

# ZENTITY

---

A repository and a platform for building repositories

Investigations into “Semantic Computing”

**Savas Parastatidis**

Software Architect/Philosopher

Bing (previously with MSR)

blog: <http://savas.me>

twitter: <http://twitter.com/savasp>

facebook: <http://facebook.com/savas>

email: [savas@parastatidis.name](mailto:savas@parastatidis.name)

# AGENDA

Introductions

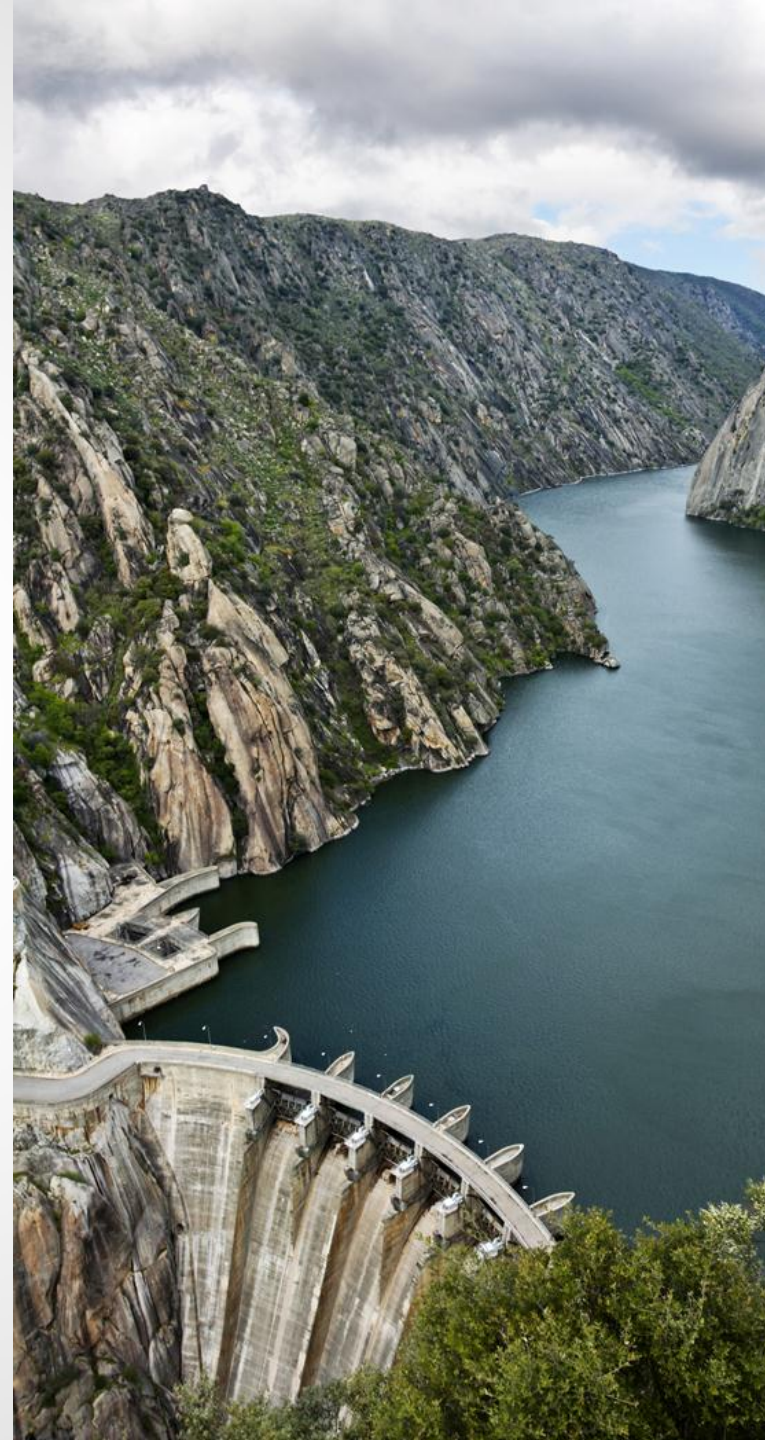
Context of Research Repositories

Subjective view of Semantic Computing concepts

Zentity Architecture

Samples/Demos

But I love to improvise so please ask  
questions along the way!





# DO ASK QUESTIONS

Unfortunately, I have to leave tonight to go to this...



# WHO ARE WE?

## **Microsoft Research**

### **External Research**

Tony Hey (Corporate Vice President)

Daron Green (Senior Director)

### **Education and Scholarly Communications**

Lee Dirks (Director)

Alex Wade (Director)

Pablo Fernicola (Group Program Manager)



# MISSION

**Optimize and extend Microsoft software to meet the specific needs of the academic community**

Our approach:

Conduct applied projects to enhance academic productivity by evolving Microsoft's scholarly communication offerings

Microsoft External Research is uniquely positioned to drive this initiative across Microsoft





Community and Geographic Outreach

Core Computer  
Science



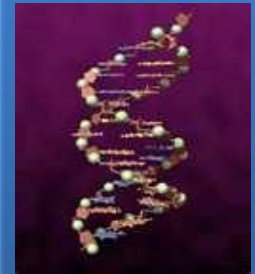
Earth, Energy and  
Environment



Education &  
Scholarly  
Communications



Health &  
Wellbeing



Advanced Research Tools and Services

# WHO WE WORK WITH

Academics / Scholars

Researchers / Scientists

Libraries / Archives

Scholarly Publishers & Societies

Governments / Policy Organizations





# OUR BUSINESS MODEL?

Dual benefit



**WHO AM I?**



# WHO AM I?

## Background

Parallel/distributed computing

eScience

Web Services

Web

Grid

Cloud

Semantic Computing



# WHO AM I?

Now with Bing (doesn't it look good? :-)) to work on  
discovery of information using semantics

Previously an Architect in External Research

**Responsible for projects like:**

Ontology plugin for Word

Chem4Word

Creative Commons plugin

and of course



# **RESEARCH OUTPUT REPOSITORIES**

---



## **Journal articles**

✓ peer-review, indexed, archived

✗ timeliness, cost, access, format limits

## **Digital Repositories**

### **Subject Repositories**

arXiv.org (Physics, Math, CS)

PubMed Central (Biomedical)

### **Institutional Repositories**

### **Data repositories**

Data sets, presentations, workflows, etc.





**SEMANTIC COMPUTING**

**DOES NOT IMPLY**

**RDF, OWL, ETC.**

---





# SEMANTICS

Term used to refer to the concept of “meaning”

The linguistics, AI, Natural Language Processing, etc. communities have been working on “meaning” and “knowledge” related technologies for decades

Emergence of a new breed of technologies to capture meaning (RDF, OWL, etc.)



# WHAT IS SEMANTIC COMPUTING?

Set of concepts/technologies/approaches

Data modeling, structured data, relationships, ontologies

Machine learning (entity extraction)

Inference, reasoning

...





# SET OF TECHNOLOGIES TO...

Model data and their connections  
e.g. RDF, Topic Maps, Unified Content Descriptors

Capture concepts and their relationships  
e.g. OWL

Query data and produce information  
e.g. SPARQL

Reason about data, concepts, information  
e.g. Pellet

Extract structured information (machine learning)  
e.g. Live Labs entity extraction

# TODAY...

Computers are  
great **tools** in



huge amounts  
of **data**

For example, Google and Microsoft both have copies of the Web for  
indexing purposes

# TOMORROW...

Computers are  
great **tools** in



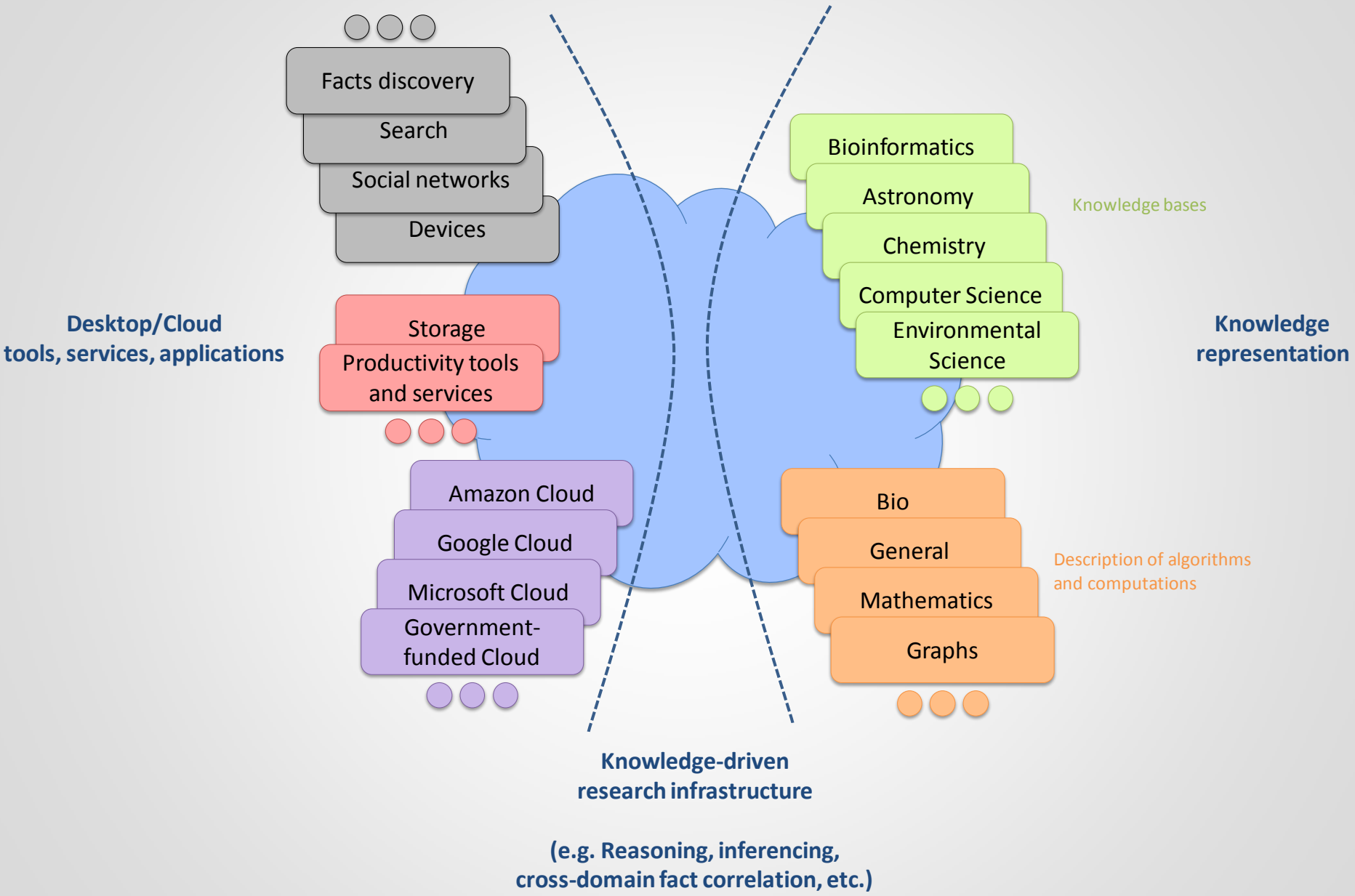
huge amounts  
of **data**

We would like  
computers to also  
help with the  
**automatic**



of the world's  
**information**





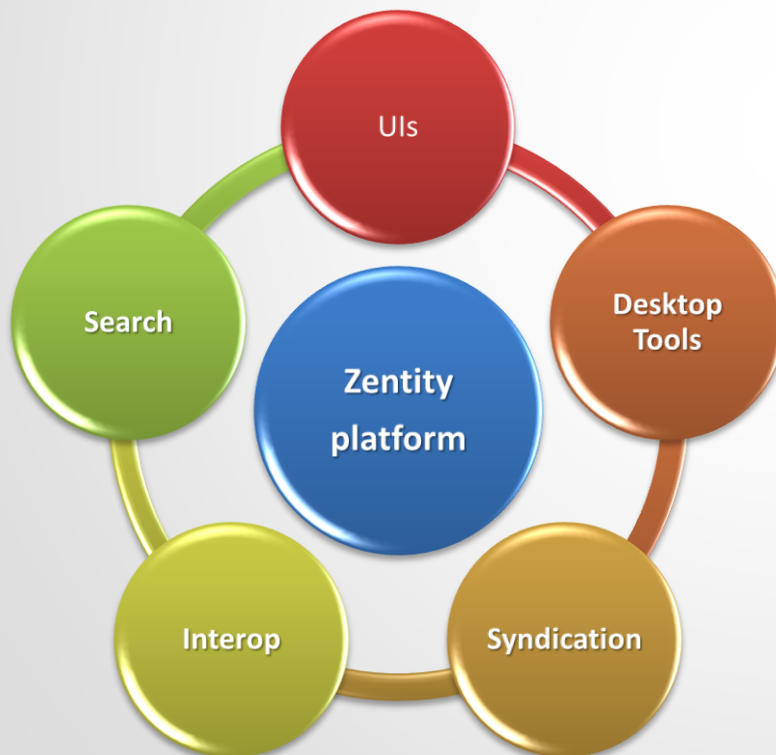
Zentity

---

A platform for building services and tools for research  
output repositories

Papers, Videos, Presentations, Lectures, References,  
Data, Code, etc.

### **Relationships between stored entities**





# GOALS

Enable a tools and services ecosystem for “research output” repositories on MS technologies

Investigate concepts related to semantic computing



# NON-GOALS

Support the lifecycle of publications

Compete with existing repository solutions





## **Modeling**

RDFs – RDF Schema

Custom vocabularies (e.g. XML-based)

Direct use of extensibility API

## **Syndication and Re-Use**

RSS/Atom

OAI-PMH – Protocol for Metadata Harvesting

OAI-ORE – Object Re-Use and Exchange

## **Ingest & publishing protocols**

SWORD – Simple Web-service Offering Repository

Deposit

AtomPub

BibTex



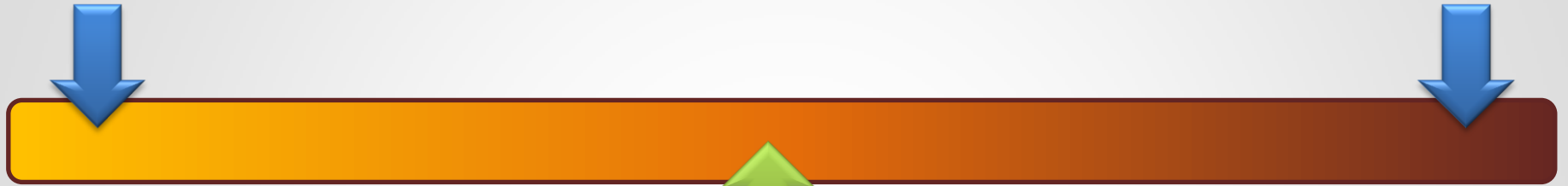
# HYBRID APPROACH

## Triple stores

- Evolution friendly
- Poor performance
- No need to model everything in advance
- Semantic interpretation at the application level

## Relational schema

- Evolution not so easy
- Great opportunities for optimization
- Model everything in advance



## Zentity Platform

- Maintain a balance
- Try to model the frequently used entities in our app domain
- Try to capture the frequently used relationships
- Allow for extensibility (Relationships, Properties)



# KEY DESIGN DECISIONS

Focus on “Resources”

Surface “Relationships” as first-class entities through our API

Model few relationships explicitly for performance reasons (e.g. “contains”, “author”, etc.) for our Scholarly Communications data model

Model key entities explicitly for our Scholarly Communications model

Expose the same functionality to arbitrary data models





# ARCHITECTURE GOALS

Create a platform for building “research output” repositories

An ecosystem of services and tools

Build an easy-to-install collection of basic services and tools

Extensibility

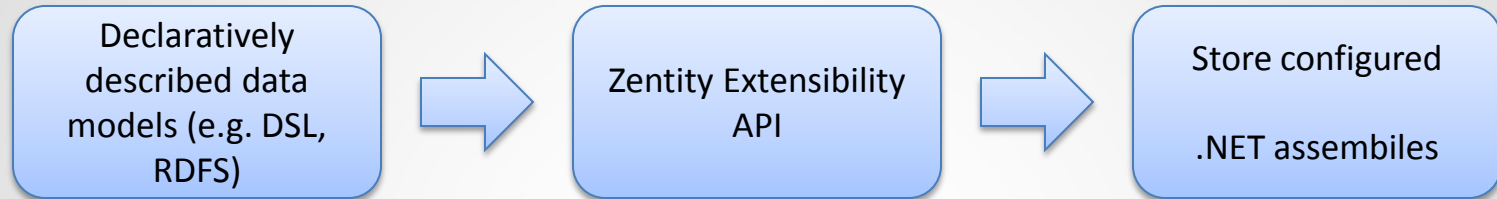
3<sup>rd</sup>-party services, tools,  
applications

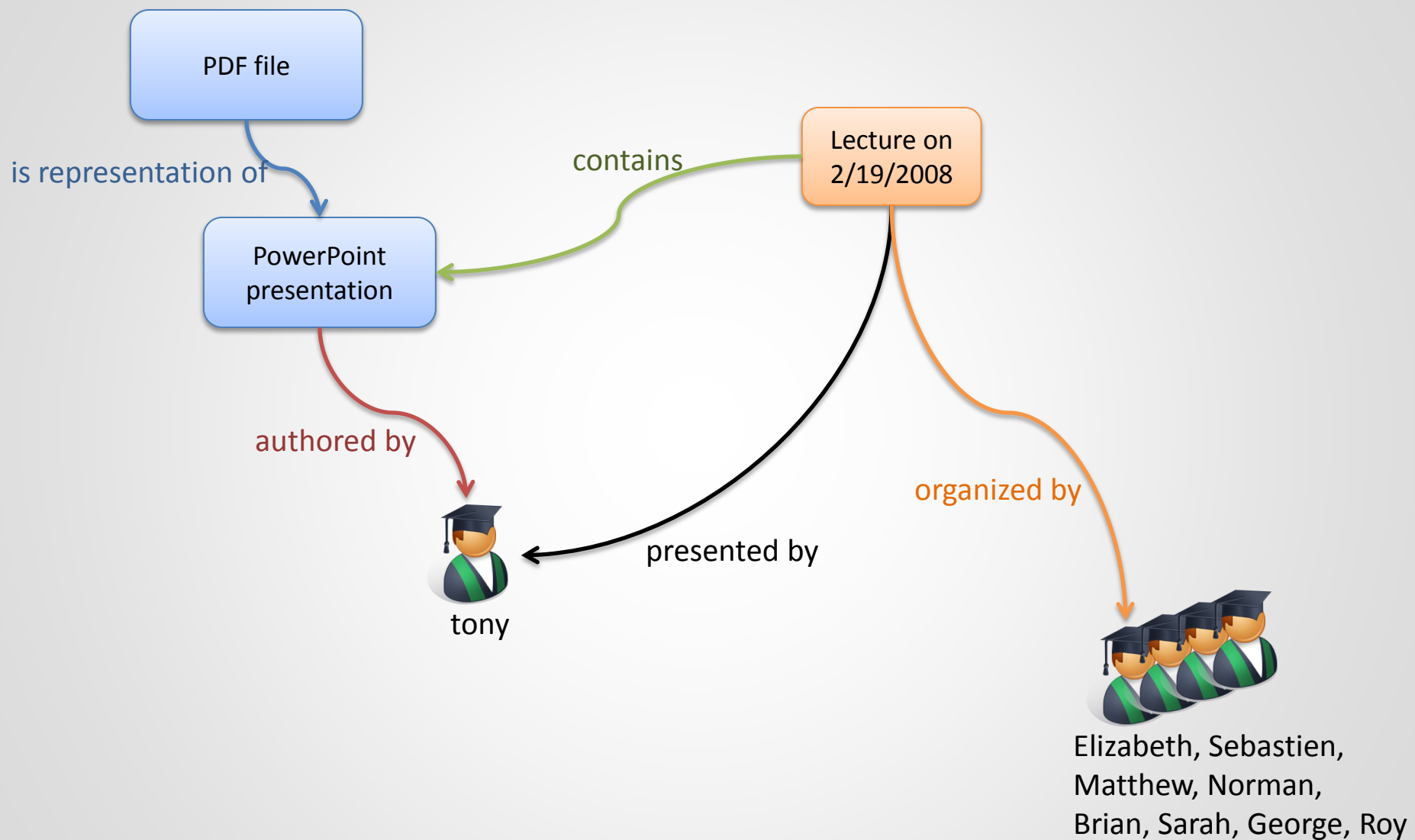
Zentity services, Web  
site, interoperability

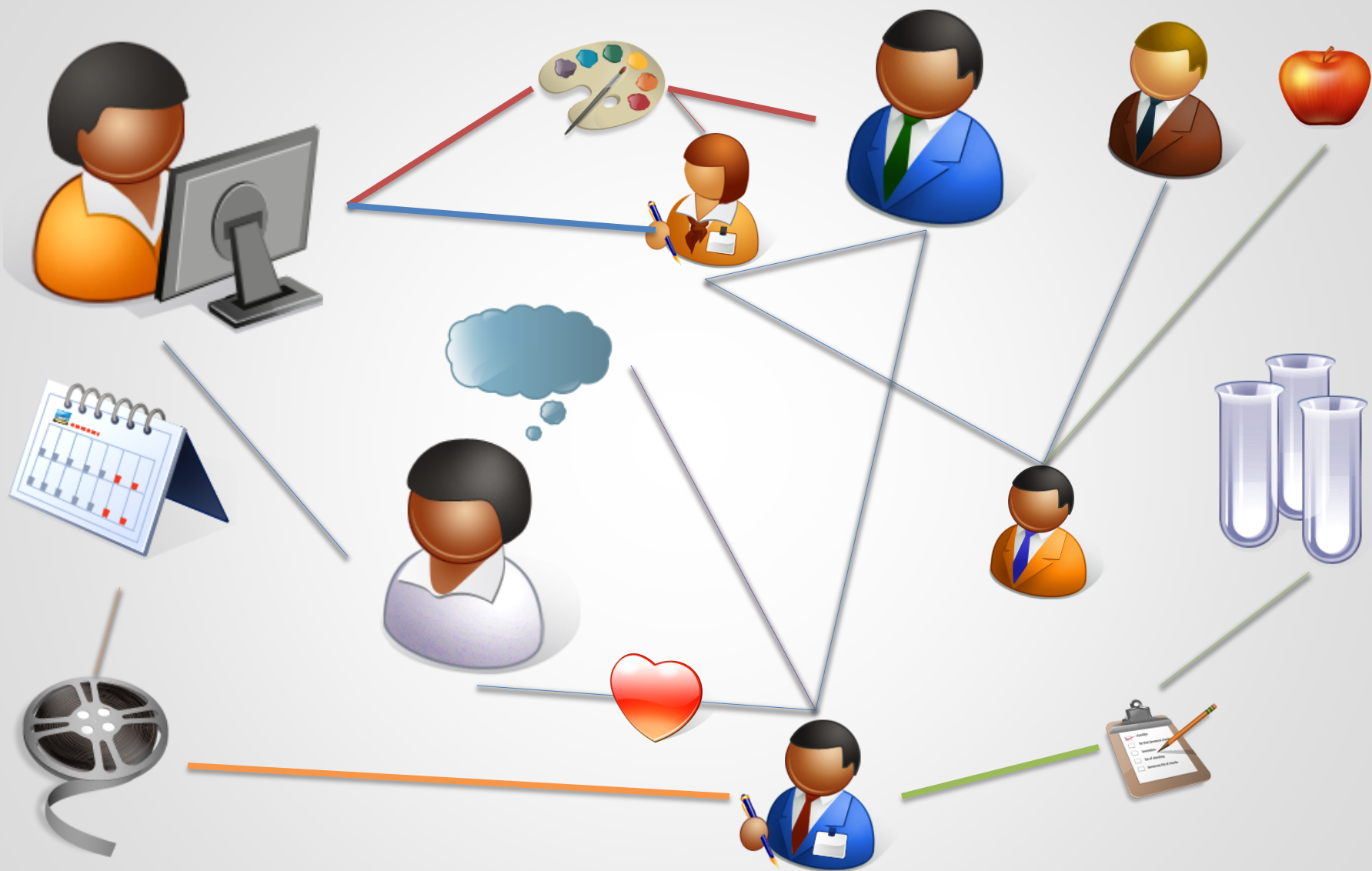
Zentity Platform  
(based on the Entity Framework + semantic data model)

SQL Server 2008, MS data storage technologies, Entity  
Framework, .NET 3.5, LINQ

# EXTENSIBILITY











# Programming Experience

```
Person tony = new Person();
```

```
Publication pub1 = new Publication();  
pub1.Title = "Title1";
```

```
Publication pub2 = new Publication();  
pub2.Title = "Title2";
```

```
pub1.Cites.Add(pub2);  
pub1.Authors.Add(tony);
```

```
Tag tag = new Tag();  
tag.Name = "keyword";  
pub1.Tags.Add(tag);
```

# ROADMAP

CTP1 – Summer 2008

Beta 1 – Autumn 2008

RC1 – Winter 2008-2009

V1 – Spring 2009

Released in May 2009 at the Open Repositories  
conference!!!

Next step: v1.1 and open source  
(provisionally: summer 2009)



# DEMOS

---

On a good day, 50% will NOT break :-)

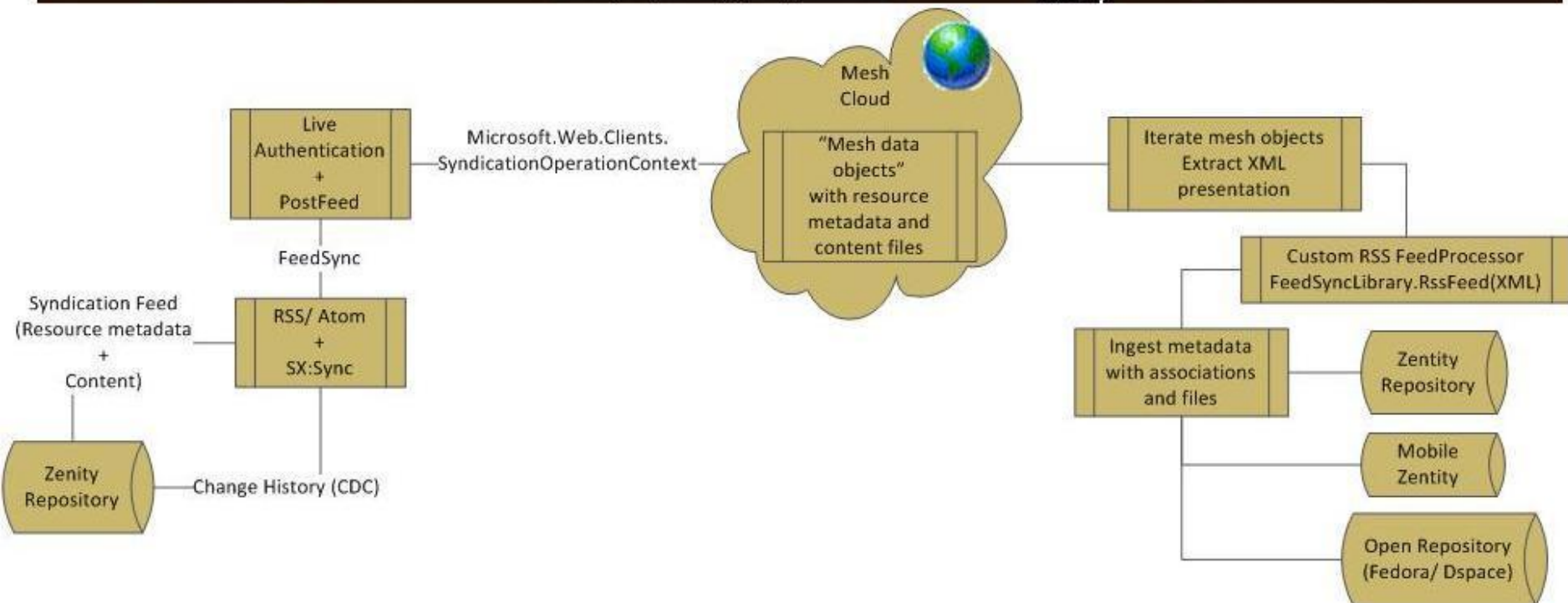
# ***Microsoft***<sup>®</sup>

*Your potential. Our passion.*<sup>™</sup>









Data  
Acquisition &  
Modeling

Collaboration  
and  
Visualization

Analysis and  
Data Mining

Disseminate  
& Share

Archiving and  
Preservation

## Article Authoring Add-in for Word 2007

Document3 - Microsoft Word

Home Insert Page Layout References Mailings Review View Addins Insert Special Chemistry Creative Commons Ontologies

Optional Sections

Title Abstract Introduction Reports Conclusions

Sections

Import Export Journal Panel Author Panel Notes Panel Upload Apply New Template Settings

References Tools

**Single-dose oral naproxen for acute postoperative pain: a quantitative systematic review**

**Abstract**

Naproxen and naproxen sodium are non-steroidal anti-inflammatory drugs used in a variety of painful conditions, including the treatment of postoperative pain. This review aims to assess the efficacy, safety and duration of action of a single oral dose of naproxen/naproxen sodium for moderate to severe acute postoperative pain in adults, compared with placebo.

**Methods**

The Cochrane Library (issue 4 2002), EMBASE, PubMed, MEDLINE and an in-house database were searched for randomised, double blind, placebo controlled trials of a single dose of orally administered naproxen or naproxen sodium in adults with acute postoperative pain. Pain relief or pain intensity data were extracted and converted into dichotomous information to give the number of patients with at least 50% pain relief over 4 to 6 hours. Relative benefit and number-needed-to-treat were then calculated. The percentage of patients with any adverse event, number-needed-to-harm, and time to remedication were also calculated.

**Results**

**Author Panel**

Authors

Lorna Mason  
(lorna.mason@pru.ox.ac.uk); Jayne P. Edwards

Additional Information

Author Notes

CorrespondenceDetails

Correspondence Information

Information regarding single dose Naproxen trials

In one US study, postsurgical pain was the leading...

New Details...

Footnotes

Footnote	Footnote ...	Footnote Type
Collabor...	a	Participating-Researchers

New Footnote ...

Page: 1 of 10 Words: 3,912

90%

# CREATIVE COMMONS ADD-IN FOR OFFICE 2007



Integration with the Creative Commons Web API so that new licenses can be created

Insert Creative Commons licenses from any Office 2007 application

Create Creative Commons License

**Allow commercial uses of your work?**

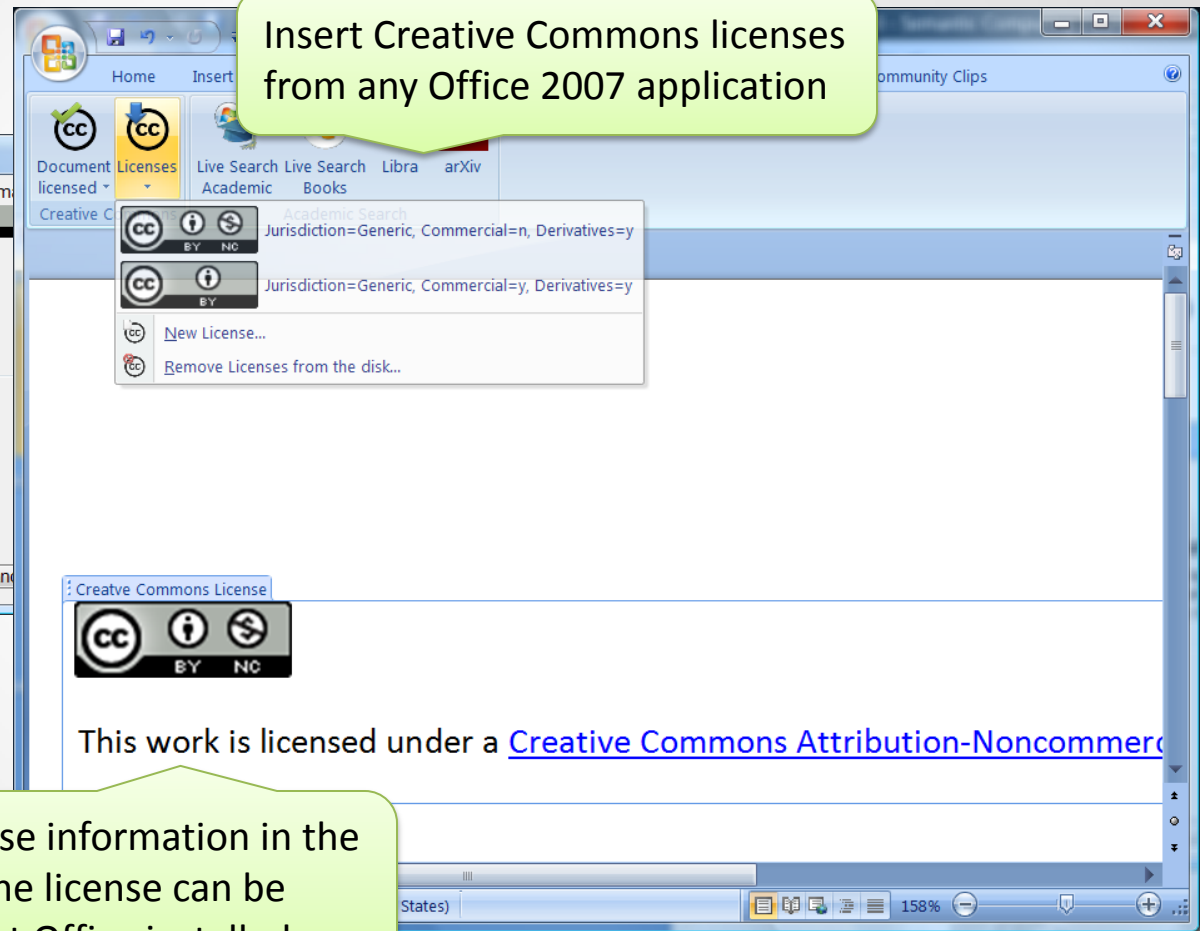
The licensor permits others to copy, distribute and transmit the work. In return, licensees may not use the work for commercial purposes — unless they get the licensor's permission.

Options

☒ Yes

☐ No

< Back   Next >   Cancel



Incorporate license information in the OOXML so that the license can be read even without Office installed



# Ontology Add-In for Word 2007

## Execution

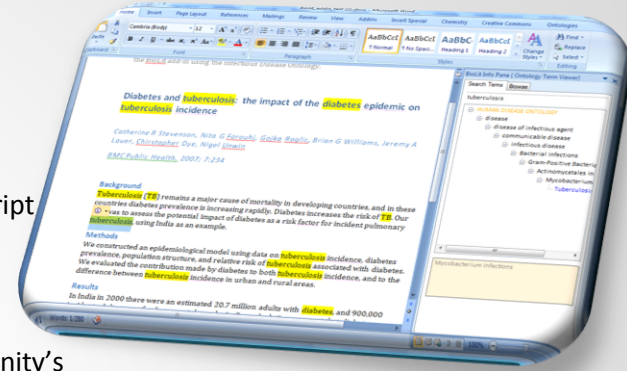
- [Joint announcement](#) between Microsoft Research & Creative Commons at O'Reilly eTech
- Binary and source code available on CodePlex as of 3/11/2009 (<http://ucsdbiolit.codeplex.com/release/>)
- Based on a research project with Dr. Phil Bourne at University of California-San Diego (2008)

## Goals

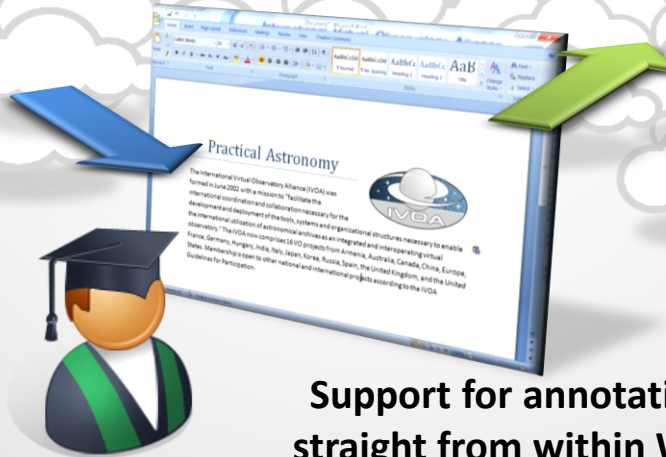
- Facilitate semantic mark-up using ontologies and controlled vocabularies
- Facilitate/automate referencing to PDB, NCBO and other bio-related resources from manuscript

## Scenario

- Authors do not need to be aware of the use of semantic technologies
- A domain-specific ontology is downloaded and made available from within Microsoft Word
- Authors can record their intention, the meaning of the terms they use based on their community's agreed vocabulary



Domain-specific  
ontology



Annotations travel  
with the document

*Can be used to improve domain-specific discovery of information, cross-linking, etc.*



Support for annotations  
straight from within Word





Data  
Acquisition &  
Modeling

Collaboration

Analysis

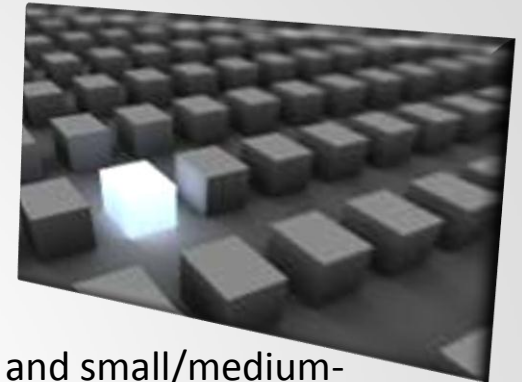
Disseminate  
& Share

Archiving

## eJournal Publishing Service

Hosted Offering for Scholarly Community

- With Microsoft Live Labs & Interoperability



### Goals

- Provide a hosted full-service solution to support scholarly societies and small/medium-sized publishers—simplify self-publishing of workshop/conference proceedings and small journals, as well as online collaboration between authors
- Make broader use of the peer review engine from MSR's Conference Management Tool (CMT) – assisting the review process
- Engage with the publishing ecosystem and enable third-party services (editorial, print-on-demand, etc.)
- Demonstrate interaction between CMT, eJournal Service, and Famulus (and other repositories)

### Execution

- Initial focus on Science / Medicine
- Built on Word 2007, SharePoint, and Office Live
- Beta release in July 2008 – broad range of participants (10+)

Data  
Acquisition &  
Modeling

Collaboration

Analysis

Disseminate  
& Share

Archiving

# Zentity

## Research Output Repository Platform

**A platform for building services and tools for research output repositories:**

- Papers, Videos, Presentations, Lectures, References, Data, Code, etc.
- Relationships between stored entities

### Goals

- Support the MSR publishing and dissemination platform for all researcher outputs
- Enable a tools and services ecosystem for “research output” repositories on MS technologies

### Execution

- Utilizing OAI-ORE, SWORD, and other community protocols
- In development, deployment within MSR in early Q4
- Release to the community in late Q4
- Built on SQL Server 2008 + Entity Framework
  - Using WPF and Silverlight for UI

