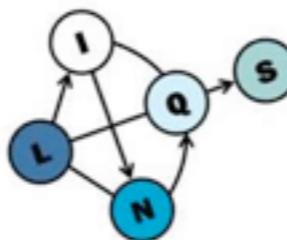


Graph Summarization in Annotated Data using Probabilistic Soft Logic

Alex Memory¹, Angelika Kimmig^{1,2}, Stephen H. Bach¹,
Louiqa Raschid¹ and Lise Getoor¹

¹University of Maryland

²KU Leuven

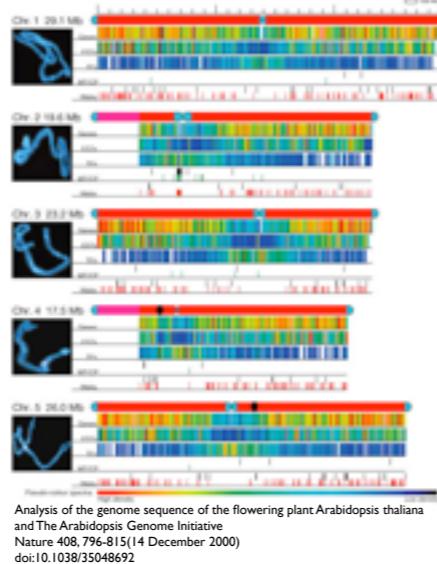


KU LEUVEN

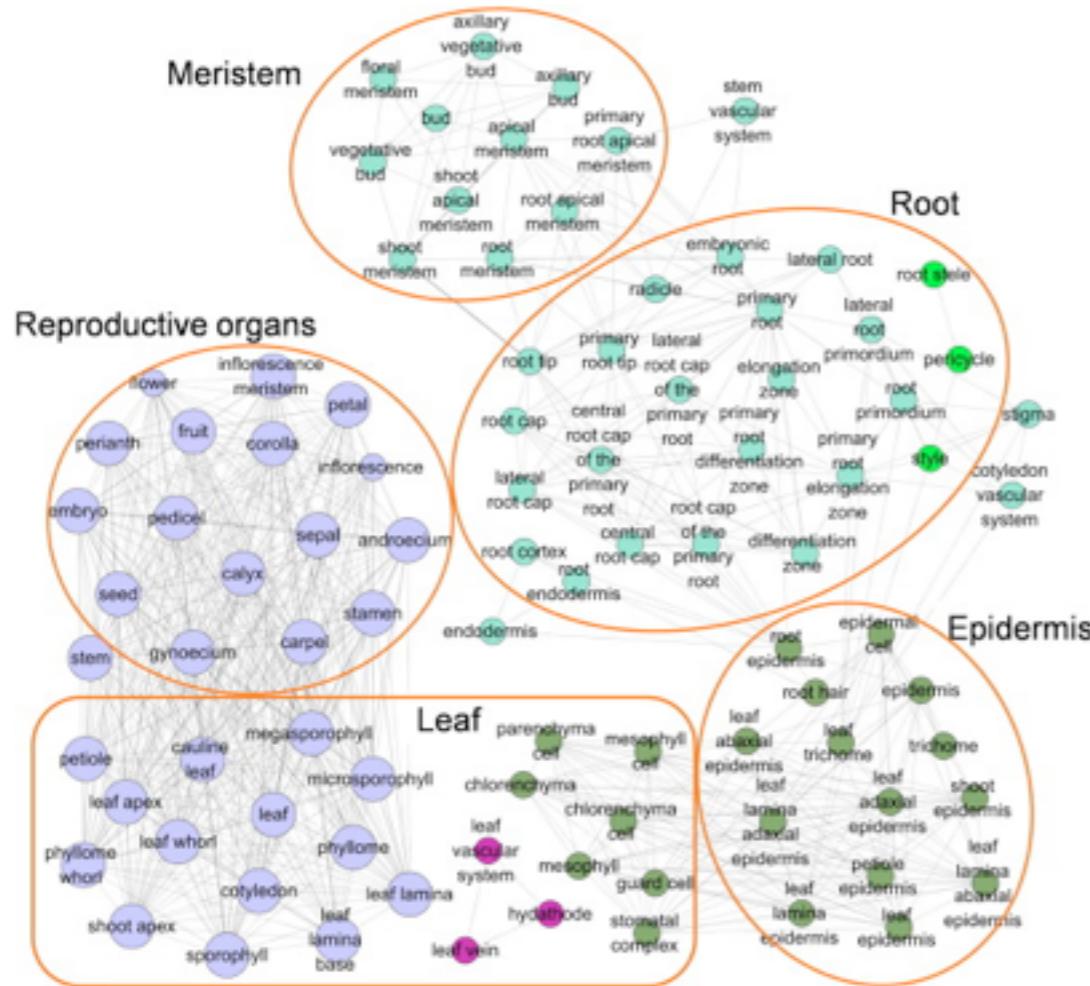
Agenda

- Annotation Graphs and Graph Summarization
- Heuristics for Graph Summarization
- Probabilistic Soft Logic
- Evaluation and Results

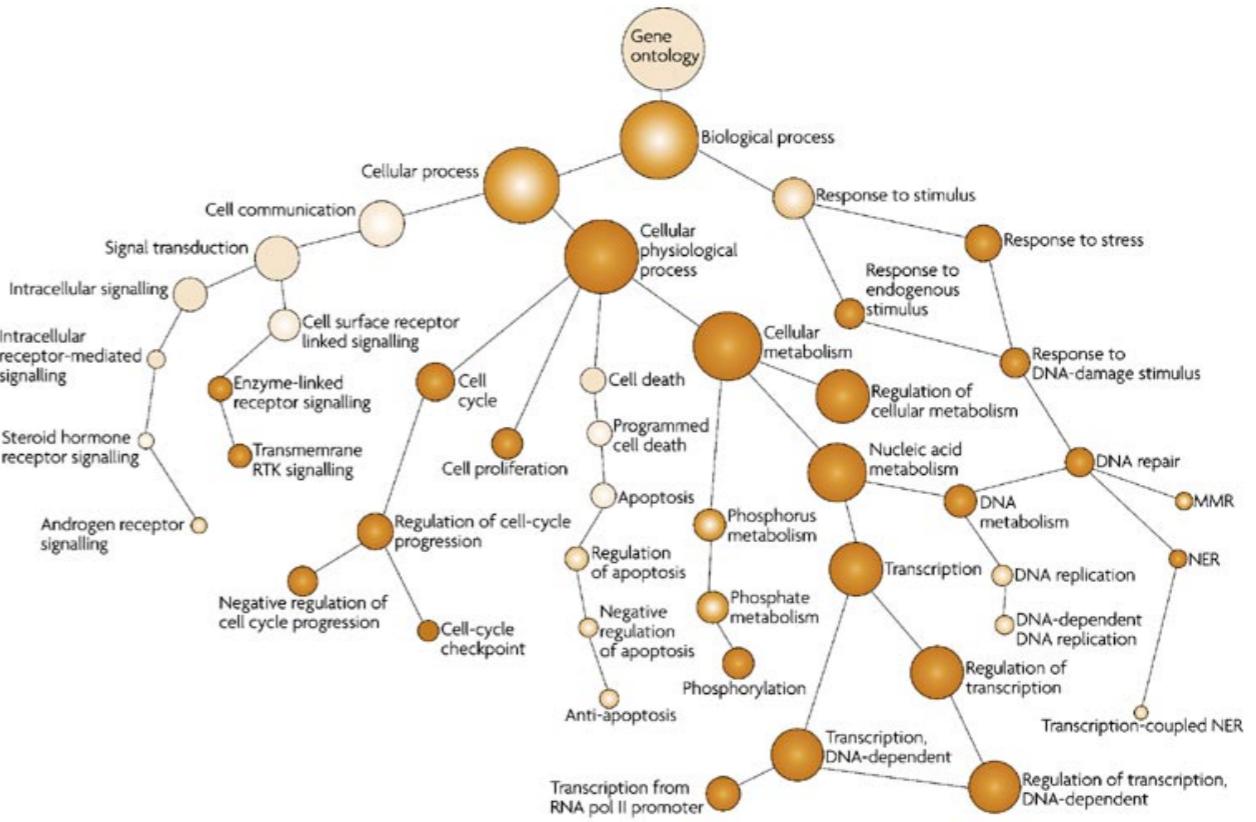
Annotated Data: the *Arabidopsis thaliana* Genome



Plant Ontology



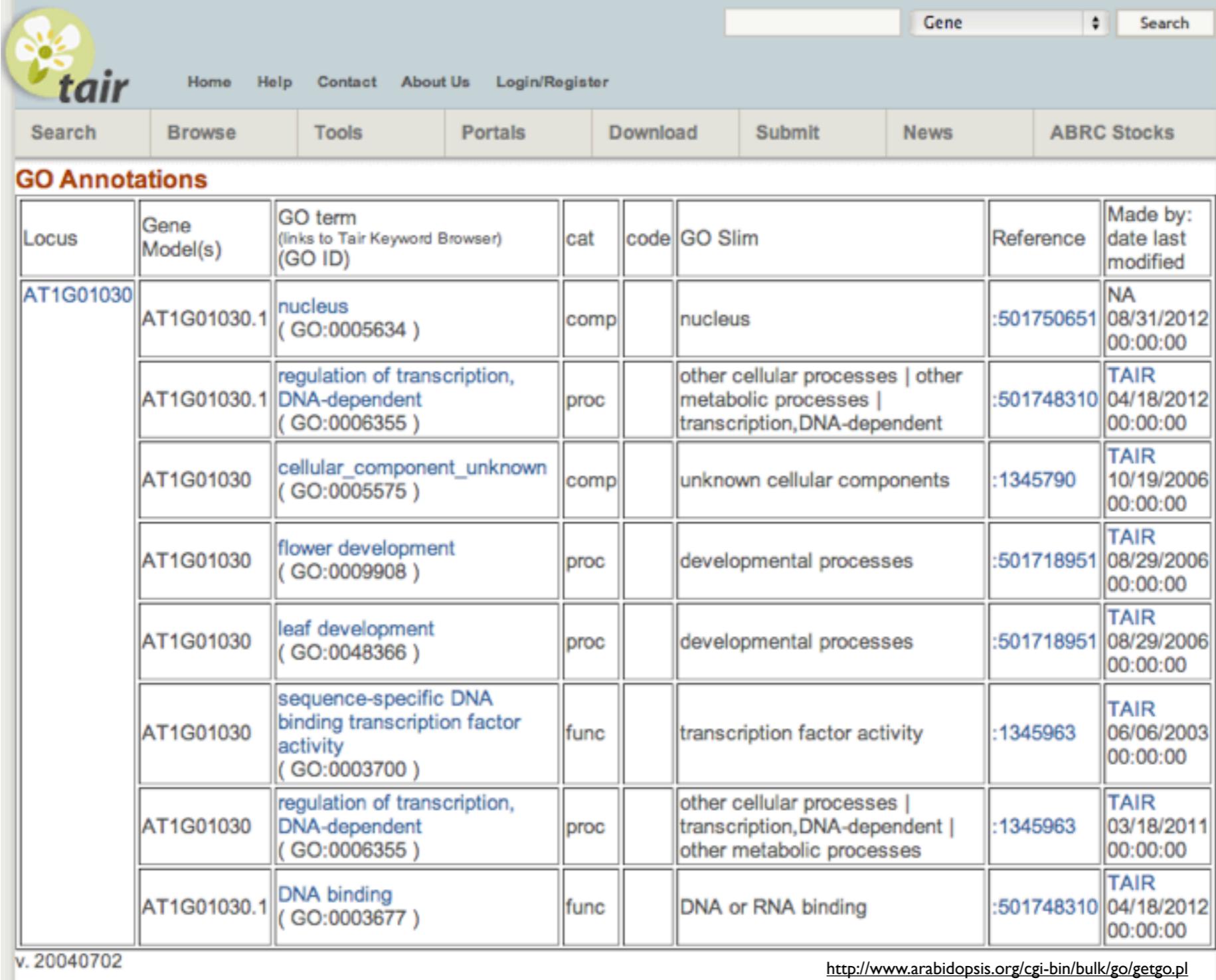
Gene Ontology



Hu et al. *Nature Reviews Cancer* 7, 23–34 (January 2007) | doi:10.1038/nrc2036

Xu Wang, Yangyang Bianb, Kai Chengb, Li-Fei Gua, Mingliang Yeb, Hanfa Zoub, Jun-Xian Hea,, A large-scale protein phosphorylation analysis reveals novel phosphorylation motifs and phosphoregulatory networks in *Arabidopsis*, *Journal of Proteomics*, Oct 2012.

The Arabidopsis Information Resource



The screenshot shows the TAIR (Arabidopsis Information Resource) website interface. At the top, there is a navigation bar with links for Home, Help, Contact, About Us, and Login/Register. Below the navigation bar is a search bar with dropdown options for "Gene" and "Search". The main content area is titled "GO Annotations" and displays a table of gene annotations. The table has columns for Locus, Gene Model(s), GO term (links to Tair Keyword Browser) (GO ID), cat, code, GO Slim, Reference, and Made by: date last modified. The data in the table is as follows:

Locus	Gene Model(s)	GO term (links to Tair Keyword Browser) (GO ID)	cat	code	GO Slim	Reference	Made by: date last modified
AT1G01030	AT1G01030.1	nucleus (GO:0005634)	comp		nucleus	:501750651	NA 08/31/2012 00:00:00
	AT1G01030.1	regulation of transcription, DNA-dependent (GO:0006355)	proc		other cellular processes other metabolic processes transcription,DNA-dependent	:501748310	TAIR 04/18/2012 00:00:00
	AT1G01030	cellular_component_unknown (GO:0005575)	comp		unknown cellular components	:1345790	TAIR 10/19/2006 00:00:00
	AT1G01030	flower development (GO:0009908)	proc		developmental processes	:501718951	TAIR 08/29/2006 00:00:00
	AT1G01030	leaf development (GO:0048366)	proc		developmental processes	:501718951	TAIR 08/29/2006 00:00:00
	AT1G01030	sequence-specific DNA binding transcription factor activity (GO:0003700)	func		transcription factor activity	:1345963	TAIR 06/06/2003 00:00:00
	AT1G01030	regulation of transcription, DNA-dependent (GO:0006355)	proc		other cellular processes transcription,DNA-dependent other metabolic processes	:1345963	TAIR 03/18/2011 00:00:00
	AT1G01030.1	DNA binding (GO:0003677)	func		DNA or RNA binding	:501748310	TAIR 04/18/2012 00:00:00

v. 20040702 <http://www.arabidopsis.org/cgi-bin/bulk/go/getgo.pl>

PO

Gene

GO

[cauline leaf](#)

[shoot apex](#)

[inflorescence meristem](#)

[leaf lamina base](#)

[plant embryo](#)

[vascular leaf](#)

[sepal](#)

[petal](#)

[cotyledon](#)

[petiole](#)

[leaf apex](#)

CRY2

PHOT1

CIB5

COP1

[vacuole \(ebi\)](#)

[response to blue light \(ebi\)](#)

[response to water deprivation \(ebi\)](#)

[regulation of flower development \(ebi\)](#)

[chromatin remodeling \(ebi\)](#)

[phototropism \(ebi\)](#)

[stomatal movement \(ebi\)](#)

[positive regulation of flower development \(ebi\)](#)

[blue light photoreceptor activity \(ebi\)](#)

[regulation of meristem growth \(ebi\)](#)

[chloroplast avoidance movement \(ebi\)](#)

[cytoplasm \(ebi\)](#)

[protein serine/threonine kinase activity \(ebi\)](#)

[FMN binding \(ebi\)](#)

[negative regulation of anion channel activity \(ebi\)](#)

[regulation of stomatal movement \(ebi\)](#)

[photomorphogenesis \(ebi\)](#)

[photoperiodism, flowering \(ebi\)](#)

[anthocyanin metabolic process \(ebi\)](#)

[nuclear ubiquitin ligase complex \(ebi\)](#)

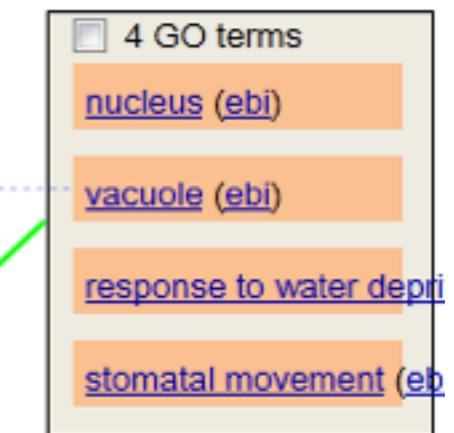
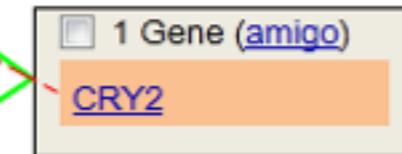
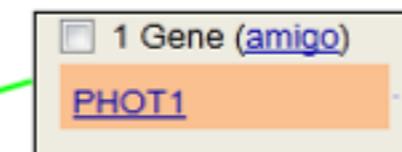
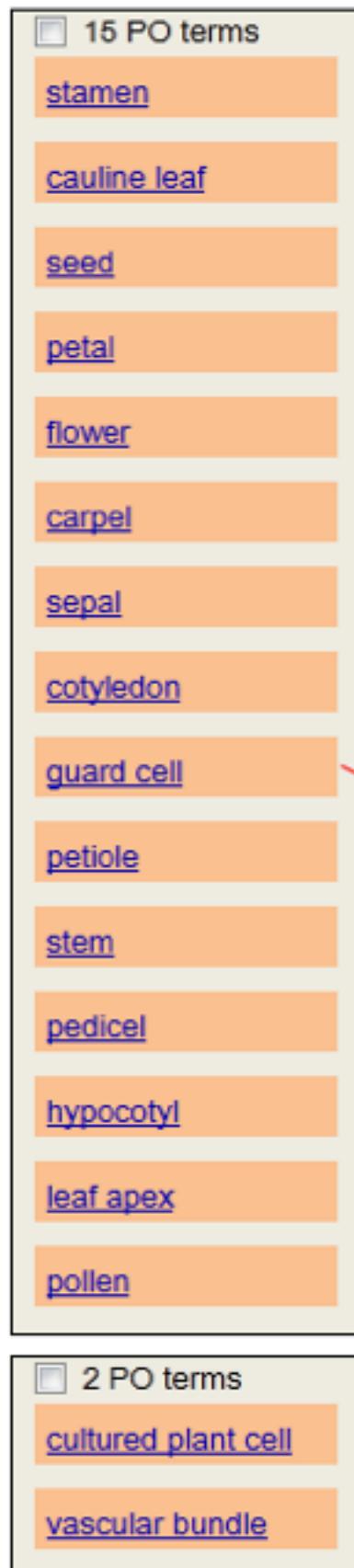
[skotomorphogenesis \(ebi\)](#)

Graph Summarization Example I

PO

Gene

GO



Legend: — Superedge - - - Deletion Addition

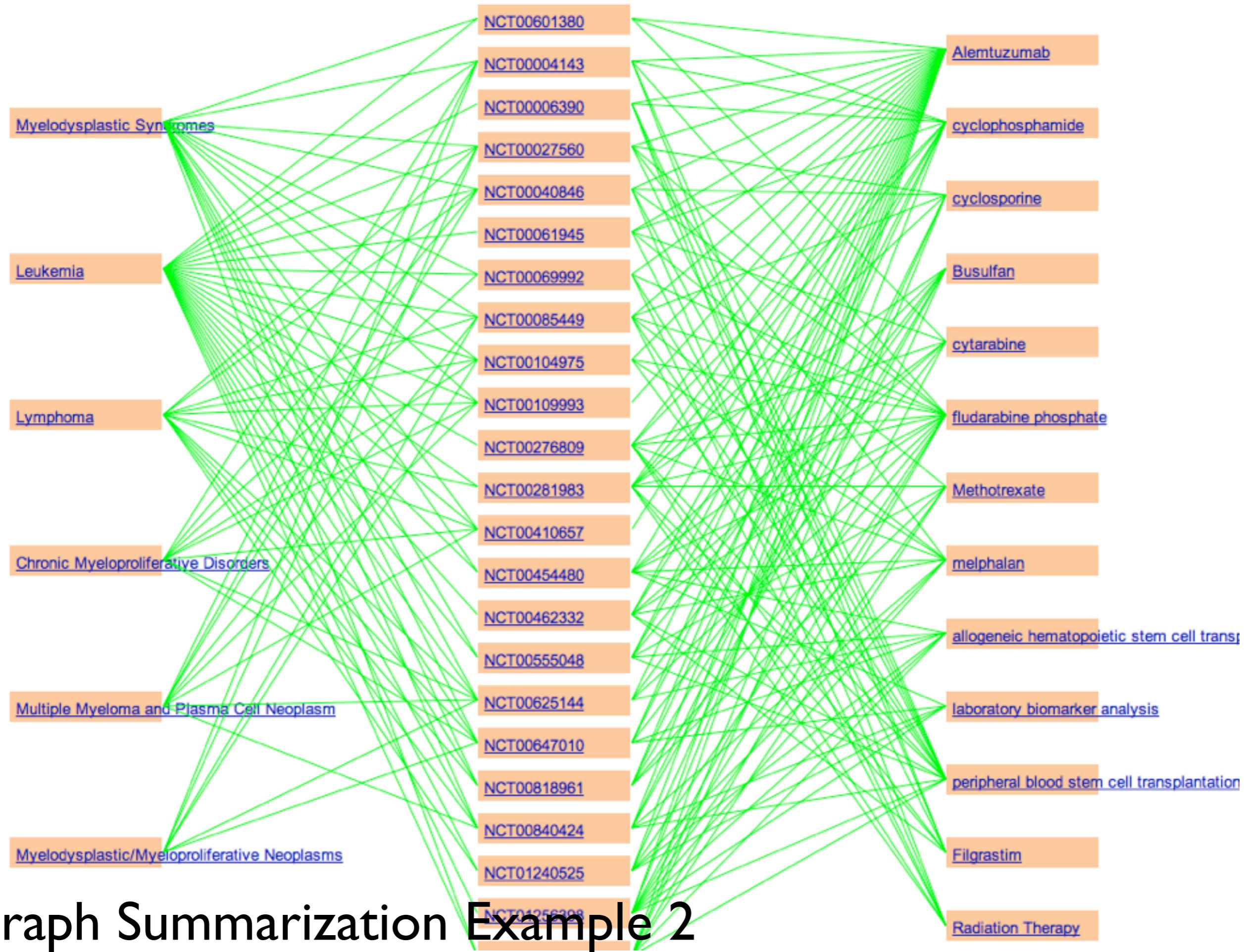
Thor, A., Anderson, P., Raschid, L., Navlakha, S., Saha, B., Khuller, S., Zhang, X.N.: Link prediction for annotation graphs using graph summarization. In: International Semantic Web Conference (ISWC). (2011)

Graph Summarization Example I

Condition

Clinical Trial

Intervention



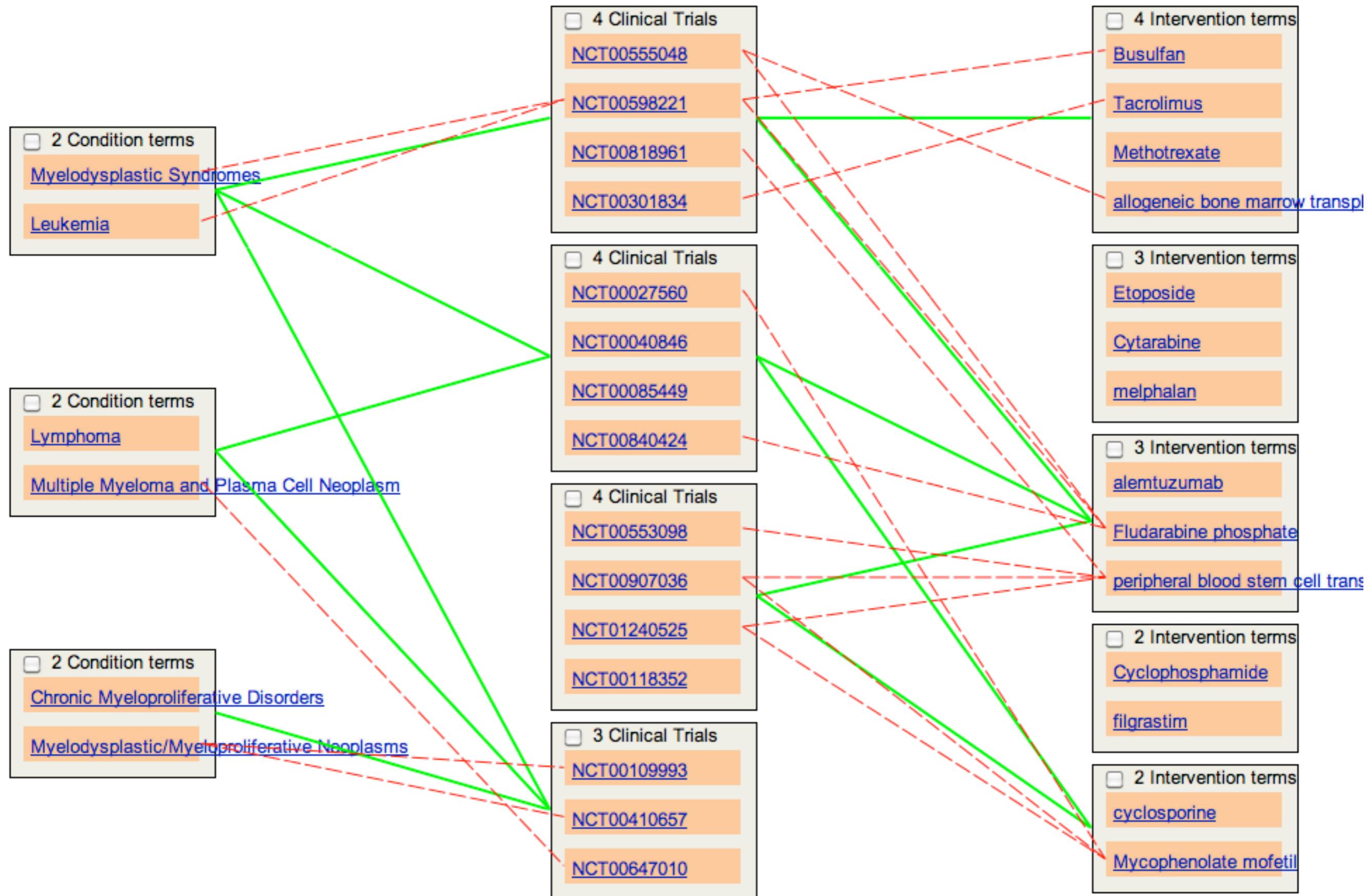
Graph Summarization Example 2

Legend:

Superedge

Deletion

Addition

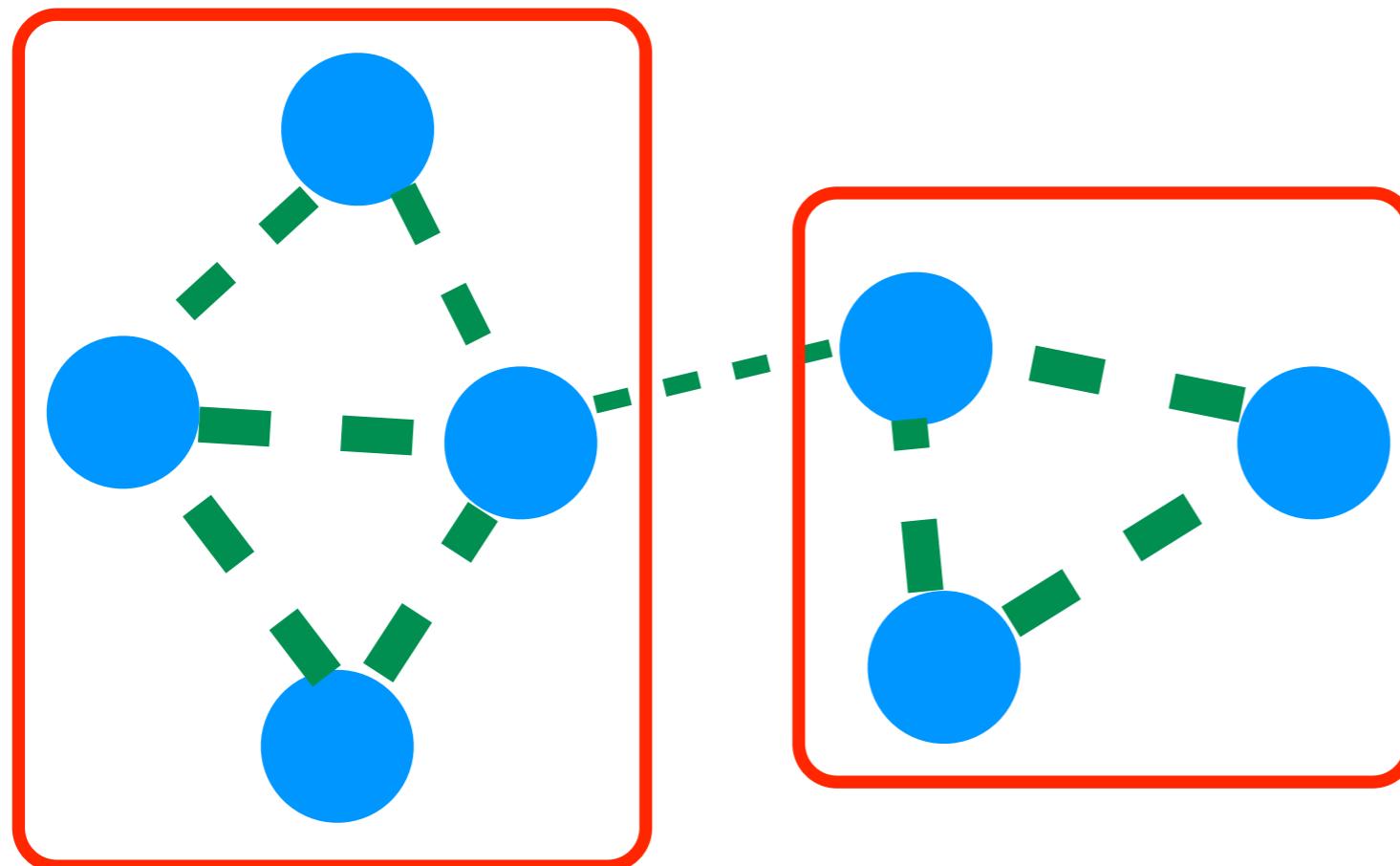


Graph Summarization Example 2

Heuristics for Graph Summarization

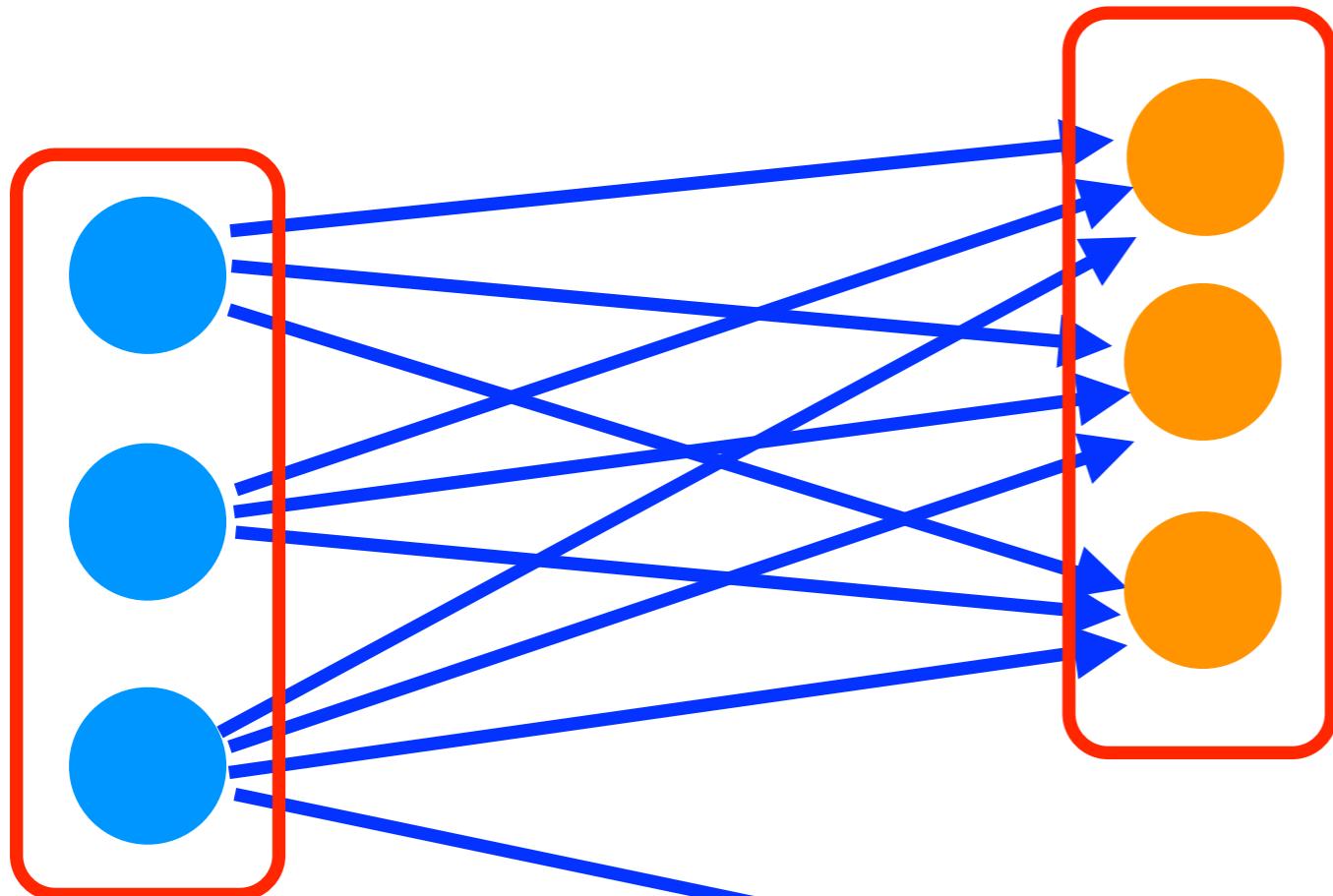
- Similarity
- Annotation Links
- Neighbor Sets

Similarity Heuristic

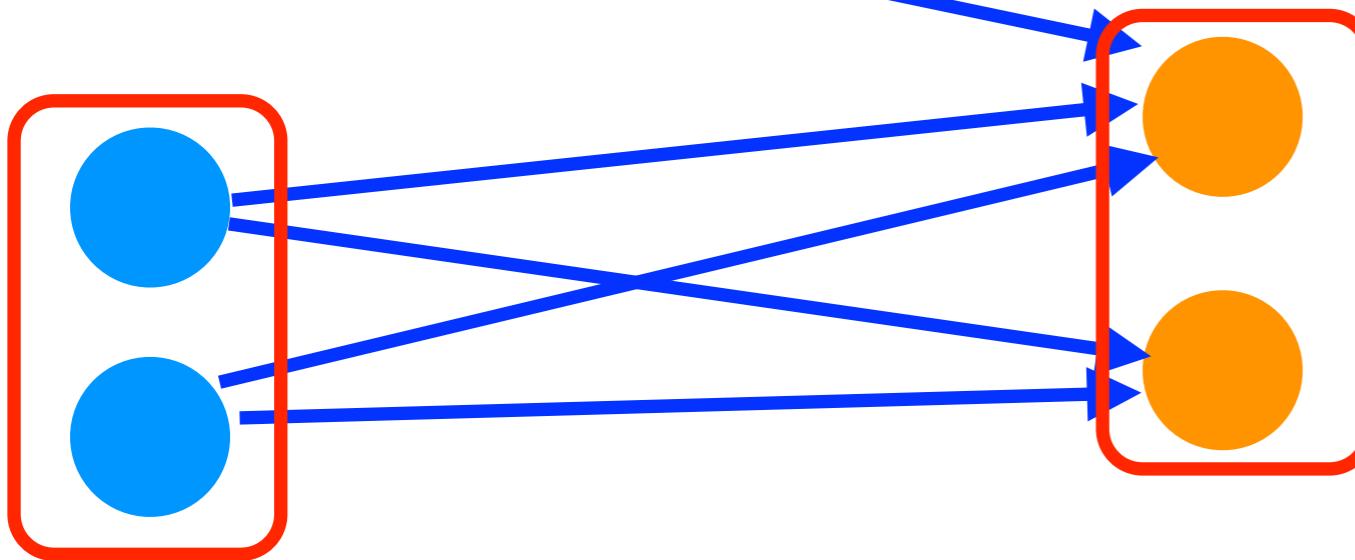


Nodes in same
cluster should
be similar

Annotation Link Heuristic

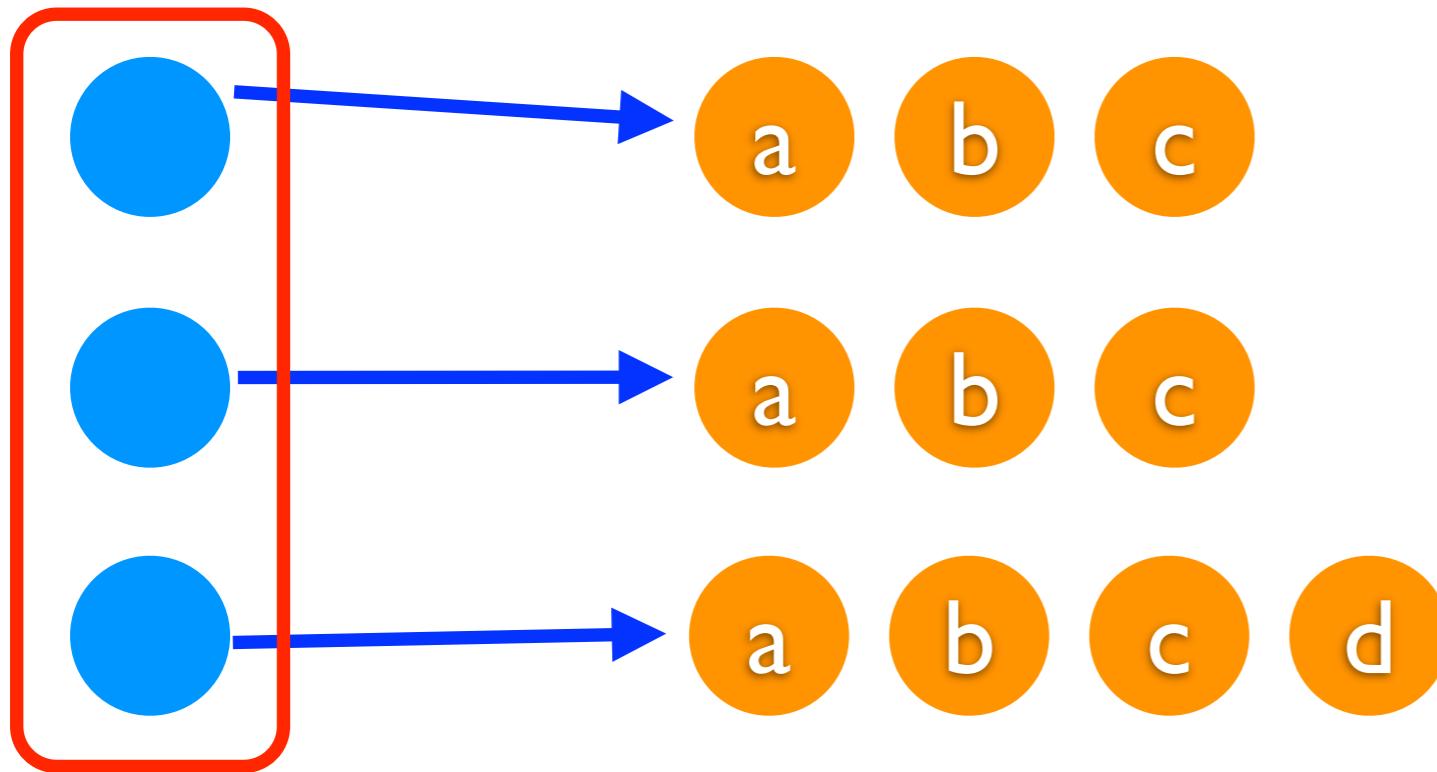


Shared neighbor →
same cluster

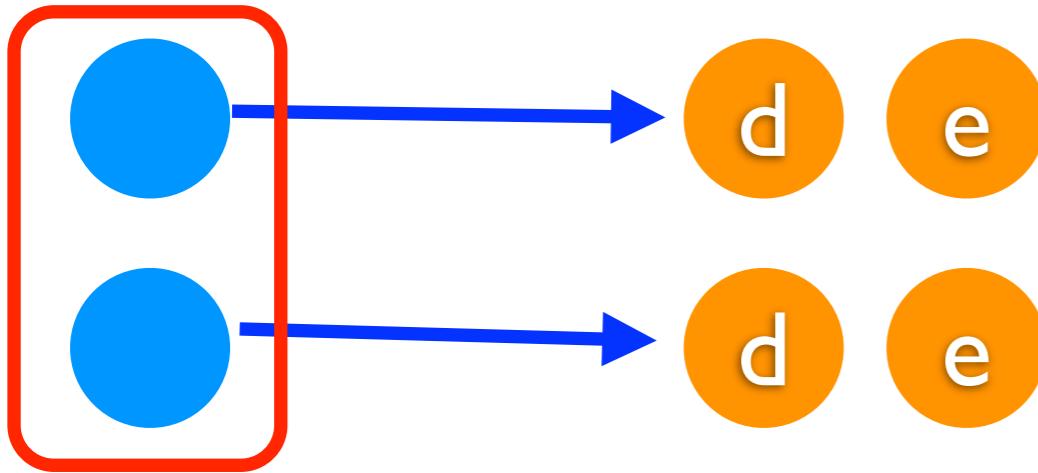


Neighbor not shared
→ different clusters

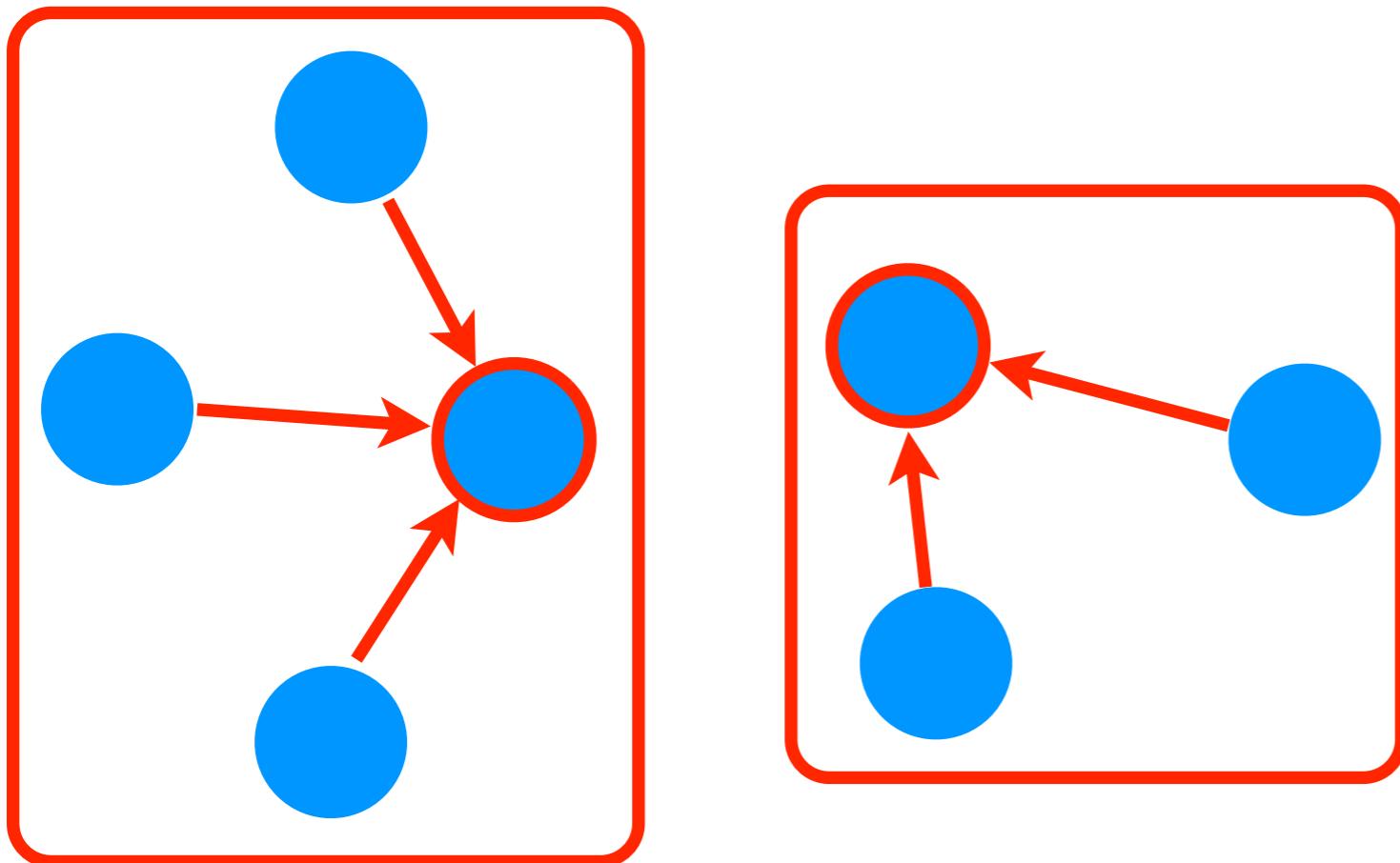
Neighbor Sets Heuristic



Similar set of
neighbors
→ same cluster



Affinity Propagation



- one exemplar per cluster
- every node has to choose an exemplar

Probabilistic Reasoning on Relational Data

- Fuzzy Logic [Lee, 1972]
- Undirected Graphical Models, FOL template language [Richardson et al, 2006]
- Probabilistic Ontologies [Costa et al, 2006]

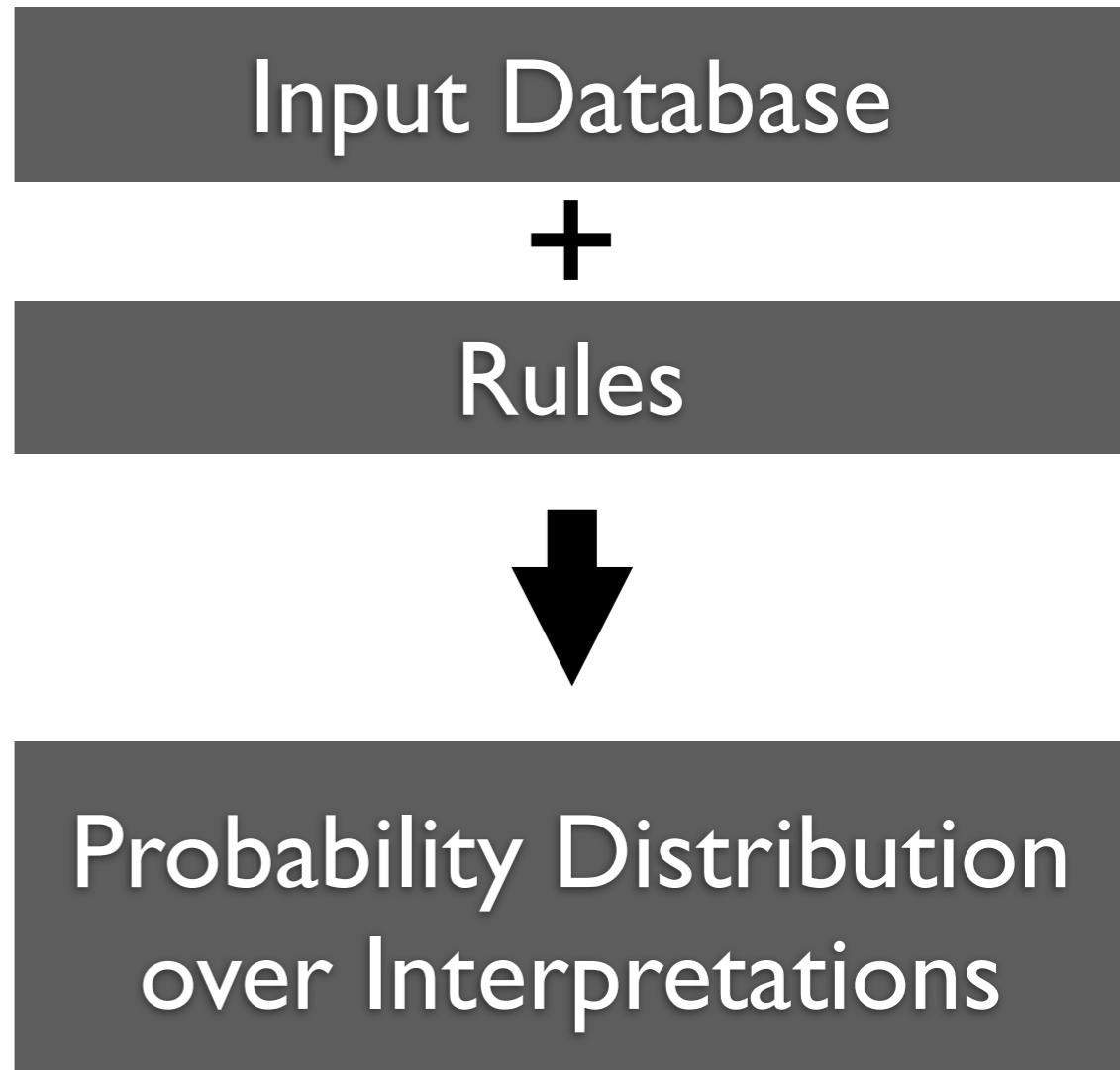
Costa, P. C. G., and K. B. Laskey. "PR-OWL: A Framework for Probabilistic Ontologies." *Frontiers in Artificial Intelligence and Applications* 150 (2006): 237.

Richardson, M., Domingos, P.: Markov logic networks. *Machine Learning* 62(1-2) (2006) 107–136

Lee, R.C.T.:Fuzzy logic and the resolution principle. *J.ACMI* 19(1)(1972)109–119

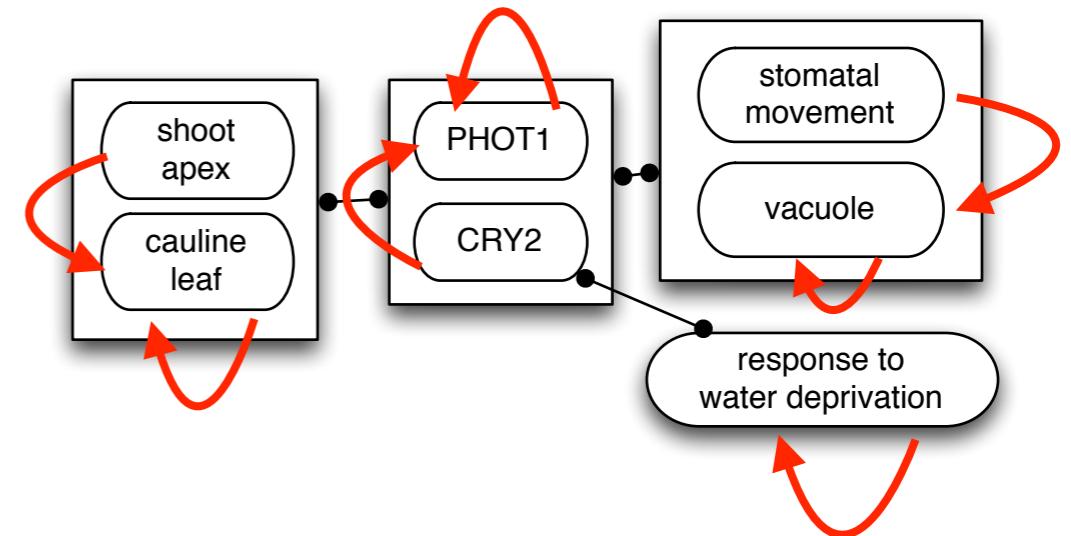
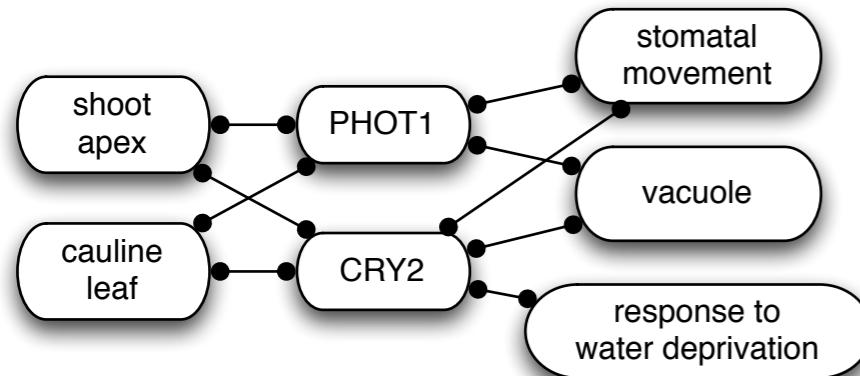
Probabilistic Soft Logic (PSL)

- Declarative language for probabilistic reasoning
 - First order rules
 - Soft truth values
 - Similarities
 - Set similarity
- Efficient inference
- <http://psl.umiacs.umd.edu/>



Broeckeler, M., Mihalkova, L., Getoor, L.: Probabilistic similarity logic. In: Conference on Uncertainty in Artificial Intelligence (UAI). (2010)

Graphs and Clusters in PSL



```

link(sa,p1)= 1
link(sa,c2) = 1
link(cl,p1) = 1
...
link(c2,v) = 1
link(c2,rw) = 1

similar(sa,cl) = .9
similar(p1,c2) = .5
similar(sm,v) = .1
...
  
```

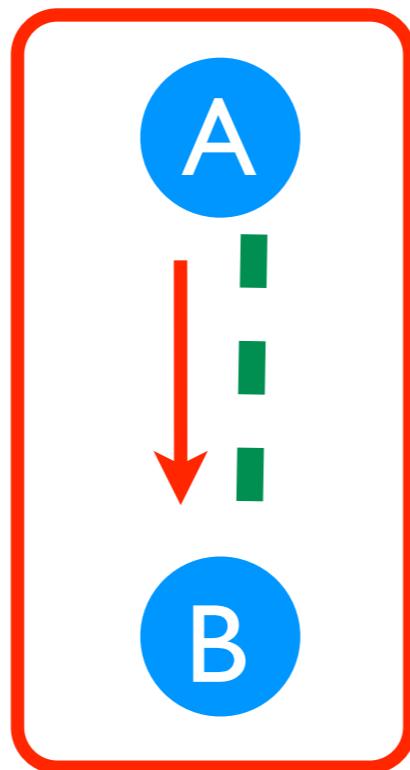
```

exemplar(sa,c1)
exemplar(cl,c1)
exemplar(p1,p1)
exemplar(c2,p1)
exemplar(sm,v)
exemplar(v,v)
exemplar(rq,rw)
  
```

GS Heuristics in PSL

- Similarity Heuristic

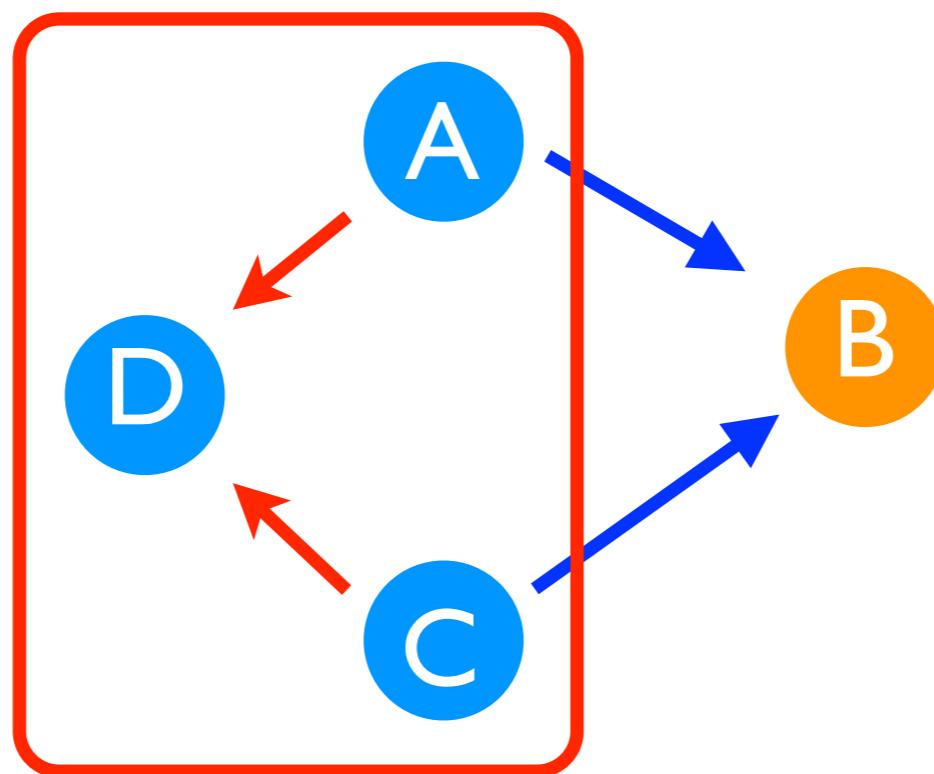
1 : **exemplar(A,B)** -> **similar(A,B)**



GS Heuristics in PSL

- Annotation Link Heuristic

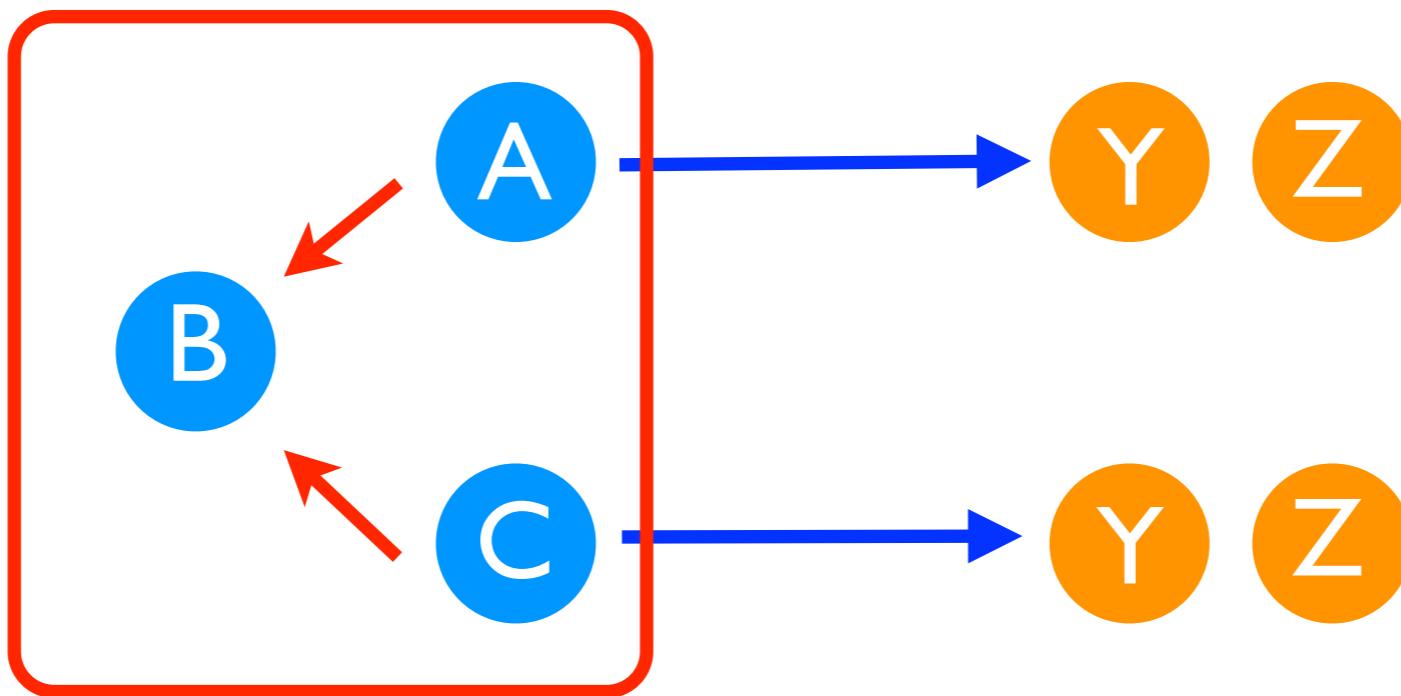
```
1 : link(A,B) & link(C,B) & exemplar(A,D) -> exemplar(C,D)  
1 : link(A,B) & exemplar(A,C) & exemplar(D,C) -> link(C,B)
```



GS Heuristics in PSL

- Neighbor Sets Heuristic

1 : `sameNeighbors({A.link}, {C.link}) & exemplar(A,B)`
 $\rightarrow \text{exemplar}(C,B)$

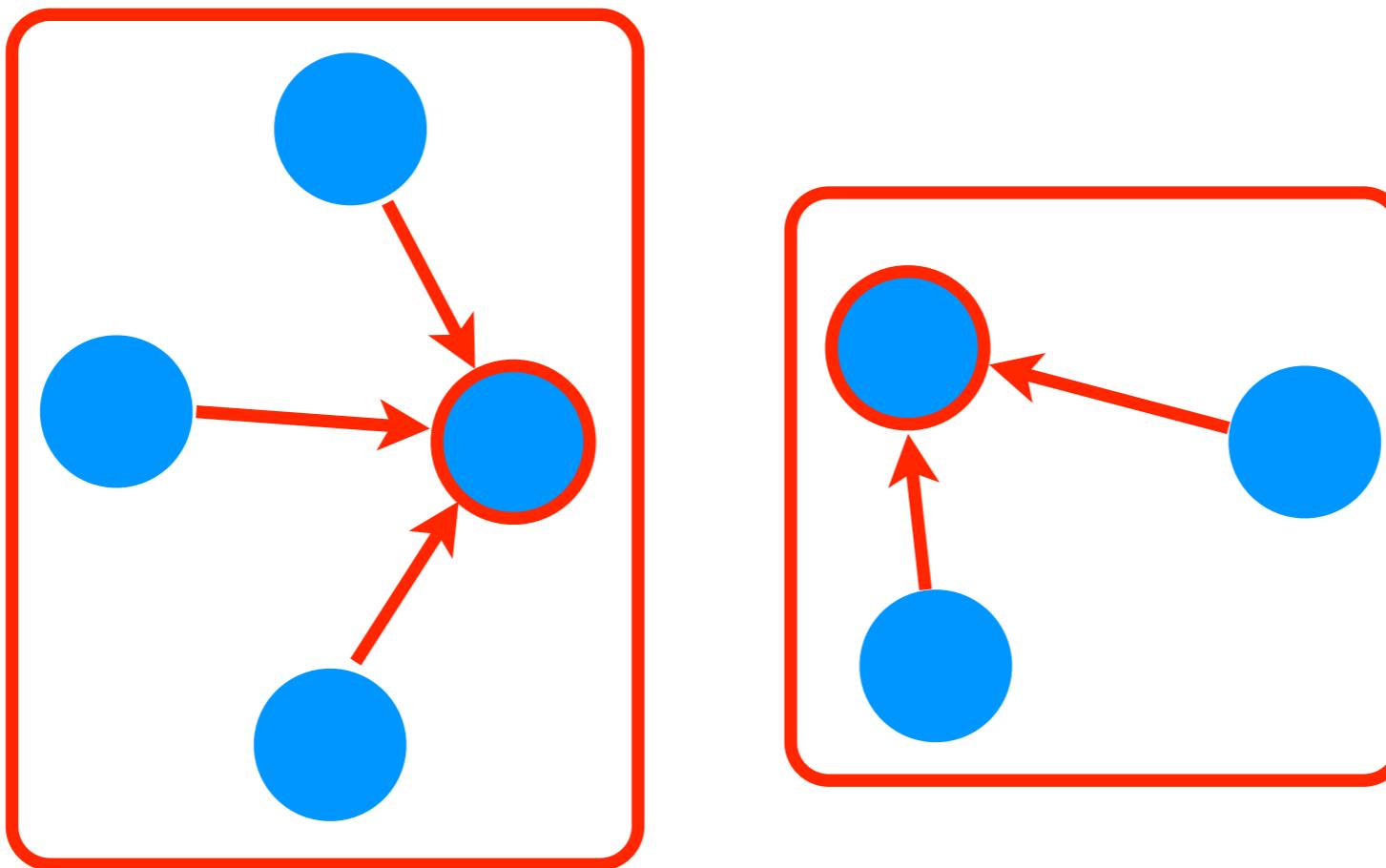


GS Heuristics in PSL

- Affinity Propagation

1000 : **exemplar(A,B)** \rightarrow **exemplar(B,B)**

Functional: **exemplar**



GS Heuristics in PSL

- Similarity Heuristic

1 : **exemplar(A,B)** → **similar(A,B)**

- Annotation Link Heuristic

1 : **link(A,B) & link(C,B) & exemplar(A,D)** → **exemplar(C,D)**

1 : **link(A,B) & exemplar(A,C) & exemplar(D,C)** → **link(D,B)**

- Neighbor Sets Heuristic

1 : **sameNeighbors({A.link}, {C.link}) & exemplar(A,B)**
→ **exemplar(C,B)**

- Affinity Propagation

1000 : **exemplar(A,B)** → **exemplar(B,B)**

Functional: **exemplar**

Distance to Satisfaction

rule satisfied \Leftrightarrow

truth value of body \leq truth value of head

in1(a,b) \rightarrow out1(a,b)		distance to satisfaction
1.0	0.0	1.0
0.9	0.6	0.3
0.7	0.6	0.1
0.5	0.6	0

Distance to Satisfaction

rule satisfied \Leftrightarrow

truth value of body \leq truth value of head

out1(a,b) -> in1(a,b)		distance to satisfaction
1.0	0.0	1.0
0.9	0.6	0.3
0.7	0.6	0.1
0.5	0.6	0

Probabilistic Model

Rule's distance to satisfaction given I

$$f(I) = \frac{1}{Z} \exp\left[- \sum_{r \in R} \lambda_r (d_r(I))^p\right]$$

Interpretation

Normalization constant

Rule's weight

Set of rule groundings

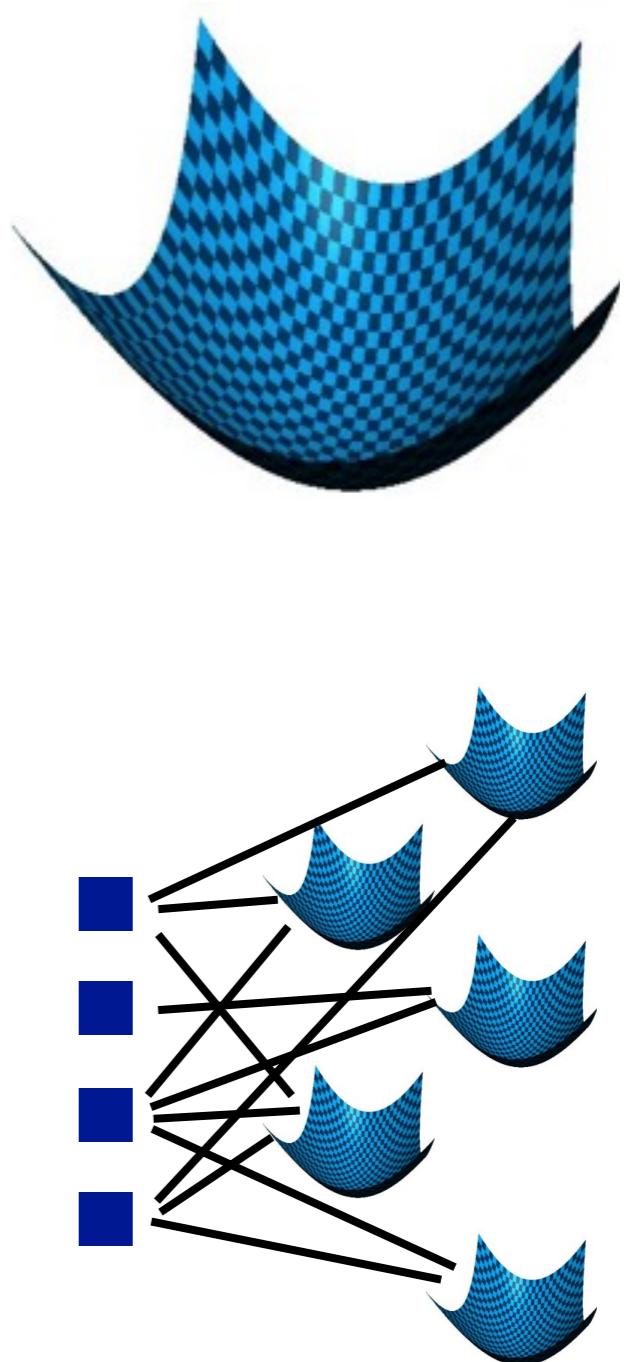
$\in \{1, 2\}$

The diagram illustrates the components of the probabilistic model equation. It shows the function $f(I)$, the normalization constant Z , the rule's weight λ_r , and the set of rule groundings R . Arrows indicate the relationship between the labels and the corresponding terms in the equation.

$$Z = \int_I \exp\left[- \sum_{r \in R} \lambda_r (d_r(I))^p\right]$$

MPE Inference in PSL

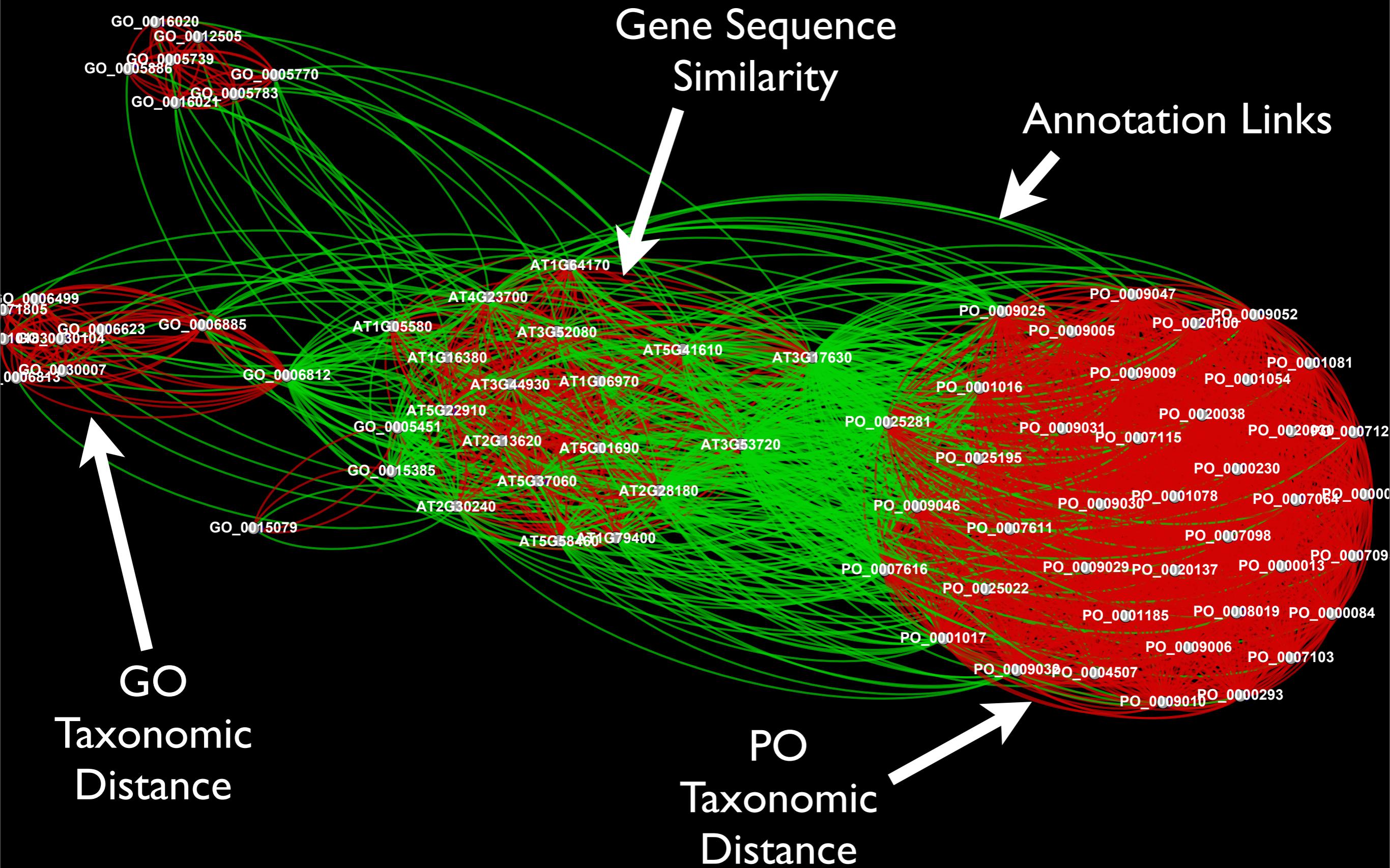
- Convex optimization problem
- Exact inference in polynomial time by transformation to second order cone program
- New solver based on consensus optimization linear time in practice [Bach et al, NIPS 12]



Bach, S., Broeckeler, M., Getoor, L., O'Leary, D.. "Scaling MPE Inference for Constrained Continuous Markov Random Fields with Consensus Optimization" Advances in Neural Information Processing Systems (NIPS) - 2012

Evaluation

Input Graph



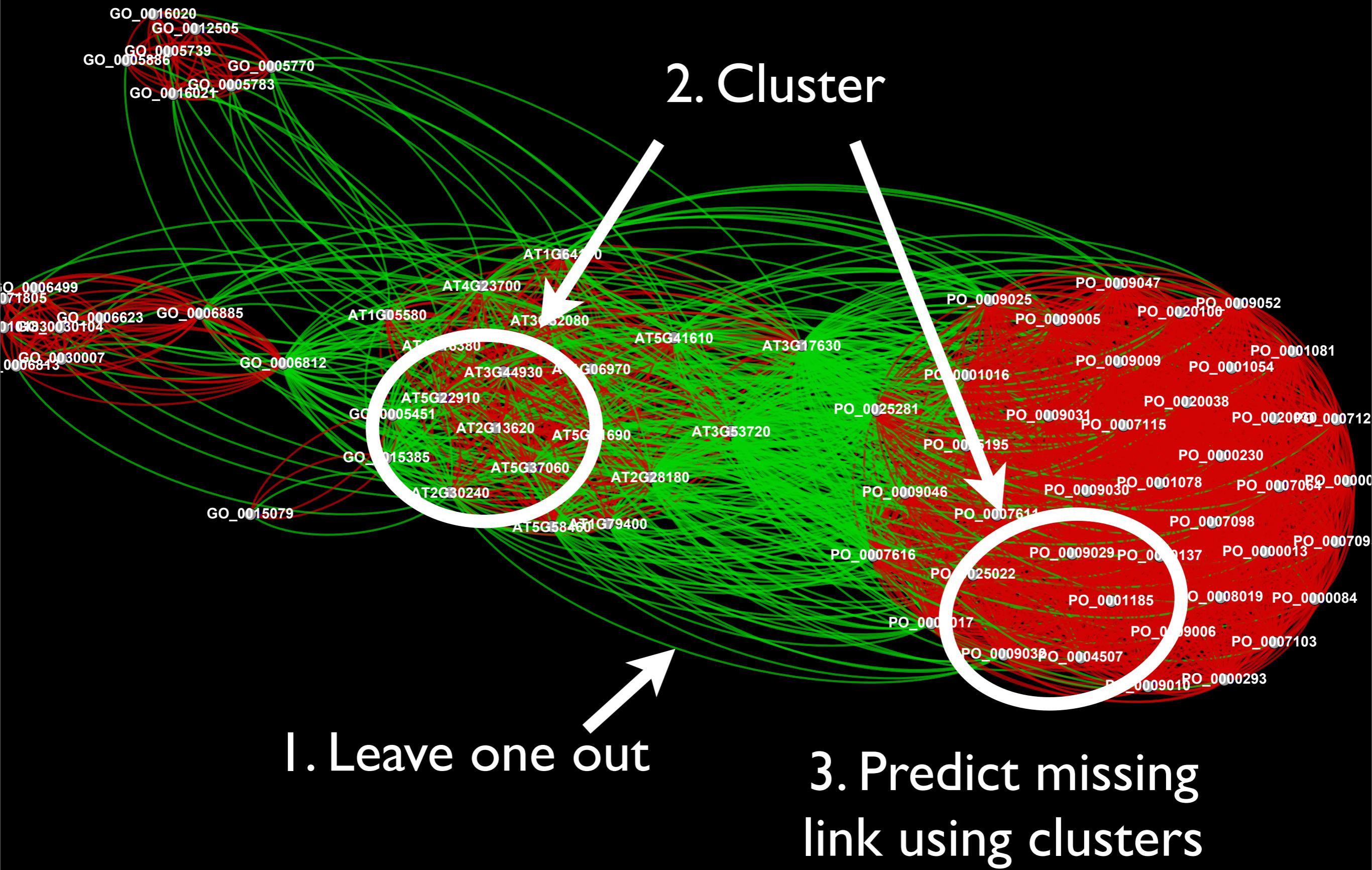
GO
Taxonomic
Distance

Gene Sequence
Similarity

Annotation Links

PO
Taxonomic
Distance

Leave-one-out Link Prediction



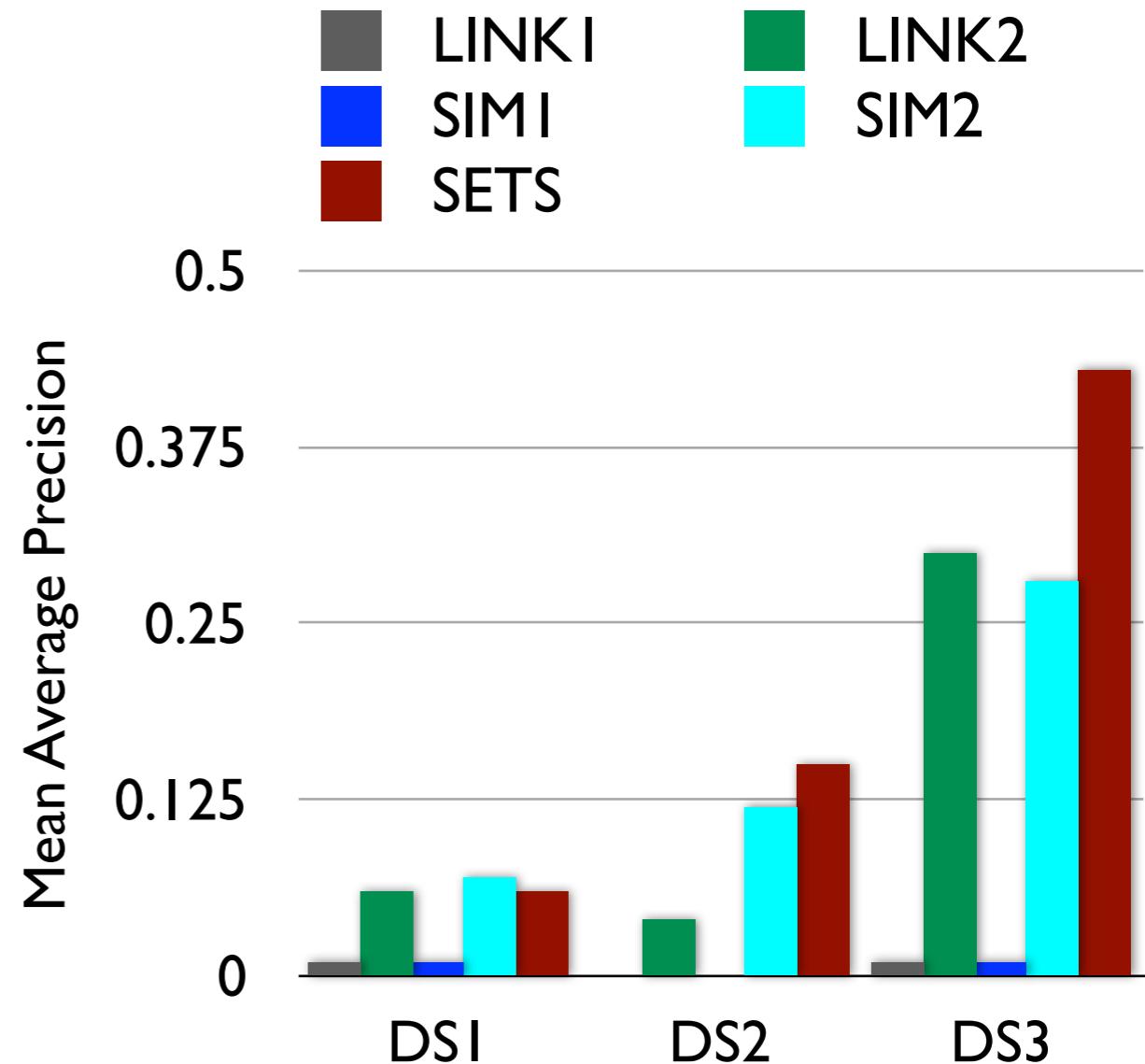
Arabidopsis thaliana

Data Sets

	DS1	DS2	DS3
Genes	10	10	18
PO Terms	53	48	40
GO Terms	44	31	19
PO-Gene	255	255	218
GO-Gene	157	157	92

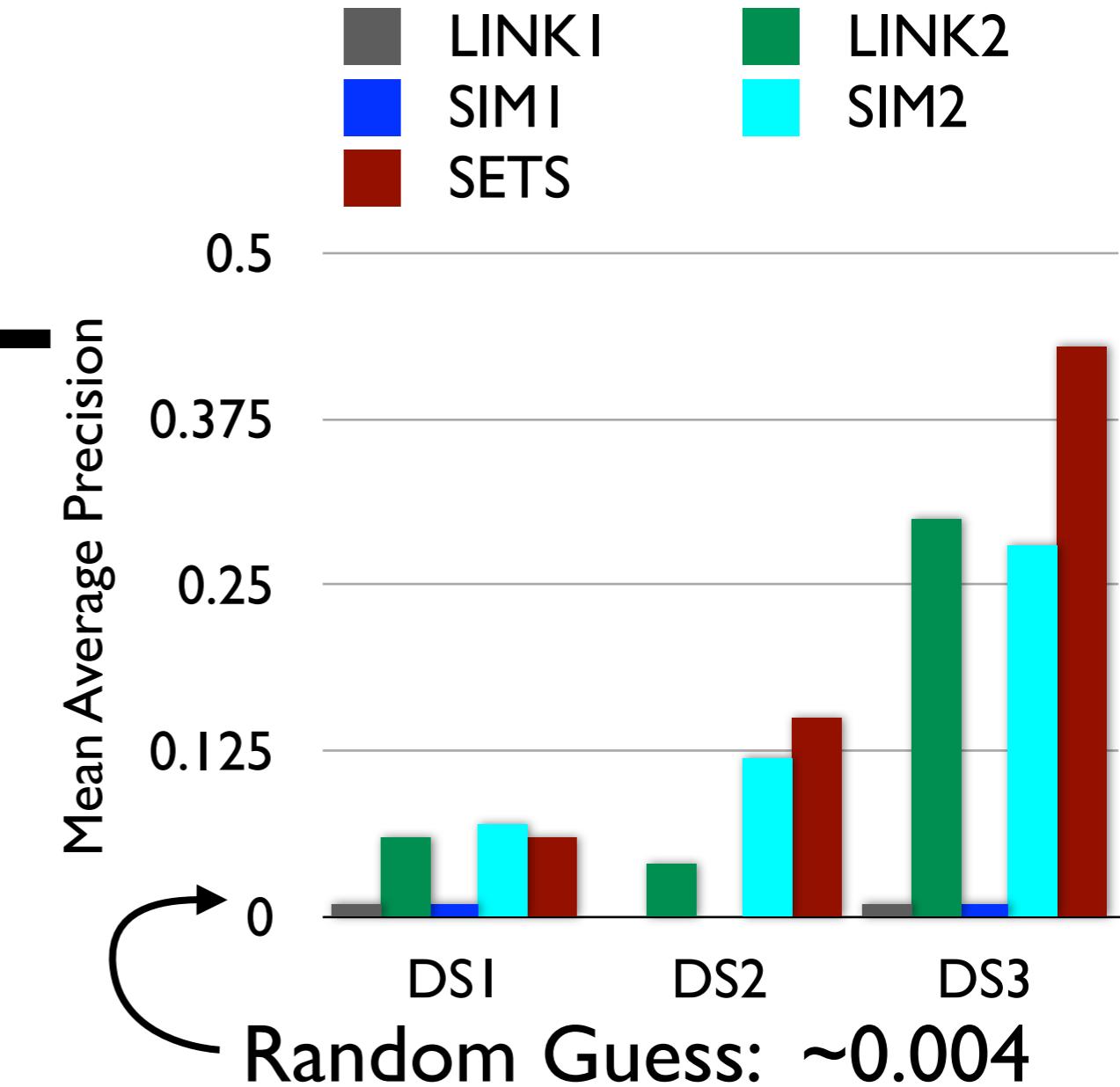
Results from Combining Heuristics

	Cluster	Link Prediction
LINK1	Similarity, Shared Neigh.	Similarity, Shared Neigh.
LINK2	Similarity, Neigh. \neg Shared	Similarity, Neigh. \neg Shared
SIM1	Similarity	Similarity, Shared Neigh.
SIM2	Similarity	Similarity, Neigh. \neg Shared
SETS	Similarity, Neigh. Sets	Similarity, Neigh. \neg Shared



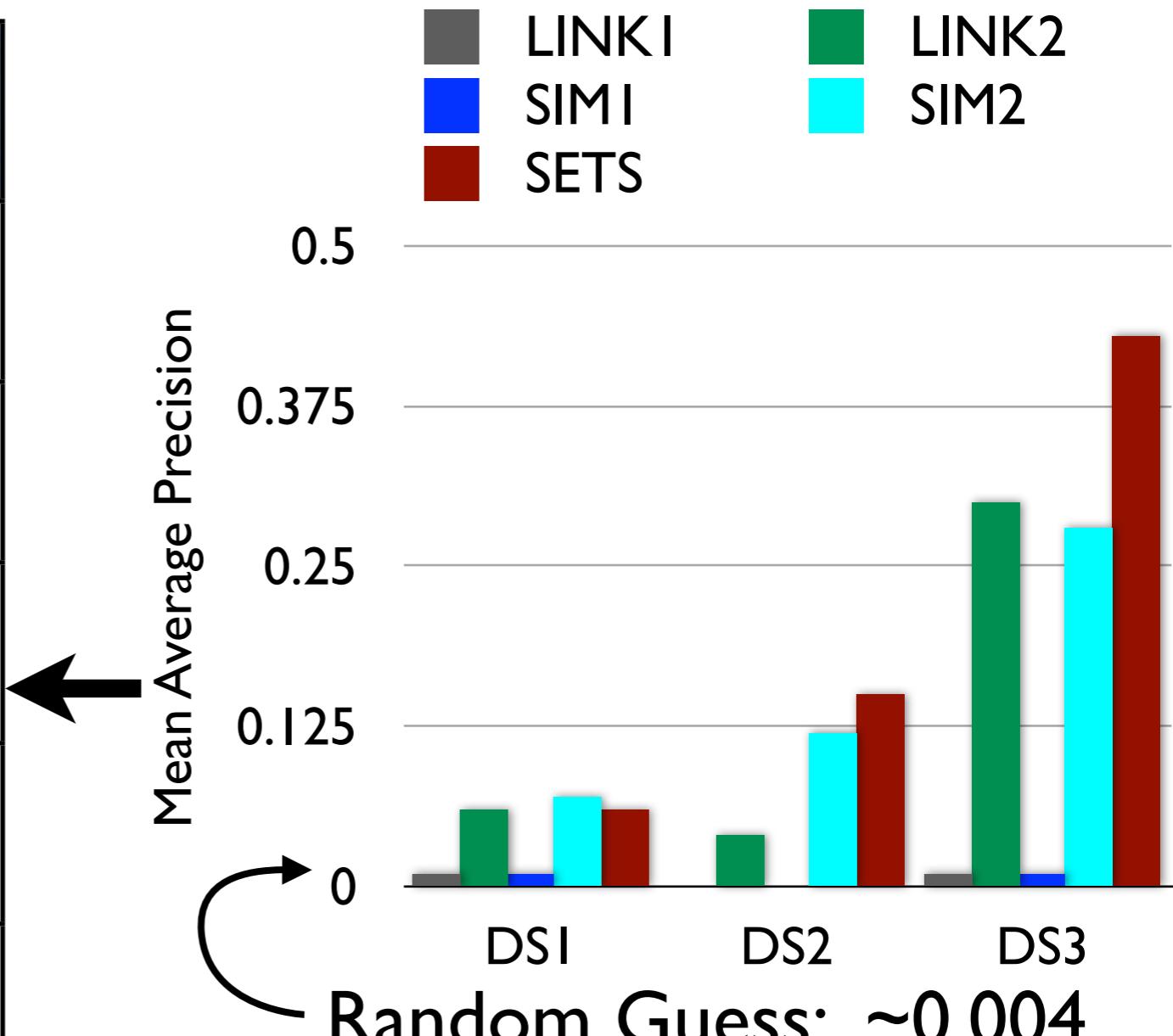
Results from Combining Heuristics

	Cluster	Link Prediction
LINK1	Similarity, Shared Neigh.	Similarity, Shared Neigh.
LINK2	Similarity, Neigh. \neg Shared	Similarity, Neigh. \neg Shared
SIM1	Similarity	Similarity, Shared Neigh.
SIM2	Similarity	Similarity, Neigh. \neg Shared
SETS	Similarity, Neigh. Sets	Similarity, Neigh. \neg Shared



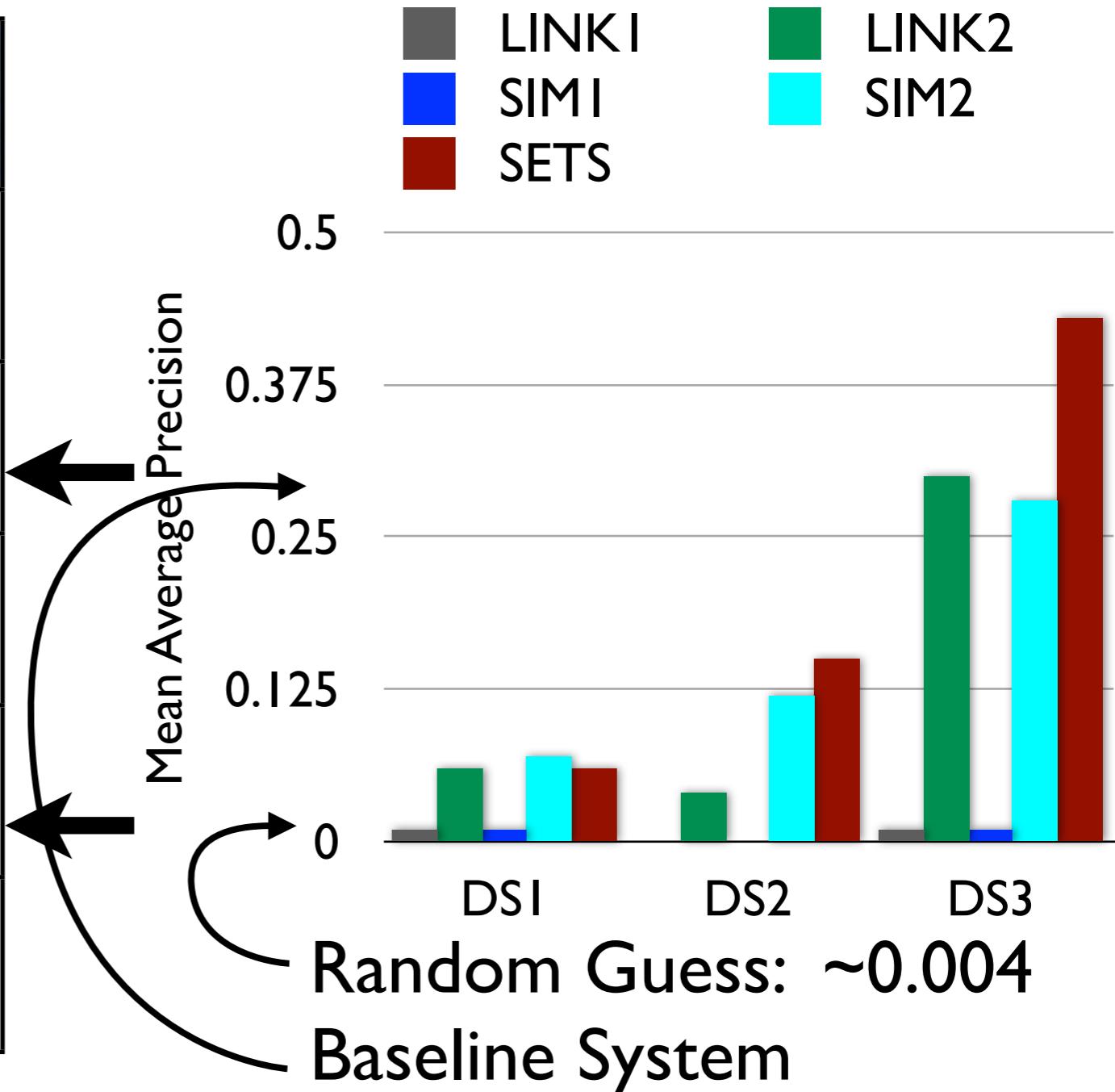
Results from Combining Heuristics

	Cluster	Link Prediction
LINK1	Similarity, Shared Neigh.	Similarity, Shared Neigh.
LINK2	Similarity, Neigh. \neg Shared	Similarity, Neigh. \neg Shared
SIM1	Similarity	Similarity, Shared Neigh.
SIM2	Similarity	Similarity, Neigh. \neg Shared
SETS	Similarity, Neigh. Sets	Similarity, Neigh. \neg Shared



Results from Combining Heuristics

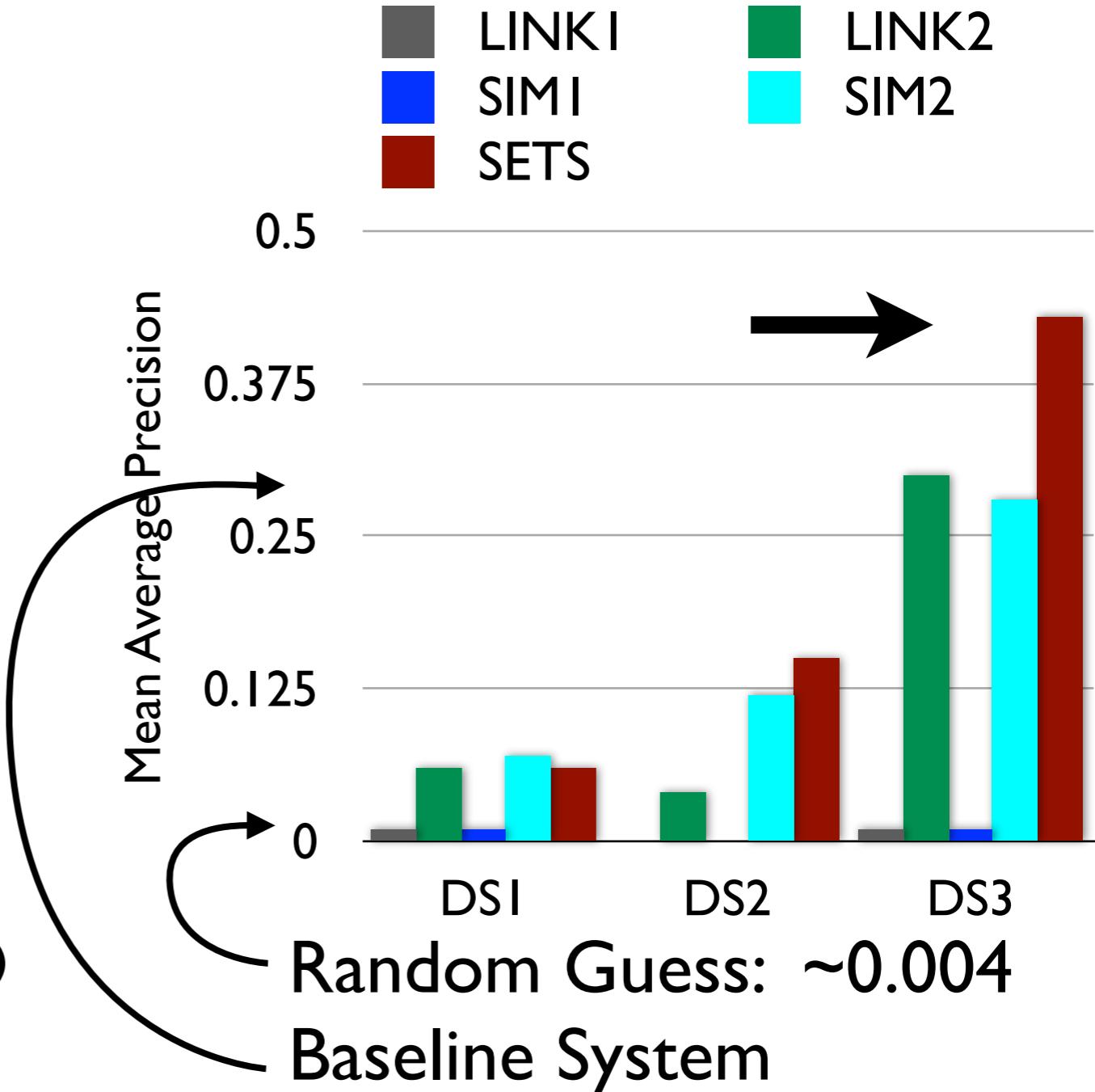
	Cluster	Link Prediction
LINK1	Similarity, Shared Neigh.	Similarity, Shared Neigh.
LINK2	Similarity, Neigh. \neg Shared	Similarity, Neigh. \neg Shared
SIM1	Similarity	Similarity, Shared Neigh.
SIM2	Similarity	Similarity, Neigh. \neg Shared
SETS	Similarity, Neigh. Sets	Similarity, Neigh. \neg Shared



Thor, A., Anderson, P., Raschid, L., Navlakha, S., Saha, B., Khuller, S., Zhang, X.N.: Link prediction for annotation graphs using graph summarization. In: International Semantic Web Conference (ISWC). (2011)

Results from Combining Heuristics

	Cluster	Link Prediction
LINK1	Similarity, Shared Neigh.	Similarity, Shared Neigh.
LINK2	Similarity, Neigh. \neg Shared	Similarity, Neigh. \neg Shared
SIM1	Similarity	Similarity, Shared Neigh.
SIM2	Similarity	Similarity, Neigh. \neg Shared
SETS	Similarity, Neigh. Sets	Similarity, Neigh. \neg Shared



Thor, A., Anderson, P., Raschid, L., Navlakha, S., Saha, B., Khuller, S., Zhang, X.N.: Link prediction for annotation graphs using graph summarization. In: International Semantic Web Conference (ISWC). (2011)

Conclusion

- Simple heuristics → competitive results
- Probabilistic Soft Logic for general reasoning over uncertainty

Thank you!

<http://psl.umiacs.umd.edu>

Backup

Distance to Satisfaction

$$d_r(I) = \max\{0, I(r_{body}) - I(r_{head})\}$$

using Lukasiewicz norms:

$$I(v_1 \wedge v_2) = \max\{0, I(v_1) + I(v_2) - 1\}$$

$$I(v_1 \vee v_2) = \min\{I(v_1) + I(v_2), 1\}$$

$$I(\neg l_1) = 1 - I(l_1)$$

1 : link(a,b) & exemplar(a,a) & exemplar(c,a) -> link(c,b)

1 0.9 0.8 0

$$\max\{0, 1+0.9+0.8-2\}=0.7 \quad d=0.7$$

1 0.5 0.3 0

$$\max\{0, 1+0.5+0.3-2\}=0 \quad d=0$$