



DESIGNING DISCOVERY

Hypergraph Models of Innovation for Science & Technology

KNOWLEDGE
LAB

James Evans
Feng (Bill) Shi

KNOWLEDGE

LAB

Big Data, Machine Learning and Intelligent Crowdsourcing enables us to:

1. Trace

2. Understand

3. Discover

4. Improve...the scientific and scholarly process

*computationally enhanced
Science of Science*

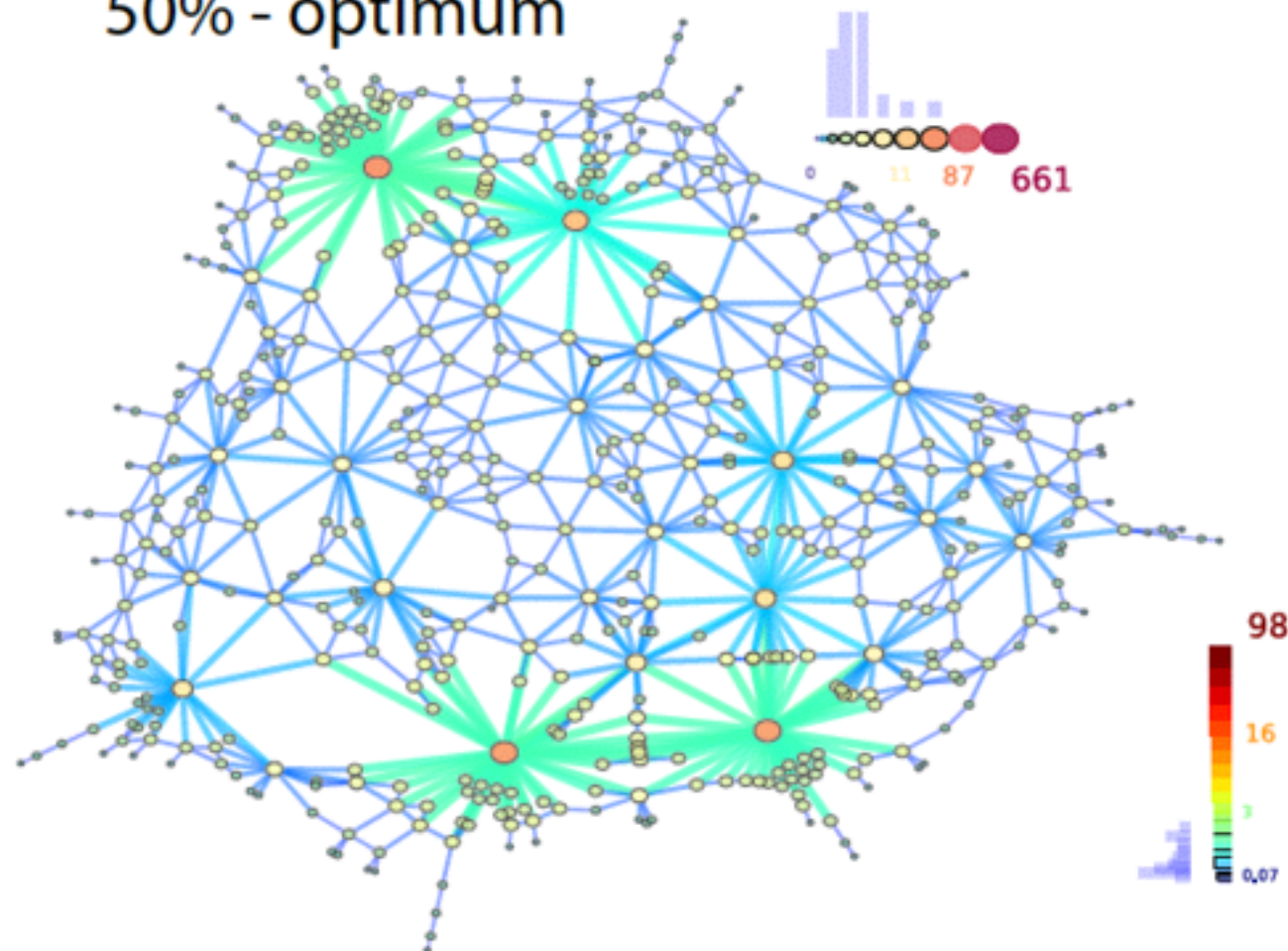
KNOWLEDGE LAB



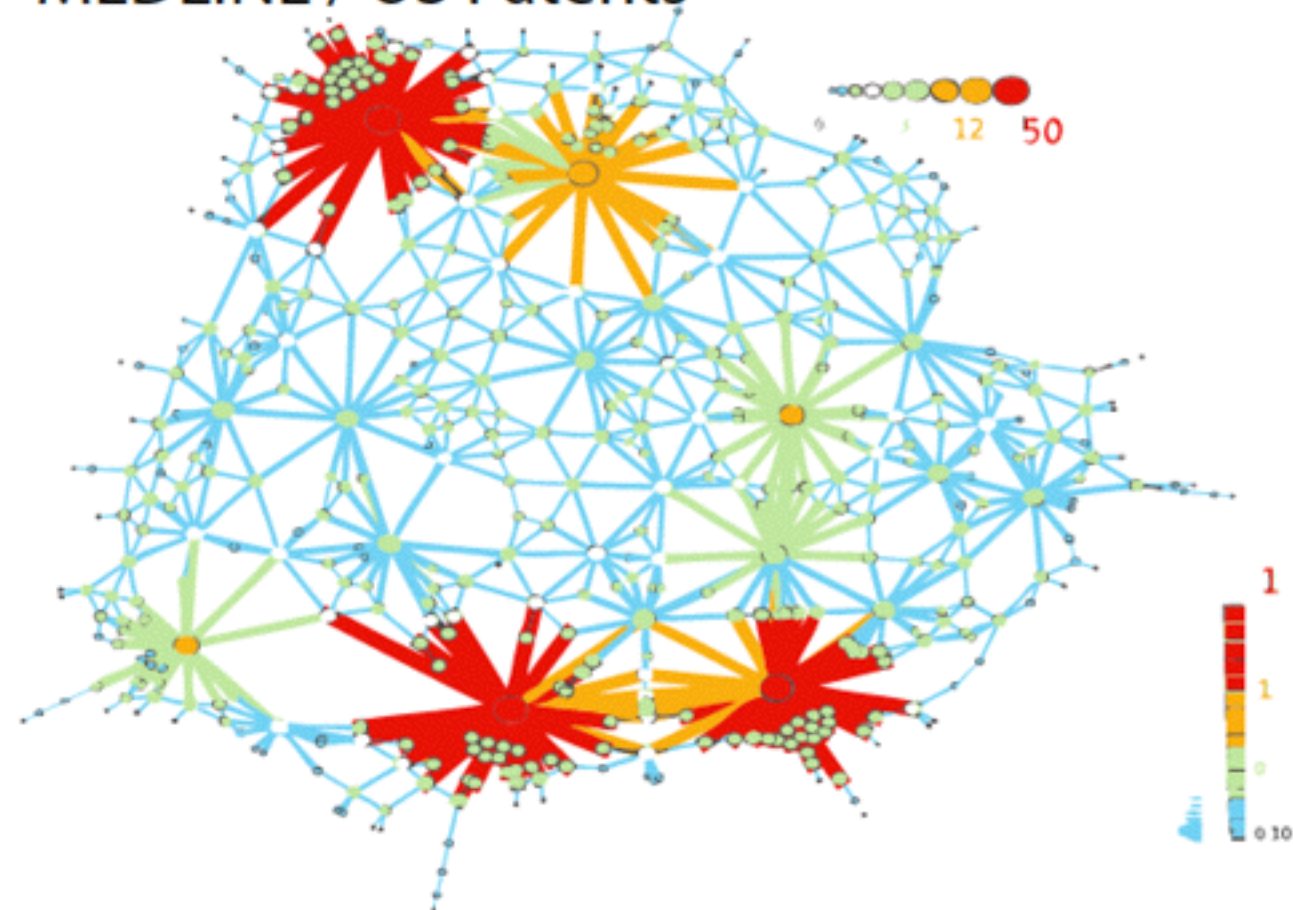
How Science Thinks

jevans@uchicago.edu

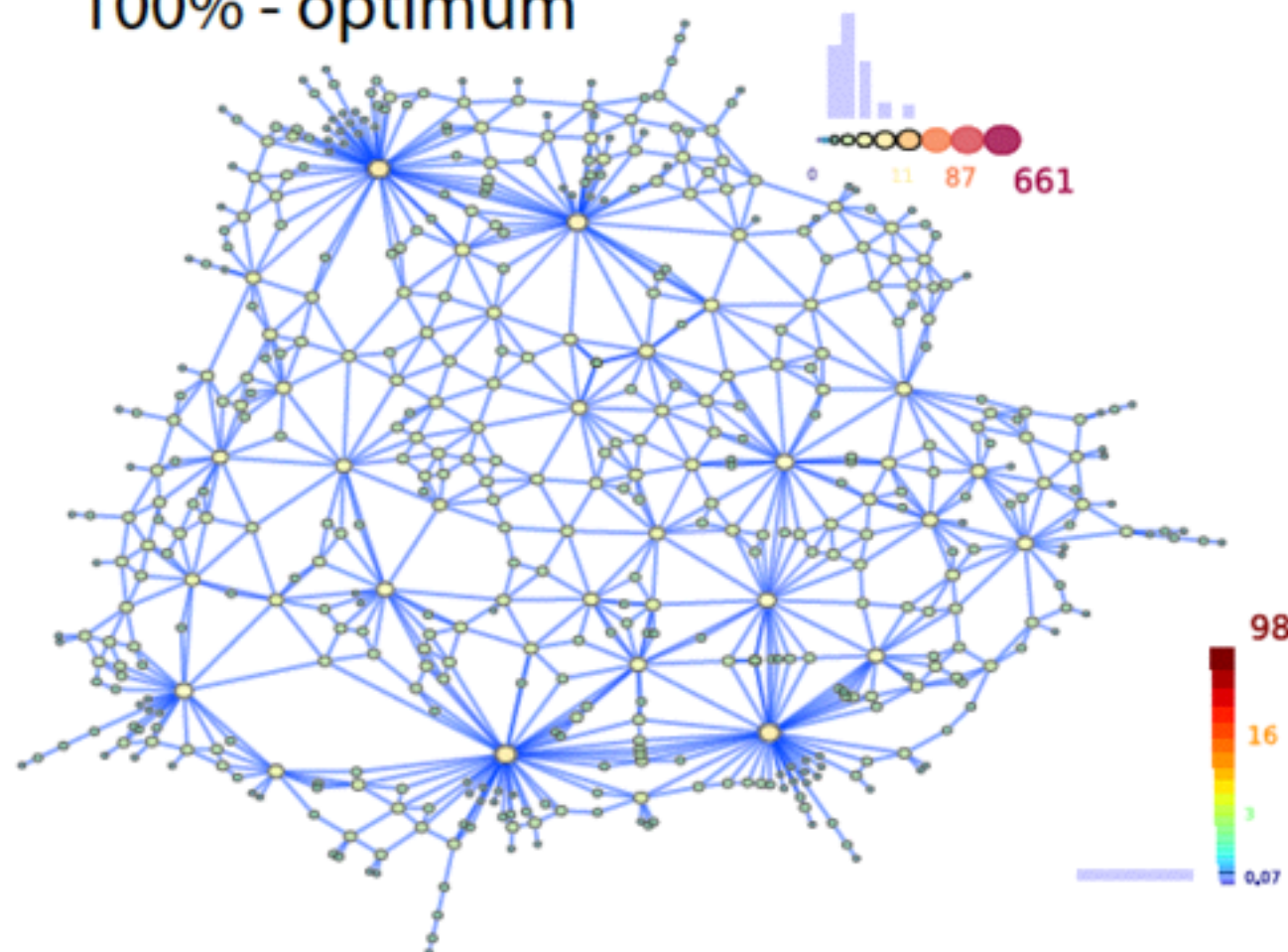
50% - optimum



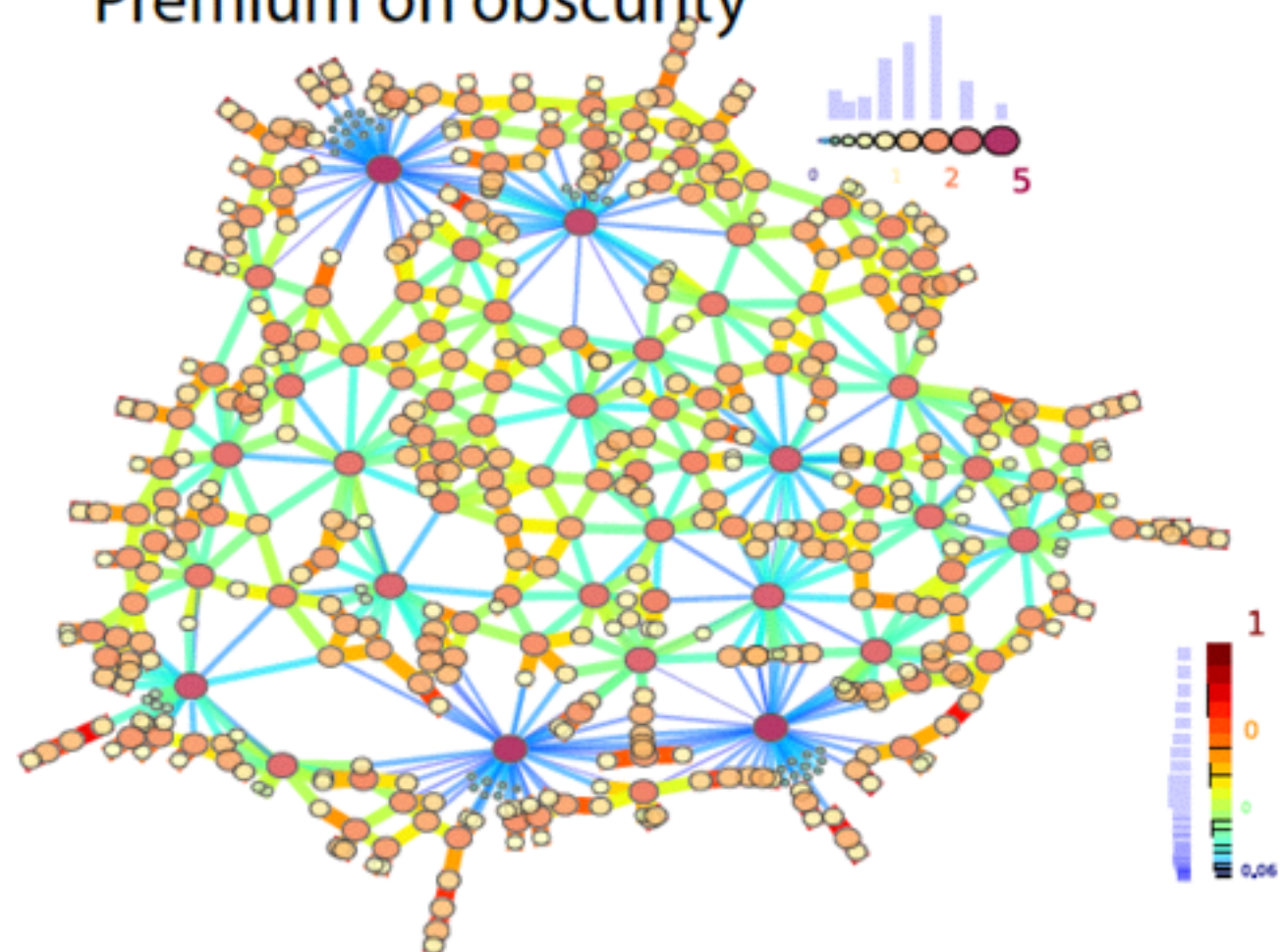
MEDLINE / US Patents

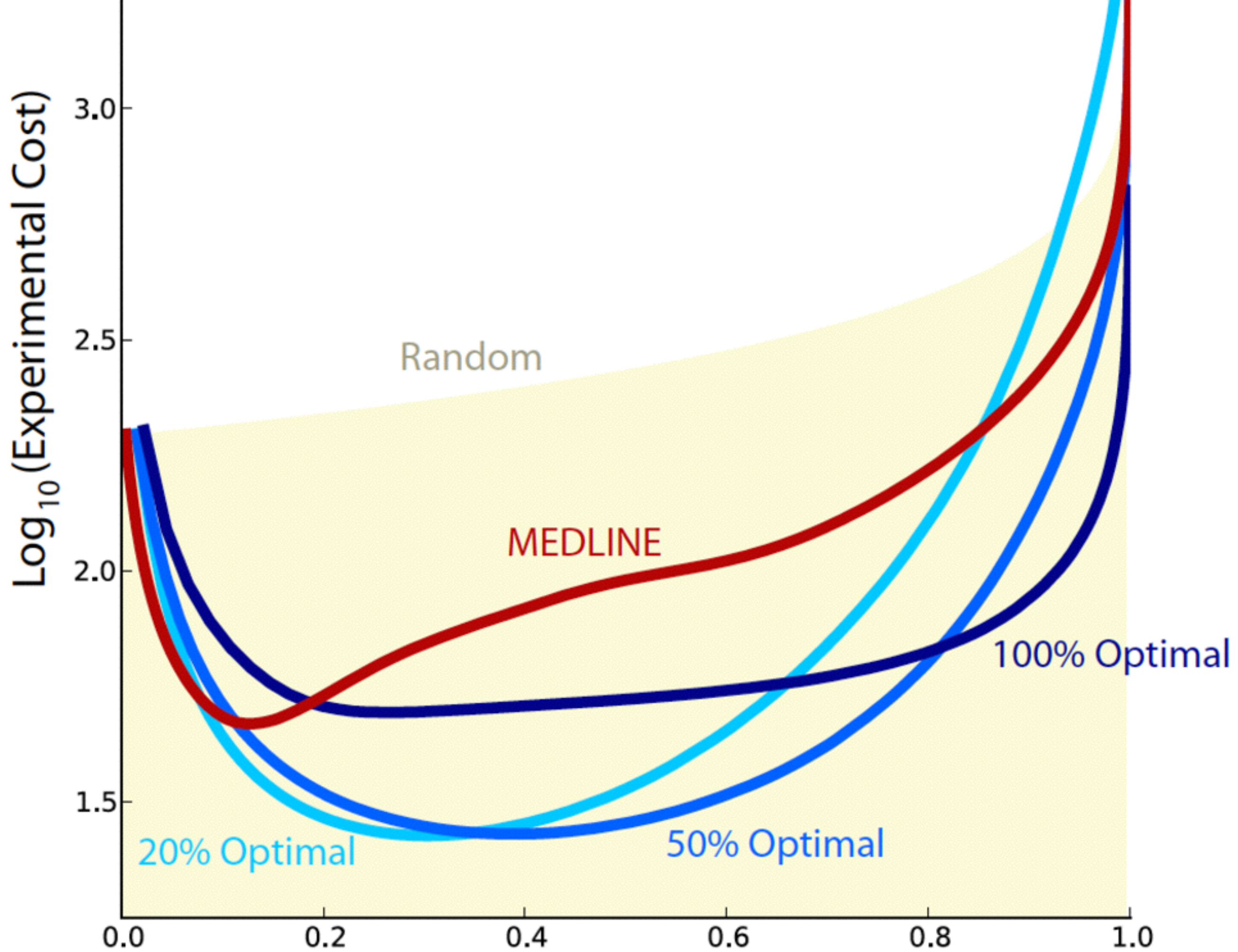


100% - optimum

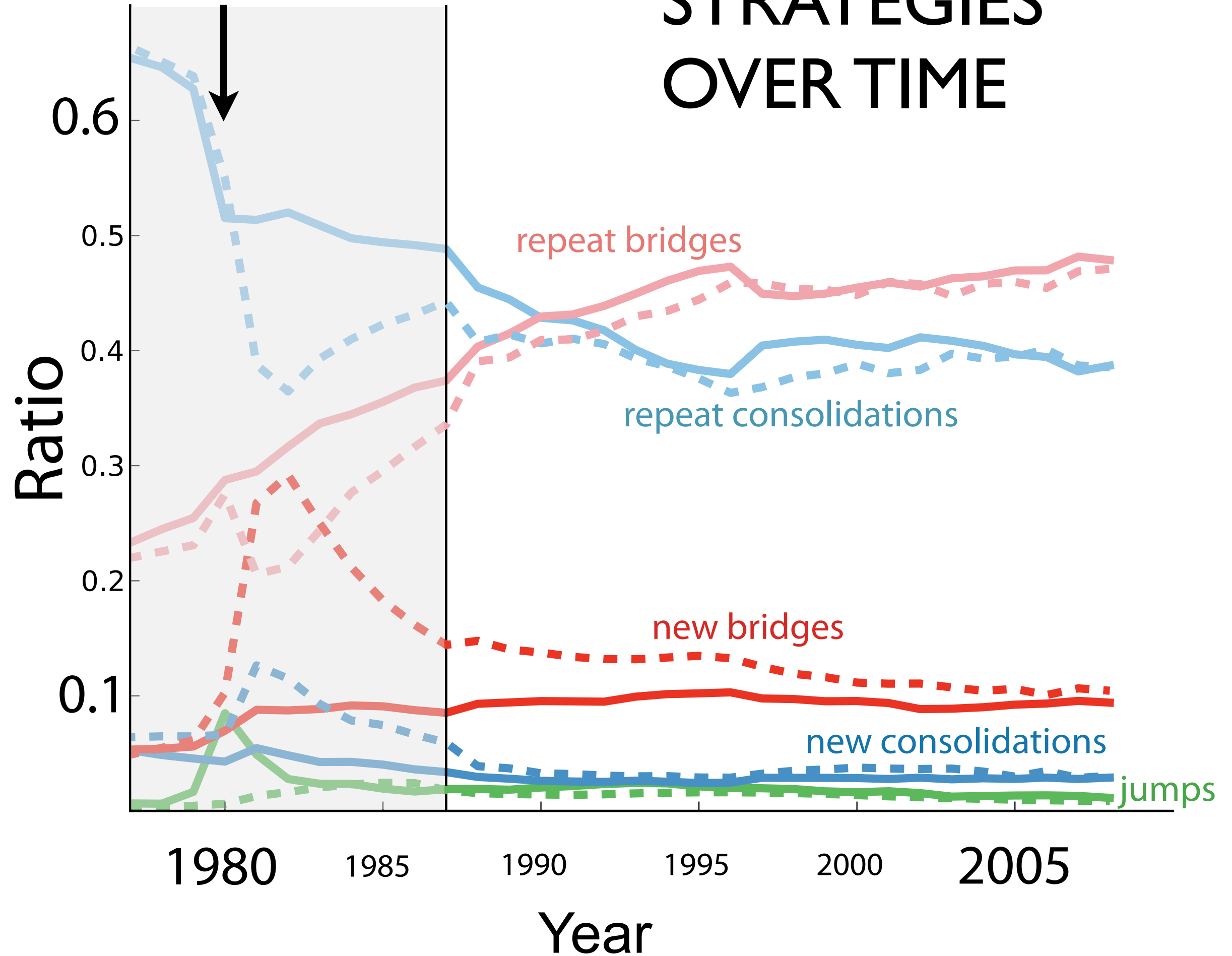


Premium on obscurity

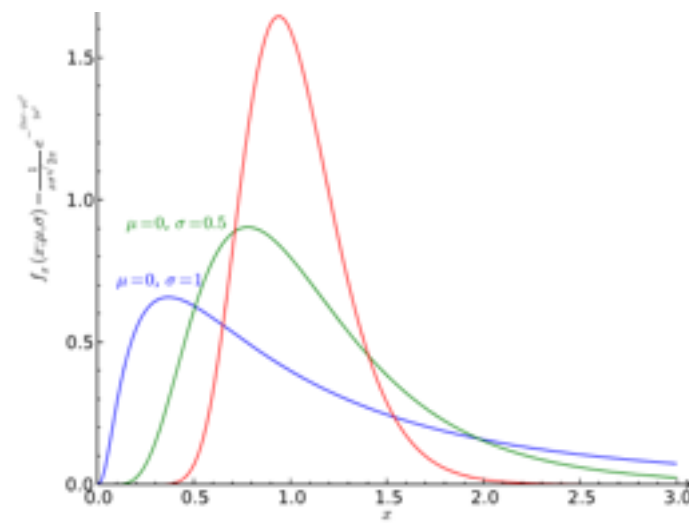
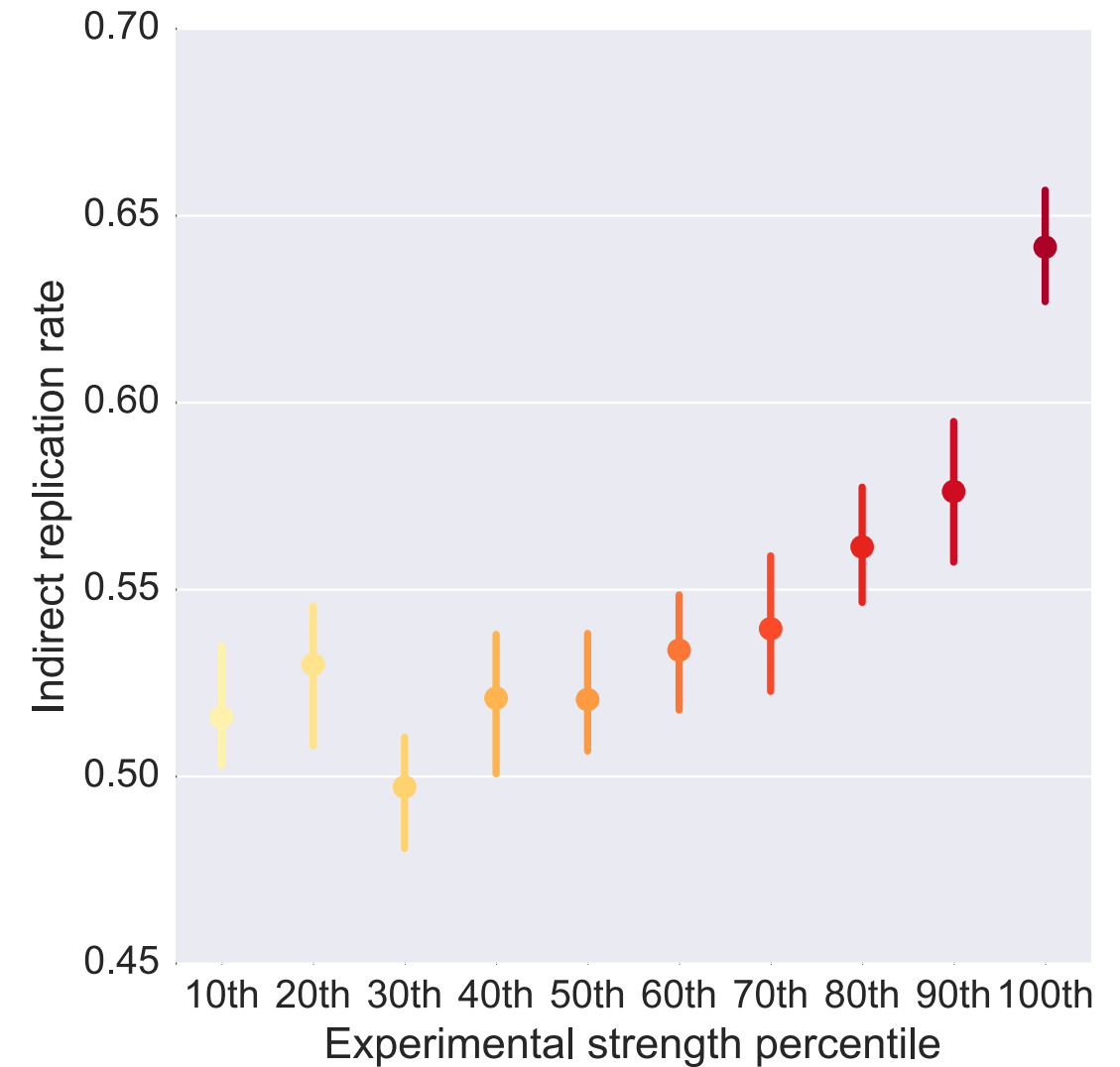
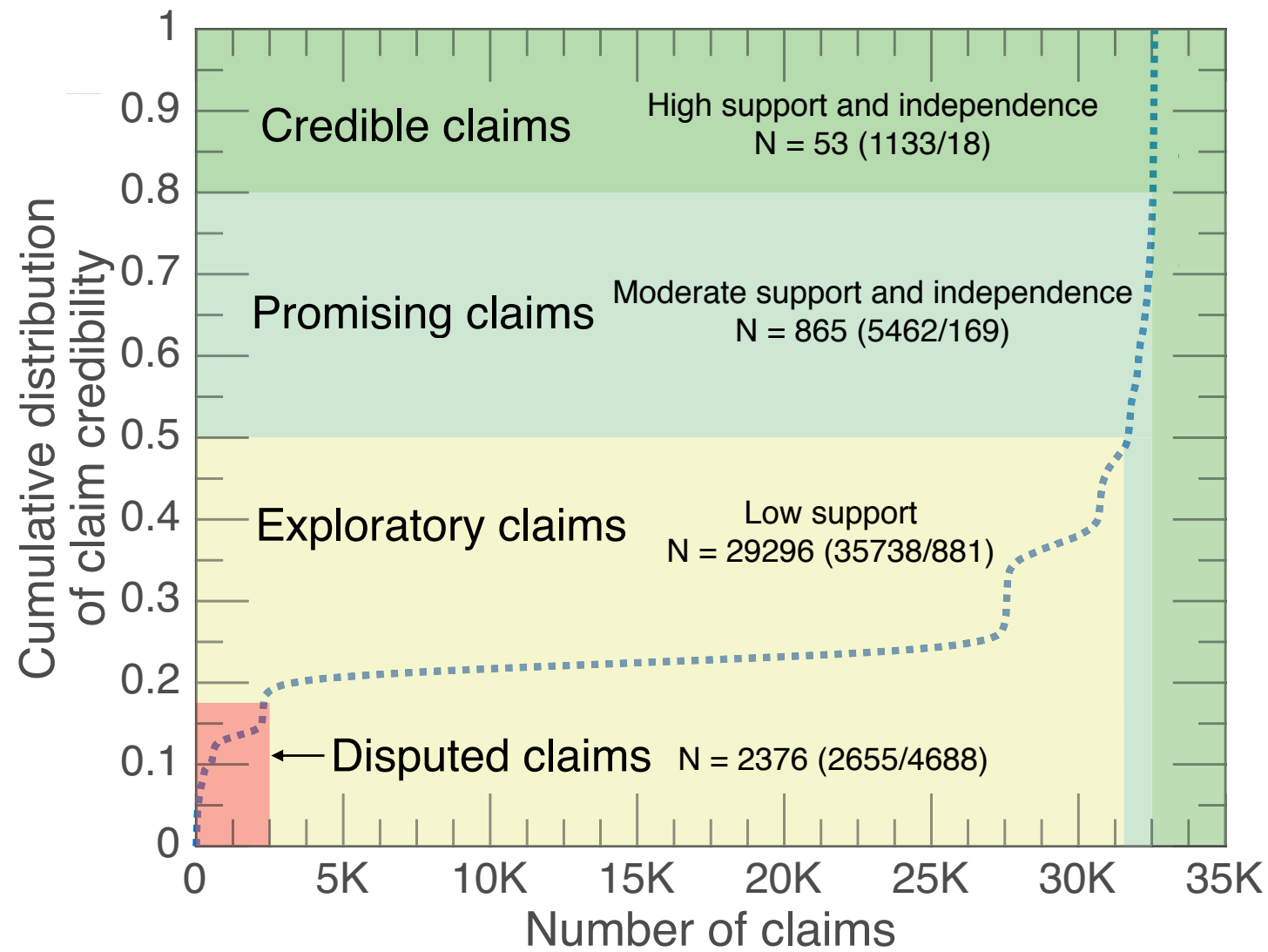




STRATEGIES OVER TIME



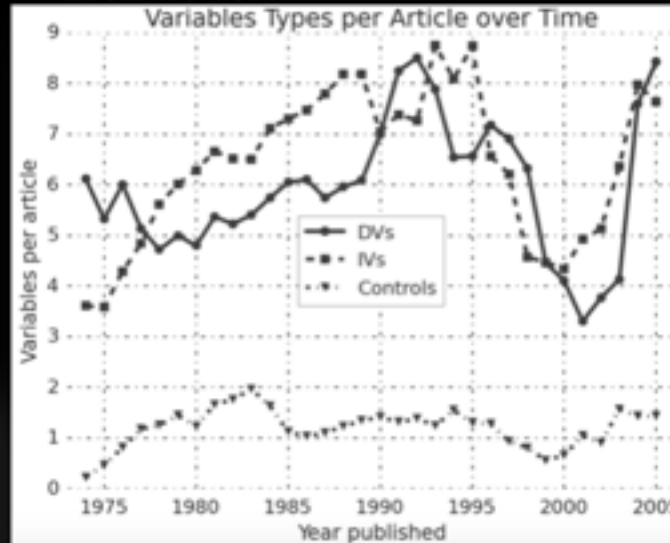
High-throughput quasi-replication



...and in the social sciences

Estimating original models

- Articles typically have
 - A few dependent variables
 - A few independent variables



Approximate Estimation of Original Models

$$\begin{aligned} Y_{1,t} &= X_t \beta + \epsilon \\ Y_{2,t} &= X_t \beta + \epsilon \\ &\vdots \\ Y_{f,t} &= X_t \beta + \epsilon \end{aligned}$$

$$X'_t = \begin{pmatrix} 1 & \dots & 1 \\ x_{1,1,t}^* & \dots & x_{1,z,t}^* \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{pmatrix}$$

arzhetsk@medicine.bsd.uchicago.edu

"Race, Sex and Feminist Outlooks" (Ransford and Miller 1983)

$$\begin{aligned} \text{FEHOME}_{1974} &= X_t \beta + \epsilon \\ \text{FEWORK}_{1974} &= X_t \beta + \epsilon \\ \text{FEPRES}_{1974} &= X_t \beta + \epsilon \\ \text{FEPOL}_{1974} &= X_t \beta + \epsilon \end{aligned} \quad X'_t = \begin{pmatrix} 1 & \dots \\ \text{MAWORK}_{1,t}^* & \dots \\ \text{OCC}_{1,t}^* & \dots \\ \text{EDUC}_{1,t}^* & \dots \\ \text{GOVAID}_{1,t}^* & \dots \\ \text{FINRELA}_{1,t}^* & \dots \\ \text{INCOM16}_{1,t}^* & \dots \\ \text{CLASS}_{1,t}^* & \dots \end{pmatrix}$$

$t \in \{1974 - 1978\}$

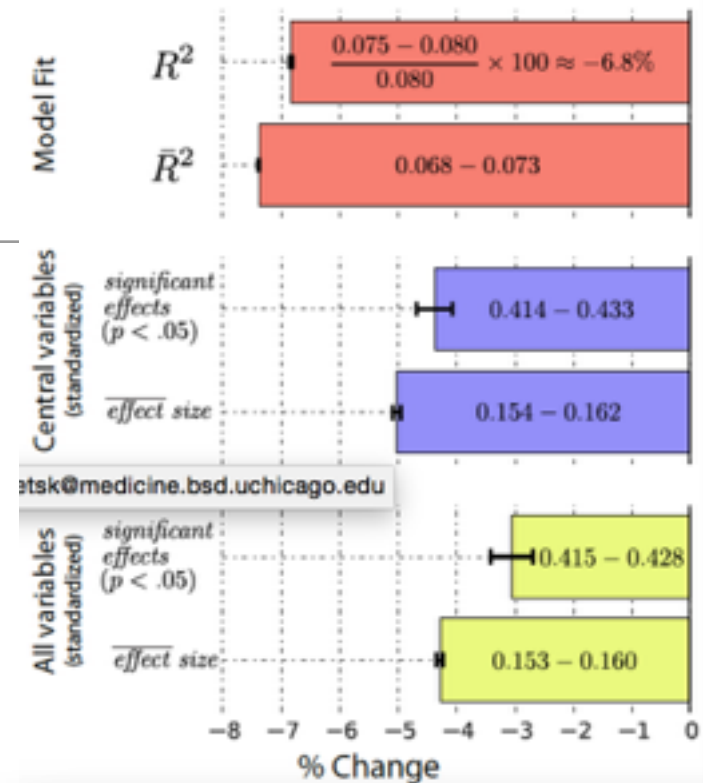
"Confidence in Science: The Gender Gap" (Fox and Firebaugh 1992)

$$\begin{aligned} \text{CONSCI}_t &= X_t \beta + \epsilon \\ \text{CONFINAN}_t &= X_t \beta + \epsilon \\ \text{CONBUS}_t &= X_t \beta + \epsilon \\ \text{CONCLERG}_t &= X_t \beta + \epsilon \\ \text{CONEDUC}_t &= X_t \beta + \epsilon \\ &\vdots \\ \text{CONPRESS}_t &= X_t \beta + \epsilon \end{aligned} \quad X'_t = \begin{pmatrix} 1 & \dots \\ \text{SEX}_{1,t}^* & \dots \\ \text{OCC}_{1,t}^* & \dots \\ \text{EDUC}_{1,t}^* & \dots \\ \text{PRESTIGE}_{1,t}^* & \dots \\ \text{WRKSTAT}_{1,t}^* & \dots \\ \vdots & \ddots \\ \text{RELIG}_{1,t}^* & \dots \end{pmatrix}$$

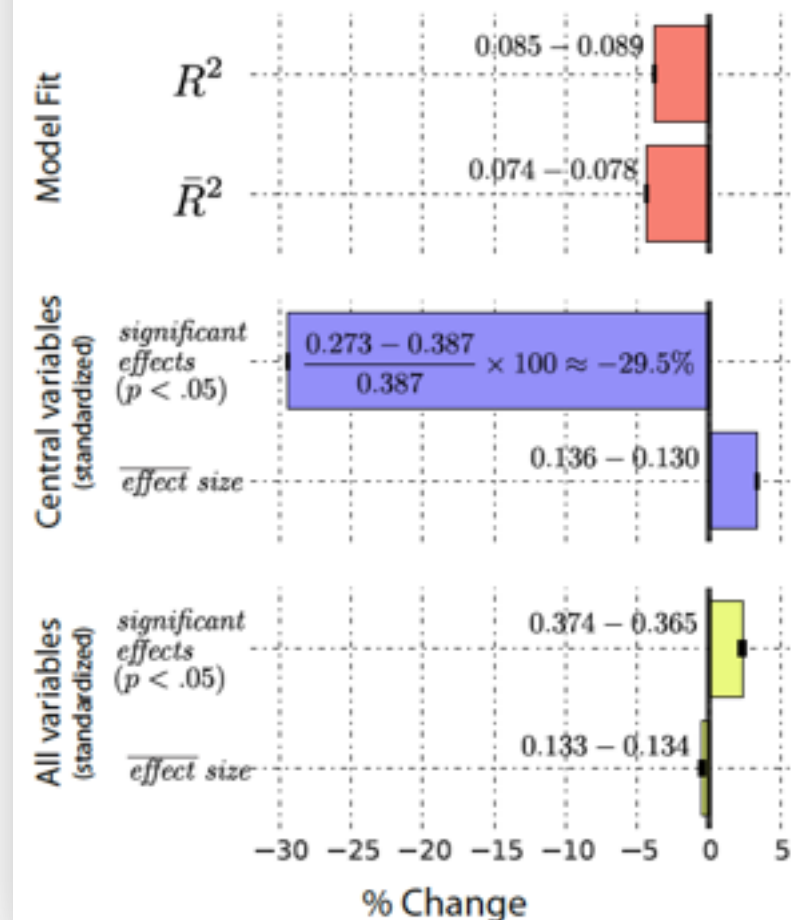
$t \in \{1973 - 1989\}$

40

The Effect of Data Substitution: Next Year Data minus Original Models



The Effect of Substituting One Variable: Perturbed minus Original models



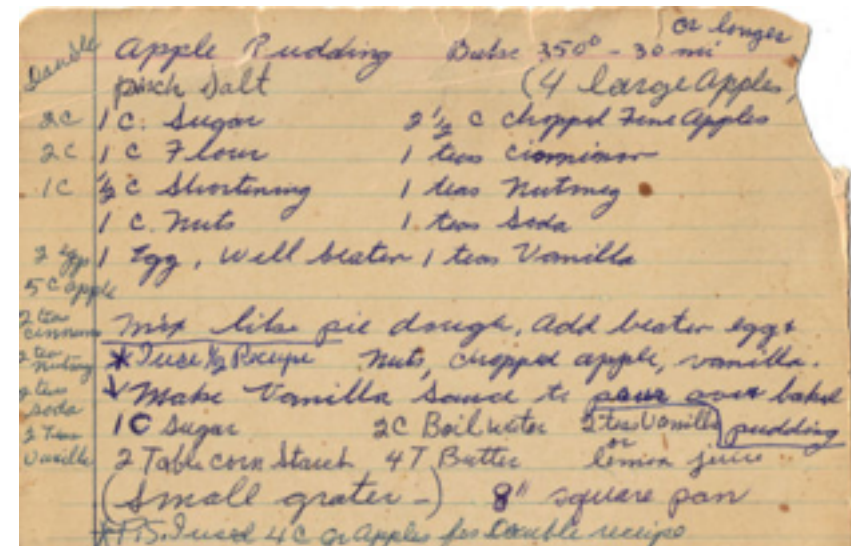
Active Learning for Intelligent Survey Design

Which place looks safer ?



PREDICTING & GENERATING SCIENTIFIC SUCCESS

- Predict **combination of concepts** in future discoveries & inventions



- Predict **level of impact** for future discoveries



KNOWLEDGE REPRESENTATION

Representation

Extraction

Inference

- Collocation Network /
Adjacency Matrix

Inexpensive

Over



$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & & & A_{2n} \\ \vdots & & & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix}$$

KNOWLEDGE REPRESENTATION

Representation

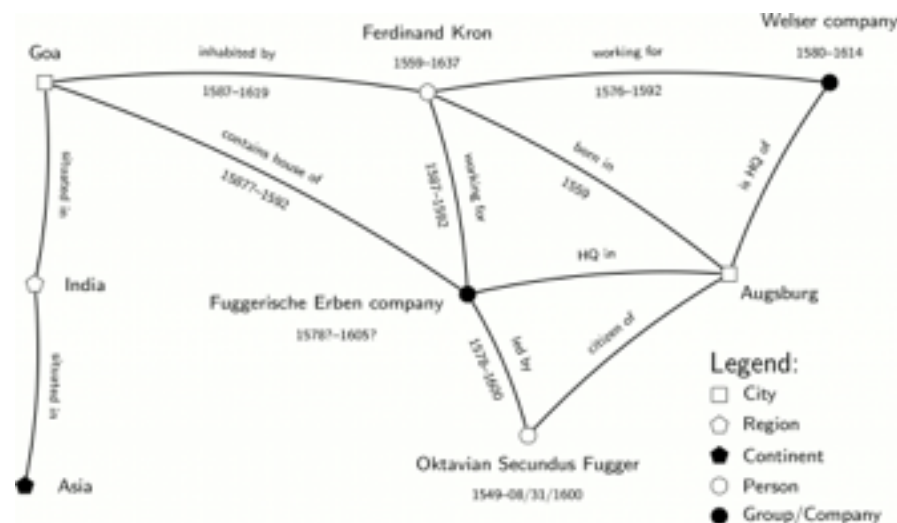
Extraction

Inference

- Collocation Network / Adjacency Matrix

Inexpensive

Over



Expensive

Under

- Semantic Graph or Hypergraph

KNOWLEDGE REPRESENTATION

Representation

Extraction

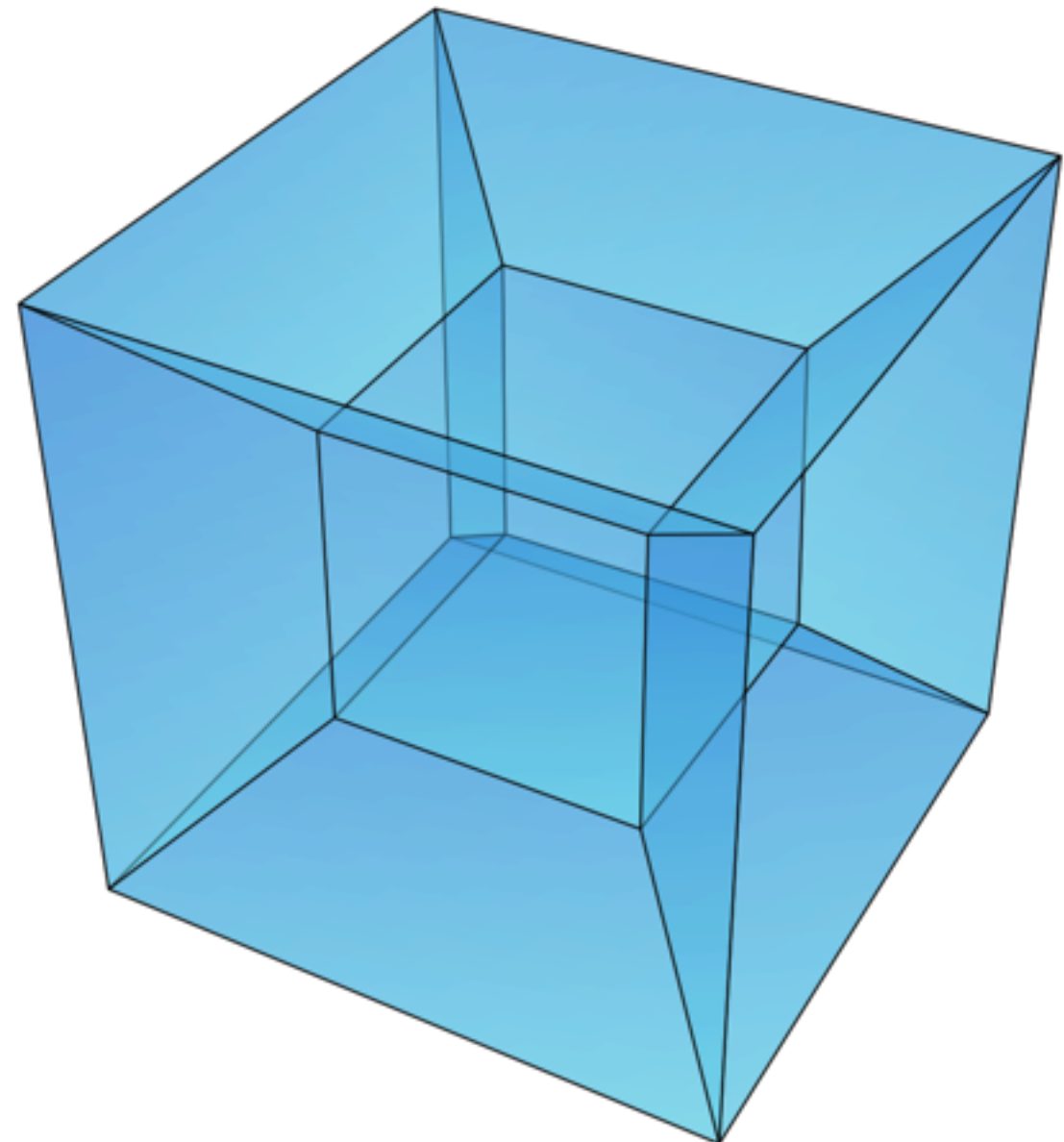
Inference

- Collocation Network /
Adjacency Matrix *Inexpensive* *Over*
- **Collocation Hypergraph /
Adjacency Tensor** *Inexpensive* *Exact*
- Semantic Graph or Hypergraph *Expensive* *Under*

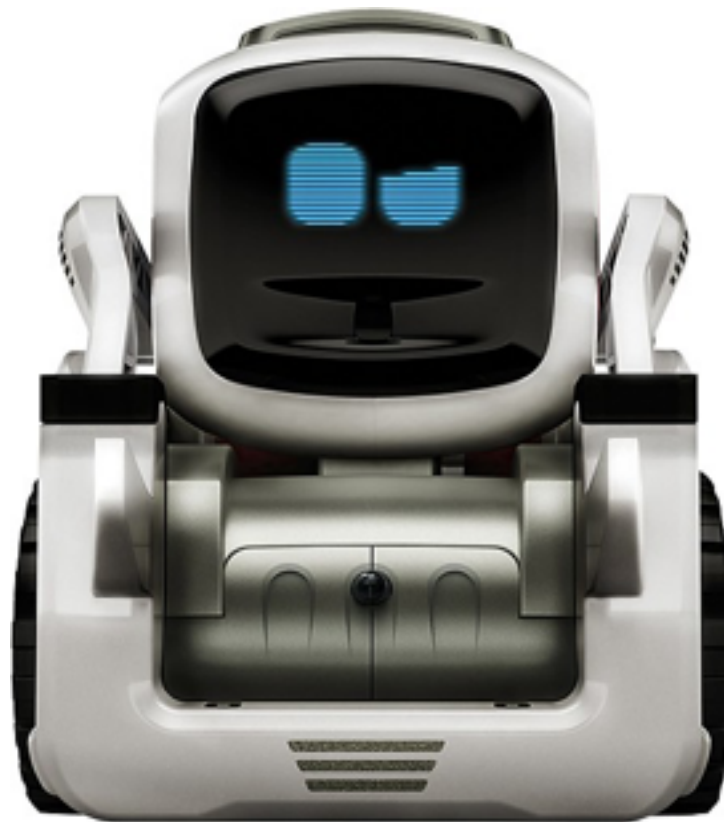
Graph vs. Hypergraph Matrix vs. Tensor

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & & & A_{2n} \\ \vdots & & & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nn} \end{bmatrix}$$

- Hypergraph CAN ALSO be rendered as a simple, 2-mode matrix where each set connects to each concept, but this removes all informational geometry
- Hypergraph rendered as a hyper-matrix or tensor retains its geometry, convexity, etc.



MODELING OPPORTUNITY



Automatically generate promising discoveries

(or feed recipes to scientist chefs)

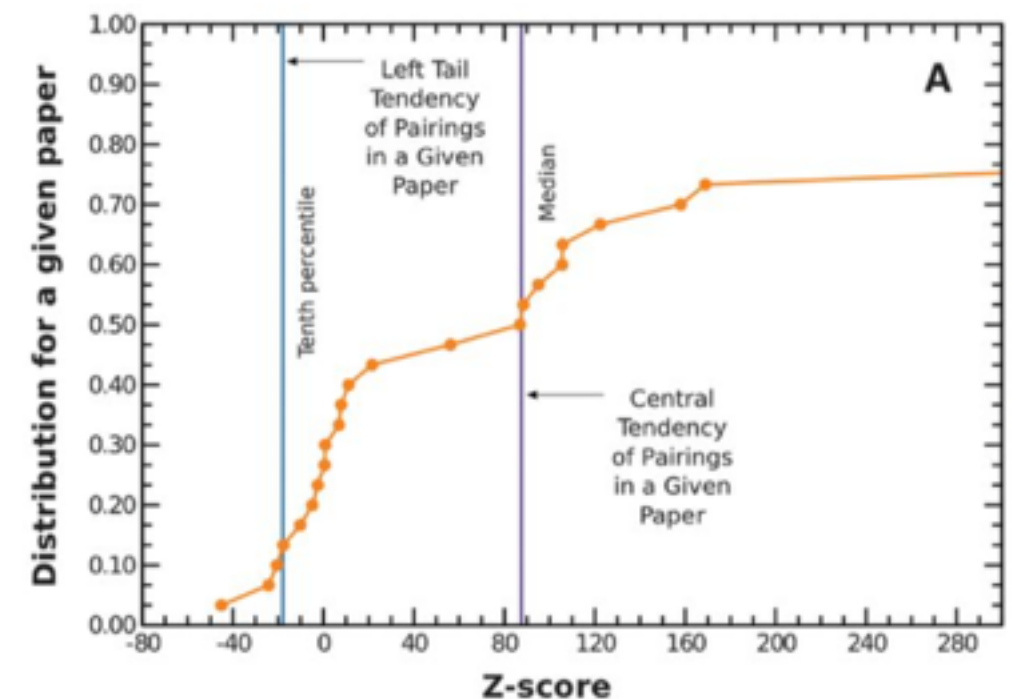
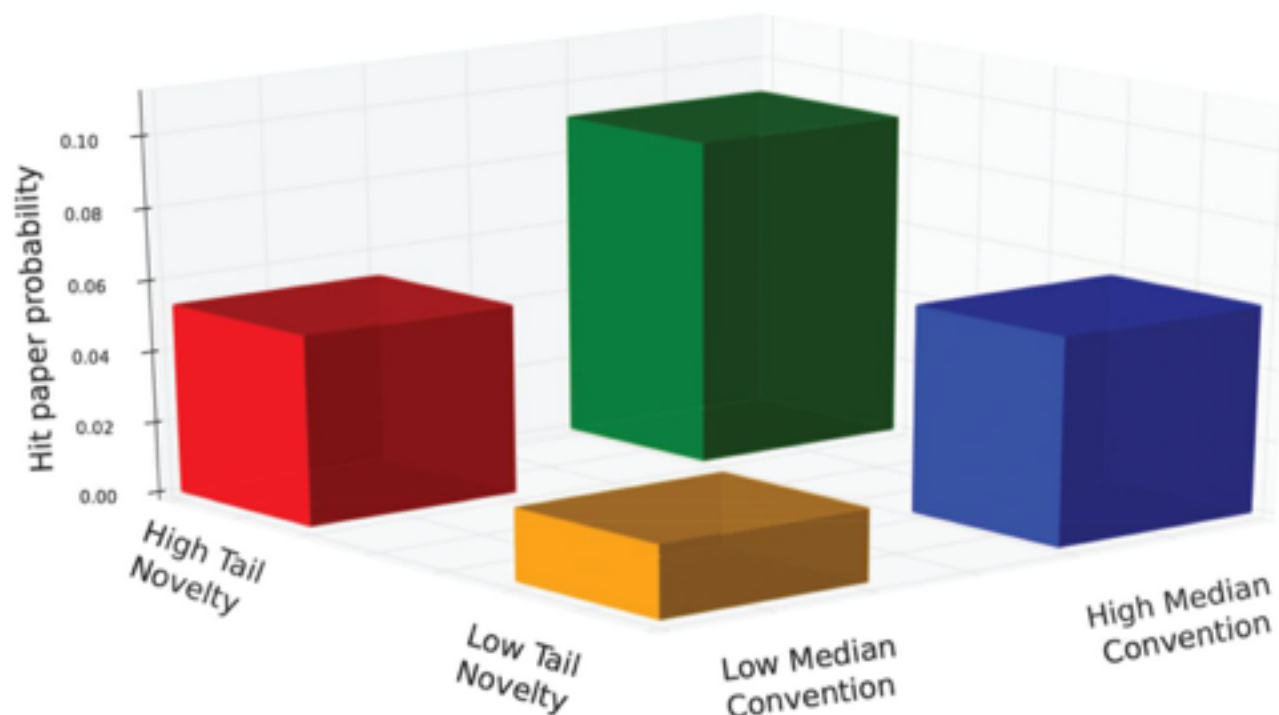
RELATED APPROACHES



REPORT

Atypical Combinations and Scientific Impact

Brian Uzzi^{1,2}, Satyam Mukherjee^{1,2}, Michael Stringer^{2,3}, Ben Jones^{1,4,*}



CONCEPTS = CITED JOURNALS

COMBINATIONS = PAIRWISE
FREQUENCY DISTRIBUTION

PREDICTION < 10%

OUR PROJECT

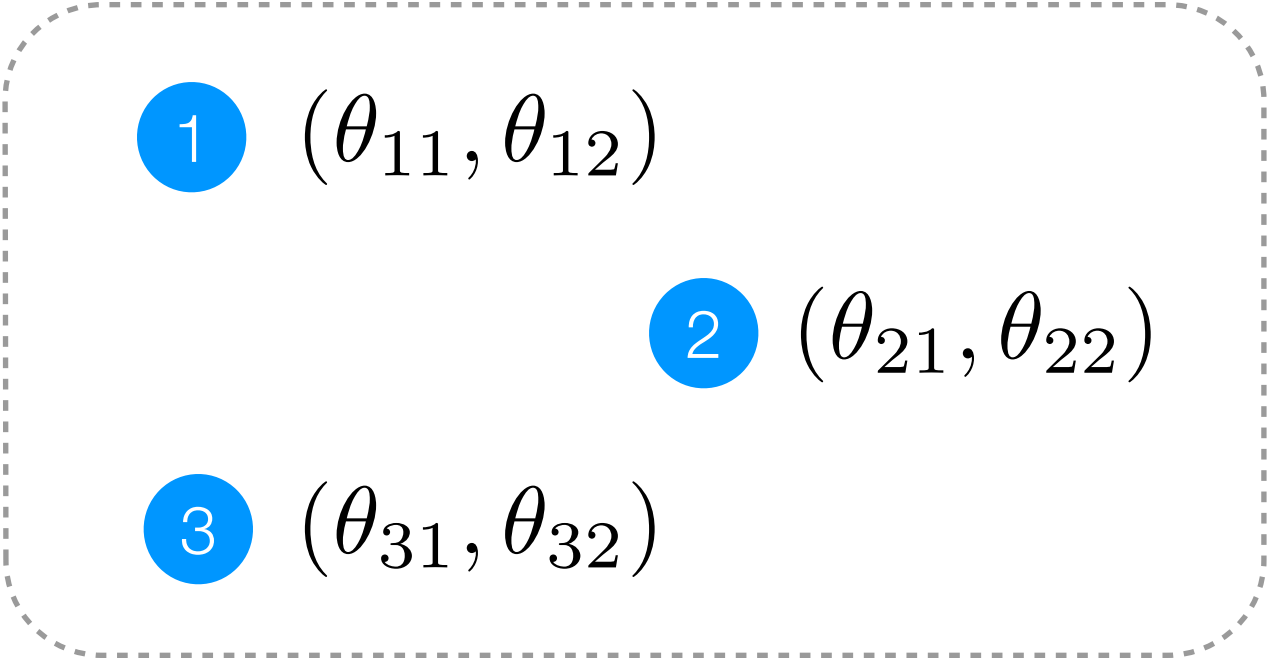
CONCEPTS = CONCEPTS

CONTEXTS = JOURNALS

COMBINATIONS = COMPLETE
COMBINATION

PREDICTION > 40%

Mixed-Membership, High-Dimensional Block Model



1 $(\theta_{11}, \theta_{12})$

2 $(\theta_{21}, \theta_{22})$

3 $(\theta_{31}, \theta_{32})$

Propensity that this combination will turn into a paper:

$$\lambda = (\theta_{11}\theta_{21}\theta_{31} + \theta_{12}\theta_{22}\theta_{32})r_1r_2r_3$$

popularity of node i

Number of papers on this combination: $X \sim \text{Poisson}(\lambda)$

Generative Model for Hypergraph

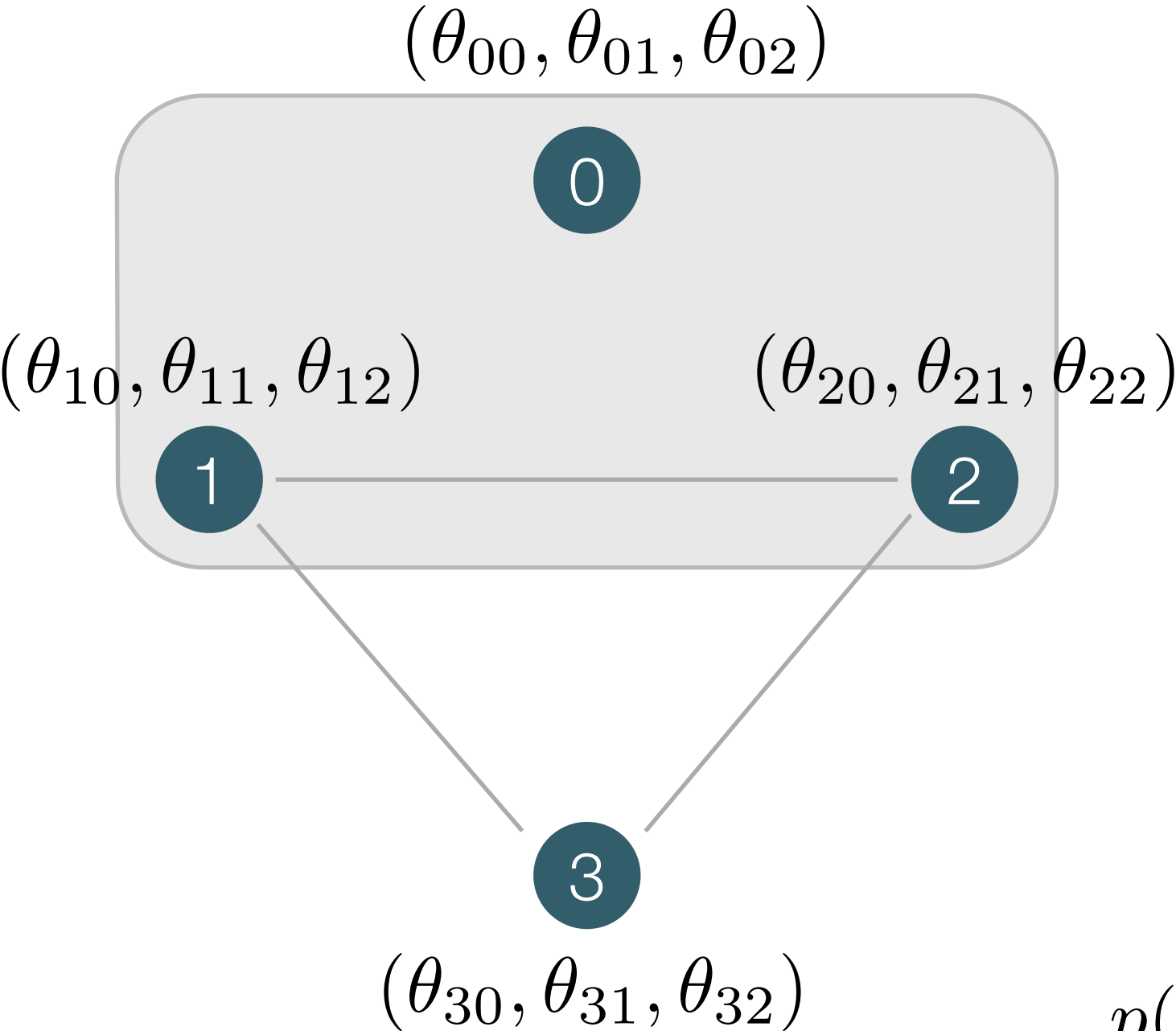
- For any combination h of nodes
- Calculate propensity $\lambda_h = \sum_k \prod_{i \in h} r_i \theta_{ik}$
- Draw the number of hyperedges of h from

$$X_h \sim \text{Poisson}(\lambda_h)$$

- Likelihood to generate the hypergraph

$$p(G|\theta, r) = \prod_{h \in H} p(x_h|\theta, r)$$

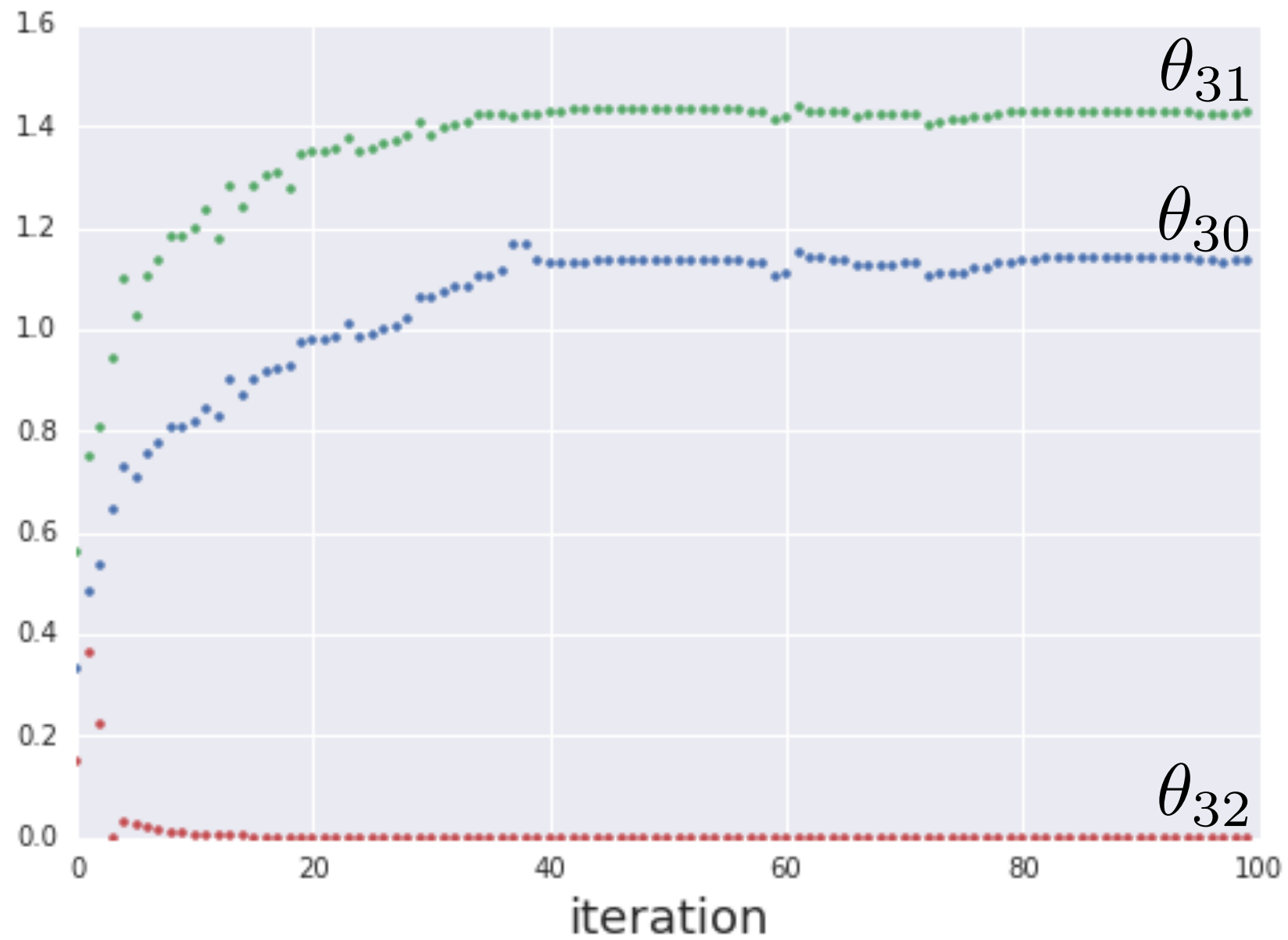
Toy Example



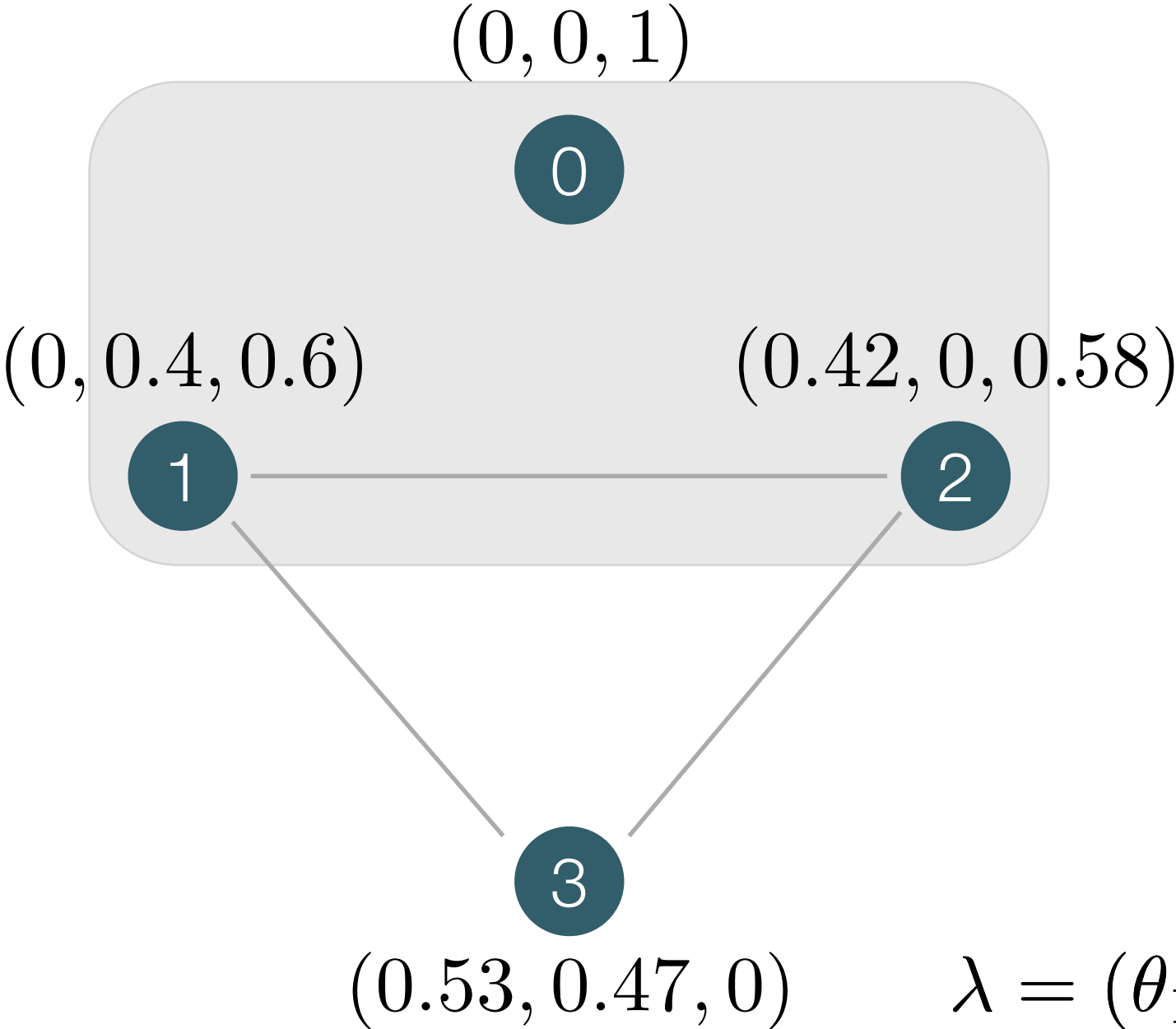
h	X_h
$\{0, 1, 2\}$	1
$\{1, 2\}$	1
$\{1, 3\}$	1
$\{2, 3\}$	1
...	0

$$p(G|\theta, r) = \prod_{h \in H} p(x_h|\theta, r)$$

Toy Example—Convergence



Toy Example

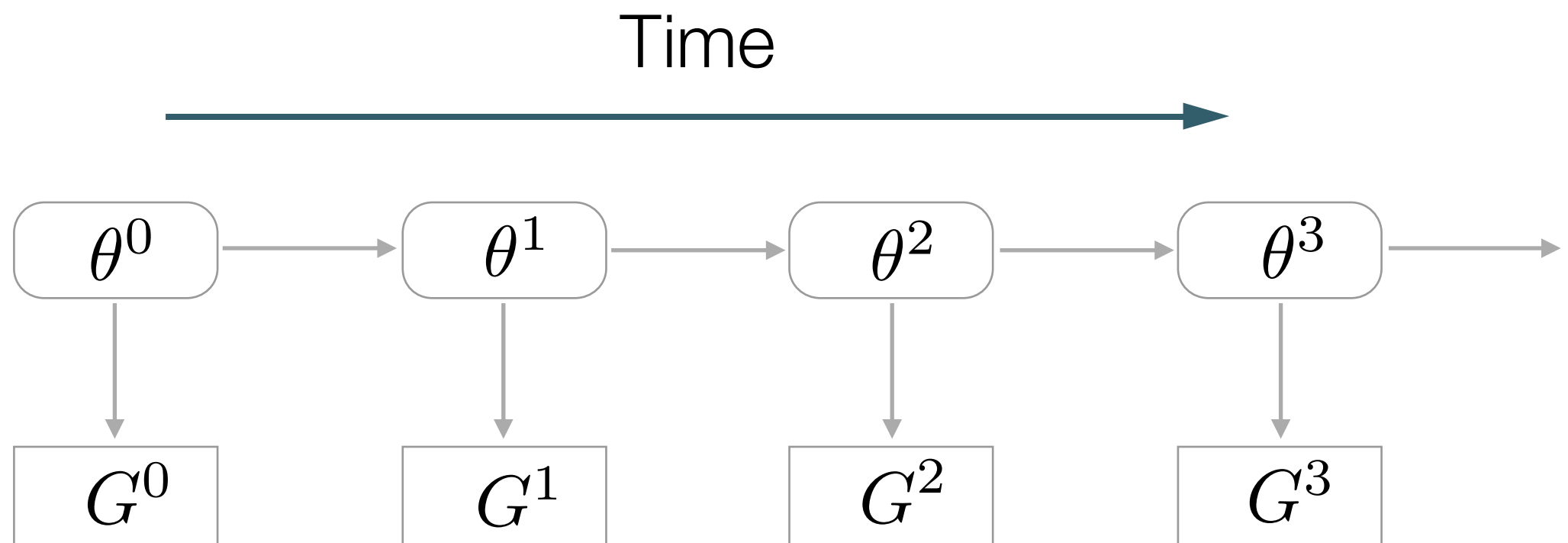


h	X_h
$\{0, 1, 2\}$	1
$\{1, 2\}$	1
$\{1, 3\}$	1
$\{2, 3\}$	1
...	0

$$\lambda = (\theta_{11}\theta_{21}\theta_{31} + \theta_{12}\theta_{22}\theta_{32})r_1r_2r_3$$

Evolution of the Network

Hidden Markov Model



G^t : observed network at time t

θ^t : latent positions of the elements at time t

Complete Model

- Log-likelihood function

$$l(\theta_1, \dots, \theta_T) = \log P(G_1, \dots, G_T | \theta_1, \dots, \theta_T)$$

$$= \sum_{t=1}^T [\log P(\theta^t | \theta^{t-1}) + \log P(G^t | \theta^t)]$$

$$= \sum_{t=1}^T \left[\sum_i \sum_k (\theta_{ik}^t - \theta_{ik}^{t-1})^2 / 2\sigma^2 + \sum_{h \in G^t} (x_h \log \sum_k \prod_{i \in h} \theta_{ik}^t - \sum_k \prod_{i \in h} \theta_{ik}^t) \right]$$

- Impossible to optimize!

Incomputable

2^N possible combinations

Maximal Likelihood Estimate

Algorithm

- Generate t from $1, \dots, T$ uniformly at random
- For $d = 2, \dots, D$, pick a random set H_d^t of combinations of order d from G^t .
- Calculate $S_d^t = \sum_{h \in H_d^t} [x_h \log \sum \prod \theta_{ik}^t - \sum \prod \theta_{ik}^t]$
- Approximate $\nabla l(\theta)$ by $\nabla(\sum_d S_d^t)$
- Update $\hat{\theta} = \hat{\theta} + \eta \nabla l(\hat{\theta})$

Maximal Likelihood Estimate

Theorem

Let $f(\theta, t) = \sum_{d=2}^D S_d^t$ and $t \sim \text{randint}(1, T)$, then

$$E[\nabla f(\theta, t)] = \nabla l(\theta)$$

Corollary

$\hat{\theta}$ will converge to the maximal likelihood estimate.





PubMed

APS
physicsTM



PubMed

APS
physicsTM

UNITED STATES
PATENT AND TRADEMARK OFFICE

uspto

Datasets

20M PubMed articles (1865 to 2015)

15,000 MeSH term Concepts (e.g., PCR, hypertension, DNA, testosterone)

.5M APS articles (1880-2015)

80,000 PACS code Concepts (e.g., neutron star core, lie algebras, polarization)

1.5M US Patents

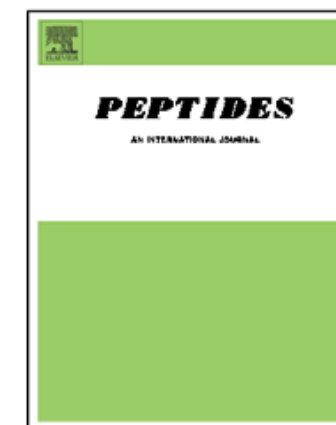
45,000 USPC subclasses (e.g., arc lamp, electrolytic condenser, paper, button)



available at www.sciencedirect.com



journal homepage: www.elsevier.com/locate/peptides



NAP protects hippocampal neurons against multiple toxins

Ilona Zemlyak^{a,b}, Nathan Manley^b, Robert Sapolsky^b, Ilana Gozes^{a,*}

^a Department of Human Molecular Genetics and Biochemistry, Sackler Faculty of Medicine, Tel Aviv University, Israel

^b Department of Biological Sciences, Stanford University, Stanford, USA

ARTICLE INFO

Article history:

Received 1 July 2007

Received in revised form

6 August 2007

Accepted 7 August 2007

Published on line 11 August 2007

Keywords:

NAP

Kainic acid

Oxygen-glucose deprivation

Sodium cyanide

Neuronal death

ABSTRACT

The femtomolar-acting protective peptide NAP (NAPVSIPQ), derived from activity-dependent neuroprotective protein (ADNP), is broadly neuroprotective in vivo and in vitro in cerebral cortical cultures and a variety of cell lines. In the present study, we have extended previous results and examined the protective potential of NAP in primary rat hippocampal cultures, using microtubule-associated protein 2 (MAP2) as a measure for neuroprotection. Results showed that NAP, at femtomolar concentrations, completely protected against oxygen-glucose deprivation, and cyanide poisoning. Furthermore, NAP partially protected against kainic acid excitotoxicity. In summary, we have significantly expanded previous findings in demonstrating here direct neuroprotective effects for NAP on vital hippocampal neurons that are key participants in cognitive function in vivo.

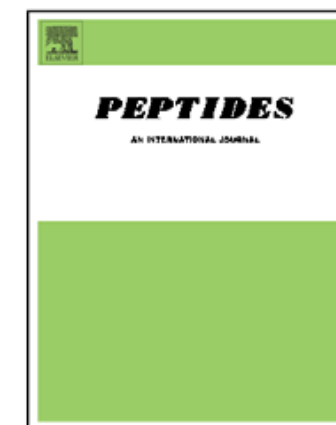
© 2007 Elsevier Inc. All rights reserved.



available at www.sciencedirect.com



journal homepage: www.elsevier.com/locate/peptides



NAP protects hippocampal neurons against multiple toxins

Ilona Zemlyak^{a,b}, Nathan Manley^b, Robert Sapolsky^b, Ilana Gozes^{a,*}

^a Department of Human Molecular Genetics and Biochemistry, Sackler Faculty of Medicine, Tel Aviv University, Israel

^b Department of Biological Sciences, Stanford University, Stanford, USA

ARTICLE INFO

Article history:

Received 1 July 2007

Received in revised form

6 August 2007

Accepted 7 August 2007

Published on line 11 August 2007

Keywords:

NAP

Kainic acid

Oxygen-glucose deprivation

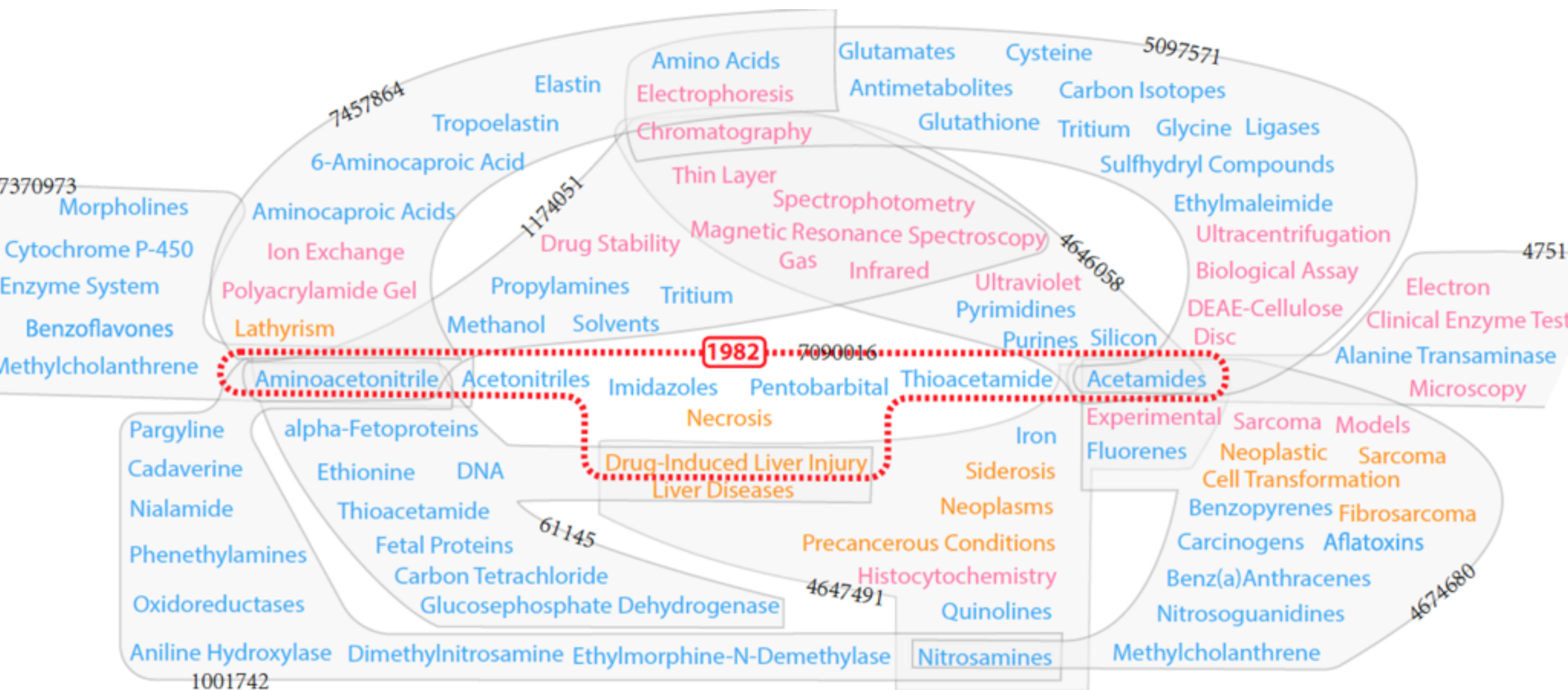
Sodium cyanide

Neuronal death

ABSTRACT

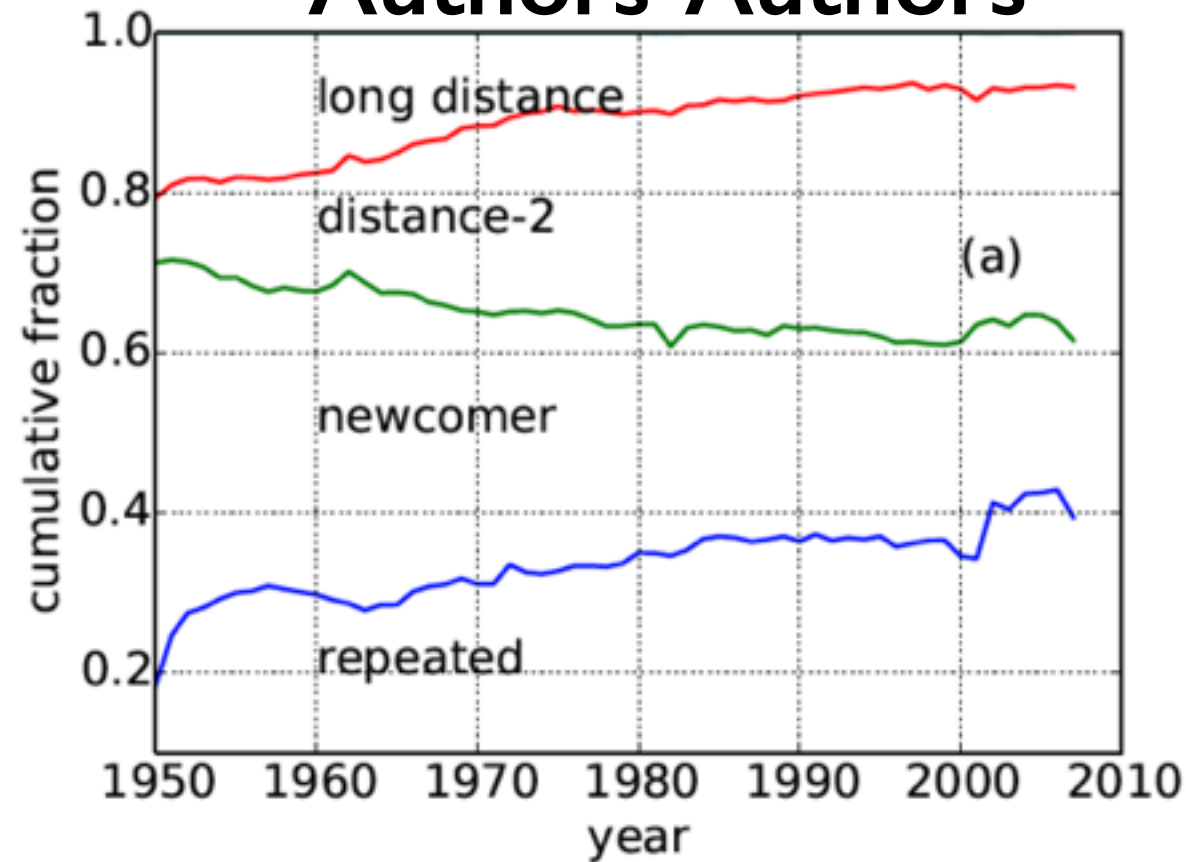
The femtomolar-acting protective peptide NAP (NAPVSIPQ), derived from activity-dependent neuroprotective protein (ADNP), is broadly neuroprotective in vivo and in vitro in cerebral cortical cultures and a variety of cell lines. In the present study, we have extended previous results and examined the protective potential of NAP in primary rat hippocampal cultures, using microtubule-associated protein 2 (MAP2) as a measure for neuroprotection. Results showed that NAP, at femtomolar concentrations, completely protected against oxygen-glucose deprivation, and cyanide poisoning. Furthermore, NAP partially protected against kainic acid excitotoxicity. In summary, we have significantly expanded previous findings in demonstrating here direct neuroprotective effects for NAP on vital hippocampal neurons that are key participants in cognitive function in vivo.

© 2007 Elsevier Inc. All rights reserved.

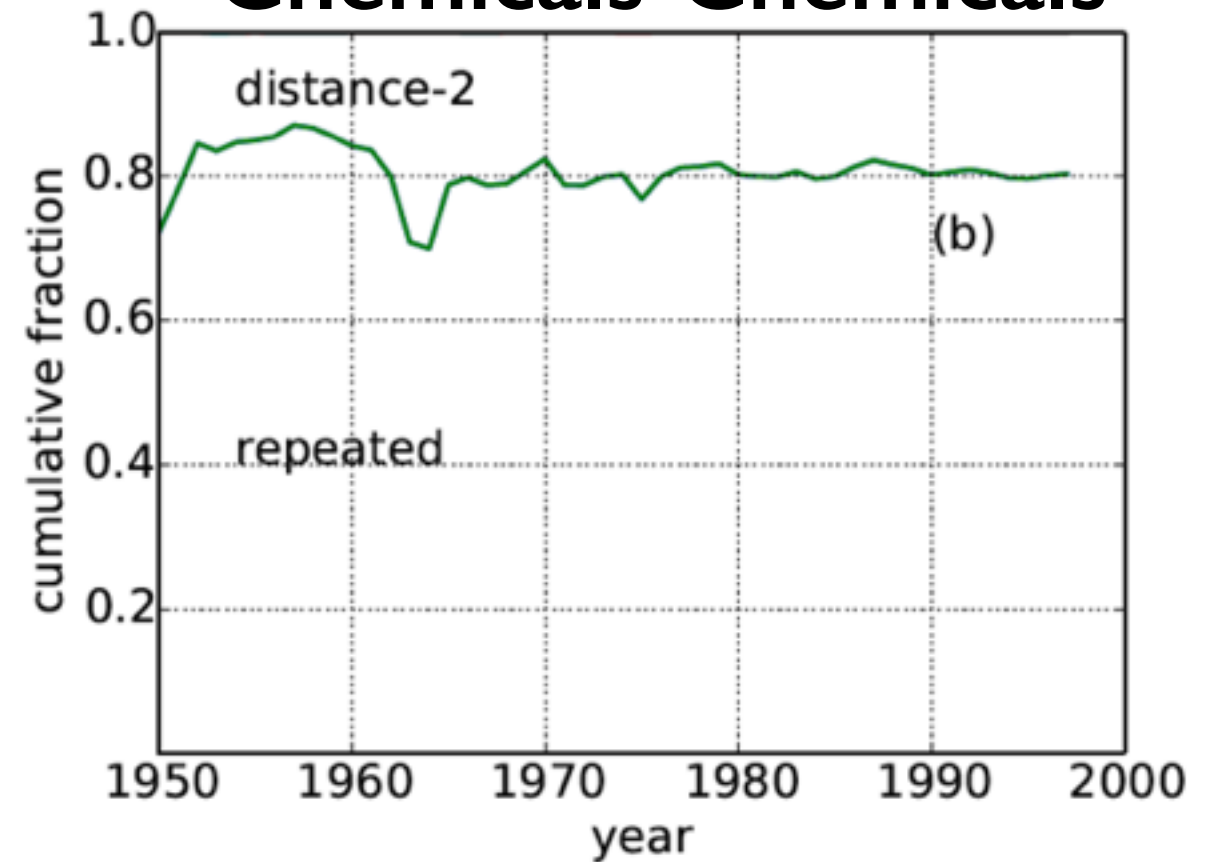


Network Distance

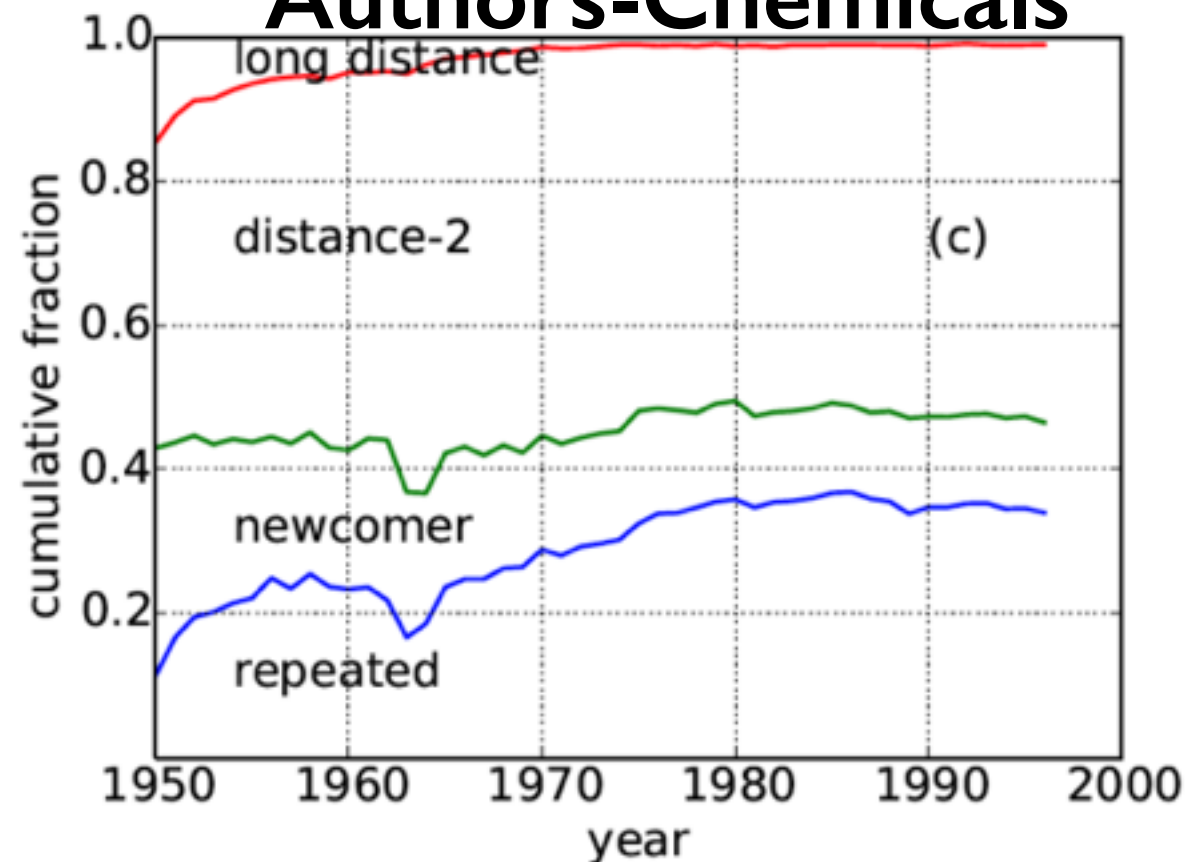
Authors-Authors



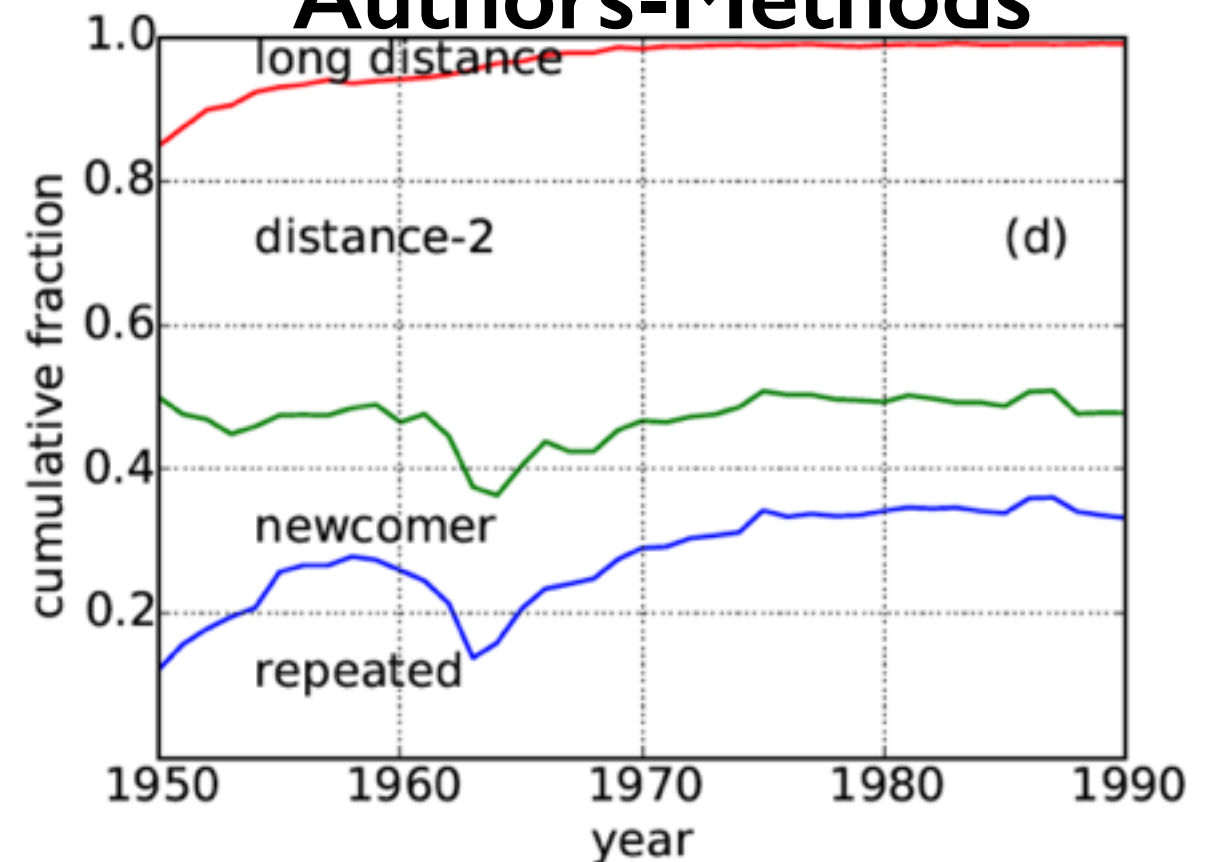
Chemicals-Chemicals



Authors-Chemicals



Authors-Methods



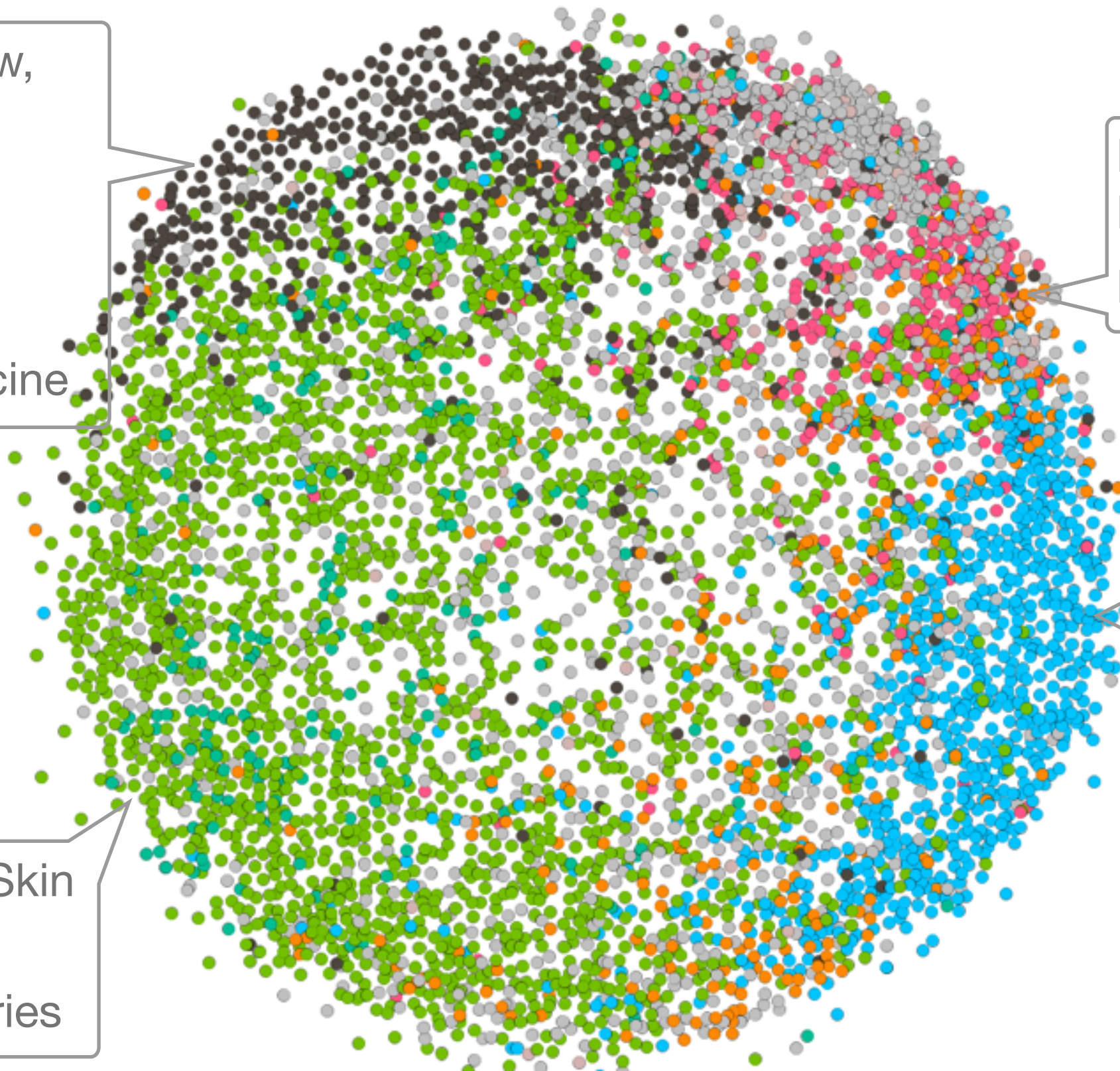
Community Structure

Mouth, Jaw,
Tooth
Diseases,
Dental
Materials,
Oral Medicine

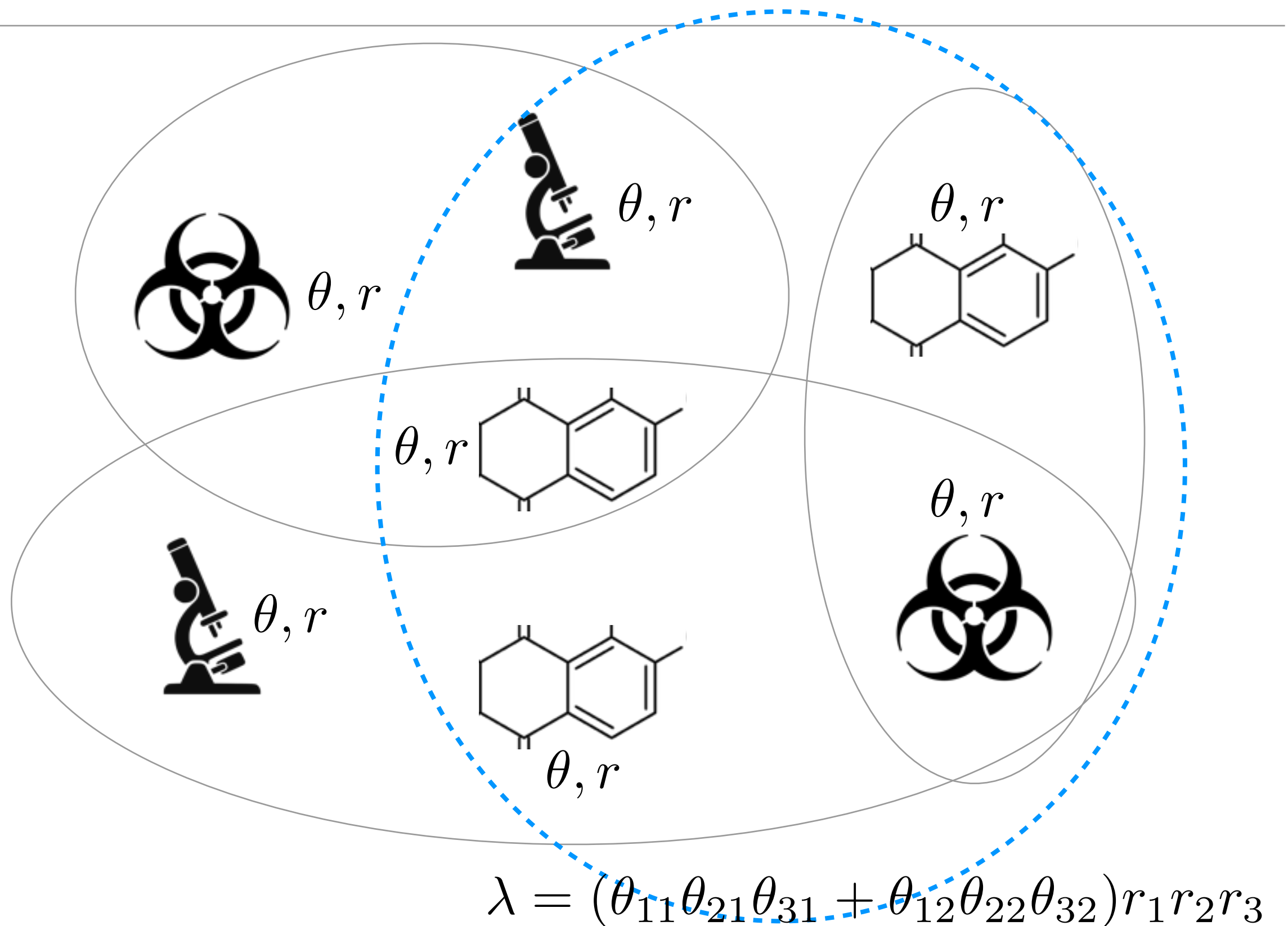
Parasitic
Diseases,
Helminthiasis

Marijuana Abuse,
Phencyclidine
Abuse,
Alkaloids,
Decompression,
Substance Abuse
Detection

Eye, Ear, Skin
Diseases,
Body Injuries

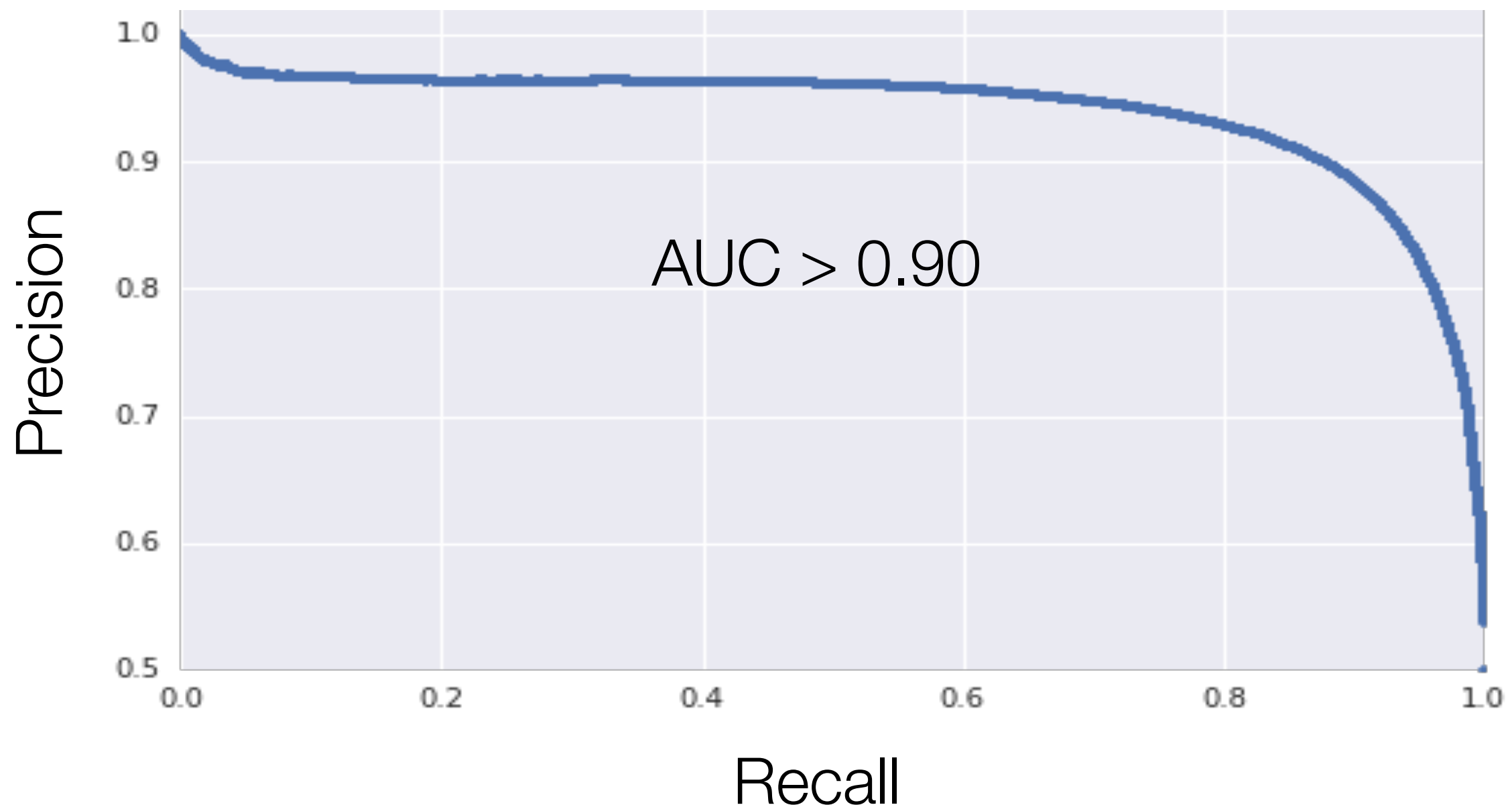


Predict New Hyperedges

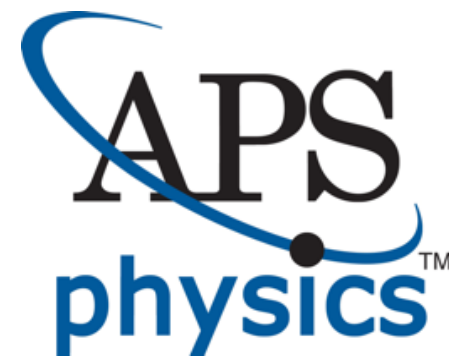


Predicting Papers

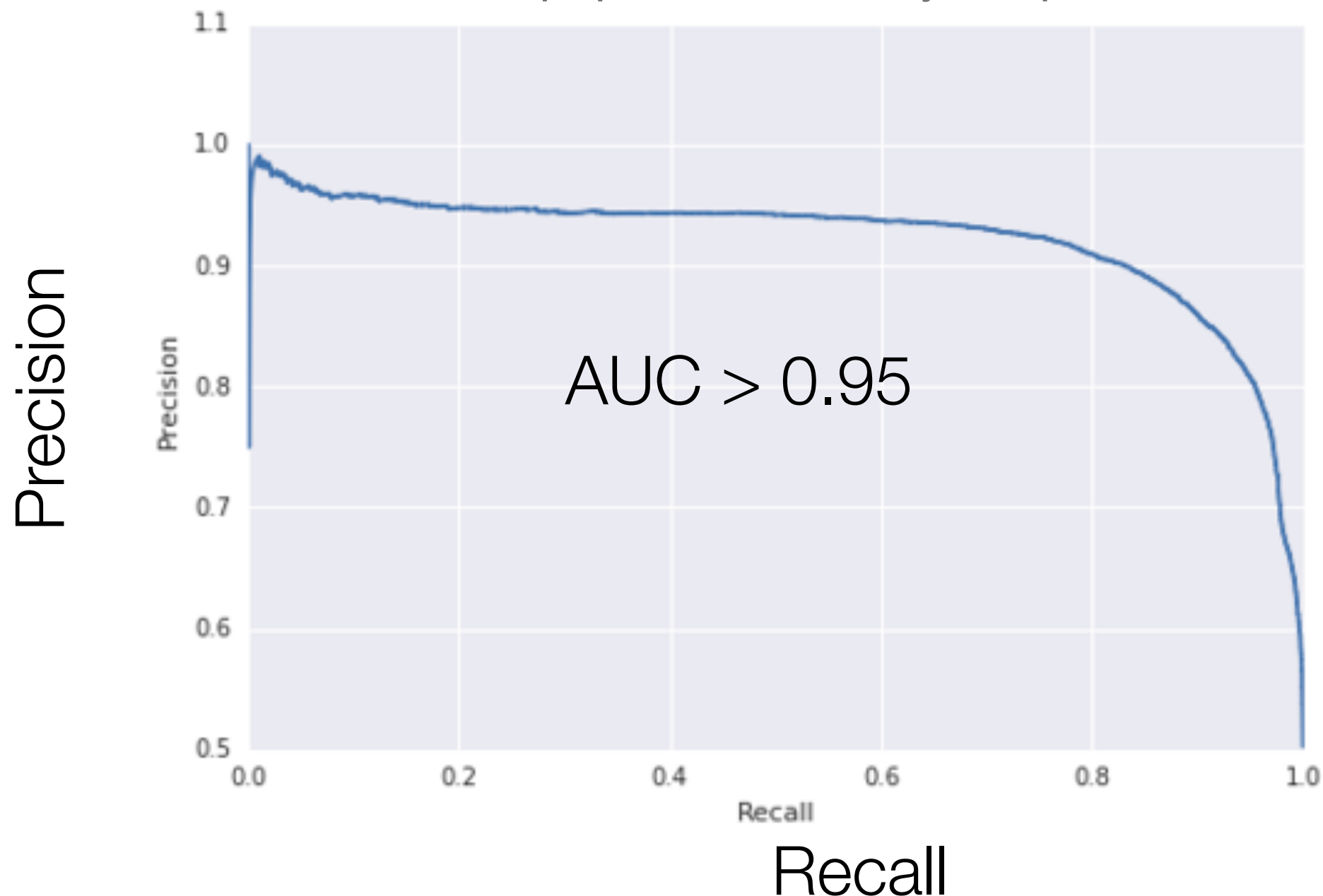
- Precision: out of all predicted papers, how many actually happen
- Recall: out of all future papers, how many are predicted to happen



Predicting Papers



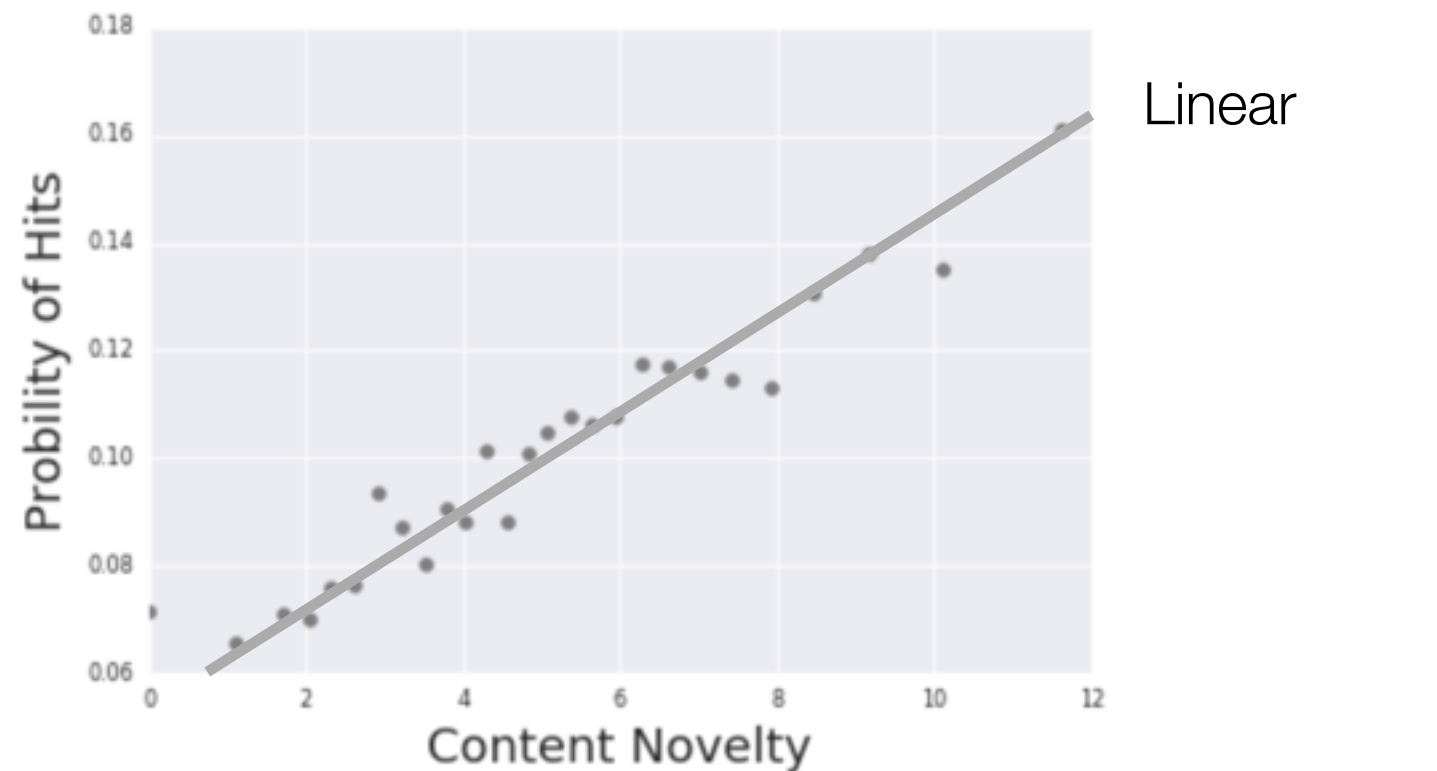
- Precision: out of all predicted papers, how many actually happen
- Recall: out of all future papers, how many are predicted to happen



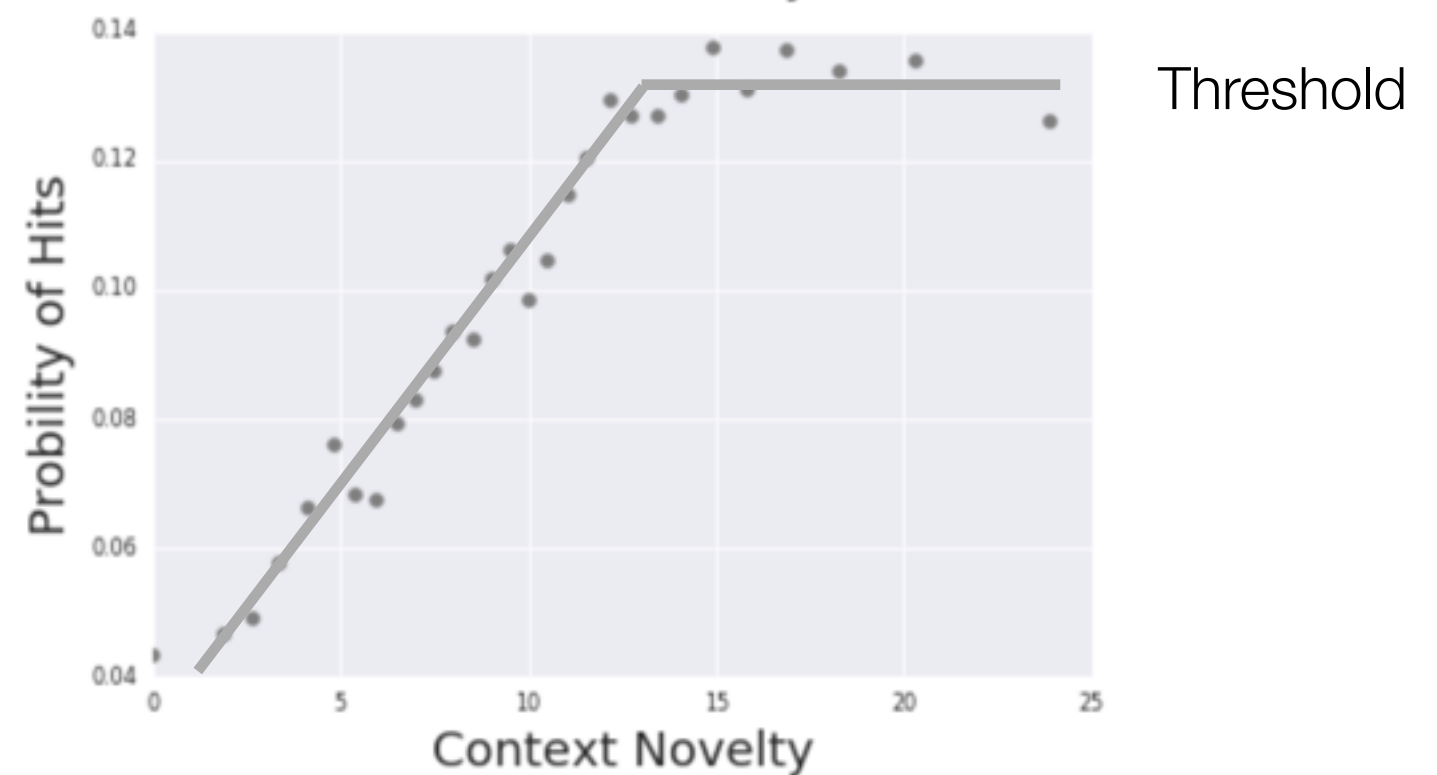
Novelty and Impact



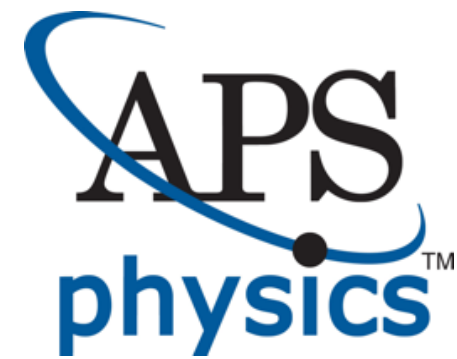
Searching Broadly



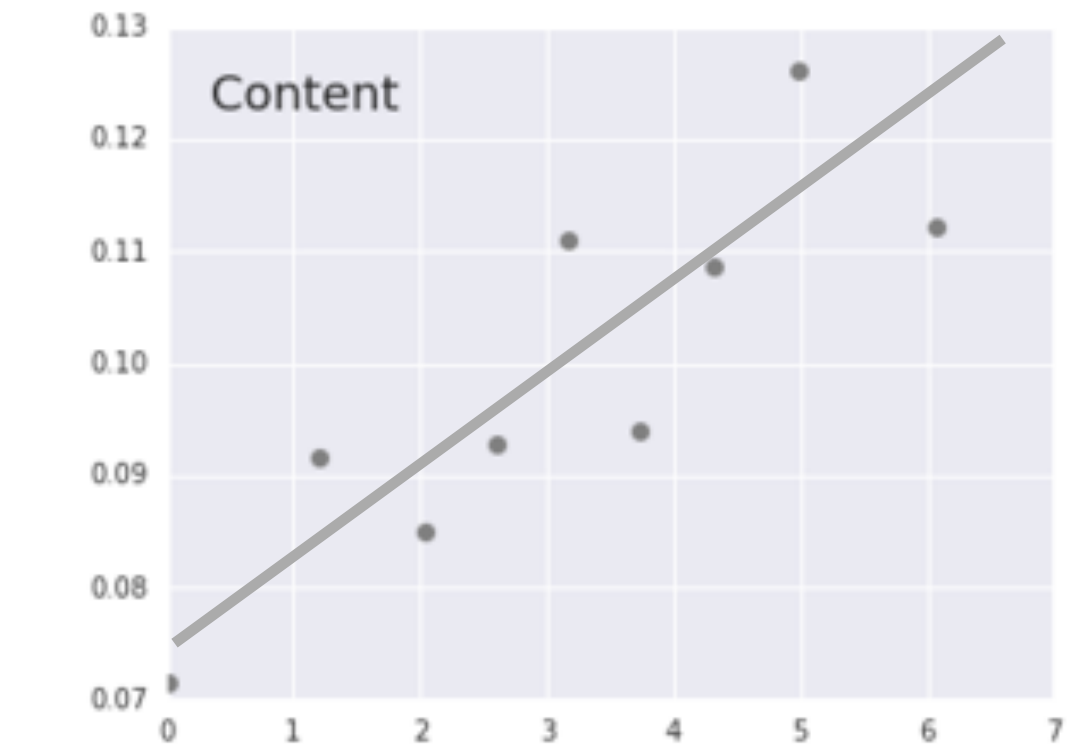
Not Citing too Broadly



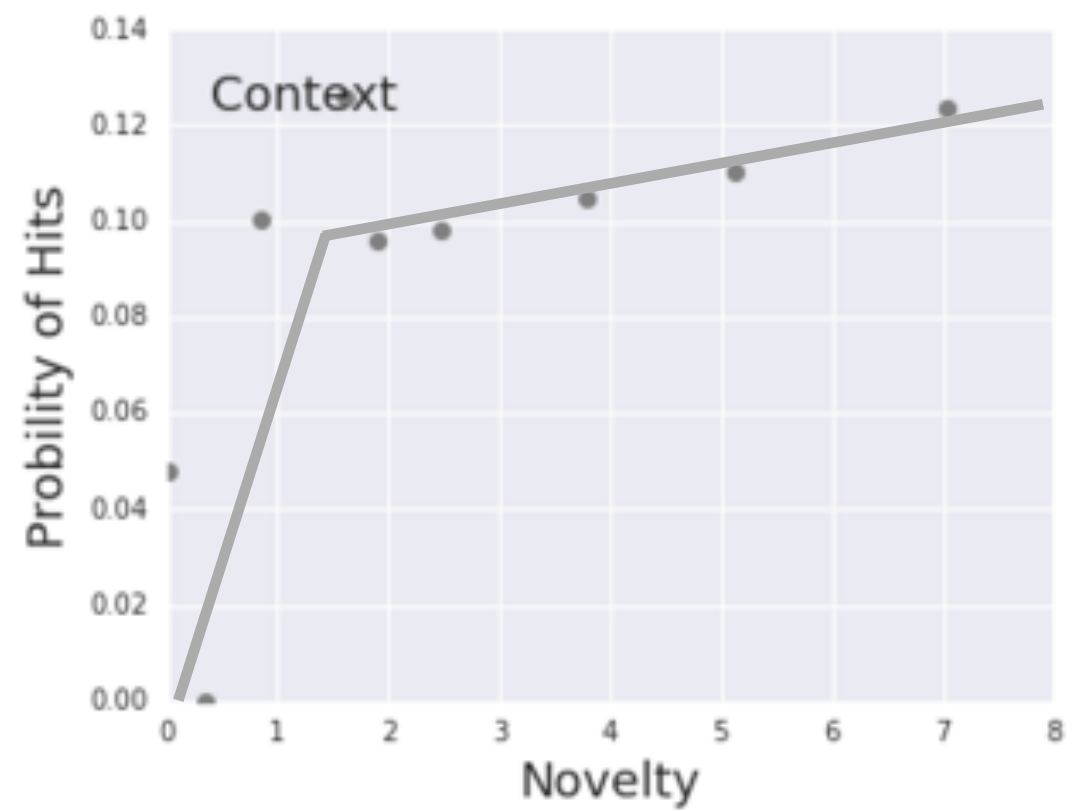
Novelty and Impact



Searching Broadly



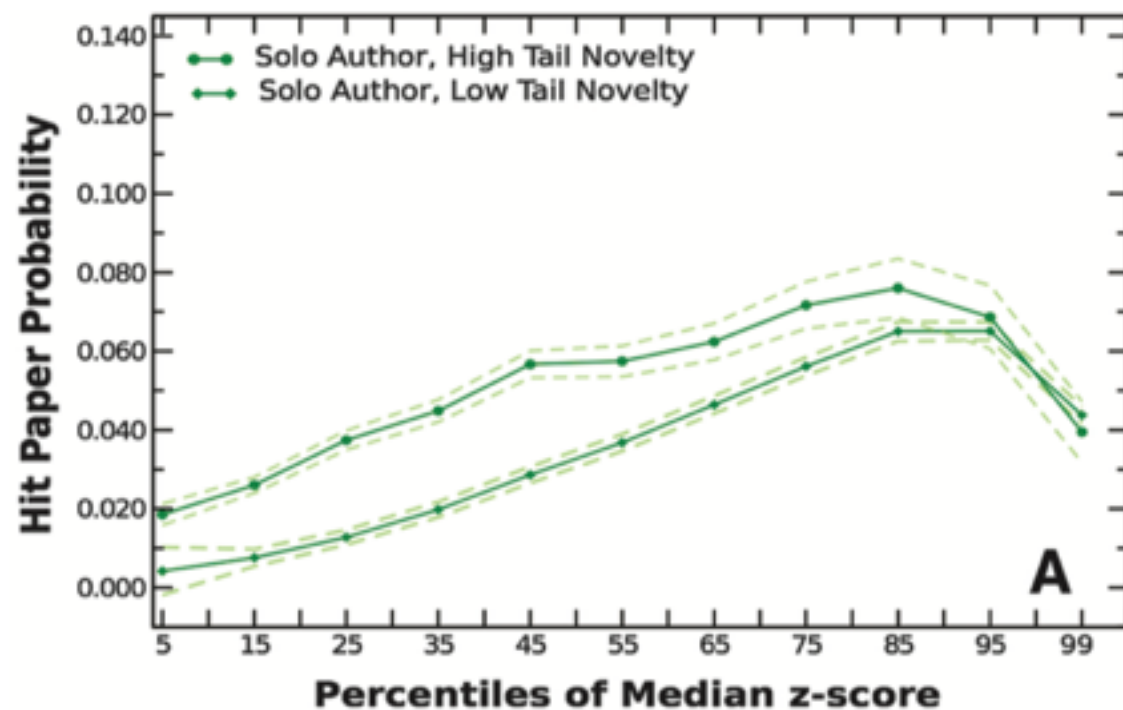
Not Citing too Broadly



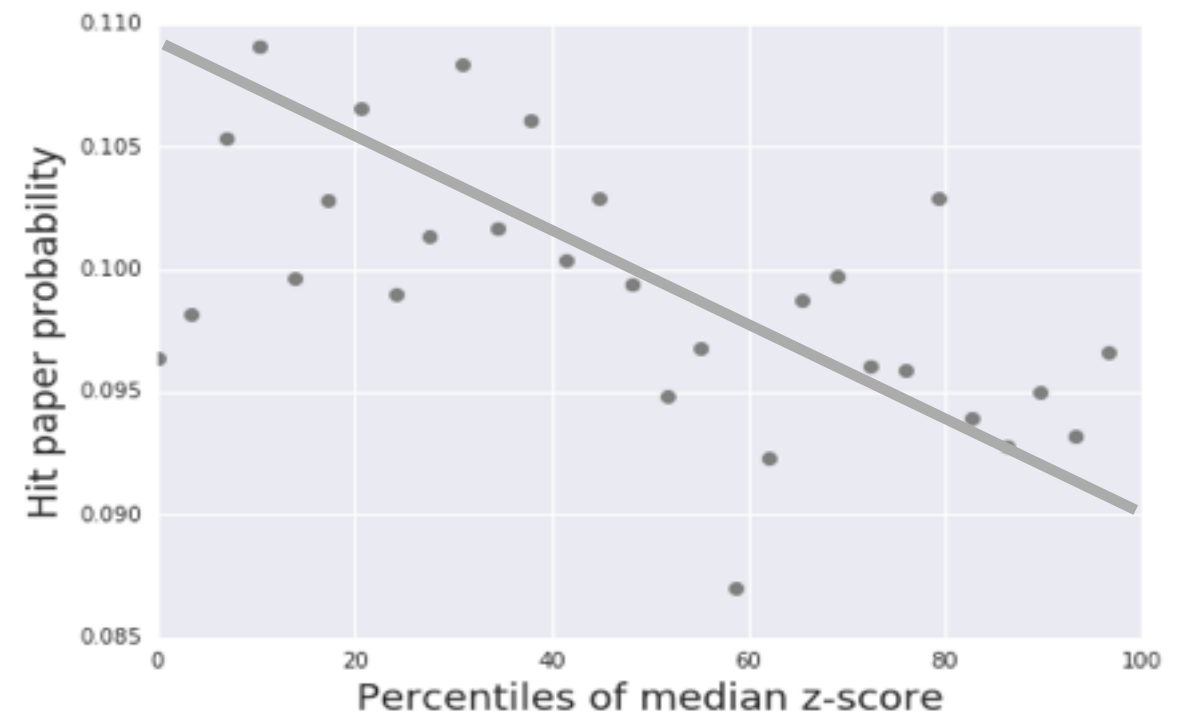
Content vs. Context

Content Correlates at **<.1** with Context
Context does NOT proxy for Content

Context vs. Content



Impact increases with conventionality of journal combinations [Uzzi, et al.]



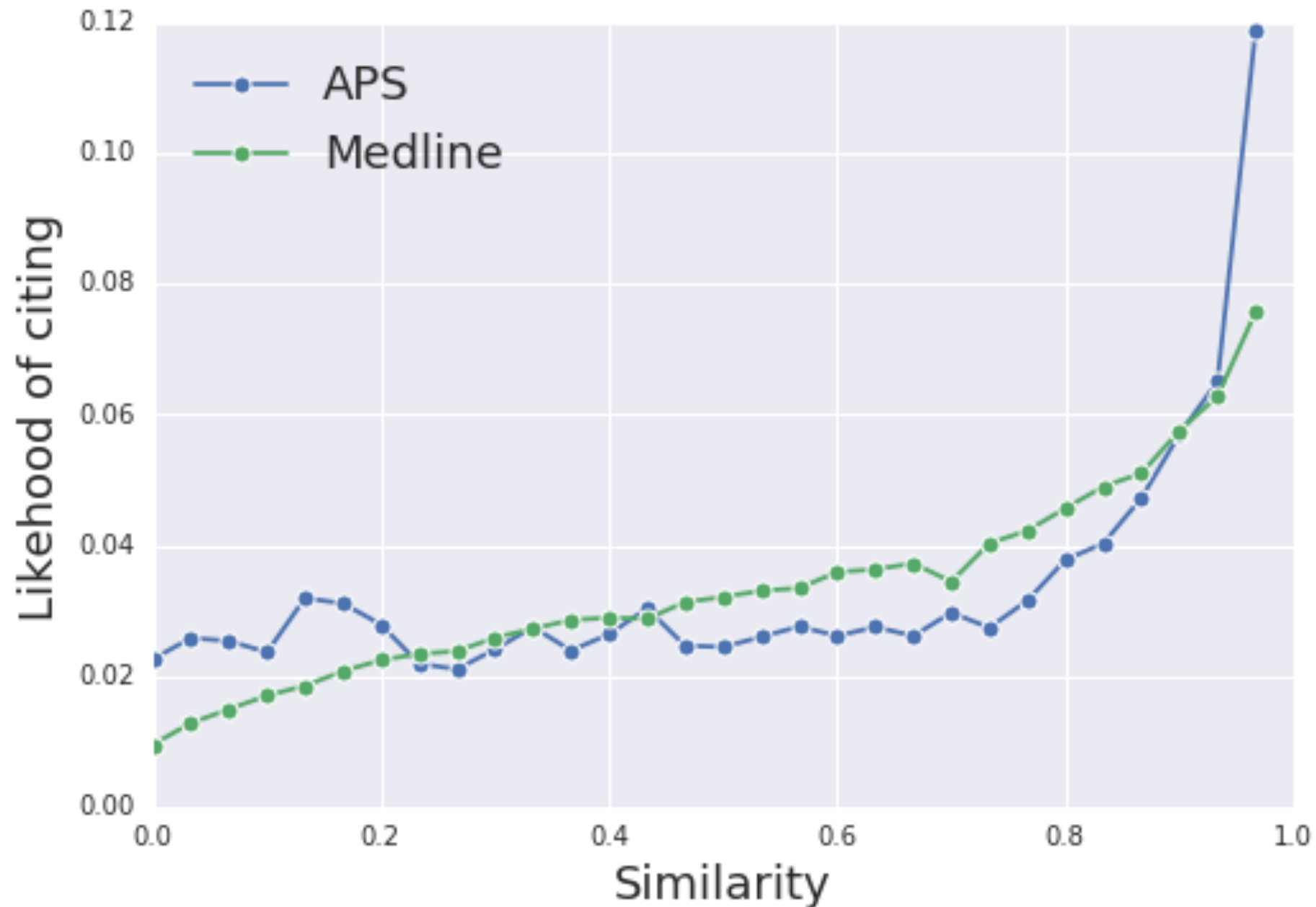
Impact (weakly) decreases with conventionality of MeSH term combinations.

Same methodology [Uzzi, et al.], but contradictory results?

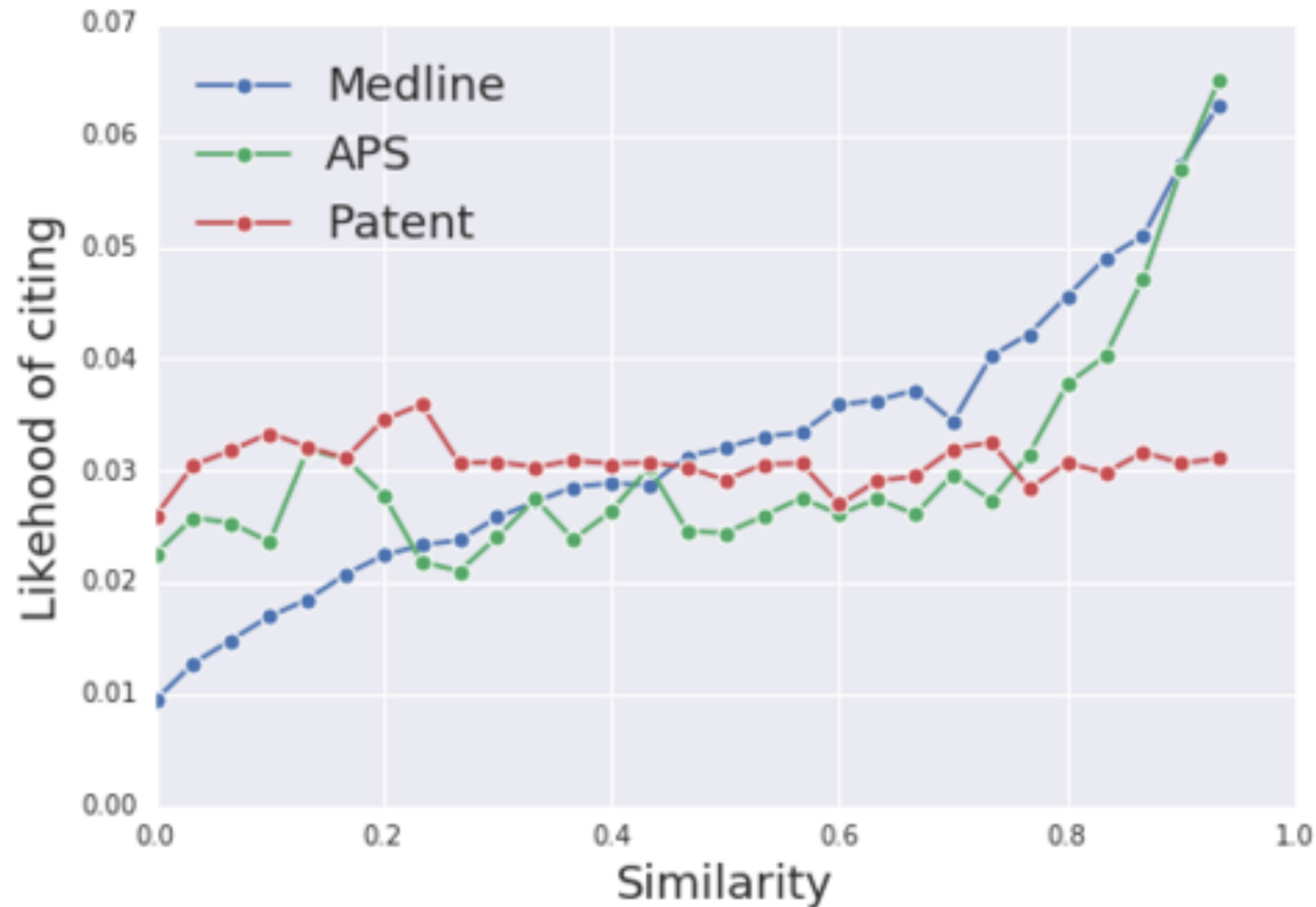
1. Technically: Not enough MeSH terms per paper to calculate median and tail.
2. Conceptually: Reference list is intended to situate a paper in the literature
3. Mismatch between content and context?

Humble Innovation:

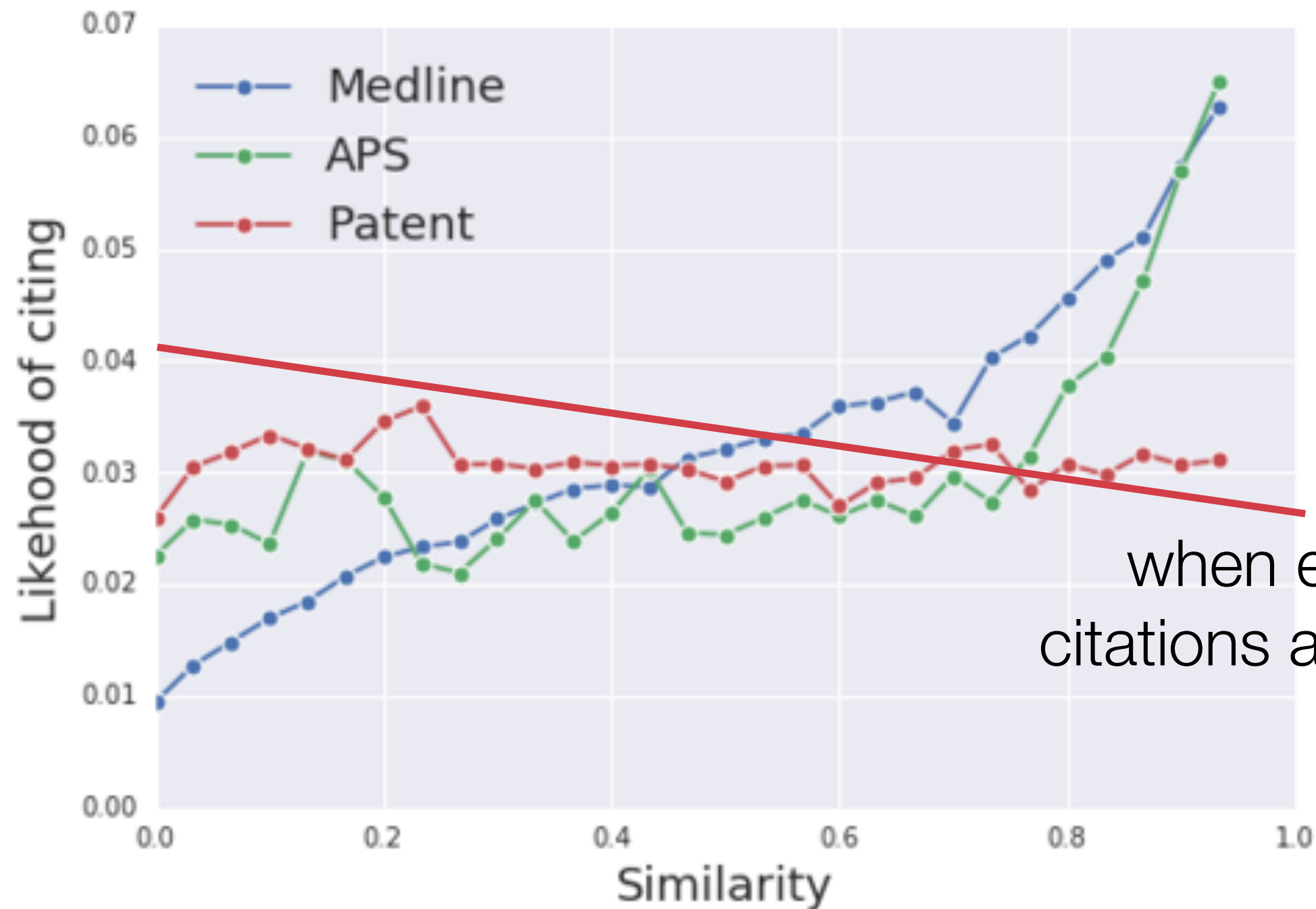
Novel Exploration and Conservative Claims



Audacious Invention: Novel Exploration and Outsized Claims



Audacious Invention: Novel Exploration and Outsized Claims



when examiner
citations are removed

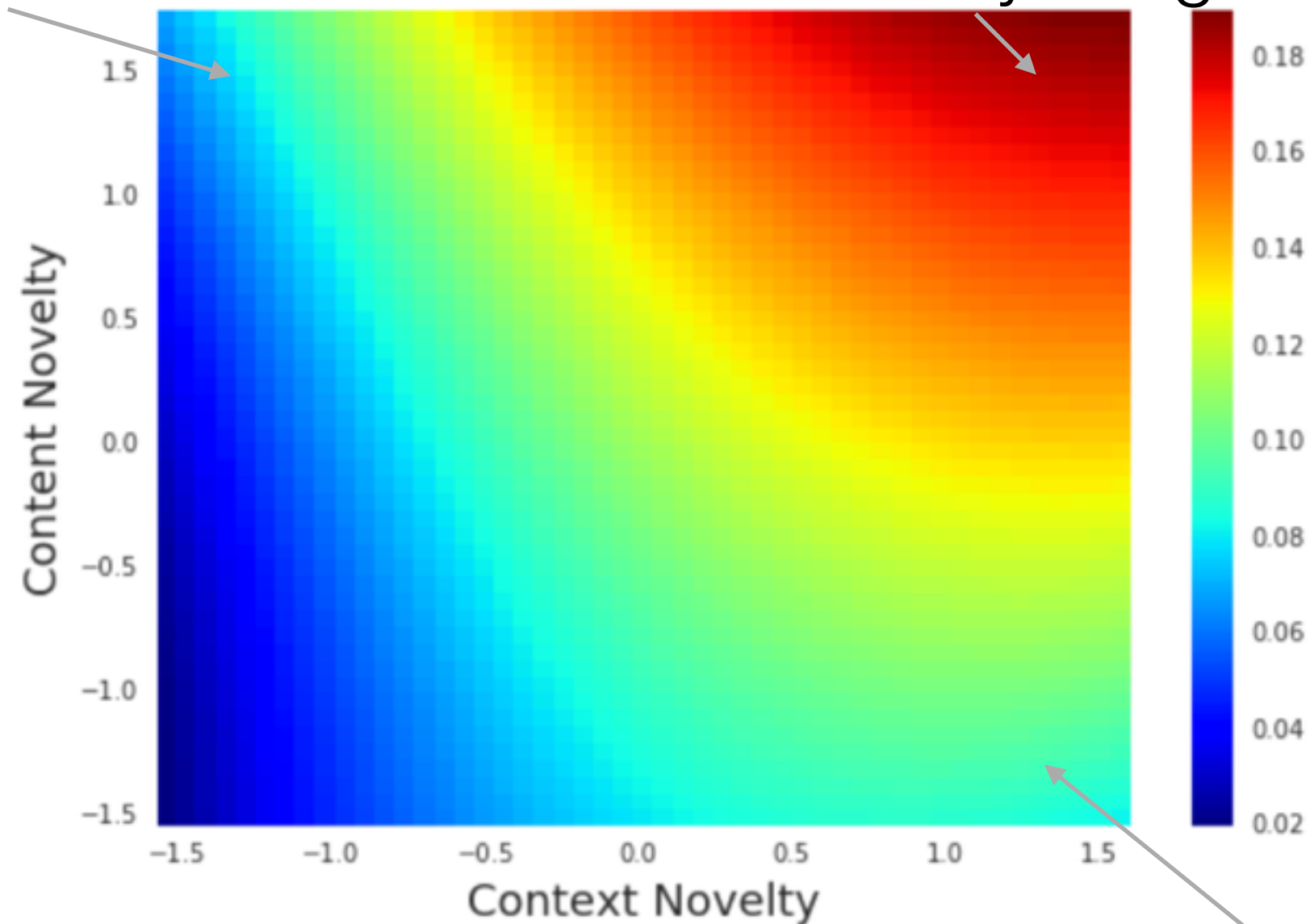
Content vs. Context

Content Correlates at $\sim \mathbf{.1}$ with Context
Context does NOT proxy for Content

Novelty and Impact

Novel Knowledge from diverse fields
unlikely imagined

Novel
but obvious

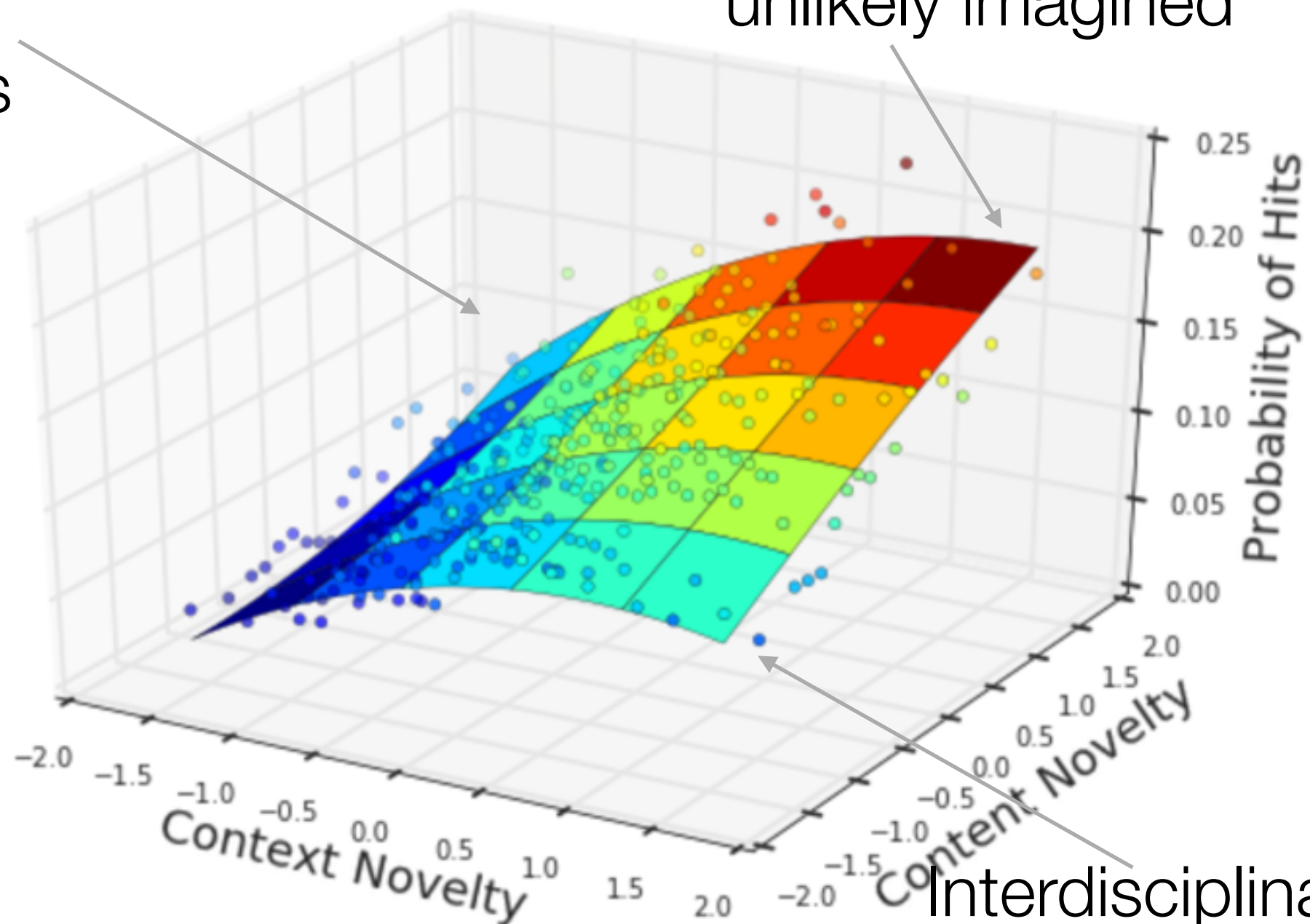


Interdisciplinary
but focused

Novelty and Impact

Novel Knowledge from diverse fields
unlikely imagined

Novel
but obvious



Interdisciplinary
but focused

Doubling Sensitivity with High Dimensionality

- 12-13% - pairwise context (Uzzi's method)
- **25%** - high dimensional context (our method)
- 10% - pairwise content (Uzzi's method)
- **20%** - high dimensional content (our method)

**Critical for estimating the effects for sparse signal—
content similarity**

Content & Context

- ~ 0.0 correlation between content and context novelty
- **30%** of hit probability captured by hypergraph of context + context
- Not linearly additive, but SUBSTANTIAL marginal effect



**Scientists think through the
complex network of content
conditional on context**

Science Think **Differently** from **Scientists** Think

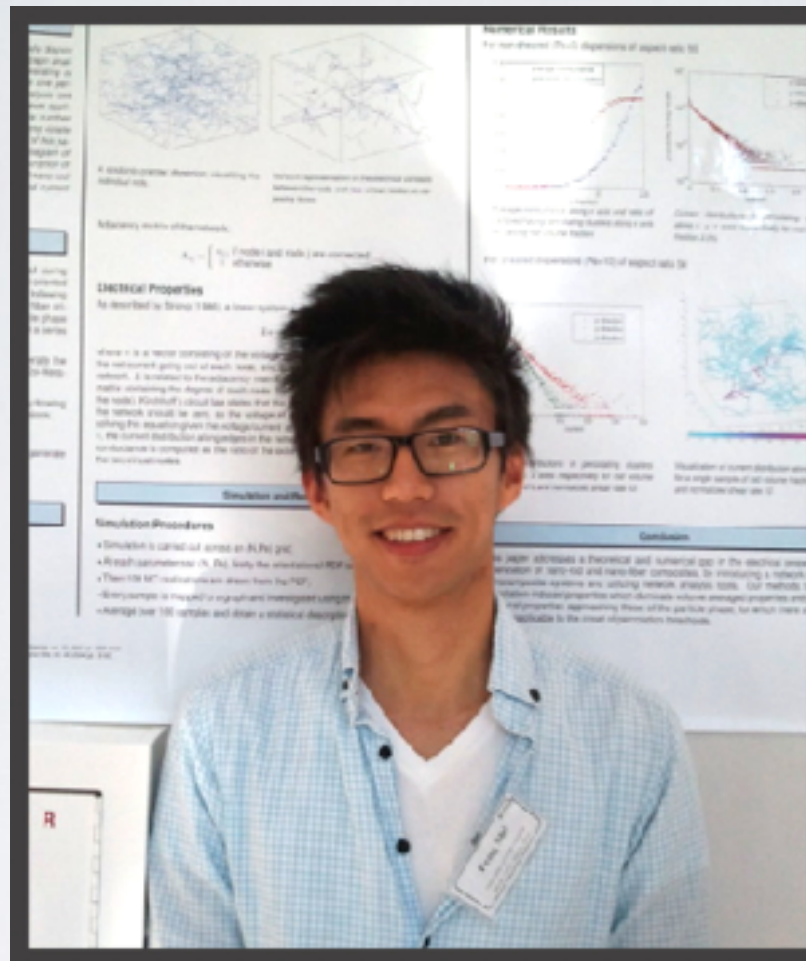
Science thinks like a **Global Bayesian**

...by conditioning success/impact on affirmation of global priors

Scientists have *much* weaker priors

...but succeed by **appearing to build on the shoulders of their audience**

Negative crowd-sourcing - finding combinations unlikely to have been imagined nearly doubled the likelihood of success



Bill Shi