

# semimetric knowledge networks



informatics  
luis rocha 2007

applications to information retrieval and social data mining

**luis m. rocha**

**Indiana university**

school of informatics

1900 East Tenth Street, Bloomington IN 47406  
and

**Instituto Gulbenkian de Ciência**  
Apartado 14, 2781-901 Oeiras, Portugal

rocha@indiana.edu  
<http://informatics.indiana.edu/rocha>

**FLAD**  
Computational Biology Collaboratorium



**INDIANA**  
UNIVERSITY

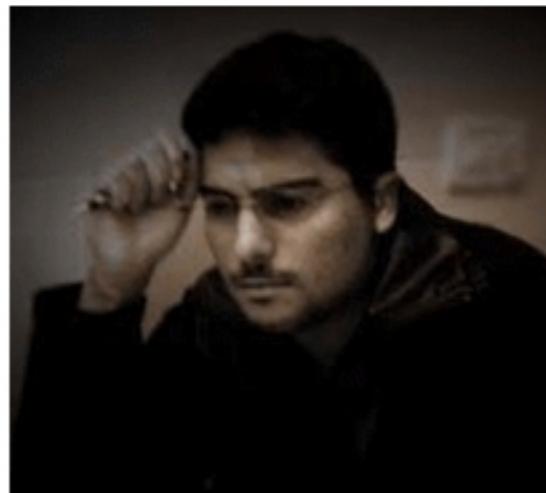


# agent-based model of genotype editing

CASCI Team: <http://casci.informatics.indiana.edu>



informatics  
luis rocha 2007



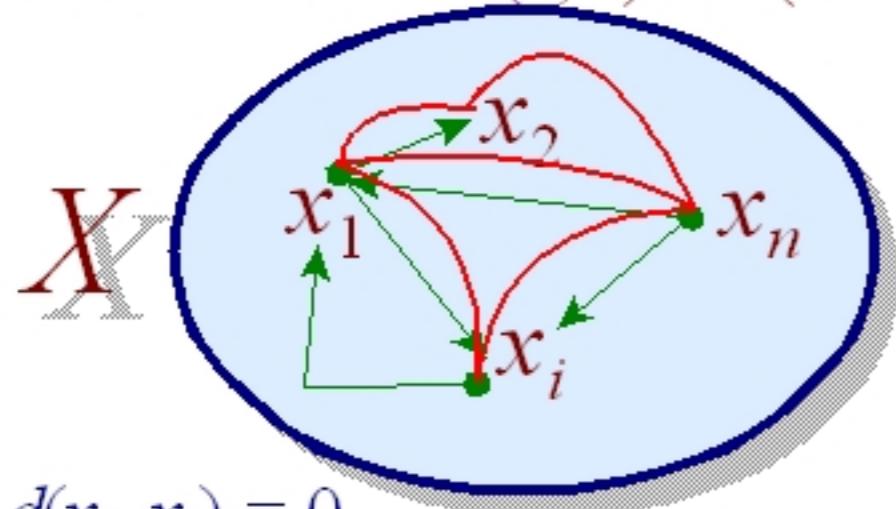
**Alaa Abi-Haidar    jasleen kaur    Tiago Simas    Ana Maguitman**

acknowledgements

Luis Bettencourt, Bharat Dravid, Mariella DiGiacomo, Fil Menczer, Pedrag Ravidojac, Andreas Retchsteiner, Elliot Smith, Karin Verspoor, Paul Wang.

measured from associative “knowledge” graphs

$d$  is a distance function on set  $X$  if it is a nonnegative, symmetric, real-valued function such that  $d(x, x) = 0$  (Shore & Sawyer 1993)



$$d(x_i, x_i) = 0$$

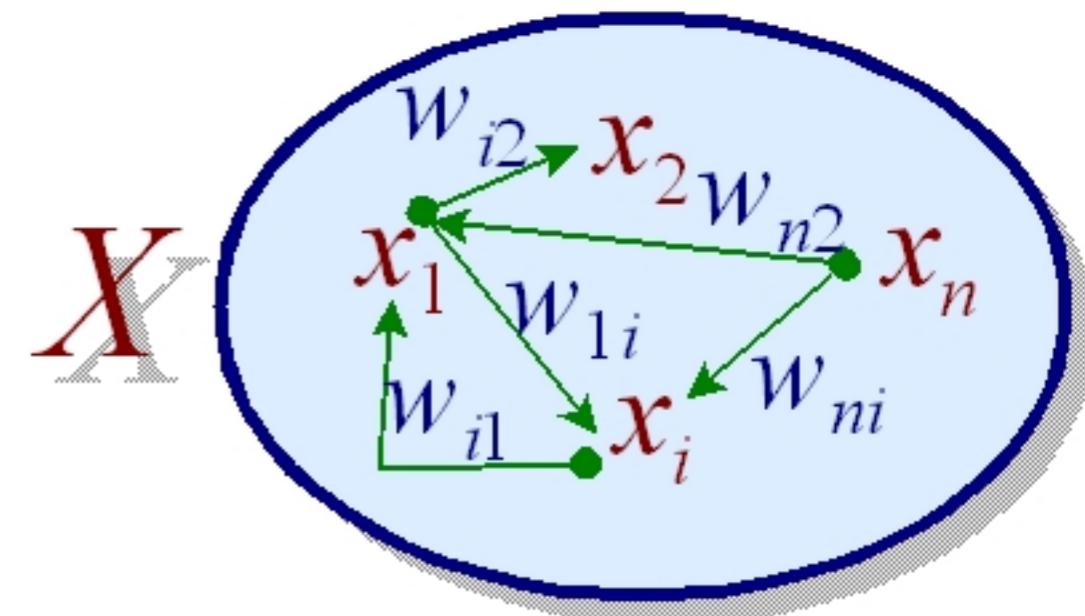
$d(x_i, x_j) = 1$ , if there is an edge

$d(x_i, x_k) = d(x_i, x_j) + \dots + d(x_l, x_k) \leq 1$ ,  
if there is a path

Due to the symmetry requirement,  
distance functions yield non-directed  
distance graphs

$$d(x_1, x_2) \leq d(x_1, x_3) + d(x_3, x_2)$$

**Metric:** the smallest distance between  
nodes is always the most direct path



In real-valued weighted graphs, derived  
distance functions can be semi-metric

$$d(k_1, k_2) \geq d(k_1, k_3) + d(k_3, k_2)$$

Semi-metric

In graphs used to store  
“knowledge”, what does  
it mean?

## operations

$$A(x) : X \rightarrow [0, 1]$$

## Standard Fuzzy Operations

$$\bar{A}(x) = 1 - A(x)$$

$$(A \cap B)(x) = \min[A(x), B(x)]$$

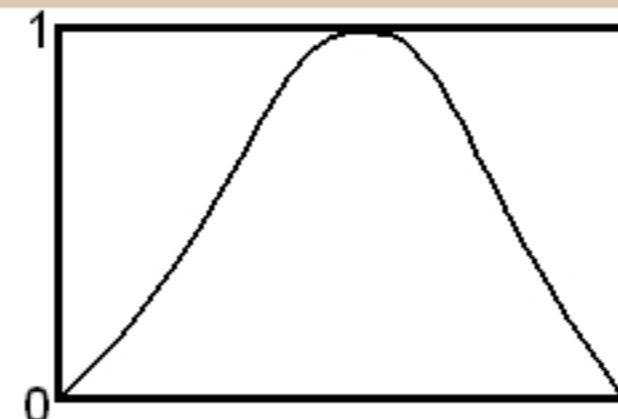
$$(A \cup B)(x) = \max[A(x), B(x)]$$

## Follows

- ▶ Involution, commutativity, associativity, distributivity, Identity, De Morgan's Laws, etc

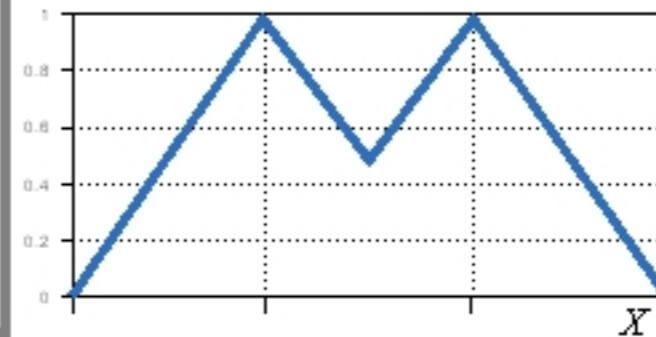
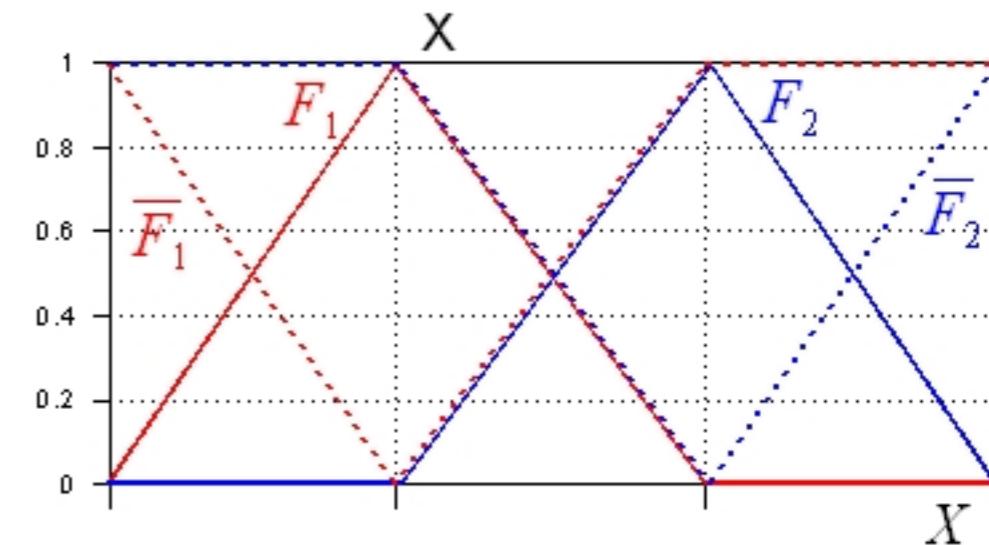
## Does not Follow

- ▶ Laws of contradiction and excluded middle

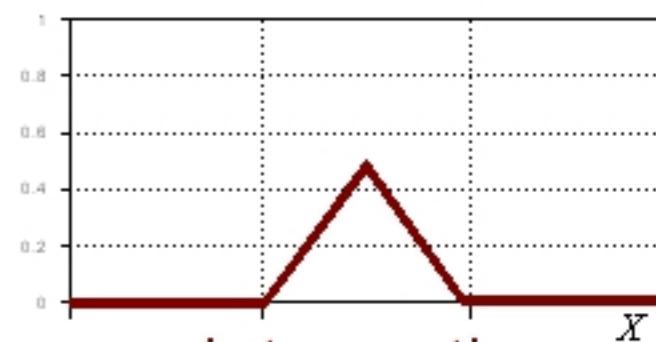


## Fuzziness

Degree of Membership/Truth



## Union



## Intersection

$$A \cap \bar{A} \neq \emptyset$$

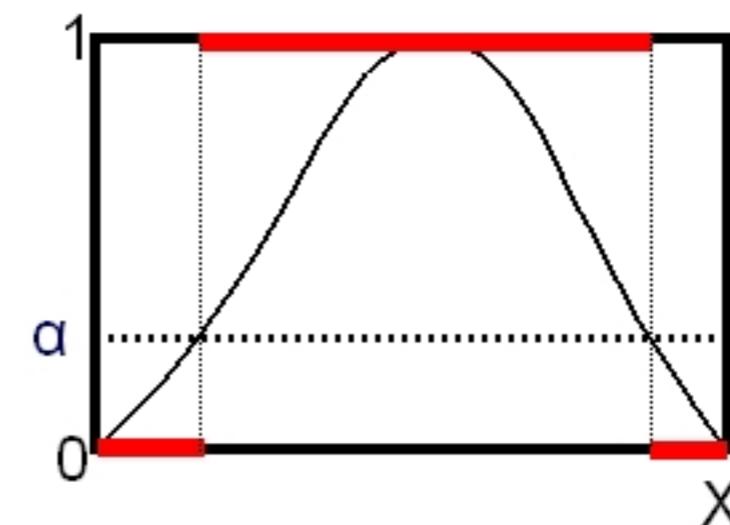
$$A \cup \bar{A} \neq X$$



$\alpha$ -cut

Fuzzy Sets:

$$A(x) : X \rightarrow [0, 1]$$



$\alpha$ -cut: Crisp  
set at  
threshold  $\alpha$

## De Morgan's Laws

$$\overline{A \cap B} = \overline{A} \cup \overline{B}$$

$$\overline{A \cup B} = \overline{A} \cap \overline{B}$$

$$i(a,b) = ab$$

$$i(a,b) = \frac{ab}{a+b-ab}$$

$$u(a,b) = a + b - ab$$

$$u(a,b) = \frac{a+b-2ab}{1-ab}$$

■ Complement, Intersection and Union must follow De Morgan's Laws plus:

► Complement

- Boundary Conditions:  $c(0)=1$  and  $c(1)=0$
- Monotonicity: if  $a \leq b$  then  $c(a) \geq c(b)$
- Continuity
- **Involutive:  $c(c(a)) = a$**

► Intersection (T-Norm)

- Boundary condition:  $i(a, 1) = a$
- Monotonicity: if  $b \leq d$  then  $i(a,b) \leq i(a,d)$
- Commutativity:  $i(a,b) = i(b,a)$
- Associativity:  $i(a,i(b,d)) = i(i(a,b),d)$
- Continuity
- Strict Monotonicity: if  $a_1 < a_2$  and  $b_1 < b_2$  then  $i(a_1,b_1) < i(a_2,b_2)$
- Subidempotency:  $i(a,a) \leq a$

► Union (T-Conorm)

- Boundary condition:  $u(a, 0) = a$
- Monotonicity: if  $b \leq d$  then  $u(a,b) \leq u(a,d)$
- Commutativity:  $u(a,b) = u(b,a)$
- Associativity:  $u(a,u(b,d)) = u(u(a,b),d)$
- Continuity
- Strict Monotonicity: if  $a_1 < a_2$  and  $b_1 < b_2$  then  $u(a_1,b_1) < u(a_2,b_2)$
- Superidempotency:  $u(a,a) \geq a$

$$c_\lambda(a) = \frac{1-a}{1+\lambda a}$$

Sugeno  
Complement:  
 $\lambda \in (-1, \infty)$



## properties

■ Reflexive

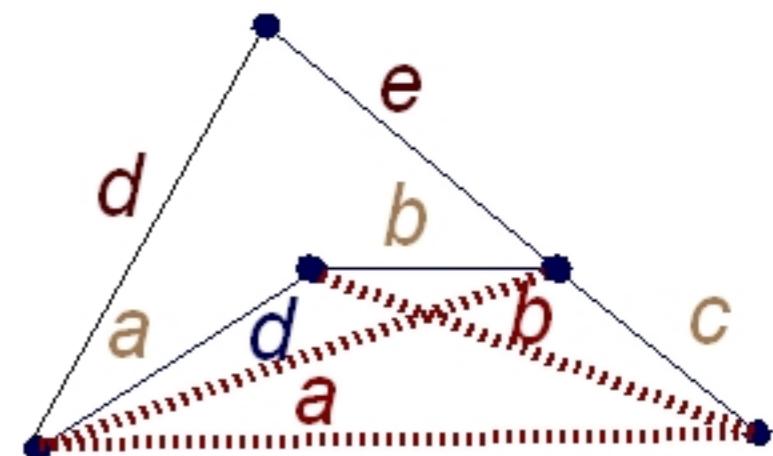
- ▶ iff  $R(x, x) = 1$  for all  $x \in X$ 
  - every element of  $X$  is maximally associated with itself

■ Symmetric

- ▶ iff  $R(x, y) = R(y, x)$  for all  $x, y \in X$ 
  - Matrices require only  $(n^2-n)/2$  elements to be defined

■ (Max-Min) Transitive

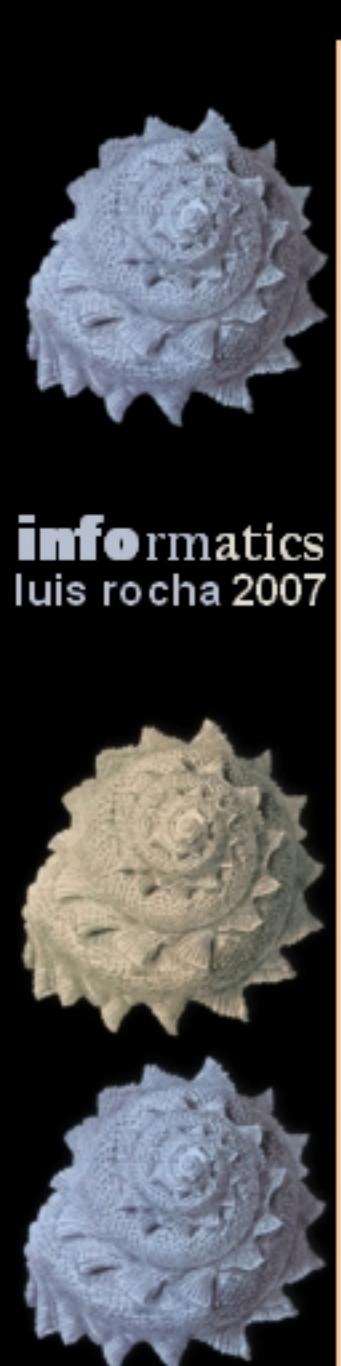
- ▶ iff  $R(x, z) \geq \max_{y \in X} \min[R(x, y), R(y, z)]$  for all  $x, z \in X$ 
  - For each indirect connection between  $x$  and  $z$  through some  $y$ , the weight of the connection is the smallest of each connection ( $x$  to  $y$  and  $y$  to  $z$ ). Finally, the weight of the connection between  $x$  and  $z$ , is the largest of all indirect connections through all  $y$  (strongest path defined by weakest link)



Max-Min Transitivity

$$a < b < c$$

$$a < d < e$$



**Max-Min Composition:**  $R \circ R = \max_k \min(r_{ik}, r_{kj}) = r'_{ij}$

where  $r_{ij}$  denotes  $R(x_i, x_j)$

The max-min composition of matrices is performed in the same way as the numerical counterpart, except that *multiplication* and *summation* are substituted by the  $\cap$  (e.g. Min) and  $\cup$  (e.g. Max) operations respectively.

■ Transitive closure of a relation  $R(X, X)$

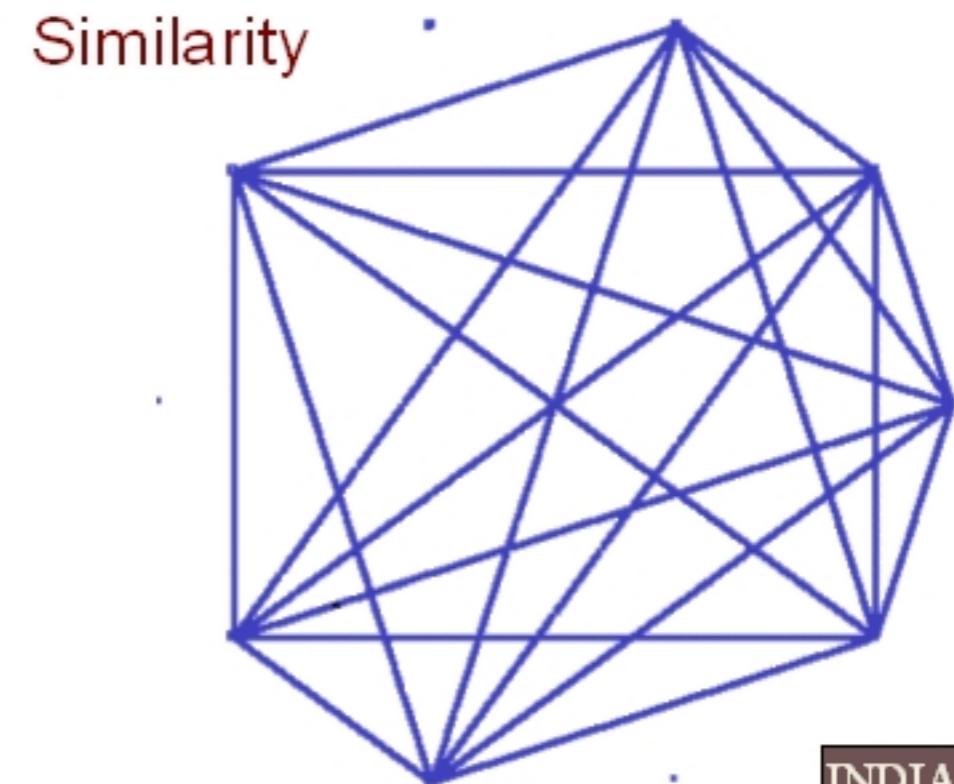
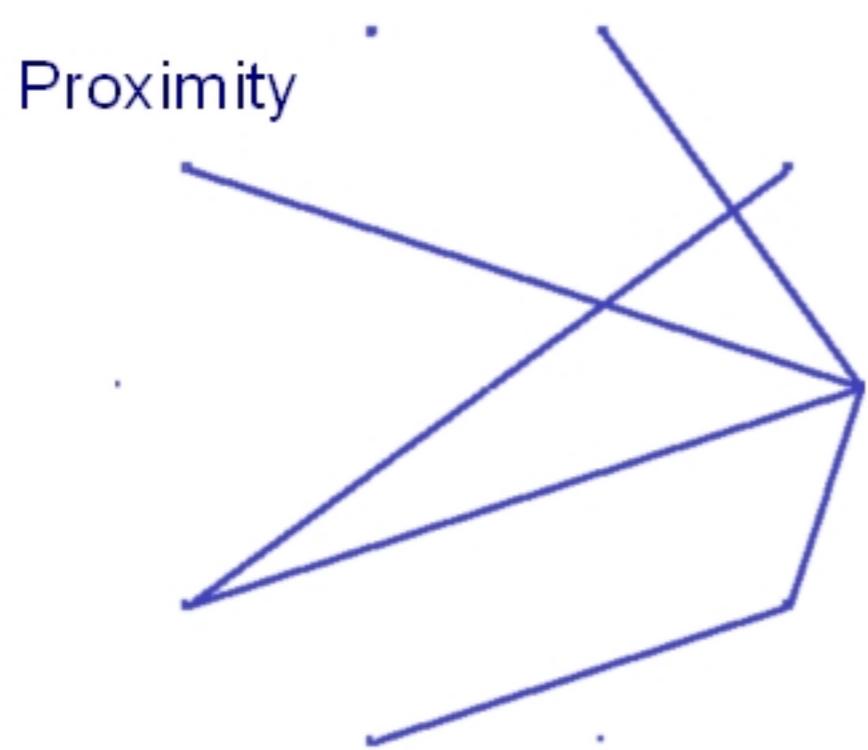
- ▶ The relation that is transitive, contains  $R(X, X)$ , and whose elements have the smallest possible membership weights that still allow the first two requirements.
  - It yields a relation where all pairs of elements which were directly or indirectly related in the original relation, are now directly related

**Generic Composition:**

$$R \circ R = \bigcup_k \cap(r_{ik}, r_{kj}) = r'_{ij}$$



- **Similarity Relation**
  - ▶ A reflexive, symmetric, and transitive binary fuzzy relation
    - Also known as an equivalence relation.
- **Proximity Relation**
  - ▶ A reflexive and symmetric binary fuzzy relation
    - Also known as a compatibility relation
    - The transitive closure of a proximity relation is a similarity relation.



## from document relations

- Document × Keyterms
  - ▶ Keyterm Co-Occurrence
- Document × Document
  - ▶ Co-Citation or Hyperlink structure
- Document × Author
  - ▶ Co-Authorship (Collaboration Network)
- Bio-entities × Keyterms
  - ▶ Gene/MeSH keyterm Co-Occurrence

$X$ (Keywords)
$R:X \times Y$
$Y$ (Documents)

Given a binary relation  $R$  between sets  $X$  and  $Y$  we extract two proximity relations:  $XYP(x_i, x_j)$  is the probability that both  $x_i$  and  $x_j$  are related in  $R$  to the same element  $y \in Y$ . Conversely,  $YXP(y_i, y_j)$  is the probability that both  $y_i$  and  $y_j$  are related in  $R$  to the same element  $x \in X$ .

$$XYP(x_i, x_j) = \frac{\sum_{k=1}^m (r_{i,k} \wedge r_{j,k})}{\sum_{k=1}^m (r_{i,k} \vee r_{j,k})}; \quad YXP(y_i, y_j) = \frac{\sum_{k=1}^n (r_{k,i} \wedge r_{k,j})}{\sum_{k=1}^n (r_{k,i} \vee r_{k,j})}$$

With some support constraint



**informatics**  
luis rocha 2007



## proximity measures

produce associative (probabilistic) networks

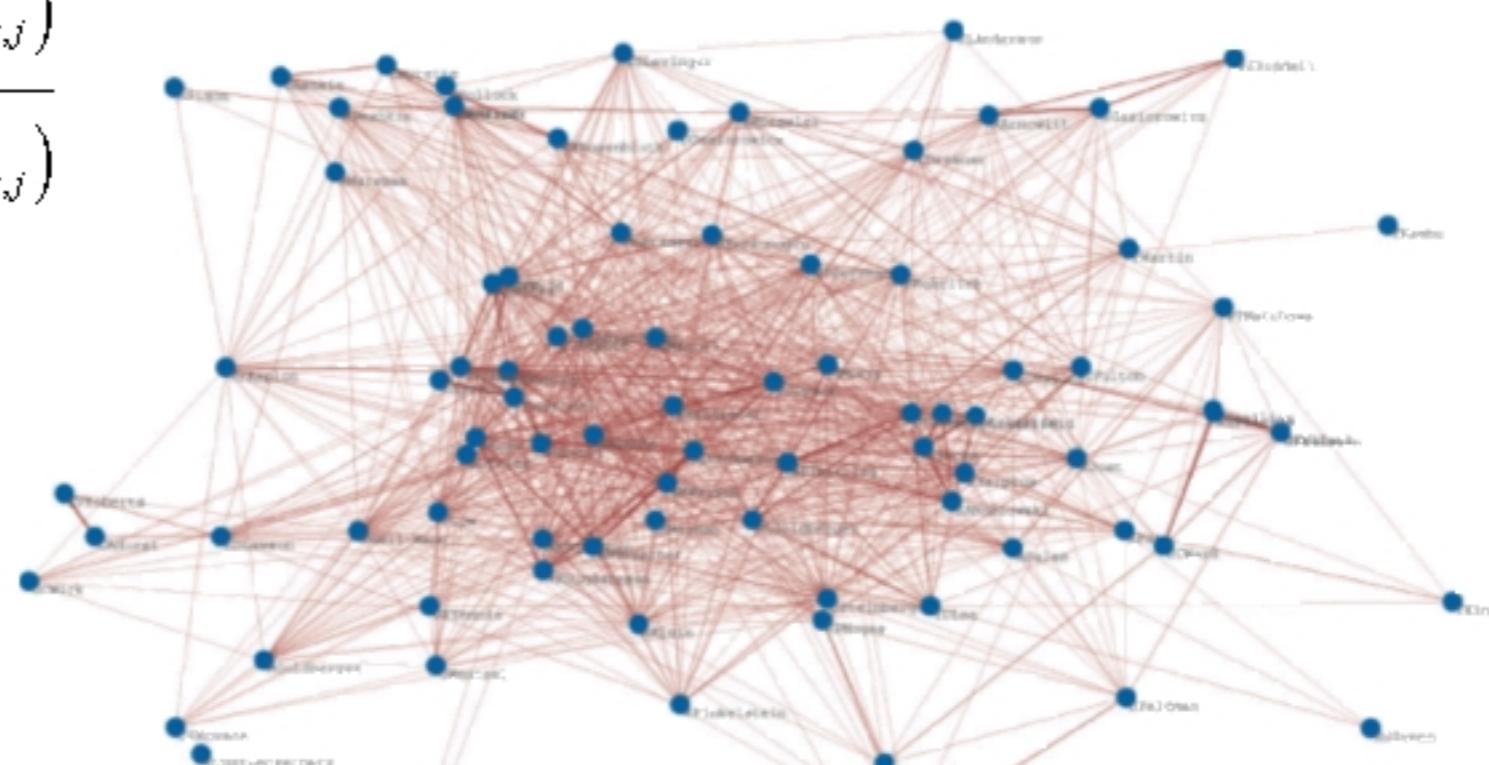
	$X$ (Keywords)
$Y$ (Documents)	$R:X \times Y$

$$XYP(x_i, x_j) = \frac{\sum_{k=1}^m (r_{i,k} \wedge r_{j,k})}{\sum_{k=1}^m (r_{i,k} \vee r_{j,k})}$$

$X$ (Keywords)	$X$ (Keywords)
$XYP:X \times X$	• Spiral classifier

$$YXP(y_i, y_j) = \frac{\sum_{k=1}^n (r_{kj} \wedge r_{kj})}{\sum_{k=1}^n (r_{kj} \vee r_{kj})}$$

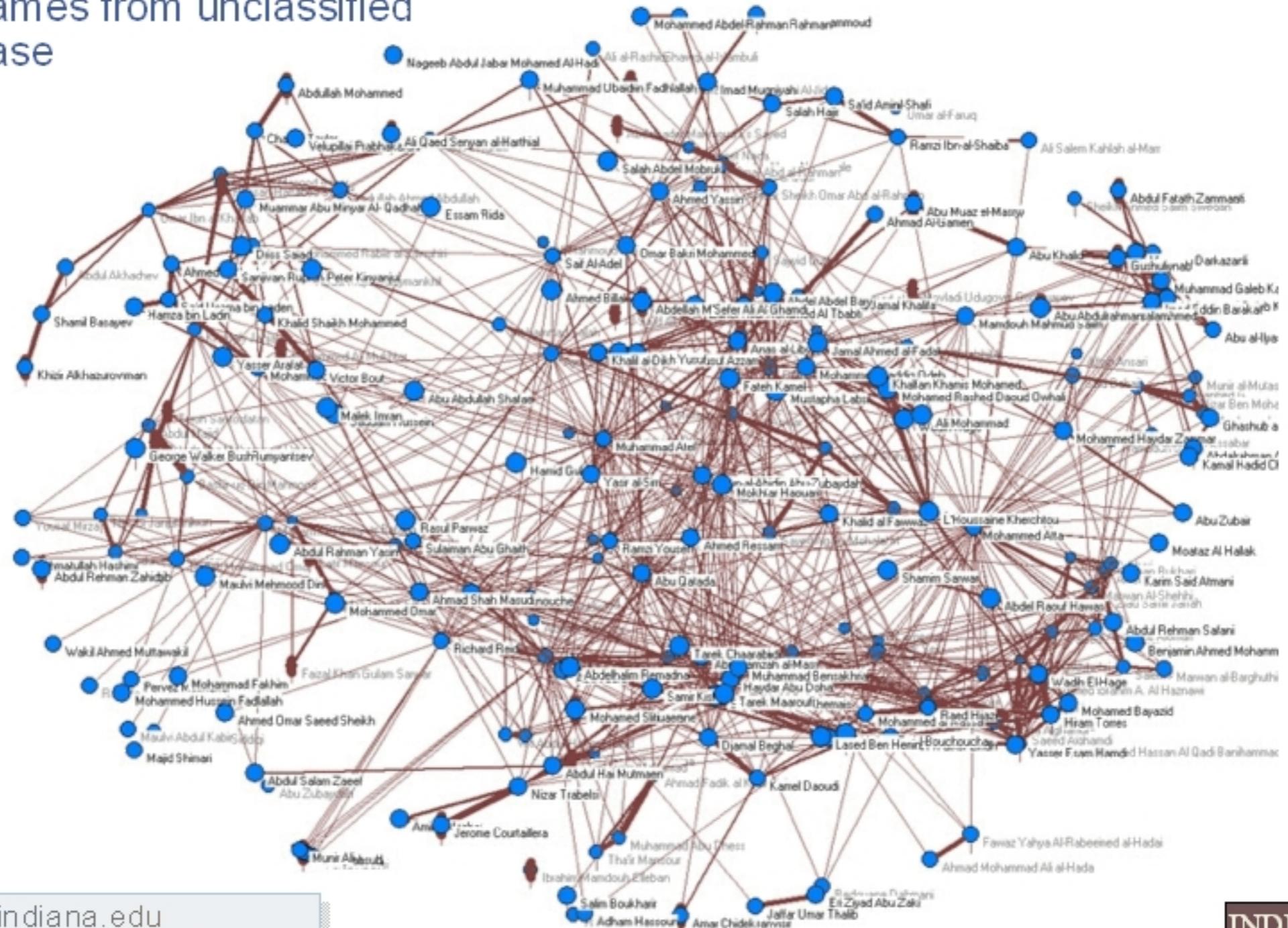
	$Y$ (Documents)
$Y$ (Documents)	$YXP:Y \times Y$



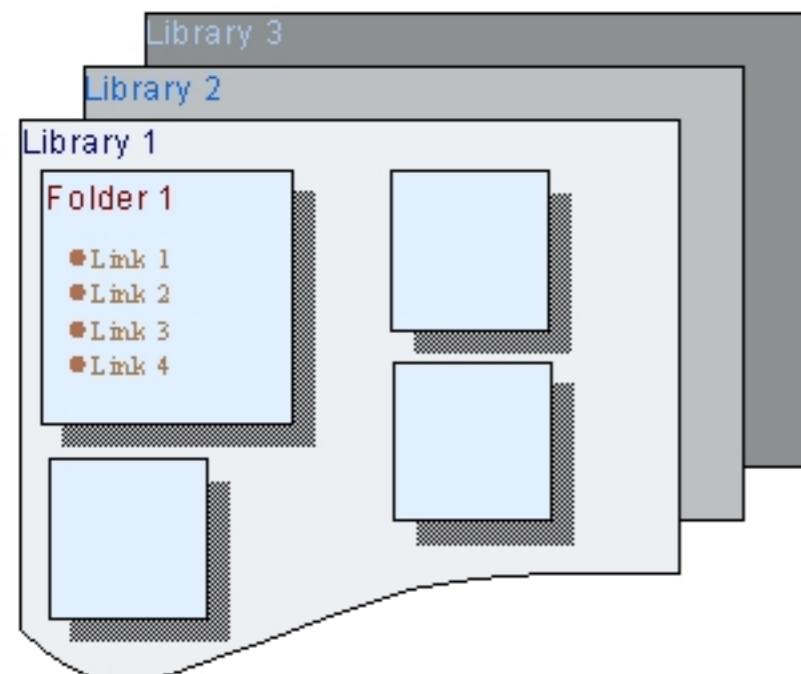
represents knowledge in  
associative manner

## PDP2

318 names from unclassified  
database



# digital library web service



## ■ Three nested entities:

### ► Libraries $\Rightarrow$ Folders $\Rightarrow$ Links

- A *library/personality* is associated with a given area of interest and consists of one or more folders.
- A *folder* contains related types of links within a library
- A *link* is a URL (typically scientific articles)

rocha@indiana.edu  
<http://informatics.indiana.edu/rocha>

Rocha, L.M., T. Simas, A. Rechtsteiner, M. DiGiacomo, R. Luce [2005]. "MyLibrary@LANL: Proximity and Semi-metric Networks for a Collaborative and Recommender Web Service". In: Proc. 2005 IEEE/WIC/ACM International Conference on Web Intelligence (WI05), IEEE Press, pp. 565-571.

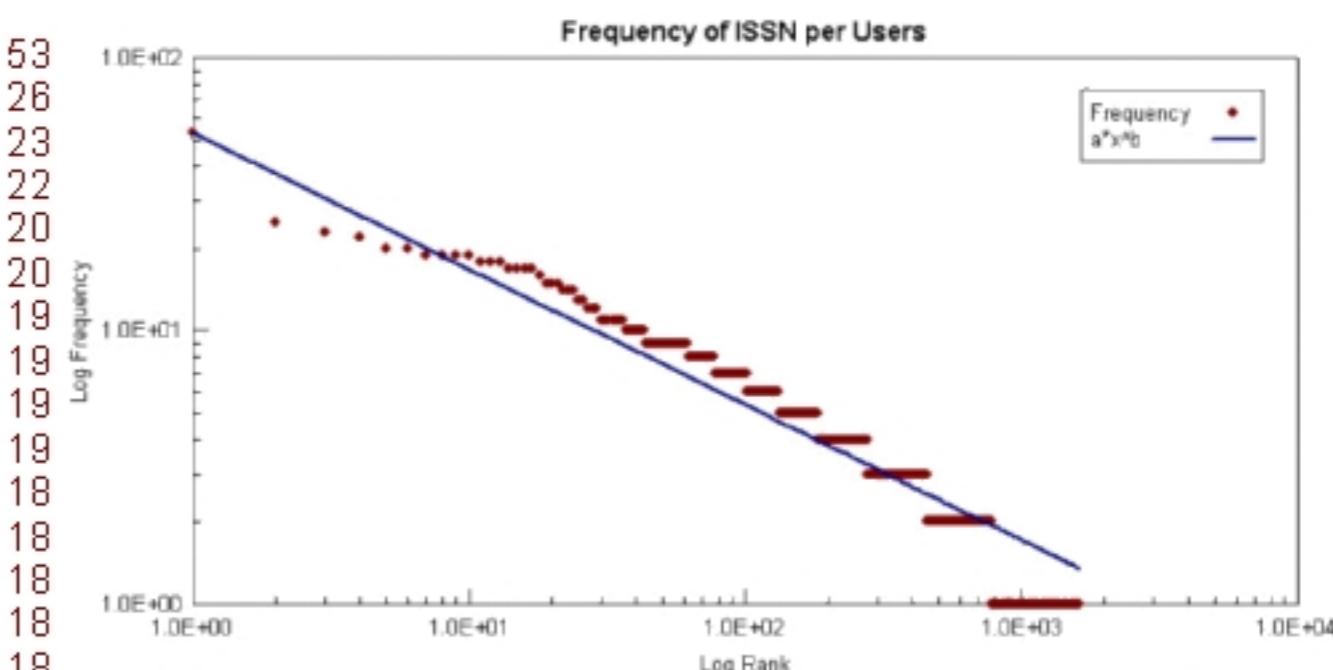
The screenshot shows a web page with a navigation bar at the top. The navigation bar includes categories like "Ethics", "Evolutionary Systems", "Fuzzy Math", "Immunology.shared", "Mathematics", and "Other". Below the navigation bar is a message box: "Messages from the library: Search Who's Who for biographical data on trial until Sept. 1 - send feedback to [steam@lanl.gov](mailto:steam@lanl.gov)". The main content area contains several sections with links:
 

- Distance Functions**: Distance Functions and Topologies, Euclidian space and grouping of biological objects, DISTANCE FUNCTIONS AND TOPOLOGIES, EXPLICIT METRIZATION
- Mathematics Web Resources**: Hyperstat Online
- Databases**: Current index to statistics, FlashPoint, INSPEC® at LANL, Jahrbuch über die Fortschritte der Mathematik, MathSciNet, SciSearch® at LANL, Zentralblatt MATH database
- Graph Theory**: A duplication growth model of gene expression networks, A Stochastic Model for the Evolution of the Web, Accelerated growth of networks, Curvature of co-links uncovers hidden thematic layers in the World Wide Web, Dynamical small-world behavior in an epidemical model of mobile individuals, Friends and Neighbors on the Web, Graph structure in the web, Intentional Walks on Scale Free Small Worlds, Intentional Walks on Scale Free Small Worlds, Local Search in Unstructured Networks, Modeling the Internet's large-scale topology, Models of the Small World: A Review, Optimization in complex networks, Random graph models of social networks, The Probability of Collective Choice with Shared Knowledge Structures, The structure and function of complex

# most frequent ISSN

## occurrences in user personalities

Physical review letters	53
Physical Review B	26
Physical review E	23
Physical review A General physics	22
Journal of physical chemistry B	22
Computers & geosciences	20
Scientific American	20
Journal of the American Chemical Society	19
Journal of Chemical Physics	19
Reviews of modern physics	19
Bioinformatics	19
IEEE trans. on geoscience and remote sensing	18
PNAS	18
Journal of computational physics	18
Advances in water resources	18
Journal of applied geophysics	17
Applied geochemistry	17
APL	17
Journal of physical chemistry A	17
Phil. mag. B Physics of condensed matter	17
Bul. of Environmental Contamination and Toxicology	16
Journal of applied physics	15
American journal of physics	14
Analytical chemistry	14
DLib	14
Chemical physics	13
NIM	13
Chemical physics letters	12
Physics reports	12
Physical review A	12



Nature	11
Physics of plasmas	11
Accounts of chemical research	11
Advanced Materials	11
Journal of physics Condensed matter	11
Journal of computational biology	11
Biochemical journal	11
Computer physics communications	11
ACM transactions on modeling and computer simulation	10
Materials science & engineering A	10
Physical review C Nuclear physics	10
Advances in physics	10
Inorganic chemistry	10
Trends in BioTechnology	10
Review of Scientific Instruments	10
Science	10

from co-occurrence in mylibrary.lanl.gov

	ISSN
Personality	$A:P \times I$

392 personalities with at least two ISSN  
 253 users with at least two ISSN  
 1702 unique ISSN occurring at least twice

Given a binary relation  $A$  between sets of Personalities  $P$  and ISSN  $I$  we extract two proximity relations:  $PIP(p_s, p_t)$  is the probability that both personalities  $p_s$  and  $p_t$  link to the same ISSN  $i \in I$ . Conversely,  $IPP(i_s, i_t)$  is the probability that both ISSN  $i_s$  and  $i_t$  co-occur in the same personality (given that one of them occurs)  $p \in P$ .

$$PIP(p_s, p_t) = \frac{\sum_{k=1}^m (a_{i,k} \wedge a_{j,k})}{\sum_{k=1}^m (a_{i,k} \vee a_{j,k})} = \frac{N_{\cap}(p_s, p_t)}{N_{\cup}(p_s, p_t)}$$

(Personality ISSN Proximity)

$$IPP(i_s, i_t) = \frac{\sum_{k=1}^m (a_{i,k} \wedge a_{j,k})}{\sum_{k=1}^m (a_{i,k} \vee a_{j,k})} = \frac{N_{\cap}(i_s, i_t)}{N_{\cup}(i_s, i_t)}$$

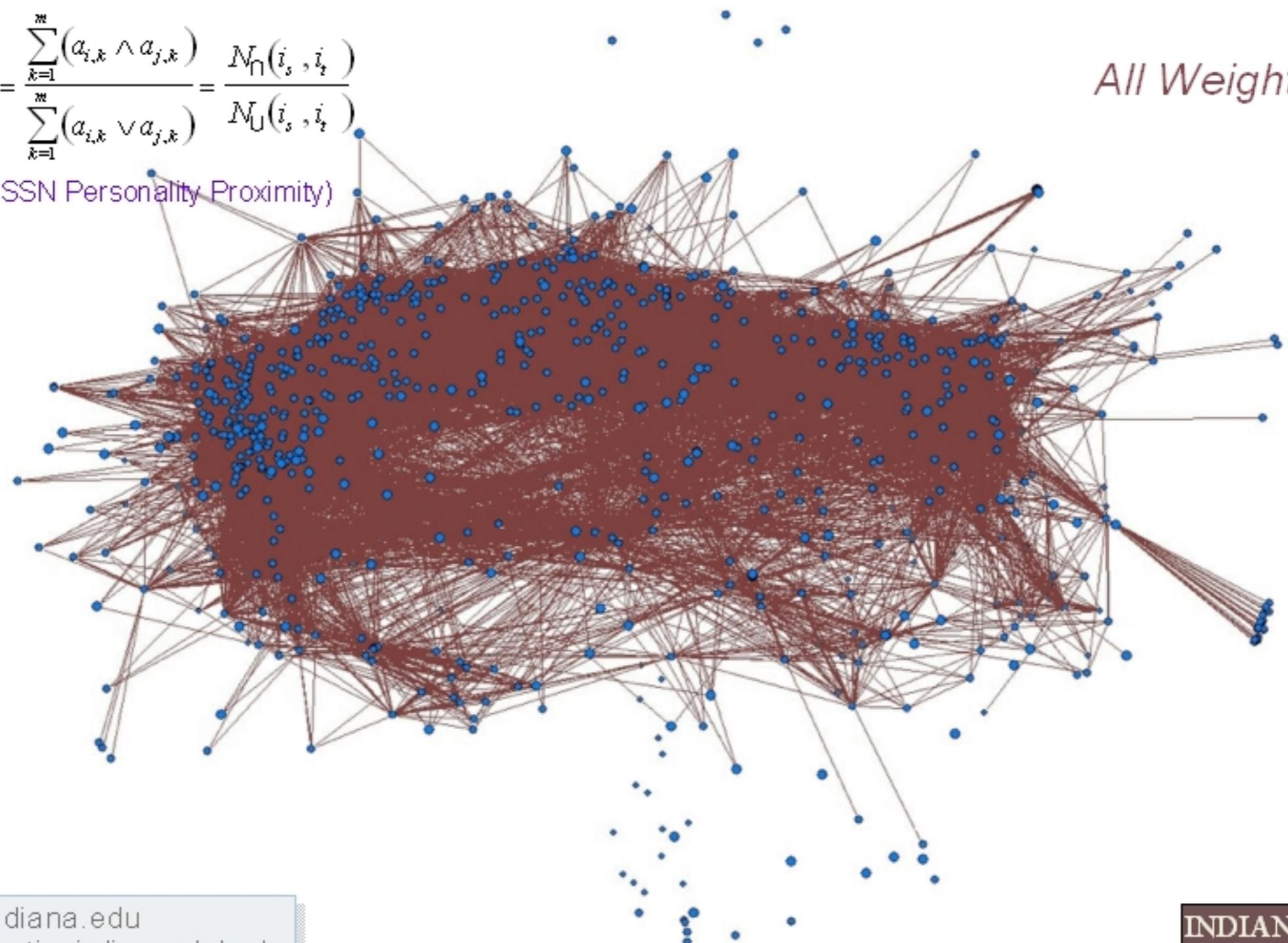
(ISSN Personality Proximity)

from co-occurrence in user personalities in mylibrary.lanl.gov: IPP

$$ipp(i_s, i_t) = \frac{\sum_{k=1}^m (a_{i,k} \wedge a_{j,k})}{\sum_{k=1}^m (a_{i,k} \vee a_{j,k})} = \frac{N_{\cap}(i_s, i_t)}{N_{\cup}(i_s, i_t)}$$

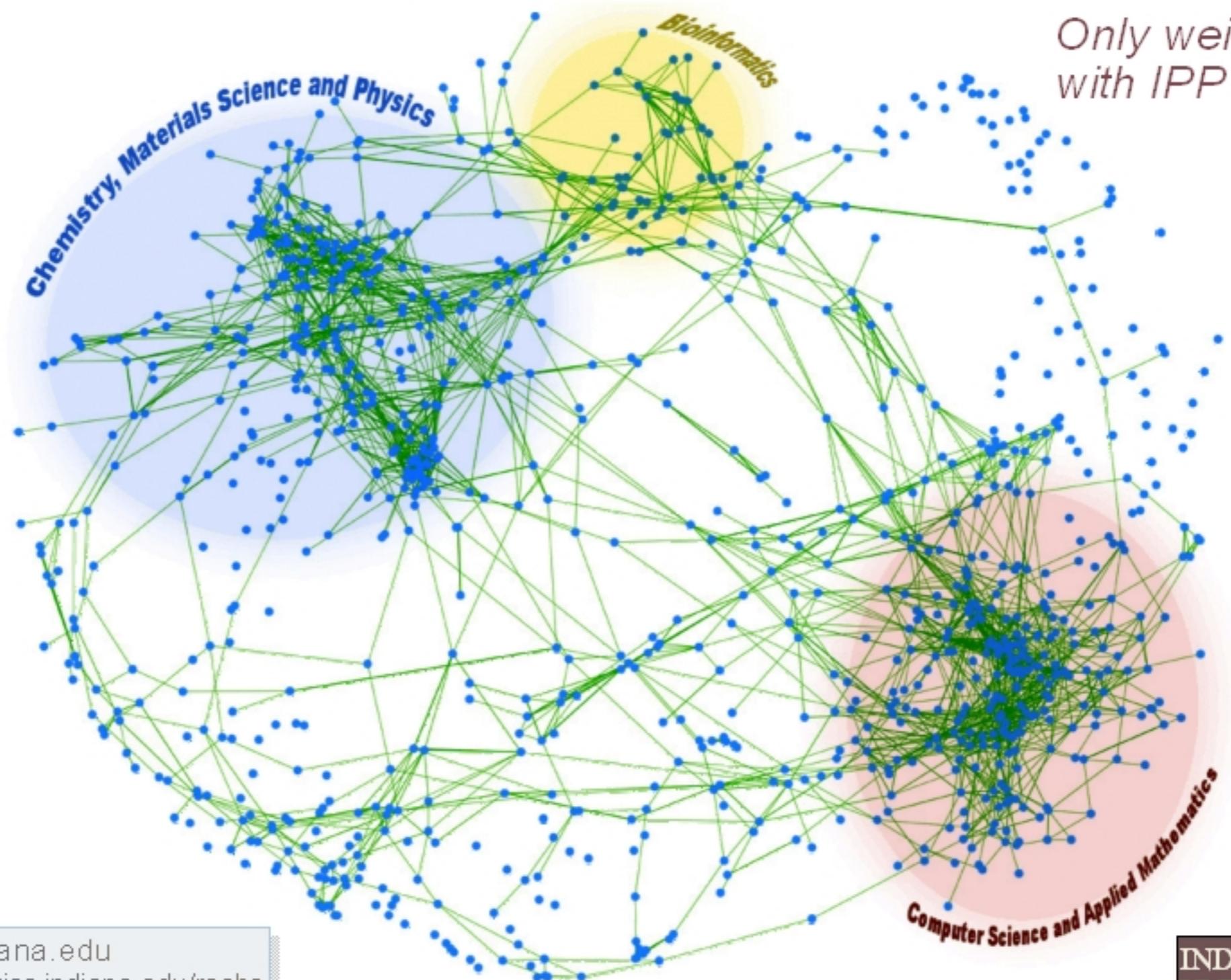
(ISSN Personality Proximity)

All Weights



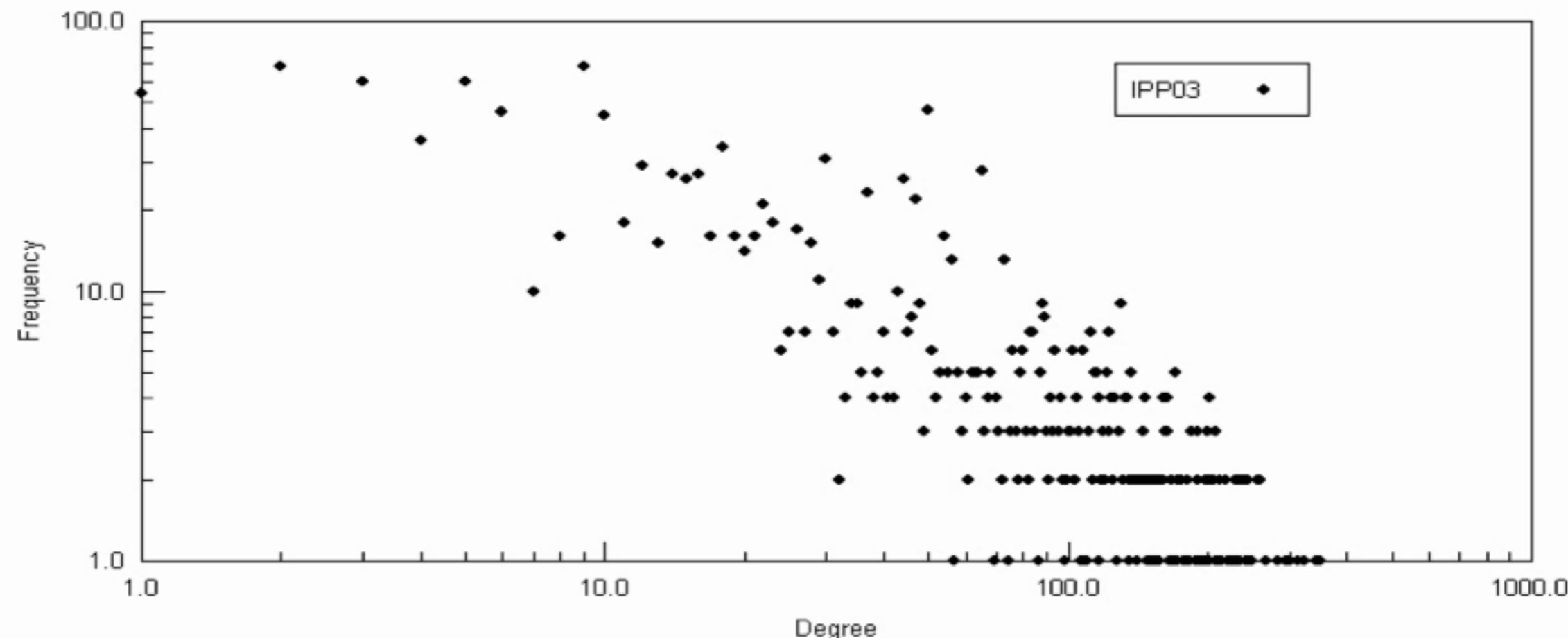
# IPP (journal) network

from co-occurrence in user personalities



# cumulative degree distribution

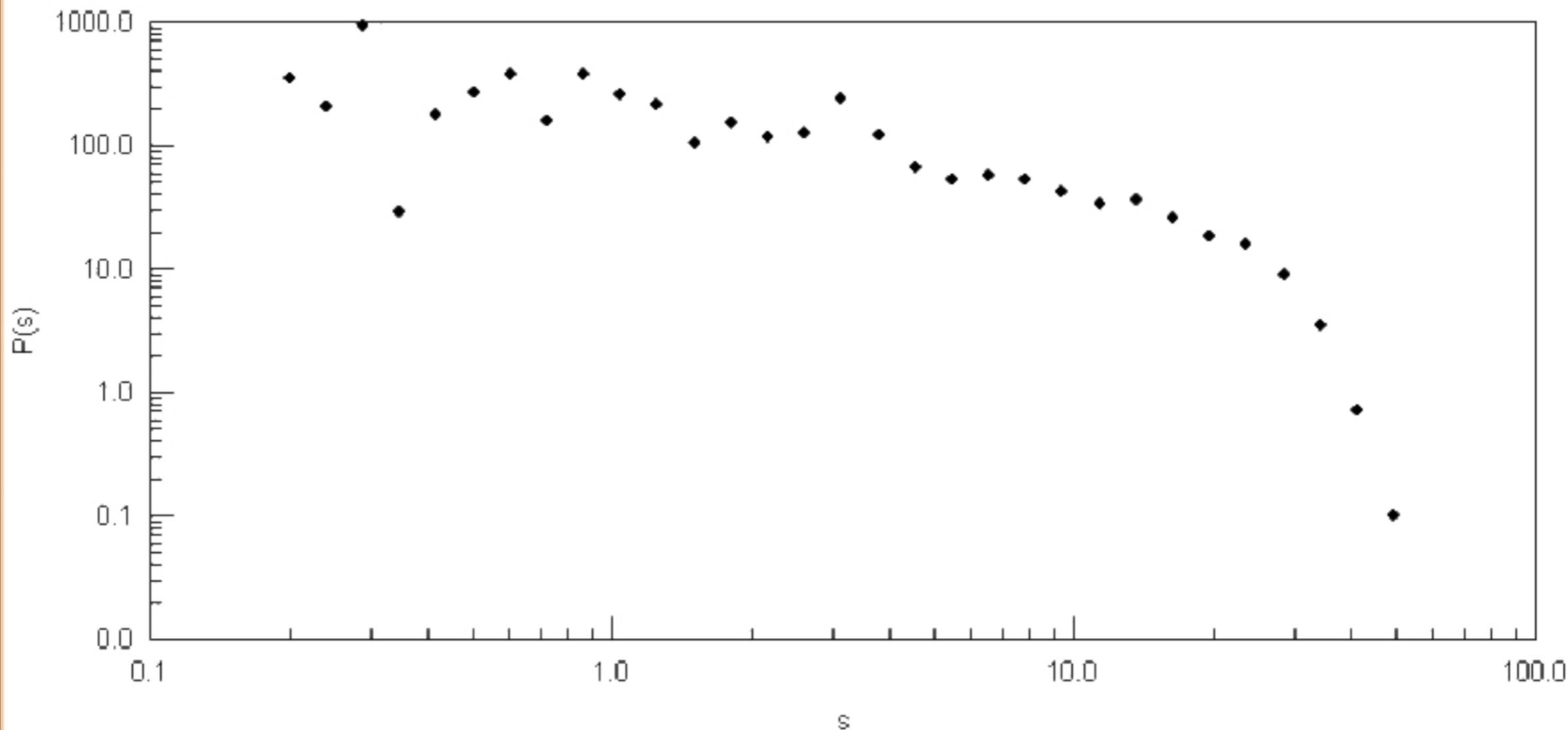
IPP: all weights





## cumulative strength distribution (binned)

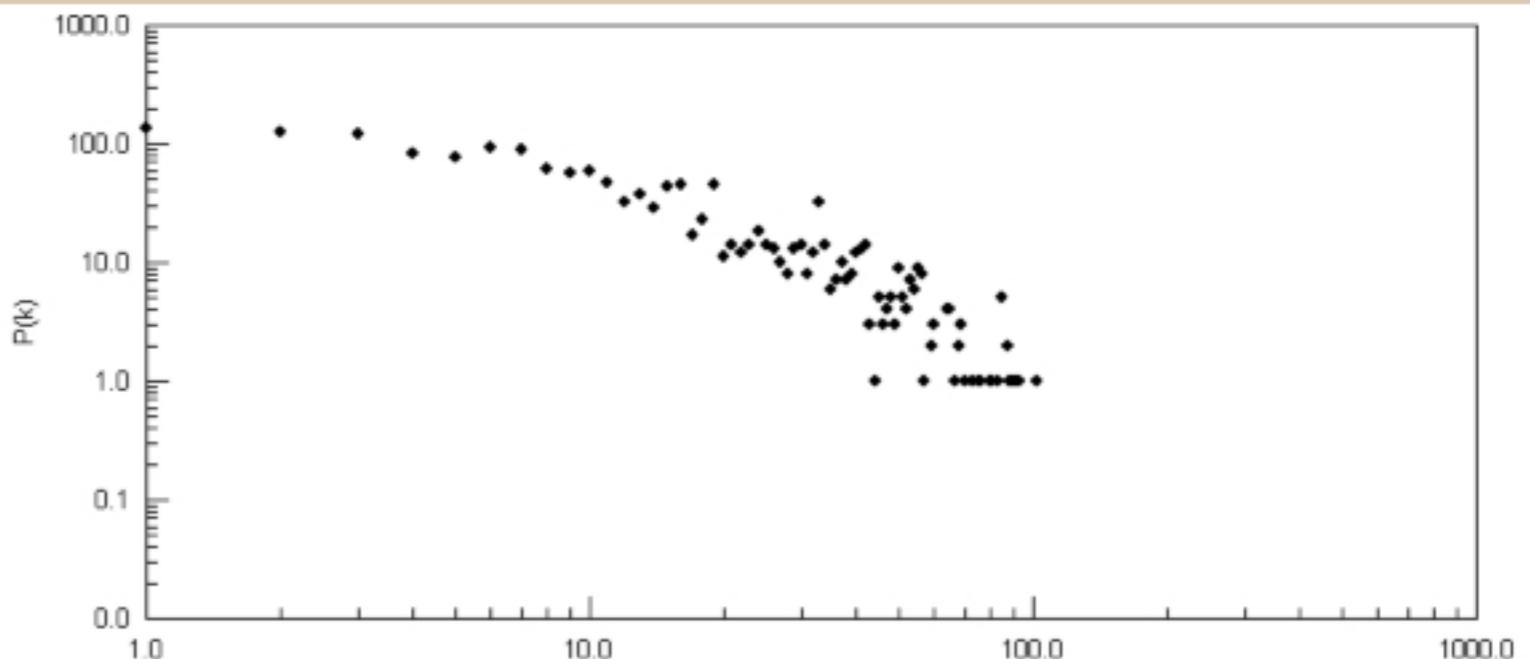
IPP: all weights



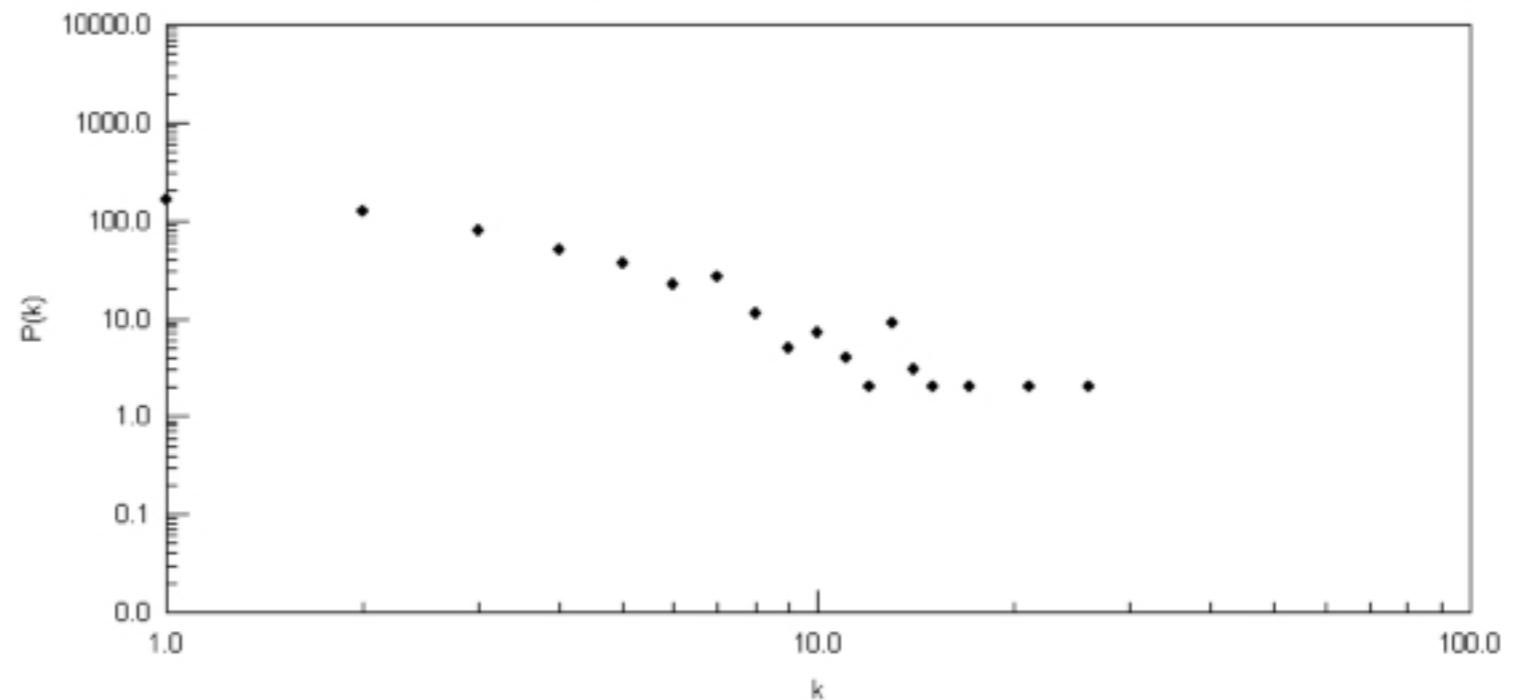
# cumulative degree distribution

IPP

$\alpha$ -cut = 0.2



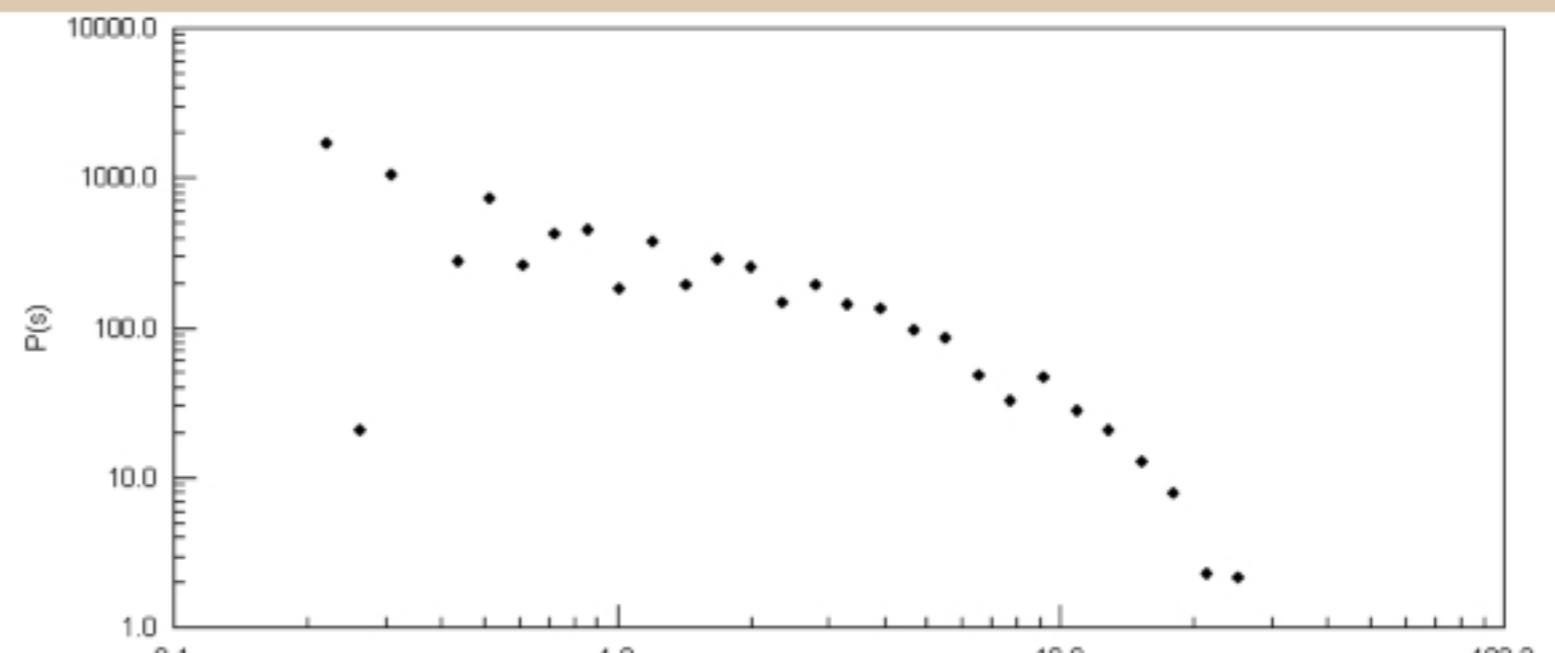
$\alpha$ -cut = 0.4



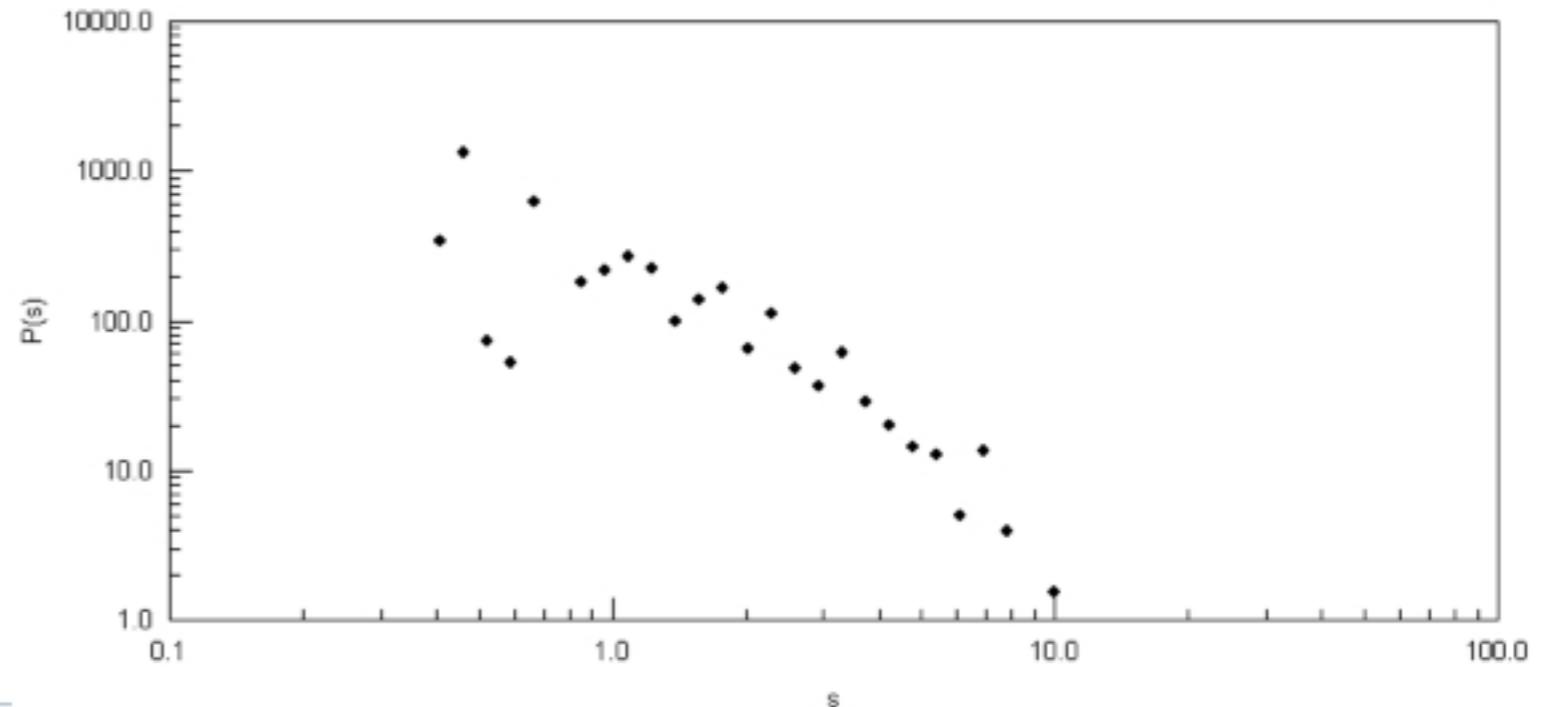
# cumulative strength distribution (binned)

IPP

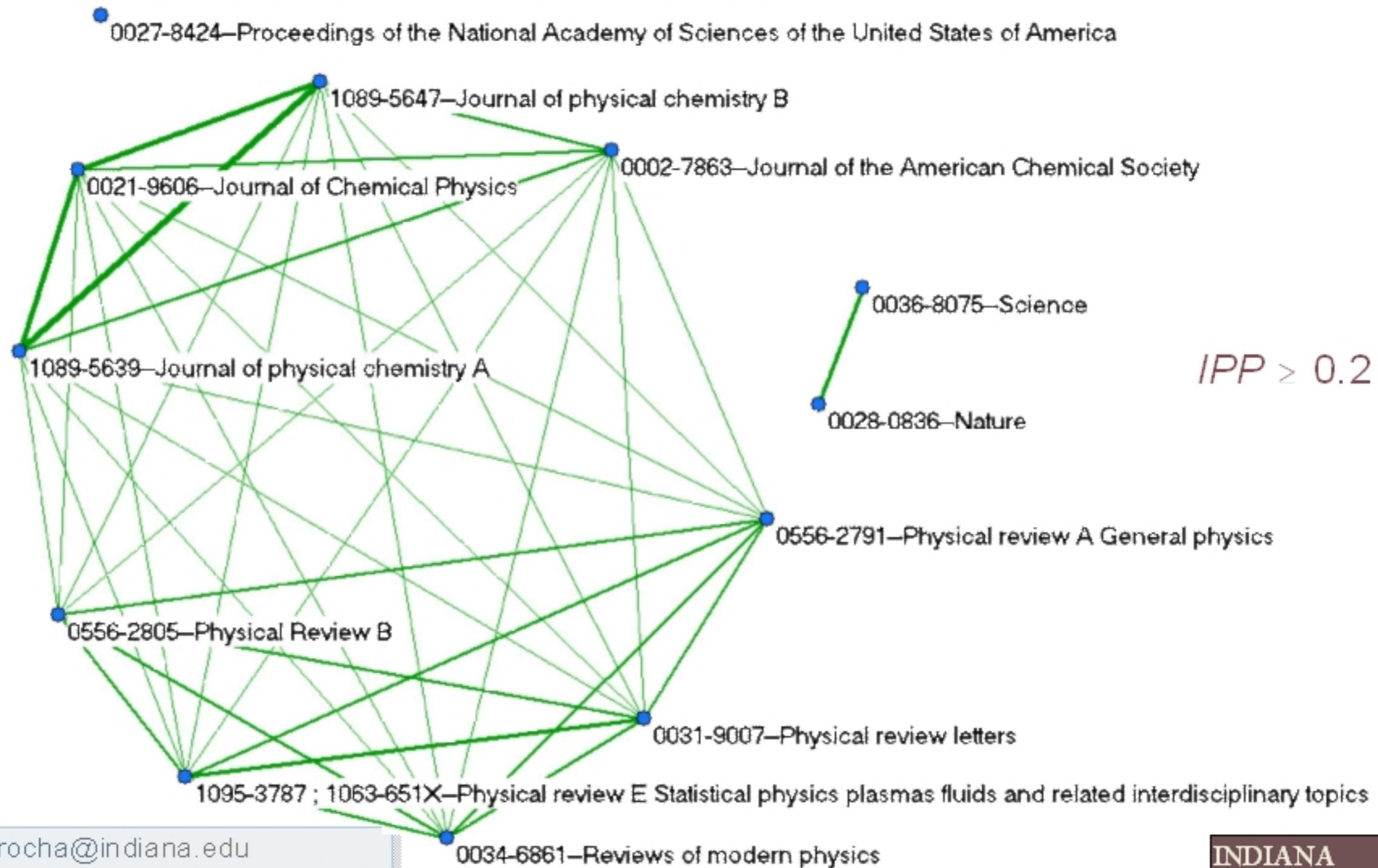
$\alpha$ -cut = 0.2



$\alpha$ -cut = 0.4



## sub-graph with top 12 most frequent journals



## integrating proximity edges

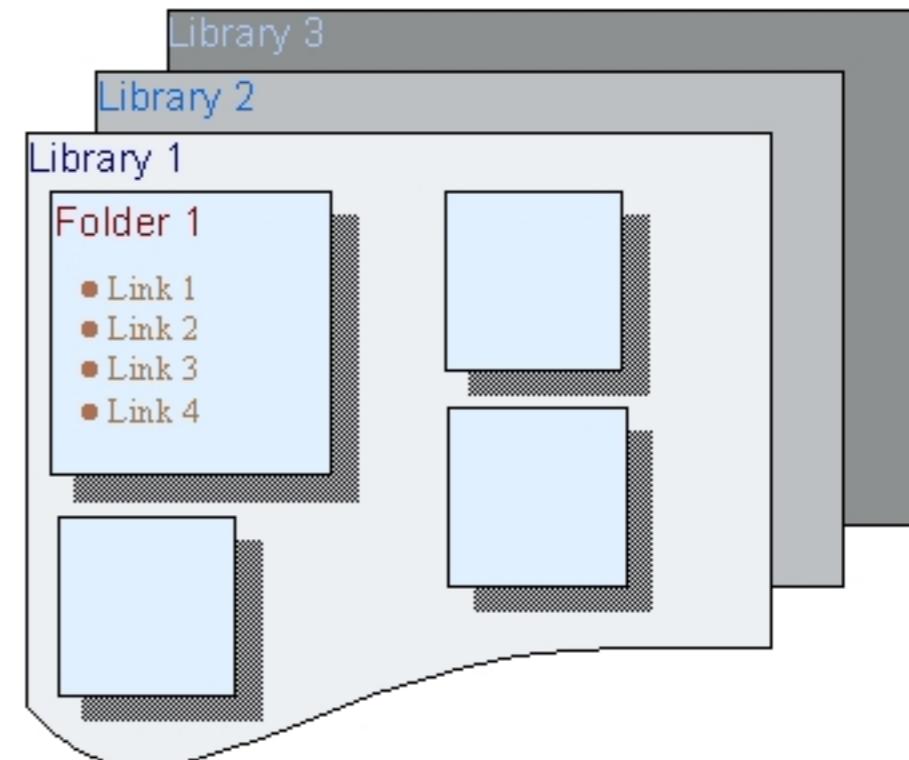
$$I_R = \left\{ i_t : \underset{\forall i_s \in I(p_u)}{MAX} (IPP(i_s, i_t)) \geq \alpha \right\}$$

Returns all journals in the network related to at least one of the journals in personality with proximity  $\geq \alpha$

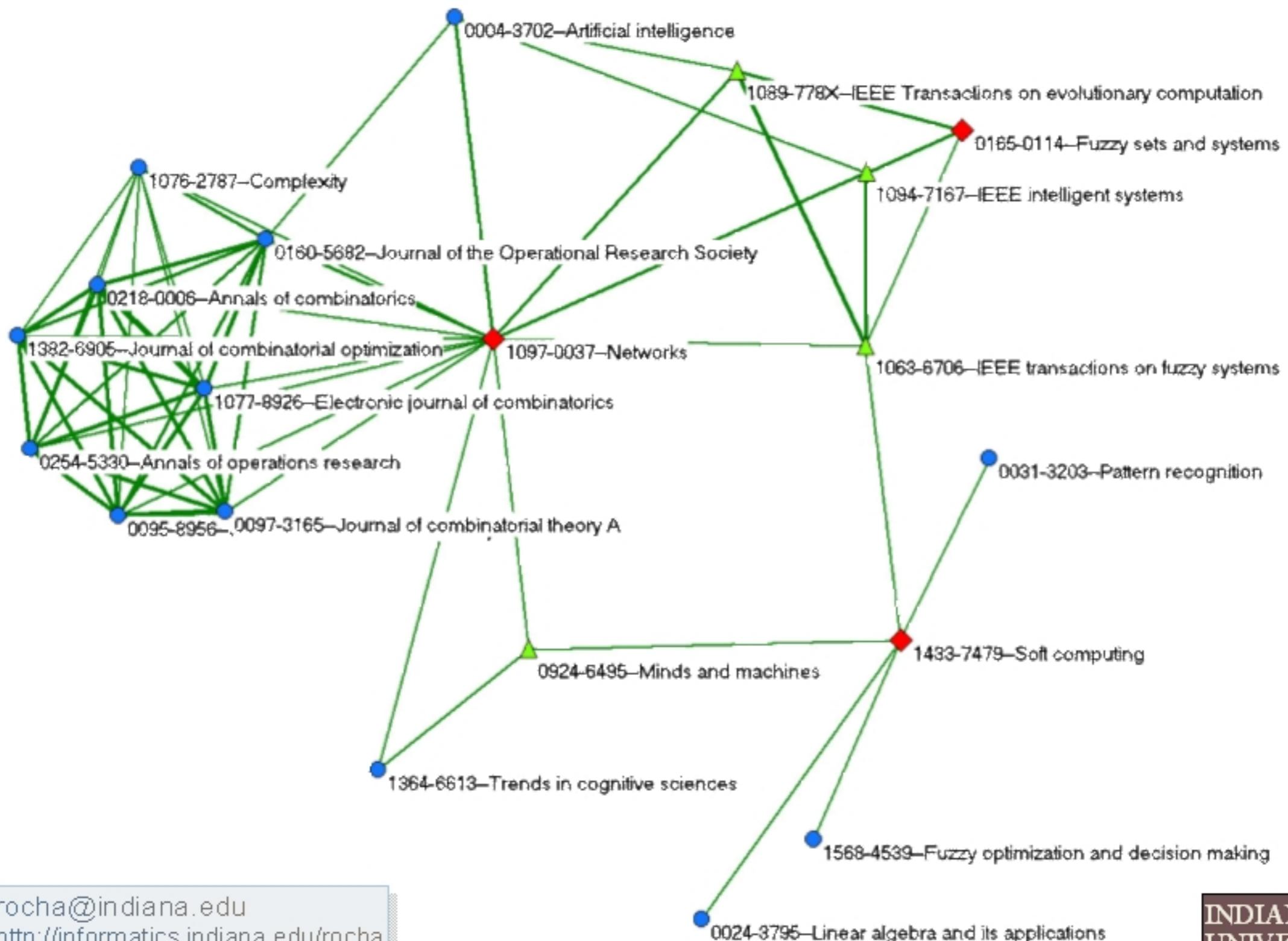
$$I_R = \left\{ i_t : \underset{\forall i_s \in I(p_u)}{MIN} (IPP(i_s, i_t)) \geq \alpha \right\}$$

Returns all journals in the network related to every journal in personality with proximity at least  $\geq \alpha$

$$I_R = \left\{ i_t : \underset{\forall i_s \in I(p_u)}{AVG} (IPP(i_s, i_t)) \geq \alpha \right\}$$



# recommendation example



mylibrary.lanl.gov

informatics  
luis rocha 2007

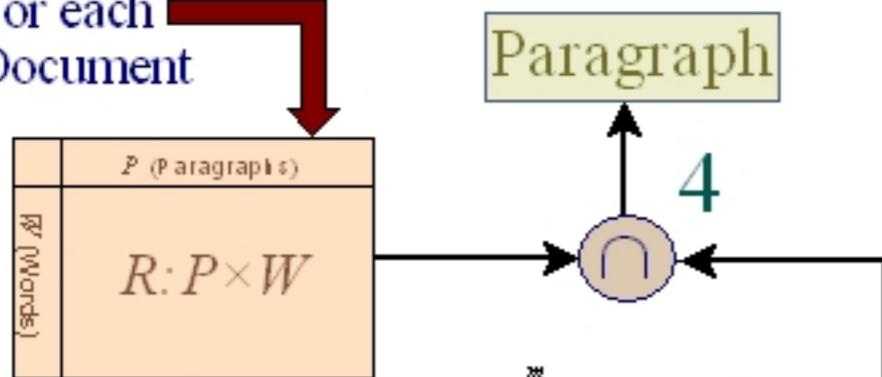
- IPP
  - ▶ Recommendations of ISSN based on co-occurrence in Personalities
    - Users who linked to this journal, also linked to...
- PIP
  - ▶ Recommendations of other users' personalities: collaboration
    - These personalities are similar to yours
  - ▶ Recommendations of specific links in close personalities
    - Users who read many of the same journals where interested in these links

# Biocreative competition (EMBO Workshop)

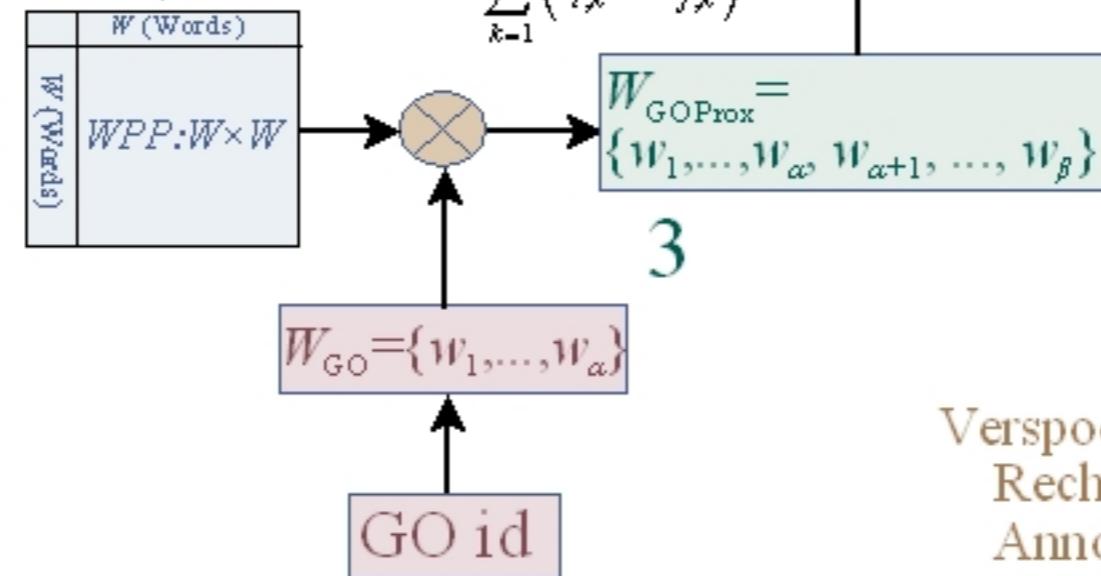
## a critical assessment of text mining methods in molecular biology

For each Document

1



2



- Task 2: Given a document, discover the portion of text most appropriate to annotate the protein's function, and produce appropriate Gene Ontology node for annotation
  - Learning set: triplets (protein, document, GO id)
  - Test set: documents

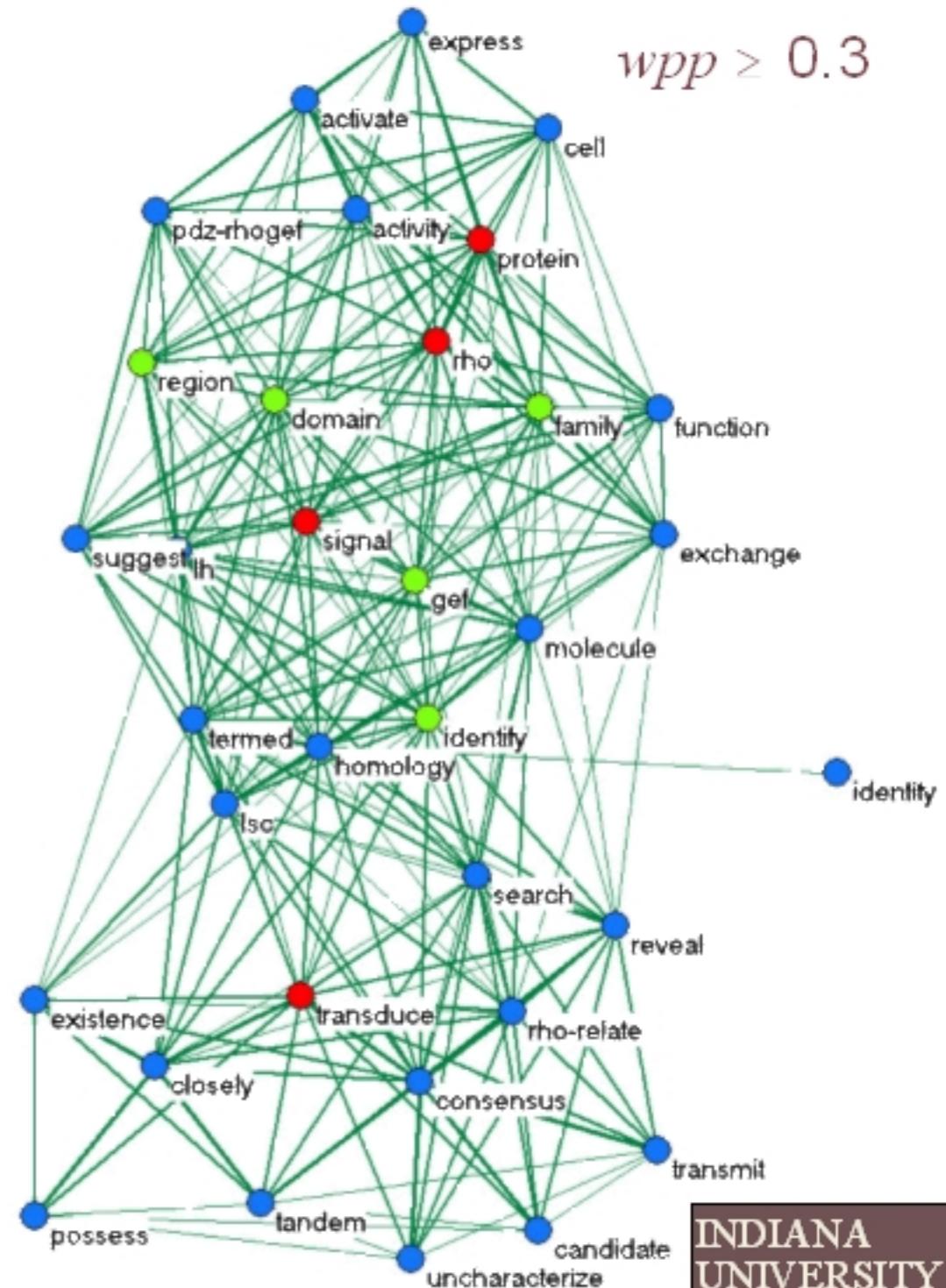
Verspoor, K., J. Cohn, C. Joslyn, S. Mniszewski, A. Rechtsteiner, L.M. Rocha, T. Simas [2005]. "Protein Annotation as Term Categorization in the Gene Ontology using Word Proximity Networks". *BMC Bioinformatics*, 6(Suppl 1):S20. doi:10.1186/1471-2105-6-S1-S20



## example document

- document bc005868
  - ▶ WPP contains 1102 words
  - ▶ Subgraph of 34 words
    - Red nodes: words removed from the respective GO annotation (0007266): Rho, protein, signal, transducer).
    - Blue nodes: words that co-occur very frequently ( $wpp > 0.5$ ) with at least one of the red nodes
    - Green nodes: additional words recommended with largest average proximity to all input words (red nodes)

Verspoor, K., J. Cohn, C. Joslyn, S. Mniszewski, A. Rechtsteiner, L.M. Rocha, T. Simas [2005]. "Protein Annotation as Term Categorization in the Gene Ontology using Word Proximity Networks". *BMC Bioinformatics*, 6(Suppl 1):S20. doi:10.1186/1471-2105-6-S1-S20



## Task 2.1 Results



Proximity-based run

User, Run	"perfect"	"generally"	cumulative
7, 1	<b>25.28%</b>	<b>14.31%</b>	<b>39.59%</b>
14, 1	<b>28.16%</b>	<b>6.41%</b>	<b>34.57%</b>
20, 1	<b>27.97%</b>	<b>5.30%</b>	<b>33.27%</b>
4, 1	<b>24.91%</b>	<b>6.88%</b>	<b>31.78%</b>
20, 2	<b>26.02%</b>	<b>5.58%</b>	<b>31.60%</b>
20, 3	<b>22.21%</b>	<b>5.48%</b>	<b>27.79%</b>
5, 2	<b>15.43%</b>	<b>8.36%</b>	<b>23.79%</b>
5, 1	<b>15.43%</b>	<b>7.16%</b>	<b>22.58%</b>
5, 3	<b>14.31%</b>	<b>7.99%</b>	<b>22.39%</b>
15, 2	<b>11.62%</b>	<b>6.41%</b>	<b>18.93%</b>
9, 1	<b>11.62%</b>	<b>1.21%</b>	<b>12.83%</b>
7, 3	<b>6.13%</b>	<b>3.72%</b>	<b>9.85%</b>
17, 1	<b>7.71%</b>	<b>1.77%</b>	<b>9.48%</b>
15, 1	<b>5.48%</b>	<b>2.60%</b>	<b>8.09%</b>
7, 2	<b>4.99%</b>	<b>3.72%</b>	<b>7.71%</b>
10, 3	<b>4.65%</b>	<b>0.37%</b>	<b>5.02%</b>
9, 3	<b>3.81%</b>	<b>0.65%</b>	<b>4.46%</b>
10, 2	<b>4.18%</b>	<b>0.19%</b>	<b>4.37%</b>
10, 1	<b>3.35%</b>	<b>0.28%</b>	<b>3.62%</b>
9, 2	<b>3.07%</b>	<b>0.46%</b>	<b>3.53%</b>
17, 2	<b>0.65%</b>	<b>0.00%</b>	<b>0.65%</b>

Verspoor, K., et al [2005]. *BMC Bioinformatics*, 6(Suppl 1):S20.  
doi:10.1186/1471-2105-6-S1-S20

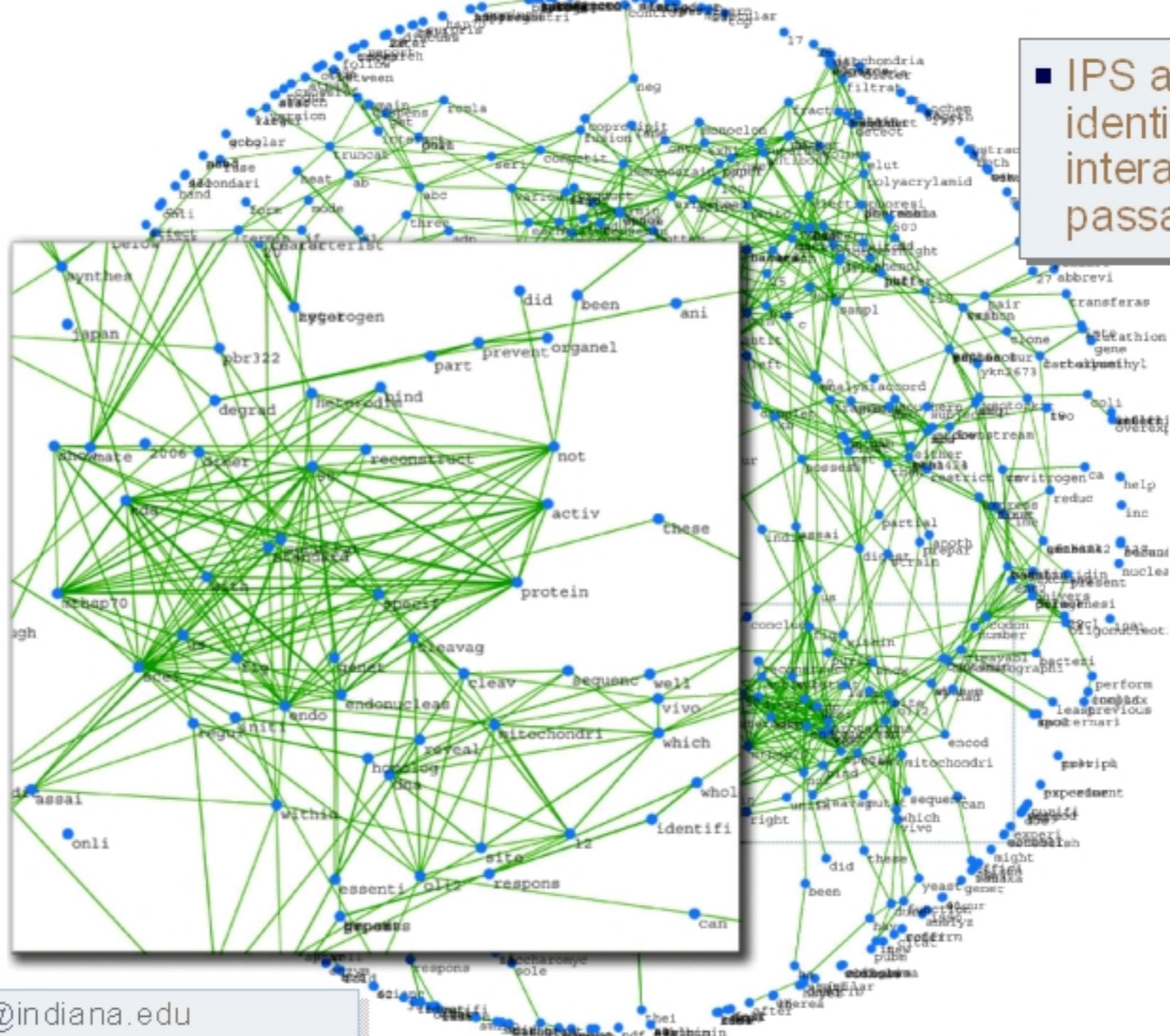


**informatics**  
luis rocha 2007



## proximity networks in Biocreative 2

Abi-Haidar et al [2008]. *Genome Biology*. In Press.



rocha@indiana.edu  
<http://informatics.indiana.edu/rocha>



INDIANA  
UNIVERSITY

## With Movielens data

- Movielens data
  - ▶ <http://www.grouplens.org/node/73>
- Dual Networks for each dataset
  - ▶ Movie (item) and User Proximity
- Results (F1 Measure)
  - ▶ Top 10 Movies
    - User-based: 0.213. Item-based: 0.184. Vector/Cosine: 0.212. LSI: 0.245
  - ▶ Arbitrary Number
    - User-based: 0.305. Item-based: 0.294. Vector/Cosine: 0.306. LSI: 0.328
- Results (Fouss' variant of Sommers D Measure)
  - ▶ Top 10 Movies
    - User-based: 88.2%. Item-based: 89.53%. Vector/Cosine: 86.02%. LSI: 91.69%
  - ▶ Arbitrary Number
    - User-based: 86.21%. Item-based: 87.95%. Vector/Cosine: 85.18%. LSI: 89.01

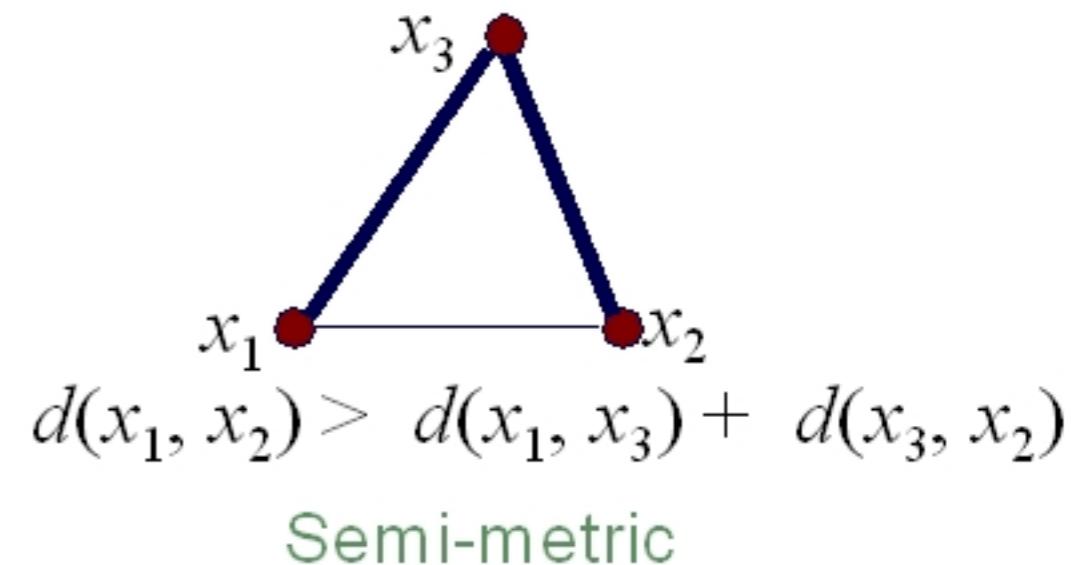
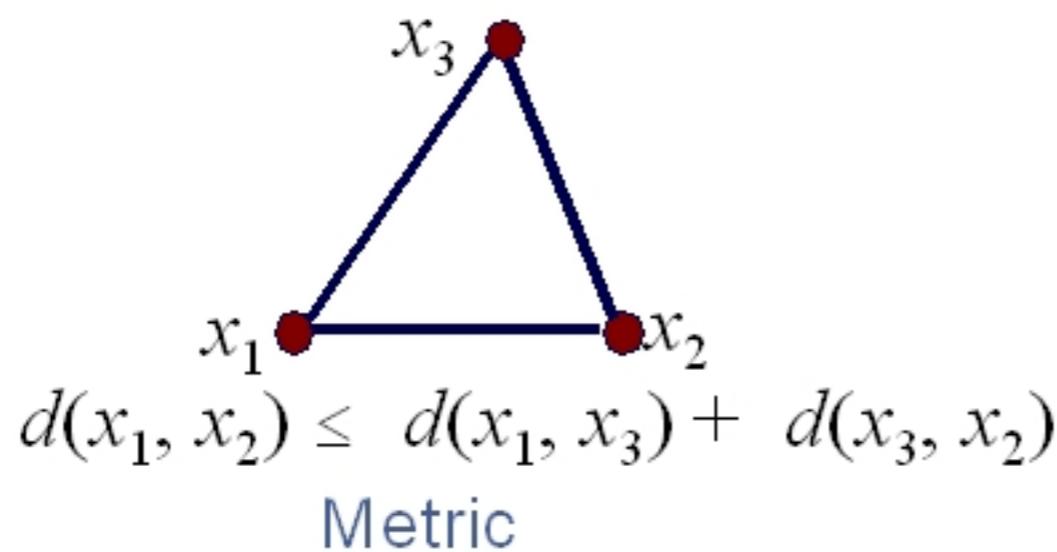


## identification of implicit associations in networks

### semi-metric behavior

$$d_X(x_i, x_j) = \frac{1}{XYP(x_i, x_j)} - 1; \quad d_Y(y_i, y_j) = \frac{1}{YXP(y_i, y_j)} - 1$$

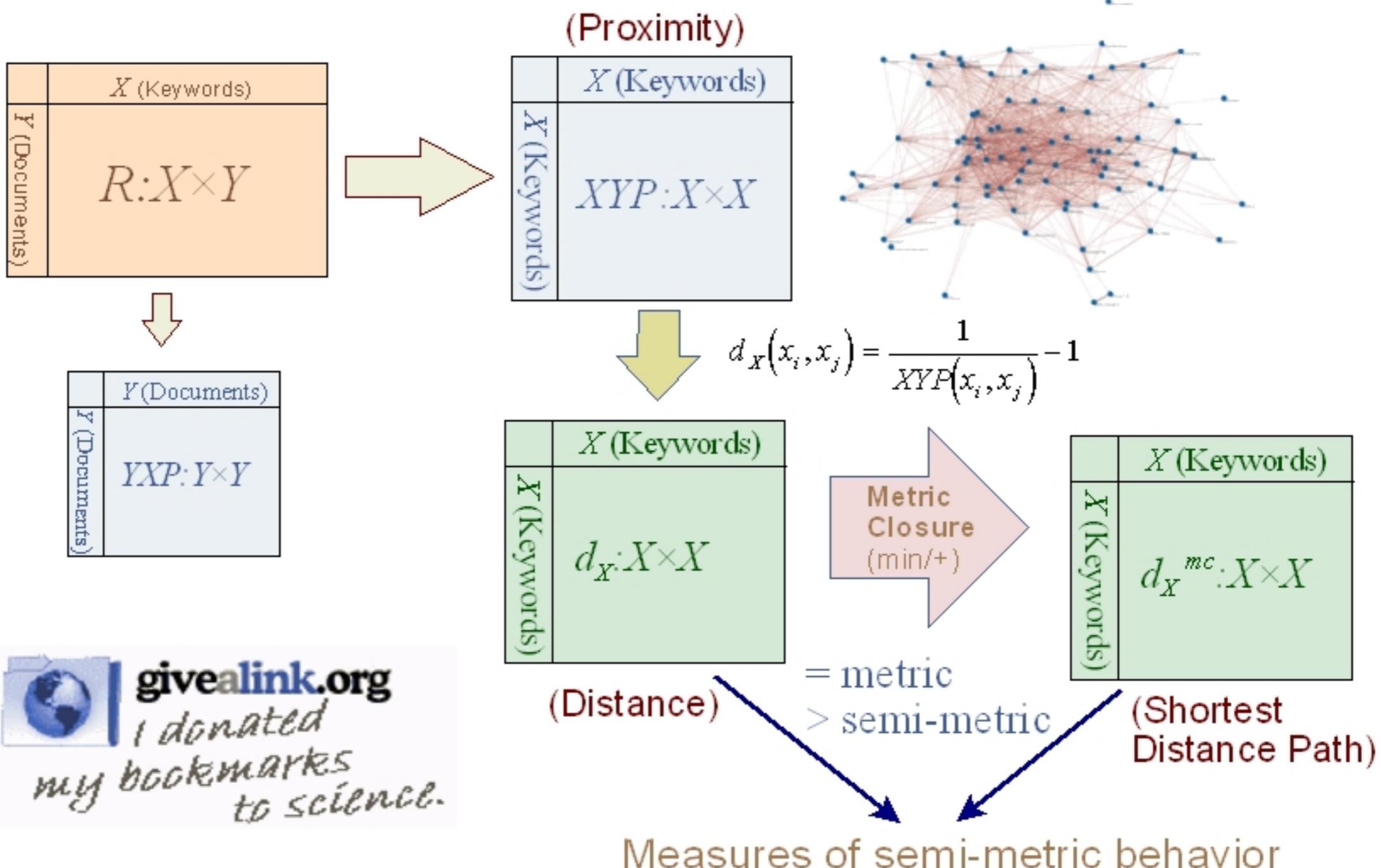
$d$  is a distance function because it is a nonnegative, symmetric, real-valued function such that  $d(k, k) = 0$



# computing semi-metric behavior



informatics  
luis rocha 2007





## Semi-metric Measures

### ■ Semi-metric ratio

- ▶ Absolute measure of indirect distance reduction

$$s(x_i, x_j) = \frac{d_{direct}(x_i, x_j)}{d_{shortest}(x_i, x_j)}$$

### ■ Relative Semi-metric ratio

- ▶ Distance reduction against maximum contraction

$$rs(x_i, x_j) = \frac{d_{direct}(x_i, x_j) - d_{shortest}(x_i, x_j)}{d_{max} - d_{min}}$$

### ■ Below Average Ratio

- ▶ Captures semi-metric distance reductions which contract to below the average distance for a given node. Captures some of the cases of initial  $\infty$  distance

$$b(x_i, x_j) = \frac{\overline{d}_{x_i}}{d_{shortest}(x_i, x_j)}$$

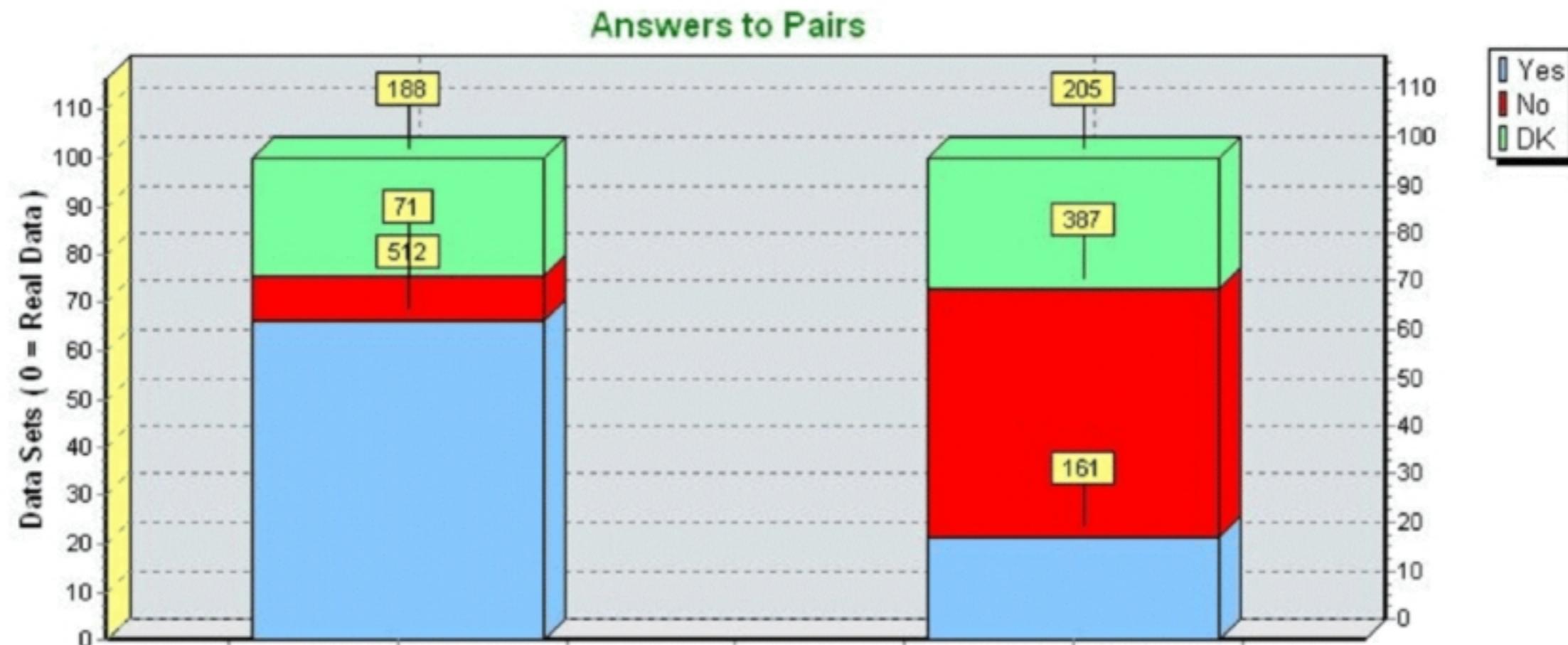
## catching strong indirect associations in mylibrary.lanl.gov

0020-1669--Inorganic chemistry 0031-9007--Physical review letters  
 0031-9007--Physical review letters 0743-7463--Langmuir  
 0003-2700--Analytical chemistry 0031-9007--Physical review letters  
 0096-3003--Applied mathematics and computation 0031-9007--Physical review letters  
 0031-9007--Physical review letters 0022-3115--Journal of nuclear materials  
 1049-3301--ACM transactions on modeling and computer simulation 0031-9007--Physical review letters  
 1364-548X--Chemical communications 0031-9007--Physical review letters  
 1064-8275--SIAM journal on scientific computing 0031-9007--Physical review letters  
 0965-5425--Computational mathematics and mathematical physics 0031-9007--Physical review letters  
 0031-9007--Physical review letters 1359-6454--Acta materialia  
 0003-7028--Applied spectroscopy 0031-9007--Physical review letters  
 0031-9007--Physical review letters 0022-2461--Journal of materials science  
 0031-9007--Physical review letters 1359-6462--Scripta materialia  
 0031-9007--Physical review letters 0022-4596--Journal of solid state chemistry  
 0031-9007--Physical review letters 0021-8898--Journal of applied crystallography  
 1097-6256--Nature neuroscience 1065-9471--Human Brain MAPPING  
 1097-6256--Nature neuroscience 0278-0062--IEEE transactions on medical imaging  
 1097-6256--Nature neuroscience 1053-8119--NeuroImage  
 1063-7796--Physics of particles and nuclei 0218-3013--International journal of modern physics E Nuclear physics  
 1053-8119--NeuroImage 1065-9471--Human Brain MAPPING  
 0031-9007--Physical review letters 0743-7463--Langmuir  
 0031-9007--Physical review letters 0020-1669--Inorganic chemistry  
 0031-9007--Physical review letters 0141-1594--Phase transitions  
 0031-9007--Physical review letters 0928-1045--Journal of computeraided materials design  
 0031-9007--Physical review letters 0042-207X--Vacuum  
 1097-6256--Nature neuroscience 0031-9155--Physics in medicine & biology  
 1097-6256--Nature neuroscience 0096-3518--IEEE transactions on acoustics speech and signal processing  
 1097-6256--Nature neuroscience 0740-7467--IEEE ASSP magazine  
 1097-6256--Nature neuroscience 1070-9908--IEEE signal processing letters  
 0022-5355--Journal of vacuum science and technology 0734-2101--Journal of vacuum science & technology A Vacuum surfaces and films

IPP\_3: parameter *b*

rocha@indiana.edu  
<http://informatics.indiana.edu/rocha>

## population from STB-RL and CCS-3



- 50% Pairs of ISSN with high  $rs$  and  $b$
- 50% Random ISSN Pairs
- 1524 answers

## what do semi-metric edges imply?

- Pairs with larger semi-metric behavior denote a *latent association*
  - ▶ Not grounded on direct evidence provided by the relation  $R$ , but rather implied by the overall network of associations in this relation.
  - ▶ Meaning depends on the semantics of the application
    - In graphs of keyword co-occurrence in documents: associated with novelty and can be used to identify trends.
    - In social networks it may identify pairs of people, groups, etc. for which we do not have direct evidence, in the available documents, that a real association exists, but who could easily be indirectly associated.

Rocha, Luis M. [2002]. "Semi-metric Behavior in Document Networks and its Application to Recommendation Systems". In: *Soft Computing Agents: A New Perspective for Dynamic Information Systems*. V. Loia (Ed.) International Series Frontiers in Artificial Intelligence and Applications. IOS Press, pp. 137-163.

Rocha, Luis M. [2002]. "Combination of Evidence in Recommendation Systems Characterized by Distance Functions". In: *Proceedings of the 2002 World Congress on Computational Intelligence: FUZZ-IEEE'02*. Honolulu, Hawaii, May 2002. IEEE Press, pp. 203-208.

Rocha, L.M., T. Simas, A. Rechtsteiner, M. DiGiacomo, R. Luce [2005]. 'MyLibrary@LANL: Proximity and Semi-metric Networks for a Collaborative and Recommender Web Service'. In: *Proc. 2005 IEEE/WIC/ACM International Conference on Web Intelligence (WT'05)*, IEEE Press, pp. 565-571.



# scientific community working on feynman diagrams

as published in *Physical Review*, 1949-54

$P(\text{author names})$	$P(\text{author names})$
$C:P \times P$	

- Collaboration Relation:  $C$ 
  - ▶ Who wrote a paper with whom
- Acknowledgment Relation:  $A$ 
  - ▶ Who acknowledged, or informally received information from whom

$P(\text{author names})$	$P(\text{author names})$
	$A:P \times P$

$$CP(p_i, p_j) = \frac{\sum_{k=1}^m (c_{i,k} \wedge c_{j,k})}{\sum_{k=1}^m (c_{i,k} \vee c_{j,k})}$$

76 Authors

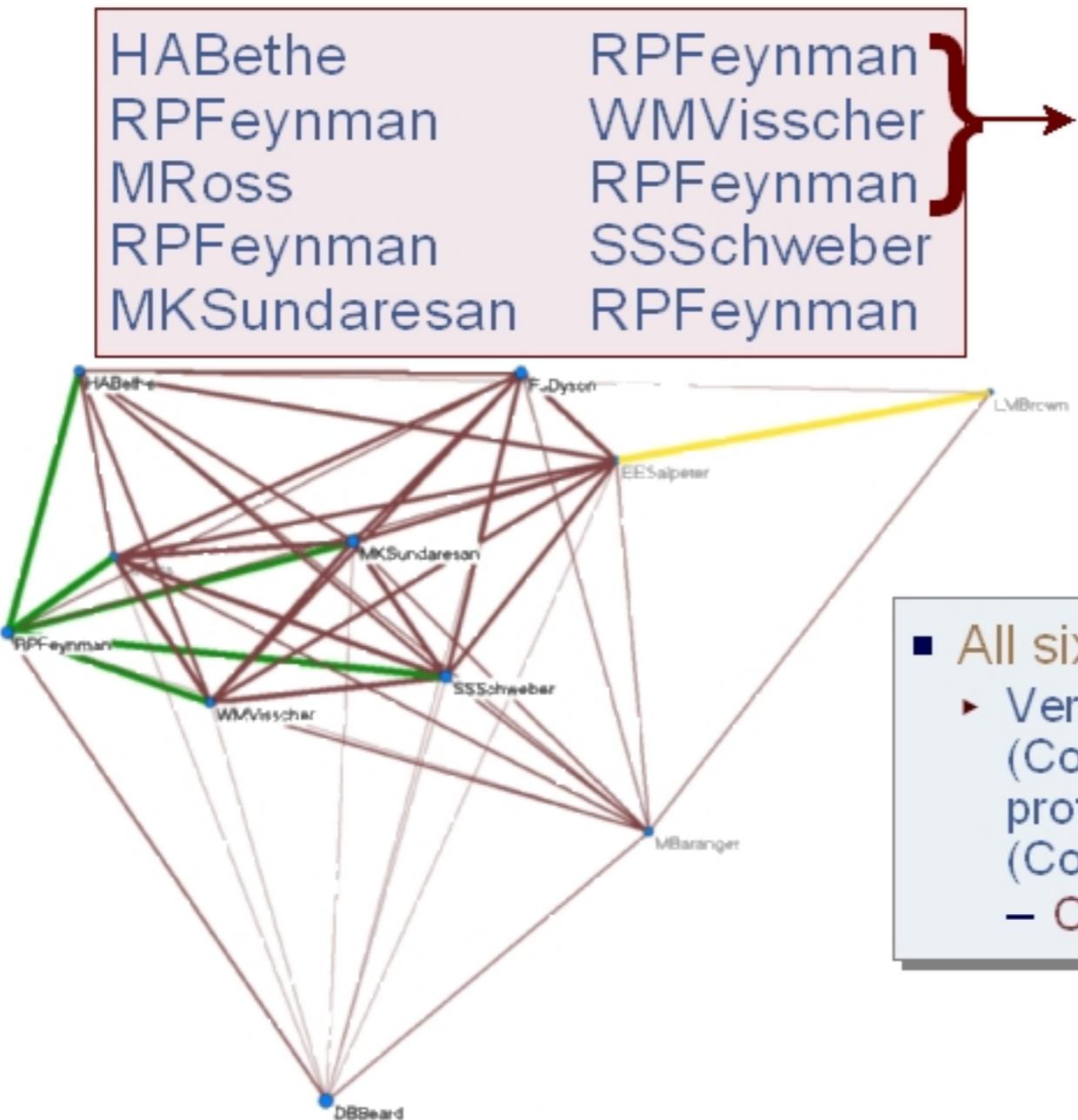
$CP(p_i, p_j)$  is a **co-collaboration probability**: the probability that two authors have collaborated with the same authors

$$AP(p_i, p_j) = \frac{\sum_{k=1}^m (a_{i,k} \wedge a_{j,k})}{\sum_{k=1}^m (a_{i,k} \vee a_{j,k})}$$

91 Authors

$AP(p_i, p_j)$  is a **co-acknowledgment probability**: the probability that two authors have acknowledged or have been acknowledged by the same authors

## 5 most semi-metric pairs (rs and b parameters)

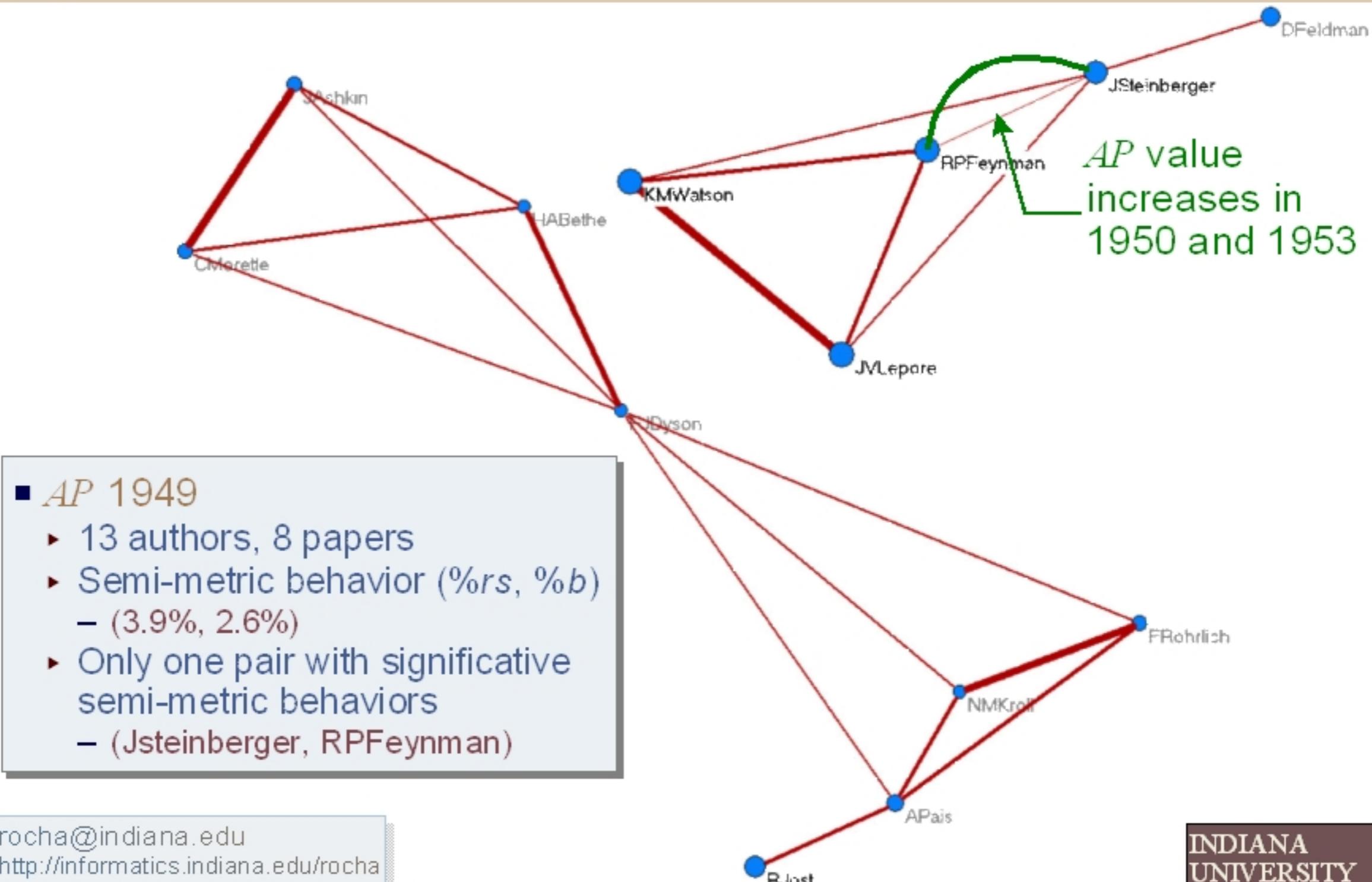


- Weak co-collaboration (*CP*), but strong co-acknowledgment (*AP*).
  - While they have not co-collaborated much in PR, they have acknowledged or have been acknowledged by many of the same people

- All six authors in top pairs
  - Very strong proximities to *EESalpeter* (Cornell) and *KMWatson* (postdoc at IAS, prof at Indiana) in *AP* and with to *EESalpeter* (Cornell) and *FJDyson* (Cornell, IAS) in *CP*
    - Cornell and Institute of Advanced Studies

# dynamics of co-acknowledgment network

1949



## NSF HSD - Dynamics of Human Behavior

- Collected Facebook data
  - ▶ Subset of IUB volunteers
    - **Questionnaires** about themselves, their relationship partner, their use of facebook, et.
    - **Ratings.** show elements of their profile pages to raters, who can rate physical attractiveness of the photo, etc.
  - ▶ Network data
    - 2 time instances with user (333 and 260) profiles as well as friends at IU, and on those who made the most recent 10 posts on their "wall."
    - Other linked users scrambled from one time measurement to the other
  - ▶ Testing
    - Semi-metric links in co-friendship network as predictor of new friends
    - Keyword content as predictor of communities, attractiveness, etc.

With Elliot Smith and Rob Goldstone

## shortest paths or weakest links

## Proximity

	$X(\text{Keywords})$
$X(\text{Keywords})$	$XYP: X \times X$

For any monotonic  
increasing distance function

$$d_X(x_i, x_j) = \frac{1}{XYP(x_i, x_j)} - 1$$

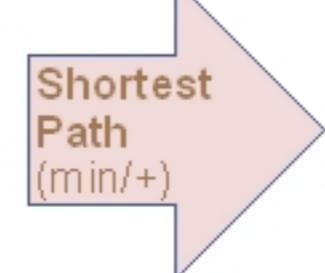
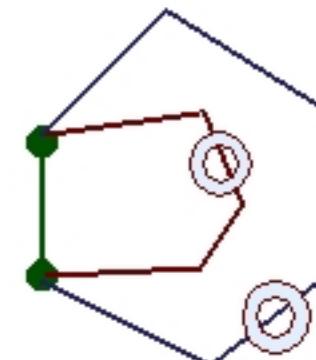
	$X(\text{Keywords})$
$X(\text{Keywords})$	$d_X: X \times X$

Distance  
semi-metric

Edges: largest of the  
weakest links in all paths:

## Similarity

	$X(\text{Keywords})$
$X(\text{Keywords})$	$XYS: X \times X$



	$X(\text{Keywords})$
$X(\text{Keywords})$	$d^*_X: X \times X$

Edges: shortest path  
(sums all edges)  
metric



	$X(\text{Keywords})$
$X(\text{Keywords})$	$d^{**}_X: X \times X$

Edges: Smallest of the  
largest edges in each path  
ultra-metric:  
immune to indirect  
path length

# effect of typical transitive closure

IPP: cumulative strength distribution (all weights)

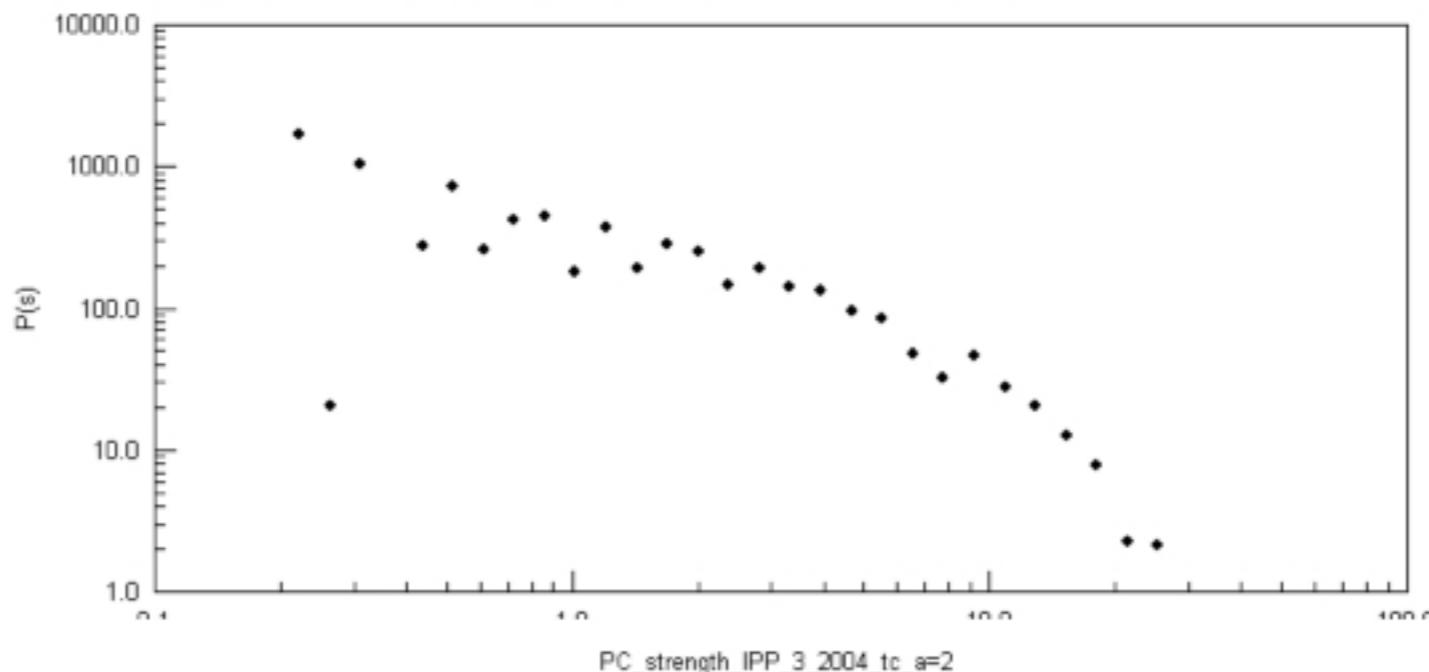
Proximity

	$X$ (Keywords)
$X$ (Keywords)	$XYP: X \times X$

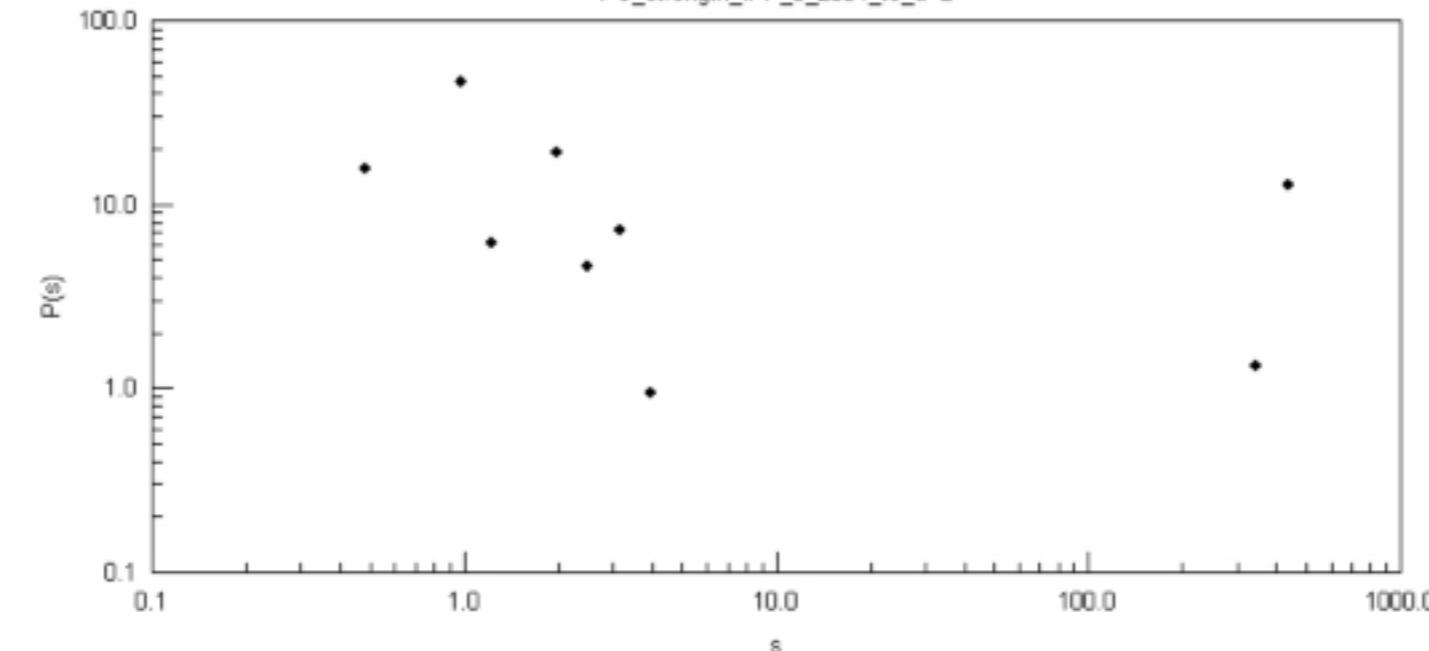
Transitive Closure  
(max/min)

Similarity

	$X$ (Keywords)
$X$ (Keywords)	$XYS: X \times X$



$\alpha$ -cut = 0.2



Destroys  
scale  
free  
behavior

# comparison of the two closures

## IPP: cumulative strength distribution

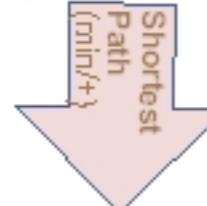
Proximity	
X (Keywords)	X (Keywords)
	$XYP: X \times X$



$$d_X(x_i, x_j) = \frac{1}{XYP(x_i, x_j)} - 1$$

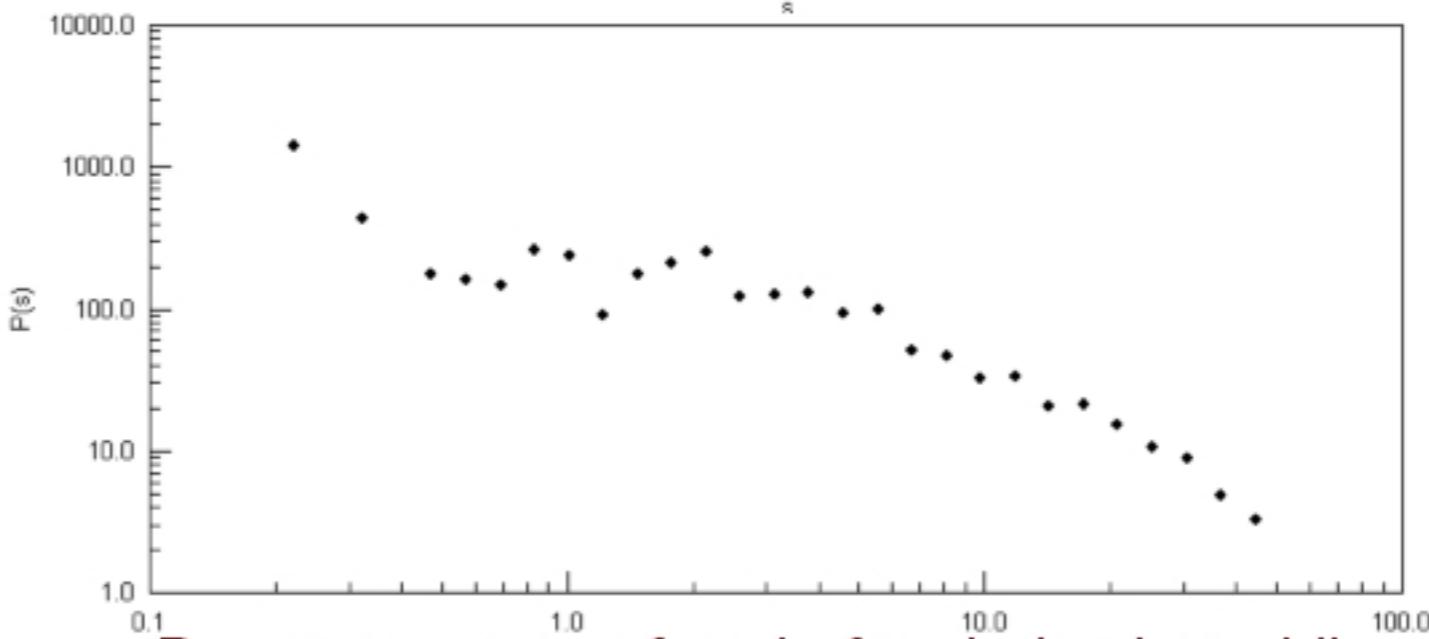
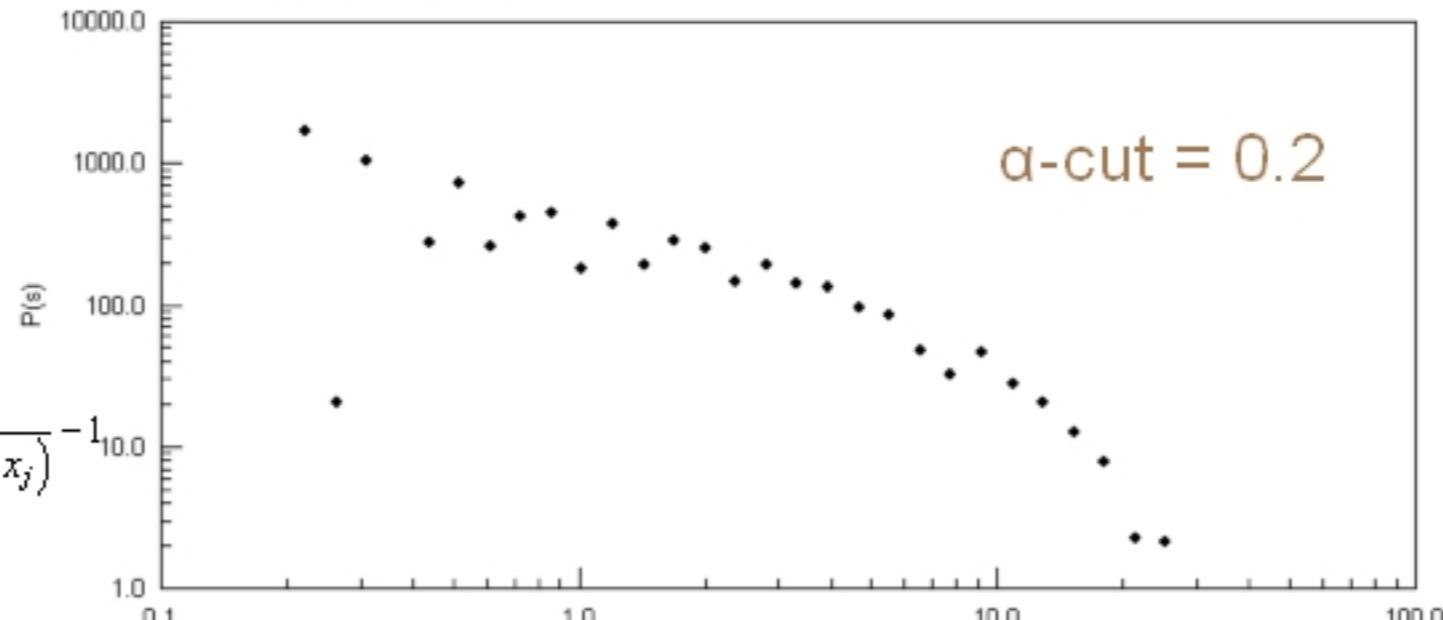
X (Keywords)	X (Keywords)
	$d_X: X \times X$

Distance



Distance  
Closure

X (Keywords)	X (Keywords)
	$d^*_X: X \times X$



Preserves more of scale free behavior, while still capturing indirect (transitive/latent) associations in the data

## Hammacher function

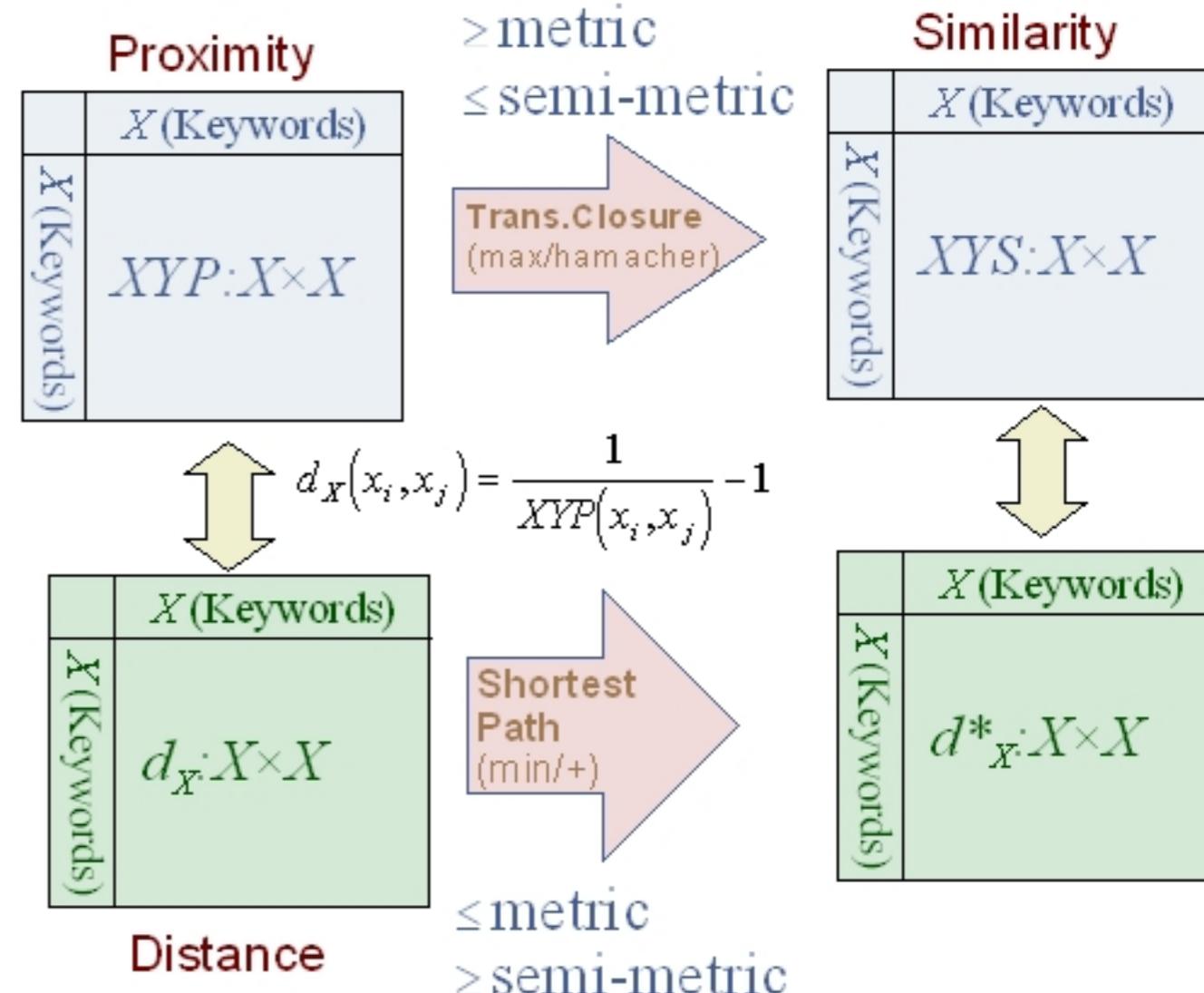
If we use :

$$i(a, b) = \frac{ab}{a + b - ab}$$

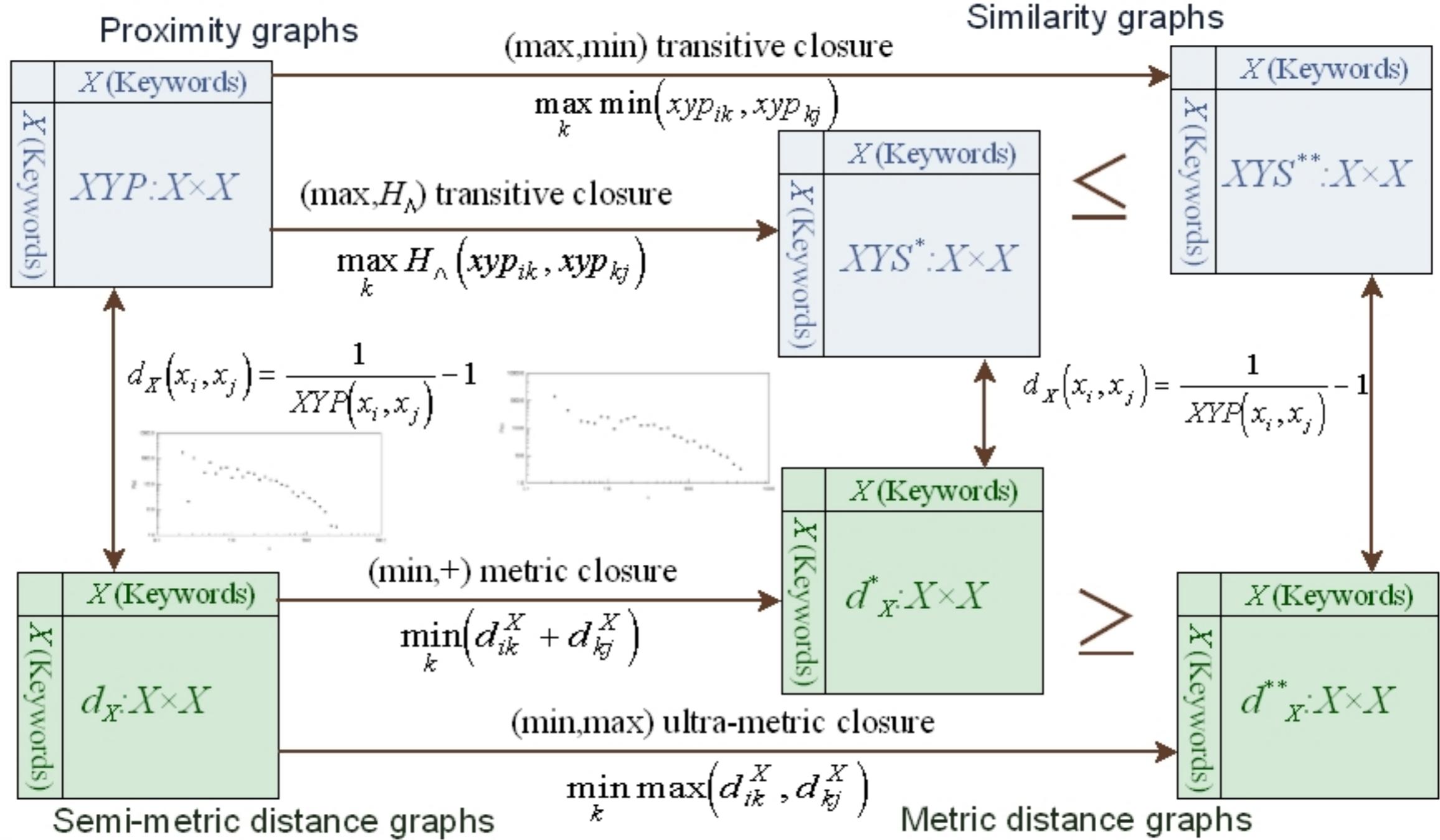
Hammacher Intersection

$$u(a, b) = \max[a, b]$$

With Tiago Simas



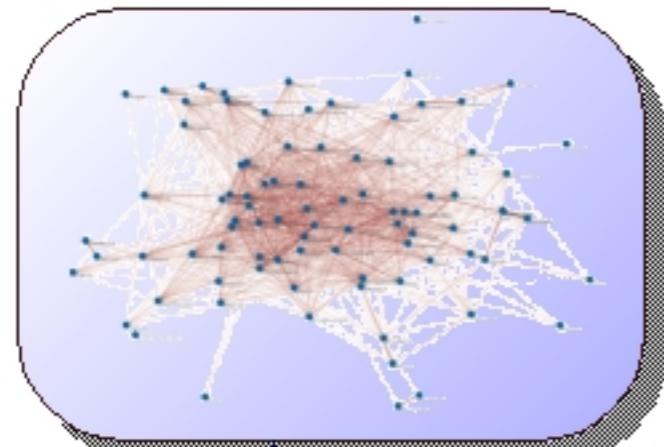
## metric and ultra-metric closures on distance graphs



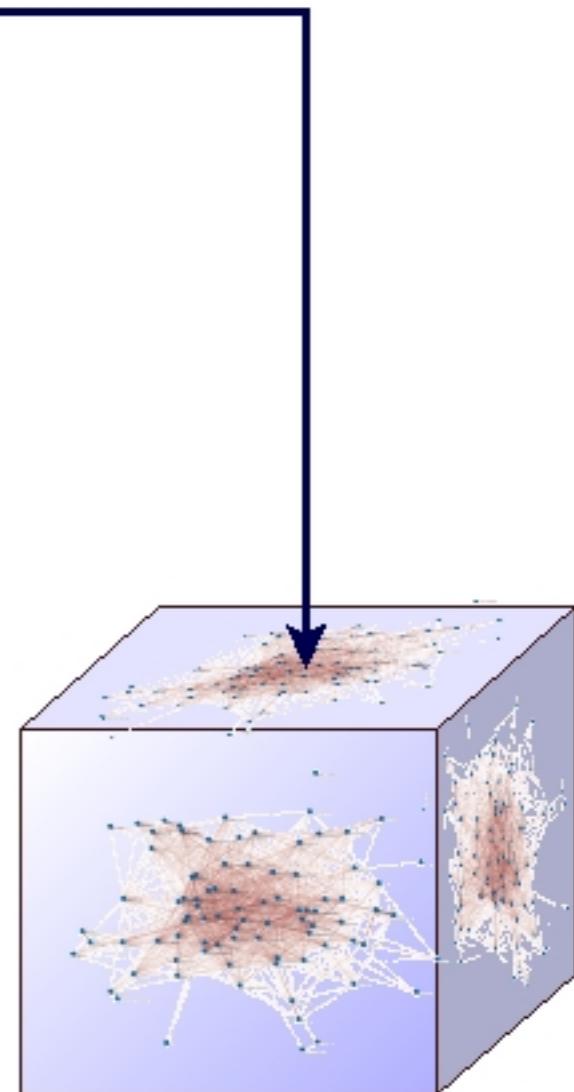
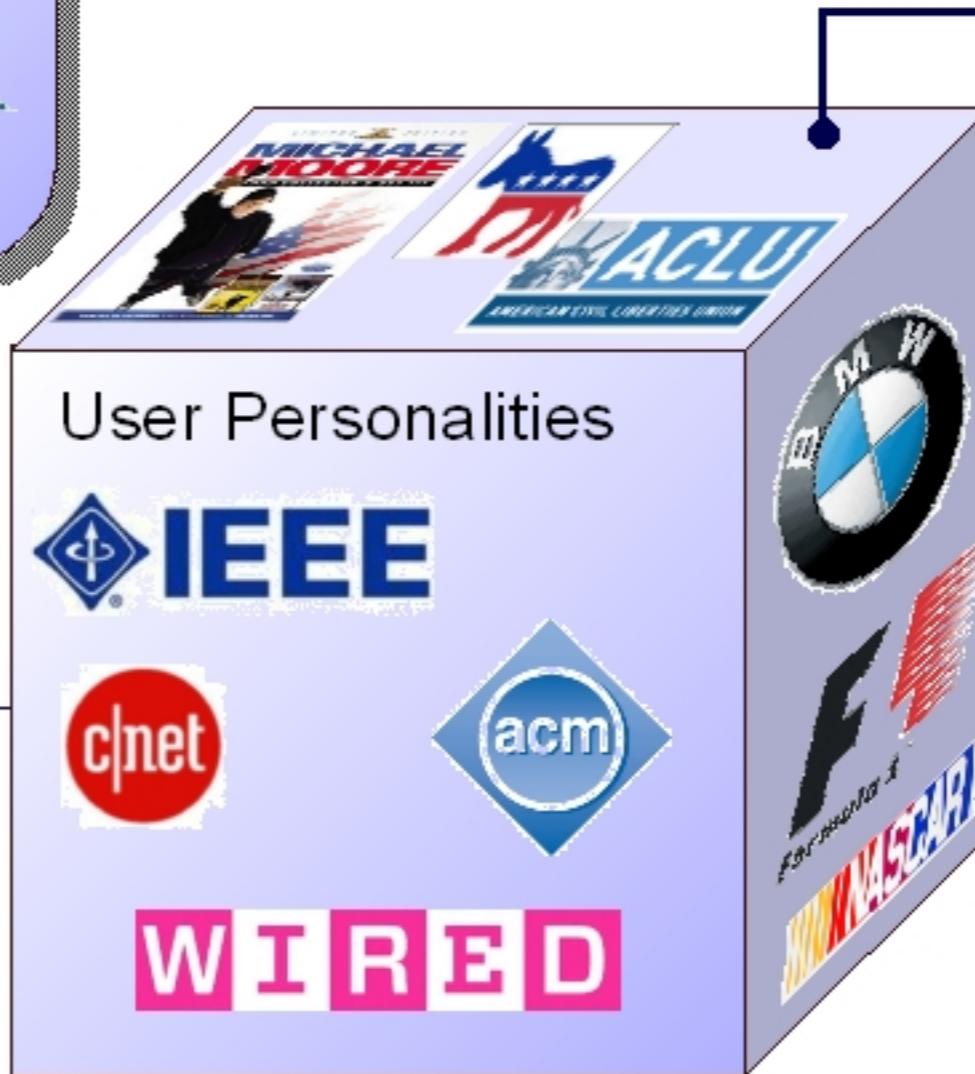


collecting data from users with multiple personalities

from single to multifaceted knowledge representations



Need for effective integration of various  
knowledge networks



What are the best ways to combine graph knowledge representations?

### Hammacher Intersection

$$i(a, b) = \frac{ab}{a + b - ab}$$

Metric closure is not dual: no complement can satisfy de Morgan's laws!

$$u(a, b) = \max[a, b]$$

- What fuzzy intersection/union comes closest to the metric closure of distance graphs?
  - ▶ Can we have a fuzzy union/intersection that follows De Morgan's laws but make it as close as possible to a metric closure?
- What captures the semi-metric behavior best?
  - ▶ The metric closure or some other transitive closure?

## Dombi family

$$u(a,b) = a \cup b = \frac{1}{1 + \left[ \left( \frac{1}{a} - 1 \right)^{-\lambda} + \left( \frac{1}{b} - 1 \right)^{-\lambda} \right]^{-\frac{1}{\lambda}}}$$

$$\lambda \rightarrow 0 \Rightarrow a \cup b \rightarrow 1$$

$$\lambda = 1 \Rightarrow a \cup b = \frac{a + b - 2ab}{1 - ab}$$

Hammacher union

$$\lambda \rightarrow \infty \Rightarrow a \cup b \rightarrow \max(a, b)$$

## De Morgan Laws and Involutive Complement

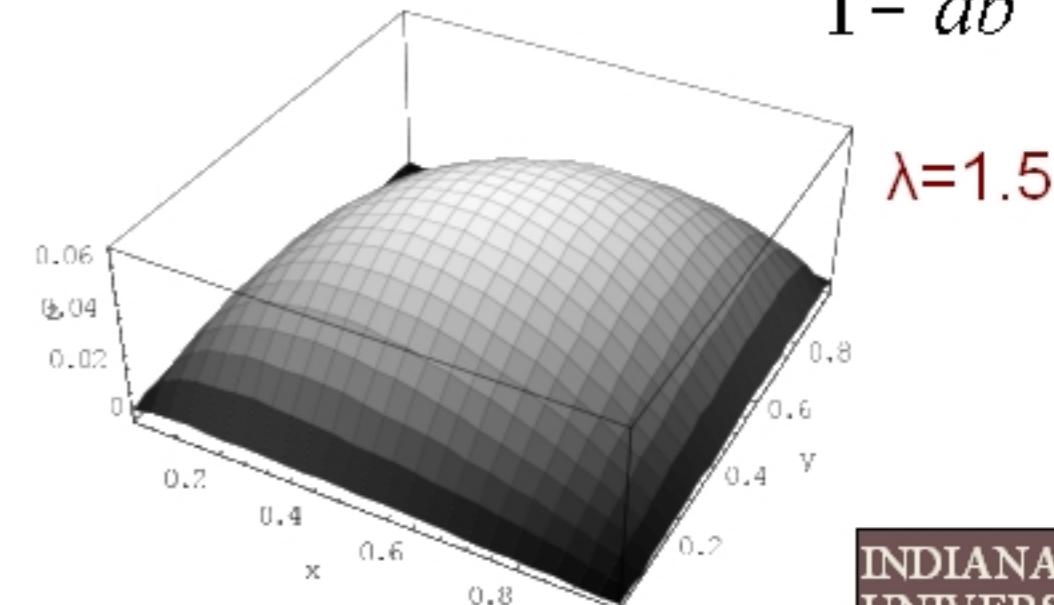
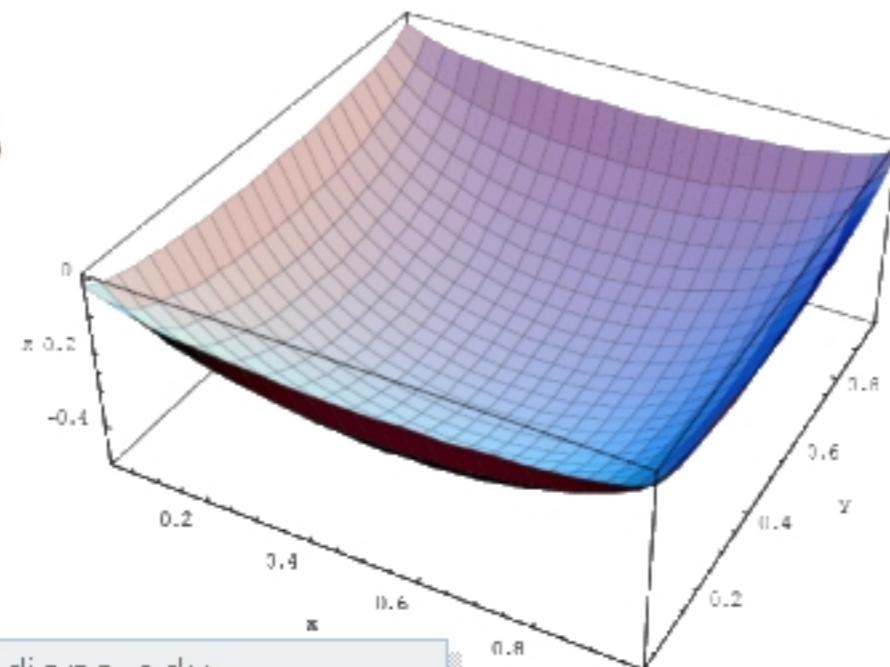
$$i(a,b) = \frac{ab}{a+b-ab} \quad u(a,b) = \frac{1}{1 + \left[ \left( \frac{1}{a} - 1 \right)^{-\lambda} + \left( \frac{1}{b} - 1 \right)^{-\lambda} \right]^{-\frac{1}{\lambda}}}$$

Requirement for desired axiomatics

$$-ab \left[ \left( \frac{1}{a} - 1 \right)^\lambda + \left( \frac{1}{b} - 1 \right)^\lambda \right]^{\frac{1}{\lambda}} + a + b - 2ab = 0$$

Unique solution for  $\lambda=1$ 

$$\lambda = 1 \Rightarrow u(a,b) = \frac{a+b-2ab}{1-ab}$$

 $\lambda=0.5$ 



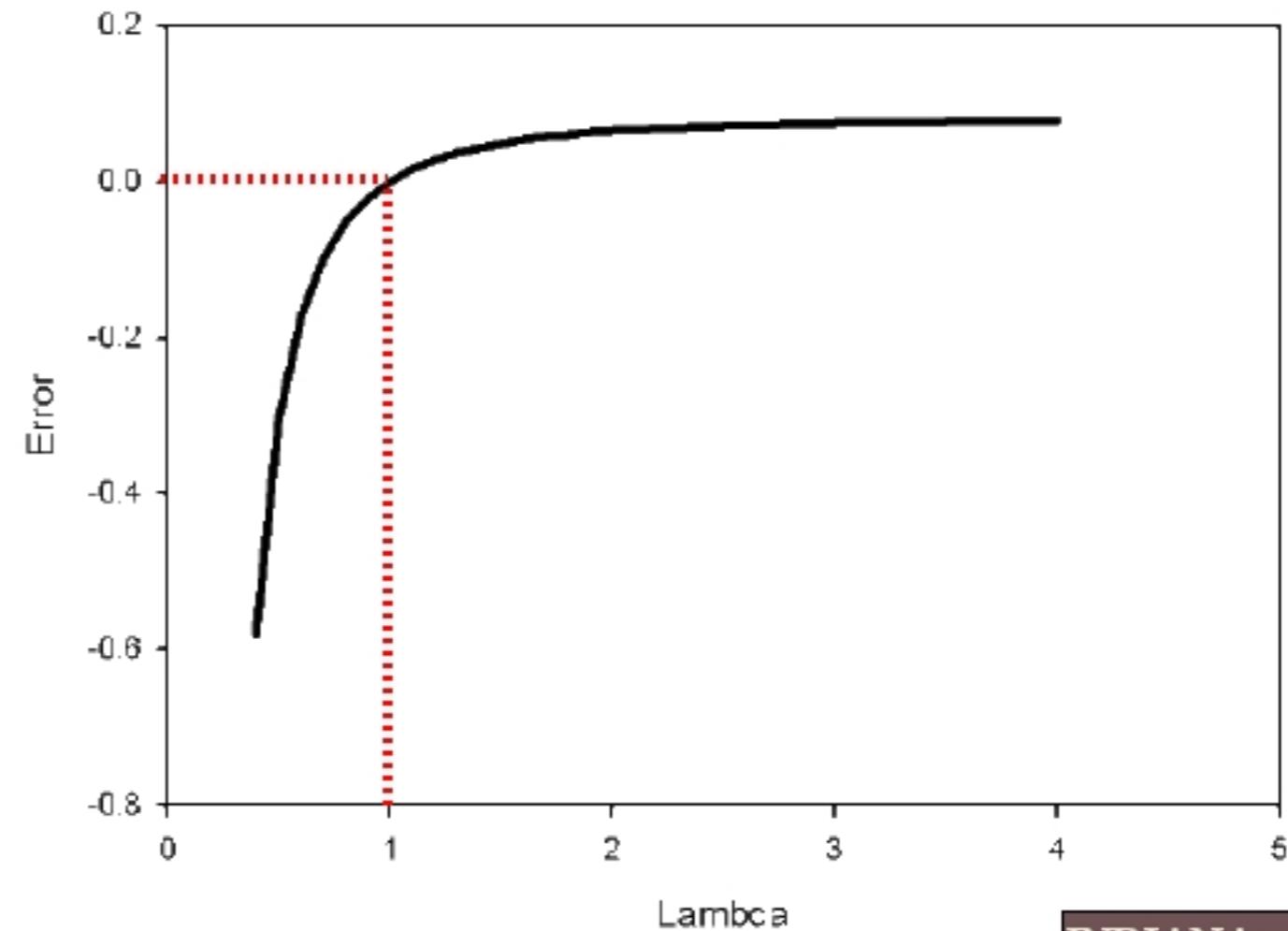
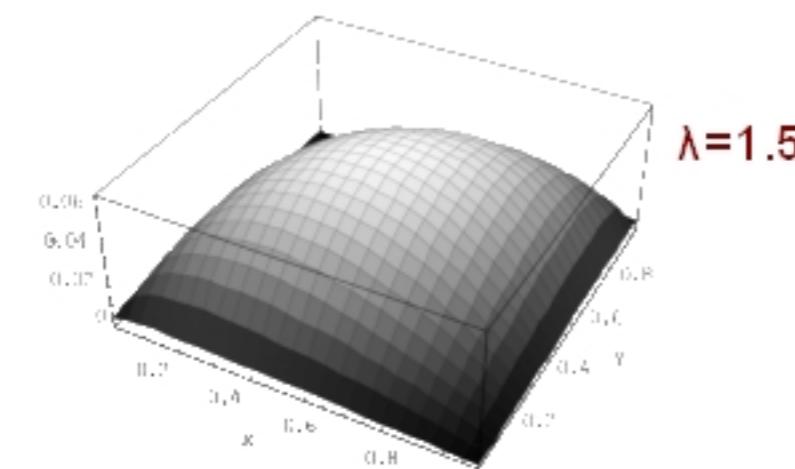
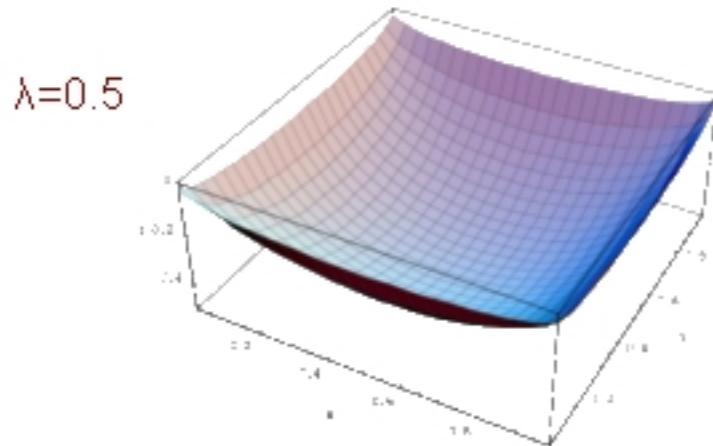
informatics  
luis rocha 2007



## De Morgan Laws and Involutive Complement

Error for other  $\lambda$

$$F(\lambda) = \int_0^1 \int_0^1 \left( -xy \left[ \left( \frac{1}{x} - 1 \right)^\lambda + \left( \frac{1}{y} - 1 \right)^\lambda \right]^{\frac{1}{\lambda}} + x + y - 2xy \right) dx dy$$



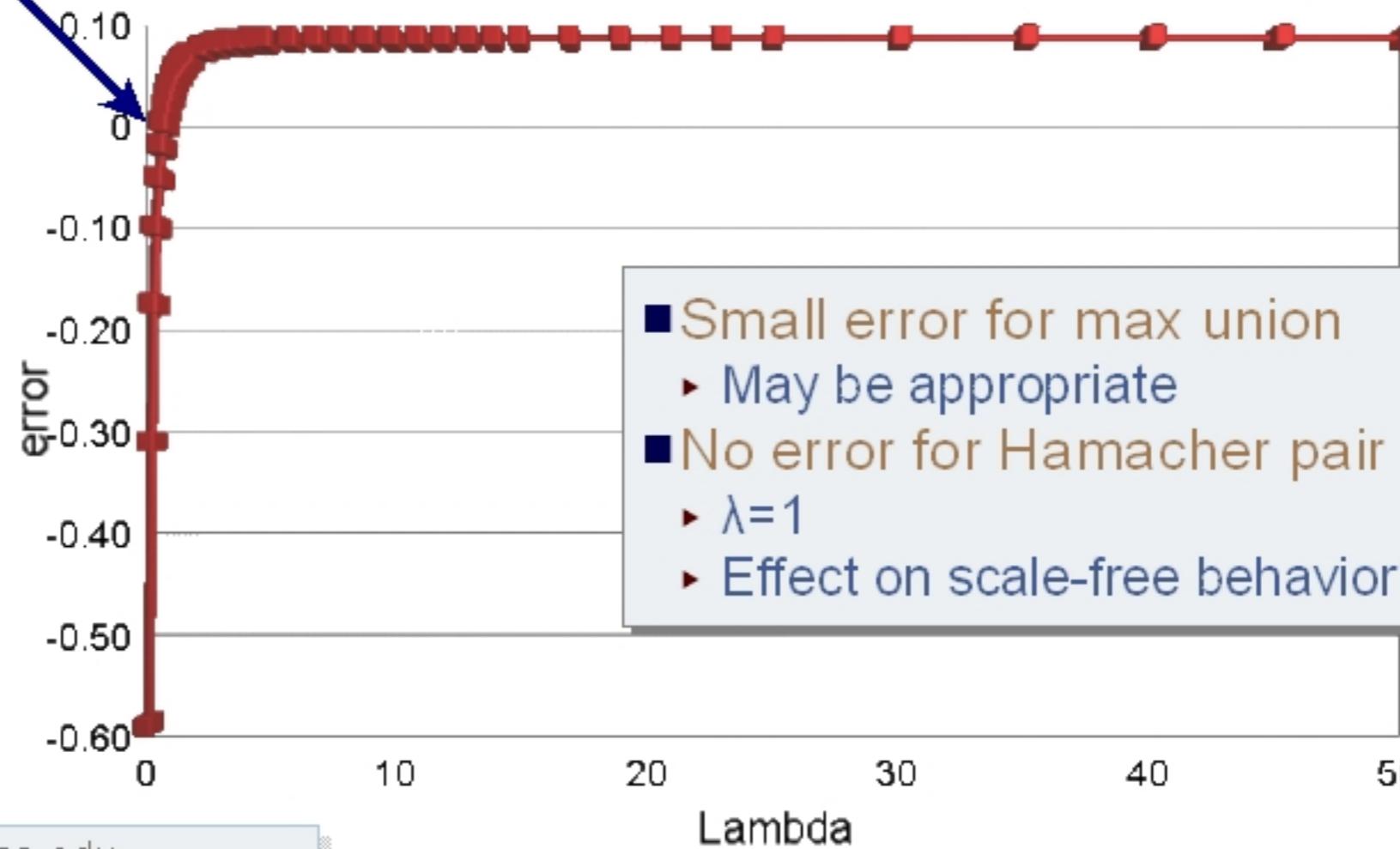


informatics  
luis rocha 2007



$$\lambda = 1 \Rightarrow a \cup b = \frac{a + b - 2ab}{1 - ab}$$

$$\lambda \rightarrow \infty \Rightarrow a \cup b \rightarrow \max(a, b)$$



## metric and ultra-metric closures on distance graphs

