# The Fourth Paradigm:
# Data-Intensive Scientific Discovery

Tony Hey
Vice President
Microsoft Research

Microsoft®
**Research** Connections

# Tony Hey – An Introduction

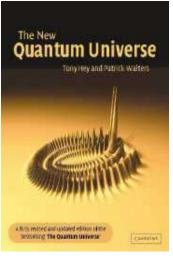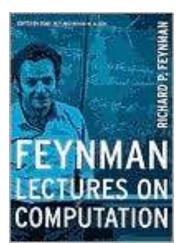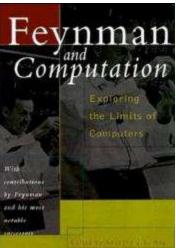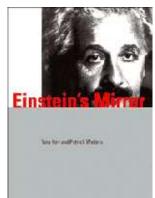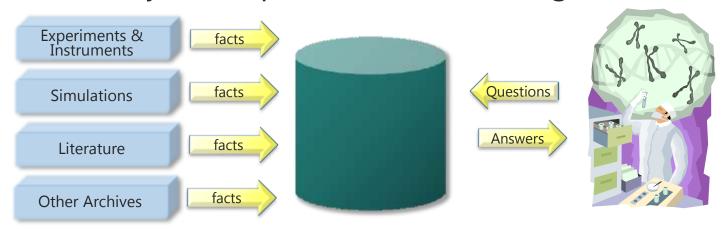# Outline

- The Fourth Paradigm and eScience
- Examples of Data-Intensive Science
- Supporting the Data Life Cycle
- Open Data and Open Science
- The Future: Semantic Computing and the Cloud

# A Tidal Wave of Scientific Data

# X-Info and Comp-X

- The evolution of X-Info and Comp-X for each discipline X
- How to codify and represent our knowledge



| Experiments & Instruments | → facts → | |
| Simulations | → facts → | |
| Literature | → facts → | |
| Other Archives | → facts → | |

← Questions

Answers →

## The Generic Problems

- Data ingest
- Managing a petabyte
- Common schema
- How to organize it
- How to *re*organize it
- How to share with others

- Query and Vis tools
- Building and executing models
- Integrating data and Literature
- Documenting experiments
- Curation and long-term preservation

*(With thanks to Jim Gray)*

# Emergence of a Fourth Research Paradigm

Thousand years ago – **Experimental Science**
- Description of natural phenomena

Last few hundred years – **Theoretical Science**
- Newton's Laws, Maxwell's Equations...
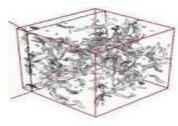
Last few decades – **Computational Science**
- Simulation of complex phenomena

Today – **Data-Intensive Science**
- Scientists overwhelmed with data sets from many different sources
  - Captured by instruments
  - Generated by simulations
  - Generated by sensor networks

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - \mathrm{K}\frac{c^2}{a^2}$$

**eScience is the set of tools and technologies to support data federation and collaboration**
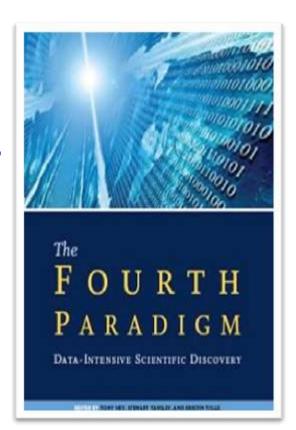- **For analysis and data mining**
- **For data visualization and exploration**
- **For scholarly communication and dissemination**

*(With thanks to Jim Gray)*

# Changing Nature of Discovery
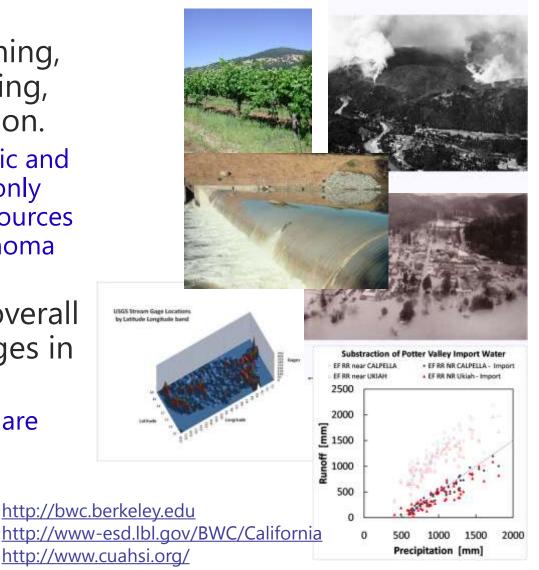
- Complex models
  - Multidisciplinary interactions
  - Wide temporal and spatial scales
- Large multidisciplinary data
  - Real-time steams
  - Structured and unstructured
- Distributed communities
  - Virtual organizations
  - Socialization and management

**http://fourthparadigm.org**



The
F O U R T H
P A R A D I G M
DATA-INTENSIVE SCIENTIFIC DISCOVERY

Microsoft Research Connections

# Digital Watersheds

- Russian River watershed challenges: forestry, farming, urbanization, gravel mining, and fish habitat restoration.
  - Can we understand historic and on-going changes using only publically available data sources such as USGS, NOAA, Sonoma Ecology Center, etc?
- Early studies examined overall water balance and changes in suspended sediment
  - Scientific data "mashups" are leading to new and useful results.

**James Hunt, BWC**

http://bwc.berkeley.edu
http://www-esd.lbl.gov/BWC/California
http://www.cuahsi.org/

# Data from Multiple Sources



Runoff [mm] — USGS

Precipitation [mm] — NOAA

# National Database for Autism Research

Federated sources of data, tools, & specimens from major US autism research funders and investigators
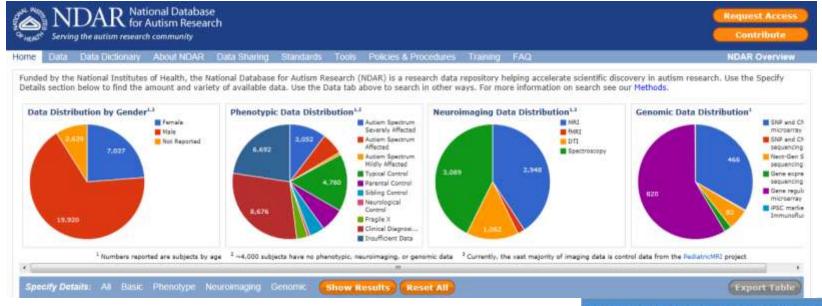
➢ partners adopting NDAR standards, e.g.:
- Global Unique Identifier (GUID)
- Data dictionary (29,000 elements defined)
- Data definition & validation tools
- Authentication scheme

➢ Data from > 85,000 subjects federated

"**Community Science**"

Science@**Microsoft**
THE FOURTH PARADIGM IN PRACTICE

## 10 Years of Inventing the Future

Share: [f] [t] [reddit] [rss]

In the 20 years since it was founded, Microsoft Research has grown from a small group of researchers to approximately 1,000 computer scientists at research labs on four continents. During this growth, the mission of Microsoft Research has remained persistent: to advance the state of the art of computer science and transfer key new technologies into Microsoft products.

In addition to making advances in technologies that can contribute to better products, research at Microsoft also affords an opportunity to use our computer science technologies to help scientists make progress on some of the great problems facing our society. This was the vision for the technical computing initiative we began in 2005, and now we have Microsoft researchers working in collaboration with leading academic researchers throughout the world on a wide range of problems related to health and the environment.

This collection of Science@Microsoft vignettes illustrates some of the progress that has been made in a number of disciplines and describes the technologies that have been deployed to gain these new insights. As can be seen, researchers are effectively applying computer science and technical computing research to fields far removed from computing. With such multidisciplinary research collaborations, Microsoft Research is reducing the time to insight for researchers and accelerating the pace of scientific discovery.

*Senior editors: Tony Hey, David Heckerman, Stephen Emmott*
*Editors: Yan Xu, Kenji Takeda*

**Download**

Science@Microsoft - The Fourth Paradigm in Practice Book (PDF, 10 MB)

**Related Sites**

Microsoft Research 20th Anniversary

Our Research

Microsoft Research Connections

Microsoft Research Connections on Microsoft.com

Microsoft Research

http://www.microsoft.com/en-us/researchconnections/science/stories/

# All Scientific Data Online

- Many disciplines overlap and use data from other sciences.

- Internet can unify all literature and data

- Go from literature *to* computation *to* data *back to* literature.

- Information at your fingertips – For everyone, everywhere

- Increase Scientific Information Velocity

- Huge increase in Science Productivity

**Literature**

**Derived and recombined data**

**Raw Data**

*(From Jim Gray's last talk)*

Microsoft Research Connections

The Data-Intensive Research Lifecycle

# Acquisition & Modeling

# The Cloud - Options



DEDICATED CLOUD

PUBLIC CLOUD

Secure Cloud
Federation

PRIVATE CLOUD

INTERNAL
IT

ENTERPRISE

Microsoft Research Connections

# AzureBLAST

## Seamless Experience

- Evaluate data and invoke computational models from Excel.
- Computationally heavy analysis done close to large database of curated data.
- Scalable for large, surge computationally heavy analysis.
- Test local, run on the cloud.





Scalability of AzureBlast

# AzureMODIS – *Remote Sensing Geoscience*



**5 TB (~600K files) upload of 9 different imagery products from 15 different locations (~6 days of download)**

**4 TB reprojected harmonized imagery ~35000 cpu hours**

**50 GB reduced science variable results ~18000 cpu hours (~14 hour download)**
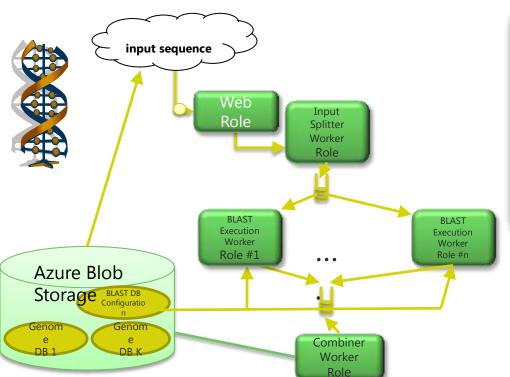
**50 GB additional reduced science analysis results ~18000 cpu hours (~14 hour download)**

Source Imagery Download Sites

Source Metadata

Request Queue

Data Collection Stage

Reprojection Queue

AzureMODIS Service Web Role Portal

Reprojection Stage

Reduction Queue

Scientist

Analysis/Reduction Stage

Scientific Results

# Project JUNIOR

- **Newcastle University, UK**
  **Paul Watson**
  - Investigating applicability of commercial clouds for scientific research
  - Build a working prototype for use-cases in chemo-informatics
  - Uses Microsoft technologies to build science-related services (Windows Azure, Silverlight…)
  - Exploits Azure and Amazon Clouds
- **Built initial proof-of-concept**
  - Silverlight UI for basic Quantitative Structure-Analysis Relationship (QSAR) modeling
  - Demonstrated ability to scale QSAR computations in Windows Azure

# Collaboration & Visualization

# World Wide Telescope
## www.worldwidetelescope.org



# Big data requires new types of data visualization tools

Collaborators:

- Alyssa Goodman; Harvard University
- Alex Szalay; Johns Hopkins University
- Curtis Wong, Jonathan Fay; Microsoft Research

- Integration of data sets and one-click contextual access
- Easy access and use

- Over 4M unique users
- Average number of WWT users is over 8K per day

**See TED talk by Roy Gould and Curtis Wong**
**http://www.youtube.com/watch?v=NPu2j3JVmnw&feature=related**



Microsoft Research Connections

# Layerscape

- Community Site for 'Data Tours'
- Data sharing by groups

# Layerscape with Seismic Data

# Visualization of Models



Tohoku events (shallow) and subduction slab

# Analysis & Data Mining

# Machine Learning and eScience

*Tackling societal challenges*

**Identifying genetic and environmental causes of disease**

**Fighting HIV/AIDS**

**Increasing energy yield of sugar cane through genome assembly**

# Fighting HIV with ML and HPC

- PhyloD.Net is a Bayes-net-based tool that deciphers evolution of HIV within a patient

- Developed by eScience research group and published in *Science*, March 2007

- Used by dozens of HIV research groups

- Analysis requires HPC to do tens of thousands of independent computations

- Integrated into .NET Bio open source library



PhyloD.Net on cover of *PLoS Comp Bio*, Nov 2008
Carlson, Kadie, & Heckerman et al.

*"Web protocol for querying and updating data"*



- Based on HTTP/ATOM
  + Metadata
  + Query options
  + URI conventions
  + JSON representation
- Uses REST semantics
- Open specification
- Submitted to OASIS

https://api.datamarket.azure.com/DataGovUK/MetOfficeWeatherOpenData/

www.odata.org

# Dissemination & Sharing

# The Berlin Declaration 2003

- 'To promote the Internet as a functional instrument for a global scientific knowledge base and for human reflection'

- Defines open access contributions as including:

  - 'original scientific research results, raw data and metadata, source materials, digital representations of pictorial and graphical materials and scholarly multimedia material'

# NSF Data Sharing Policy 2010

*"Investigators are expected to share with other researchers, at no more than incremental cost and within a reasonable time, the primary data, samples, physical collections and other supporting materials created or gathered in the course of work under NSF grants. Grantees are expected to encourage and facilitate such sharing."*

All future grant proposals now require a two-page Data Management Plan that addresses the above requirement and the Plan will be subject to peer review.

# Datacite and ORCID

**DataCite**

- International consortium to establish easier access to scientific research data
- Increase acceptance of research data as legitimate, citable contributions to the scientific record
- Support data archiving that will permit results to be verified and re-purposed for future study.


**ORCID** - Open Research & Contributor ID

- Aims to solve the author/contributor name ambiguity problem in scholarly communications
- Central registry of unique identifiers for individual researchers
- Open and transparent linking mechanism between ORCID and other current author ID schemes.
- Identifiers can be linked to the researcher's output to enhance the scientific discovery process

# Microsoft Academic Search

Microsoft Academic Search is a free academic search engine developed by Microsoft Research Asia, which also serves as a test-bed for our object-level vertical search vertical search research.

- Easily search the top papers, authors, conferences, and journals for a topic.

- See details about a specific paper, author, conference, journal or organization.

- Quickly find relationships between authors.

- Discover influential papers, authors, conferences, journals and organizations within a certain field.

- Get the latest Call for Papers.

**http://academic.research.microsoft.com/**

# Top 10 Computer Science Organizations



Microsoft® Academic Search (Beta)

[ search box ]   🔍   Advanced Search

Author »
Publication »
Conference »
Journal »
Organization »
Keyword »

Academic > Top organizations in Computer Science                    1 - 100 of 5,730 results

[ Computer Science ▼ ]  [ Overall for Computer Science ▼ ]  [ Last 5 Years ▼ ]  [ All Continents ▼ ]

| Organization | Publications | Citations |
|---|---|---|
| Microsoft (H-Index: 285) | 9846 | 37983 |
| Stanford University (H-Index: 365) | 6371 | 26084 |
| Massachusetts Institute of Technology (H-Index: 362) | 6977 | 23939 |
| Carnegie Mellon University (H-Index: 279) | 8379 | 23145 |
| University of California Berkeley (H-Index: 349) | 5804 | 21467 |
| IBM (H-Index: 244) | 7326 | 17166 |
| University of Illinois Urbana Champaign (H-Index: 221) | 6684 | 16700 |
| Georgia Institute of Technology (H-Index: 176) | 5685 | 12749 |
| The french National Institute for Research in Computer science and Control (H-Index: 134) | 4794 | 12358 |
| University of Maryland (H-Index: 210) | 4435 | 11647 |

Microsoft Research Connections

# Public API for Academic Search

- **Application Programming Interface**
  - Supports queries against all academic entities and their basic information
- **With the API, you can**
  - Work with others to share information
  - Help users to build useful clients
- **Openly available to everyone**
  - Targeting the academic community
  - API is available for non-commercial use only

API details at:

**http://academic.research.microsoft.com/About/Help.htm#5**

# The Eigenfactor Project



**http://mas.eigenfactor.org/**

- The *Eigenfactor Project™* is a non-commercial academic research project in the Department of Biology at the University of Washington.
- We aim to use recent advances in network analysis and information theory to develop novel methods for evaluating the influence of scholarly periodicals and for mapping the structure of academic research.
- We are committed to broadly disseminating our research findings and technological developments, while respecting the confidentiality of the data sources we use.
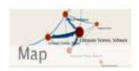


Home | Recommend | Map: Journals | Map: Papers | Explore | Rank | Categorize | About

eigenFACTOR.org

**Recommend** — By uncovering the hierarchical structure of scholarly citation, we can identify key papers pertaining to any search query. For a reader new to the field we can find the classic and foundational papers; for an expert we can find the latest innovations.

**Map** — From patterns of scholarly citation, we use Rosvall and Bergstrom's map equation to chart the topography of science and the relations among fields and subfields. [journal map] [paper map]

**Explore** — By integrating a hierarchical clustering of citation networks with semantic analysis, we develop a scalable map of scientific fields and the key research terms and topics therein.

**Rank** — Scientific influence is often quantified using simple citation counts, but the structure of a citation network provides far more information than can be revealed by these simple counts. This is principle behind the Eigenfactor metrics; we can better rank the importance of scientific journals or papers by viewing them in the context of the full citation network.

APPL PHYS LETT
PHYS REV LETT
ASTRON ASTROPHYS
ASTROPHYS J
MON NOT R ASTRON SOC
APPL ENVIRON MICROB
PLANT CELL
PLANT PHYSIOL
EARTH PLANET SC LETT
GEOPHYS RES LETT
CIRCULATION
NEW ENGL J MED
LANCET
NAT MED
GENE DEV
CELL
J BIOL CHEM
NATURE
SCIENCE
P NATL ACAD SCI USA
NAT GENET
J EXP MED
IMMUNITY
CURR OPIN CELL BIOL
NEURON

# Archiving & Preservation

# Key driver from a UK Research Council

EPSRC Policy Framework on research data (May 2011)

- "all institutions in receipt of their funding should develop a clear roadmap for research data management, which should be implemented by May 1st 2015"

- "organisations will ensure that EPSRC-funded research data is securely preserved for a minimum of 10 years"
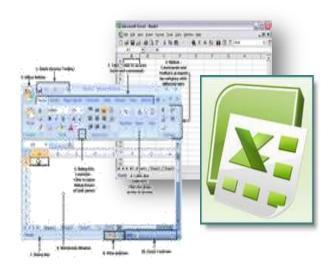
# Data Curation for Excel (DataUp)

- Excel add-in and HTML5 browser app
- Microsoft Research with California Digital Library's Curation Center
- Part of the DataONE (an NSF DataNet Project)

  - Goal to facilitate data management, sharing, and archiving for scientists.
  - Open-source add-in and service for Microsoft Excel to assist in documenting and preparing Excel data for archiving and sharing.
  - Targeting environmental scientists but should be useful for wide audiences
  - Define archiving as movement or storage of data with metadata, to a repository for long-term retention.

**http://dcxl.cdlib.org/**

- Ensure long-term access to Europe's cultural and scientific heritage
  - Improve decision-making about long term preservation
  - Ensure long-term access to valued digital content
  - Control the costs through automation, scalable infrastructure
  - Ensure wide adoption across the user community
  - Establish market place for preservation services and tools
- Build practical solutions
  - Integrate existing expertise, designs and tools
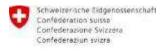  - Share and build

planets

**P**reservation and

**L**ong-term

**A**ccess through

**NET**worked

**S**ervices

# PLANETS Partners

**The British Library**
**National Library, Netherlands**
**Austrian National Library**
**State and University Library, Denmark**
**Royal Library, Denmark**

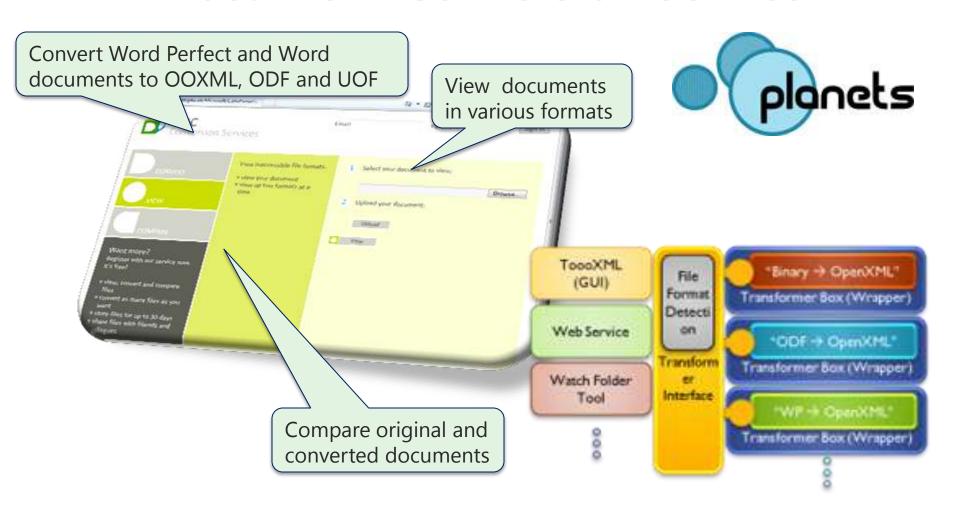**National Archives, UK**
**Swiss Federal Archives**
**National Archives, Netherlands**

**Hatii at University of Glasgow**
**University of Freiburg**
**Technical University of Vienna**
**University at Cologne**

**Tessella Plc**
**IBM Netherlands**
**Microsoft Research, Cambridge**
**ARC Seibersdorf research**

# Archiving and Preservation:
## *A Document Conversion Service*

Convert Word Perfect and Word documents to OOXML, ODF and UOF

View documents in various formats

Compare original and converted documents

planets

ToooXML (GUI)

Web Service

Watch Folder Tool

File Format Detection

Transformer Interface

"Binary → OpenXML" Transformer Box (Wrapper)

"ODF → OpenXML" Transformer Box (Wrapper)

"WP → OpenXML" Transformer Box (Wrapper)
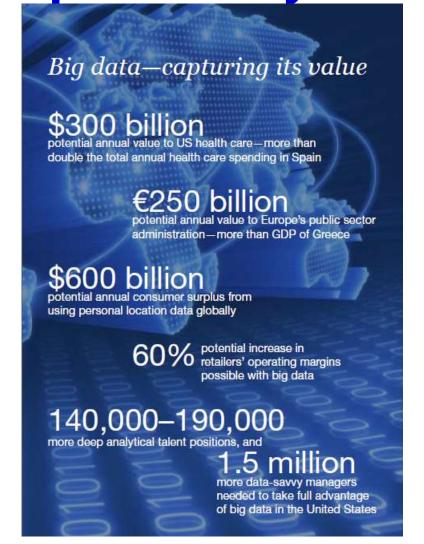
http://odf-converter.sourceforge.net/

# Big data: The next frontier for innovation, competition, and productivity

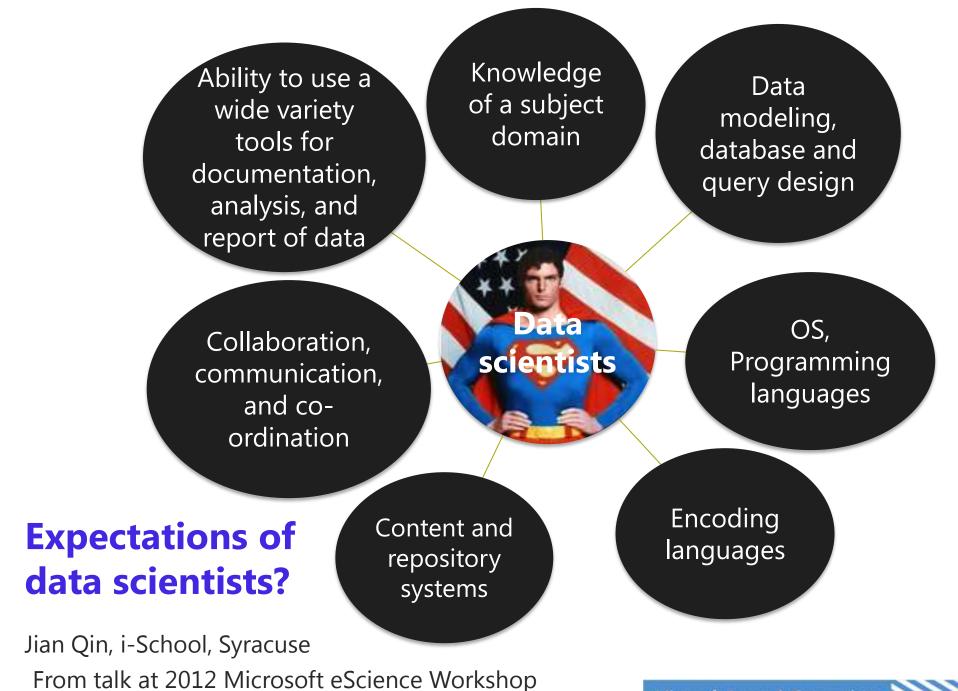

McKinsey Global Institute Report
May 2011

# Educating a New Breed of Data Scientists for Scientific Data Management

**Jian Qin**

**School of Information Studies**
**Syracuse University**

**Microsoft eScience Workshop, Chicago, October 9, 2012**

**Expectations of data scientists?**

Jian Qin, i-School, Syracuse

From talk at 2012 Microsoft eScience Workshop

Microsoft Research Connections

# Semantic Computing?

computers are great **tools** for

| | |
|---|---|
| storing | computing |
| managing | indexing |

huge amounts of **data**

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

| | |
|---|---|
| acquisition | discovery |
| aggregation | organization |
| correlation | analysis |
| interpretation | inference |

of the world's **information** and **knowledge**

**bing**

Community      Home    Blogs    Forums    Media    Events    Toolbox

Site Blogs » Search Blog » Introducing Schema.org: Bing, Google and Yahoo Unite to Build the Web of Objects

## Search Blog

# Introducing Schema.org: Bing, Google and Yahoo Unite to Build the Web of Objects

The Bing Team  6/2/2011 10:01 AM  Comments (3)
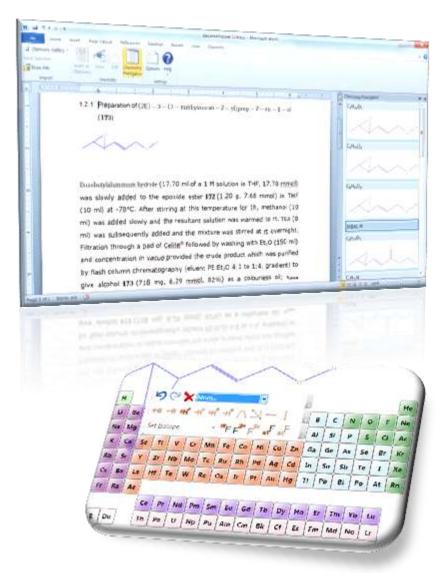
RATE THIS
★★★★★

We've been talking for a while about the need to rethink the search experience to better reflect both the changing web and advancing user habits.

One of the biggest challenges and opportunities we see is to literally create a high-definition proxy of the physical world inside of Bing. In other words, we want to be able to model the world in which we all live to the level that search can actually help you make decisions and get things done in real life by understanding all the options the world presents.

We've made great progress on the technical front to begin to model the real world from the messy bits of data scattered across the web. Things like movies have benefitted from this work. We're now able to understand "Casablanca" is a movie and literally mine the web to re-assemble information about that movie from millions of sites.
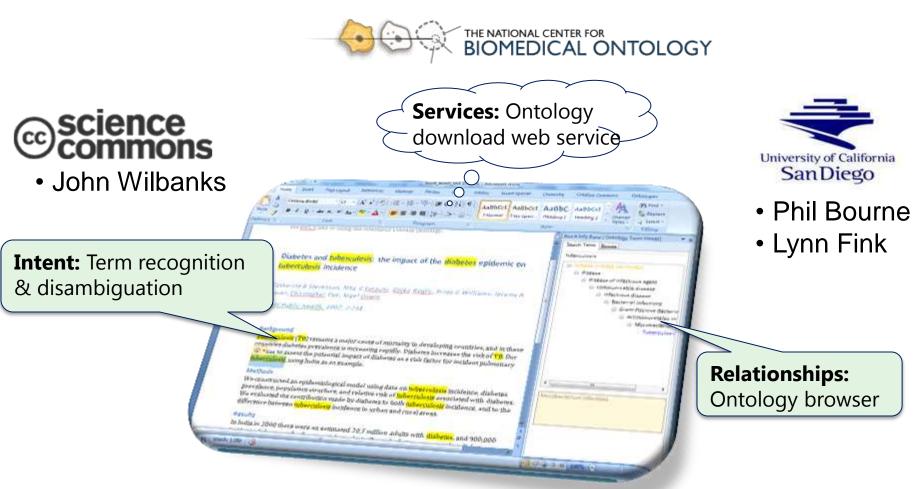
But we think we can do better. We want to enable publishers to give us hints about what things they are describing on their sites. Rather than rely solely on machine learning and other AI techniques, we asked "what if we could enable publishers to have a single schema they could use to describe their sites that all search engines could understand?"

Microsoft Research Connections

# Semantic Chemistry Add-in for Word

- Authoring and rendering of semantic-rich chemical information ([CML](#))

- In partnership with the University of Cambridge

- Support for Office 2007 and Office 2010

- Available under Apache 2.0

- Over **360K** [downloads](#) since March 22nd, 2010

# Ontology Add-in for Word

THE NATIONAL CENTER FOR
**BIOMEDICAL ONTOLOGY**

**Services:** Ontology download web service

science commons

• John Wilbanks

University of California
**San Diego**

• Phil Bourne
• Lynn Fink

**Intent:** Term recognition & disambiguation

**Relationships:** Ontology browser

Source code + binary:
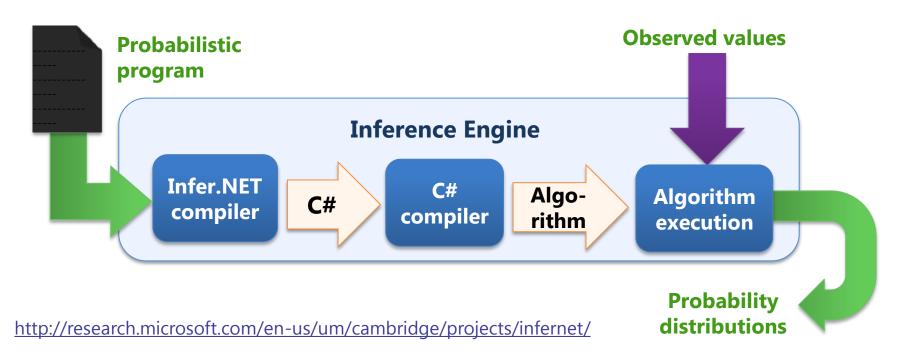http://research.microsoft.com/ontology/

# InnerEye: Semantic Understanding of Medical Images

- InnerEye focuses on the analysis of patient scans using machine learning techniques for automatic detection and segmentation of healthy anatomy as well as anomalies.

- In this image, we see that InnerEye can separate a carotid artery visually from adjacent parts of a human body

# Infer.NET
# Machine Learning  as a Service?

- Provides a probabilistic programming language allowing rapid development of new models
- Goal is to provide a platform for modern applications of Bayesian inference



http://research.microsoft.com/en-us/um/cambridge/projects/infernet/

# Towards a Semantic Cyberinfrastructure?

Future Research Infrastructure will use semantic knowledge services on Client + Cloud



visualization and analysis services

scholarly communications

domain-specific services

search books citations

blogs & social networking

Reference management

Project management

instant messaging

identity

mail

notification

document store

storage/data services

knowledge management

compute services virtualization

knowledge discovery

# Some Resources

- Microsoft Research
  - http://research.microsoft.com
  - Microsoft Research downloads: http://research.microsoft.com/research/downloads
- Microsoft Research Connections
  - http://research.microsoft.com/en-us/collaboration/
- Science at Microsoft
  - http://www.microsoft.com/science
- Scholarly Communications
  - http://www.microsoft.com/scholarlycomm
- CodePlex
  - http://www.codeplex.com
- Outercurve Foundation
  - http://www.outercurve.org/
- Tony Hey on eScience
  - http://tonyhey.net/