

Thoughts on the Future of Scientific Dissemination

Philip E. Bourne
University of California San Diego
pbourne@ucsd.edu

eResearch Australia June 2007

“So Much Done and So Much to Do”

Winston Churchill, upon being
accused of being slightly tipsy by
a member of the temperance
league

eResearch Australia June 2007

“So Much Done and So Much to Do”

Winston Churchill, upon being
accused of being slightly tipsy by
a member of the temperance
league

eResearch Australia June 2007

Motivational Questions

- Have publishers have yet to realize the power of the Internet?
- Is a database really different from a journal? *PLoS Comp. Biol.* 2005 1(3), e34
- Does a scientific society meet the needs of its members?
- Are scientists taking full advantage of the digital world?

eResearch Australia June 2007

Motivational Questions

- Have publishers have yet to realize the power of the Internet? **NO**
- Is a database really different from a journal? **NO** *PLoS Comp. Biol.* 2005 1(3), e34
- Does a scientific society meet the needs of its members? **NO**
- Are scientists taking full advantage of the digital world? **NO**

eResearch Australia June 2007

Background

- As a computational biologist this is a view from the trenches
- This is a view from the life sciences – while disciplines have similarities there are also differences
- I do not actively follow developments in eResearch
- As an EIC of a major biological journal and a developer of a major biological database I have a certain perspective (bias)

eResearch Australia June 2007

Agenda

- Introduce a research vision
- Diagnose that vision
- Illustrate what we are doing to address this vision
- Invite comment

eResearch Australia June 2007

The Vision...

Prior to leaving home the graduate student syncs her IPOL with the latest video papers delivered overnight by the journal via RSS feed. On the bus she reviews the stream, selecting a paper close to her interest in HIV-1 proteases. The data shows apparent anomalies with her own work. She notices that her colleague has also discovered the same paper and they IM annotating the results. By the time the bus stops she has recomputed the results, proven the anomaly and written a rebuttal, included a podcast as a letter to the Editor and sent it to the journal

eResearch Australia June 2007

Science Fiction?

- Five years ago Yes... Today No...
- Five years ago the idea of downloading data on a bus would have been absurd – not today
- Five years ago an IPOL would be absurd - not today
- Journals are providing RSS feeds today
- IM is prevalent but not for scientific discourse
- Video and podcasting is prevalent but not for scientific discourse
- Why should the way we do science not change in the next five years?

eResearch Australia June 2007

What is Missing to Make the Vision a Reality?

1. Seamless integration between the data and the publication upon which that data are based
2. Seamless integration of the authoring and publishing process
3. Notion of traditional publications being associated with podcasts and video
4. Professional networking akin to social networking

eResearch Australia June 2007

What are the Catalysts for Change?

- Open access publishing
- The emerging generation of digital scientists
- The increased ease of working with digital media, notably sound and video

eResearch Australia June 2007

What is Missing to Make the Vision a Reality?

1. Seamless integration between the data and the publication upon which that data are based
2. Seamless integration of the authoring and publishing process
3. Notion of traditional publications being associated with podcasts and video
4. Professional networking akin to social networking

eResearch Australia June 2007

1. Database and Journal Integration- The Test Bed



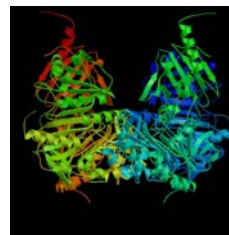
Journals

<http://www ww pdb.org/>

RCSB **PDB**
PROTEIN DATA BANK

WORLDWIDE
ww PDB
PROTEIN DATA BANK

eResearch Australia June 2007



Database

The PLoS Corpus



- Established in 2000
- Identified as a high quality publications (*PLoS Biology* impact factor 14.7)
- Currently 8 journals with healthy growth
- Open Access – free to all

It is the last point that makes this work a reality

eResearch Australia June 2007

Open Access (Creative Commons License)

1. All published materials available on-line free to all (author pays model)
2. Unrestricted access to all published material in various formats eg XML provided attribution is given to the original author(s)
3. Copyright remains with the author



eResearch Australia June 2007

Open Access (Creative Commons License)

1. All published materials available on-line free to all (author pays model)
2. **Unrestricted access to all published material in various formats eg XML provided attribution is given to the original author(s)**
3. Copyright remains with the author



The catalyst

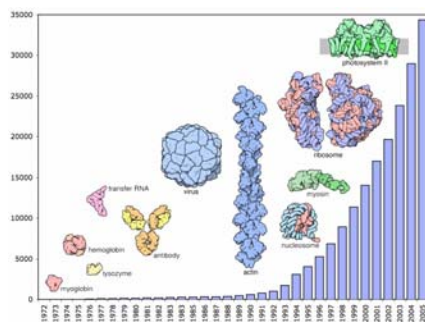
eResearch Australia June 2007

The PLoS Corpus – Under the Hood

- Conforms to the NLM DTD – little markup of content
- Parallel development of Topaz – different emphasis – manuscript and content management with reusable software and backend infrastructure

eResearch Australia June 2007

The Protein Data Bank

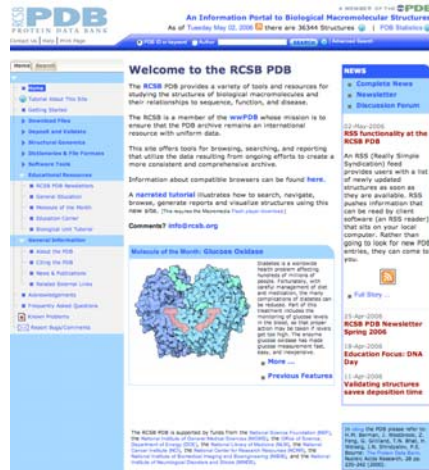


- The single worldwide repository for data on the structure of biological macromolecules
- Vital for drug discovery and the life sciences
- Over 30 years old
- Free to all

Nucleic Acid Research 2000 **28(1)**, 235 – 242 – 5000 citations

eResearch Australia June 2007

The Protein Data Bank



<http://www.pdb.org>

eResearch Australia June 2007

- Paper not published unless data are deposited – strong data to literature correspondence
- Highly structured data conforming to an extensive ontology
- DOI's assigned to every structure

1. Seamless Integration between Data and the Literature – What Does That Imply?

- Improving semantic consistency in the literature – best done at the point of authoring
- Post processing to establish semantic content
- New forms of visualization and interaction at the presentation layer

eResearch Australia June 2007

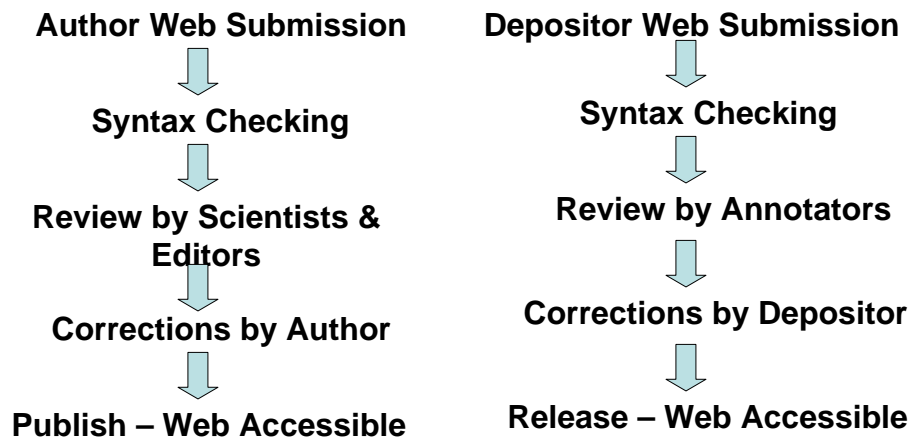
I argue:

This really should not be that difficult since journal processes and content are not that different to database processes and content

eResearch Australia June 2007



Similar Processes Lead to Similar Resources



eResearch Australia June 2007

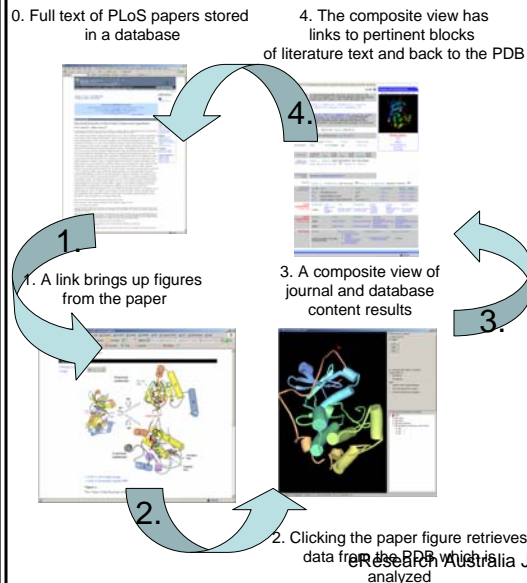
Can it Happen?

- The data repositories have been available for a long time
- With open access the knowledge repositories are available – its more than just abstracts
- If the perception of the difference between data and knowledge is lowered
- The technologies are there

eResearch Australia June 2007

BioLit: Tools for New Modes of Scientific Dissemination

The Knowledge and Data Cycle



- **Biolit integrates biological literature and biological databases and includes:**
 - A database of journal text
 - Authoring tools to facilitate database storage of journal text
 - Tools to make static tables and figures interactive

eResearch Australia June 2007

PDB-PLoS Journal Integration – What has Been Done so Far?



Lynn Fink

- Retrieve full XML files of all articles in PubMed Central
- Parse into database
 - Extract DOIs, captions for figures
 - Extract PDB IDs from article and figure text
 - Use NLP algorithms to find PDB structure properties (i.e., mutations)
 - Extract MeSH terms, EC numbers



←Research Australia June 2007→



ALERT: Our data files are changing soon. Please see <http://www wwpdb.org> for more details.

1mu2 DOI 10.2210/pdb/1mu2/pdb

Title CRYSTAL STRUCTURE OF HIV-2 REVERSE TRANSCRIPTASE

Authors Ran, J., Bird, L.E., Chamberlain, P.P., Stewart-Jones, G.B., Stuart, D.I., Stammers, D.K.

Primary Citation Ran, J., Bird, L.E., Chamberlain, P.P., Stewart-Jones, G.B., Stuart, D.I., Stammers, D.K. Structure of HIV-2 reverse transcriptase at 2.35-Å resolution and the mechanism of resistance to non-nucleoside inhibitors. *Proc Natl Acad Sci USA* 99 pp. 14410-14415, 2002. [Abstract]

Additional Literature Show figures from secondary citations

History Deposition 2002-09-23 Release 2002-10-30

Experimental Method Type X-RAY DIFFRACTION Data [EDS]

| Parameters | Resolution(Å) | R-value | R-Free | Space Group |
|------------|---------------|--------------|--------|--|
| | 2.35 | 0.192 (work) | 0.240 | P 2 ₁ 2 ₁ 2 ₁ |

Display Options
RMSD
3mol
WebMol
My Stuff
QuickPDB
All Images

eResearch Australia June 2007

[CONTACT US](#) | [HELP](#) | [PRINT PAGE](#) | | | |

ALERT: Our data files are changing soon. Please see <http://www wwptdb.org> for more details.

[Home](#) | [Search](#) | [Structure](#) | [Quickstart](#) | [PDB ID or keyword](#) | [Author](#) | [Site Search](#) | [Advanced Search](#)

As of Tuesday Jun 19, 2007 there are 44191 Structures | [PDB Statistics](#)

[1mu2](#) | [Download Files](#) | [FASTA Sequence](#) | [DOI 10.2210/pdb/1mu2/pdb](#) | [Images and Visualization](#)

[Biological Molecule / Asymmetric Unit](#)

Why Do HIV-1 and HIV-2 Use Different Pathways to Develop AZT Resistance?
 Royer P, Santalucia DJ, Clark PK, Arnold E, Hughes MD
 PLoS Pathog. 2006 Feb;2(2):e110. Epub 2006 Feb 17.
[Full paper](#)

Figure 6:
 In the HIV-1 RT complex (PDB code 1HQJ), shown in yellow, the 3' end of the primer is in the N site, unliganded HIV-2 RT (PDB code 1MU2) is shown in magenta. The residues used for the superposition were 107-112 and 155-215. The Van der Waal radii for the atoms in the amino acid residues at positions 117 and 214 are shown as spheres: yellow for HIV-1 RT and blue for HIV-2 RT. The superposition highlights two major differences between the two enzymes: a) the interaction between the N terminus and residue Ser117 of p66 in HIV-1 RT (red dotted line) is absent in HIV-2 RT because in HIV-2 RT the N terminus has moved away from residue 117 (cyan dotted line); b) in HIV-1 RT, Ser117 interacts with Leu214 in a way that differs from the interaction of Ser117 and the bulkier Phe214 in HIV-2 RT.

Figure 7:
 Unliganded HIV-2 RT (shown in magenta; PDB code 1MU2) and HIV-1 RT/DNA/dTTP (shown in white with the dTTP and DNA omitted for clarity; PDB code 1RTD), is superimposed on the Cε protein backbone of HIV-1 RT/DNA/tamofovir-diphosphate (shown in cyan; PDB code 1TOS). The residues used for the superposition were 107-112 and 155-215. The alignment shows that K220 of HIV-2 RT is the residue structurally equivalent to K219 in HIV-1 RT.

Figure 8:
 Unliganded HIV-2 RT (PDB code 1MU2) is shown in yellow, unliganded HIV-1RT (PDB code 1DLO) is shown in cyan, and HIV-1 RT/DNA/dTTP (PDB code 1RTD) is shown in red. The superposition shows that a Phe116Tyr mutation in HIV-1 RT would affect the interactions of position 116 with the main-chain carbonyl of Lys73 in the fingers subdomain of HIV-1 RT (dotted line). However, the same mutation in HIV-2 RT is not likely to cause similar interactions of Tyr116 with the fingers subdomain of HIV-2 because the differences in the position of the small helix that carries residues 115-117 in HIV-2 RT make the interaction unlikely.

Additional Literatures: Structures that have been published in the same paper as this structure:
 1HQJ
 1HQO
 1RTD
 1TOS

eResearch Australia June 2007

What is Missing to Make the Vision a Reality?

1. Seamless integration between the data and the publication upon which that data are based
2. Seamless integration of the authoring and publishing process
3. Notion of traditional publications being associated with podcasts and video
4. Professional networking akin to social networking

Microsoft

- Collaboration with Microsoft to develop Word 2007 plug-in
 - Semantic mark-up using ontologies and controlled vocabularies
 - Facilitate/automate referencing to PDB (and other resources) from manuscript
 - Conversion of manuscript to NLM DTD for direct submission to publisher

eResearch Australia June 2007

Semantic Mark up Using Ontologies

"Also, it has been shown that APH(3')-IIIa has some protein kinase activity, providing a functional link between the APHs and TPK."

gene_ontology.obo

```
id: GO:0004672
name: protein kinase activity
namespace: molecular_function

def: "Catalysis of the transfer of a phosphate group,
usually from ATP, to a protein substrate." [GOC:j]
xref_analog: EC:2.7.1.37
is_a: GO:0016301 ! kinase activity
is_a: GO:0016773 ! phosphotransferase activity, alcohol
group as acceptor
```

Structural Evolution of the Protein Kinase-Like Superfamily, PLoS Comput Biol. 2005 Oct;1(5):e49, Scheeff ED, Bourne PE.

Link to Data Repository

Human PEX5 [PDB:<ext-link ext-link-type="pdb" xlink:href="1FCH">1FCH</ext-link>]

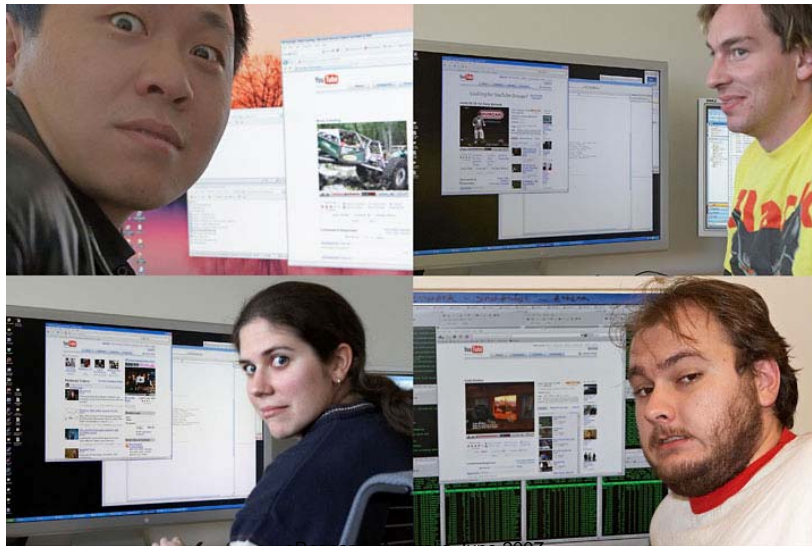
eResearch Australia June 2007

What is Missing to Make the Vision a Reality?

1. Seamless integration between the data and the publication upon which that data are based
2. Seamless integration of the authoring and publishing process
3. **Notion of traditional publications being associated with podcasts and video**
4. Professional networking akin to social networking

eResearch Australia June 2007

Motivation



Motivation

- Scientific understanding requires we digest ever more and diverse information – **16,000** papers are being added to PubMed every week
- The graduate students and postdocs of today are the leading scientists of tomorrow
- This generation is comfortable with the short video clip and sound bite format
- Scientific publishers have been slow to adopt and disseminate sound and video

eResearch Australia June 2007



What is Wrong with Text Only?

- Often times text does not capture the excitement of the work i.e., impersonal, boring
- Time consuming to digest – the abstract is too limiting and the full text too consuming
- Video conveys information not possible with text e.g. a simulation or experimental

eResearch Australia June 2007



Concepts

- Use video and podcasts to enhance dissemination and comprehension
- Leverage existing technology
- Cater to the YouTube generation
- Partner with a respectable scientific publisher known for innovation
- Have the author present only what is in the paper

eResearch Australia June 2007

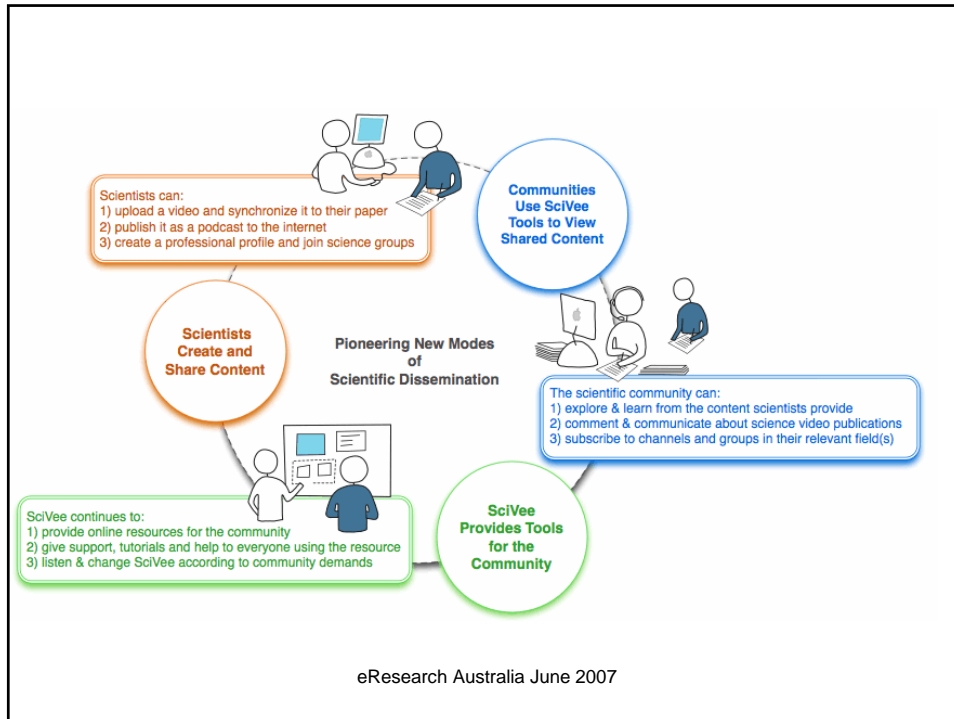


Product

- Video clips and podcasts of presentations by authors of peer reviewed scientific papers – duration 3, 5, 7, 10 minutes
- Later – Other forms of scientific video material – conferences, K12 education etc.

eResearch Australia June 2007



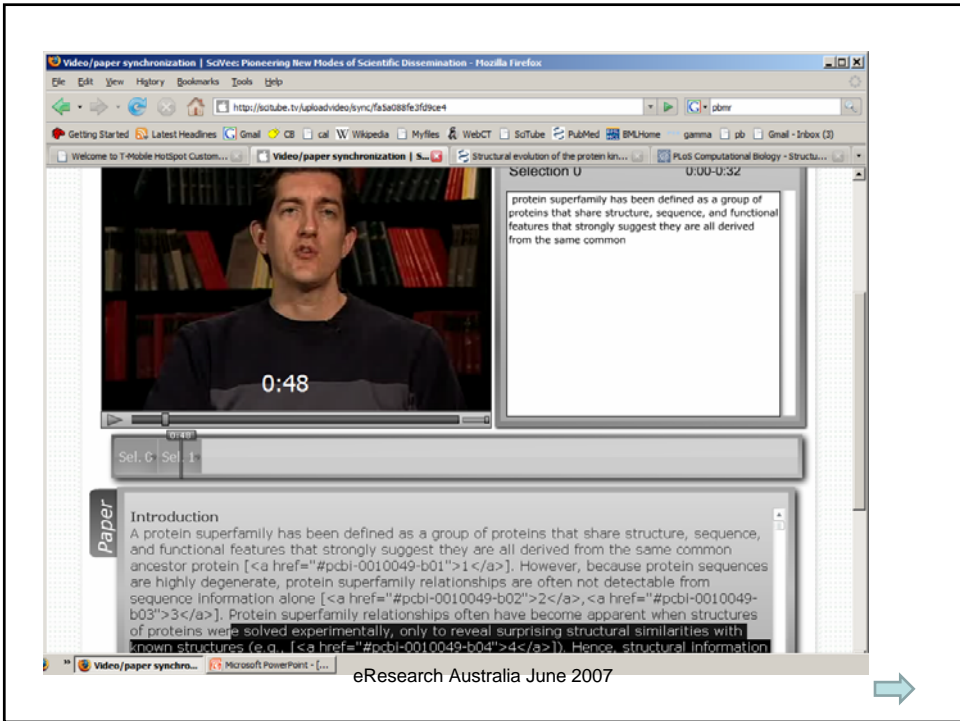
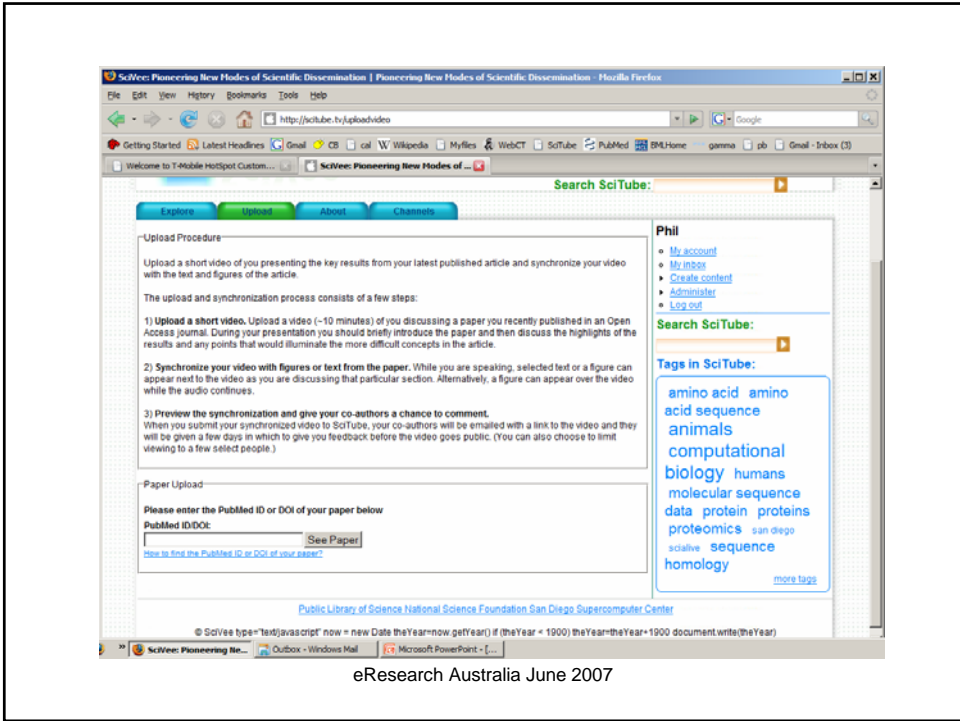


Developmental Phases

- **Phase I** (One Year) – Invite authors of papers published in PLoS journals to upload a video or podcast to *SciVee.tv* describing the motivation, key results and major conclusions of the published study. Establish linkage between literature and video – source of metadata etc. – **September 2007**
- **Phase II** (Years 2 - 3) Scrape PubMed on a daily basis and extend the invitation to authors of all papers in the life sciences; develop video authoring server; provide ratings and virtual community comment
- **Phase III** (Year 4)- Extend to other scientific disciplines







What is Missing to Make the Vision a Reality?

1. Seamless integration between the data and the publication upon which that data are based
2. Seamless integration of the authoring and publishing process
3. Notion of traditional publications being associated with podcasts and video
4. Professional networking akin to social networking

eResearch Australia June 2007

Issues and Solutions

- Science is hierarchical – let the graduate students and post docs be heard – if not be their peers by each other
- How to get profiles to start? – PLoS author/reviewer/editor base
- How to organize? – around labs, organizations, topic areas

eResearch Australia June 2007

Acknowledgements



PLoS for being PLoS



BioLit

Lynn Fink

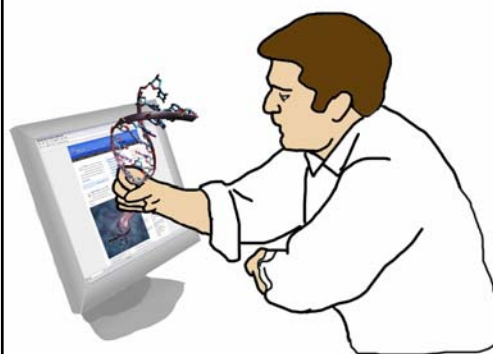


- SciVee Team

- Apryl Bailey
- Diane Baxter
- Leo Chalupa
- Louie Coffman
- Caroline Drakely
- John Matherly
- Alex Ramos
- Willy Suwanto
- Kevin Walsh



eResearch Australia June 2007



Questions?





- Underlying Postgres relational database
 - Fedora (Flexible Extensible Digital Object Repository Architecture).
 - Support for RDF metadata
 - Ajax front end
 - applications will access the repository's data by means of the four APIs by which Fedora is exposed: management, access, search (exposed via HTTP or SOAP) and the OAI provider API (exposed via HTTP).
- 3 Tier Architecture
 - SQL-92 Database (MySQL)
 - J2EE Application Server (JBoss, hibernate)
 - Web Client, Web Services Client, Webworks
 - Lucene indexing engine
 - Standards compliant using available open source tools. Intention is to make the site more mobile while offering easy access and integration to the diverse community.

Deshpande et al., Nucleic Acids Research. 2005 33: D233-D237