHAUL, T.; SHILLINGFORD, B.; AND DE FREITAS, N., 2016. **Learning to Learn by Gradient Descent by Gradient Descent**. In *the Advances in Neural Information Processing Systems*, 3981–3989. (cited on pages [1](#), [4](#), [6](#), [7](#), [9](#), [41](#), [63](#), [64](#), [65](#), [66](#), [67](#), [68](#), [89](#), [90](#), [91](#), [96](#), [103](#), [107](#), [113](#), [114](#), [117](#), [118](#), [119](#), [125](#), [126](#), [127](#), [129](#), [134](#), [137](#), [140](#), [147](#), [172](#), [176](#), [177](#), and [178](#))
- ATKINSON, K. E., 2008. *An Introduction to Numerical Analysis*. John Wiley & Sons. (cited on page [59](#))
- AYADI, I. AND TURINICI, G., 2021. **Stochastic Runge-Kutta Methods and Adaptive SGD-G2 Stochastic Gradient Descent**. In *the International Conference on Pattern Recognition*, 8220–8227. (cited on page [60](#))
- BA, J.; HINTON, G. E.; MNIH, V.; LEIBO, J. Z.; AND IONESCU, C., 2016a. **Using Fast Weights to Attend to the Recent Past**. In *the Advances in Neural Information Processing Systems*, 4331–4339. (cited on page [62](#))
- BA, J.; KIROS, J.; AND HINTON, G. E., 2016b. **Layer Normalisation**. In *the Neural Information Processing Systems Deep Learning Symposium*. (cited on pages [20](#) and [219](#))
- BAHDANAU, D.; CHO, K.; AND BENGIO, Y., 2015. **Neural machine translation by jointly learning to align and translate**. In *the International Conference on Learning Representations*. (cited on pages [5](#) and [76](#))
- BAI, S.; KOLTER, J. Z.; AND KOLTUN, V., 2018. **An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modelling**. *arXiv preprint arXiv:1803.01271*, (2018). (cited on pages [5](#), [53](#), [89](#), [118](#), and [120](#))

- BALDUZZI, D. AND GHIFARY, M., 2016. **Strongly-Typed Recurrent Neural Networks.** In *the International Conference on Machine Learning*, 1292–1300. (cited on pages [54](#), [96](#), [101](#), [102](#), [105](#), and [176](#))
- BARAKAT, A. AND BIANCHI, P., 2018. **Convergence and Dynamical Behaviour of the Adam Algorithm for Non-Convex Stochastic Optimisation.** *arXiv preprint arXiv:1810.02263*, (2018). (cited on pages [35](#) and [222](#))
- BAVELIER, D.; LEVI, D. M.; LI, R. W.; DAN, Y.; AND HENSCH, T. K., 2010. **Removing Brakes on Adult Brain Plasticity: From Molecular to Behavioural Interventions.** *Journal of Neuroscience*, (2010), 14964–14971. (cited on page [72](#))
- BAYER, J.; WIERSTRA, D.; TOGELIUS, J.; AND SCHMIDHUBER, J., 2009. **Evolving Memory Cell Structures for Sequence Learning.** In *the International Conference on Artificial Neural Networks*, 755–764. (cited on pages [50](#), [51](#), and [55](#))
- BENGIO, Y., 2009. *Learning Deep Architectures for AI*. Now Publishers Inc. (cited on page [19](#))
- BENGIO, Y.; COURVILLE, A.; AND VINCENT, P., 2013. **Representation Learning: A Review and New Perspectives.** In *the IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1798–1828. (cited on page [16](#))
- BENGIO, Y. AND DELALLEAU, O., 2011. **On the Expressive Power of Deep Architectures.** In *the International Conference on Algorithmic Learning Theory*, 18–36. (cited on page [18](#))
- BENGIO, Y.; DELALLEAU, O.; ROUX, N. L.; PAIEMENT, J.-F.; VINCENT, P.; AND OUIMET, M., 2004. **Learning Eigenfunctions Links Spectral Embedding and Kernel PCA.** *Neural Computation*, (2004), 2197–2219. (cited on page [16](#))
- BENGIO, Y.; DELALLEAU, O.; AND SIMARD, C., 2010. **Decision Trees Do Not Generalise to New Variations.** *Computational Intelligence*, (2010), 449–467. (cited on page [17](#))
- BENGIO, Y.; LAMBLIN, P.; POPOVICI, D.; AND LAROCHELLE, H., 2007. **Greedy Layer-Wise Training of Deep Networks.** In *the Advances in Neural Information Processing Systems*, 153–160. (cited on page [19](#))
- BENGIO, Y. AND MONPERRUS, M., 2005. **Non-Local Manifold Tangent Learning.** In *the Advances in Neural Information Processing Systems*, 129–136. (cited on page [16](#))
- BMBF, 2014. **Project of the Future: Industry 4.0** by the German Federal Ministry of Education and Research (Bundesministerium für Bildung und Forschung). *Internal Market, Industry, Entrepreneurship and SMEs, the European Commission*, (2014). (cited on page [1](#))
- BOEING, G., 2016. **Visual Analysis of Nonlinear Dynamical Systems: Chaos, Fractals, Self-Similarity and the Limits of Prediction.** *Systems*, (2016), 37. (cited on page [32](#))

-
- BRADBURY, J.; MERITY, S.; XIONG, C.; AND SOCHER, R., 2017. **Quasi-Recurrent Neural Networks**. In *the International Conference on Representation Learning*. (cited on page 53)
- BREIMAN, L.; FRIEDMAN, J.; STONE, C. J.; AND OLSHEN, R. A., 1984. *Classification and Regression Trees*. CRC press. (cited on page 17)
- BROMLEY, J.; BENTZ, J. W.; BOTTOU, L.; GUYON, I.; LECUN, Y.; MOORE, C.; SÄCKINGER, E.; AND SHAH, R., 1993. **Signature Verification using a “Siamese” Time Delay Neural Network**. *International Journal of Pattern Recognition and Artificial Intelligence*, (1993), 669–688. (cited on page 60)
- BROWN, P. F.; DELLA PIETRA, S. A.; DELLA PIETRA, V. J.; LAI, J. C.; AND MERCER, R. L., 1992. **An Estimate of an Upper Bound for the Entropy of English**. *Computational Linguistics*, (1992), 31–40. (cited on pages 104 and 215)
- BROWN, T.; MANN, B.; RYDER, N.; SUBBIAH, M.; KAPLAN, J. D.; DHARIWAL, P.; NEELAKANTAN, A.; SHYAM, P.; SASTRY, G.; ASKELL, A.; AGARWAL, S.; HERBERT-VOSS, A.; KRUEGER, G.; HENIGHAN, T.; CHILD, R.; RAMESH, A.; ZIEGLER, D.; WU, J.; WINTER, C.; HESSE, C.; CHEN, M.; SIGLER, E.; LITWIN, M.; GRAY, S.; CHESS, B.; CLARK, J.; BERNER, C.; MCCANDLISH, S.; RADFORD, A.; SUTSKEVER, I.; AND AMODEI, D., 2020. **Language Models are Few-Shot Learners**. In *the Advances in Neural Information Processing Systems, 1877–1901*. (cited on page 5)
- BURGER, M. AND NEUBAUER, A., 2003. **Analysis of Tikhonov Regularisation for Function Approximation by Neural Networks**. *Neural Networks*, (2003), 79–90. (cited on page 89)
- BUTCHER, J. C., 1963. **Coefficients for the Study of Runge-Kutta Integration Processes**. *Journal of the Australian Mathematical Society*, (1963), 185–201. (cited on page 37)
- CARREIRA-PERPINÁN, M. A. AND IDELBAYEV, Y., 2018. **“Learning-Compression” Algorithms for Neural Net Pruning**. In *the IEEE Conference on Computer Vision and Pattern Recognition*, 8532–8541. (cited on page 138)
- CARUANA, R., 1997. **Multitask learning**. *Machine Learning*, (1997), 41–75. (cited on page 71)
- CHA, S.; HSU, H.; CALMON, F. P.; AND MOON, T., 2020. **CPR: Classifier-Projection Regularization for Continual Learning**. In *the International Conference on Learning Representations*. (cited on page 75)
- CHANG, S.; ZHANG, Y.; HAN, W.; YU, M.; GUO, X.; TAN, W.; CUI, X.; WITBROCK, M.; HASEGAWA-JOHNSON, M. A.; AND HUANG, T. S., 2017. **Dilated Recurrent Neural Networks**. In *the Advances in Neural Information Processing Systems*, 77–87. (cited on page 89)

- CHAUDHARI, P.; CHOROMANSKA, A.; SOATTO, S.; LECUN, Y.; BALDASSI, C.; BORGS, C.; CHAYES, J.; SAGUN, L.; AND ZECCHINA, R., 2019. **Entropy-SGD: Biasing Gradient Descent into Wide Valleys**. *Journal of Statistical Mechanics: Theory and Experiment*, (2019), 124018. (cited on page 29)
- CHAUDHRY, A.; RANZATO, M.; ROHRBACH, M.; AND ELHOSEINY, M., 2018. **Efficient Lifelong Learning with A-GEM**. In *the International Conference on Learning Representations*. (cited on pages 77 and 143)
- CHAUDHRY, A.; ROHRBACH, M.; ELHOSEINY, M.; AJANTHAN, T.; DOKANIA, P. K.; TORR, P. H.; AND RANZATO, M., 2019. **On Tiny Episodic Memories in Continual Learning**. In *the International Conference on Machine Learning Workshop on Multi-Task and Lifelong Reinforcement Learning*. (cited on pages 77, 78, 141, 143, 144, 145, 148, 149, 151, 155, 180, and 181)
- CHE, Z.; PURUSHOTHAM, S.; CHO, K.; SONTAG, D.; AND LIU, Y., 2018. **Recurrent Neural Networks for Multivariate Time Series with Missing Values**. *Scientific Reports*, (2018), 1–12. (cited on page 18)
- CHEN, R. T.; RUBANOVA, Y.; BETTENCOURT, J.; AND DUVENAUD, D. K., 2018. **Neural Ordinary Differential Equations**. In *the Advances in Neural Information Processing Systems*, 6571–6583. (cited on pages 37 and 160)
- CHEN, Y.; HOFFMAN, M. W.; COLMENAREJO, S. G.; DENIL, M.; LILLICRAP, T. P.; BOTVINICK, M.; AND DE FREITAS, N., 2017. **Learning to Learn Without Gradient Descent by Gradient Descent**. In *the International Conference on Machine Learning*, 748–756. (cited on pages 126 and 177)
- CHO, K.; RAIKO, T.; ILIN, A.; AND KARHUNEN, J., 2013. **A Two-stage Pretraining Algorithm for Deep Boltzmann Machines**. In *the International Conference on Artificial Neural Networks*, 106–113. (cited on page 19)
- CHO, K.; VAN MERRIËNBOER, B.; GULCEHRE, C.; BAHDANAU, D.; BOUGARES, F.; SCHWENK, H.; AND BENGIO, Y., 2014. **Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation**. In *the Conference on Empirical Methods in Natural Language Processing*, 1724–1734. (cited on pages 8, 41, 47, 48, 49, 50, 96, 102, and 172)
- CHOPRA, S.; HADSELL, R.; AND LECUN, Y., 2005. **Learning a Similarity Metric Discriminatively, with Application to Face Verification**. In *the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 539–546. (cited on pages 60 and 61)
- CHUNG, J.; GULCEHRE, C.; CHO, K.; AND BENGIO, Y., 2014. **Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modelling**. In *the Workshop on Deep Learning of the Advances in Neural Information Processing Systems*. (cited on page 47)

-
- CLANUWAT, T.; BOBER-IRIZAR, M.; KITAMOTO, A.; LAMB, A.; YAMAMOTO, K.; AND HA, D., 2018. **Deep Learning for Classical Japanese Literature**. In *the Advances in Neural Information Processing Systems Workshop on Machine Learning for Creativity and Design*. (cited on pages [124](#), [161](#), [177](#), and [214](#))
- COATES, A.; NG, A.; AND LEE, H., 2011. **An Analysis of Single-Layer Networks in Unsupervised Feature Learning**. In *the International Conference on Artificial Intelligence and Statistics*, 215–223. (cited on pages [134](#), [140](#), and [214](#))
- COOIJMANS, T.; BALLAS, N.; LAURENT, C.; GÜLÇEHRE, Ç.; AND COURVILLE, A., 2017. **Recurrent Batch Normalisation**. In *the International Conference on Learning Representations*. (cited on page [89](#))
- DANILUK, M.; ROCKTÄSCHEL, T.; WELBL, J.; AND RIEDEL, S., 2017. **Frustratingly Short Attention Spans in Neural Language Modeling**. In *the International Conference on Machine Learning*. (cited on page [53](#))
- DAVIS, R. A.; LIU, K.-S.; AND POLITIS, D. N., 2011. **Remarks on Some Non-Parametric Estimates of a Density Function**. In *Selected Works of Murray Rosenblatt*, 95–100. Springer. (cited on page [60](#))
- DELALLEAU, O. AND BENGIO, Y., 2011. **Shallow vs. Deep Sum-Product Networks**. In *the Advances in Neural Information Processing Systems*, 666–674. (cited on page [18](#))
- DESAI, S.; ZHAN, H.; AND ALY, A., 2019. **Evaluating Lottery Tickets Under Distributional Shifts**. In *the Workshop on Deep Learning Approaches for Low-Resource NLP of the Association for Computational Linguistics*. (cited on page [185](#))
- DEVLIN, J.; CHANG, M.-W.; LEE, K.; AND TOUTANOVA, K., 2019. **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**. In *the Association for Computational Linguistics: Human Language Technologies*. (cited on page [5](#))
- DINH, L.; PASCANU, R.; BENGIO, S.; AND BENGIO, Y., 2017. **Sharp Minima can Generalize for Deep Nets**. In *the International Conference on Machine Learning*, 1019–1028. (cited on pages [28](#) and [220](#))
- DOUILLARD, A.; CORD, M.; OLLION, C.; ROBERT, T.; AND VALLE, E., 2020. **Podnet: Pooled Outputs Distillation for Small-Tasks Incremental Learning**. In *the European Conference on Computer Vision*, 86–102. (cited on pages [74](#) and [162](#))
- DRAELOS, T. J.; MINER, N. E.; LAMB, C. C.; COX, J. A.; VINEYARD, C. M.; CARLSON, K. D.; SEVERA, W. M.; JAMES, C. D.; AND AIMONE, J. B., 2017. **Neurogenesis Deep Learning: Extending Deep Networks to Accommodate New Classes**. In *the International Joint Conference on Neural Networks*, 526–533. (cited on page [78](#))
- DYRE, J. C., 2006. **Colloquium: The Glass Transition and Elastic Models of Glass-Forming Liquids**. *Reviews of Modern Physics*, (2006), 953. (cited on page [97](#))

- EHRET, B.; HENNING, C.; CERVERA, M. R.; MEULEMANS, A.; VON OSWALD, J.; AND GREWE, B. F., 2021. **Continual Learning in Recurrent Neural Networks with Hypernetworks**. In *the International Conference on Learning Representations*. (cited on page 131)
- ELMAN, J. L., 1990. **Finding Structure in Time**. *Cognitive Science*, (1990), 179–211. (cited on page 43)
- ELSKEN, T.; METZEN, J. H.; AND HUTTER, F., 2019. **Neural Architecture Search: A Survey**. *Journal of Machine Learning Research*, (2019), 1–21. (cited on pages 2 and 127)
- FEURER, M. AND HUTTER, F., 2019. **Hyperparameter Optimisation**. In *Hutter et al. [2019]*, 3–38. (cited on page 2)
- FINN, C.; ABBEEL, P.; AND LEVINE, S., 2017. **Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks**. In *the International Conference on Machine Learning*, 1126–1135. (cited on pages 29, 59, and 155)
- FITZHUGH, R., 1955. **Mathematical Models of Threshold Phenomena in the Nerve Membrane**. *Bulletin of Mathematical Biophysics*, (1955), 257–278. (cited on page 37)
- FIX, E., 1951. *Discriminatory analysis: nonparametric discrimination, consistency properties*. USAF School of Aviation Medicine. (cited on page 60)
- FOX, C. W. AND ROBERTS, S. J., 2012. **A Tutorial on Variational Bayesian Inference**. *Artificial Intelligence Review*, (2012), 85–95. (cited on page 17)
- FRANKLE, J. AND CARBIN, M., 2019. **The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks**. In *the International Conference on Learning Representations*. (cited on pages 138, 179, and 181)
- FRENCH, R. M., 1991. **Using Semi-Distributed Representations to Overcome Catastrophic Forgetting in Connectionist Networks**. In *the Annual Cognitive Science Society Conference*, 173–178. (cited on pages 76, 154, and 180)
- FRENCH, R. M., 1997. **Pseudo-Recurrent Connectionist Networks: An Approach to the ‘Sensitivity-Stability’ Dilemma**. *Connection Science*, (1997), 353–380. (cited on page 77)
- FRENCH, R. M., 1999. **Catastrophic Forgetting in Connectionist Networks**. *Trends in Cognitive Sciences*, (1999), 128–135. (cited on page 71)
- GAROFOLO, J. S.; LAMEL, L. F.; FISHER, W. M.; FISCUS, J. G.; AND PALLETT, D. S., 1993. **DARPA TIMIT Acoustic-Phonetic Continuous Speech Corpus CD-ROM. NIST Speech Disc 1-1.1**. *STIN*, (1993), 27403. (cited on page 28)
- GASTALDI, X., 2017. **Shake-Shake Regularisation**. In *the Workshop of International Conference on Learning Representations*. (cited on pages 37 and 223)

-
- GATYS, L. A.; ECKER, A. S.; AND BETHGE, M., 2016. **Image Style Transfer using Convolutional Neural Networks**. In *the IEEE Conference on Computer Vision and Pattern Recognition*, 2414–2423. (cited on page 21)
- GE, R.; HUANG, F.; JIN, C.; AND YUAN, Y., 2015. **Escaping from Saddle Points — On-line Stochastic Gradient for Tensor Decomposition**. In *the Conference on Learning Theory*, 797–842. (cited on page 27)
- GERS, F. A. AND SCHMIDHUBER, J., 2000. **Recurrent Nets that Time and Count**. In *the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, 189–194. (cited on page 47)
- GERS, F. A.; URGEN SCHMIDHUBER, J.; AND CUMMINS, F., 1999. **Learning to Forget: Continual Prediction with LSTM**. *Technical Report of IDSIA*, (1999). (cited on pages 46, 47, 49, and 88)
- GLENDINNING, P., 1994. *Stability, Instability, and Chaos: An Introduction to the Theory of Nonlinear Differential Equations*, vol. 11. Cambridge University Press. (cited on pages 6, 31, 35, 83, 97, 109, 120, 131, 158, 171, 175, and 178)
- GLOROT, X. AND BENGIO, Y., 2010. **Understanding the Difficulty of Training Deep Feedforward Neural Networks**. In *the International Conference on Artificial Intelligence and Statistics*, 249–256. (cited on pages 19, 133, and 218)
- GOLMANT, N.; VEMURI, N.; YAO, Z.; FEINBERG, V.; GHOLAMI, A.; ROTHAUGE, K.; MAHONEY, M. W.; AND GONZALEZ, J., 2018. **On the Computational Inefficiency of Large Batch Sizes for Stochastic Gradient Descent**. *arXiv preprint arXiv:1811.12941*, (2018). (cited on page 29)
- GOODFELLOW, I.; BENGIO, Y.; COURVILLE, A.; AND BENGIO, Y., 2016. *Deep learning*. MIT press Cambridge. (cited on page 213)
- GOODFELLOW, I.; LEE, H.; LE, Q. V.; SAXE, A.; AND NG, A. Y., 2009. **Measuring Invariances in Deep Networks**. In *the Advances in Neural Information Processing Systems*, 646–654. (cited on page 18)
- GOODFELLOW, I. J.; BULATOV, Y.; IBARZ, J.; ARNOUD, S.; AND SHET, V., 2014. **Multi-Digit Number Recognition from Street View Imagery using Deep Convolutional Neural Networks**. In *the International Conference on Learning Representations*. (cited on page 106)
- GOULD, S.; HARTLEY, R.; AND CAMPBELL, D. J., 2021. **Deep Declarative Networks**. (cited on page 183)
- GRANT, E.; FINN, C.; LEVINE, S.; DARRELL, T.; AND GRIFFITHS, T., 2018. **Recasting Gradient-Based Meta-Learning as Hierarchical Bayes**. In *the International Conference on Learning Representations*. (cited on page 59)

-
- GRAVES, A., 2013. **Generating Sequences with Recurrent Neural Networks**. *arXiv preprint arXiv:1308.0850*, (2013). (cited on page 217)
- GRAVES, A.; WAYNE, G.; AND DANIHELKA, I., 2014. **Neural Turing Machines**. *arXiv preprint arXiv:1410.5401*, (2014). (cited on page 63)
- GREFF, K.; SRIVASTAVA, R. K.; KOUTNÍK, J.; STEUNEBRINK, B. R.; AND SCHMIDHUBER, J., 2016. **LSTM: A Search Space Odyssey**. In *the IEEE Transactions on Neural Networks and Learning Systems*, 2222–2232. (cited on pages 6, 50, 52, 55, 72, 107, 108, 173, 175, 176, and 183)
- GREGOR, K.; DANIHELKA, I.; GRAVES, A.; REZENDE, D.; AND WIERSTRA, D., 2015. **DRAW: A Recurrent Neural Network For Image Generation**. In *the International Conference on Machine Learning*, 1462–1471. (cited on pages 6, 8, 96, 103, 106, and 176)
- GURBUZBALABAN, M.; SIMSEKLI, U.; AND ZHU, L., 2021. **The Heavy-Tail Phenomenon in SGD**. In *the International Conference on Machine Learning*, 3964–3975. (cited on page 28)
- HA, D.; DAI, A.; AND LE, Q. V., 2017. **Hypernetworks**. In *the International Conference on Learning Representations*. (cited on pages 46, 61, 62, 120, 140, and 177)
- HAN, S.; POOL, J.; TRAN, J.; AND DALLY, W., 2015. **Learning Both Weights and Connections for Efficient Neural Network**. In *the Advances in Neural Information Processing Systems*, 1135–1143. (cited on pages 9, 129, 137, 138, 172, and 179)
- HASSIBI, B. AND STORK, D. G., 1993. **Second Order Derivatives for Network Pruning: Optimal Brain Surgeon**. In *the Advances in Neural Information Processing Systems*, 164–171. (cited on page 27)
- HASTIE, T.; TIBSHIRANI, R.; AND FRIEDMAN, J., 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Science & Business Media. (cited on page 18)
- HAYES, T. L.; KRISHNAN, G. P.; BAZHENOV, M.; SIEGELMANN, H. T.; SEJNOWSKI, T. J.; AND KANAN, C., 2021. **Replay in Deep Learning: Current Approaches and Missing Biological Elements**. *Neural Computation*, (2021). (cited on page 72)
- HE, K.; ZHANG, X.; REN, S.; AND SUN, J., 2015. **Delving Deep into Rectifiers: Surpassing Human-Level Performance on Imagenet Classification**. In *the IEEE International Conference on Computer Vision*, 1026–1034. (cited on pages 19 and 218)
- HE, K.; ZHANG, X.; REN, S.; AND SUN, J., 2016a. **Deep Residual Learning for Image Recognition**. In *the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778. (cited on pages 3, 5, 14, 15, 20, 21, 29, 31, 34, 36, 39, 71, 100, 107, 118, 129, 131, 134, 140, 148, 160, 163, 167, 171, and 182)

-
- HE, K.; ZHANG, X.; REN, S.; AND SUN, J., 2016b. **Identity Mappings in Deep Residual Networks**. In *the European Conference on Computer Vision*, 630–645. (cited on page [36](#))
- HE, Z.; GAO, S.; XIAO, L.; LIU, D.; HE, H.; AND BARBER, D., 2017. **Wider and Deeper, Cheaper and Faster: Tensorised LSTMs for Sequence Learning**. In *the Advances in Neural Information Processing Systems*, 1–11. (cited on page [89](#))
- HEAFIELD, K.; POUZYREVSKY, I.; CLARK, J. H.; AND KOEHN, P., 2013. **Scalable Modified Kneser-Ney Language Model Estimation**. In *the Association for Computational Linguistics*, 690–696. (cited on page [52](#))
- HENSCH, T. K., 2004. **Critical Period Regulation**. *Annu. Rev. Neurosci.*, (2004), 549–579. (cited on page [72](#))
- HIGGINS, I.; MATTHEY, L.; PAL, A.; BURGESS, C.; GLOROT, X.; BOTVINICK, M.; MOHAMED, S.; AND LERCHNER, A., 2016. **Beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework**. In *the International Conference on Learning Representations*. (cited on page [183](#))
- HINTON, G.; VINYALS, O.; AND DEAN, J., 2015. **Distilling the Knowledge in a Neural Network**. *University of Toronto, Geoffrey E. Hinton Abstracts*, (2015). (cited on pages [19](#), [72](#), and [75](#))
- HINTON, G. E., 2009. **Deep Belief Networks**. *Scholarpedia*, (2009), 5947. (cited on page [19](#))
- HINTON, G. E.; OSINDERO, S.; AND TEH, Y.-W., 2006. **A Fast Learning Algorithm for Deep Belief Nets**. *Neural Computation*, (2006), 1527–1554. (cited on page [19](#))
- HINTON, G. E.; SRIVASTAVA, N.; KRIZHEVSKY, A.; SUTSKEVER, I.; AND SALAKHUTDINOV, R. R., 2012. **Improving Neural Networks by Preventing Co-Adaptation of Feature Detectors**. *University of Toronto Computer Science*, (2012). (cited on page [168](#))
- HOCHREITER, S. AND SCHMIDHUBER, J., 1995. **Simplifying Neural Nets by Discovering Flat Minima**. In *the Advances in Neural Information Processing Systems*, 529–536. (cited on pages [26](#), [27](#), [28](#), [29](#), and [220](#))
- HOCHREITER, S. AND SCHMIDHUBER, J., 1997a. **Flat Minima**. *Neural Computation*, (1997), 1–42. (cited on pages [23](#) and [171](#))
- HOCHREITER, S. AND SCHMIDHUBER, J., 1997b. **Long Short-Term Memory**. *Neural Computation*, (1997), 1735–1780. (cited on pages [1](#), [6](#), [7](#), [16](#), [39](#), [41](#), [44](#), [45](#), [48](#), [49](#), [53](#), [61](#), [64](#), [72](#), [81](#), [102](#), [111](#), [113](#), [129](#), [161](#), [171](#), and [217](#))
- HOCHREITER, S.; YOUNGER, A. S.; AND CONWELL, P. R., 2001. **Learning to Learn Using Gradient Descent**. In *the International Conference on Artificial Neural Networks*, 87–94. (cited on page [62](#))

- HODGKIN, A. L. AND HUXLEY, A. F., 1952. **Currents Carried by Sodium and Potassium Ions through the Membrane of the Giant Axon of Loligo.** *Journal of Physiology*, (1952), 449. (cited on page 37)
- HOPFIELD, J. J., 1982. **Neural Networks and Physical Systems with Emergent Collective Computational Abilities.** In *the Proceedings of the National Academy of Sciences*, 2554–2558. (cited on pages 8, 96, 97, and 175)
- HOSPEDALES, T.; ANTONIOU, A.; MICAELLI, P.; AND STORKEY, A., 2020. **Meta-Learning in Neural Networks: A Survey.** *arXiv preprint arXiv:2004.05439*, (2020). (cited on page 59)
- HOWARD, J. AND RUDER, S., 2018. **Universal Language Model Fine-tuning for Text Classification.** In *the Association for Computational Linguistics*, 328–339. (cited on page 131)
- HUANG, G.; SUN, Y.; LIU, Z.; SEDRA, D.; AND WEINBERGER, K. Q., 2016. **Deep networks with stochastic depth.** In *the European Conference on Computer Vision*, 646–661. (cited on pages 36 and 37)
- HUTTER, F.; KOTTHOFF, L.; AND VANSCHOREN, J., 2019. **Automated Machine Learning: Methods, Systems, Challenges.** Springer Nature. (cited on pages 1, 2, 127, 171, 192, and 209)
- IOFFE, S. AND SZEGEDY, C., 2015a. **Batch Normalisation: Accelerating Deep Network Training by Reducing Internal Covariate Shift.** In *the International Conference on Machine Learning*, 448–456. (cited on pages 19 and 218)
- IOFFE, S. AND SZEGEDY, C., 2015b. **Batch Normalisation: Accelerating Deep Network Training by Reducing Internal Covariate Shift.** In *the International Conference on Machine Learning*, 448–456. (cited on pages 90 and 181)
- ISHIKAWA, M., 1996. **Structural Learning with Forgetting.** *Neural Networks*, (1996), 509–521. (cited on page 138)
- JAEGER, H., 2001. **The “Echo State” Approach to Analysing and Training Recurrent Neural Networks – with an Erratum Note.** *Bonn, Germany: German National Research Centre for Information Technology GMD Technical Report*, (2001), 13. (cited on pages 55, 110, and 175)
- JASTRZEBSKI, S.; KENTON, Z.; ARPIT, D.; BALLAS, N.; FISCHER, A.; BENGIO, Y.; AND STORKEY, A., 2017. **Three Factors Influencing Minima in SGD.** *arXiv preprint arXiv:1711.04623*, (2017). (cited on page 28)
- JASTRZEBSKI, S.; KENTON, Z.; BALLAS, N.; FISCHER, A.; BENGIO, Y.; AND STORKEY, A., 2019. **On the Relation between the Sharpest Directions of DNN Loss and the SGD Step Length.** In *International Conference on Learning Representations*. (cited on pages 28 and 220)

-
- JAVED, K. AND WHITE, M., 2019. **Meta-Learning Representations for Continual Learning**. In *the Advances in Neural Information Processing Systems*. (cited on page 155)
- JOZEFOWICZ, R.; ZAREMBA, W.; AND SUTSKEVER, I., 2015. **An Empirical Exploration of Recurrent Network Architectures**. In *the International Conference on Machine Learning*, 2342–2350. (cited on pages 6, 50, 51, 55, 108, and 173)
- JUNG, H.; JU, J.; JUNG, M.; AND KIM, J., 2016. **Less-Forgetting Learning in Deep Neural Networks**. *arXiv preprint arXiv:1607.00122*, (2016). (cited on page 75)
- KAMRA, N.; GUPTA, U.; AND LIU, Y., 2017. **Deep Generative Dual Memory Network for Continual Learning**. *arXiv preprint arXiv:1710.10368*, (2017). (cited on pages 77 and 78)
- KANG, J.; ZHANG, W.-Q.; LIU, W.-W.; LIU, J.; AND JOHNSON, M. T., 2018. **Advanced Recurrent Network-Based Hybrid Acoustic Models for Low Resource Speech Recognition**. *EURASIP Journal on Audio, Speech, and Music Processing*, (2018), 6. (cited on page 47)
- KARPATHY, A.; JOHNSON, J.; AND FEI-FEI, L., 2016. **Visualising and Understanding Recurrent Networks**. In *the International Conference on Learning Representations*. (cited on pages 6, 52, 53, 54, 55, and 173)
- KEENER, J. P. AND SNEYD, J., 1998. *Mathematical Physiology*. Springer. (cited on pages 6 and 31)
- KEMKER, R. AND KANAN, C., 2018. **FearNet: Brain-Inspired Model for Incremental Learning**. In *the International Conference on Learning Representations*. (cited on page 76)
- KEMKER, R.; MCCLURE, M.; ABITINO, A.; HAYES, T.; AND KANAN, C., 2018. **Measuring Catastrophic Forgetting in Neural Networks**. In *the AAAI Conference on Artificial Intelligence*. (cited on page 71)
- KESKAR, N. S.; MUDIGERE, D.; NOCEDAL, J.; SMELYANSKIY, M.; AND TANG, P. T. P., 2017. **On Large-Batch Training for Deep Learning: Generalisation Gap and Sharp Minima**. In *the International Conference on Learning Representations*. (cited on pages 27, 28, and 220)
- KIM, J.; EL-KHAMY, M.; AND LEE, J., 2017. **Residual LSTM: Design of a Deep Recurrent Architecture for Distant Speech Recognition**. In *the International Speech Communication Association*. (cited on page 54)
- KINGMA, D. P. AND BA, J., 2015. **Adam: A Method for Stochastic Optimization**. In *the International Conference on Learning Representations*. (cited on pages 35, 64, and 221)
- KINGMA, D. P. AND WELING, M., 2014. **Auto-Encoding Variational Bayes**. In *the International Conference on Learning Representations*. (cited on page 17)

- KIRKPATRICK, J.; PASCANU, R.; RABINOWITZ, N.; VENESS, J.; DESJARDINS, G.; RUSU, A. A.; MILAN, K.; QUAN, J.; RAMALHO, T.; GRABSKA-BARWINSKA, A.; HASSABIS, D.; CLOPATH, C.; KUMARAN, D.; AND HADSELL, R., 2017. **Overcoming Catastrophic Forgetting in Neural Networks**. In *the Proceedings of the National Academy of Sciences*, 3521–3526. (cited on pages [10](#), [75](#), [143](#), [148](#), [149](#), [155](#), [162](#), [164](#), [168](#), [180](#), [182](#), and [183](#))
- KOCH, G.; ZEMEL, R.; AND SALAKHUTDINOV, R., 2015. **Siamese Neural Networks for One-Shot Image Recognition**. In *the International Conference on Machine Learning Deep Learning Workshop*. (cited on page [60](#))
- KOUTNIK, J.; GREFF, K.; GOMEZ, F.; AND SCHMIDHUBER, J., 2014. **A Clockwork RNN**. In *the International Conference on Machine Learning*, 1863–1871. (cited on pages [55](#) and [110](#))
- KRAUSKOPF, B.; OSINGA, H. M.; DOEDEL, E. J.; HENDERSON, M. E.; GUCKENHEIMER, J.; VLADIMIRSKY, A.; DELLNITZ, M.; AND JUNGE, O., 2006. **A Survey of Methods for Computing (Un)Stable Manifolds of Vector Fields**. In *Modelling And Computations In Dynamical Systems: In Commemoration of the 100th Anniversary of the Birth of John von Neumann*, 67–95. World Scientific. (cited on page [185](#))
- KRIZHEVSKY, A., 2009. **Learning Multiple Layers of Features from Tiny Images**. In *Tech Report of the University of Toronto*. (cited on pages [28](#), [103](#), [107](#), [119](#), [124](#), [162](#), [176](#), and [214](#))
- KRIZHEVSKY, A.; SUTSKEVER, I.; AND HINTON, G. E., 2012. **Imagenet Classification with Deep Convolutional Neural Networks**. In *the Advances in Neural Information Processing Systems*, 1097–1105. (cited on pages [1](#), [5](#), [16](#), and [71](#))
- KRUEGER, D.; MAHARAJ, T.; KRAMÁR, J.; PEZESHKI, M.; BALLAS, N.; KE, N. R.; GOYAL, A.; BENGIO, Y.; COURVILLE, A.; AND PAL, C., 2017. **Zoneout: Regularising RNNs by Randomly Preserving Hidden Activations**. In *the International Conference on Learning Representations*. (cited on pages [88](#) and [89](#))
- KUCHAIEV, O. AND GINSBURG, B., 2017. **Factorisation Tricks for LSTM Networks**. In *the Workshop Track of International Conference on Representation Learning*. (cited on page [54](#))
- KULLBACK, S. AND LEIBLER, R. A., 1951. **On Information and Sufficiency**. *The Annals of Mathematical Statistics*, 22, 1 (1951), 79–86. (cited on page [182](#))
- KUNITA, H., 2010. **Itô’s Stochastic Calculus: Its Surprising Power for Applications**. *Stochastic Processes and their Applications*, (2010), 622–652. (cited on page [37](#))
- KUO, N.; HARANDI, M.; FOURRIER, N.; WALDER, C.; FERRARO, G.; AND SUOMINEN, H., 2020a. **M2SGD: Learning to Learn Important Weights**. In *the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop on Continual Learning in Computer Vision*, 957–964. (cited on pages [9](#), [113](#), and [129](#))

-
- KUO, N.; HARANDI, M.; FOURRIER, N.; WALDER, C.; FERRARO, G.; AND SUOMINEN, H., 2021a. **Plastic and Stable Gated Classifiers for Continual Learning**. In *the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop on Continual Learning in Computer Vision*. (cited on pages 10, 141, and 158)
- KUO, N. I., 2017. **Understanding and Modifying Dynamical Hopfield Neural Networks for Generating Multiple Coherent Patterns**. *Masters Thesis of the University of Auckland*, (2017). (cited on page 175)
- KUO, N. I.; HARANDI, M.; FOURRIER, N.; WALDER, C.; FERRARO, G.; AND SUOMINEN, H., 2020b. **An Input Residual Connection for Simplifying Gated Recurrent Neural Networks**. In *the International Joint Conference on Neural Networks*, 1–8. (cited on pages 8, 81, 96, 102, and 110)
- KUO, N. I.; HARANDI, M.; FOURRIER, N.; WALDER, C.; FERRARO, G.; SUOMINEN, H.; ET AL., 2020c. **MTL2L: A Context Aware Neural Optimiser**. In *the Workshop on Automated Machine Learning of the International Conference on Machine Learning*. (cited on pages 9, 113, and 117)
- KUO, N. I.; HARANDI, M. T.; SUOMINEN, H.; FOURRIER, N.; WALDER, C.; AND FERRARO, G., 2018. **DecayNet: A Study on the Cell States of Long Short Term Memories**. (2018). (cited on pages 8, 81, 83, 90, and 110)
- KUO, N. I.; HARANDI, M. T.; SUOMINEN, H.; FOURRIER, N.; WALDER, C.; AND FERRARO, G., 2021b. **Learning to Continually Learn Rapidly from Few and Noisy Data**. In *the Meta-Learning and Co-Hosted Competition of the AAAI Conference on Artificial Intelligence*. (cited on pages 10, 141, and 143)
- KURTZ, M.; KOPINSKY, J.; GELASHVILI, R.; MATVEEV, A.; CARR, J.; GOIN, M.; LEISERSON, W.; MOORE, S.; NELL, B.; SHAVIT, N.; AND ALISTARH, D., 2020. **Inducing and Exploiting Activation Sparsity for Fast Neural Network Inference**. In *the International Conference on Machine Learning*. (cited on pages 18 and 168)
- KUSUPATI, A.; SINGH, M.; BHATIA, K.; KUMAR, A.; JAIN, P.; AND VARMA, M., 2018. **FastGRNN: A Fast, Accurate, Stable and Tiny Kilobyte Sized Gated Recurrent Neural Network**. In *the Advances in Neural Information Processing Systems*, 9017–9028. (cited on pages 54, 96, 101, 102, 105, and 176)
- LA SALLE, J.; LEFSCHETZ, S.; AND ALVERSON, R., 1962. **Stability by Lyapunov’s Direct Method with Applications**. *PhT*, (1962), 59. (cited on page 35)
- LAMME, V. A., 2006. **Towards a True Neural Stance on Consciousness**. *Trends in Cognitive Sciences*, (2006), 494–501. (cited on page 42)
- LANCZOS, C., 1950. **An Iteration Method for the Solution of the Eigenvalue Problem of Linear Differential and Integral Operators**. United States Governm. Press Office Los Angeles, CA. (cited on page 28)

- LANG, K. J.; WAIBEL, A. H.; AND HINTON, G. E., 1990. **A Time-Delay Neural Network Architecture for Isolated Word Recognition.** *Neural Networks*, (1990), 23–43. (cited on page 60)
- LARSSON, G.; MAIRE, M.; AND SHAKHAROVICH, G., 2017. **FractalNet: Ultra-Deep Neural Networks without Residuals.** In *the International Conference on Learning Representations*. (cited on page 37)
- LE, Q. V.; JAITLEY, N.; AND HINTON, G. E., 2015. **A Simple Way to Initialise Recurrent Networks of Rectified Linear Units.** *arXiv preprint arXiv:1504.00941*, (2015). (cited on pages 55, 83, 89, and 174)
- LECUN, Y.; BOTTOU, L.; BENGIO, Y.; AND HAFFNER, P., 1998. **Gradient-Based Learning Applied to Document Recognition.** *Proceedings of the IEEE*, (1998), 2278–2324. (cited on pages 89, 119, 124, 134, 161, 162, 176, and 213)
- LECUN, Y.; DENKER, J. S.; AND SOLLA, S. A., 1990. **Optimal Brain Damage.** In *the Advances in Neural Information Processing Systems*, 598–605. (cited on pages 138 and 139)
- LECUN, Y. A.; BOTTOU, L.; ORR, G. B.; AND MÜLLER, K.-R., 2012. **Efficient Backprop.** In *Neural Networks: Tricks of the Trade*, 9–48. Springer. (cited on page 27)
- LEE, C.-Y.; XIE, S.; GALLAGHER, P.; ZHANG, Z.; AND TU, Z., 2015. **Deeply-Supervised Nets.** In *the International Conference on Artificial Intelligence and Statistics*, 562–570. (cited on page 19)
- LEGG, S. AND HUTTER, M., 2007. **Universal Intelligence: A Definition of Machine Intelligence.** *Minds and Machines*, (2007), 391–444. (cited on pages 10 and 140)
- LEI, T.; ZHANG, Y.; WANG, S. I.; DAI, H.; AND ARTZI, Y., 2018a. **Simple Recurrent Units for Highly Parallelisable Recurrence.** In *the Conference on Empirical Methods in Natural Language Processing*, 4470–4481. (cited on pages 8, 96, 101, 102, 105, 120, and 176)
- LEI, T.; ZHANG, Y.; WANG, S. I.; DAI, H.; AND ARTZI, Y., 2018b. **Simple Recurrent Units for Highly Parallelisable Recurrence.** In *the Conference on Empirical Methods in Natural Language Processing*, 4470–4481. (cited on page 53)
- LI, H.; XU, Z.; TAYLOR, G.; STUDER, C.; AND GOLDSTEIN, T., 2018. **Visualising the Loss Landscape of Neural Nets.** In *the Advances in Neural Information Processing Systems*, 6389–6399. (cited on pages 29 and 221)
- LI, Z. AND HOIEM, D., 2017. **Learning without Forgetting.** In *the IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2935–2947. (cited on pages 72, 75, and 155)

-
- LI, Z.; ZHOU, F.; CHEN, F.; AND LI, H., 2017. **Meta-SGD: Learning to Learn Quickly for Few-Shot Learning.** *arXiv preprint arXiv:1707.09835*, (2017). (cited on pages [9](#), [126](#), [129](#), [130](#), [136](#), [137](#), [140](#), [144](#), [145](#), and [178](#))
- LIM, J. J.; SALAKHUTDINOV, R. R.; AND TORRALBA, A., 2011. **Transfer Learning by Borrowing Examples for Multiclass Object Detection.** In *the Advances in Neural Information Processing Systems*, 118–126. (cited on page [60](#))
- LIN, M.; CHEN, Q.; AND YAN, S., 2014. **Network in Network.** In *the International Conference on Learning Representations*. (cited on page [38](#))
- LING, Z. AND QIU, R. C., 2019. **Spectrum Concentration in Deep Residual Learning: A Free Probability Approach.** *IEEE Access*, (2019), 105212–105223. (cited on page [15](#))
- LOMONACO, V. AND MALTONI, D., 2017. **Core50: A New Dataset and Benchmark for Continuous Object Recognition.** In *the Conference on Robot Learning*, 17–26. (cited on pages [73](#) and [76](#))
- LOMONACO, V.; PELLEGRINI, L.; RODRIGUEZ, P.; CACCIA, M.; SHE, Q.; CHEN, Y.; JODELET, Q.; WANG, R.; MAI, Z.; VAZQUEZ, D.; ET AL., 2020. **CVPR 2020 Continual Learning in Computer Vision Competition: Approaches, Results, Current Challenges and Future Directions.** *arXiv preprint arXiv:2009.09929*, (2020). (cited on pages [71](#), [73](#), and [78](#))
- LOPEZ-PAZ, D. AND RANZATO, M., 2017. **Gradient Episodic Memory for Continual Learning.** In *the Advances in Neural Information Processing Systems*, 6467–6476. (cited on pages [1](#), [4](#), [7](#), [41](#), [77](#), [78](#), [140](#), [141](#), [143](#), [144](#), [148](#), [149](#), [163](#), [164](#), [168](#), [172](#), [173](#), and [180](#))
- LU, Y.; ZHONG, A.; LI, Q.; AND DONG, B., 2018. **Beyond Finite Layer Neural Networks: Bridging Deep Architectures and Numerical Differential Equations.** In *the International Conference on Machine Learning*, 3276–3285. (cited on page [37](#))
- LUKOŠEVIČIUS, M., 2012. **A Practical Guide to Applying Echo State Networks.** In *Neural Networks: Tricks of The Trade*, 659–686. Springer. (cited on page [98](#))
- LUO, Y., 2017. **Recurrent Neural Networks for Classifying Relations in Clinical Notes.** *Journal of Biomedical Informatics*, (2017), 85–95. (cited on page [125](#))
- LV, K.; JIANG, S.; AND LI, J., 2017. **Learning Gradient Descent: Better Generalisation and Longer Horizons.** In *the International Conference on Machine Learning*. (cited on page [125](#))
- MAI, Z.; KIM, H.; JEONG, J.; AND SANNER, S., 2020. **Batch-level Experience Replay with Review for Continual Learning.** In *the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshop on Continual Learning in Computer Vision*. (cited on pages [148](#), [151](#), and [163](#))

- MARCUS, M.; SANTORINI, B.; AND MARCINKIEWICZ, M. A., 1993. **Building a Large Annotated Corpus of English: The Penn Treebank.** *Computational Linguistics*, (1993), 313–330. (cited on pages [104](#) and [215](#))
- MARTENS, J. AND GROSSE, R., 2015. **Optimising Neural Networks with Kronecker-Factored Approximate Curvature.** In *the International Conference on Machine Learning*, 2408–2417. (cited on page [184](#))
- MCCLELLAND, J. L.; MCNAUGHTON, B. L.; AND O'REILLY, R. C., 1995. **Why there are Complementary Learning Systems in the Hippocampus and Neocortex: Insights from the Successes and Failures of Connectionist Models of Learning and Memory.** *Psychological Review*, (1995), 419. (cited on pages [76](#) and [78](#))
- MCCLOSKEY, M. AND COHEN, N. J., 1989. **Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem.** In *Psychology of Learning and Motivation*, 109–165. (cited on pages [10](#), [71](#), [141](#), and [172](#))
- MERITY, S.; KESKAR, N. S.; AND SOCHER, R., 2018. **Regularising and Optimising LSTM Language Models.** In *the International Conference on Learning Representations*. (cited on pages [6](#), [50](#), [76](#), [96](#), [103](#), [105](#), and [176](#))
- MERMILLOD, M.; BUGAJSKA, A.; AND BONIN, P., 2013. **The Stability-Plasticity Dilemma: Investigating the Continuum from Catastrophic Forgetting to Age-Limited Learning Effects.** *Frontiers in Psychology*, (2013), 504. (cited on pages [72](#) and [172](#))
- MIAO, Y.; LI, J.; WANG, Y.; ZHANG, S.-X.; AND GONG, Y., 2016. **Simplifying Long Short-Term Memory Acoustic Models for Fast Training and Decoding.** In *the IEEE International Conference on Acoustics, Speech and Signal Processing*, 2284–2288. (cited on page [54](#))
- MIRZADEH, S. I.; FARAJTABAR, M.; AND GHASEMZADEH, H., 2020a. **Dropout as an Implicit Gating Mechanism for Continual Learning.** In *the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops of Continual Learning in Vision*. (cited on page [168](#))
- MIRZADEH, S. I.; FARAJTABAR, M.; PASCANU, R.; AND GHASEMZADEH, H., 2020b. **Understanding the Role of Training Regimes in Continual Learning.** In *the International Conference on Machine Learning Workshop on Continual Learning*. (cited on pages [73](#) and [164](#))
- MONTUFAR, G. F.; PASCANU, R.; CHO, K.; AND BENGIO, Y., 2014. **On the Number of Linear Regions of Deep Neural Networks.** In *the Advances in Neural Information Processing Systems*, 2924–2932. (cited on page [18](#))
- MUNKHDALAI, T. AND YU, H., 2017. **Meta Networks.** *Proceedings of Machine Learning Research*, (2017), 2554. (cited on page [63](#))

-
- NAGUMO, J.; ARIMOTO, S.; AND YOSHIZAWA, S., 1962. **An Active Pulse Transmission Line Simulating Nerve Axon.** *Proceedings of the IRE*, (1962), 2061–2070. (cited on page [37](#))
- NAIR, V. AND HINTON, G. E., 2010. **Rectified Linear Units Improve Restricted Boltzmann Machines.** In *the International Conference on Machine Learning*. (cited on page [18](#))
- NAKKIRAN, P.; KAPLUN, G.; BANSAL, Y.; YANG, T.; BARAK, B.; AND SUTSKEVER, I., 2020. **Deep Double Descent: Where Bigger Models and More Data Hurt.** In *the International Conference on Learning Representations*. (cited on pages [17](#) and [217](#))
- NARANG, S.; ELSÉN, E.; DIAMOS, G.; AND SENGUPTA, S., 2017. **Exploring Sparsity in Recurrent Neural Networks.** In *the International Conference on Learning Representations*. (cited on page [138](#))
- NEISHTADT, A. I., 1984. **The Separation of Motions in Systems with Rapidly Rotating Phase.** *Journal of Applied Mathematics and Mechanics*, (1984), 133–139. (cited on pages [6](#) and [31](#))
- NESTEROV, Y. E., 1983. **A Method for Solving the Convex Programming Problem with Convergence Rate $O(1/k^2)$.** In *Dokl. akad. nauk Sssr*, 543–547. (cited on pages [35](#), [64](#), [123](#), and [221](#))
- NGUYEN, C. V.; LI, Y.; BUI, T. D.; AND TURNER, R. E., 2018. **Variational Continual Learning.** In *the International Conference on Learning Representations*. (cited on page [75](#))
- NICHOL, A.; ACHIAM, J.; AND SCHULMAN, J., 2018. **On First-Order Meta-Learning Algorithms.** *arXiv preprint*, (2018). (cited on page [155](#))
- NOROUZI, M.; MIKOLOV, T.; BENGIO, S.; SINGER, Y.; SHLENS, J.; FROME, A.; CORRADO, G. S.; AND DEAN, J., 2013. **Zero-Shot Learning by Convex Combination of Semantic Embeddings.** In *the International Conference on Learning Representations*. (cited on page [2](#))
- PAN, S. J. AND YANG, Q., 2009. **A Survey on Transfer Learning.** In *the IEEE Transactions on Knowledge and Data Engineering*, 1345–1359. (cited on pages [71](#) and [131](#))
- PARISI, G. I.; KEMKER, R.; PART, J. L.; KANAN, C.; AND WERMTER, S., 2019. **Continual Lifelong Learning with Neural Networks: A Review.** *Neural Networks*, (2019), 54–71. (cited on page [2](#))
- PARISI, G. I.; TANI, J.; WEBER, C.; AND WERMTER, S., 2017. **Lifelong Learning of Human Actions with Deep Neural Network Self-Organisation.** *Neural Networks*, (2017), 137–149. (cited on page [76](#))

- PARISI, G. I.; TANI, J.; WEBER, C.; AND WERMTER, S., 2018. **Lifelong Learning of Spatiotemporal Representations with Dual-Memory Recurrent Self-Organisation.** *Frontiers in Neurorobotics*, (2018), 78. (cited on page 76)
- PASCANU, R.; GULCEHRE, C.; CHO, K.; AND BENGIO, Y., 2014. **How to Construct Deep Recurrent Neural Networks.** In *the International Conference on Learning Representations*. (cited on pages 50, 51, and 55)
- PASCANU, R.; MIKOLOV, T.; AND BENGIO, Y., 2012. **Understanding the Exploding Gradient Problem.** *CoRR, abs/1211.5063*, (2012). (cited on page 89)
- PASCANU, R.; MIKOLOV, T.; AND BENGIO, Y., 2013. **On the Difficulty of Training Recurrent Neural Networks.** In *the International Conference on Machine Learning*, 1310–1318. (cited on pages 6, 15, 43, 44, and 120)
- PEREYRA, G.; TUCKER, G.; CHOROWSKI, J.; KAISER, Ł.; AND HINTON, G., 2017. **Regularising Neural Networks by Penalising Confident Output Distributions.** In *the Workshop Track of the International Conference on Learning Representations*. (cited on page 75)
- PHAM, H.; GUAN, M.; ZOPH, B.; LE, Q.; AND DEAN, J., 2018a. **Efficient Neural Architecture Search via Parameters Sharing.** In *the International Conference on Machine Learning*, 4095–4104. (cited on page 184)
- PHAM, H.; GUAN, M. Y.; ZOPH, B.; LE, Q. V.; AND DEAN, J., 2018b. **Faster Discovery of Neural Architectures by Searching for Paths in a Large Model.** In *the International Conference on Machine Learning*. (cited on page 175)
- QIAN, N., 1999. **On the Momentum Term in Gradient Descent Learning Algorithms.** *Neural Networks*, (1999), 145–151. (cited on pages 35, 63, 132, 178, and 222)
- RAGHU, A.; RAGHU, M.; BENGIO, S.; AND VINYALS, O., 2019. **Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML.** In *the International Conference on Learning Representations*. (cited on page 60)
- RAIKO, T.; VALPOLA, H.; AND LECUN, Y., 2012. **Deep Learning Made Easier by Linear Transformations in Perceptrons.** In *the International Conference on Artificial Intelligence and Statistics*, 924–932. (cited on page 217)
- RAVI, S. AND LAROCHELLE, H., 2017. **Optimisation as a Model for Few-Shot Learning.** In *the International Conference on Learning Representations*. (cited on page 130)
- REBUFFI, S.-A.; KOLESNIKOV, A.; SPERL, G.; AND LAMPERT, C. H., 2017. **iCaRL: Incremental Classifier and Representation Learning.** In *the IEEE conference on Computer Vision and Pattern Recognition*, 2001–2010. (cited on pages 77, 143, and 162)
- RENDELL, L. A.; SHESHU, R.; AND TCHENG, D. K., 1987. **Layered Concept-Learning and Dynamically Variable Bias Management.** In *the International Joint Conferences on Artificial Intelligence*, 308–314. (cited on page 118)

-
- RIEMER, M.; CASES, I.; AJEMIAN, R.; LIU, M.; RISH, I.; TU, Y.; AND TESAURO, G., 2018. **Learning to Learn Without Forgetting by Maximizing Transfer and Minimizing Interference.** In *the International Conference on Learning Representations*. (cited on pages [147](#) and [155](#))
- ROBBINS, H. AND MONRO, S., 1951. **A Stochastic Approximation Method.** *The Annals of Mathematical Statistics*, (1951), 400–407. (cited on pages [4](#), [63](#), [113](#), [115](#), and [172](#))
- ROBINS, A., 1995. **Catastrophic Forgetting, Rehearsal and Pseudorehearsal.** *Connection Science*, (1995). (cited on pages [77](#) and [143](#))
- ROSSUM, G. V., 1995. **Python Tutorial.** Technical Report CS-R9526, Centrum voor Wiskunde en Informatica, Amsterdam. (cited on page [122](#))
- ROWEIS, S. T. AND SAUL, L. K., 2000. **Nonlinear Dimensionality Reduction by Locally Linear Embedding.** *Science*, (2000), 2323–2326. (cited on page [16](#))
- RUMELHART, D. E.; HINTON, G. E.; AND WILLIAMS, R. J., 1986. **Learning Representations by Back-Propagating Errors.** *Nature*, (1986), 533–536. (cited on pages [1](#), [4](#), [35](#), [172](#), and [221](#))
- RUSAK, Z.; WANG, S.; XU, L.; AND TAYLOR, S., 2012. **On the Global Nonlinear Stability of a Near-Critical Swirling Flow in a Long Finite-Length Pipe and the Path to Vortex Breakdown.** *Journal of Fluid Mechanics*, (2012), 295. (cited on pages [6](#) and [31](#))
- RUSU, A. A.; RABINOWITZ, N. C.; DESJARDINS, G.; SOYER, H.; KIRKPATRICK, J.; KAVUKCUOGLU, K.; PASCANU, R.; AND HADSELL, R., 2016. **Progressive Neural Networks.** In *the Deep Learning Symposium of Advances in Neural Information Processing Systems*. (cited on pages [75](#), [167](#), and [180](#))
- SALAKHUTDINOV, R. AND HINTON, G., 2007. **Learning a Nonlinear Embedding by Preserving Class Neighbourhood Structure.** In *the International Conference on Artificial Intelligence and Statistics*, 412–419. (cited on page [61](#))
- SALIMANS, T. AND KINGMA, D. P., 2016. **Weight Normalisation: A Simple Reparameterisation to Accelerate Training of Deep Neural Networks.** In *the Advances in Neural Information Processing Systems*, 901–909. (cited on pages [20](#), [181](#), and [219](#))
- SANTORO, A.; BARTUNOV, S.; BOTVINICK, M.; WIERSTRA, D.; AND LILICRAP, T., 2016. **Meta-Learning with Memory-Augmented Neural Networks.** In *the International Conference on Machine Learning*, 1842–1850. (cited on page [62](#))
- SANTURKAR, S.; TSIPRAS, D.; ILYAS, A.; AND MADRY, A., 2018. **How does Batch Normalisation Help Optimisation?** In *the Advances in Neural Information Processing Systems*, 2483–2493. (cited on pages [20](#), [181](#), and [219](#))
- SCHMIDHUBER, J., 1987. *Evolutionary Principles in Self-Referential Learning, or on Learning How to Learn: The Meta-Meta-... Hook.* Ph.D. thesis, Technische Universität München. (cited on pages [45](#), [72](#), and [73](#))

- SCHMIDHUBER, J., 1992. **Learning to Control Fast-Weight Memories: An Alternative to Dynamic Recurrent Networks.** *Neural Computation*, (1992), 131–139. (cited on pages [46](#) and [61](#))
- SCHMIDHUBER, J., 1993. **A Neural Network that Embeds Its Own Meta-Levels.** In *the IEEE International Conference on Neural Networks*, 407–412. (cited on page [46](#))
- SCHRAUDOLPH, N., 1998. **Accelerated Gradient Descent by Factor-Centring Decomposition.** *Technical Report/IDSIA*, (1998). (cited on pages [5](#), [14](#), [16](#), [36](#), [81](#), [171](#), and [217](#))
- SCHROPP, J. AND SINGER, I., 2000. **A Dynamical Systems Approach to Constrained Minimisation.** *Numerical Functional Analysis and Optimisation*, (2000), 537–551. (cited on page [35](#))
- SCHULTZ, M. AND JOACHIMS, T., 2003. **Learning a Distance Metric from Relative Comparisons.** In *the Advances in Neural Information Processing Systems*, 41–48. (cited on page [60](#))
- SEPLIARSKAIA, A.; KISELEVA, J.; AND DE RIJKE, M., 2021. **How to Not Measure Disentanglement.** (cited on page [184](#))
- SERRA, J.; SURIS, D.; MIRON, M.; AND KARATZOGLOU, A., 2018. **Overcoming Catastrophic Forgetting with Hard Attention to the Task.** In *the International Conference on Machine Learning*, 4548–4557. (cited on pages [76](#), [149](#), and [164](#))
- SHAN, S.; WILLSON, E.; WANG, B.; LI, B.; ZHENG, H.; AND ZHAO, B. Y., 2019. **Gotta Catch’Em All: Using Concealed Trapdoors to Detect Adversarial Attacks on Neural Networks.** *arXiv preprint arXiv:1904.08554*, (2019). (cited on page [134](#))
- SHANNON, C. E., 2001. **A Mathematical Theory of Communication.** *ACM SIGMOBILE Mobile Computing and Communications Review*, (2001), 3–55. (cited on pages [18](#) and [26](#))
- SHARIF RAZAVIAN, A.; AZIZPOUR, H.; SULLIVAN, J.; AND CARLSSON, S., 2014. **CNN Features Off-The-Shelf: An Astounding Baseline for Recognition.** In *the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 806–813. (cited on page [2](#))
- SHIMODAIRA, H., 2000. **Improving Predictive Inference under Covariate Shift by Weighting the Log-likelihood Function.** *Journal of Statistical Planning and Inference*, (2000), 227–244. (cited on page [20](#))
- SHIN, H.; LEE, J. K.; KIM, J.; AND KIM, J., 2017. **Continual Learning with Deep Generative Replay.** In *the Advances in Neural Information Processing Systems*, 2994–3003. (cited on pages [78](#) and [143](#))

-
- SIMONYAN, K. AND ZISSERMAN, A., 2015. **Very Deep Convolutional Networks for Large-Scale Image Recognition**. In *the International Conference on Learning Representations*. (cited on page [129](#))
- ŞİMŞEKLI, U.; SENER, O.; DELIGIANNIDIS, G.; AND ERDOGDU, M. A., 2020. **Hausdorff Dimension, Stochastic Differential Equations, and Generalization in Neural Networks**. In *the Advances in Neural Information Processing Systems*. (cited on page [28](#))
- SNELL, J.; SWERSKY, K.; AND ZEMEL, R., 2017. **Prototypical Networks for Few-Shot Learning**. In *the Advances in Neural Information Processing Systems*, 4077–4087. (cited on page [61](#))
- SNEYD, J.; HAN, J. M.; WANG, L.; CHEN, J.; YANG, X.; TANIMURA, A.; SANDERSON, M. J.; KIRK, V.; AND YULE, D. I., 2017. **On the Dynamical Structure of Calcium Oscillations**. In *the Proceedings of the National Academy of Sciences*, 1456–1461. (cited on pages [1](#), [6](#), and [31](#))
- SOCHER, R.; GANJOO, M.; MANNING, C. D.; AND NG, A., 2013. **Zero-Shot Learning through Cross-Modal Transfer**. In *Advances in Neural Information Processing Systems*, 935–943. (cited on page [163](#))
- SOMPOLINSKY, H.; CRISANTI, A.; AND SOMMERS, H.-J., 1988. **Chaos in Random Neural Networks**. *Physical Review Letters*, (1988), 259. (cited on pages [97](#) and [175](#))
- SRIVASTAVA, N.; HINTON, G.; KRIZHEVSKY, A.; SUTSKEVER, I.; AND SALAKHUTDINOV, R., 2014. **Dropout: A Simple Way to Prevent Neural Networks from Overfitting**. *Journal of Machine Learning Research*, (2014), 1929–1958. (cited on pages [20](#), [88](#), and [174](#))
- SRIVASTAVA, R. K.; GREFF, K.; AND SCHMIDHUBER, J., 2015. **Highway Networks**. In *the Workshop of International Conference on Machine Learning*. (cited on pages [3](#), [5](#), [14](#), [15](#), [21](#), [34](#), [39](#), [43](#), [81](#), [103](#), [111](#), [158](#), [159](#), [171](#), [182](#), and [217](#))
- SU, W.; BOYD, S.; AND CANDÉS, E., 2014. **A Differential Equation for Modelling Nesterov’s Accelerated Gradient Method: Theory and Insights**. In *the Advances in Neural Information Processing Systems*, 2510–2518. (cited on pages [35](#) and [222](#))
- SÜLI, E. AND MAYERS, D. F., 2003. *An Introduction to Numerical Analysis*. Cambridge University Press. (cited on page [59](#))
- SUN, Q.; LIU, Y.; CHUA, T.-S.; AND SCHIELE, B., 2019. **Meta-Transfer Learning for Few-Shot Learning**. In *the IEEE Conference on Computer Vision and Pattern Recognition*, 403–412. (cited on page [59](#))
- SUN, X.; ZHANG, Z.; REN, X.; LUO, R.; AND LI, L., 2021. **Exploring the Vulnerability of Deep Neural Networks: A Study of Parameter Corruption**. In *the AAAI Conference on Artificial Intelligence*. (cited on page [118](#))

- SUNG, F.; YANG, Y.; ZHANG, L.; XIANG, T.; TORR, P. H.; AND HOSPEDALES, T. M., 2018. **Learning to Compare: Relation Network for Few-Shot Learning.** In *the IEEE Conference on Computer Vision and Pattern Recognition*, 1199–1208. (cited on page 61)
- SUSSILLO, D. AND ABBOTT, L. F., 2009. **Generating Coherent Patterns of Activity from Chaotic Neural Networks.** *Neuron*, (2009), 544–557. (cited on pages 8, 97, and 175)
- SUTSKEVER, I.; MARTENS, J.; DAHL, G.; AND HINTON, G., 2013. **On the Importance of Initialisation and Momentum in Deep Learning.** In *the International Conference on Machine Learning*, 1139–1147. (cited on page 35)
- SUTSKEVER, I.; VINYALS, O.; AND LE, Q. V., 2014. **Sequence to Sequence Learning with Neural Networks.** In *the Advances in Neural Information Processing Systems*, 3104–3112. (cited on page 48)
- SUTTON, R. S. AND BARTO, A. G., 2018. *Reinforcement Learning: An Introduction.* MIT press. (cited on page 2)
- SZEGEDY, C.; LIU, W.; JIA, Y.; SERMANET, P.; REED, S.; ANGUELOV, D.; ERHAN, D.; VANHOUCHE, V.; AND RABINOVICH, A., 2015. **Going Deeper with Convolutions.** In *the IEEE conference on Computer Vision and Pattern Recognition*, 1–9. (cited on pages 3, 5, 15, 17, 37, 38, 147, and 159)
- TAKESIAN, A. E. AND HENSCH, T. K., 2013. **Balancing Plasticity/Stability Across Brain Development.** *Progress in Brain Research*, (2013), 3–34. (cited on page 72)
- TAKEUCHI, K., 1976. **Distribution of Information Statistic and Validity Criterion of Models.** *Mathematical Science*, (1976), 12–18. (cited on page 26)
- TALLEC, C. AND OLLIVIER, Y., 2018. **Can Recurrent Neural Networks Warp Time?** In *the International Conference on Learning Representations*. (cited on page 88)
- TENENBAUM, J. B.; DE SILVA, V.; AND LANGFORD, J. C., 2000. **A Global Geometric Framework for Nonlinear Dimensionality Reduction.** *Science*, (2000), 2319–2323. (cited on page 16)
- TENNYSON, J., 1989. **Flip-Flop Modes in Symmetric and Asymmetric Colliding-Beam Storage Rings.** Technical report, Lawrence Berkeley National Lab.(LBNL), Berkeley, CA (United States). (cited on page 33)
- TIELEMAN, T. AND HINTON, G., 2012. **Lecture 6.5-RMSProp: Divide the Gradient by a Running Average of its Recent Magnitude.** *COURSERA: Neural Networks for Machine Learning*, (2012), 26–31. (cited on pages 3, 64, and 89)
- TSANEVA-ATANASOVA, K.; OSINGA, H. M.; RIESS, T.; AND SHERMAN, A., 2010. **Full System Bifurcation Analysis of Endocrine Bursting Models.** *Journal of Theoretical Biology*, 264, 4 (2010), 1133–1146. (cited on page 185)

-
- ULYANOV, D.; VEDALDI, A.; AND LEMPITSKY, V., 2016. **Instance Normalisation: The Missing Ingredient for Fast Stylisation**. *arXiv preprint arXiv:1607.08022*, (2016). (cited on pages [21](#) and [220](#))
- VAN DE VEN, G. M. AND TOLIAS, A. S., 2018. **Three Scenarios for Continual Learning**. In *the Continual Learning Workshop of Advances in Neural Information Processing Systems*. (cited on pages [73](#), [74](#), [78](#), [148](#), and [180](#))
- VAN DEN OORD, A.; DIELEMAN, S.; ZEN, H.; SIMONYAN, K.; VINYALS, O.; GRAVES, A.; KALCHBRENNER, N.; SENIOR, A.; AND KAVUKCUOGLU, K., 2016. **WaveNet: A Generative Model for Raw Audio**. In *the ISCA Speech Synthesis Workshop*, 125–125. (cited on page [5](#))
- VAN DER POL, B., 1926. **On “Relaxation-Oscillations”**. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, (1926), 978–992. (cited on page [37](#))
- VANSCHOREN, J., 2019. **Meta-Learning**. In [Hutter et al. \[2019\]](#), 39–68. (cited on page [2](#))
- VASWANI, A.; SHAZEER, N.; PARMAR, N.; USZKOREIT, J.; JONES, L.; GOMEZ, A. N.; KAISER, Ł.; AND POLOSUKHIN, I., 2017. **Attention is All You Need**. In *the Advances in Neural Information Processing Systems*, 5998–6008. (cited on pages [5](#) and [53](#))
- VEIT, A.; WILBER, M. J.; AND BELONGIE, S., 2016. **Residual Networks Behave Like Ensembles of Relatively Shallow Networks**. In *the Advances in Neural Information Processing Systems*, 550–558. (cited on page [36](#))
- VILALTA, R.; BLIX, G.; AND RENDELL, L., 1997. **Global Data Analysis and the Fragmentation Problem in Decision Tree Induction**. In *the European Conference on Machine Learning*, 312–326. (cited on page [17](#))
- VINYALS, O.; BENGIO, S.; AND KUDLUR, M., 2016a. **Order Matters: Sequence to Sequence for Sets**. In *the International Conference of Learning Representations*. (cited on page [61](#))
- VINYALS, O.; BLUNDELL, C.; LILLICRAP, T.; WIERSTRA, D.; AND KAVUKCUOGLU, K., 2016b. **Matching Networks for One Shot Learning**. In *the Advances in Neural Information Processing Systems*, 3630–3638. (cited on pages [2](#), [59](#), and [61](#))
- VIRTANEN, P.; GOMMERS, R.; OLIPHANT, T. E.; HABERLAND, M.; REDDY, T.; COURNAPEAU, D.; BUROVSKI, E.; PETERSON, P.; WECKESSER, W.; BRIGHT, J.; VAN DER WALT, S. J.; BRETT, M.; WILSON, J.; JARROD MILLMAN, K.; MAYOROV, N.; NELSON, A. R. J.; JONES, E.; KERN, R.; LARSON, E.; CAREY, C.; POLAT, İ.; FENG, Y.; MOORE, E. W.; VANDERPLAS, J.; LAXALDE, D.; PERKTOLD, J.; CIMRMAN, R.; HENRIKSEN, I.; QUINTERO, E. A.; HARRIS, C. R.; ARCHIBALD, A. M.; RIBEIRO, A. H.; PEDREGOSA, F.; VAN MULBREGT, P.; AND CONTRIBUTORS, S. . ., 2020. **SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python**. *Nature Methods*, (2020), 261–272. (cited on page [122](#))

- VON HELMHOLTZ, H., 1978/1903. *Die Tatsachen in der Wahrnehmung*. Vorträge und Reden (Braunschweig: Vieweg). (cited on page 70)
- VON WEIZSÄCKER, C. F., 1985. *Aufbau der Physik*. Hauser. (cited on pages 45, 46, 61, and 73)
- VUORIO, R.; SUN, S.-H.; HU, H.; AND LIM, J. J., 2019. **Multimodal Model-Agnostic Meta-Learning via Task-Aware Modulation**. In *the Advances in Neural Information Processing Systems*, 1–12. (cited on page 59)
- WAH, C.; BRANSON, S.; WELINDER, P.; PERONA, P.; AND BELONGIE, S., 2011. **The Caltech-UCSD Birds-200-2011 Dataset**. *California Institute of Technology*, (2011). (cited on pages 162 and 214)
- WAN, L.; ZEILER, M.; ZHANG, S.; LE CUN, Y.; AND FERGUS, R., 2013. **Regularisation of Neural Networks Using DropConnect**. In *the International Conference on Machine Learning*, 1058–1066. (cited on page 88)
- WANG, J. X.; KURTH-NELSON, Z.; TIRUMALA, D.; SOYER, H.; LEIBO, J. Z.; MUNOS, R.; BLUNDELL, C.; KUMARAN, D.; AND BOTVINICK, M., 2016. **Learning to Reinforce Learn**. *CogSci*, (2016). (cited on pages 126 and 177)
- WEIGEND, A. S.; RUMELHART, D. E.; AND HUBERMAN, B. A., 1991. **Generalization by Weight-Elimination with Application to Forecasting**. In *the Advances in Neural Information Processing Systems*, 875–882. (cited on page 27)
- WEINBERGER, K. Q.; BLITZER, J.; AND SAUL, L. K., 2006. **Distance Metric Learning for Large Margin Nearest Neighbour Classification**. In *the Advances in Neural Information Processing Systems*, 1473–1480. (cited on page 60)
- WESTHEIMER, G., 2008. **Was Helmholtz a Bayesian?** *Perception*, (2008), 642–650. (cited on page 70)
- WESTON, J.; CHOPRA, S.; AND BORDES, A., 2015. **Memory networks**. In *the International Conference on Learning Representations*. (cited on pages 62 and 63)
- WESTON, J.; RATLE, F.; MOBAHI, H.; AND COLLOBERT, R., 2012. **Deep Learning via Semi-Supervised Embedding**. In *Neural Networks: Tricks of the Trade*, 639–655. Springer. (cited on page 19)
- WICHROWSKA, O.; MAHESWARANATHAN, N.; HOFFMAN, M. W.; COLMENAREJO, S. G.; DENIL, M.; DE FREITAS, N.; AND SOHL-DICKSTEIN, J., 2017. **Learned Optimisers that Scale and Generalise**. In *the International Conference on Machine Learning*. (cited on pages 125, 177, and 178)
- WILLIAMS, T., 2019. **John Duns Scotus**. In *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter edn. (cited on page 82)

-
- WITTEN, E., 1995. **String Theory Dynamics in Various Dimensions**. *Nuclear Physics B*, (1995), 85–126. (cited on page 95)
- WU, Y. AND HE, K., 2018. **Group Normalisation**. In *the European Conference on Computer Vision*, 3–19. (cited on page 21)
- WU, Z. AND KING, S., 2016. **Investigating Gated Recurrent Networks for Speech Synthesis**. In *the IEEE International Conference on Acoustics, Speech and Signal Processing*, 5140–5144. (cited on pages 52 and 55)
- XIAO, H.; RASUL, K.; AND VOLLGRAF, R., 2017. **Fashion-MNIST: A Novel Image Dataset for Benchmarking Machine Learning Algorithms**. *arXiv preprint arXiv:1708.07747*, (2017). (cited on pages 90, 124, 134, and 214)
- XU, J.; SUN, X.; ZHANG, Z.; ZHAO, G.; AND LIN, J., 2019. **Understanding and Improving Layer Normalisation**. In *the Advances in Neural Information Processing Systems*, 4381–4391. (cited on pages 21 and 219)
- YANG, Z.; LIU, Y.; BAO, C.; AND SHI, Z., 2020. **Interpolation between Residual and Non-Residual Networks**. In *the International Conference on Machine Learning*. (cited on page 167)
- YAO, K.; COHN, T.; VYLOMOVA, K.; DUH, K.; AND DYER, C., 2015. **Depth-Gated LSTM**. In *the Jelinek Summer Workshop*, 1. (cited on page 54)
- YAO, Z.; GHOLAMI, A.; LEI, Q.; KEUTZER, K.; AND MAHONEY, M. W., 2018. **Hessian-Based Analysis of Large Batch Training and Robustness to Adversaries**. In *the Advances in Neural Information Processing Systems*, 4949–4959. (cited on pages 27, 28, and 220)
- YOON, J.; YANG, E.; LEE, J.; AND HWANG, S. J., 2018. **Lifelong Learning with Dynamically Expandable Networks**. In *the International Conference on Learning Representations*. (cited on page 76)
- ZENG, J.; CAI, H.; PENG, H.; WANG, H.; ZHANG, Y.; AND AKUTSU, T., 2020. **Causalcall: Nanopore Basecalling Using a Temporal Convolutional Network**. *Frontiers in Genetics*, (2020), 1332. (cited on page 5)
- ZENKE, F.; POOLE, B.; AND GANGULI, S., 2017. **Continual Learning through Synaptic Intelligence**. In *the International Conference on Machine Learning*, 3987–3995. (cited on page 75)
- ZHANG, C.; BENGIO, S.; HARDT, M.; RECHT, B.; AND VINYALS, O., 2017a. **Understanding Deep Learning Requires Rethinking Generalisation**. In *International Conference on Learning Representations*. (cited on page 29)
- ZHANG, X.; LI, Z.; CHANGE LOY, C.; AND LIN, D., 2017b. **PolyNet: A Pursuit of Structural Diversity in Very Deep Networks**. In *the IEEE Conference on Computer Vision and Pattern Recognition*, 718–726. (cited on pages 37, 222, and 223)

- ZHAO, J.; HUANG, F.; LV, J.; DUAN, Y.; QIN, Z.; LI, G.; AND TIAN, G., 2020. **Do RNN and LSTM have Long Memory?** In *the International Conference on Machine Learning*, 11365–11375. (cited on page [53](#))
- ZHOU, G.; SOHN, K.; AND LEE, H., 2012. **Online Incremental Feature Learning with Denoising Autoencoders.** In *the International Conference on Artificial Intelligence and Statistics*, 1453–1461. (cited on page [167](#))
- ZOPH, B. AND LE, Q. V., 2017. **Neural Architecture Search with Reinforcement Learning.** In *the International Conference on Learning Representations*. (cited on pages [36](#), [126](#), [175](#), and [183](#))

A: Datasets

This appendix describes some repeatedly used datasets in this thesis. Descriptions of unlisted datasets could be found therein their respective experimental sections.

MNIST

This computer vision dataset is used in Section(s) 4.1.6.1, 4.1.6.2, 5.1.5, 6.1.5, and 6.2.5.

MNIST [LeCun et al., 1998] is a popular machine learning benchmark database of grey-scale handwritten digits of numbers between [0-9]. There are 60K and 10K images for training and testing respectively; and the images are presented on 28×28 -pixel boxes. Interested readers should refer to <http://yann.lecun.com/exdb/mnist/> for descriptions documented by the original data providers.

Training

Since MNIST is a classification task, we use a *softmax layer* [Goodfellow et al., 2016] with 10 output dimensions as the classifier in this thesis for this task. In brief, softmax is a generalisation of the logistic function for normalising neural network output as a probability distribution for predicted output classes. A softmax layer $\sigma_s : \mathbb{K} \rightarrow \mathbb{K}$ typically follows the following formulation

$$\sigma_s(\mathbf{z}) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad (1)$$

for $i = 1, \dots, K$ on some $\mathbf{z} = (z_1, \dots, z_K) \in \mathbb{K}$.

Then in order to fine-tune the network, we set our objective loss as the discrete *cross entropy loss* [Goodfellow et al., 2016]. The cross entropy loss is an information measurement on the same underlying events over two probability distributions of p and q . In addition, the measurement is defined as

$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log(q(x)) \quad (2)$$

for individual events x over a support set of \mathcal{X} .

FashionMNIST and KMNIST

These computer vision datasets are used in Section(s) 5.1.5 and 5.2.4.

FashionMNIST [Xiao et al., 2017] and KMNIST [Clanuwat et al., 2018] are both designed similarly to MNIST. They have 60K training and 10K test greyscale images with 10 classes of items on 28×28 pixels. FashionMNIST has images of fashion items and KMNIST consists kuzushiji (崩し字) versions of Hiragana (平仮名) characters from classical Japanese literature. Interested readers should refer to <https://github.com/zalandoresearch/fashion-mnist> for descriptions documented by the original FashionMNIST providers; and likely, to <https://github.com/rois-codh/kmnist> for descriptions documented by the original KMNIST providers.

Cifar10 and Cifar100

These computer vision datasets are used in Section(s) 4.2.5.3, 5.1.5, 6.1.5, and 6.2.5.

The *Canadian Institute For Advanced Research* (Cifar) collected images that are commonly used to train machine learning and computer vision algorithms. They provided two of the most widely used datasets for machine learning research. Cifar10 [Krizhevsky, 2009] is a coloured dataset with 10 classes of objects. There are 50K training and 10K test images all presented on 32×32 pixel boxes. Interested readers should refer to <https://www.cs.toronto.edu/~kriz/cifar.html> for descriptions documented by the original data providers.

STL10

This computer vision dataset is used in Section(s) 5.2.4.

The STL10 [Coates et al., 2011] dataset is inspired by Cifar10 and it consists 5K training and 8K test RGB colour images in 10 classes on 96×96 pixels. Its main difference to Cifar10 is that it also included unlabelled data for unsupervised learning. Interested readers should refer to <https://cs.stanford.edu/~acoates/stl10/> for descriptions documented by the original data providers.

CUB200

This computer vision dataset is used in Section(s) 6.2.5.

The CUB200 [Wah et al., 2011] dataset consisted 5,994 training and 5,794 test over 200 different coloured (predominately North American) birds species. The sizes of most images are roughly 500×500 pixels. Interested readers should refer to <http://www.vision.caltech.edu/visipedia/CUB-200-2011.html> for descriptions documented by the original data providers.

Penn Treebank

This natural language processing dataset is used in Section(s) 4.2.5.1.

Penn Treebank (PTB) [Marcus et al., 1993] is an English corpus dataset which consists 929K training words, 73K validation words, and 82K test words. It has a vocabulary size of 10K words.

Training

In this thesis, we use PTB for language modelling. A language model \mathcal{F} constructs the probability mapping \mathbb{P} in a sequence \mathfrak{W} of words \mathfrak{w} of length D following the Bayesian inferential process such that

$$\mathcal{F}(\mathfrak{W}) = \mathcal{F}(\mathfrak{w}_1|\mathfrak{w}_0) \times \mathcal{F}(\mathfrak{w}_2|\mathfrak{w}_0\mathfrak{w}_1) \times \cdots \times \mathcal{F}(\mathfrak{w}_D|\mathfrak{w}_0\mathfrak{w}_1 \cdots \mathfrak{w}_{D-1}) \quad (3)$$

$$= \prod_{i=1}^n \mathcal{F}(\mathfrak{w}_i|\mathfrak{w}_0\mathfrak{w}_1 \cdots \mathfrak{w}_{i-1}). \quad (4)$$

Since a language model predicts the next word token given preceding ones, it employs softmax (see Equation (1)) as its final layer.

Though both language models and image recognition models use softmax layers, the well-being of language models are not evaluated by the cross entropy loss (see Equation (2)). Instead, they use the *Perplexity Loss* (PPL) [Brown et al., 1992] as a measurement. In brief, PPL is the inverse probability of the test set of which is normalised by the number of words

$$\text{PPL}(\mathfrak{W}) = \frac{1}{\mathbb{P}(\mathfrak{w}_1, \mathfrak{w}_2, \dots, \mathfrak{w}_D)^{\frac{1}{D}}}. \quad (5)$$

However given cross entropy loss \mathfrak{H} , PPL can be calculated as

$$\text{PPL} = e^{\mathfrak{H}}. \quad (6)$$

An in-depth discussion for the derivation can be found in <https://towardsdatascience.com/the-relationship-between-perplexity-and-entropy-in-nlp-f81888775ccc>.

B: Additional Notes for the Preliminaries

As we mentioned in the foreword of this preliminary chapter, this appendix was purposely written in a more narrative style for those readers who are less familiar with our thesis prerequisites. A more analytical tone is used for the rest of the thesis. Nonetheless, some highlights and interesting notes can be found below in this in-chapter clarification (ICC) for the preliminary chapter.

For Prerequisite 1

ICC-01: Both the highway network paper [Srivastava et al., 2015] and the LSTM RNN paper [Hochreiter and Schmidhuber, 1997b] were (co-)developed by Professor Jürgen Schmidhuber and the former could be seen as a natural extension of the latter. [PDF return.](#)

ICC-02: The work of Schraudolph [1998] studied the variability of updates in gradient descent. Prior of 2015 (the creation of highway networks), this method was revisited by Raiko et al. [2012] of which the authors of that paper made explicit modifications to a MLP network for accelerating and stabilising network training. Also prior to 2015, Graves [2013] added shortcut connections to his deep predictive RNN architecture to improve the flow of information. Graves’s work was important because it was considered as one of the seminal papers that addressed sequence generation with RNNs. [PDF return.](#)

ICC-03: During the initial training phase of a statistical model, the training loss as well as the testing loss will decrease as we fine-tune the parameters. However, if we over-train the model, the parameters will be adjusted solely for explaining the training data and would thus not generalise well for unseen data. Hence after a certain amount of training, the testing loss will start to increase. Interestingly, the work of Nakkiran et al. [2020] showed that state of the art deep learning models do not follow these principles. Instead, the testing loss exhibits *double descent* – which it would experience an initial decrease, then an increase, but ultimately a second and much more sharper decrease. Interested readers should refer to Figure 1 on page 1 of the mentioned paper. [PDF return.](#)

ICC-04: By random initialisation, we are referring to the common practice of setting biases to a vector of 0s and sampling entries of the transformation matrix from a uniform distribution $U[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}]$ where n is the output dimension of the previous layer. PDF [return](#).

ICC-05: One finding of [Glorot and Bengio \[2010\]](#) was that neural network biases could be trained extremely efficiently, causing the early training phase of the network to rely heavily on the biases as the main inferential parameter. If a layer activation was given as $Wx + b$, this would cause the error associated with x to be pushed towards 0. Hence if the sigmoid function was employed as the activation function, the shape of its derivative would cause the back-propagated 0 to stuck in a local region. PDF [return](#).

ICC-06: The formal derivation of the Xavier initialisation can be found on page 5 of [Glorot and Bengio \[2010\]](#). The derivation was mainly achieved using the product relations of variance, and it replaced the uniform distribution of ICC-04 with $U[-\frac{\sqrt{6}}{\sqrt{n_j+n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j+n_{j+1}}}]$ where j referred to the input dimension of the current layer, $j+1$ referred to the output dimension of the current layer. Note, j was also equivalent to the output dimension of the previous dimension. PDF [return](#).

ICC-07: Similar to [Glorot and Bengio \[2010\]](#), [He et al. \[2015\]](#) also kept the biases as 0 but they changed the sampling of the weights to a zero-mean Gaussian distribution with standard deviation of $\sqrt{\frac{2}{n}}$. This change was made because, in addition to the input variances studied in [Glorot and Bengio](#), [He et al.](#) also studied the output variances. Refer to Equations (6) - (10) on page 4 of [He et al.](#) for the extra factors. Their novel changes provided specific highlights on the changes in the variability of output features as the transformation weights varied. PDF [return](#).

ICC-08: BN [[Ioffe and Szegedy, 2015a](#)] replaced the common layer output of $\hat{x} = x_{tijk}$ with $\hat{x} = \text{BN}_{\gamma,\beta}(x_{tijk})$ where γ and β were 2 extra learnable parameters. Here we denoted $x_{tijk} \in \mathbb{R}^{T \times C \times W \times H}$ as the input in batch B where the (k, j) -element spanned the spatial dimensions, i for the feature channel, and that t for the index of the batch. Formally, BN was given as

$$B \text{ mean: } \mu_i = \frac{1}{HWT} \sum_{t=1}^T \sum_{l=1}^W \sum_{m=1}^H x_{tijk}; \quad B \text{ variance: } \sigma_i^2 = \frac{1}{HWT} \sum_{t=1}^T \sum_{l=1}^W \sum_{m=1}^H (x_{tijk} - \mu_i)^2;$$

$$\text{Normalisation: } x_{tijk}^{\hat{}} \leftarrow \frac{x_{tijk} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}; \quad \& \text{ Scale and Shift: } \text{BN}_{\gamma,\beta}(x_{tijk}) = \gamma x_{tijk}^{\hat{}} + \beta.$$

We could see that the layer-wise neural values were normalised as according to the batch-wise feature-specific (i -wise) mean and variance; and that parameters β and γ were introduced to provide neural network with expressiveness post-normalisation.

Note that BN should only use the feature statistics that were gathered during the training data and not with those in the testing data. If BN were not froze during test time, the B mean and B variance (also known as the running-mean and running-variance) would continue to be updated with the testing data hence giving BN an

extremely powerful and unfair advantage. PDF [return](#).

ICC-09: Refer to the full analyses with the Hessian on page 8 of [Santurkar et al. \[2018\]](#). A similar analysis, albeit simpler, could also be found in the well-known blog³ by Mr. Keita Kurita of the Carnegie Mellon University.

Consider the scenario where we are optimising the loss function $f(w)$ for weight w with gradient descent and that the weight is currently parameterised as w_0 . The second order Taylor expansion around the current weights is

$$f(w) \approx f(w_0) + (w - w_0)^T g + \frac{1}{2}(w - w_0)^T H(w - w_0)$$

where g and H are the gradient and Hessian matrix of loss f at w_0 . That is, a slight update in this direction with size ϵ will result in

$$f(w_0 - \epsilon g) \approx f(w_0) - \epsilon g^T g + \frac{1}{2}\epsilon^2 g^T H g.$$

Note that both the first order term and the second order term $\frac{1}{2}\epsilon^2 g^T H g$ are semi-positive definite. Hence a large second order term could potentially increase the loss instead of decreasing it. The gradient would however most likely decrease with BN because BN indirectly controlling the scalings (like that of a step size) of the different order terms. PDF [return](#).

ICC-10: WN [[Salimans and Kingma, 2016](#)] did not require direct access to information in the batch. For the pre-activated transformation $\hat{x} = Wx + b$ with weight W and bias b , WN separately trained the *magnitude* g and *direction* v of W with the new weight $\hat{W} = \frac{g}{|v|}v$ where $|v|$ was the Euclidean norm of v . PDF [return](#).

ICC-11: Realising this potential issue, the authors of WN [[Santurkar et al., 2018](#)] proposed a complementary method to WN which they called the *Mean-Only* BN (MOBN). Unlike the standard BN, MOBN did not divide the neural values by the mini-batch standard deviation. PDF [return](#).

ICC-12: LN [[Ba et al., 2016b](#)] removed the dependency on the batch by modifying the mean and variance in ICC-08 with

$$\text{mean: } \mu_i = \frac{1}{|B|} \sum_{\xi=1}^B x_{\xi}; \quad \text{variance: } \sigma_i^2 = \frac{1}{|B|} \sum_{\xi=1}^B (x_{\xi} - \mu_i)^2.$$

Note that the variable had changed from x_{tijk} to x_{ξ} . This was because of 2 reasons: first, that LN was proposed specifically for recurrent neural information processing thus unlike CNNs, the weights in conventional RNNs did not have the additional features/colour channels; and second, that LN took all features of the same layer into account regardless of their individual uniqueness. PDF [return](#).

ICC-13: [Xu et al. \[2019\]](#) found that LN made important modifications to the backward pass. Furthermore, they showed that when the new features were formulated as $\hat{x} = g \odot \text{LN}(x) + b$ with shifts and rescales from parameters b and g , the net-

³Refer to “Intuition 2: High-Order Effects” of <https://mlexplained.com/2018/01/10/an-intuitive-explanation-of-why-batch-normalization-really-works-normalization-in-deep-learning-part-1/>

work could actually be severely over-parameterised thus leading to over-fitting. PDF [return](#).

ICC-14: Similar to LN in ICC-12, IN [[Ulyanov et al., 2016](#)] also ignored the batch-wise information of BN. It was modified as below

$$B \text{ mean: } \mu_i = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H x_{tijk}; \quad B \text{ variance: } \sigma_i^2 = \frac{1}{HW} \sum_{l=1}^W \sum_{m=1}^H (x_{tijk} - \mu_i)^2.$$

Note and compare to ICC-08, the dependency on t no longer existed. PDF [return](#).

For Prerequisite 2

ICC-15: Regarding random projection, see the definition of matrix A in page 5 of [Keskar et al. \[2017\]](#), where “For that purpose, we introduce an $n \times p$ matrix A , whose columns are randomly generated.” PDF [return](#).

ICC-16: Unfortunately, the paper of [Yao et al. \[2018\]](#) did not provide clear descriptions on how they derived their Hessian eigenvalues. PDF [return](#).

ICC-17: [Dinh et al. \[2017\]](#)’s criticism towards the the concept of sharpness was mainly due to how our imagination is largely plagued when the dimensionality increases well beyond the third dimension. PDF [return](#).

ICC-18: There are actually substantive differences on how [Hochreiter and Schmidhuber \[1995\]](#) and [Dinh et al. \[2017\]](#) conceptualise the term of “minima”. The former, see Section 2 on “Task / Architecture / Boxes” of that paper, largely described the search of flat minima as a task of measuring the differences in the change of volume of a contiguous space. [Hochreiter and Schmidhuber](#) conceptualisation of minima was hence limited to a 3 dimensional projection of the high dimensional parametric space. In contrast and in Section 4.1 “Volume ϵ -flatness” of the latter paper, [Dinh et al.](#) argued that minima could be conceptualised in terms of its volume, but not by direct measurement and instead via fitting the said volume with high dimensional spheres. PDF [return](#).

ICC-19: Following ICA-18, see Definition 4 on page 4 of [Dinh et al. \[2017\]](#) for *non-negative homogeneity*, then Definition 4 on page 4 for the concept of α -scale transformations. PDF [return](#).

ICC-20: [Dinh et al. \[2017\]](#)’s infinite volume proof was a two part proof presented in Section 4.1 “Volume ϵ -flatness”. The proof mainly played around 2 concepts: on the α -scale transformations (see ICC-19) for maximal non-identifiability and on the geometric meaning of the determinant of the Jacobian of the α -scale transformations. The first part of the proof addressed when the 2 layers of the ReLU network were of different hidden dimensionality; while the second part addressed when the 2 layers shared the same hidden dimensionality. PDF [return](#).

ICC-21: The changes in the eigenvalues were not necessarily large under [Jastrzębski et al. \[2019\]](#)’s experiments. Figure 4 of page 5 of [Jastrzębski et al. \[2019\]](#) showed that the eigenvalues could easily increase over 2 or even 3 orders of magnitudes; but their Figure 2 of page 3 also showed that sometimes the eigenvalues merely doubled

the initial values. PDF [return](#).

ICC-22: However, the student-author of this thesis personally found that several claims and motivation behind [Li et al. \[2018\]](#) were interesting but dubious at best. Several concerns regarding this paper could also be found on this website <https://openreview.net/forum?id=HkmaTz-0W>. PDF [return](#).

For Prerequisite 3

ICC-23: A comparison of optimisation schemes can be found below.

(Stochastic) gradient descent:

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}$$

where $\mathcal{L} = \mathcal{L}(\theta_t)$ is the momentum coefficient.

SGD with Momentum [[Rumelhart et al., 1986](#)]:

$$v_{t+1} = pv_t - \alpha \nabla_{\theta} \mathcal{L}$$

$$\theta_{t+1} = \theta_t + v_{t+1}$$

where $p \in [0, 1]$ is the momentum coefficient.

NAG [[Nesterov, 1983](#)]:

$$v_{t+1} = pv_t - \alpha \nabla_{\theta} \mathcal{L}(\theta_t + pv_t)$$

$$\theta_{t+1} = \theta_t + v_{t+1}$$

where $p \in [0, 1]$ is the momentum coefficient.

ADAM [[Kingma and Ba, 2015](#)]:

$$g_t \leftarrow \nabla_{\theta} \mathcal{L}$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

$$\theta_{t+1} = \theta_t - \alpha \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$$

where the 1st moment and 2nd moment vectors of m_t and v_t are initialised as 0s. Hyper-parameters β_1 and β_2 are exponential decay rates of which

[Kingma and Ba \[2015\]](#) suggested to use 0.9 and 0.999 respectively, and $\epsilon = 10^{-8}$.

PDF [return](#).

ICC-24: This ICC summarises 2 important steps in [Qian \[1999\]](#)'s paper.

Step I: Relating SGD & momentum to a Newtonian particle

See Equation (3) of that paper where the authors rearranged the discretised version of gradient descent into its differential form. The authors also replaced loss \mathcal{L} as E of which they used to describe the energy of Newtonian motion. Equation (4) is where they generalised the gradient descent update as a Newtonian particle motion. Of which they discretised in Equation (5) to establish the link between SGD with momentum to Newtonian motions.

Step II: Existence of Lyapunov function for SGD & momentum

See Section 3.1 of that paper, where the authors claimed that, since SGD with momentum could be linked to Newtonian motions, then it was of common knowledge of which the total energy was a Lyapunov function that monotonically decreased due to the presence of friction. As a side note, this was also the reason why they rewrote their Equation (4) to Equation (10) with $-H(w - w_0)$ where matrix H was positive definite. Furthermore, they used this property to analyse the speed-up provided via using momentum.

PDF [return](#).

ICC-25: See Equation (4) in page 5 of [Su et al. \[2014\]](#). Their work was mainly based on the assumption of smoothness of NAG via Lipschitz constant, then taking a Taylor expansion on NAG of which they matched the coefficients up to the second order terms. PDF [return](#).

ICC-26: The work of [Barakat and Bianchi \[2018\]](#) was based on a asymptotic analysis on the Adam parameter to treat parameteric updates as a differential process. Similar to the work of [Qian \[1999\]](#), they also sought for Lyapunov functions for proving the existence of minima. Their continuous analysis found that there did not necessarily exist a Lyapunov function for Adam updates; and their discrete time analysis found that Adam could also certainly not converge unless with the decrease of learning rates. From there, further analysis on an Adam optimiser with decreasing learning rates was conducted.

The daunting equation (2.1) is just a concatenate of x_n , m_n , and v_n of their Algorithm 2.1. Note, [Barakat and Bianchi \[2018\]](#) used γ to represent for the learning rate, and used α and β to represent for the 1st and 2nd order moments. PDF [return](#).

ICC-27: The main focus of the PolyNet paper [[Zhang et al., 2017b](#)] was dedicated towards finding different compositions of component $\mathcal{H}(n_{L(\gamma-1)}, W_{L\gamma})$ for the ResNet design of Equation (2.1). They provided the perspective such that any ResNet connection of

$$\mathbf{n}_{L\gamma} = \mathbf{n}_{L\gamma-1} + \mathcal{H}(\mathbf{n}_{L(\gamma-1)}, \mathbf{W}_{L\gamma})$$

could instead be expressed as

$$\mathbf{n}_{L\gamma} = (\mathbb{I} + \mathcal{M})\mathbf{n}_{L(\gamma-1)}$$

where \mathbb{I} denotes the identity operator and \mathcal{M} as some arbitrary transform. Under this perspective, \mathcal{M} no longer needs to be $\mathcal{H}(n_{L(\gamma-1)}, W_{L\gamma})$ of which originally represented a single sequence of non-linear transform. Instead \mathcal{M} could now represent some linear combination of second-order transformation of $\mathcal{H} + \mathcal{H}(\mathcal{H})$. The au-

thors of that paper named this new composition *(m)poly-2*. For different variants of *(m)poly-2*, refer to Figure 4 on page 4 of [Zhang et al. \[2017b\]](#). PDF [return](#).

ICC-28: The main focus of the Shake-Shake paper [[Gastaldi, 2017](#)] was to increase generalisation in ResNet via network modification. They changed the original ResNet of Equation (2.1)

$$\mathbf{n}_{L_\gamma} = \mathbf{n}_{L_{\gamma-1}} + \mathcal{H}(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma})$$

into

$$\mathbf{n}_{L_\gamma} = \mathbf{n}_{L_{\gamma-1}} + \omega_{L_\gamma} \mathcal{H}_1(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma}^{(1)}) + (1 - \omega_{L_\gamma}) \mathcal{H}_2(\mathbf{n}_{L_{(\gamma-1)}}, \mathbf{W}_{L_\gamma}^{(2)})$$

where for every layer, a variable ω_{L_γ} was to be randomly sampled from uniform distribution of the range of $[0, 1]$ during training. Then during testing, each ω_{L_γ} was then set as 0.5. PDF [return](#).