# Analyzing the Structure of Genomic Science:
# Mapping the Diffusion of Ideas & Data across Disciplines
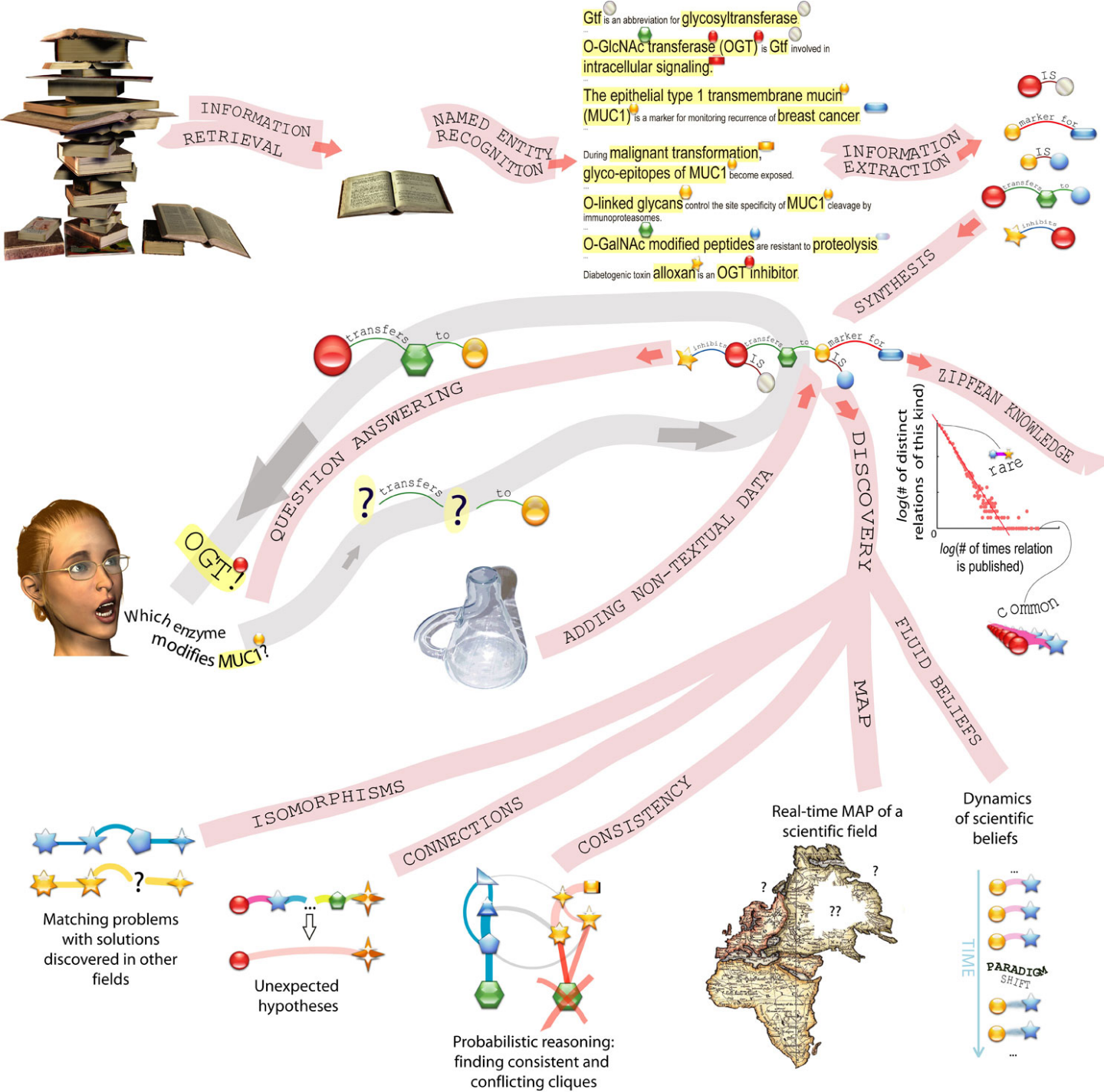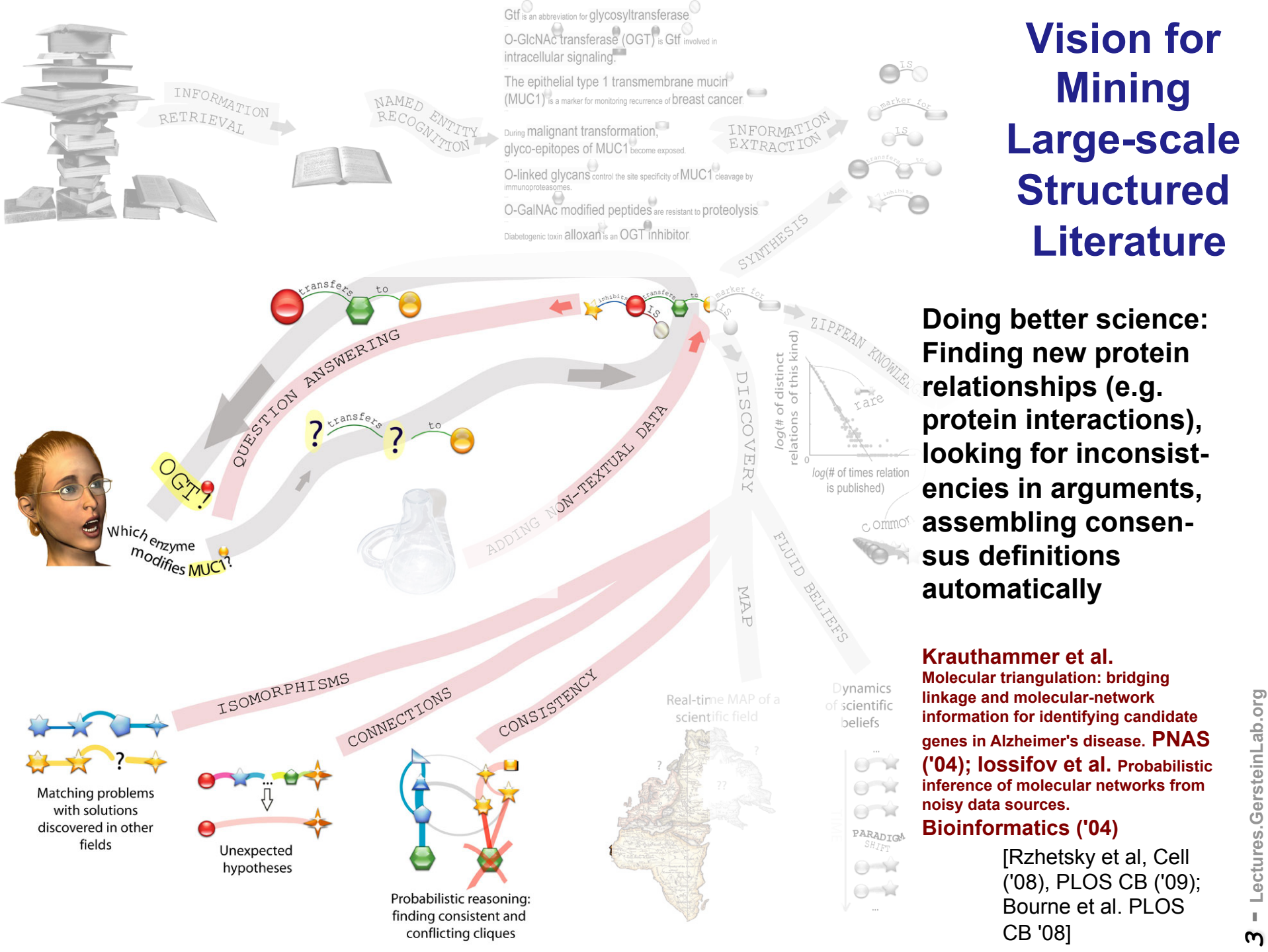
**Mark Gerstein
Yale**

# Vision for Mining Large-scale Structured Literature

Harnessing the "Data Exhaust" from large-scale efforts

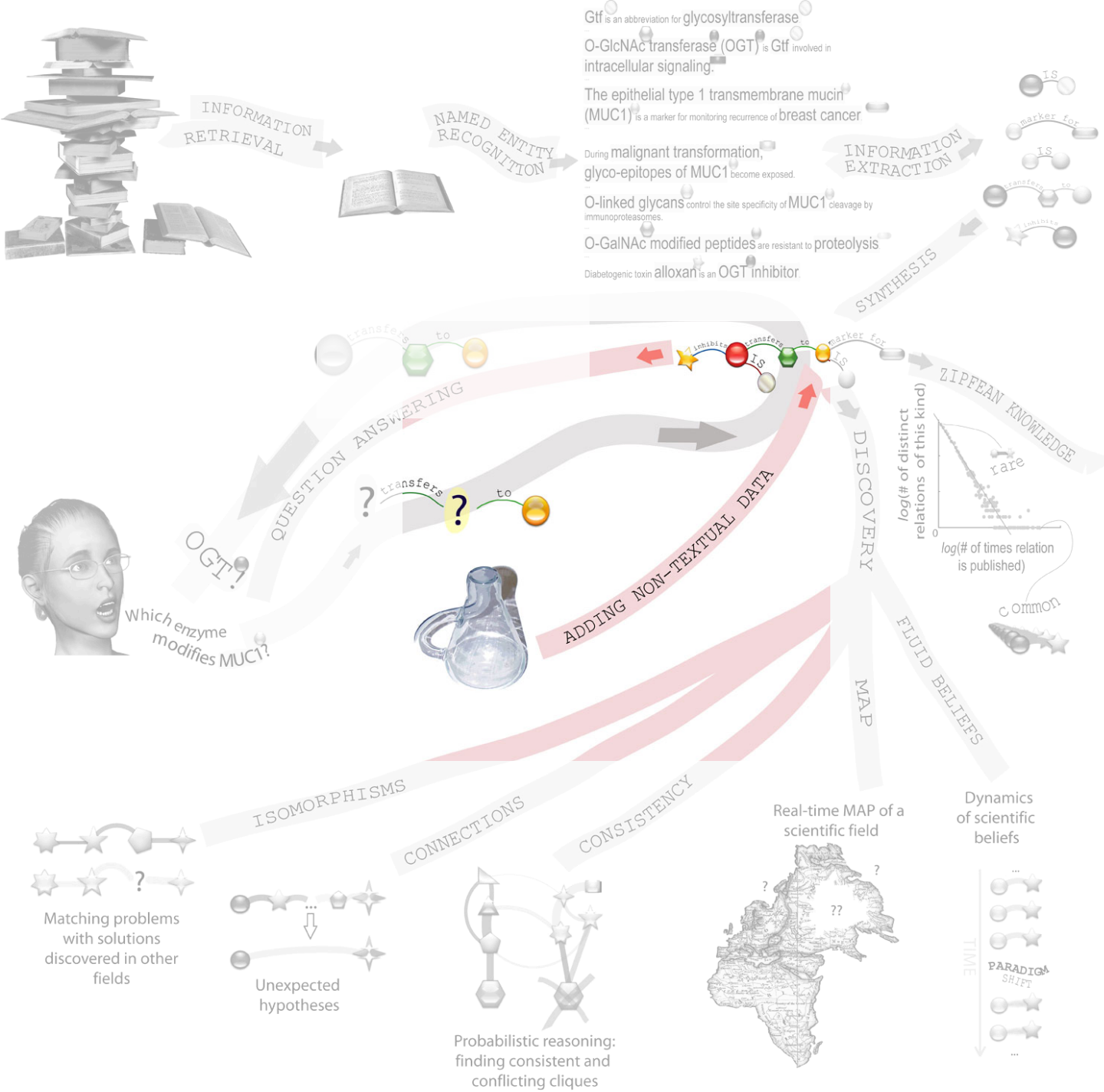[Rzhetsky et al, Cell ('08), PLOS CB ('09); Bourne et al. PLOS CB '08]

# Vision for Mining Large-scale Structured Literature

Gtf is an abbreviation for glycosyltransferase

O-GlcNAc transferase (OGT) is Gtf involved in intracellular signaling.

The epithelial type 1 transmembrane mucin (MUC1) is a marker for monitoring recurrence of breast cancer.

During malignant transformation, glyco-epitopes of MUC1 become exposed.

O-linked glycans control the site specificity of MUC1 cleavage by immunoproteasomes.

O-GalNAc modified peptides are resistant to proteolysis

Diabetogenic toxin alloxan is an OGT inhibitor.

INFORMATION RETRIEVAL

NAMED ENTITY RECOGNITION

INFORMATION EXTRACTION

SYNTHESIS

QUESTION ANSWERING

Which enzyme modifies MUC1?

OGT!

ADDING NON-TEXTUAL DATA

DISCOVERY

ZIPFEAN KNOWLEDGE

log(# of distinct relations of this kind)

rare

log(# of times relation is published)

common

FLUID BELIEFS

MAP

Real-time MAP of a scientific field

Dynamics of scientific beliefs

PARADIGM SHIFT

ISOMORPHISMS

Matching problems with solutions discovered in other fields

CONNECTIONS

Unexpected hypotheses

CONSISTENCY

Probabilistic reasoning: finding consistent and conflicting cliques

**Doing better science: Finding new protein relationships (e.g. protein interactions), looking for inconsist-encies in arguments, assembling consen-sus definitions automatically**

**Krauthammer et al.**
**Molecular triangulation: bridging linkage and molecular-network information for identifying candidate genes in Alzheimer's disease. PNAS ('04); Iossifov et al. Probabilistic inference of molecular networks from noisy data sources.**
**Bioinformatics ('04)**

[Rzhetsky et al, Cell ('08), PLOS CB ('09); Bourne et al. PLOS CB '08]

# Vision for Mining Large-scale Structured Literature
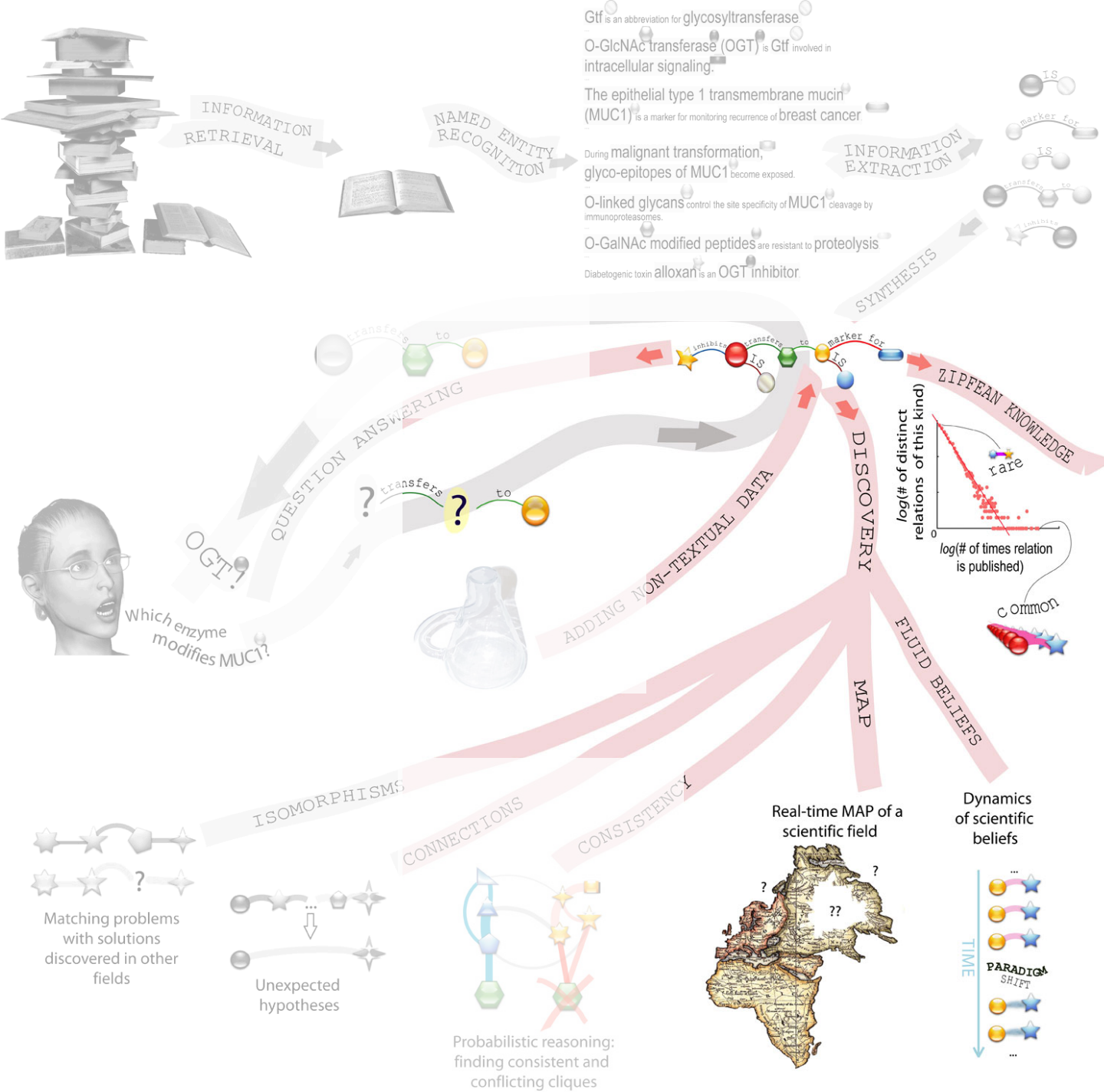
**Making it understand-able (through "mashup")**

**SciVee, podcasts**

[Rzhetsky et al, Cell ('08), PLOS CB ('09); Bourne et al. PLOS CB '08]

# Vision for Mining Large-scale Structured Literature

**Mapping Science + Studying its Dynamics & Evolution**

Gtf is an abbreviation for glycosyltransferase

O-GlcNAc transferase (OGT) is Gtf involved in intracellular signaling.

The epithelial type 1 transmembrane mucin (MUC1) is a marker for monitoring recurrence of breast cancer.

During malignant transformation, glyco-epitopes of MUC1 become exposed.

O-linked glycans control the site specificity of MUC1 cleavage by immunoproteasomes.

O-GalNAc modified peptides are resistant to proteolysis

Diabetogenic toxin alloxan is an OGT inhibitor

INFORMATION RETRIEVAL

NAMED ENTITY RECOGNITION

INFORMATION EXTRACTION

SYNTHESIS

QUESTION ANSWERING

OGT?

Which enzyme modifies MUC1?

ADDING NON-TEXTUAL DATA

DISCOVERY

ZIPFEAN KNOWLEDGE

$log$(# of distinct relations of this kind)

$log$(# of times relation is published)

rare

common

FLUID BELIEFS

MAP

ISOMORPHISMS

Matching problems with solutions discovered in other fields

Unexpected hypotheses

CONNECTIONS

CONSISTENCY

Probabilistic reasoning: finding consistent and conflicting cliques

Real-time MAP of a scientific field

Dynamics of scientific beliefs

TIME

PARADIGM SHIFT

[Rzhetsky et al, Cell ('08), PLOS CB ('09); Bourne et al. PLOS CB '08]

# Vision for Mining Large-scale Structured Literature

INFORMATION RETRIEVAL

NAMED ENTITY RECOGNITION

Gtf is an abbreviation for glycosyltransferase

O-GlcNAc transferase (OGT) is Gtf involved in intracellular signaling.

The epithelial type 1 transmembrane mucin (MUC1) is a marker for monitoring recurrence of breast cancer.

During malignant transformation, glyco-epitopes of MUC1 become exposed.

O-linked glycans control the site specificity of MUC1 cleavage by immunoproteasomes.

O-GalNAc modified peptides are resistant to proteolysis

Diabetogenic toxin alloxan is an OGT inhibitor

INFORMATION EXTRACTION

SYNTHESIS

QUESTION ANSWERING

OGT?

Which enzyme modifies MUC1?

ADDING NON-TEXTUAL DATA

DISCOVERY

MAP

FLUID BELIEFS

ZIPFEAN KNOWLEDGE

log(# of distinct relations of this kind)

rare

log(# of times relation is published)

common

ISOMORPHISMS

CONNECTIONS

CONSISTENCY

Matching problems with solutions discovered in other fields

Unexpected hypotheses

Probabilistic reasoning: finding consistent and conflicting cliques

Real-time MAP of a scientific field

Dynamics of scientific beliefs

TIME

PARADIGM SHIFT

• Revealing patterns of collaboration
• Understanding basis of terms & nomenclature
• Tracking the evolution of ideas
• Models for the evolution of science;
• Helping set policy & research directions

[Rzhetsky et al, Cell ('08), PLOS CB ('09); Bourne et al. PLOS CB '08]

# Analyzing the Structure of Genomic Science: Mapping the diffusion of ideas & data across disciplines

- Intro: Using the Data Exhaust from Consortia
- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers
- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases
- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
  - Moderated by parm. w/ 2 time regimes

# Analyzing the Structure of Genomic Science: Mapping the diffusion of ideas & data across disciplines

- Intro: Using the Data Exhaust from Consortia
- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers
- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases
- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
  - Moderated by parm. w/ 2 time regimes

# Increase in Consortium Science

9 -

# Examples Illuminating Current State of Affairs:
# Using Network Representations to Make Maps of Science -- Studying the Publication Patterns of Genomics Consortia



[Douglas et al. GenomeBiol. ('05), pubnet.gersteinlab.org]

# Different Representations of the Publication Network of a Structural Genomics Center (NESG)



a) Paper / Co-Authorship

b) Author / Co-Authorship

c) Paper / Shared MeSH term

d) Paper / Shared Location

BSGC

CESG

JCSG (45)

MCSG

NESG (19)

NYSGXRC

SECSG

SGPP

TB (7)

**Co-authorship Networks comparing the 9 NIH Structural Genomics Centers**

**Average Degree**

[Douglas et al. GenomeBiol. ('05), pubnet.gersteinlab.org]

- Intro: Using the Data Exhaust from Consortia

- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers

- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases

- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
  - Moderated by parm. w/ 2 time regimes

Timeline of genomics milestones:

**Above the timeline:**
- HumanGenome Project
- ENCODE Pilot
- ENCODE Production
- ComparativeENCODE
- Epigenome Roadmap

**Timeline markers:** 2000, 2005, 2010, 2015

**Below the timeline:**
- Worm Genome
- modENCODE
- 1000 Genomes Pilot
- 1000 Genomes Phase 3
- GTEx

- A THOUSAND GENOMES
- END OF THE BEGINNING
- Mapping our differences

**# Authors**       **Yr. ('04 to '15)**

■ non–ENCODE (papers used ENCODE data)     ■ ENCODE

With help of NHGRI, identified:
1,786 ENCODE members & 8,263 non-members
from 558 consortium papers supported by ENCODE funding &
702 community papers that used ENCODE data but were not supported by
ENCODE funding

# ENCODE co-authorship network



2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014

Legend:
- ● ENCODE member
- ● non-member
- ● ENCODE member broker
- ● non-member broker
- — co-authorship

# Network statistics highlight change in modularity with consortium rollouts (L) & importance of broker role (R)



[Wang et al., TIG ('16)]

# Similar Findings in terms of modularity & broker scientists in the modENCODE consortium as for ENCODE



2014　2013　2012　2011　2010　2009　2008　2007

**modENCODE**

Number of member neighbors
Number of non-member neighbors

Modularity
Number of co-authorship cluster
Year

- consortium member
- non–member member
- broker non–member
- broker consortium network
- non–consortium network
- random network
- co–authorship

[Wang et al., TIG ('16)]

- Intro: Using the Data Exhaust from Consortia

- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers

- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases

- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
  - Moderated by parm. w/ 2 time regimes

# Diffusion of a data type (sequenced bases) measured by occurrence in specialty journals



Cumulative bases deposited in the Sequence Read Archive by journal

Legend:
- Nature
- Science
- Cell
- Nucleic Acids Res.
- Genome Biology
- Nat. Biotech.
- ISME J.
- Proc Natl Acad Sci USA
- Nat Chem Biol.
- Molecular ecology

Number of bases

Date

[Muir et al. GenomeBiol. '16]

# RNAi: Birth of a Field in the Literature Culmin-ating in the 2006 Nobel



1998

1999

2000

2001

2002

2003

● Andrew Fire    ● Craig Mello

- Intro: Using the Data Exhaust from Consortia

- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers

- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases

- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
  - Moderated by parm. w/ 2 time regimes

# Spread of information as a diffusion process

The knowledge of a scientific publication is a piece of information
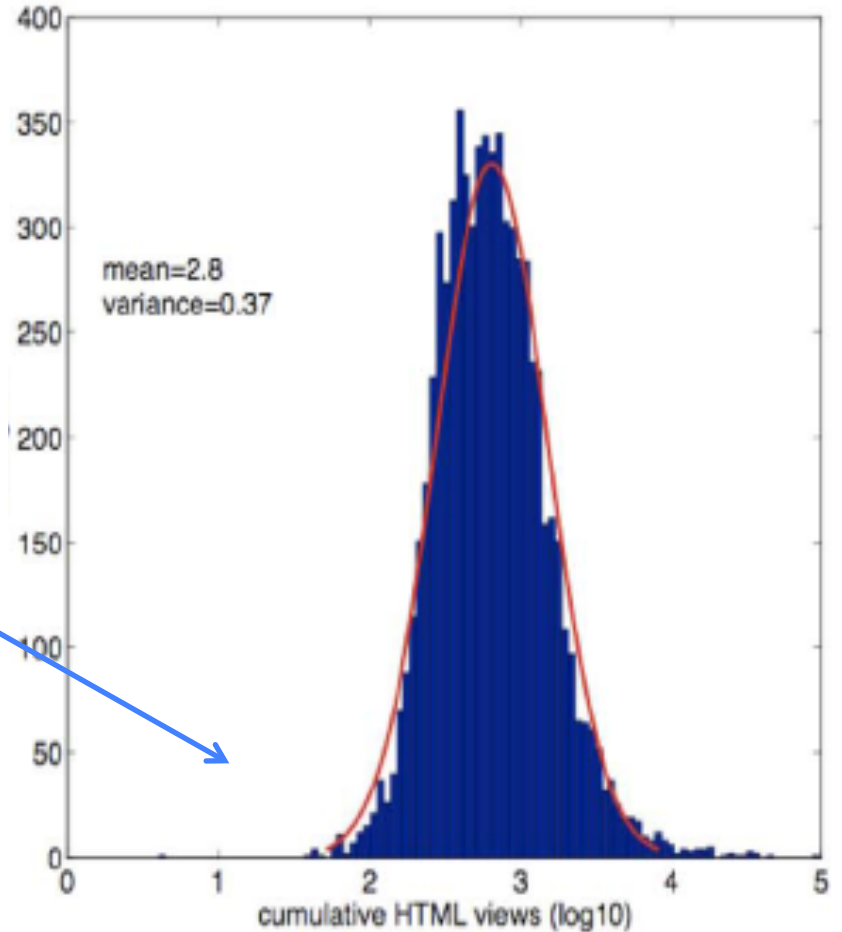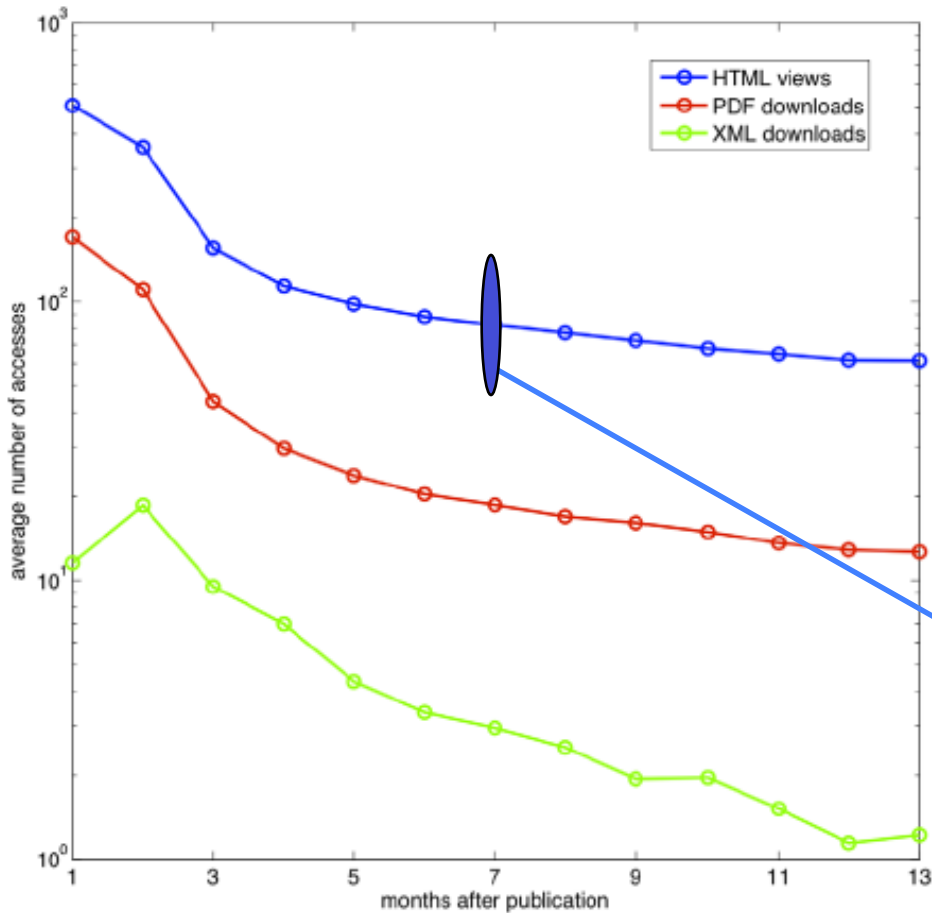


Based on PLOS ALM data for ~7000 papers

the access of different articles follows a log-normal distribution

Yan and Gerstein PLoS One 2011.
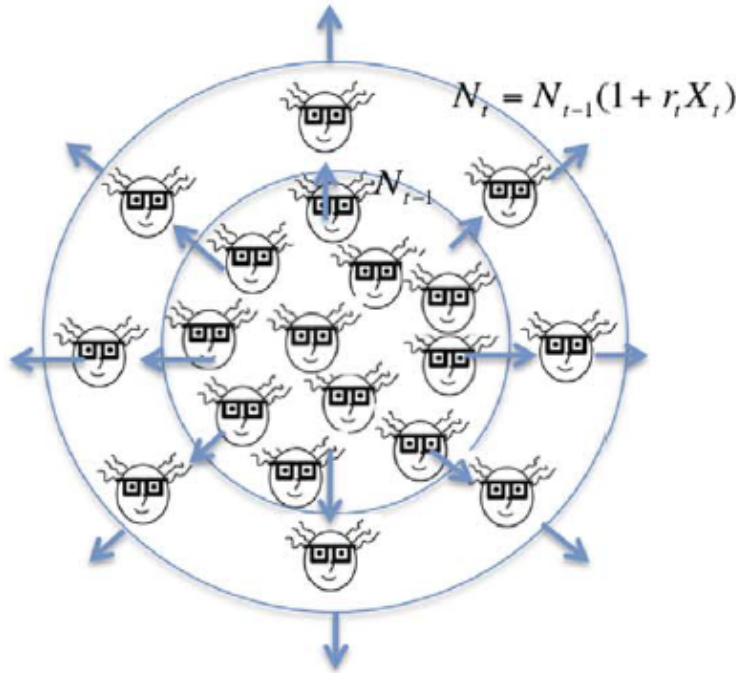
# Spread of information as a diffusion process

The knowledge of a scientific publication is a piece of information



the access of different articles follows a log-normal distribution

Yan and Gerstein PLoS One 2011.

## Modeling Information diffusion

log-normal distribution suggests a simple model,
random multiplicative process for a given paper p:

$$N_t(p) = N_{t-1}(p) \left( 1 + r_t X_t(p) \right)$$

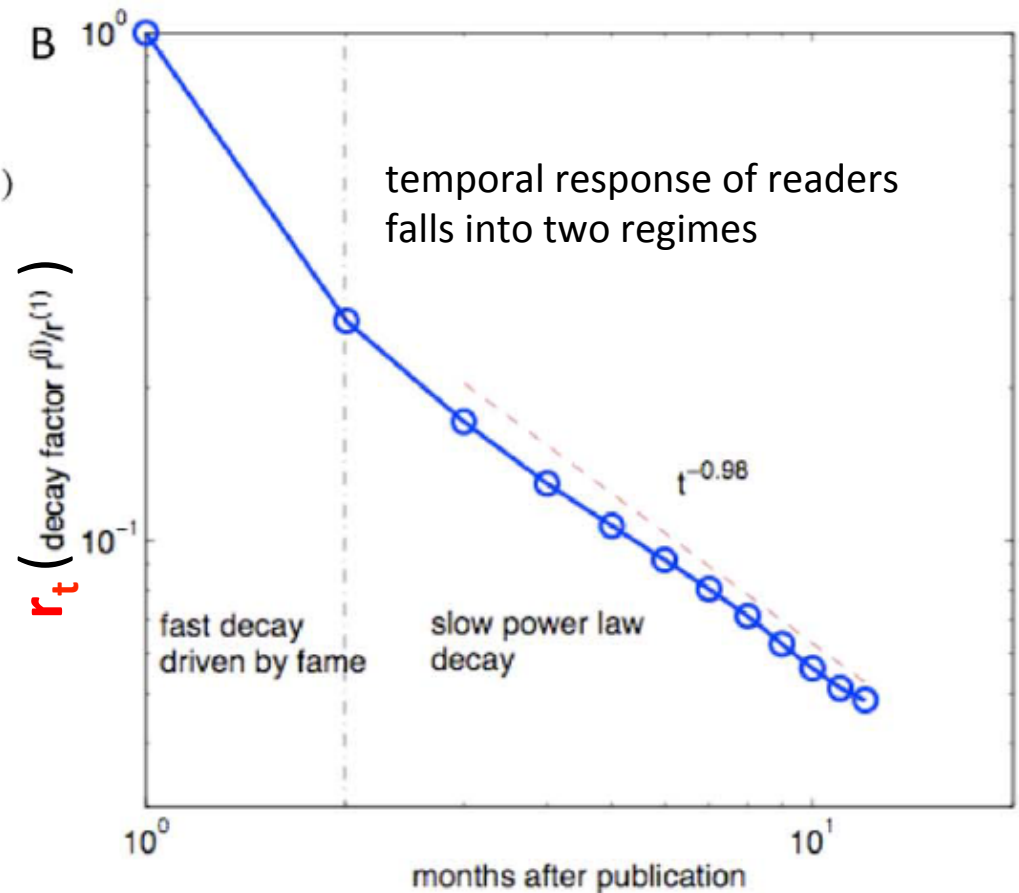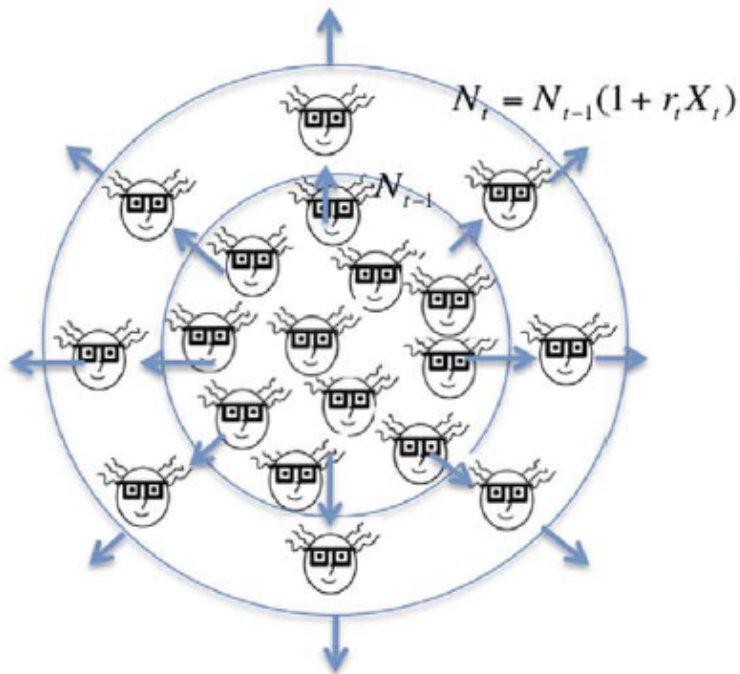share the information with friends?



$$N_t = N_{t-1}(1 + r_t X_t)$$

$N_t$ = cumulative number of accesses for a given paper up to time t

$X_t$ = iid random variable whose mean represents avg. fraction of scientists willing to "spread" the paper at t

$r_t$ = moderating parameter on how the "spreading" changes over time

Yan and Gerstein PLoS One 2011.

# Modeling Information diffusion

## Refine values of $r_t$ for different times



$$N_t = N_{t-1}(1 + r_t X_t)$$

temporal response of readers falls into two regimes

$t^{-0.98}$

fast decay driven by fame

slow power law decay

$r_t$ (decay factor $r^{(t)}/r^{(1)}$)

months after publication

Yan and Gerstein PLoS One 2011.

- Intro: Using the Data Exhaust from Consortia

- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers

- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases

- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
  - Moderated by parm. w/ 2 time regimes

- Intro: Using the Data Exhaust from Consortia

- PubNet Tool for analyzing co-publication patterns
  - Different patterns for Str. Genomics centers

- Analysis of Evolution of the ENCODE Consortium
  - Differences in "modularity" for members & users
  - Key role for brokers

- Other Examples of Idea & Data Diffusion
  - RNAi & SRA bases

- Quant. Model of Information Diffusion
  - Based on PLOS ALM
  - Random multiplicative process
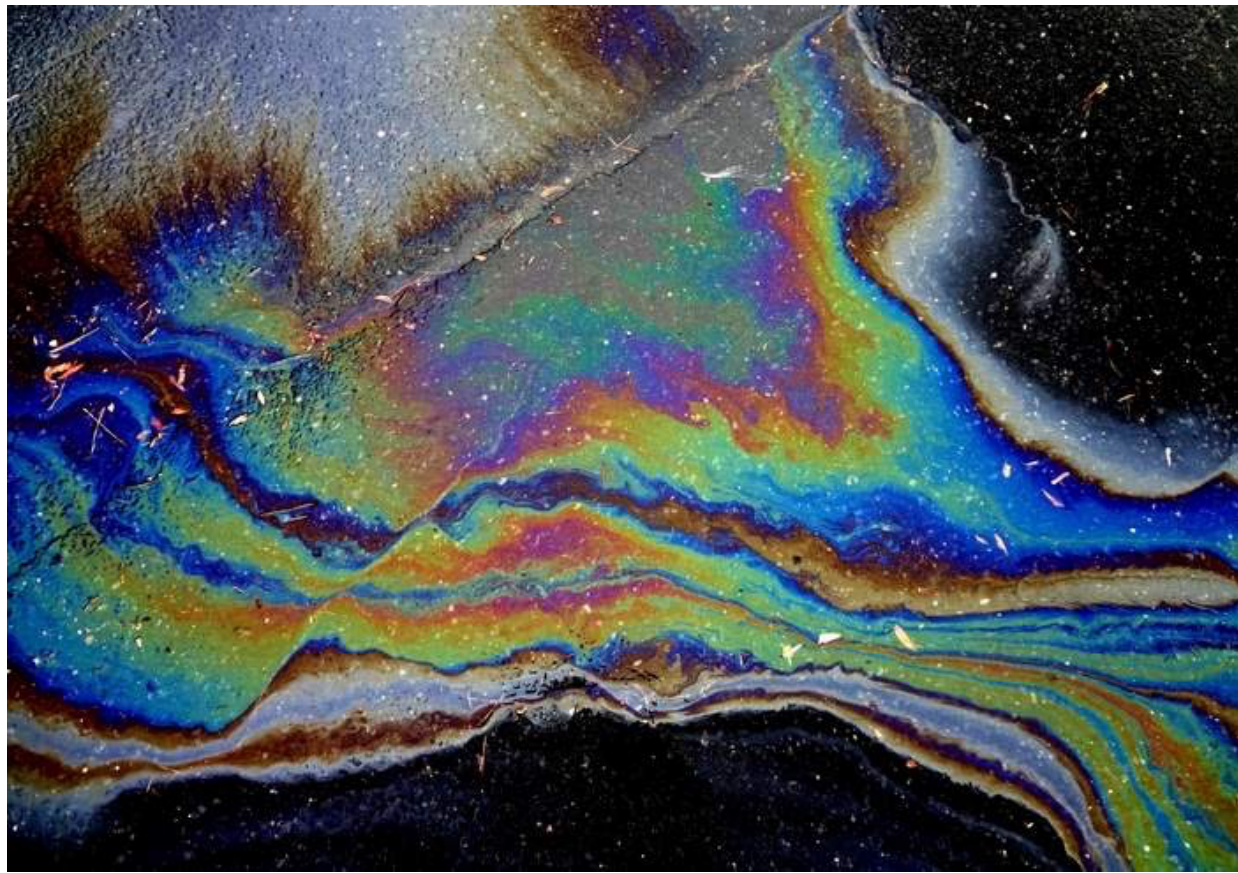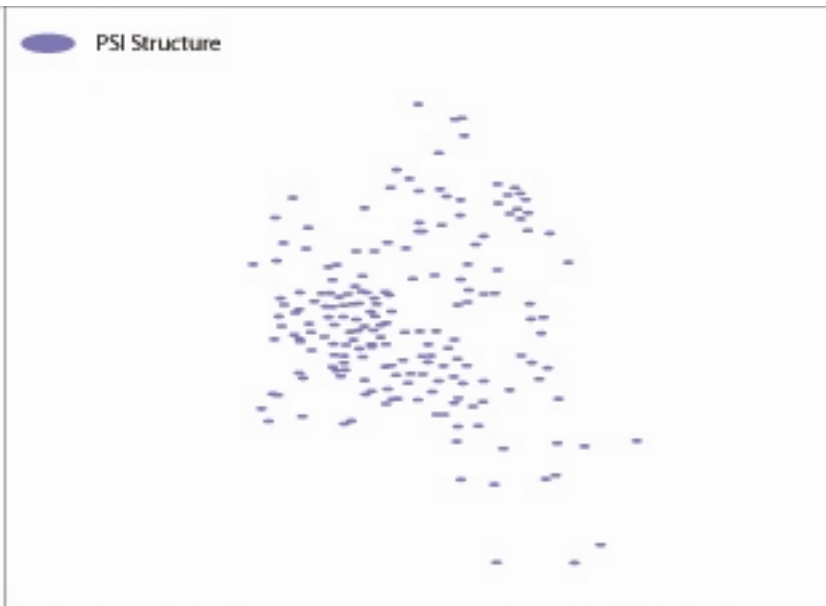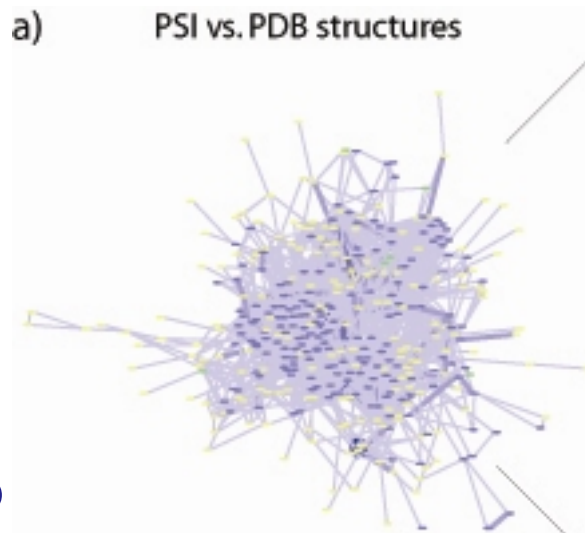  - Moderated by parm. w/ 2 time regimes

# Acknowledgments

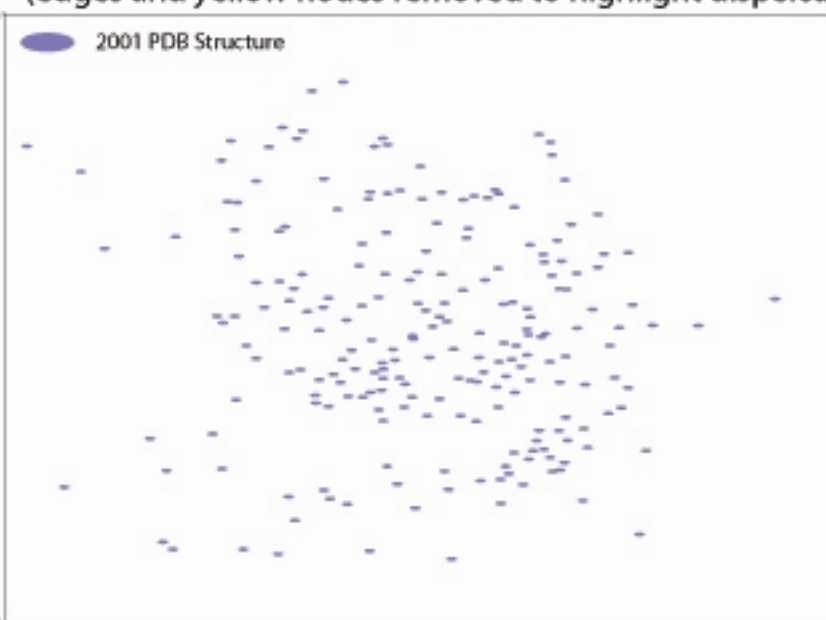**Hiring Postdocs. See gersteinlab.org/jobs**

# Extra

**Clustering structures determined by struc. genomics consortia according to functional similarity: Is there a functional bias in consortia structures?**



a) PSI vs. PDB structures — PSI Structure

b) PDB vs PDB structures (control) — 2001 PDB Structure

(edges and yellow nodes removed to highlight dispersal)

|  | Avg. Degree | Avg. Path | Clust. Coeff. | Diameter |
|---|---|---|---|---|
| PSI | 24 | 2.6 | 37% | 7 |
| PDB | 6 | 3.9 | 31% | 9 |

[Douglas et al. GenomeBiol. ('05), pubnet.gersteinlab.org]

31 – Lect

# Over-representation of crystallography among the Nobel Prizes, highlighted by the 2006 Nobels

| | MeSH term | Crystallography | Protein Conformation | Chemistry |
|---|---|---|---|---|
| **1970-2006** | Related Nobel Prizes | 7*** | 9 | 36 |
| | Fraction of All PubMed records | 0.3% | 1.1% | 9.3% |
| | Fraction of All Chemistry records | **4%** | **12%** | 100% |
| | Fraction of Available Nobel | **19%** | **25%** | 100% |
| **1996-2006** | Related Nobel Prizes | 4**** | 5 | 10 |
| | Fraction of All PubMed records | 0.6% | 2.1% | 9.0% |
| | Fraction of All Chemistry records | **7%** | **23%** | 100% |
| | Fraction of Available Nobel | **40%** | **50%** | 100% |

[Seringhaus & Gerstein, Science (2007)]

# Info about content in this slide pack

- General PERMISSIONS
  - This Presentation is copyright Mark Gerstein, Yale University, 2015.
  - Please read permissions statement at www.**gersteinlab.org/misc/permissions.html** .
  - Feel free to use slides & images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).
  - Paper references in the talk were mostly from Papers.GersteinLab.org.

- PHOTOS & IMAGES. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see http://streams.gerstein.info .
  - In particular, many of the images have particular EXIF tags, such as  kwpotppt , that can be easily queried from flickr, viz: http://www.flickr.com/photos/mbgmbg/tags/kwpotppt