

KNOWLEDGE GRAPH CONSTRUCTION

Jay Pujara

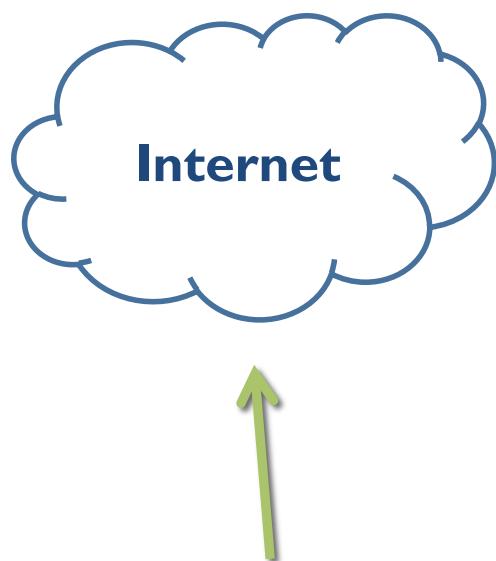
University of Maryland, College Park

Max Planck Institute

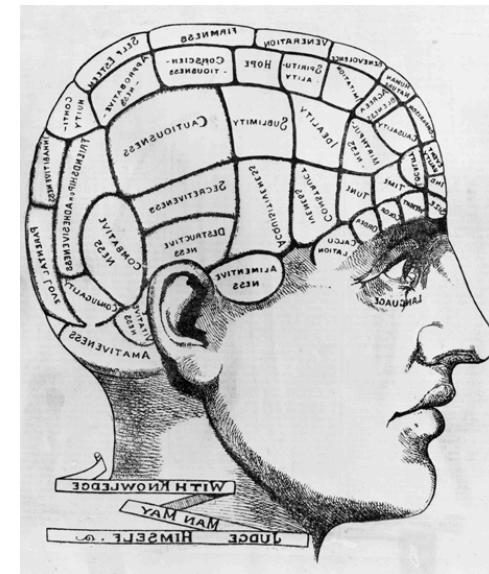
7/9/2015



Can Computers Create Knowledge?



Massive source of
publicly available
information



Knowledge

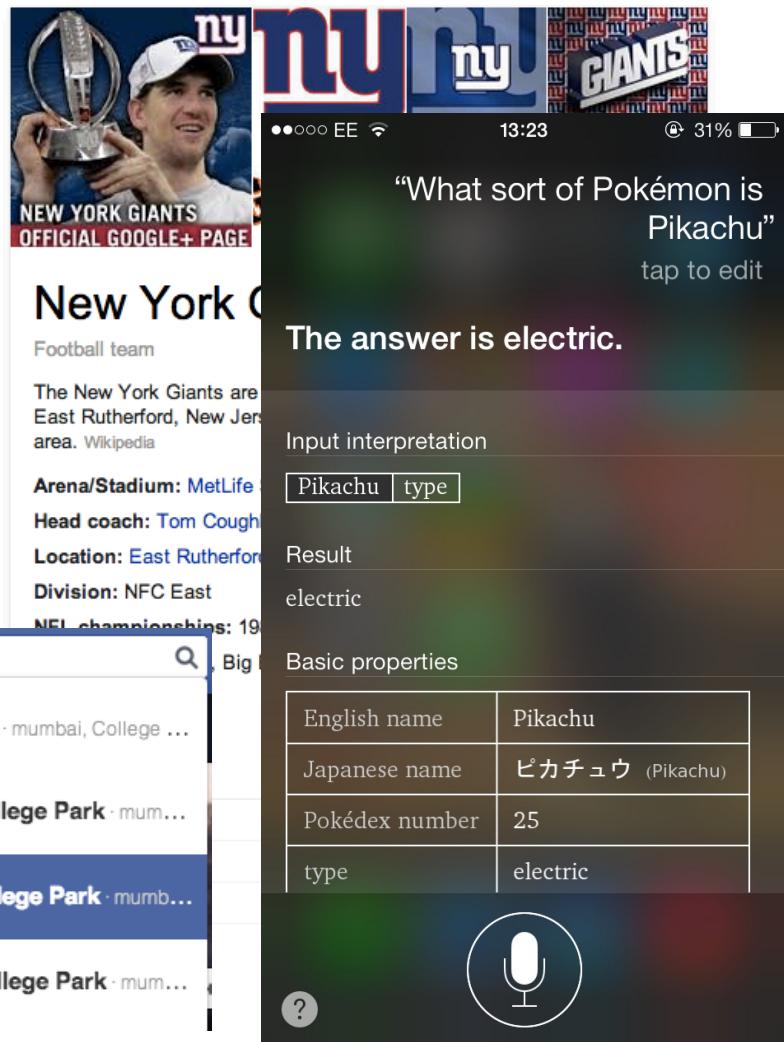
Computers + Knowledge = ❤

New York Giants
4-6, 3rd in NFC Eastern Division

Yesterday, 4:25 PM (ET)
MetLife Stadium, East Rutherford, New Jersey

| | | |
|---|-----------|--|
|  Green Bay Packers (5-5) | 13 - 27 |  New York Giants (4-6) |
| | Final | Total |
| Packers | 1 0 6 0 7 | 13 |
| Giants | 7 3 10 7 | 27 |

Sun, Nov 24 vs.  Cowboys 4:25 PM (ET)



The smartphone screen displays a search result for "What sort of Pokémon is Pikachu". The result is "The answer is electric." Below this, it shows the input interpretation "Pikachu type" and the result "electric". A table titled "Basic properties" provides detailed information about Pikachu:

| | |
|----------------|-----------------|
| English name | Pikachu |
| Japanese name | ピカチュウ (Pikachu) |
| Pokédex number | 25 |
| type | electric |

News for Giants

 People I know who studied at University of Maryland, College Park

 People I know who studied at **University of Maryland, College Park** · mumbai, College ...

 Friends of people I know who studied at **University of Maryland, College Park** · mum...
 Photos of people I know who studied at **University of Maryland, College Park** · mumb...

 Photos by people I know who studied at **University of Maryland, College Park** · mum...

What does it mean to create knowledge?
What do we mean by knowledge?

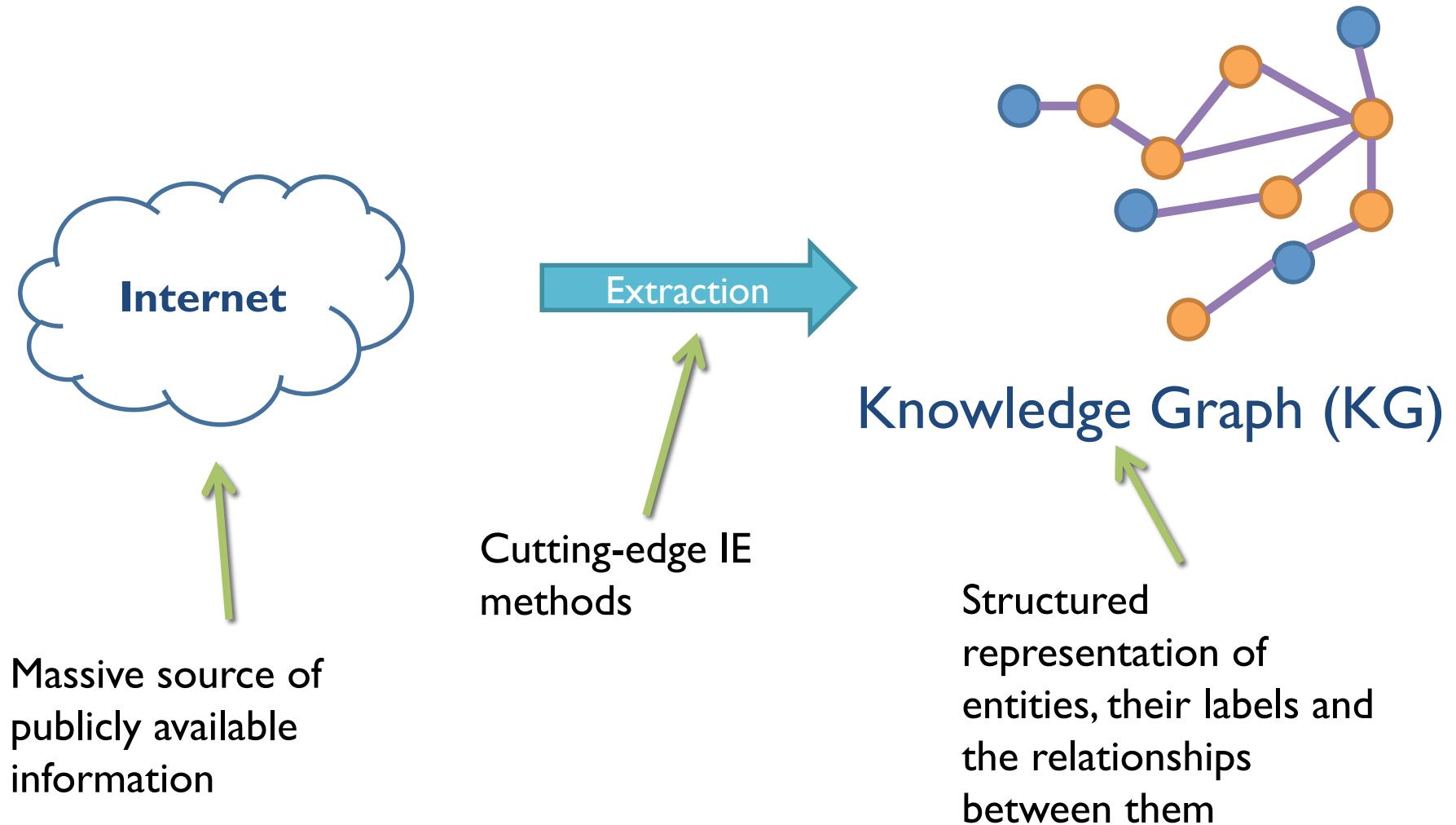
Defining the Questions

- Extraction
- Representation
- Reasoning and Inference

Defining the Questions

- Extraction
- Representation
- **Reasoning and Inference**

A Revised Knowledge-Creation Diagram



Knowledge Graphs in the wild

| New York Giants | |
|--|---------------------------|
| 4-6, 3rd in NFC Eastern Division | |
| Yesterday, 4:25 PM (ET) | |
| MetLife Stadium, East Rutherford, New Jersey | |
|  Green Bay Packers (5-5) | 13 - 27 Final |
|  New York Giants (4-6) | |
| | 1 2 3 4 Total |
| Packers | 0 6 0 7 13 |
| Giants | 7 3 10 7 27 |
| Sun, Nov 24 vs.  Cowboys | 4:25 PM (ET) |

News for Giants



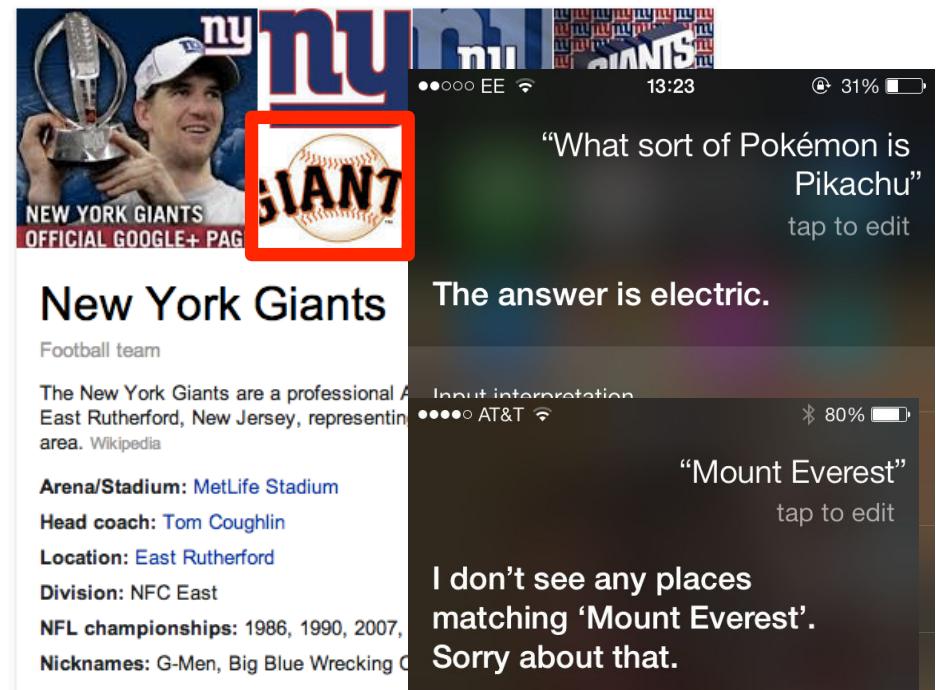
Tim Hudson San Francisco **Giants** close in on deal
USA TODAY - by Jorge Ortiz - 17 minutes ago
The San Francisco **Giants**, determined to bolster a once-proud rotation that faltered in 2013, a previous ...

Newsday

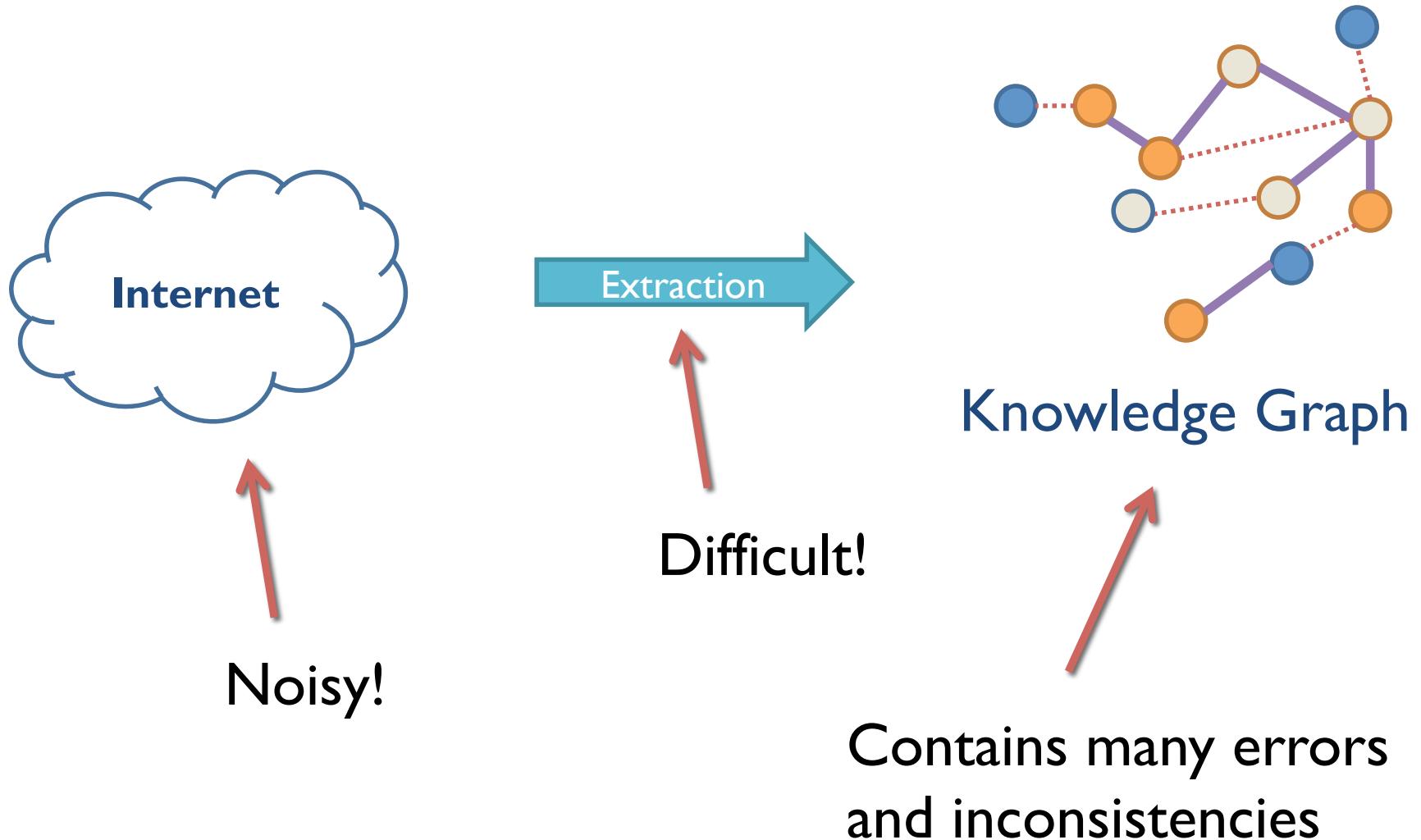
Facebook icon

People I know who s

-  People I know who studied at University of Maryland, College Park
-  People I know who studied at **University of Maryland, College Park** · mumbai, College ...
-  Friends of people I know who studied at **University of Maryland, College Park** · mum...
sc
-  Photos of people I know who studied at **University of Maryland, College Park** · mumb...
-  Photos by people I know who studied at **University of Maryland, College Park** · mum...
sc



Motivating Problem: Real Challenges



NELL: The Never-Ending Language Learner

- Large-scale IE project
(Carlson et al., AAAI10)

- Lifelong learning: aims to
“read the web”

- Ontology of known
labels and relations

- Knowledge base
contains millions of facts

- person
 - monarch
 - astronaut
 - personbylocation
 - personnorthamerica
 - personcanada
 - personus
 - politicanus
 - personmexico
 - personeurope
 - personaaustralia
 - personafrica
 - personsouthamerica
 - personasia
 - personantarctica
 - visualartist
 - model
 - scientist
 - journalist
 - female
 - actor
 - professor
 - director
 - architect
 - politician
 - politicanus
 - musician
 - athlete
 - chef
 - male
 - writer
 - ceo
 - judge
 - mlauthor
 - coach
 - celebrity
 - comedian
 - criminal

Examples of NELL errors

Entity co-reference errors

Kyrgyzstan has many variants:

- Kyrgystan
- Kyrgistan
- Kyrgyzstan
- Kyrgzstan
- Kyrgyz Republic

Saudi Cultural Days in the **Kyrgyz Republic** has concluded its activities in the capital Bishkek in the weekend in a special ceremony held on this occasion. The event was attended by Deputy Minister of Culture and Tourism of the **Kyrgyz Republic** Koulev Mirza; Kyrgyzstan's Ambassador to Saudi Arabia Jusupbek Sharipov; the Saudi Embassy Acting Chargé d'affaires to Kyrgyzstan, Mari bin Barakah Al-Derbas and members of the embassy staff, in the presence of a heavy turnout of Kyrgyz citizens.

The Days of Culture of Saudi Arabia in **Kyrgyzstan** will be held from 6 to 9 May.

Refugees are often from areas where conflict is historically embedded and marked in ideology and injustice. The Tsarnaev family emigrated from the Chechen diaspora in **Kyrgzstan**, a region Stalin deported the Chechens to in 1943. After the fall of the Berlin Wall in 1991, Chechens engaged in a battle for independence from Russia that led to the Tsarnaevs' petition for refugee status in the early

[Home](#) > [Holiday Destinations](#) > [Kyrghyzstan](#) > [Bishkek](#) > Climate Profile



Fast Forecast

Holiday Weather

Missing and spurious labels

[Erik Kleyheeg](#) has just returned from Lesvos with some new bird images. Included here are: [Common Scops-Owl](#), [Wood Warbler](#), [Spanish Sparrow](#), [Red-throated Pipit](#), [Eurasian Chiff-chaff](#), and [Cretzschmar's Bunting](#).

[Anssi Kullberg](#) has sent along some great trip reports to unusual places, including [Kyrgyzstan](#), [Pakistan](#),

Kyrgyzstan is labeled a bird and a country

Kyrgyzstan (/kɜːrgɪ'sta:n/ *kur-gi-STAHN*;^[5] Kyrgyz: Кыргызстан (IPA: [qurⱥs'tan]); Russian: Киргизия), officially the **Kyrgyz Republic** (Kyrgyz: Кыргыз Республикасы; Russian: Кыргызская Республика), is a country located in Central Asia.^[6] Landlocked and mountainous, Kyrgyzstan is bordered by Kazakhstan to the north, Uzbekistan to the west, Tajikistan to the southwest and China to the east. Its capital and largest city is Bishkek.

Missing and spurious relations

Guidance

Kazakhstan / Kyrgyzstan – Consular Fees

Organisation: Foreign & Commonwealth Office
Page history: Published 4 April 2013

Kyrgyzstan's location is ambiguous – Kazakhstan, Russia and US are included in possible locations

Kyrgyzstan U.S. Air Base Future Unclear

A Central Asian country of incredible natural beauty and proud nomadic traditions, most of Kyrgyzstan was formally annexed to Russia in 1876. The Kyrgyz staged a major revolt against the Tsarist Empire in 1916 in which almost one-sixth of the Kyrgyz population was killed. Kyrgyzstan became a Soviet republic in 1936 and

Violations of ontological knowledge

- Equivalence of co-referent entities (`sameAs`)
 - `SameEntity(Kyrgyzstan, Kyrgyz Republic)`
- Mutual exclusion (`disjointWith`) of labels
 - `MUT(bird, country)`
- Selectional preferences (domain/range) of relations
 - `RNG(countryLocation, continent)`

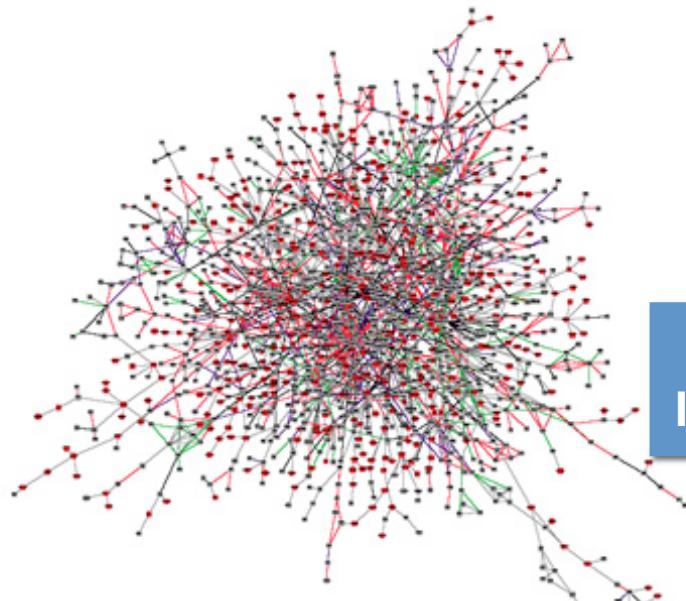
Enforcing these constraints requires **jointly** considering multiple extractions *across documents*

Examples where joint models have succeeded

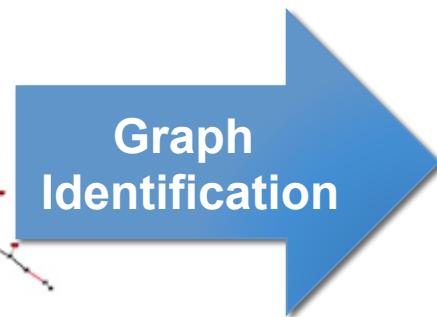
- Information extraction
 - ER+Segmentation: Poon & Domingos, AAAI07
 - SRL: Srikumar & Roth, EMNLP11
 - Within-doc extraction: Singh et al., AKBC13
- Social and communication networks
 - Fusion: Eldardiry & Neville, MLG10
 - EMailActs: Carvalho & Cohen, SIGIR05
 - GraphID: Namata et al., KDD11

GRAPH IDENTIFICATION

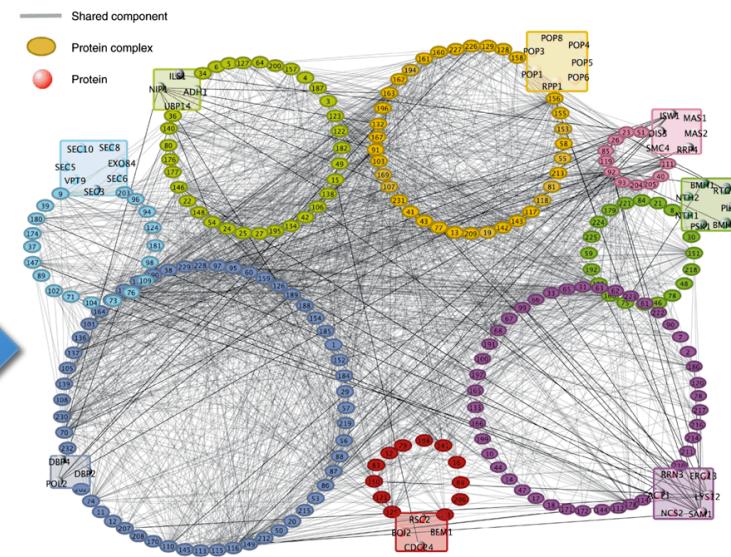
Transformation



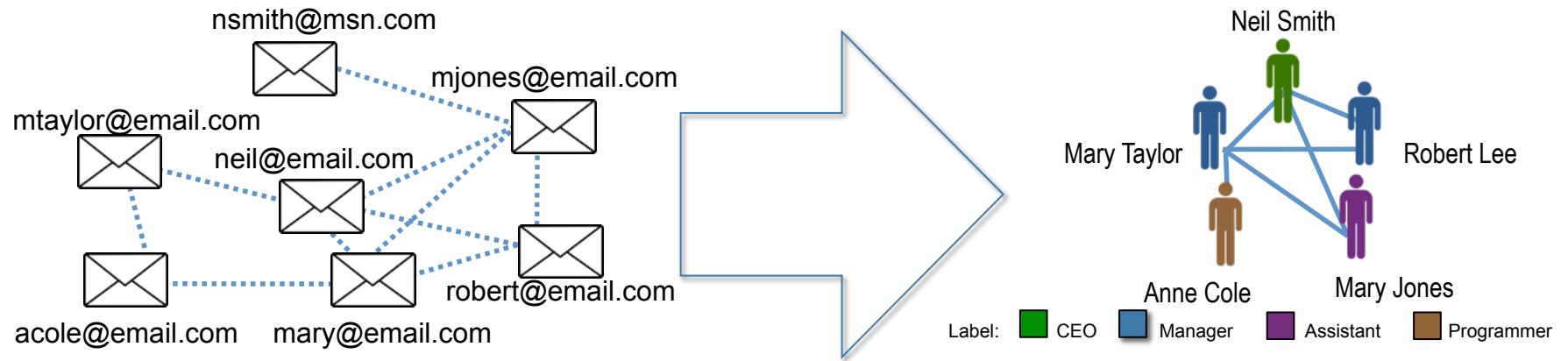
Input Graph



Graph
Identification



Motivation: Different Networks



Communication Network

Nodes: Email Address

Edges: Communication

Node Attributes: Words

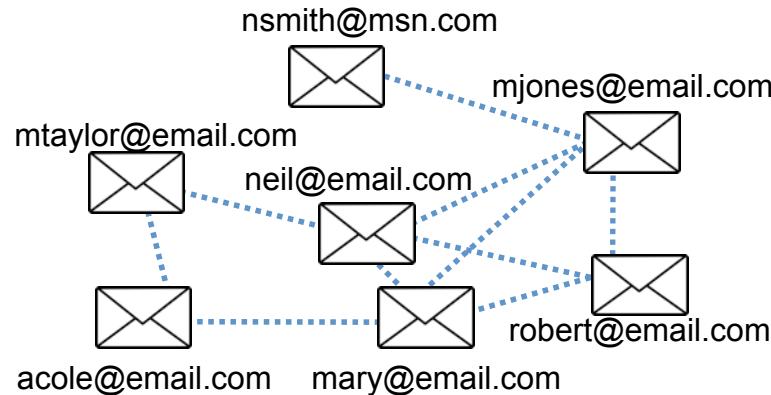
Organizational Network

Nodes: Person

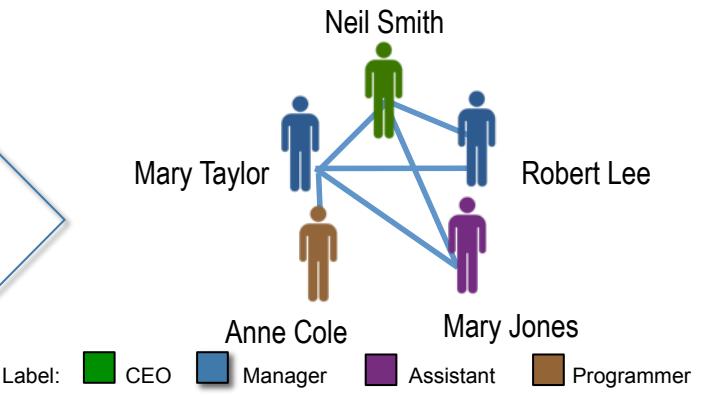
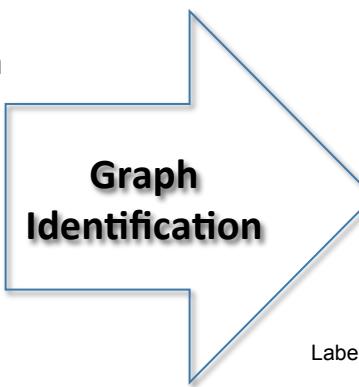
Edges: Manages

Node Labels: Title

Graph Identification

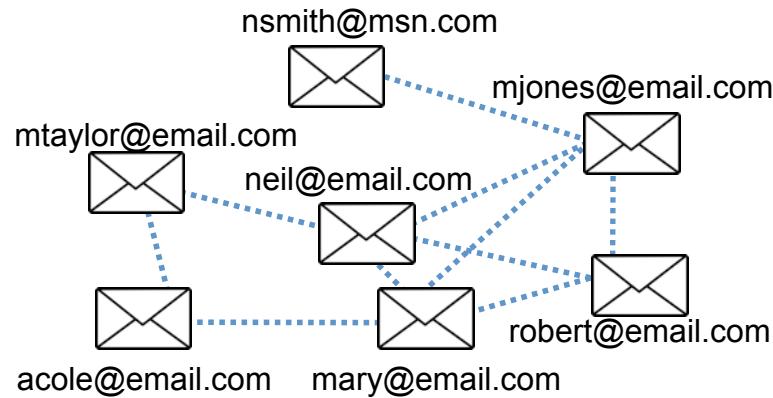


Input Graph: Email Communication Network

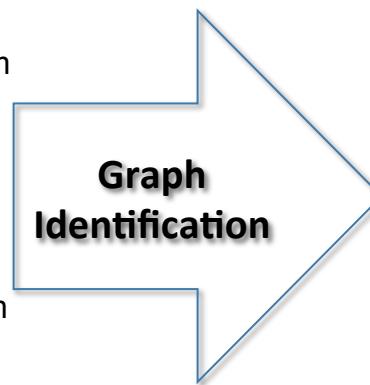


Output Graph: Social Network

Graph Identification



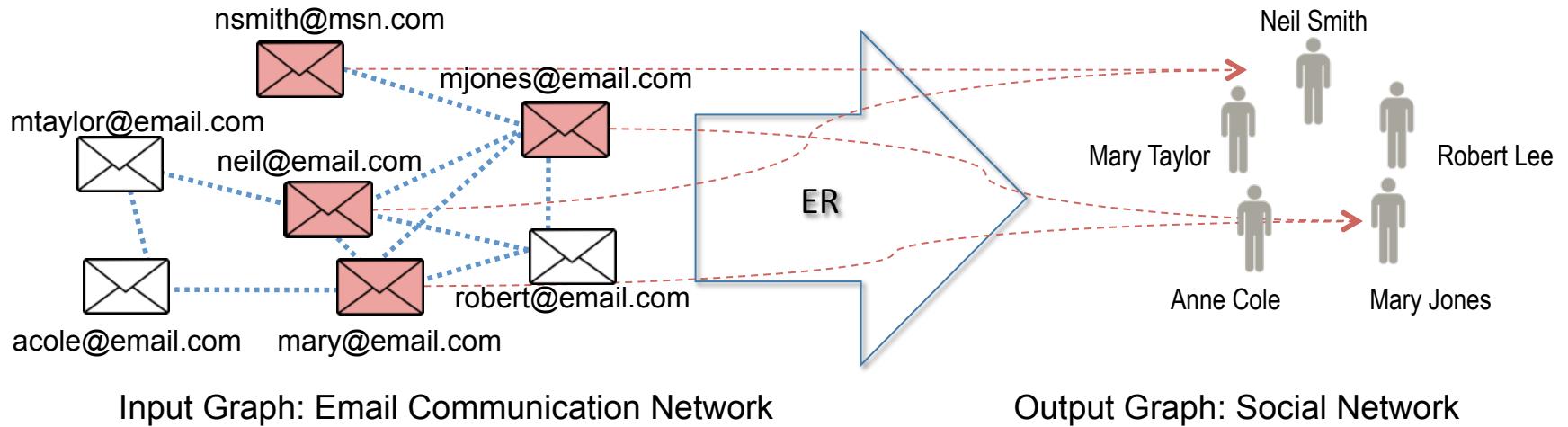
Input Graph: Email Communication Network



Output Graph: Social Network

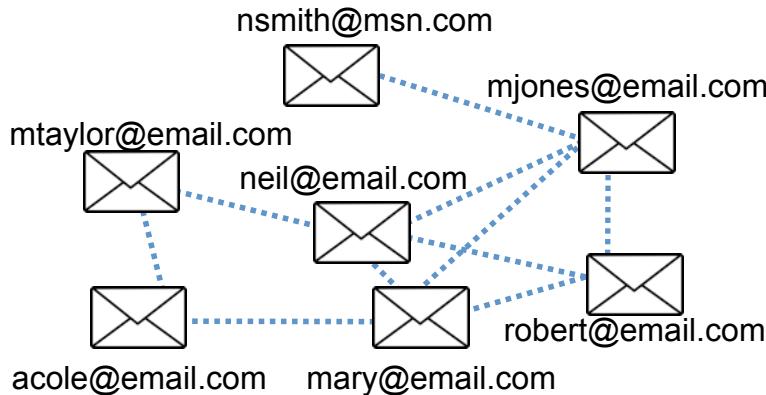
- What's involved?

Graph Identification

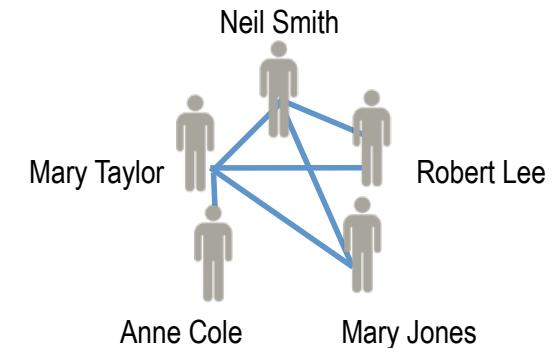
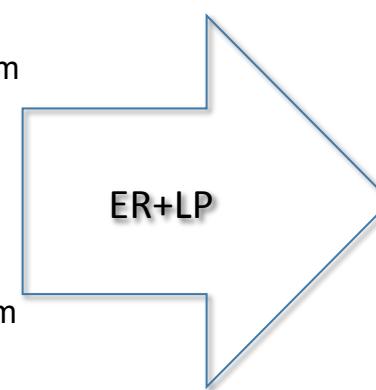


- What's involved?
 - Entity Resolution (ER): Map input graph nodes to output graph nodes

Graph Identification



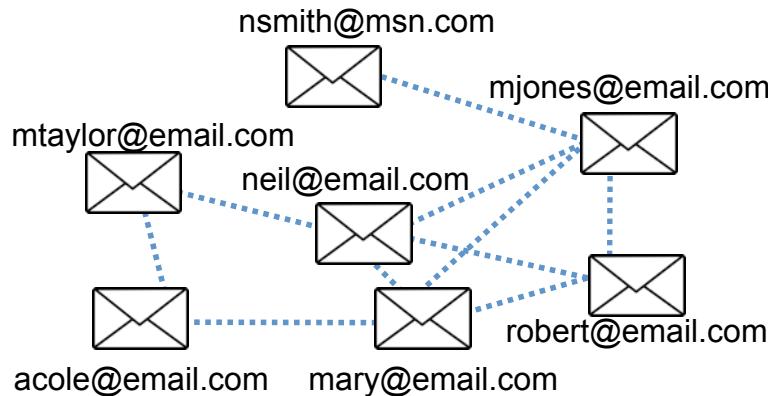
Input Graph: Email Communication Network



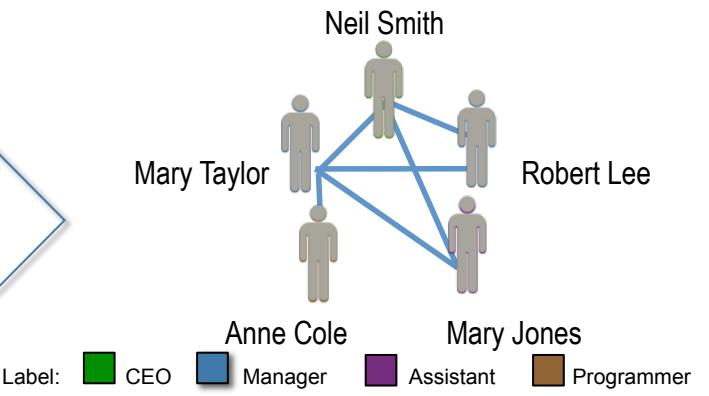
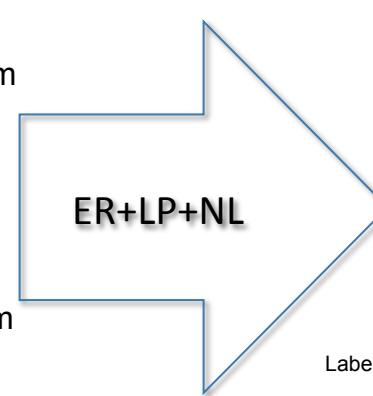
Output Graph: Social Network

- What's involved?
 - Entity Resolution (ER): Map input graph nodes to output graph nodes
 - Link Prediction (LP): Predict existence of edges in output graph

Graph Identification



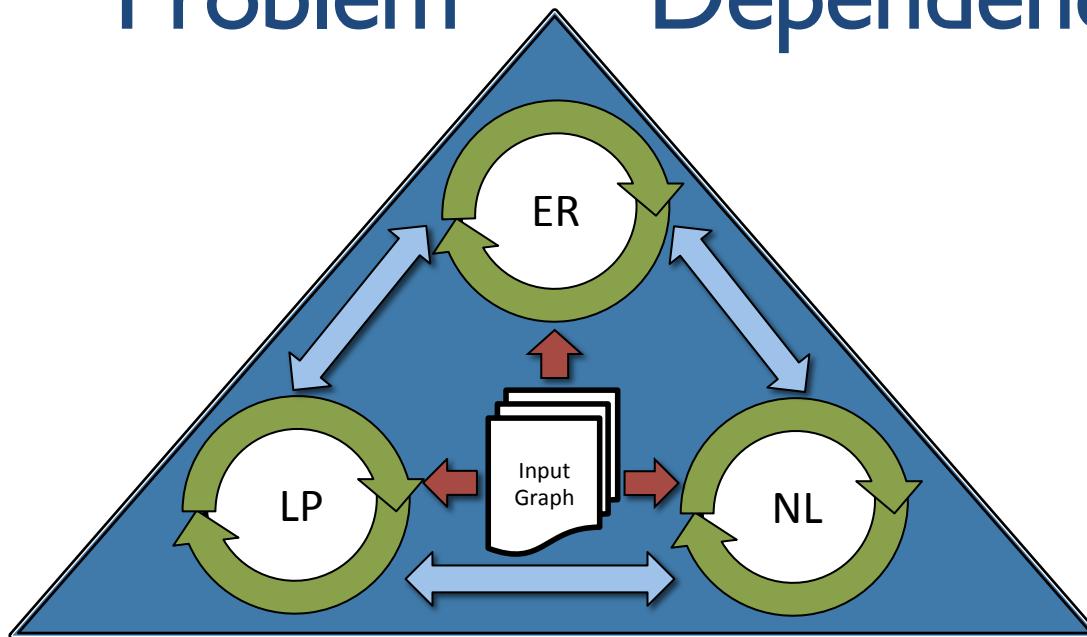
Input Graph: Email Communication Network



Output Graph: Social Network

- What's involved?
 - Entity Resolution (ER): Map input graph nodes to output graph nodes
 - Link Prediction (LP): Predict existence of edges in output graph
 - Node Labeling (NL): Infer the labels of nodes in the output graph

Problem Dependencies

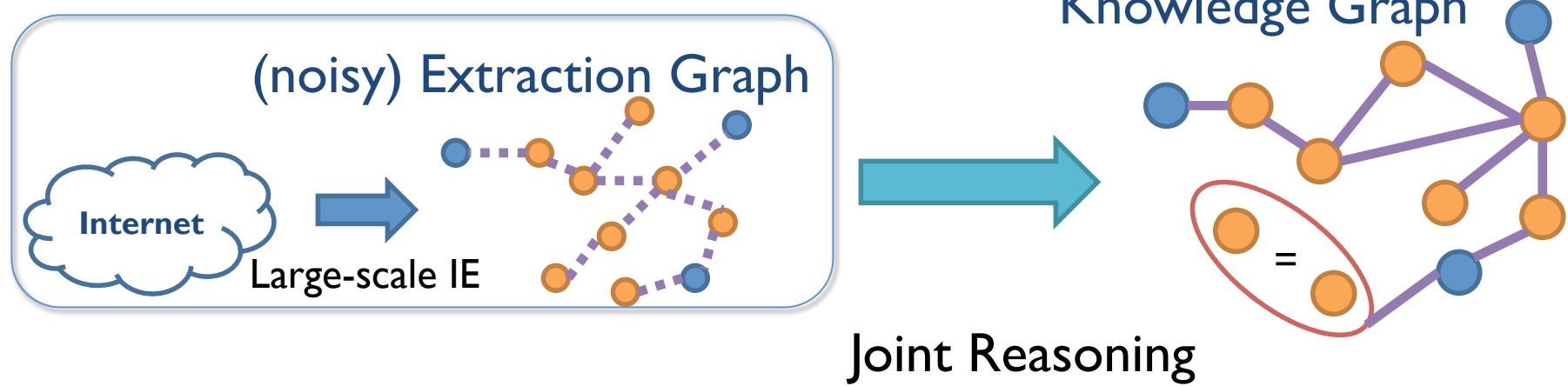


- Most work looks at these tasks in isolation
- In graph identification they are:
 - Evidence-Dependent – Inference depend on observed input graph
e.g., ER depends on input graph
 - Intra-Dependent – Inference within tasks are dependent
e.g., NL prediction depend on other NL predictions
 - Inter-Dependent – Inference across tasks are dependent
e.g., LP depend on ER and NL predictions

KNOWLEDGE GRAPH IDENTIFICATION

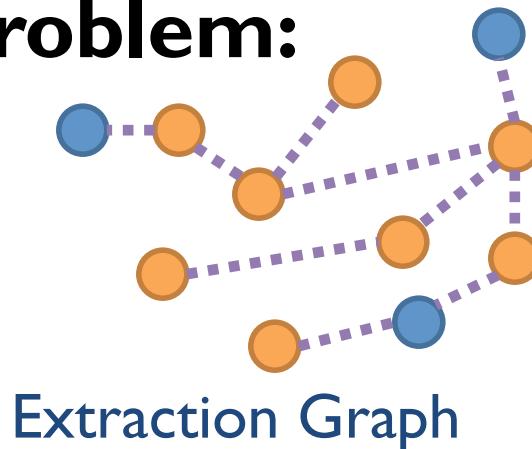
Pujara, Miao, Getoor, Cohen, ISWC 2013 (best student paper)

Motivating Problem (revised)



Knowledge Graph Identification

Problem:

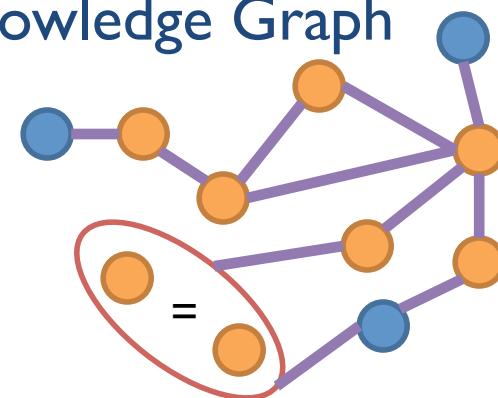


Extraction Graph



Knowledge
Graph
Identification

Knowledge Graph



Solution: Knowledge Graph Identification (KGI)

- Performs *graph identification*:
 - entity resolution
 - node labeling
 - link prediction
- Enforces *ontological constraints*
- Incorporates *multiple uncertain sources*

Illustration of KGI: Extractions

Uncertain Extractions:

- .5: Lbl(Kyrgyzstan, bird)
- .7: Lbl(Kyrgyzstan, country)
- .9: Lbl(Kyrgyz Republic, country)
- .8: Rel(Kyrgyz Republic, Bishkek, hasCapital)

Illustration of KGI: Ontology + ER

Uncertain Extractions:

- .5: Lbl(Kyrgyzstan, bird)
- .7: Lbl(Kyrgyzstan, country)
- .9: Lbl(Kyrgyz Republic, country)
- .8: Rel(Kyrgyz Republic, Bishkek, hasCapital)

Extraction Graph

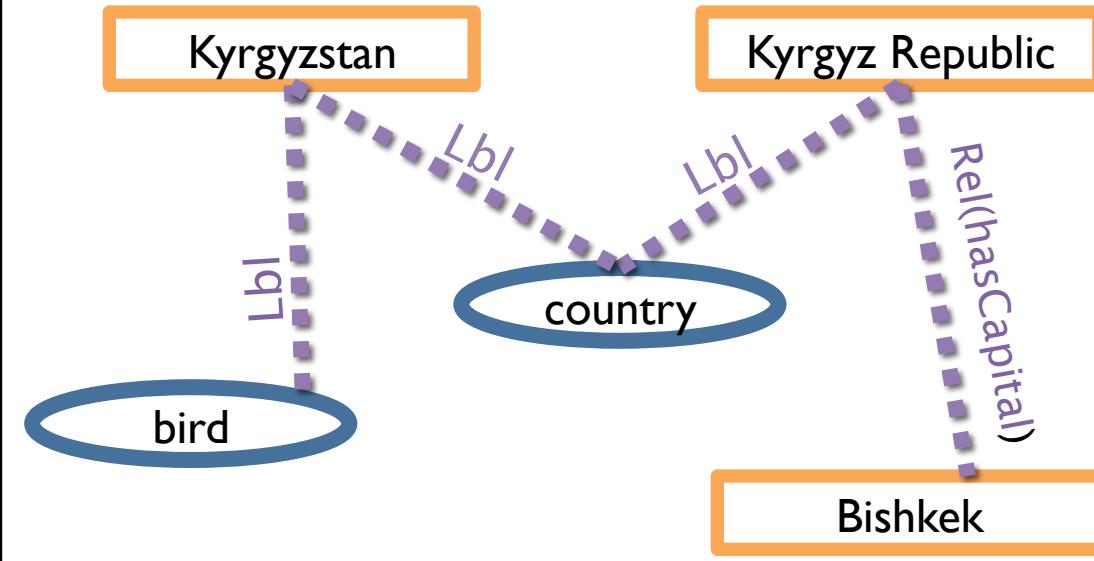


Illustration of KGI: Ontology + ER

Uncertain Extractions:

- .5: Lbl(Kyrgyzstan, bird)
- .7: Lbl(Kyrgyzstan, country)
- .9: Lbl(Kyrgyz Republic, country)
- .8: Rel(Kyrgyz Republic, Bishkek, hasCapital)

Ontology:

- Dom(hasCapital, country)
- Mut(country, bird)

Entity Resolution:

- SameEnt(Kyrgyz Republic, Kyrgyzstan)

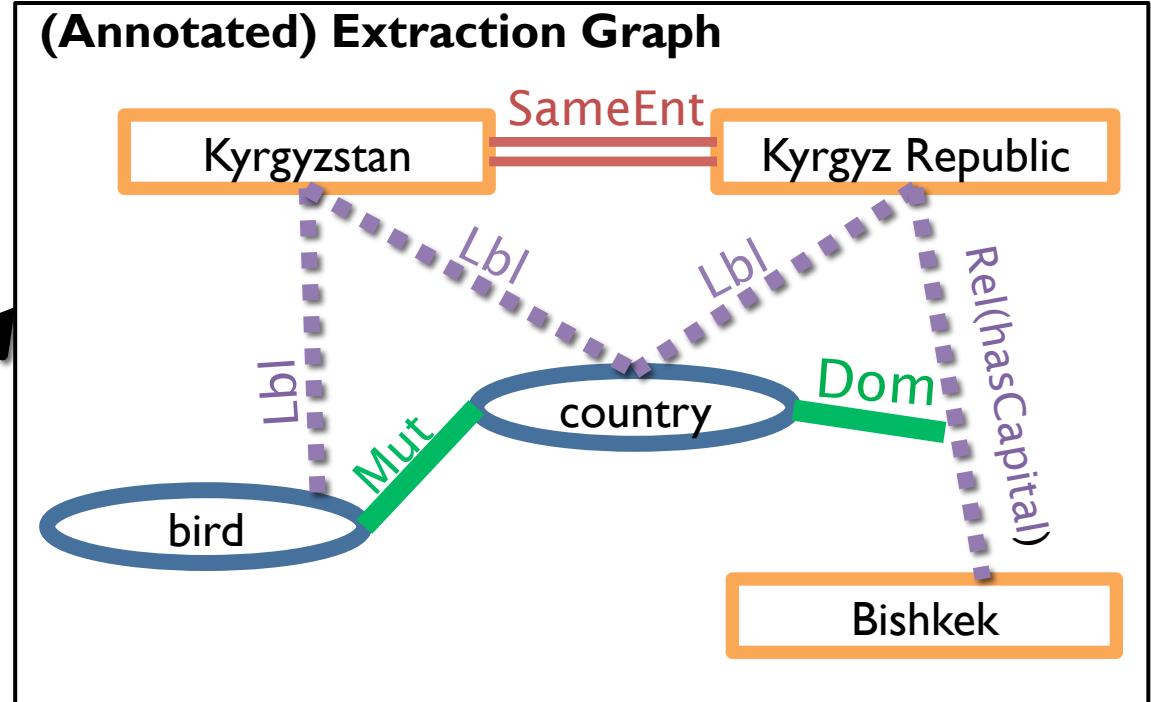


Illustration of KGI

Uncertain Extractions:

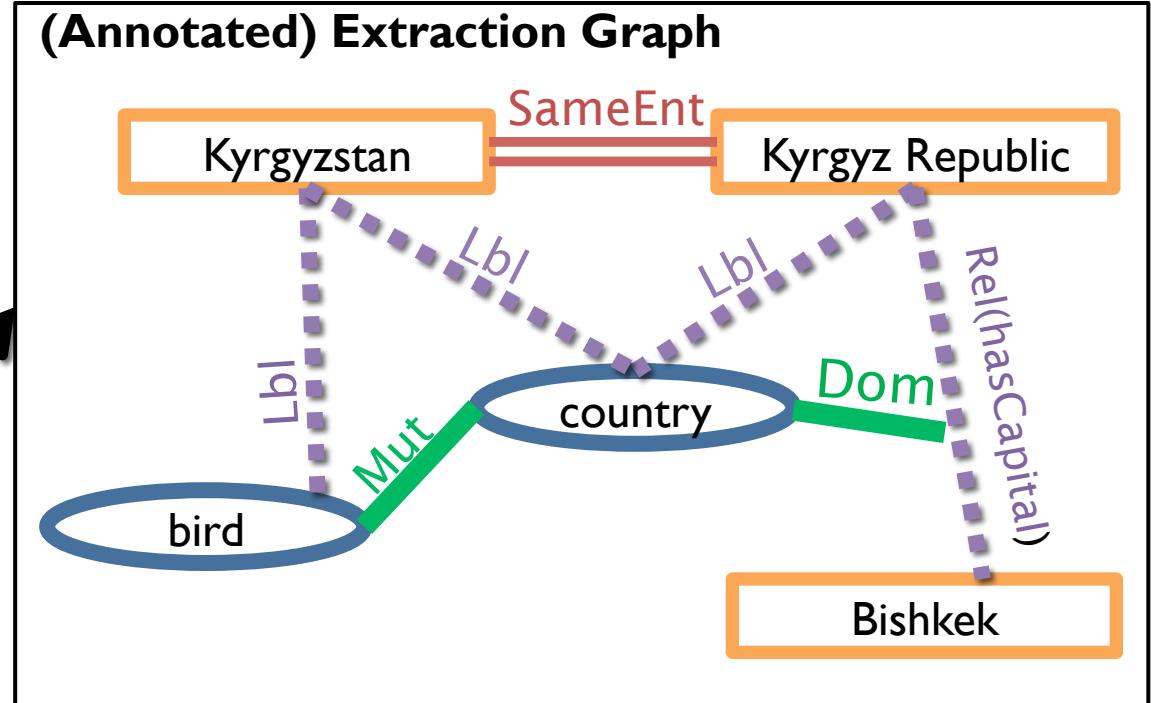
- .5: Lbl(Kyrgyzstan, bird)
- .7: Lbl(Kyrgyzstan, country)
- .9: Lbl(Kyrgyz Republic, country)
- .8: Rel(Kyrgyz Republic, Bishkek, hasCapital)

Ontology:

Dom(hasCapital, country)
Mut(country, bird)

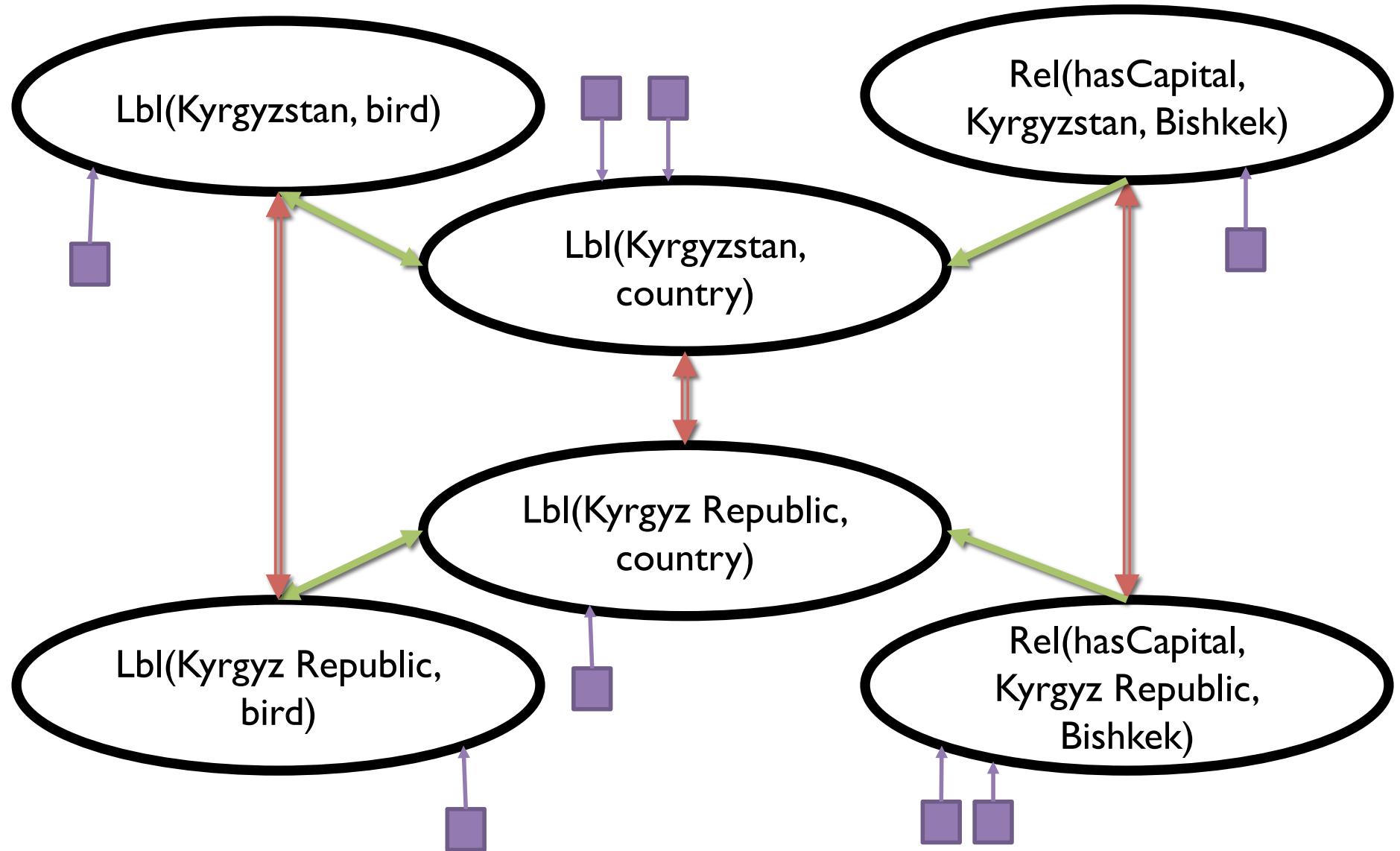
Entity Resolution:

SameEnt(Kyrgyz Republic, Kyrgyzstan)



Modeling Knowledge Graph Identification

Viewing KGI as a probabilistic graphical model



Background: Probabilistic Soft Logic (PSL)

(Broeckeler et al., UAI10; Kimmig et al., NIPS-ProbProg12)

- Templating language for hinge-loss MRFs, very scalable!
- Model specified as a collection of logical formulas

$$\text{SAMEENT}(E_1, E_2) \underset{\curvearrowleft}{\wedge} \text{LBL}(E_1, L) \Rightarrow \text{LBL}(E_2, L)$$

Uses soft-logic formulation

- Truth values of atoms relaxed to $[0,1]$ interval
- Truth values of formulas derived from Lukasiewicz t-norm

$$p \tilde{\wedge} q = \max(0, p + q - 1)$$

$$p \tilde{\vee} q = \min(1, p + q)$$

$$\tilde{\neg}p = 1 - p$$

$$p \tilde{\Rightarrow} q = \min(1, q - p + 1)$$

Soft Logic Tutorial: Rules to Groundings

- Given a database of evidence, we can convert rule templates to instances (grounding)
- Rules are *grounded* by substituting literals into formulas

$$\text{SAMEENT}(E_1, E_2) \ \tilde{\wedge} \ \text{LBL}(E_1, L) \Rightarrow \text{LBL}(E_2, L)$$

$$\text{SAMEENT}(\text{Kyrgyzstan}, \text{Kyrgyz Republic})$$
$$\tilde{\wedge} \ \text{LBL}(\text{Kyrgyzstan}, \text{country})$$
$$\Rightarrow \text{LBL}(\text{Kyrgyz Republic}, \text{country})$$

- The soft logic interpretation assigns a “satisfaction” value to each ground rule

Soft Logic Tutorial: Groundings to Satisfaction

SAMEENT(Kyrgyzstan, Kyrgyz Republic) : 0.9 $\tilde{\wedge}$
LBL(Kyrgyzstan, country) : 0.8

$$p \tilde{\vee} q = \max(0, p + q - 1)$$

$$\begin{aligned} & \text{SAMEENT(Kyrgyzstan, Kyrgyz Republic)} \tilde{\wedge} \\ & \text{LBL(Kyrgyzstan, country)} \\ &= \max(0, 0.9 + 0.8 - 1) \end{aligned}$$

Soft Logic Tutorial: Groundings to Satisfaction

(SAMEENT(Kyrgyzstan, Kyrgyz Republic)
 \wedge LBL(Kyrgyzstan, country)) : 0.7
 \Rightarrow LBL(Kyrgyz Republic, country) : 0.6

$$p \tilde{\Rightarrow} q = \min(1, q - p + 1)$$

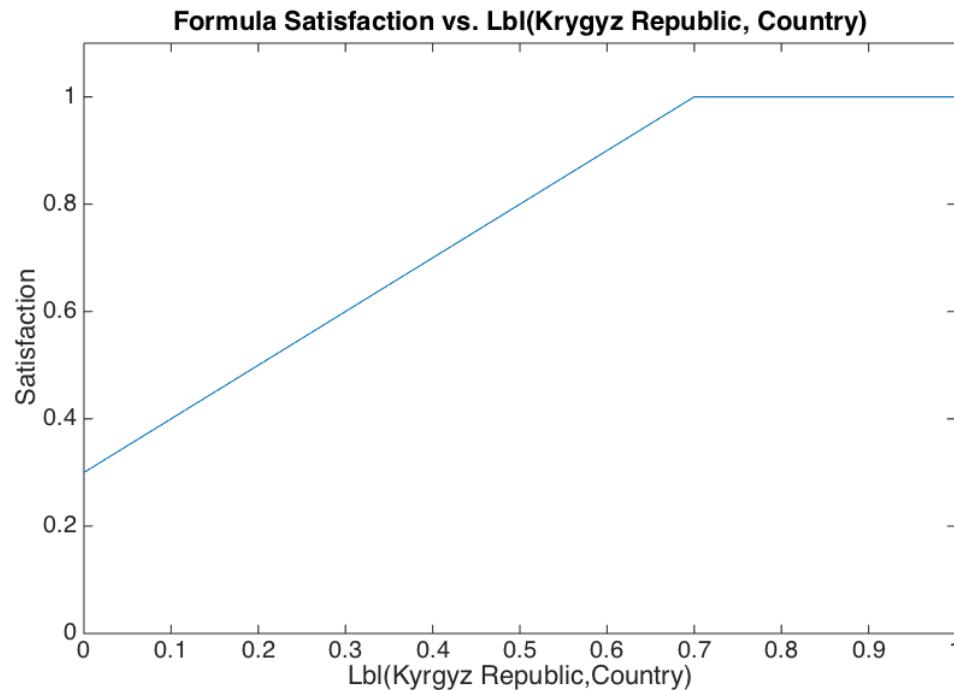
SAMEENT(Kyrgyzstan, Kyrgyz Republic)
 \wedge LBL(Kyrgyzstan, country)
 \Rightarrow LBL(Kyrgyz Republic, country)
= $\min(1, 0.6 - 0.7 + 1) = 0.9$

Soft Logic Tutorial: Inferring Satisfaction

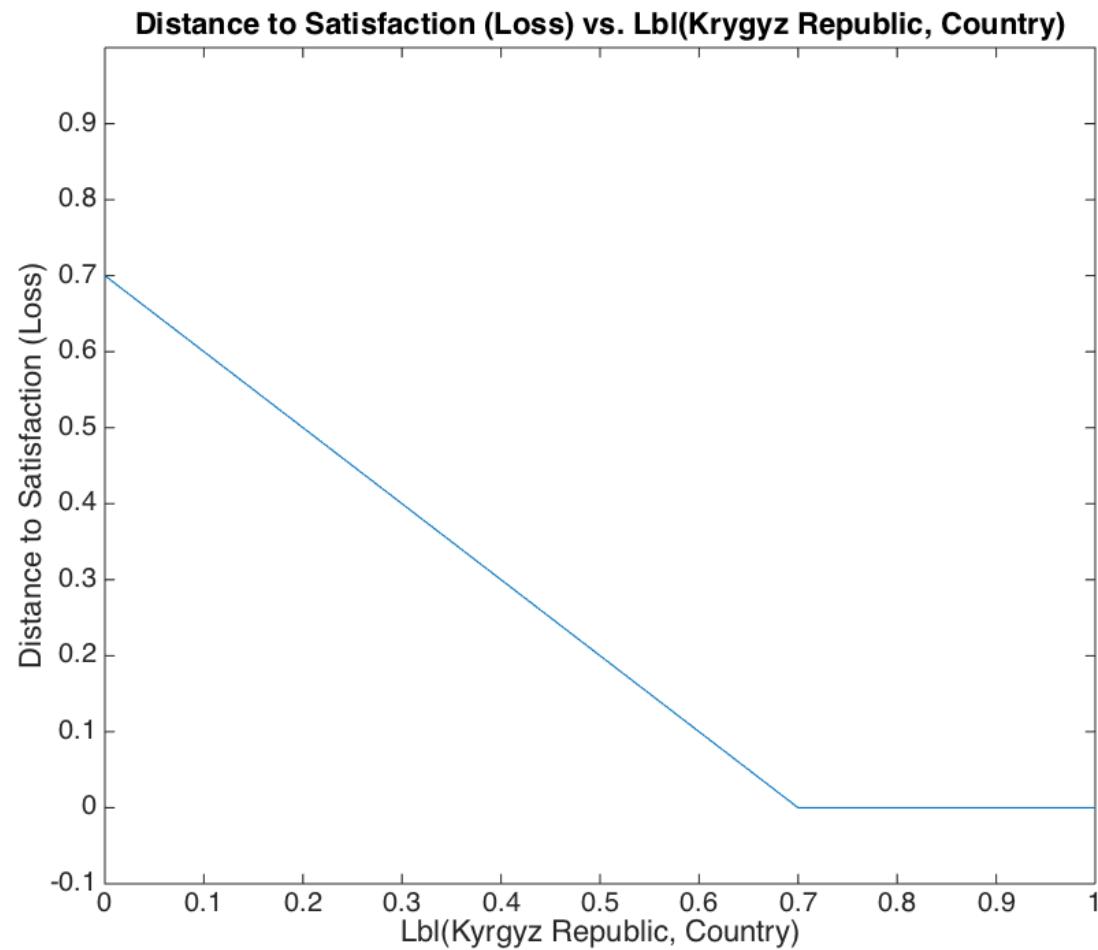
(SAMEENT(Kyrgyzstan, Kyrgyz Republic)

\wedge LBL(Kyrgyzstan, country)) : 0.7

\Rightarrow LBL(Kyrgyz Republic, country) :?



Soft Logic Tutorial: Distance to Satisfaction



Background: PSL Rules to Distributions

- Rules are *grounded* by substituting literals into formulas

$w_{EL} : \text{SAMEENT}(\text{Kyrgyzstan}, \text{Kyrgyz Republic}) \wedge$
 $\text{LBL}(\text{Kyrgyzstan}, \text{country}) \Rightarrow \text{LBL}(\text{Kyrgyz Republic}, \text{country})$

- Each ground rule has a weighted distance to satisfaction derived from the formula's truth value

$$P(G | E) = \frac{1}{Z} \exp \left[- \sum_{r \in R} w_r \varphi_r(G) \right]$$

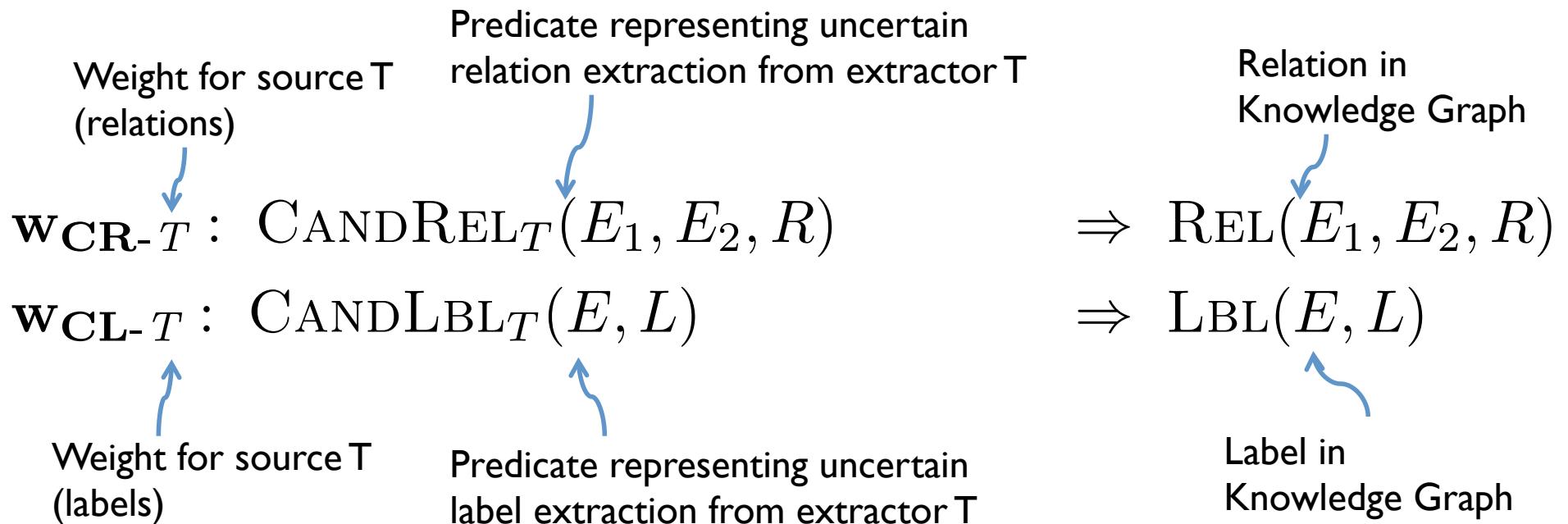
- The PSL program can be interpreted as a joint probability distribution over all variables in knowledge graph, conditioned on the extractions

Background: Finding the best knowledge graph

- MPE inference solves $\max_G P(G)$ to find the best KG
- In PSL, inference solved by convex optimization
- Efficient: running time empirically scales with $O(|R|)$
(Bach et al., NIPS12)

PSL Rules for KGI Model

PSL Rules: Uncertain Extractions



PSL Rules: Entity Resolution

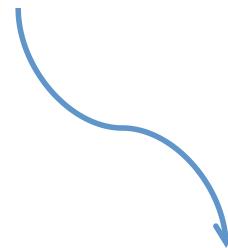
wEL : SAMEENT(E_1, E_2) $\tilde{\wedge}$ LBL(E_1, L) \Rightarrow LBL(E_2, L)

wER : SAMEENT(E_1, E_2) $\tilde{\wedge}$ REL(E_1, E, R) \Rightarrow REL(E_2, E, R)

wER : SAMEENT(E_1, E_2) $\tilde{\wedge}$ REL(E, E_1, R) \Rightarrow REL(E, E_2, R)



SameEnt predicate captures confidence that entities are co-referent



- Rules require co-referent entities to have the same labels and relations
- Creates an equivalence class of co-referent entities

PSL Rules: Ontology

Inverse:

$$\mathbf{w_O} : \text{INV}(R, S) \quad \tilde{\wedge} \quad \text{REL}(E_1, E_2, R) \quad \Rightarrow \quad \text{REL}(E_2, E_1, S)$$

Selectional Preference:

$$\mathbf{w_O} : \text{DOM}(R, L) \quad \tilde{\wedge} \quad \text{REL}(E_1, E_2, R) \quad \Rightarrow \quad \text{LBL}(E_1, L)$$

$$\mathbf{w_O} : \text{RNG}(R, L) \quad \tilde{\wedge} \quad \text{REL}(E_1, E_2, R) \quad \Rightarrow \quad \text{LBL}(E_2, L)$$

Subsumption:

$$\mathbf{w_O} : \text{SUB}(L, P) \quad \tilde{\wedge} \quad \text{LBL}(E, L) \quad \Rightarrow \quad \text{LBL}(E, P)$$

$$\mathbf{w_O} : \text{RSUB}(R, S) \quad \tilde{\wedge} \quad \text{REL}(E_1, E_2, R) \quad \Rightarrow \quad \text{REL}(E_1, E_2, S)$$

Mutual Exclusion:

$$\mathbf{w_O} : \text{MUT}(L_1, L_2) \quad \tilde{\wedge} \quad \text{LBL}(E, L_1) \quad \Rightarrow \quad \neg \text{LBL}(E, L_2)$$

$$\mathbf{w_O} : \text{RMUT}(R, S) \quad \tilde{\wedge} \quad \text{REL}(E_1, E_2, R) \quad \Rightarrow \quad \neg \text{REL}(E_1, E_2, S)$$

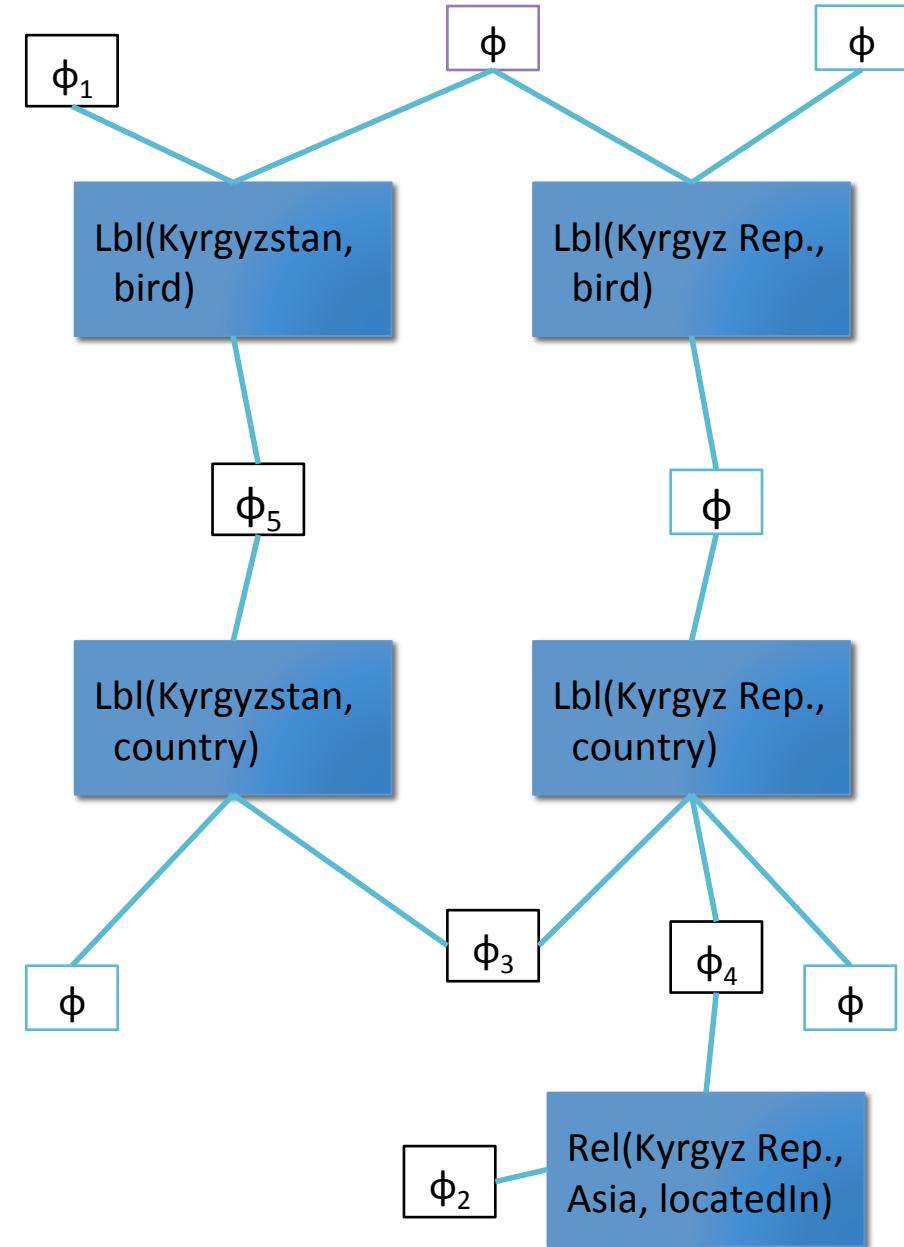
$[\phi_1] \text{ CANDLBL}_{\text{struct}}(\text{Kyrgyzstan}, \text{bird})$
 $\Rightarrow \text{LBL}(\text{Kyrgyzstan}, \text{bird})$

$[\phi_2] \text{ CANDREL}_{\text{pat}}(\text{Kyrgyz Rep.}, \text{Asia}, \text{locatedIn})$
 $\Rightarrow \text{REL}(\text{Kyrgyz Rep.}, \text{Asia}, \text{locatedIn})$

$[\phi_3] \text{ SAMEENT}(\text{Kyrgyz Rep.}, \text{Kyrgyzstan})$
 $\wedge \text{LBL}(\text{Kyrgyz Rep.}, \text{country})$
 $\Rightarrow \text{LBL}(\text{Kyrgyzstan}, \text{country})$

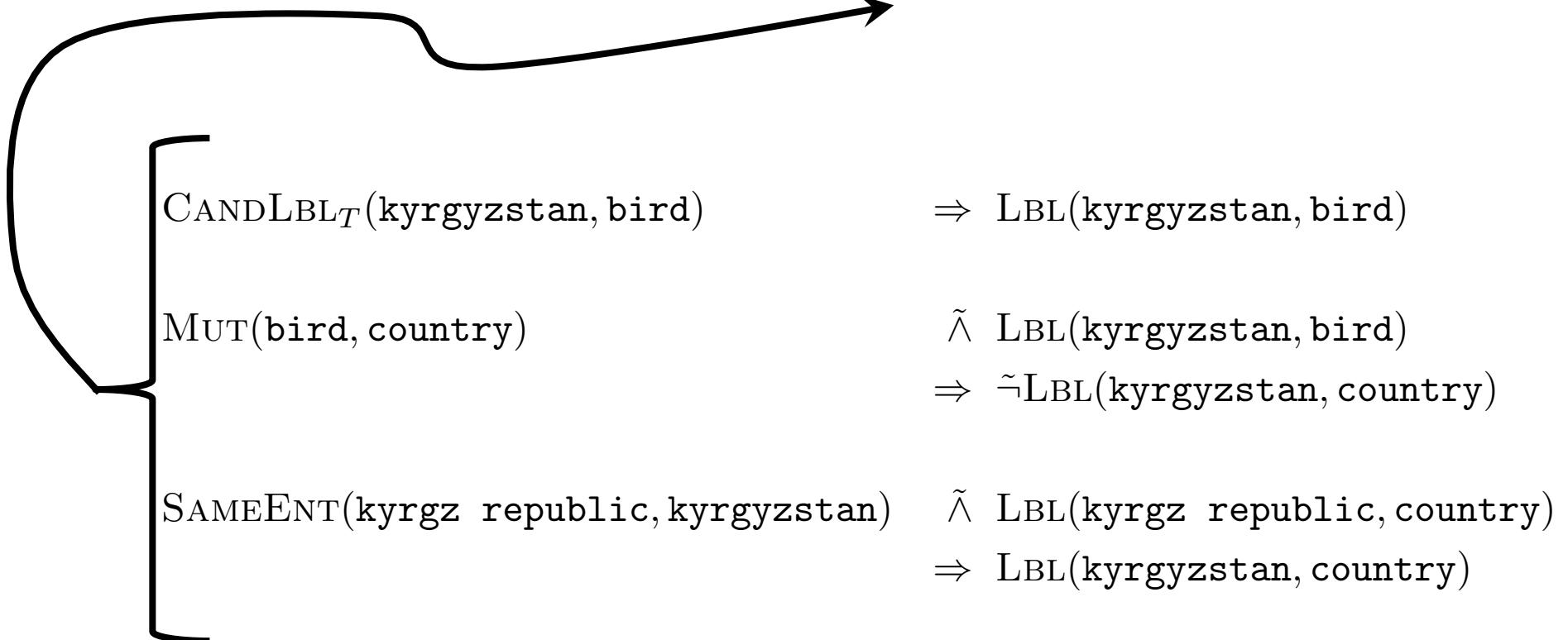
$[\phi_4] \text{ DOM}(\text{locatedIn}, \text{country})$
 $\wedge \text{REL}(\text{Kyrgyz Rep.}, \text{Asia}, \text{locatedIn})$
 $\Rightarrow \text{LBL}(\text{Kyrgyz Rep.}, \text{country})$

$[\phi_5] \text{ MUT}(\text{country}, \text{bird})$
 $\wedge \text{LBL}(\text{Kyrgyzstan}, \text{country})$
 $\Rightarrow \neg \text{LBL}(\text{Kyrgyzstan}, \text{bird})$



Probability Distribution over KGs

$$P(G | E) = \frac{1}{Z} \exp \left[- \sum_{r \in R} w_r \varphi_r(G) \right]$$



Evaluation

Two Evaluation Datasets

| | LinkedBrainz | NELL |
|-----------------------------|---|--|
| Description | Community-supplied data about musical artists, labels, and creative works | Real-world IE system extracting general facts from the WWW |
| Noise | Realistic synthetic noise | Imperfect extractors and ambiguous web pages |
| Candidate Facts | 810K | 1.3M |
| Unique Labels and Relations | 27 | 456 |
| Ontological Constraints | 49 | 67.9K |

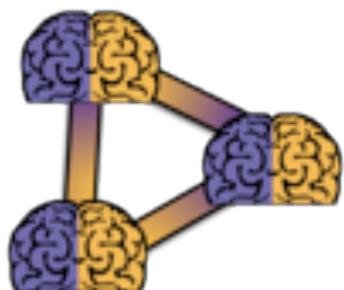
LinkedBrainz



- Open source community-driven structured database of music metadata
- Uses proprietary schema to represent data

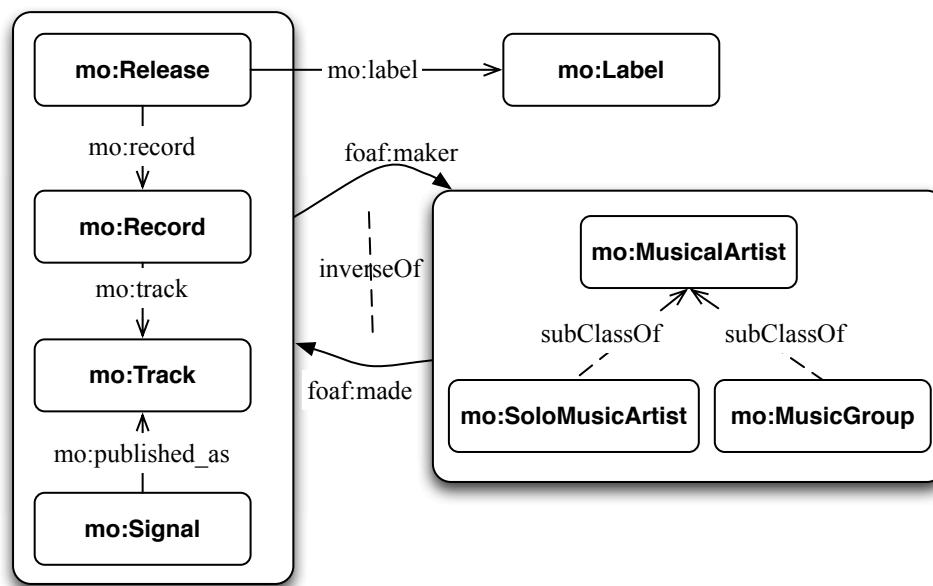


- Built on popular ontologies such as FOAF and FRBR
- Widely used for music data (e.g. BBC Music Site)



LinkedBrainz project provides an RDF mapping from MusicBrainz data to Music Ontology using the D2RQ tool

LinkedBrainz dataset for KGI



Mapping to FRBR/FOAF ontology

| | |
|------|--------------------|
| DOM | rdfs:domain |
| RNG | rdfs:range |
| INV | owl:inverseOf |
| SUB | rdfs:subClassOf |
| RSUB | rdfs:subPropertyOf |
| MUT | owl:disjointWith |

LinkedBrainz experiments

Comparisons:

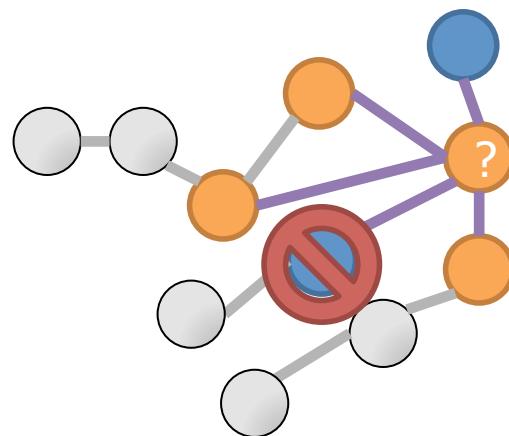
- Baseline** Use noisy truth values as fact scores
- PSL-EROnly** Only apply rules for Entity Resolution
- PSL-OntOnly** Only apply rules for Ontological reasoning
- PSL-KGI** Apply Knowledge Graph Identification model

| | AUC | Precision | Recall | F1 at .5 | Max F1 |
|-------------|------------|------------------|---------------|-----------------|---------------|
| Baseline | 0.672 | 0.946 | 0.477 | 0.634 | 0.788 |
| PSL-EROnly | 0.797 | 0.953 | 0.558 | 0.703 | 0.831 |
| PSL-OntOnly | 0.753 | 0.964 | 0.605 | 0.743 | 0.832 |
| PSL-KGI | 0.901 | 0.970 | 0.714 | 0.823 | 0.919 |

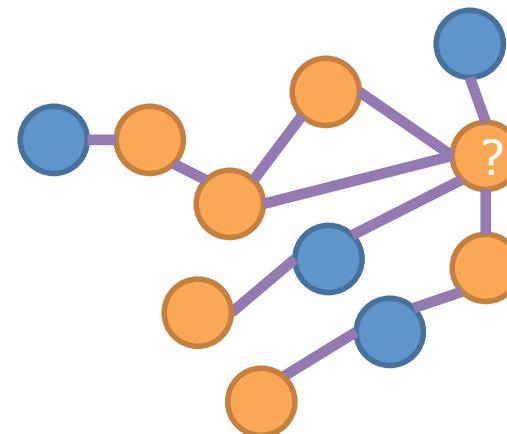
NELL Evaluation: two settings

Target Set: restrict to a subset of KG

(Jiang, ICDM12)



Complete: Infer full knowledge graph



- Closed-world model
- Uses a target set: subset of KG
- Derived from 2-hop neighborhood
- Excludes trivially satisfied variables

- Open-world model
- All possible entities, relations, labels
- Inference assigns truth value to each variable

NELL experiments: Target Set

Task: Compute truth values of a target set derived from the evaluation data

Comparisons:

Baseline Average confidences of extractors for each fact in the NELL candidates

NELL Evaluate NELL's promotions (on the full knowledge graph)

MLN Method of (Jiang, ICDM12) – estimates marginal probabilities with MC-SAT

PSL-KGI Apply full Knowledge Graph Identification model

Running Time: Inference completes in 10 seconds, values for 25K facts

| | AUC | FI |
|-----------------|------|------|
| Baseline | .873 | .828 |
| NELL | .765 | .673 |
| MLN (Jiang, 12) | .899 | .836 |
| PSL-KGI | .904 | .853 |

NELL experiments: Complete knowledge graph

Task: Compute a full knowledge graph from uncertain extractions

Comparisons:

NELL NELL's strategy: ensure ontological consistency with existing KB

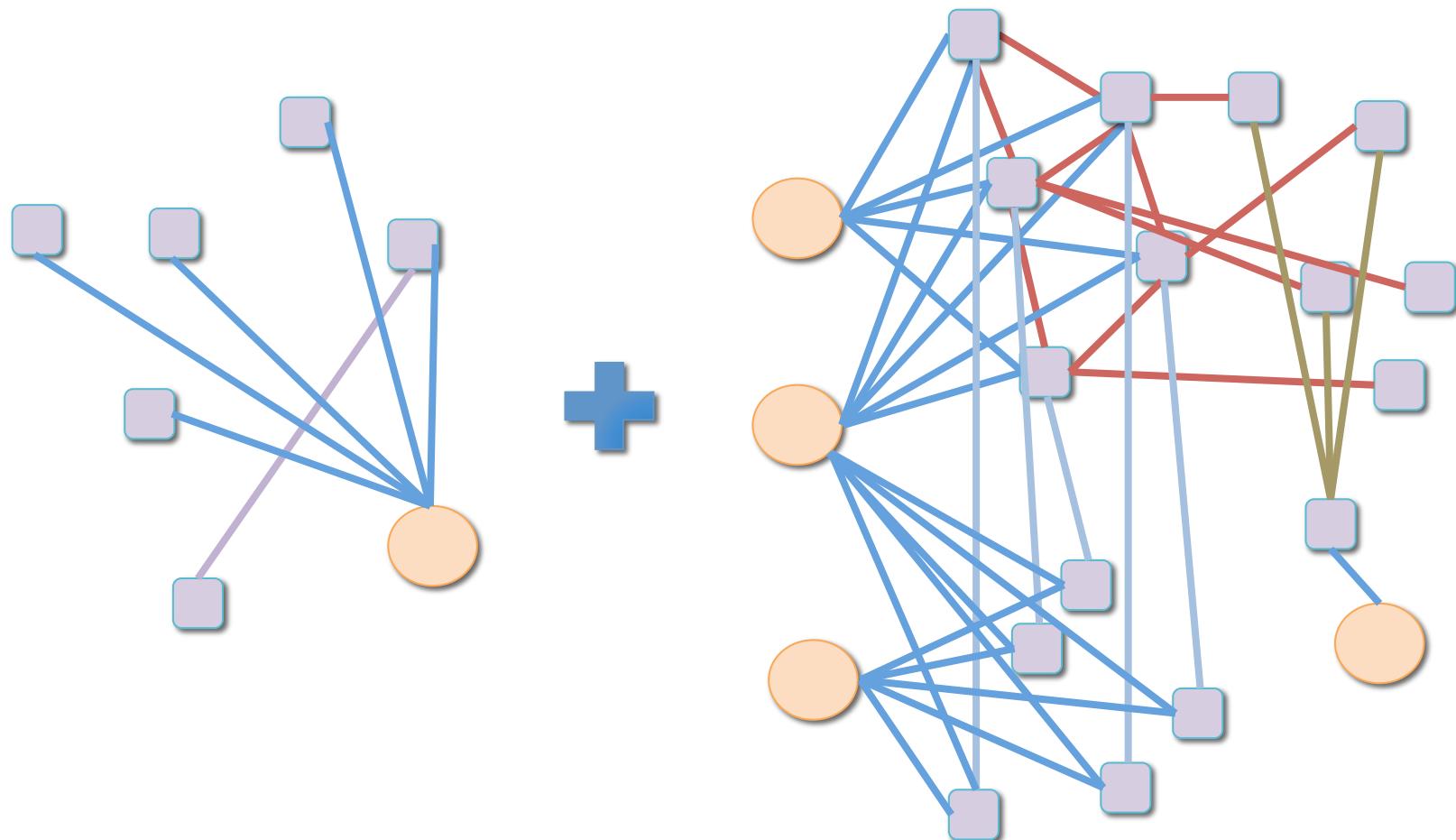
PSL-KGI Apply full Knowledge Graph Identification model

Running Time: Inference completes in 130 minutes, producing 4.3M facts

| | AUC | Precision | Recall | F1 |
|---------|-------|-----------|--------|-------|
| NELL | 0.765 | 0.801 | 0.477 | 0.634 |
| PSL-KGI | 0.892 | 0.826 | 0.871 | 0.848 |

KNOWLEDGE GRAPH ENTITY RESOLUTION

Problem: Merge domain KG to global KG



Approach: Factored Entity Resolution model

- Goal: Build a generic entity resolution model for KGs
- Build on vast amount of work on Entity Resolution
- PSL provides an easy, flexible, sophisticated models

| | Local | Collective |
|-----------------|--------------------|------------------------|
| General | String similarity | Sparsity; Transitivity |
| New Entity | New Entity prior | New Entity penalty |
| Knowledge Graph | Type compatibility | Relation compatibility |
| Domain-Specific | (Album length) | (Artist's country) |

Preliminary Results

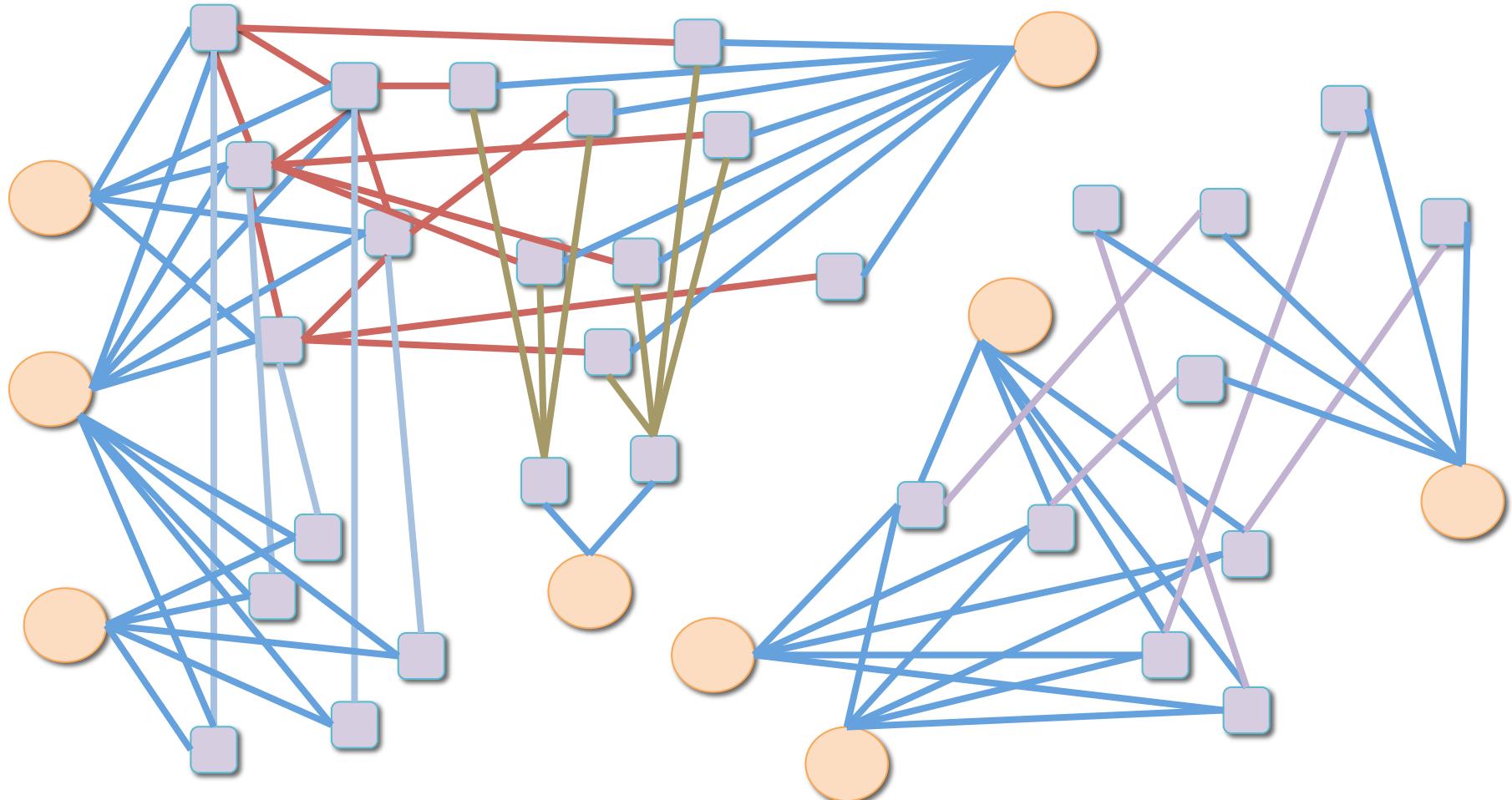
- Task: ER from MusicBrainz to Google KG
- Data:
 - 11K MusicBrainz entities (5/5-6/29/14)
 - 330K Freebase entities
 - 15.7M relations
 - 11K human labels

| Methods | F1 | AUPRC |
|-------------|-------|-------|
| General | 0.734 | 0.416 |
| +Collective | 0.805 | 0.569 |
| +NewEntity | 0.840 | 0.724 |

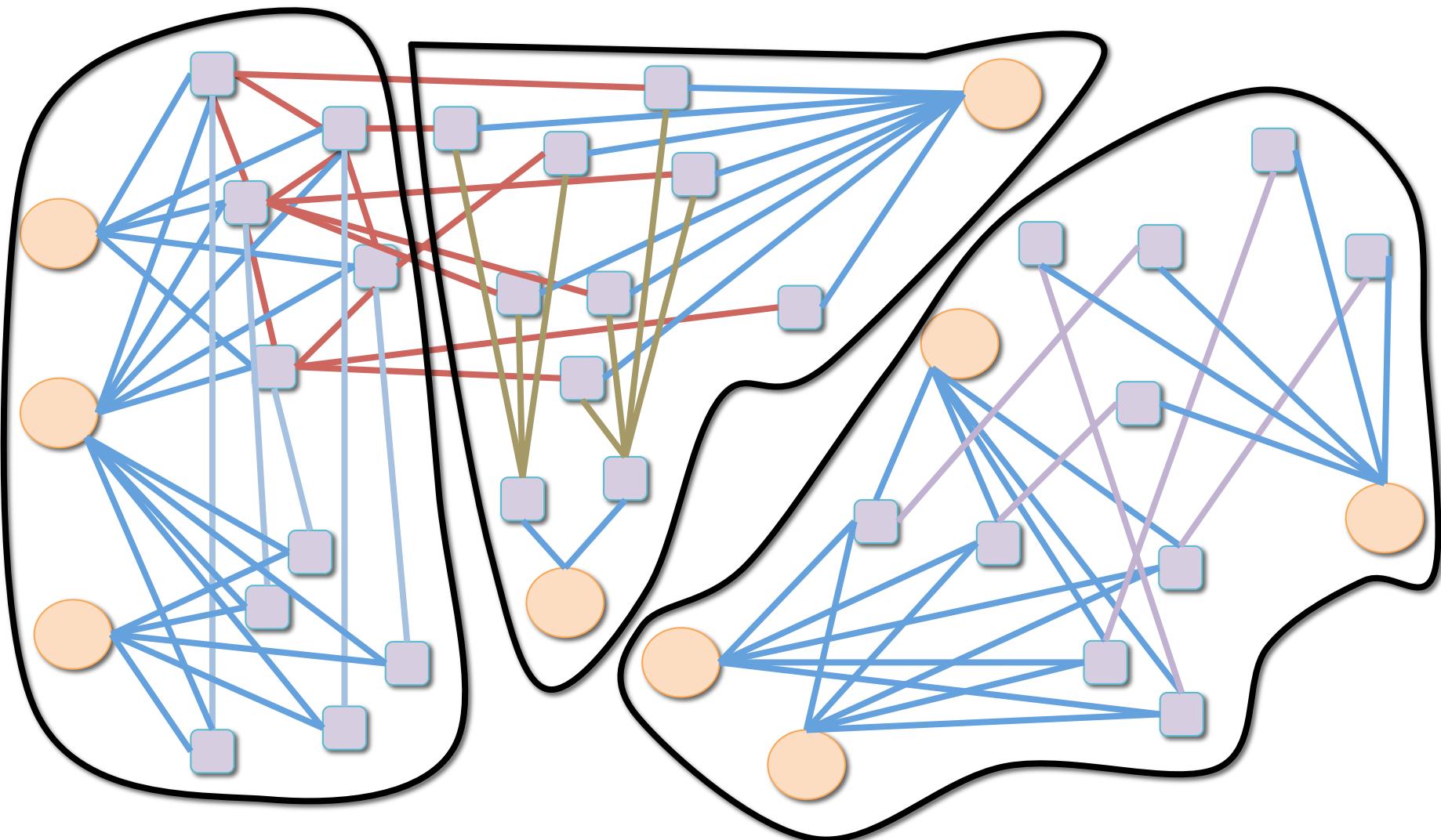
FASTER KNOWLEDGE GRAPH CONSTRUCTION

Partitioning

Problem: Knowledge Graphs are **HUGE**



Solution: Partition the Knowledge Graph

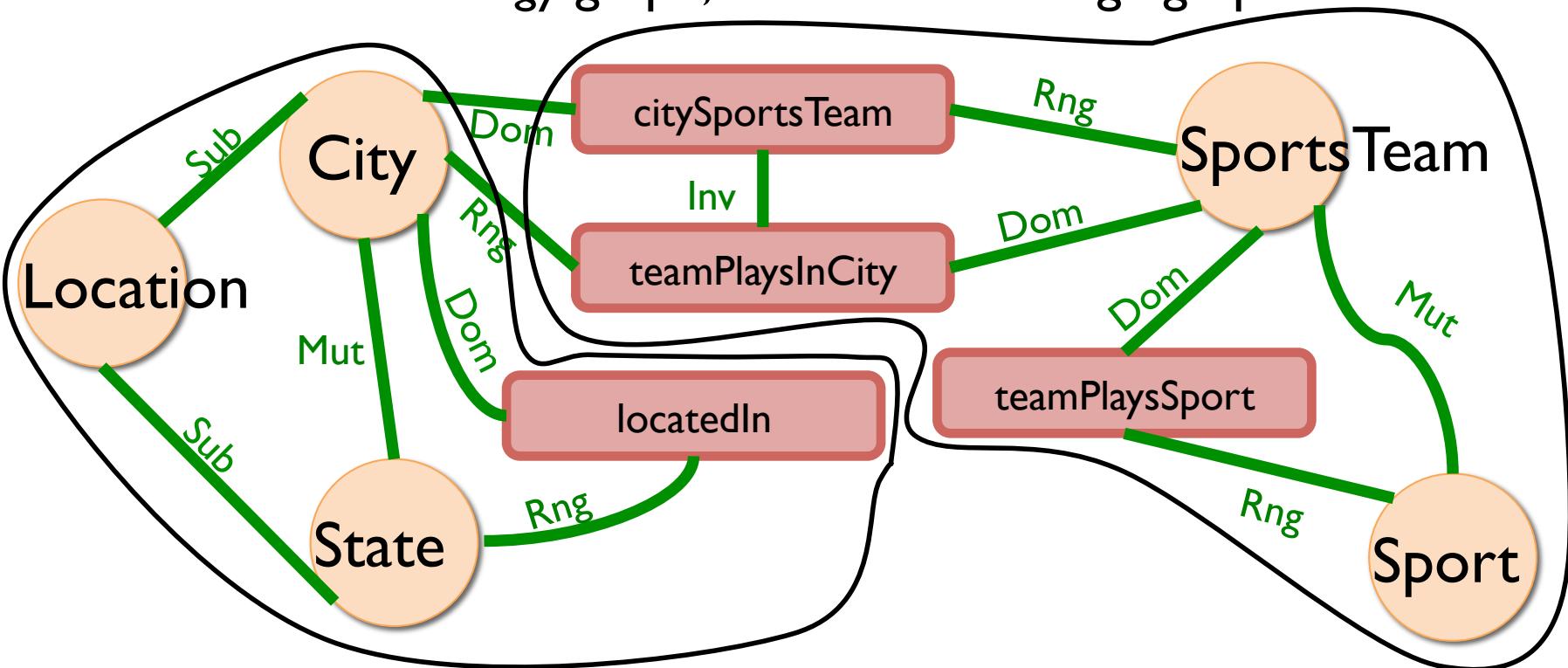


Partitioning: advantages and drawbacks

- Advantages
 - Smaller problems
 - Parallel Inference
 - Speed / Quality Tradeoff
- Drawbacks
 - Partitioning large graph time-consuming
 - Key dependencies may be lost
 - New facts require re-partitioning

Key idea: Ontology-aware partitioning

- Partition the *ontology graph*, not the knowledge graph



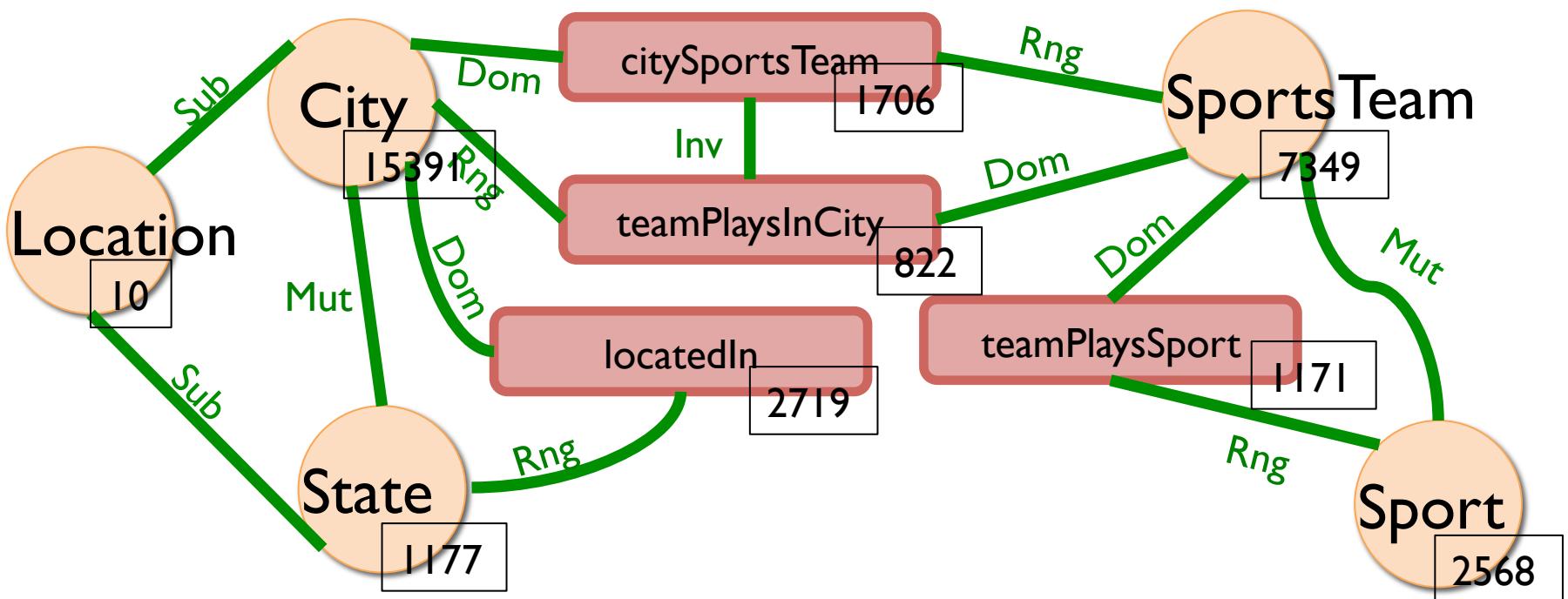
- Induce a partitioning of the knowledge graph based on the ontology partition

Considerations: Ontology-aware Partitions

- Advantages:
 - Ontology is a smaller graph
 - Ontology coupled with dependencies
 - New facts can reuse partitions
- Disadvantages:
 - Insensitive to data distribution
 - All dependencies treated equally

Refinement: include data frequency

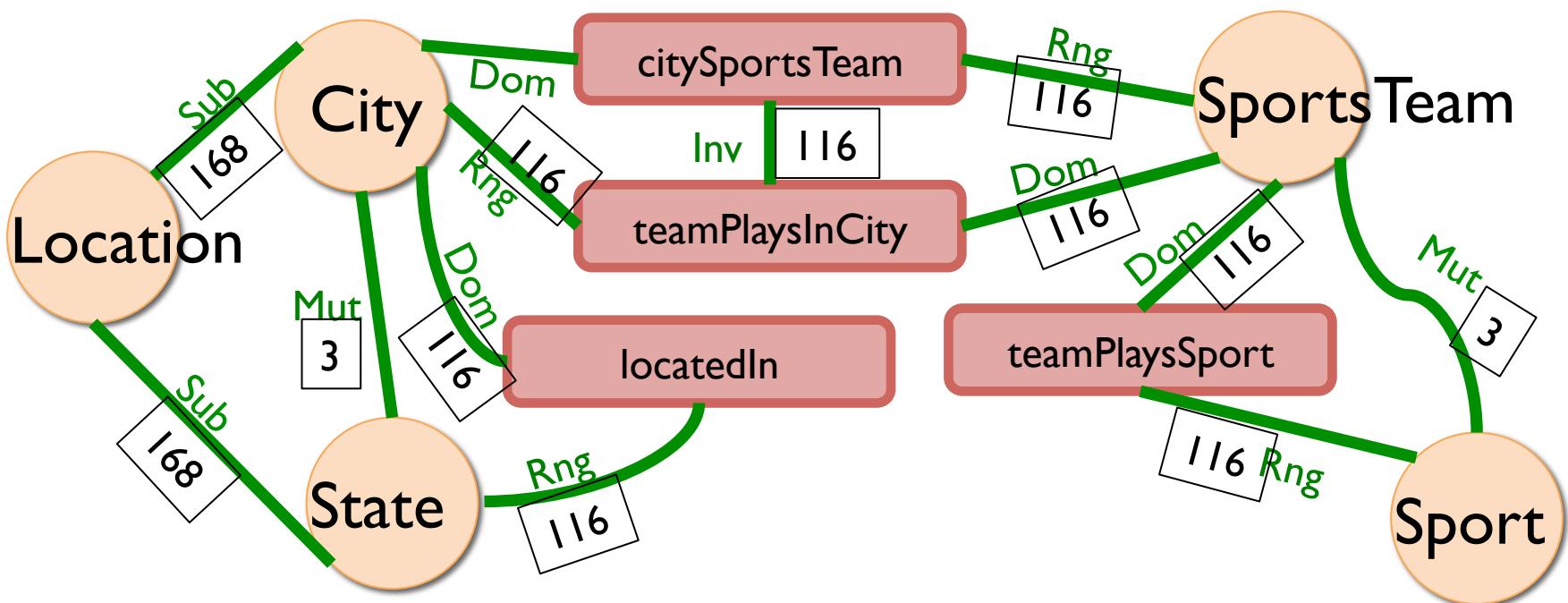
- Annotate each ontological element with its frequency



- Partition ontology with constraint of equal vertex weights

Refinement: weight edges by type

- Weight edges by their ontological importance



Experiments: Partitioning Approaches

Comparisons (6 partitions):

NELL

Default promotion strategy, no KGI

KGI

No partitioning, full knowledge graph model

baseline

KGI, Randomly assign extractions to partition

Ontology

KGI, Edge min-cut of ontology graph

O+Vertex

KGI, Weight ontology vertices by frequency

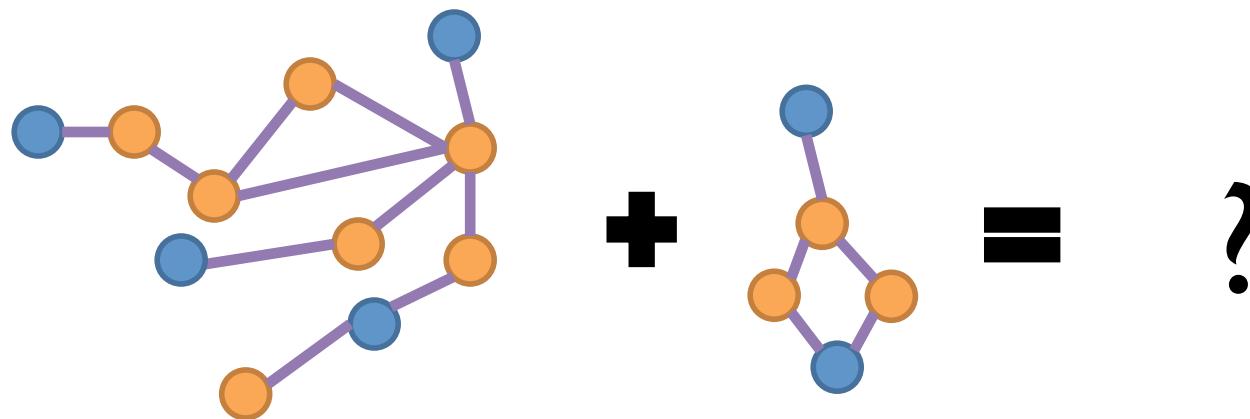
O+V+Edge

KGI, Weight ontology edges by inv. frequency

| | AUPRC | Running Time (min) | Opt. Terms |
|----------|--------------|--------------------|------------|
| NELL | 0.765 | - | |
| KGI | 0.794 | 97 | 10.9M |
| baseline | 0.780 | 31 | 3.0M |
| Ontology | 0.788 | 42 | 4.2M |
| O+Vertex | 0.791 | 31 | 3.7M |
| O+V+Edge | 0.790 | 31 | 3.7M |

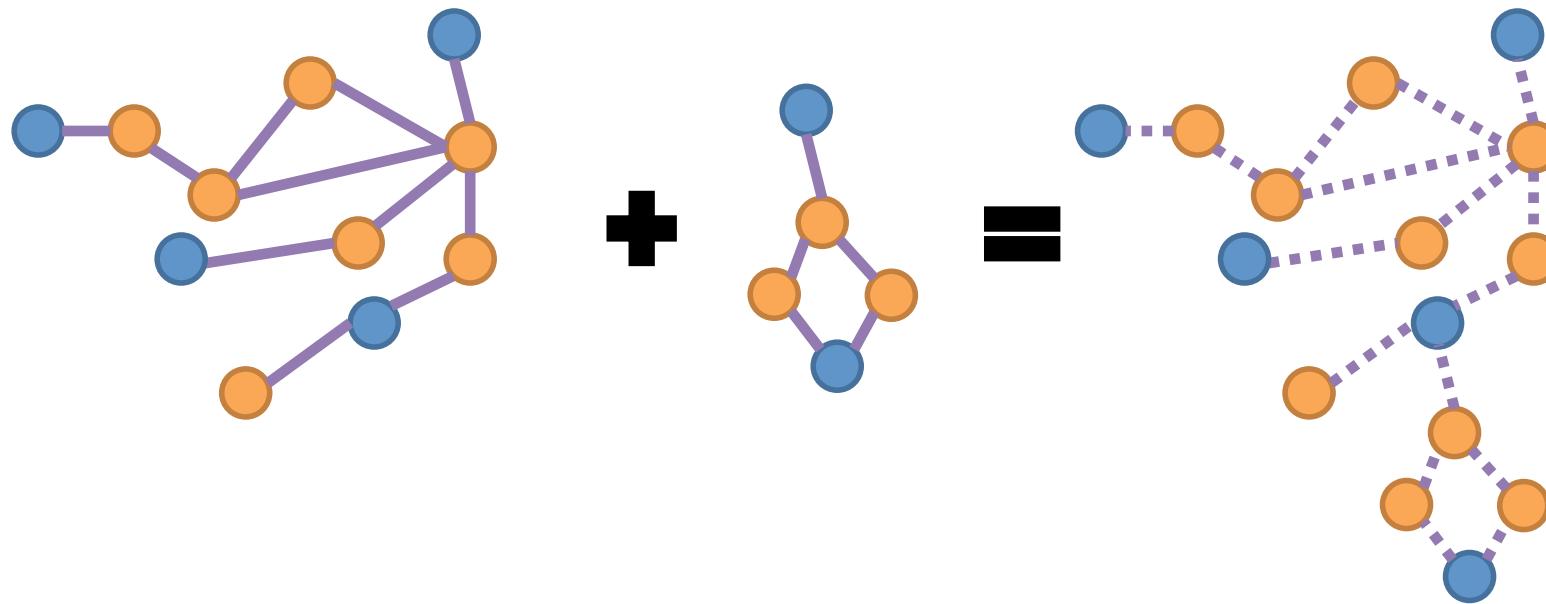
Evolving Models

Problem: Incremental Updates to KG



How do we add new extractions to the Knowledge Graph?

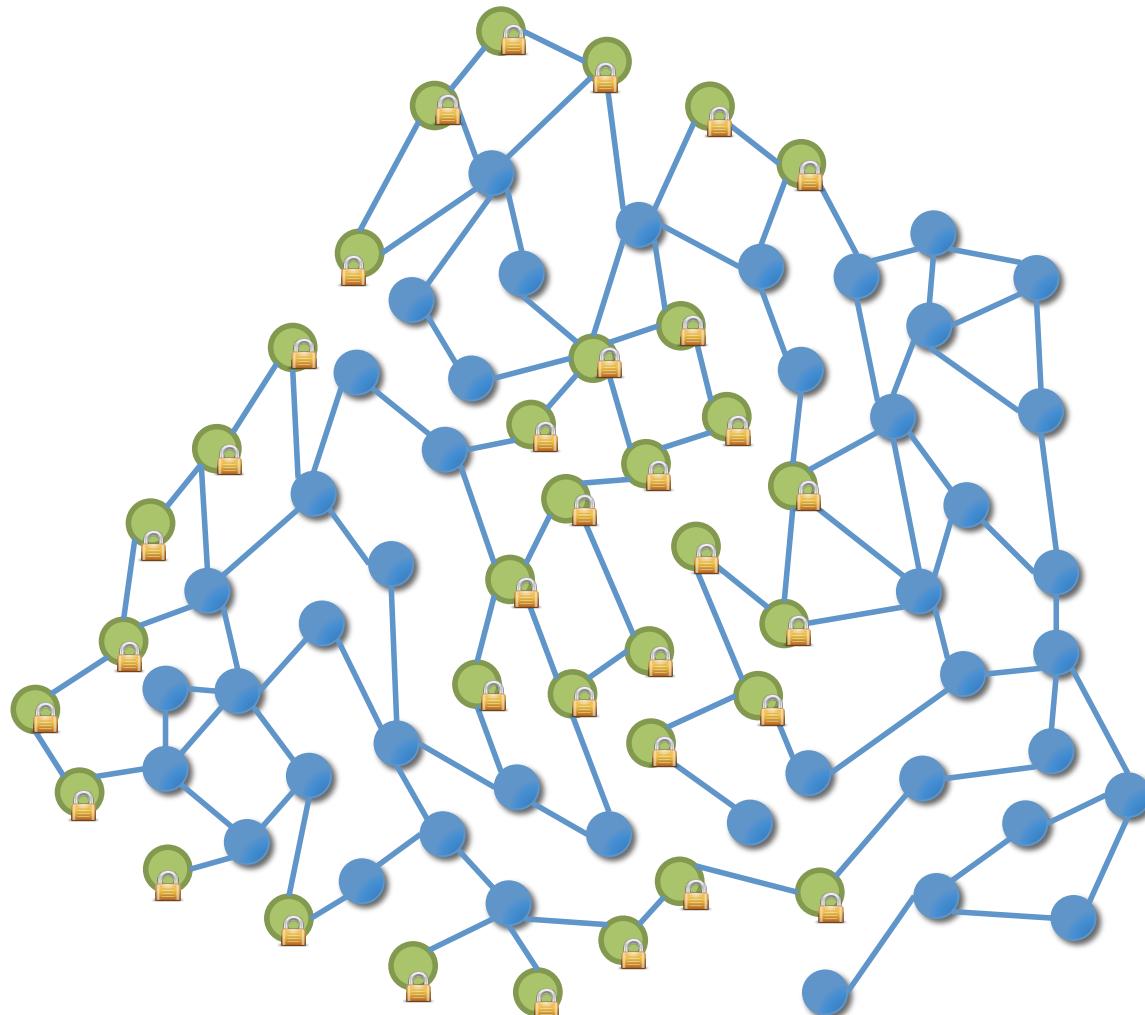
Naïve Approach: Full KGI over extractions



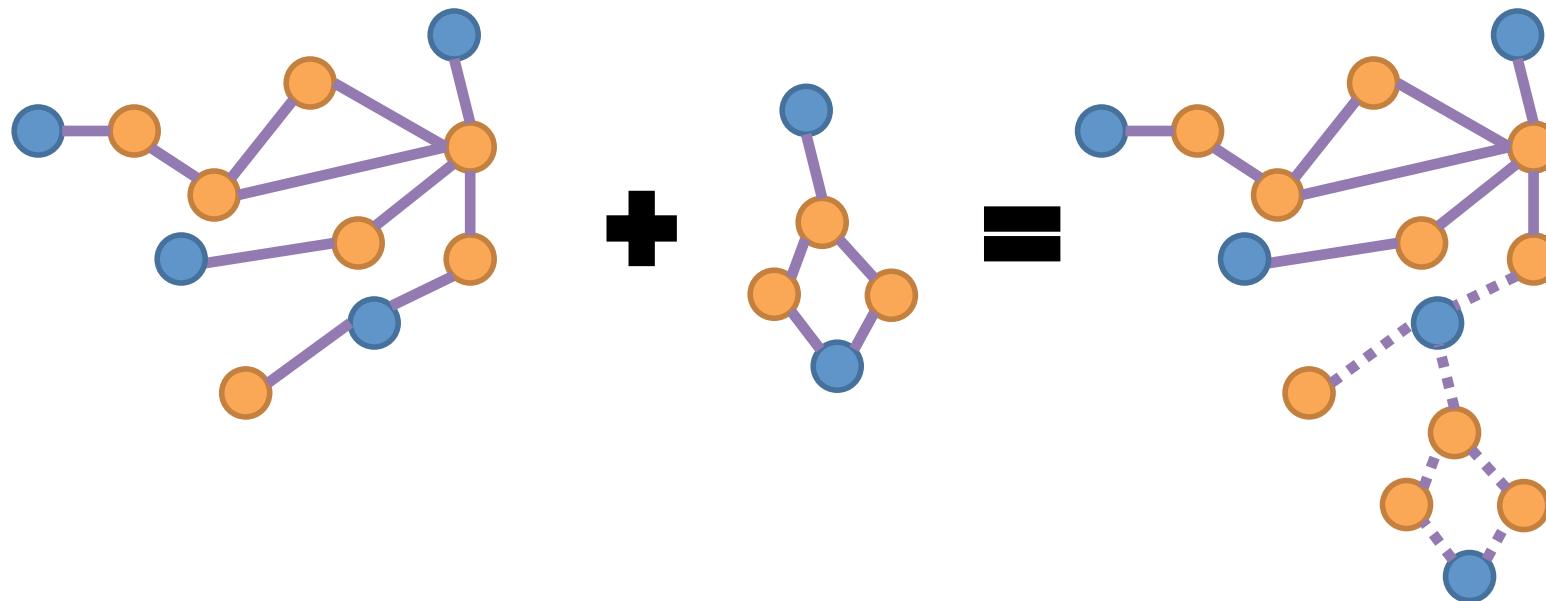
Improving the naïve approach

- **Intuition:** Much of previous KG does not change
- **Online collective inference:**
 - Selectively update the MAP state
 - Bound the *regret* of partial updates
 - Efficiently determine which variables to infer

Key Idea: fix some variables, infer others



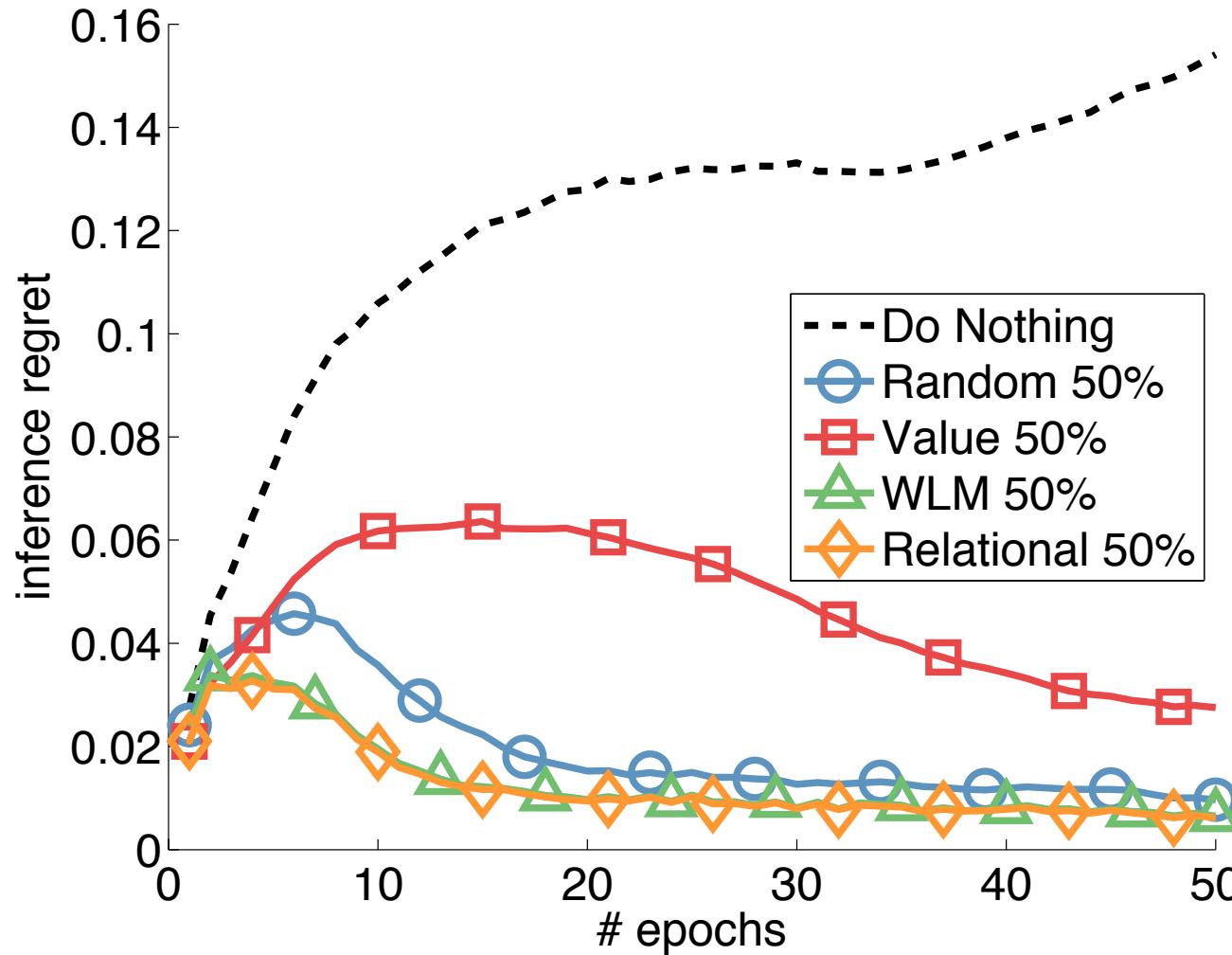
Approximation: KGI over subset of graph



Theory: Regret of approximating update

$$\mathfrak{R}_n(\mathbf{x}, \mathbf{y}_{\mathcal{S}}; \dot{\mathbf{w}}) \leq O\left(\sqrt{\frac{B\|\mathbf{w}\|_2}{n \cdot w_p} \|\mathbf{y}_{\mathcal{S}} - \hat{\mathbf{y}}_{\mathcal{S}}\|_1}\right)$$

Practice: Regret and Approximation Algo



Conclusion

- Knowledge Graph Identification is a powerful technique for producing knowledge graphs from noisy IE system output
- Using PSL we are able to enforce global ontological constraints and capture uncertainty in our model
- Unlike previous work, our approach infers complete knowledge graphs for datasets with millions of extractions

Code available on GitHub:

<https://github.com/linqs/KnowledgeGraphIdentification>

Key Collaborators

