

KNOWLEDGE GRAPH IDENTIFICATION

Jay Pujara

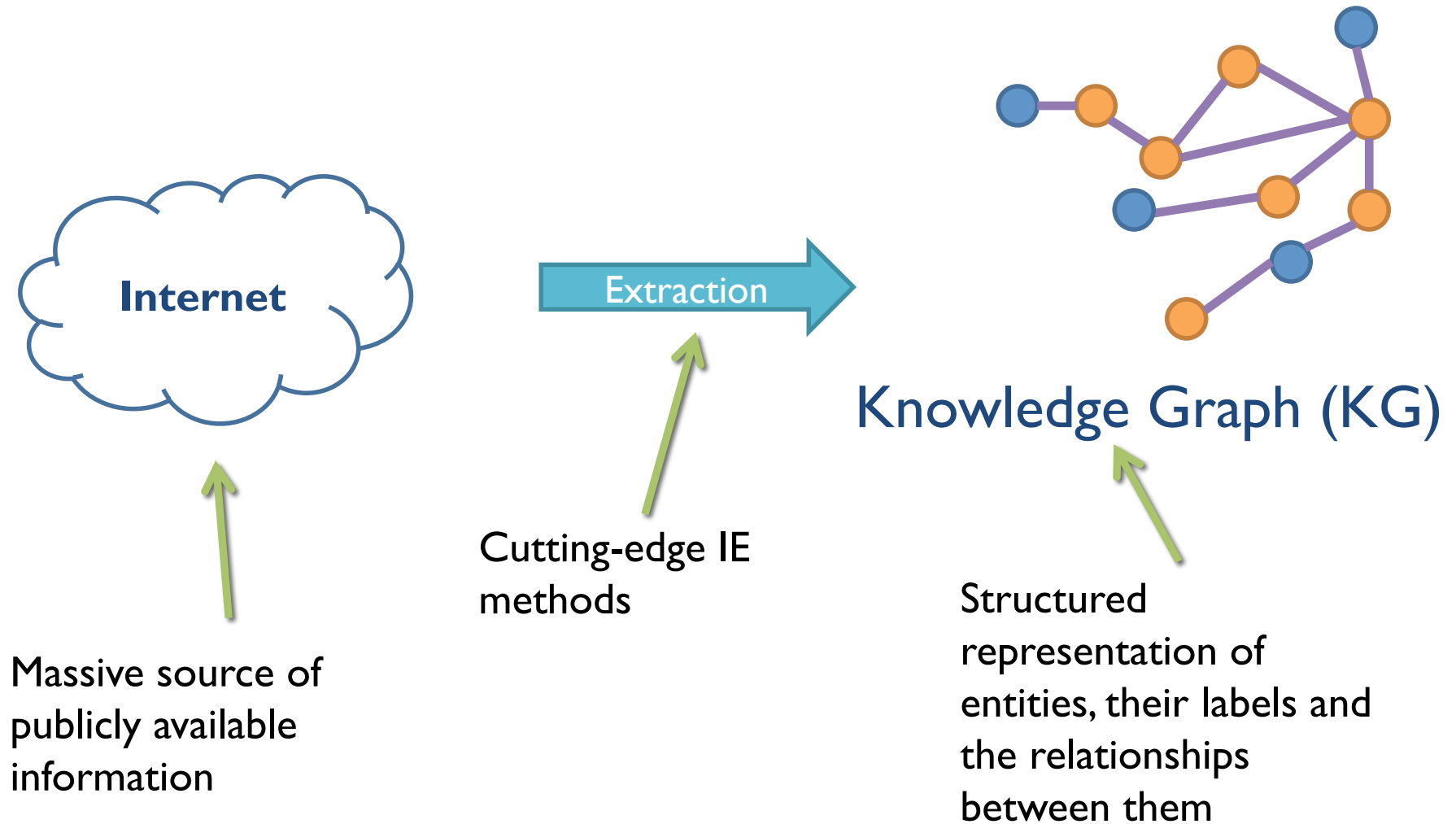
`jay@cs.umd.edu`

11/5/2014

presented at Carnegie Mellon University





Motivating Problem: Opportunities




Knowledge Graphs in the wild

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

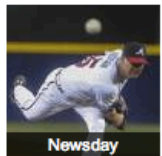
Yesterday, 4:25 PM (ET)
MetLife Stadium, East Rutherford, New Jersey

 **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 ...

Giants close to deal with Tim Hudson
ESPN - 1 hour ago

Giants close to signing veteran hurler Hudson
MLB.com - 2 hours ago

 **NEW YORK GIANTS**
OFFICIAL GOOGLE+ PAGE

New York Giants
Football team

The New York Giants are a professional American football team based in East Rutherford, New Jersey, representing the area. Wikipedia

Arena/Stadium: MetLife Stadium
Head coach: Tom Coughlin
Location: East Rutherford
Division: NFC East
NFL championships: 1986, 1990, 2007,
Nicknames: G-Men, Big Blue Wrecking Crew

13:23 31%
“What sort of Pokémon is Pikachu”
tap to edit

The answer is electric.

Input interpretation
80%
“Mount Everest”
tap to edit

I don't see any places matching 'Mount Everest'.
Sorry about that.

ピカチュウ (Pikachu)
25
electric

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

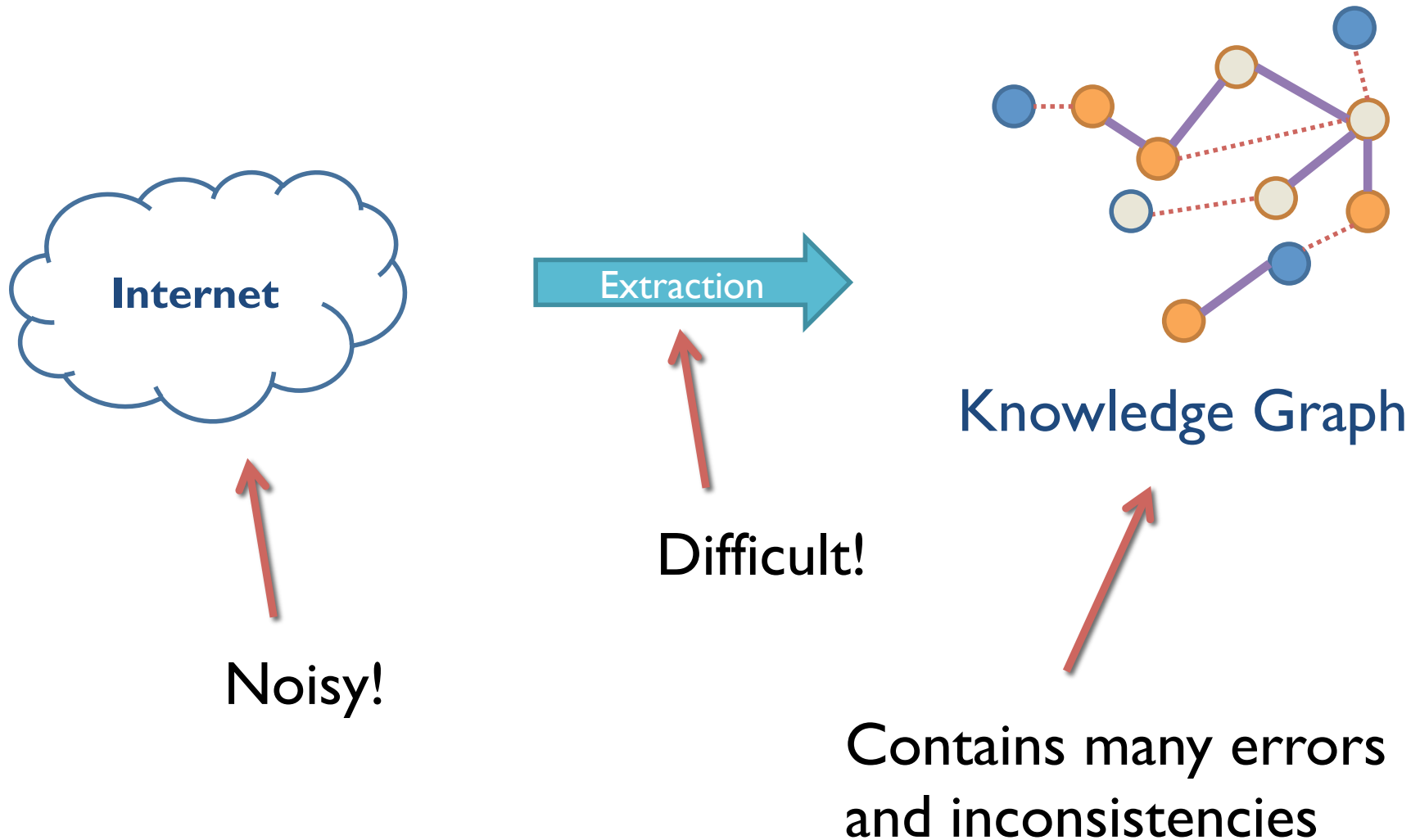
People I know who studied at **University of Maryland, College Park** · mumbai, College ...

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...

Motivating Problem: Real Challenges



Examples of NELL errors

Entity co-reference errors

Kyrgyzstan has many variants:

- Kyrgystan
- Kyrgistan
- Kyrghyzstan
- Kyrgyzstan
- 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 **Kyrgyzstan**, 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 ([/kɜrɡɪ'stɑːn/](#) *kur-gi-STAN*;^[5] [Kyrgyz](#): Кыргызстан (IPA: [\[qɯrɣɯs'stan\]](#)); [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

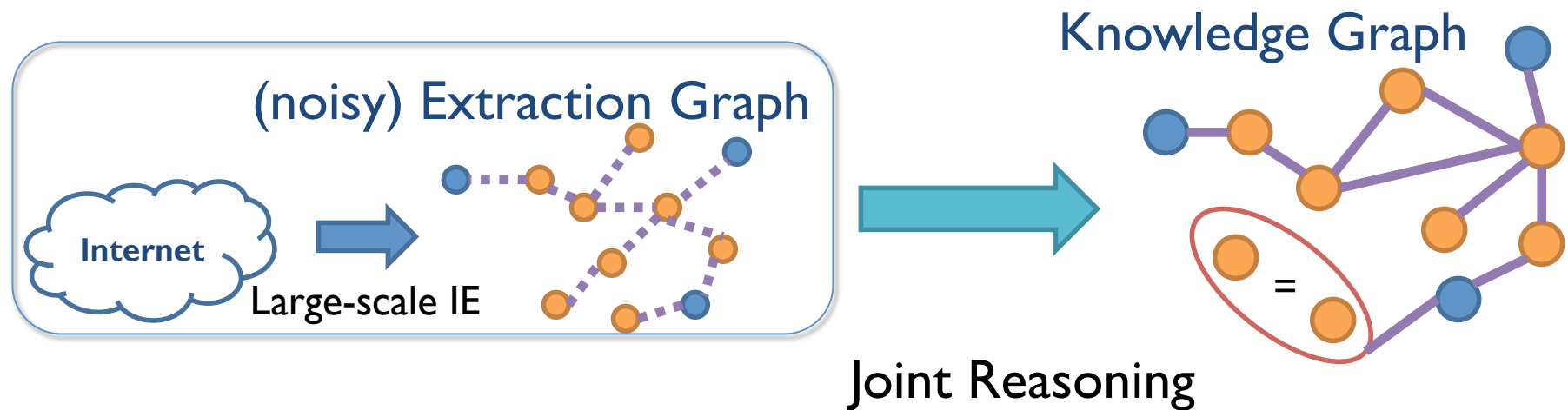
A Central Asian country of incredible natural beauty and proud nomadic traditions, most of Kyrgyzstan was formally annexed to Russia in 1876. The Kyrgyz staged a major revolt against the Tsarist Empire in 1916 in which almost one-sixth of the Kyrgyz population was killed. Kyrgyzstan became a Soviet republic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

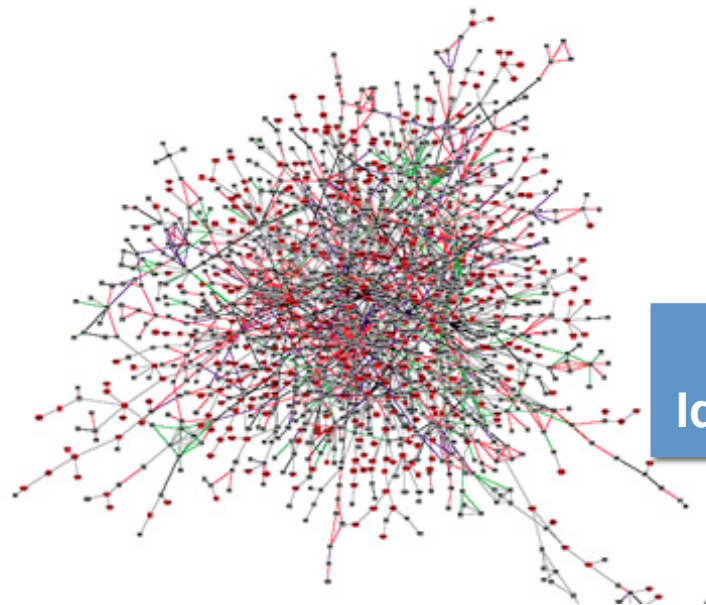
Enforcing these constraints requires **jointly** considering multiple extractions *across* documents

Motivating Problem (revised)



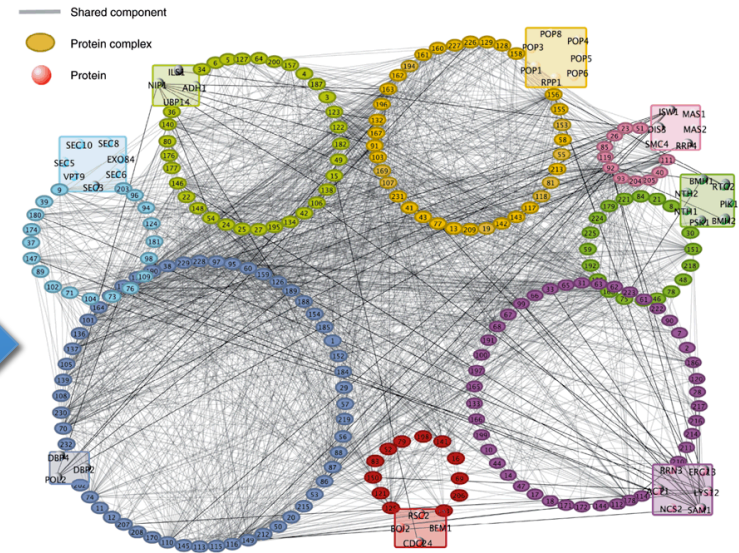
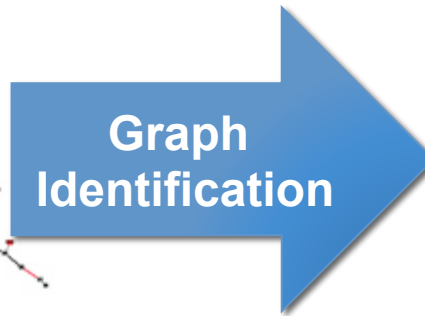
GRAPH IDENTIFICATION

Transformation



Input Graph

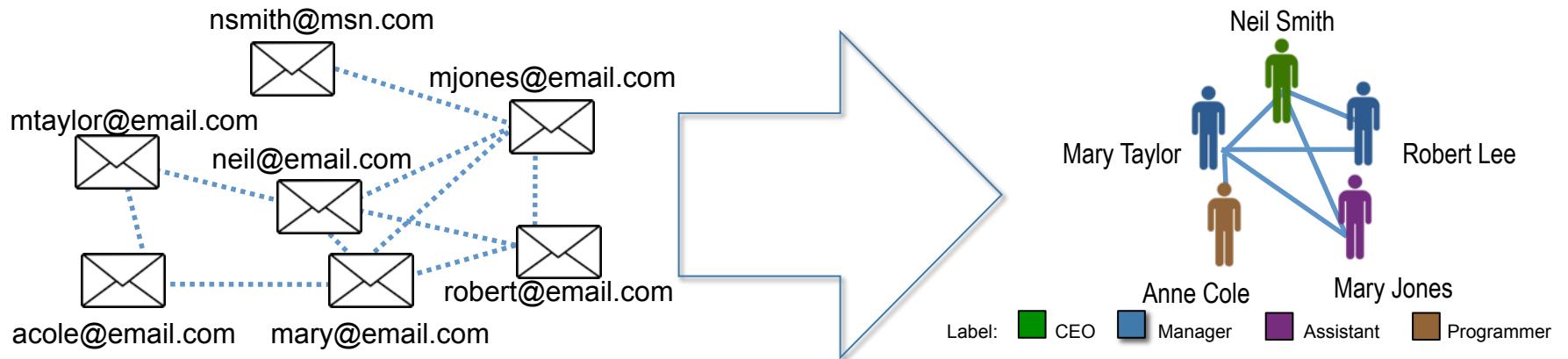
Available but inappropriate for analysis



Output Graph

Appropriate for further analysis

Motivation: Different Networks



Communication Network

Nodes: Email Address

Edges: Communication

Node Attributes: Words

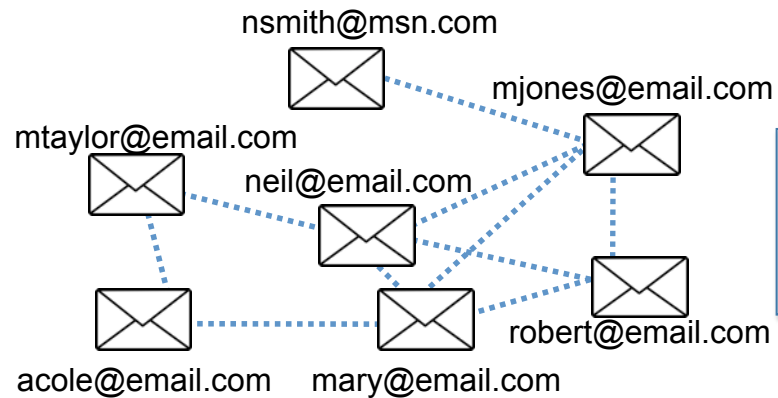
Organizational Network

Nodes: Person

Edges: Manages

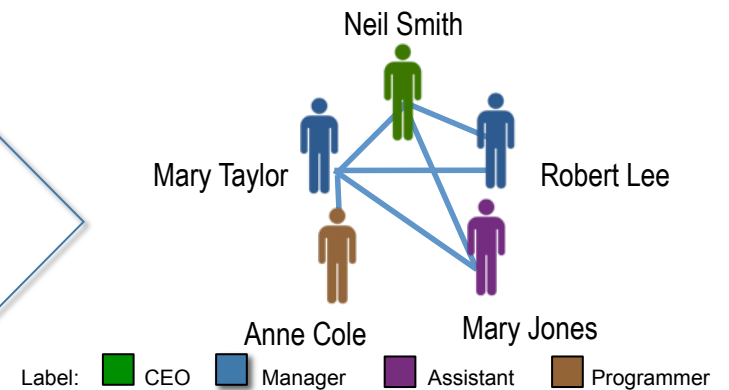
Node Labels: Title

Graph Identification



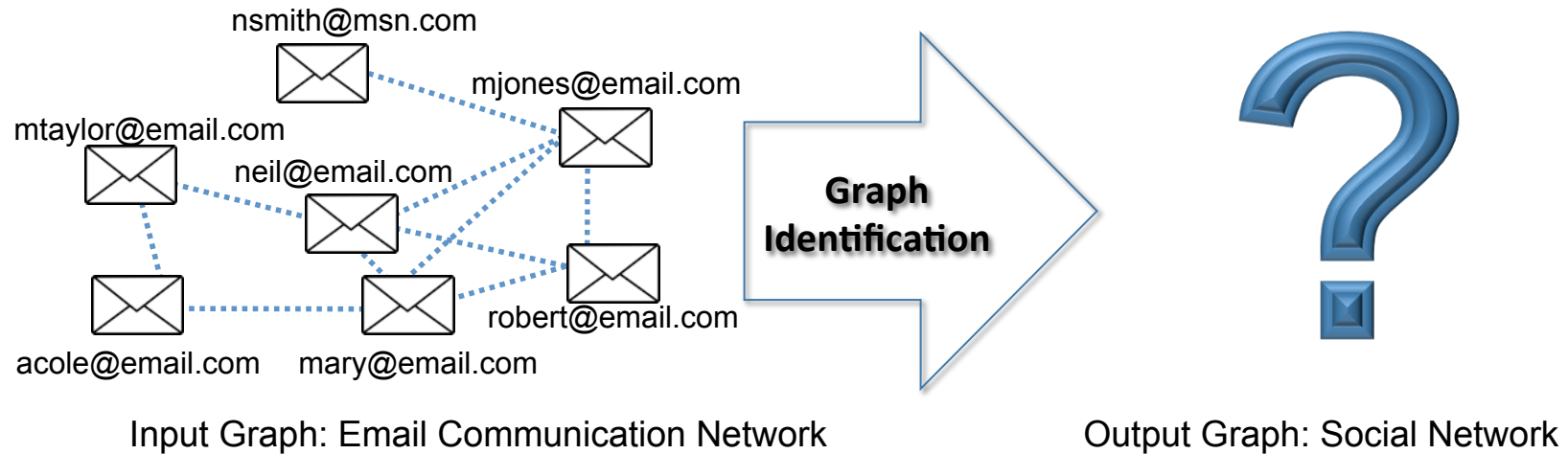
Input Graph: Email Communication Network

**Graph
Identification**



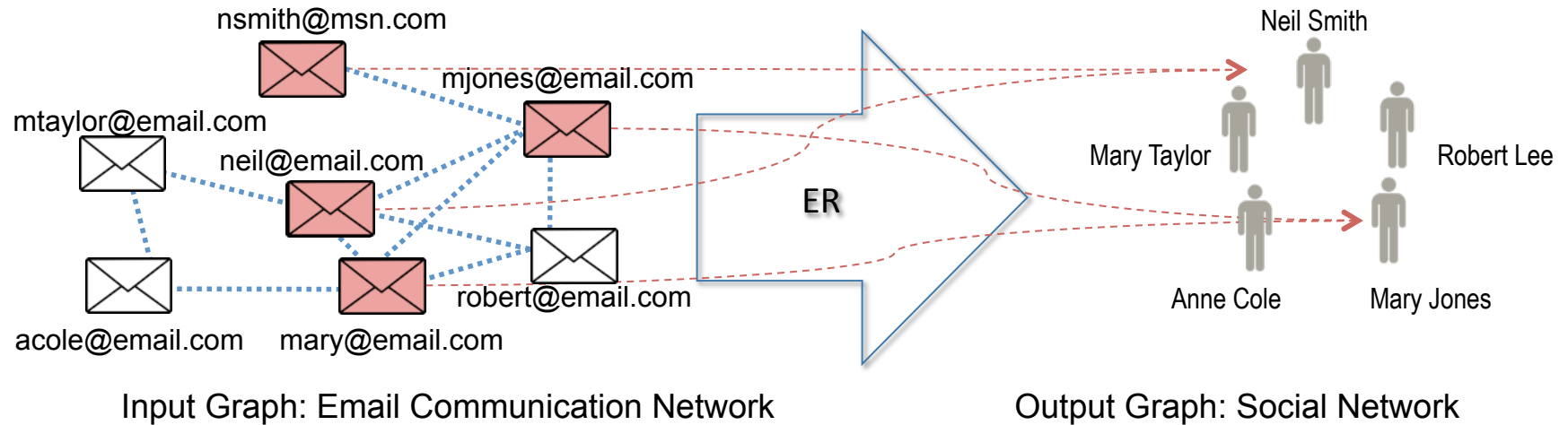
Output Graph: Social Network

Graph Identification



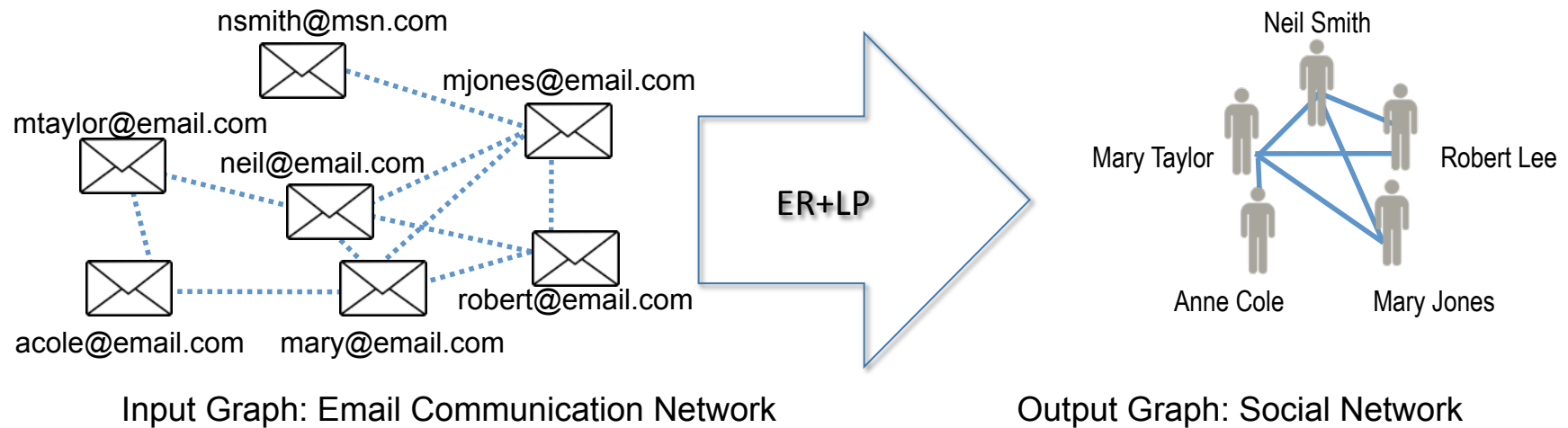
- What's involved?

Graph Identification



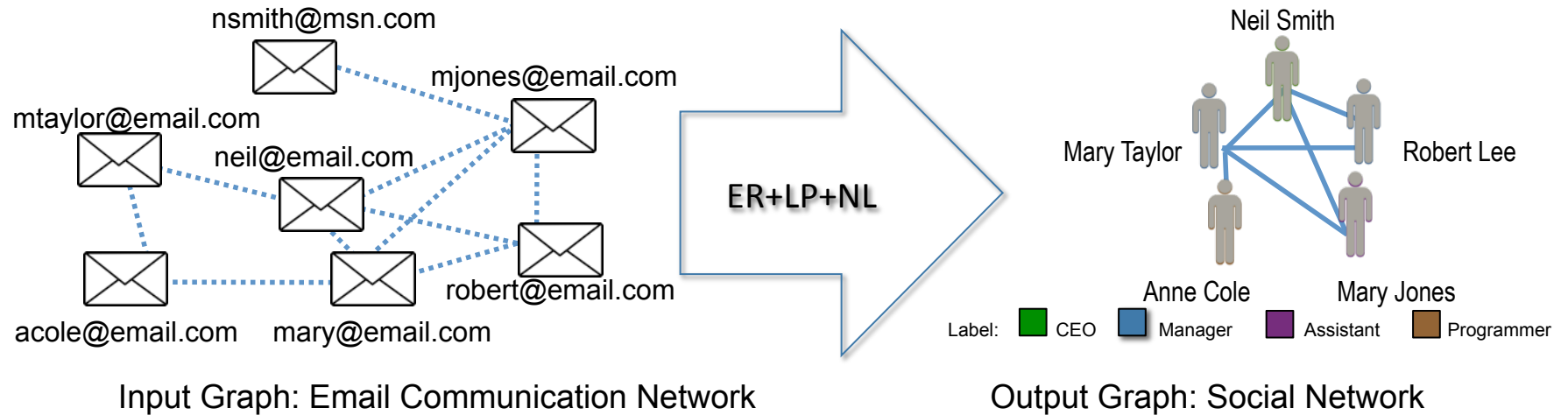
- What's involved?
 - Entity Resolution (ER): Map input graph nodes to output graph nodes

Graph Identification



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

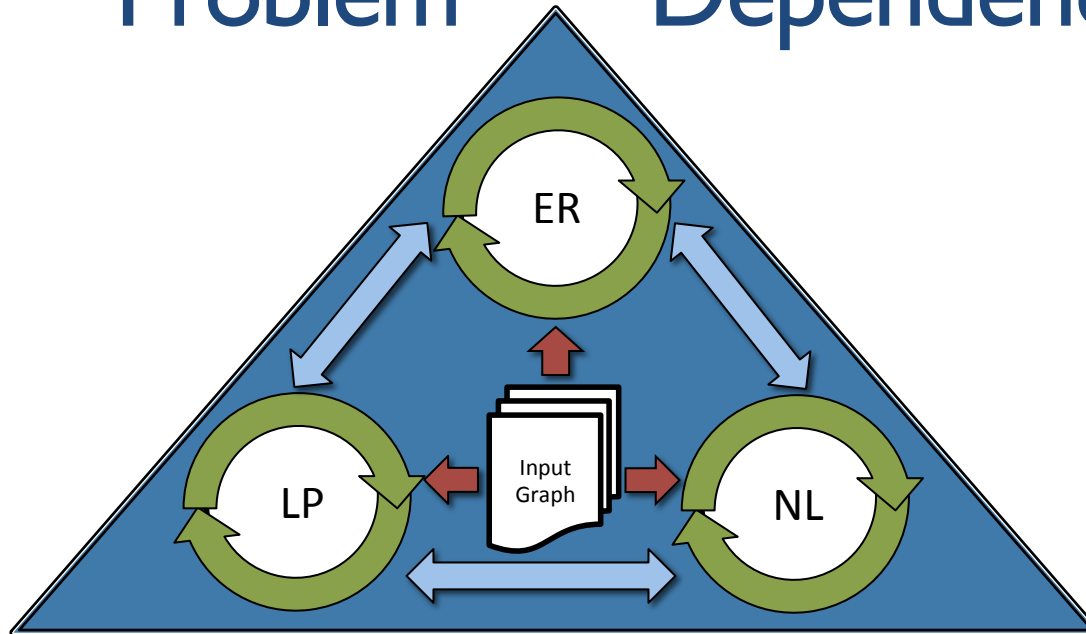
Graph Identification



- 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

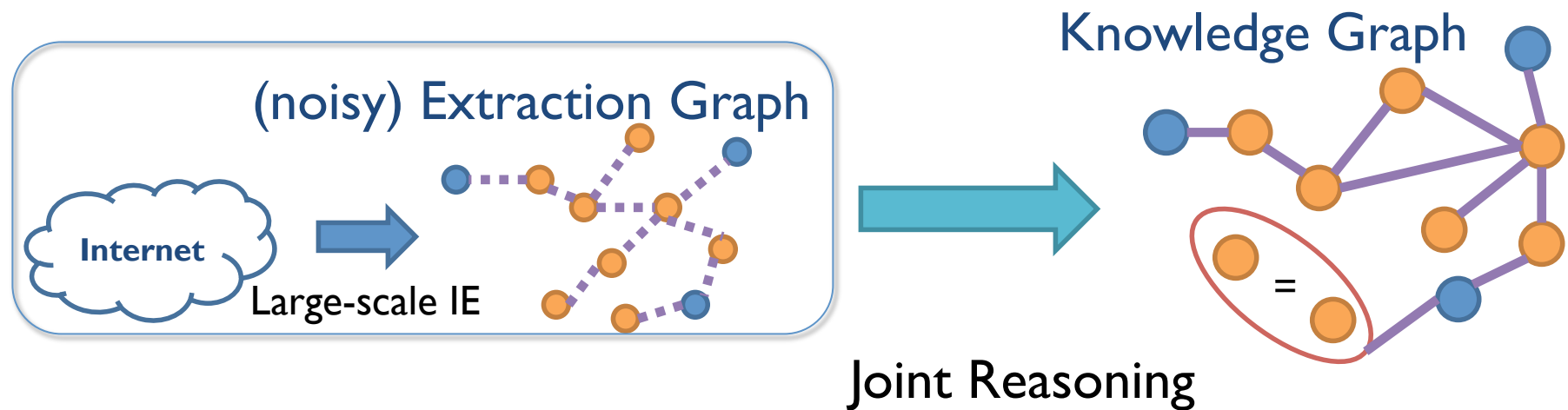


- Most work looks at these tasks in **isolation**
- In graph identification they are:
 - Evidence-Dependent – Inference depend on observed input graph
e.g., ER depends on input graph
 - Intra-Dependent – Inference within tasks are dependent
e.g., NL prediction depend on other NL predictions
 - Inter-Dependent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

KNOWLEDGE GRAPH IDENTIFICATION

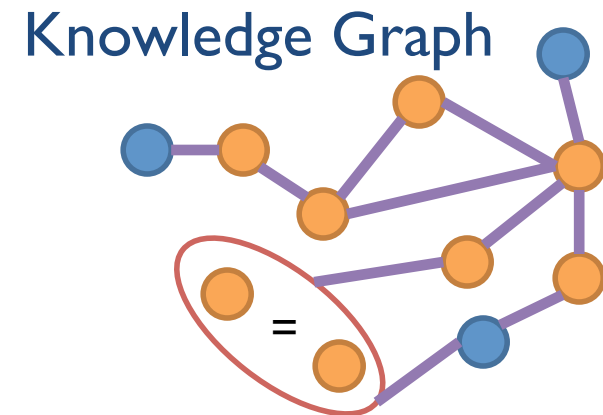
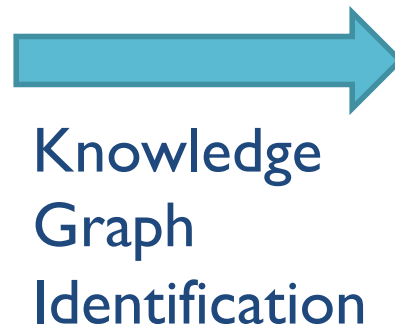
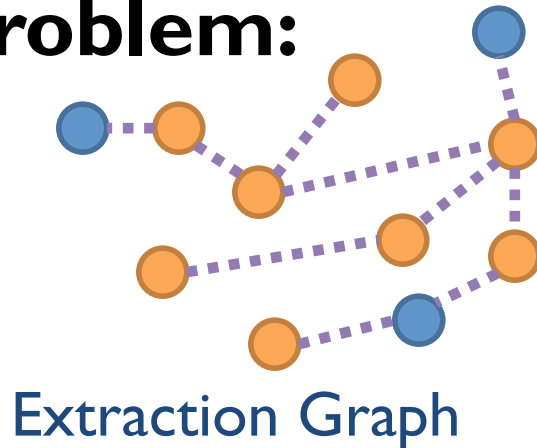
Pujara, Miao, Getoor, Cohen, ISWC 2013 (best student paper)

Motivating Problem (revised)



Knowledge Graph Identification

Problem:



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

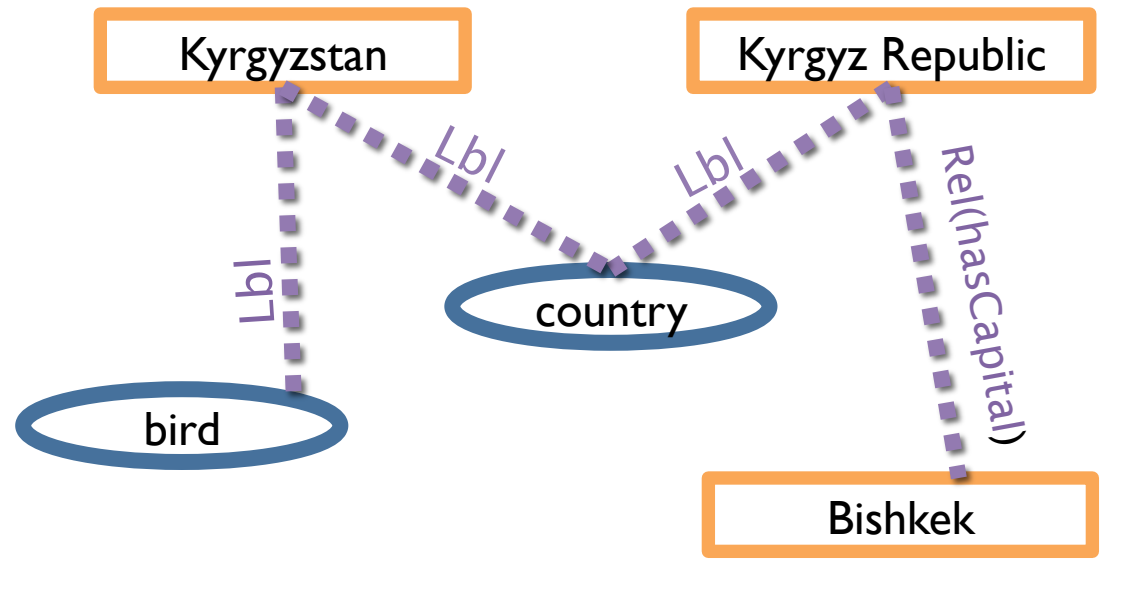


Illustration of KGI: Ontology + ER

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)

(Annotated) Extraction Graph

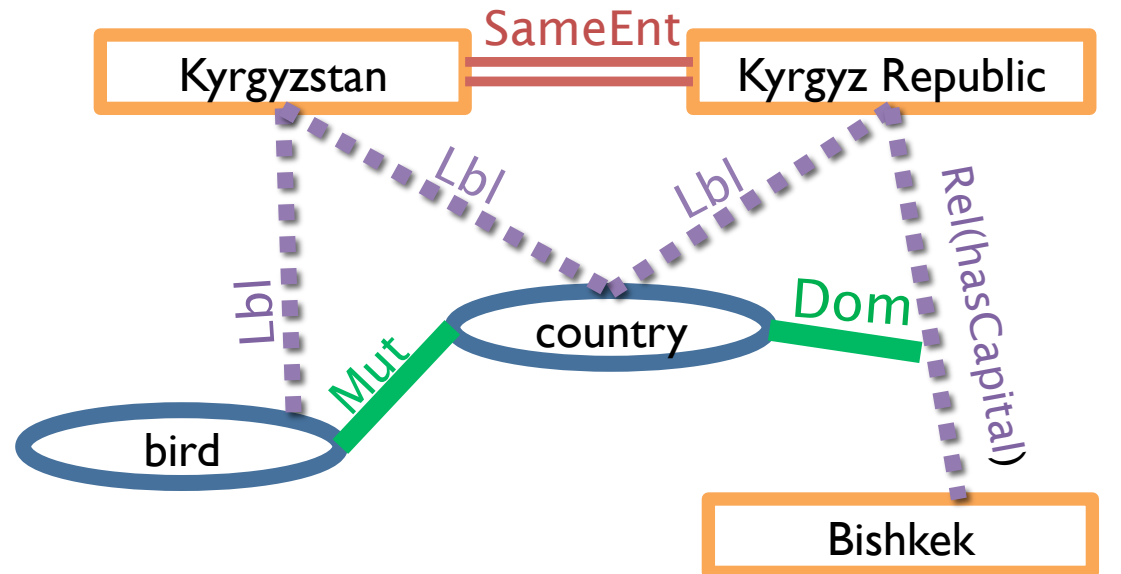


Illustration of KGI

Uncertain Extractions:

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

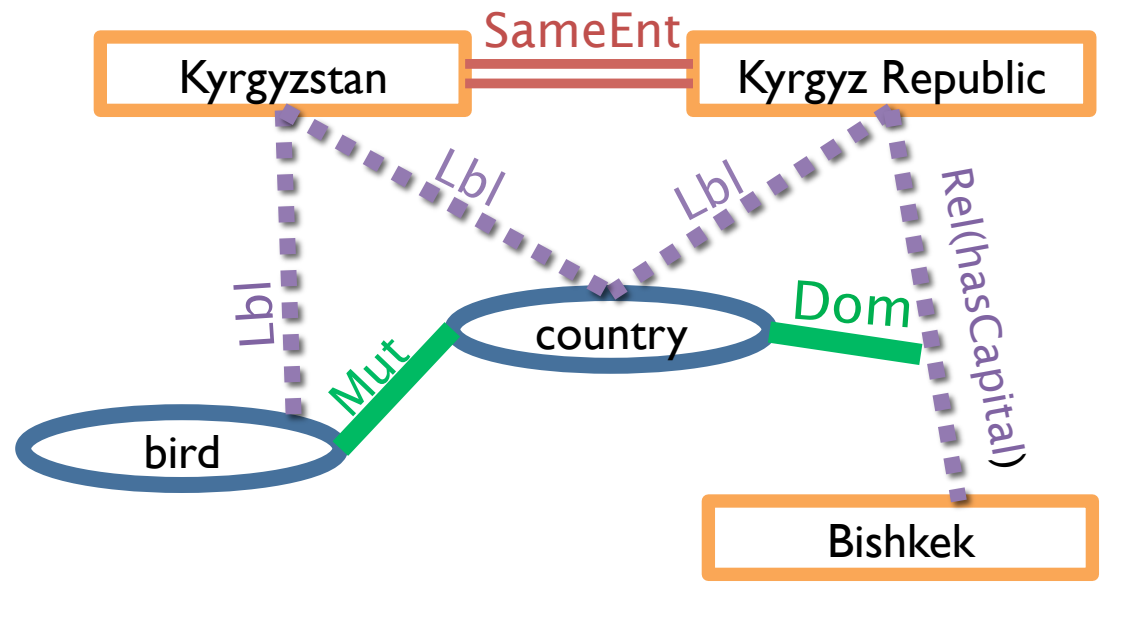
Ontology:

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

Entity Resolution:

SameEnt(Kyrgyz Republic, Kyrgyzstan)

(Annotated) Extraction Graph

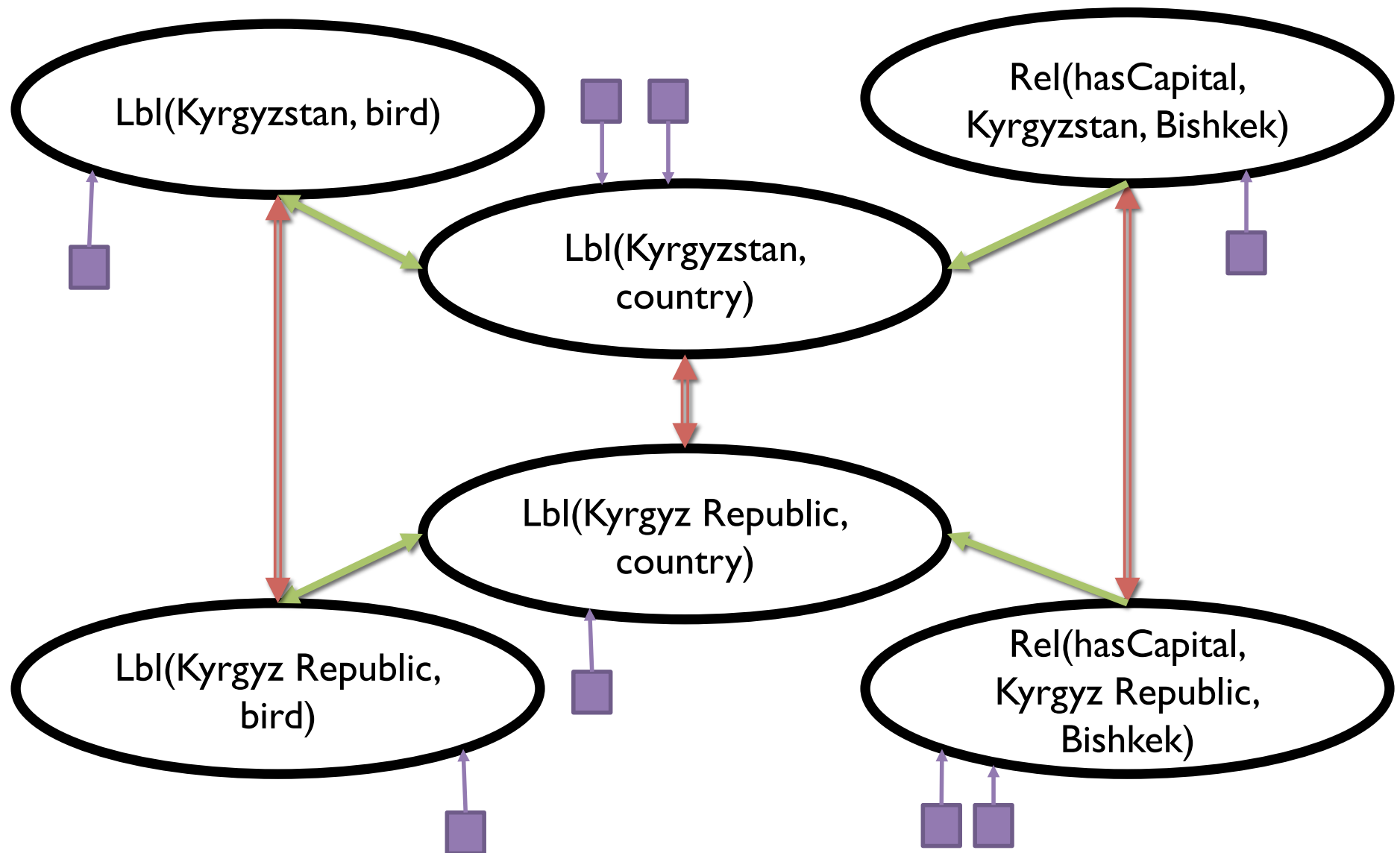


After Knowledge Graph Identification



Modeling Knowledge Graph Identification

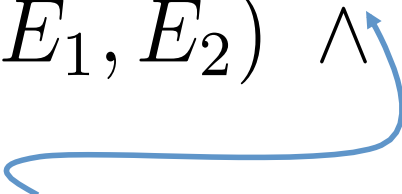
Viewing KGI as a probabilistic graphical model



Background: Probabilistic Soft Logic (PSL)

(Broecheler et al., UAI10; Kimming et al., NIPS-ProbProg12)

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

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


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

Background: PSL Rules to Distributions

- Rules are *grounded* by substituting literals into formulas

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

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

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

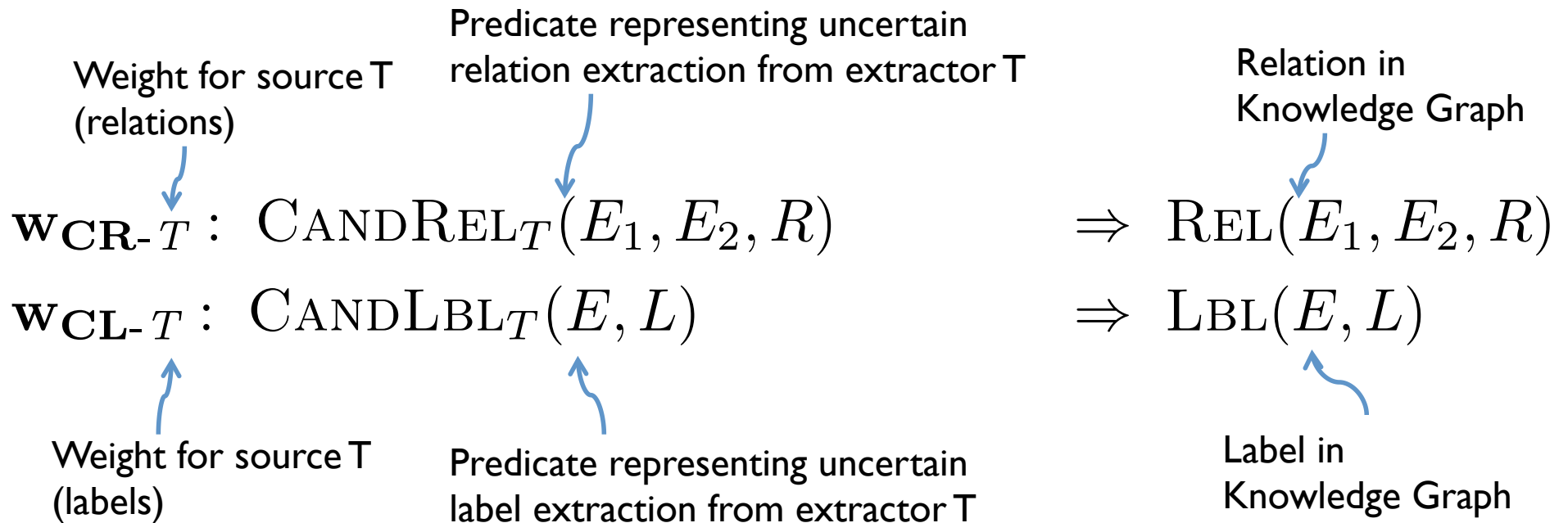
- The PSL program can be interpreted as a joint probability distribution over all variables in knowledge graph, conditioned on the extractions

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

$$\mathbf{w}_{\mathbf{EL}} : \text{SAMEENT}(E_1, E_2) \tilde{\wedge} \text{LBL}(E_1, L) \Rightarrow \text{LBL}(E_2, L)$$

$$\mathbf{w}_{\mathbf{ER}} : \text{SAMEENT}(E_1, E_2) \tilde{\wedge} \text{REL}(E_1, E, R) \Rightarrow \text{REL}(E_2, E, R)$$

$$\mathbf{w}_{\mathbf{ER}} : \text{SAMEENT}(E_1, E_2) \tilde{\wedge} \text{REL}(E, E_1, R) \Rightarrow \text{REL}(E, E_2, R)$$

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:

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

Selectional Preference:

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

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

Subsumption:

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

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

Mutual Exclusion:

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

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

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

$\Rightarrow \text{LBL}(\text{Kyrgyzstan}, \text{bird})$

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

$\Rightarrow \text{REL}(\text{Kyrgyz Rep.}, \text{Asia}, \text{locatedIn})$

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

$\wedge \text{LBL}(\text{Kyrgyz Rep.}, \text{country})$

$\Rightarrow \text{LBL}(\text{Kyrgyzstan}, \text{country})$

$[\phi_4] \text{DOM}(\text{locatedIn}, \text{country})$

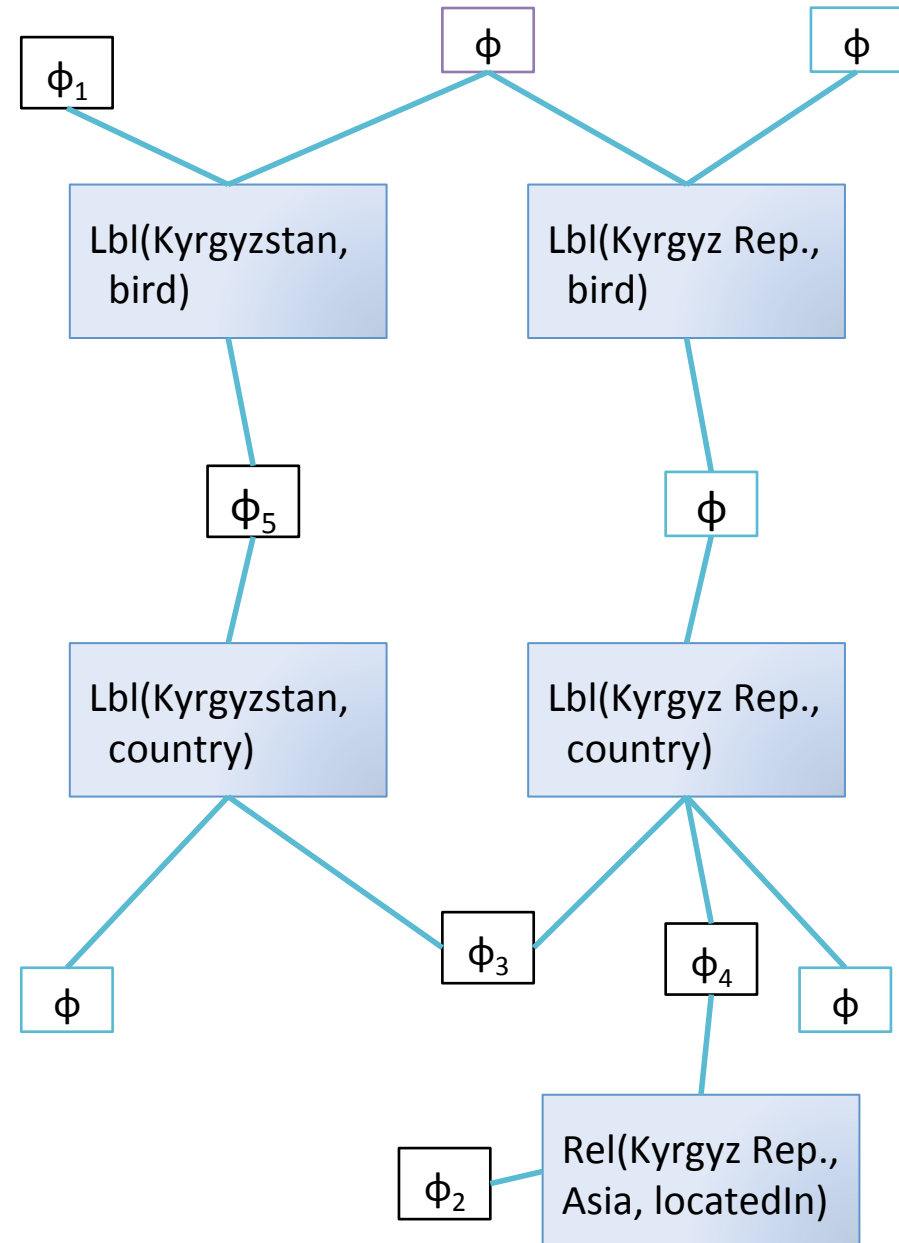
$\wedge \text{REL}(\text{Kyrgyz Rep.}, \text{Asia}, \text{locatedIn})$

$\Rightarrow \text{LBL}(\text{Kyrgyz Rep.}, \text{country})$

$[\phi_5] \text{MUT}(\text{country}, \text{bird})$

$\wedge \text{LBL}(\text{Kyrgyzstan}, \text{country})$

$\Rightarrow \neg \text{LBL}(\text{Kyrgyzstan}, \text{bird})$



Probability Distribution over KGs

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

CANDLBL_T(kyrgyzstan, bird)

⇒ LBL(kyrgyzstan, bird)

MUT(bird, country)

$\tilde{\wedge}$ LBL(kyrgyzstan, bird)

⇒ $\tilde{\neg}$ LBL(kyrgyzstan, country)

SAMEENT(kyrgyz republic, kyrgyzstan)

$\tilde{\wedge}$ LBL(kyrgyz republic, country)

⇒ LBL(kyrgyzstan, country)

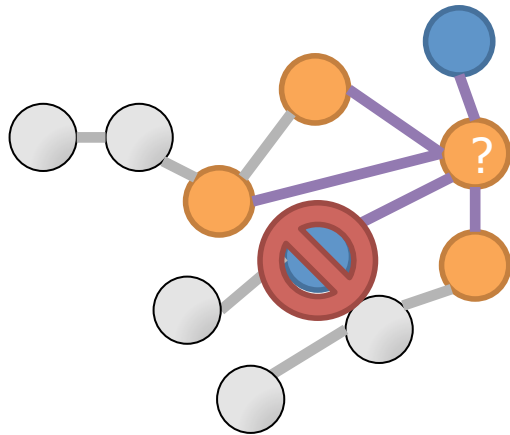
Evaluation

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

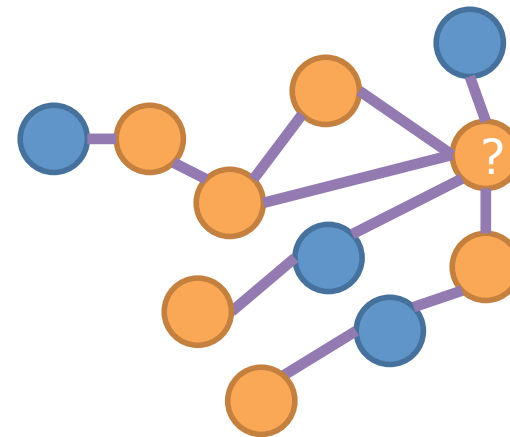
NELL Evaluation: two settings

Target Set: restrict to a subset of KG
(Jiang, ICDM12)



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

Complete: Infer full knowledge graph



- 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

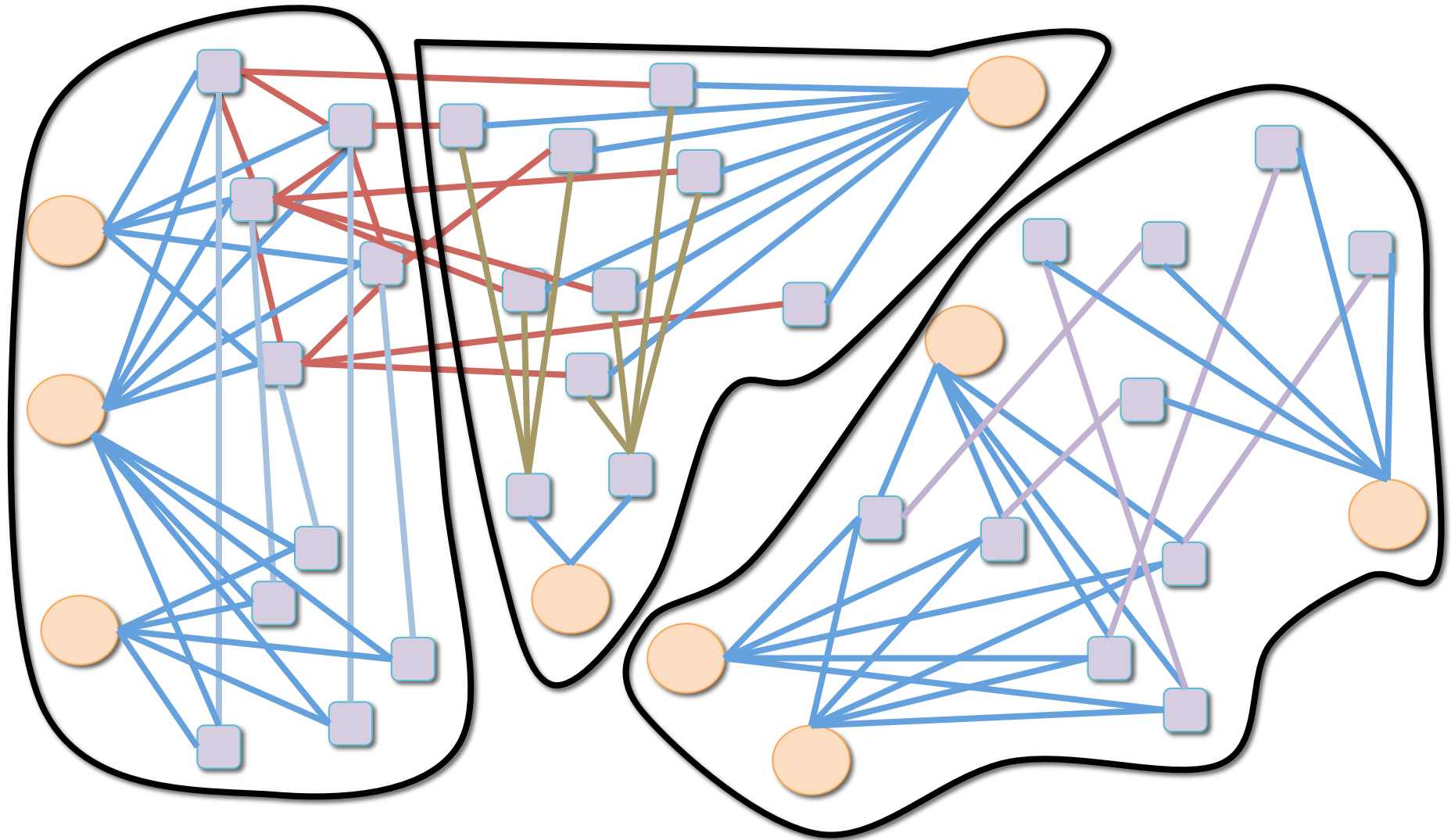
PSL-KGI Apply full Knowledge Graph Identification model

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

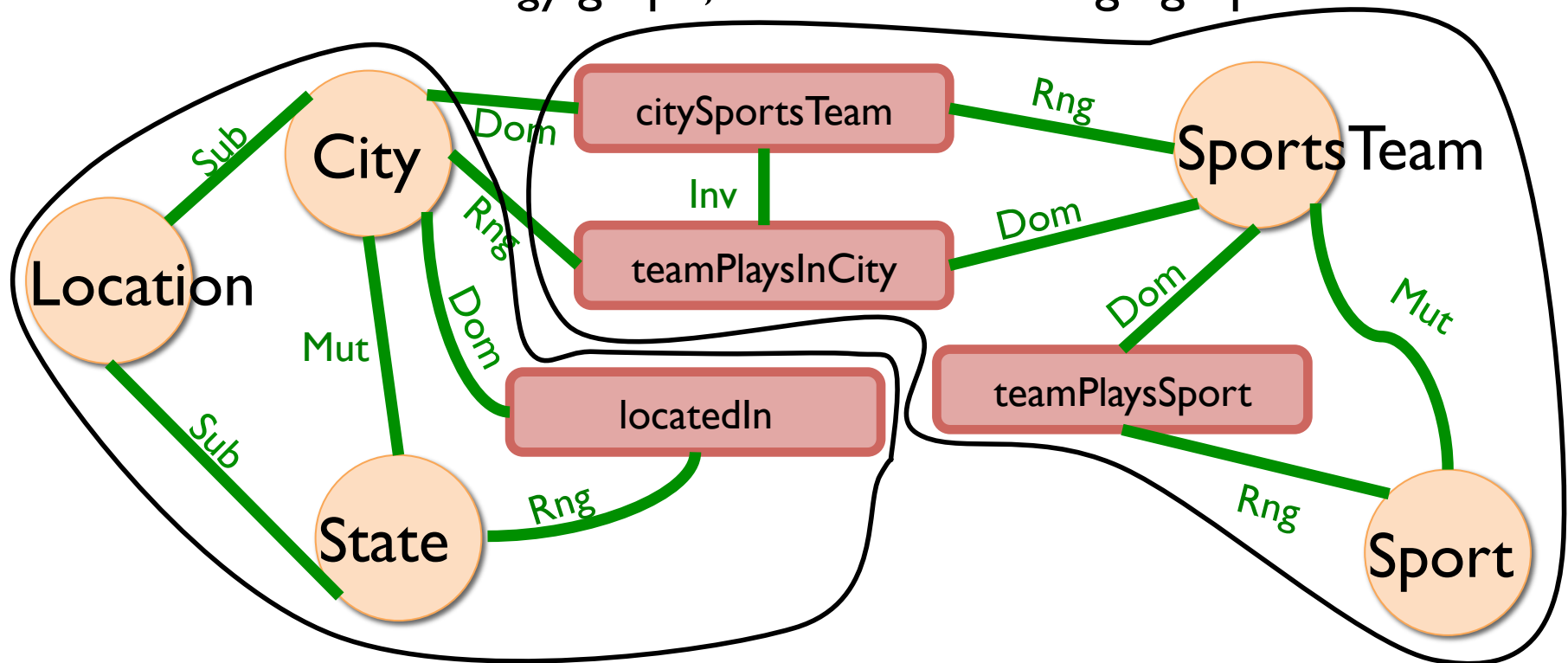
Ontology-Aware Partitioning

Problem: Partition the Knowledge Graph



Key idea: Ontology-aware partitioning

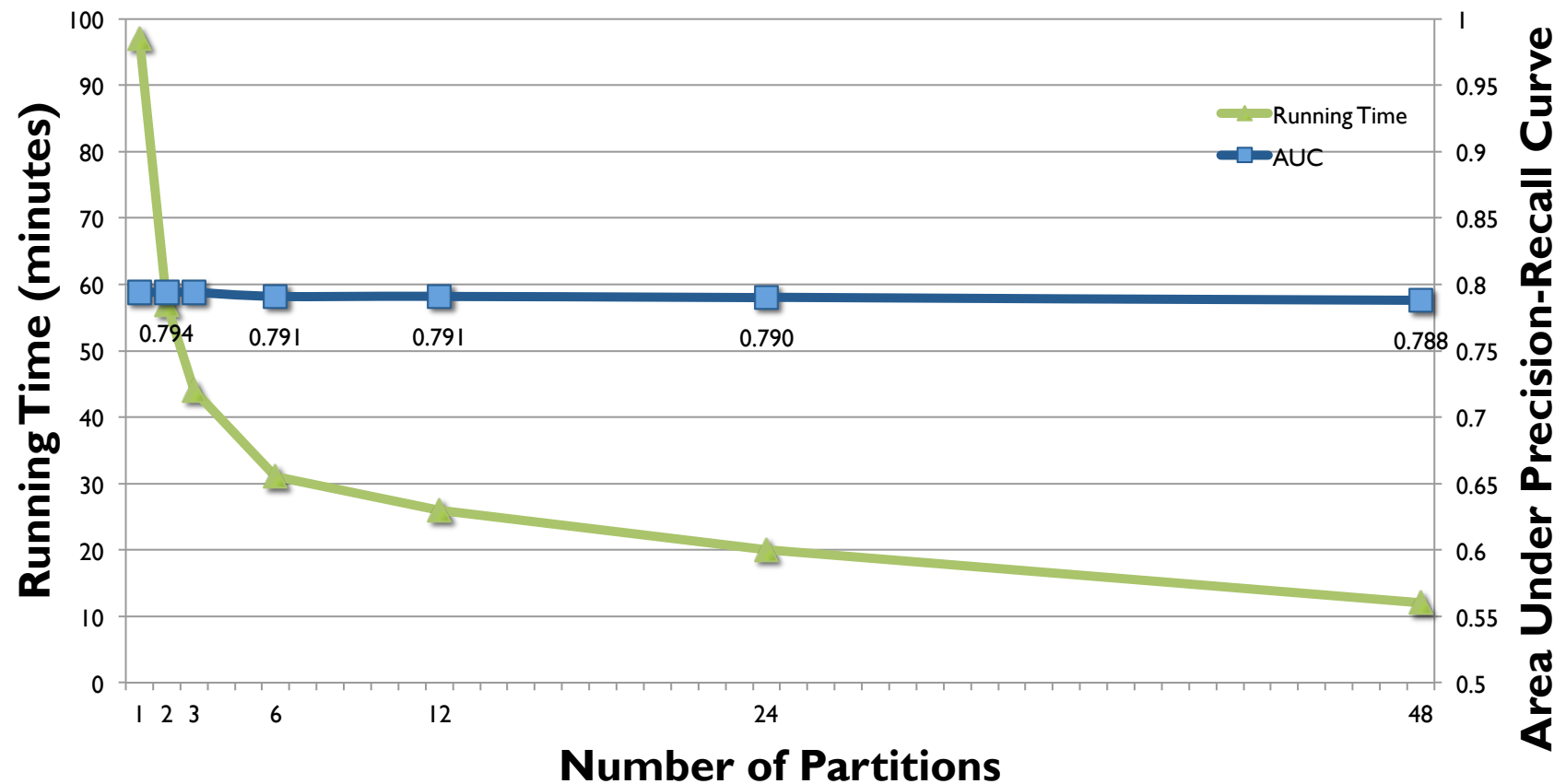
- Partition the *ontology* graph, not the knowledge graph



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

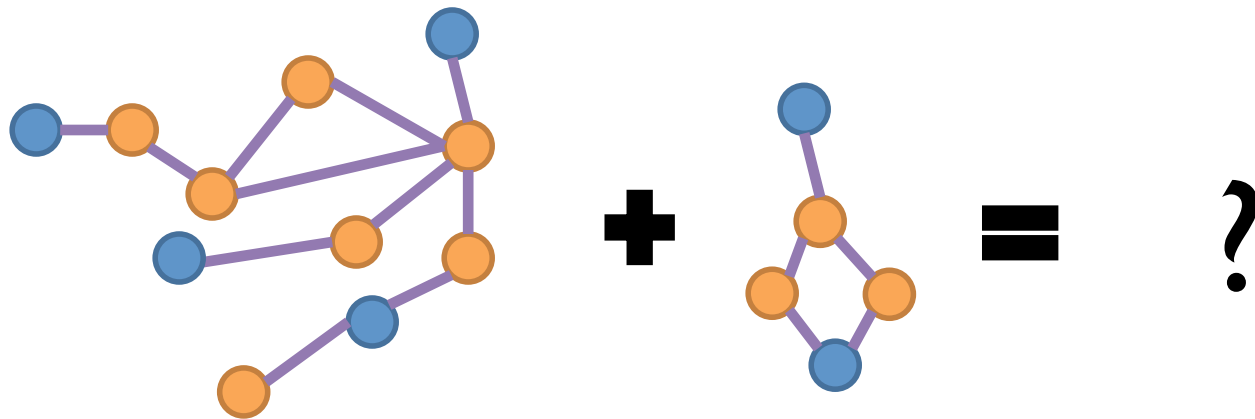
Experiments: Scalability

Partitions vs. Performance



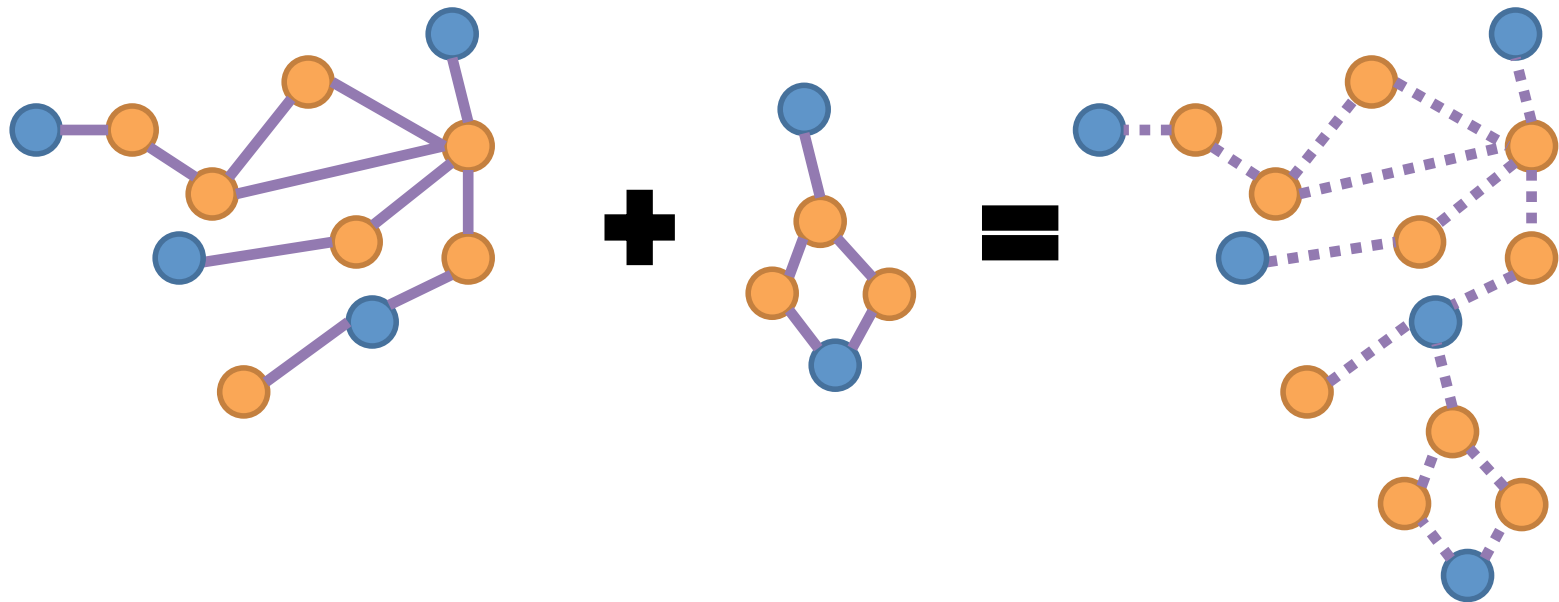
Dynamic Knowledge Graphs

Problem: Incremental Updates to KG

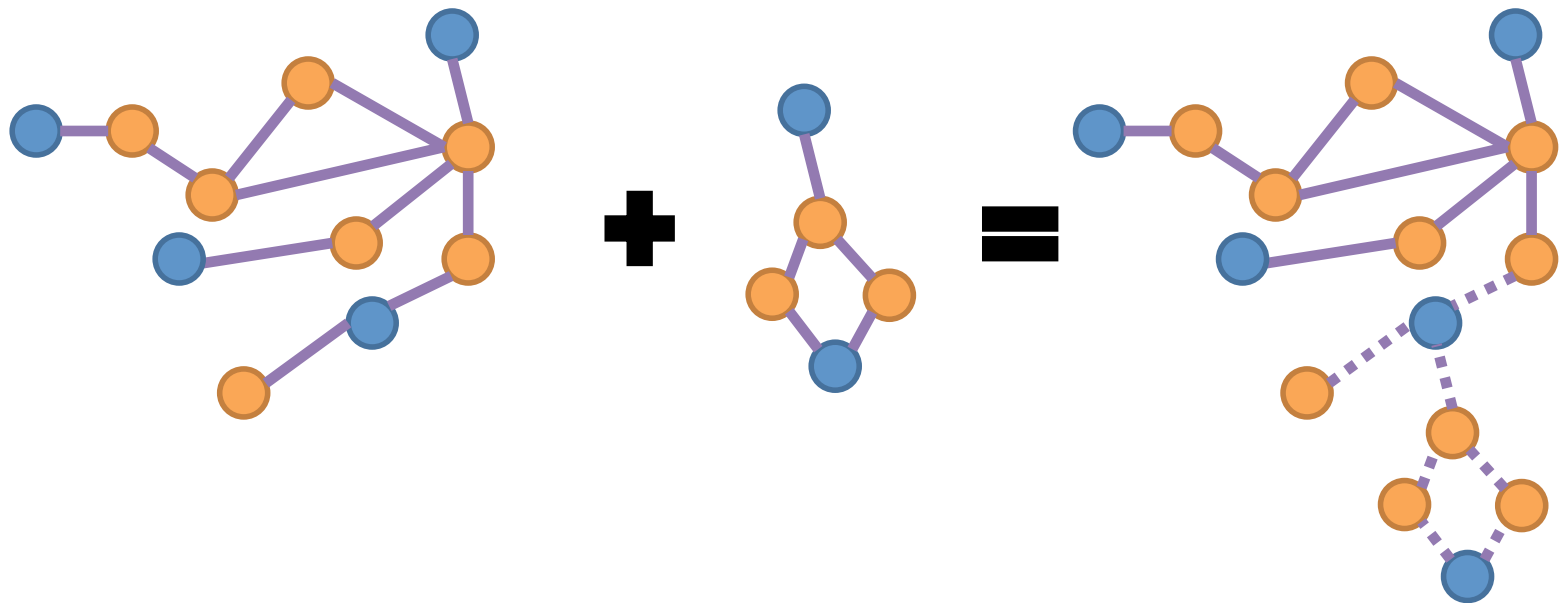


How do we add new extractions to the Knowledge Graph?

Naïve Approach: Full KGI over extractions



Approximation: KGI over subset of graph



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>

Questions?