



Penn  
UNIVERSITY OF PENNSYLVANIA

# Natural Language Understanding with Indirect Supervision

Daniel Khashabi

# Age of Big Data

- Big data:
  - Over 56 billion pages
  - Over 500 million tweets are sent every day.
  - Over 4 million blog posts are published on the Internet every day.
- Deep learning:
  - 1.5 billion parameters [Radford et al. 2019]
  - Super-human performance [Devlin et al. 2018]



# Troubling Observations

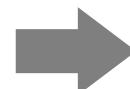
- Brittleness with respect to small changes

[K at al. 2016; Jia et al. 2017;  
Ribeiro et al. 2018; others]

**Context:** In the United States especially, several high-profile cases such as Debra LaFave, Pamela Rogers, and Mary Kay Letourneau have caused increased scrutiny on teacher misconduct.

**Question:** What has been the result of this publicity?

What's



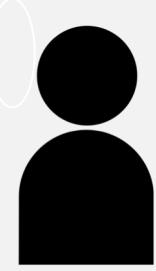
“increased scrutiny on teacher misconduct”



“teacher misconduct”

# Scenarios with Little (no?) Supervision

- Majority of our success has been on tasks w/ abundant annotations.
  - And tasks with little annotated data get the least attention.
- There will be settings where there is not “enough” direct supervision.
  - Unseen/unexpected scenarios.
  - Change of style, context, domain, etc.
  - These all result in vast space of possibilities for meanings.

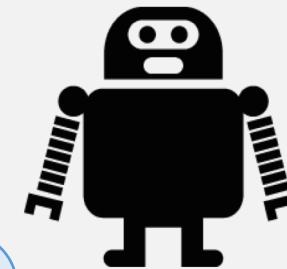


*Show me some restaurants nearby.*

*I don't like crowds.*

*Here are some options I found nearby:*

....



11:31 LTE

Show me restaurants nearby if I'm allergic to peanuts  
Tap to Edit >

Pattaya appears to serve peanuts and averages 3½ stars.

MAPS

**Pattaya**  
Thai · 800 feet  
★★★★★ (278) on Yelp · \$\$

**Bobby's Burger Palace**  
American (Traditional) · 450 feet  
★★★★★ (756) on Yelp · \$\$

**POD**  
Asian Fusion · 0.2 miles  
★★★★★ (731) on Yelp · \$\$\$

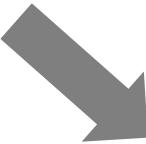
**Han Dynasty**  
Szechuan · 0.2 miles  
★★★★★ (658) on Yelp · \$\$

# Talk Statement

- It's unlikely that we will have directly “annotated” data that cover all aspects of natural language understanding.
- Data provides “hints” that exist independently of the task at hand.
- Weak signals can be amplified to produce higher quality signals.
  - Requires effective use of representation, knowledge and putting them together.



This talk



# “Supervision” vs “Data”

Minimal

$\{(x,y)\}$  →

(, "spam")

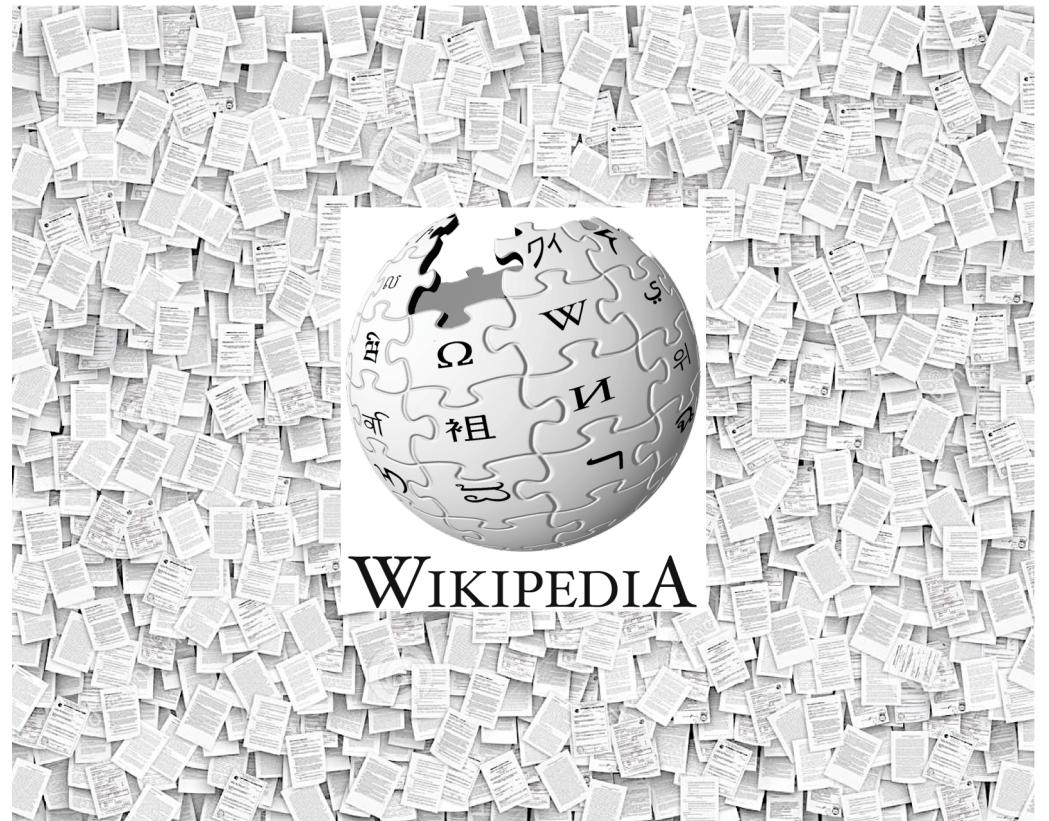
(, "ham")

(, "spam")

⋮

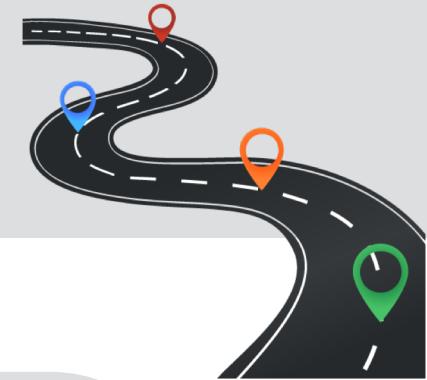
*Supervision  
(direct annotations)*

Abundant



*data*

# This talk



## ❑ Introduction

## ❑ Answering Questions

## ❑ Semantic Typing of Entities

## ❑ Future Work



*with minimal supervision*

- Representations
- Wikipedia
- Structure of the problem
- Compositionality
- Other learned models
- ...

# **ANSWERING QUESTIONS** *with minimal supervision*

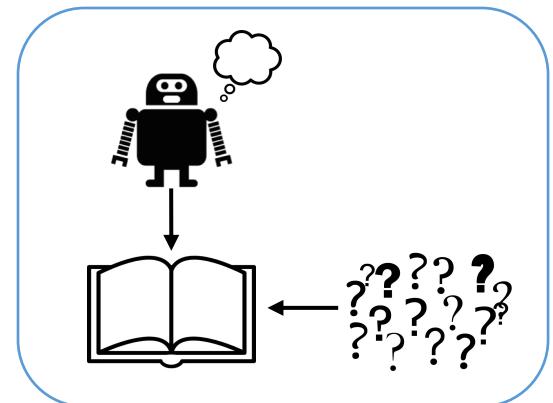
**K** et al. Question Answering as Global Reasoning over Semantic Abstractions. AAAI 18.

**K** et al. Question Answering via Integer Programming over Semi-Structured Knowledge. IJCAI 16.

Clark, EKSTTK. Combining Retrieval, Statistics, and Inference to Answer Elementary Science Questions. AAAI 16.

# Why Answering Questions?

- The grand goal: *Natural Language Understanding (NLU)*.
- Measuring progress by answering questions.
  - A system that is better at understanding language should have a higher chance of answering questions.
  - This has been used in the field for many years.  
[Winograd, 1972; Lehnert, 1977b; others]
  - Question Answering (QA), Reading Comprehension (RC), Textual Entailment (TE), etc.



# Answering Questions: The Setting

- Standardized science exams. [Clark et al. 2015]
- Simple language; machines require the ability to use the knowledge and abstract over it.



**Question:** *A bear survives winters with what structure?*



- (A) big ears
- (B) black nose
- (C) thick fur
- (D) brown eyes



Attached to each question is an evidence paragraph, potentially with the answer to the question.

# Linguistic Variability



**Question:** *A bear survives winters with what structure?*

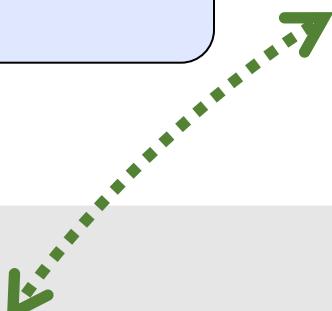


- (A) big ears
- (B) black nose
- (C) **thick fur**
- (D) brown eyes

Evidence  
paragraph



... and *a bear survives winters using its thick fur ...*



# Linguistic Variability



**Question:** *A bear survives winters with what structure?*



- (A) big ears
- (B) black nose
- (C) **thick fur**
- (D) brown eyes

Evidence  
paragraph



*... Polar bears, saved from the bitter cold by their thick fur coats, are among the animals in danger of extinction because of global warming and human activities. ...*

A given “meaning” can be phrased in many surface forms!

# Linguistic Variability



**Question:** A bear survives winters with what structure?



- (A) big ears
- (B) black nose
- (C) **thick fur**
- (D) brown eyes

**Evidence  
paragraph**

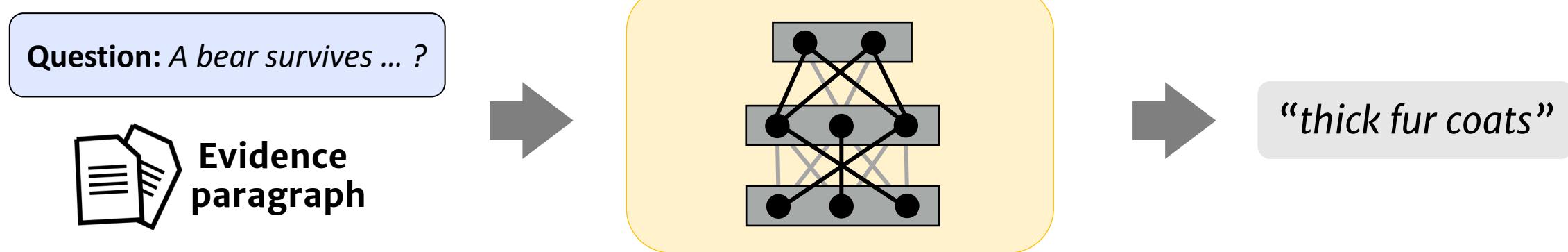


Polar bears have white fur so that they can camouflage into their environment. Their coat is so well camouflaged in Arctic environments that it can sometimes pass as a snow drift. They have a thick layer of body fat, which keeps them warm while swimming, and a double-layered coat that insulates them from the cold Arctic air.

**Polar bears, saved from the bitter cold by their thick fur coats, are among the animals in danger of extinction because of global warming and human activities.** Polar bears' lives depend wholly on the sea, their main source of food, and the place they spend most of their lives. But as the climate warms, that ice is melting, threatening polar bears. A common method of hunting by polar bears involves the bear keeping perfectly still by a seal's breathing hole, waiting for hours—or even days—for a seal to pop up for air.

# A Common Approach: Supervised Learning

- **Input:** question, an evidence paragraph.
- **Output:** predicted answer.



- Much success: Mostly with abundantly annotated data.
- Things can break down!



## Question: A bear survives winters with what structure?

Polar bears have white fur so that they can camouflage into their environment. Their coat is so well camouflaged in Arctic environments that it can sometimes pass as a snow drift. They have a thick layer of body fat, which keeps them warm while swimming, and a double-layered coat that insulates them from the cold Arctic air.

Polar bears, saved from the bitter cold by their thick fur coats, are among the animals in danger of extinction because of global warming and human activities. Polar bears' lives depend wholly on the sea, their main source of food, and the place they spend most of their lives. But as the climate warms, that ice is melting, threatening polar bears. A common method of hunting by polar bears involves the bear keeping perfectly still by a seal's breathing hole, waiting for hours—or even days—for a seal to pop up for air.

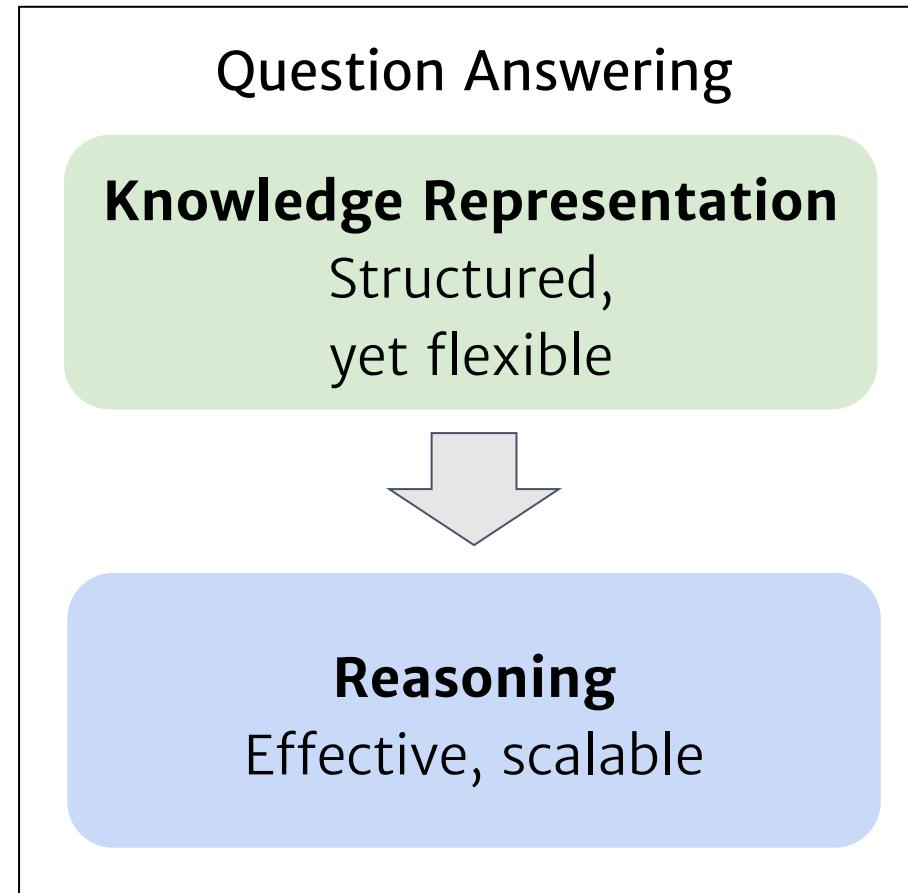
- Can we “explain” the decision?
- Can we “fix” such behaviors?

Predicted  
Answer

[Fetched on March 26, 2019]  
<https://demo.allennlp.org>  
[Seo et al, 17, Gardner et al, 18]

# Semi-Structured Inference: High-level View

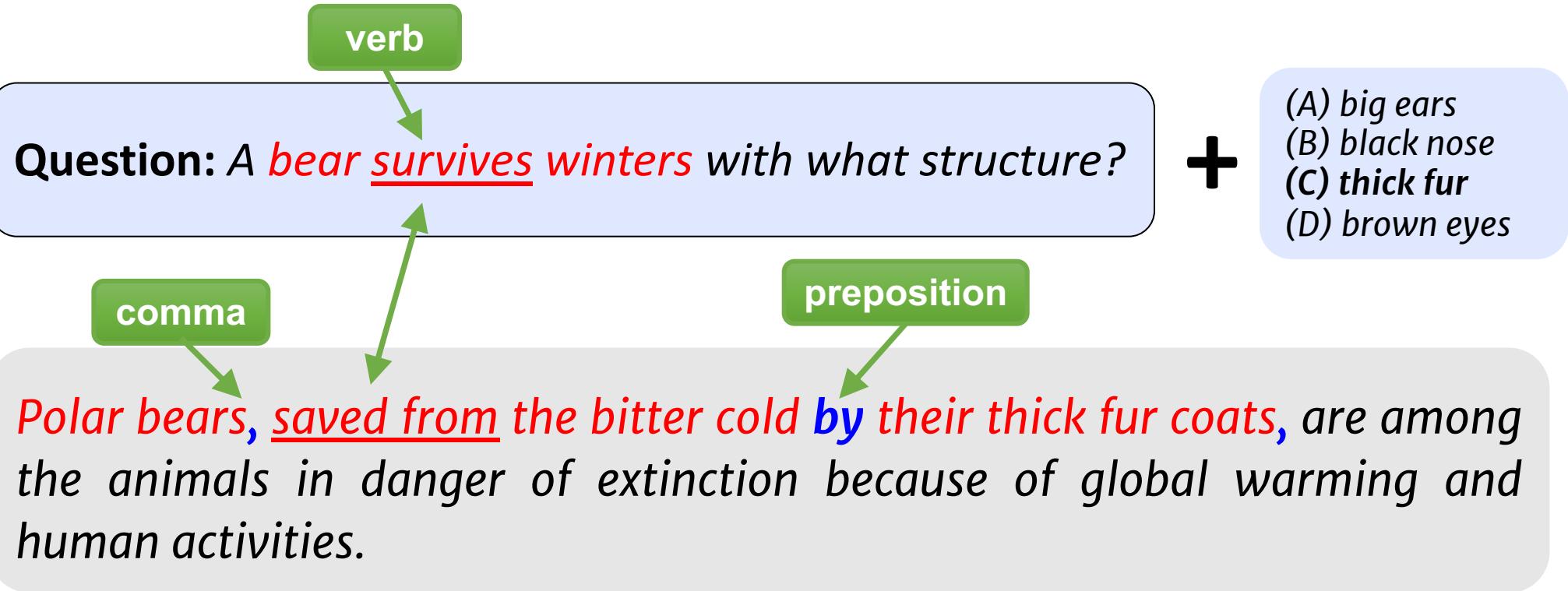
Question Answering  
as **Global Reasoning**  
over **Semi-Structured Knowledge**



# Language Understanding Phenomena



Evidence  
paragraph



QA is fundamentally a NLU problem

# “Lifting” Meaning as Semantic Graphs

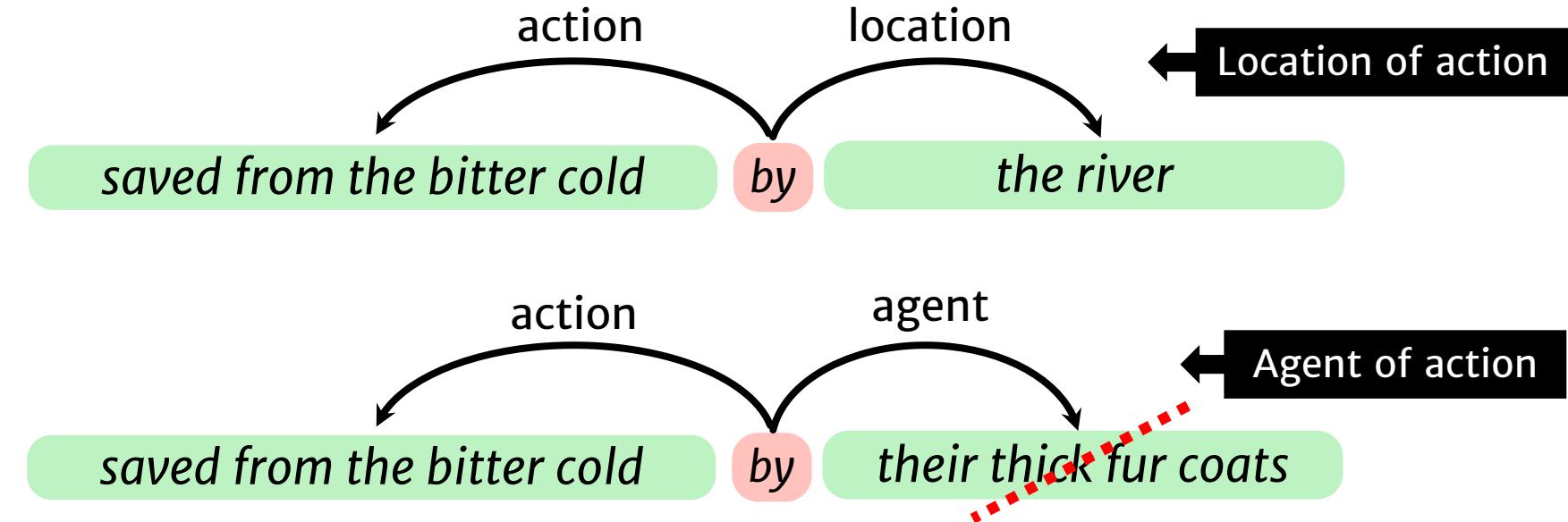
Oxford English Dictionary lists 8 primary meanings for “by”.

Disambiguation!

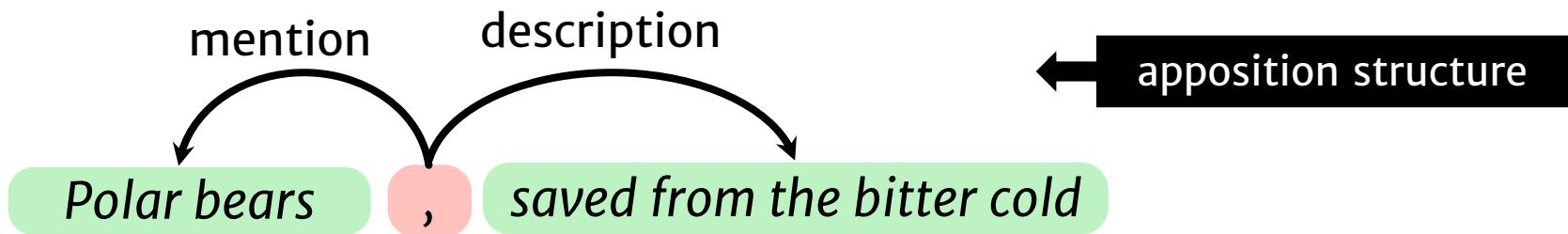
Evidence paragraph



... Polar bears, saved from the bitter cold by their thick fur coats, are among the animals in danger of extinction because of global warming and human activities.



# “Lifting” Meaning as Semantic Graphs



Evidence  
paragraph



... Polar bears , saved from the bitter cold by their thick fur coats, are among the animals in danger of extinction because of global warming and human activities.

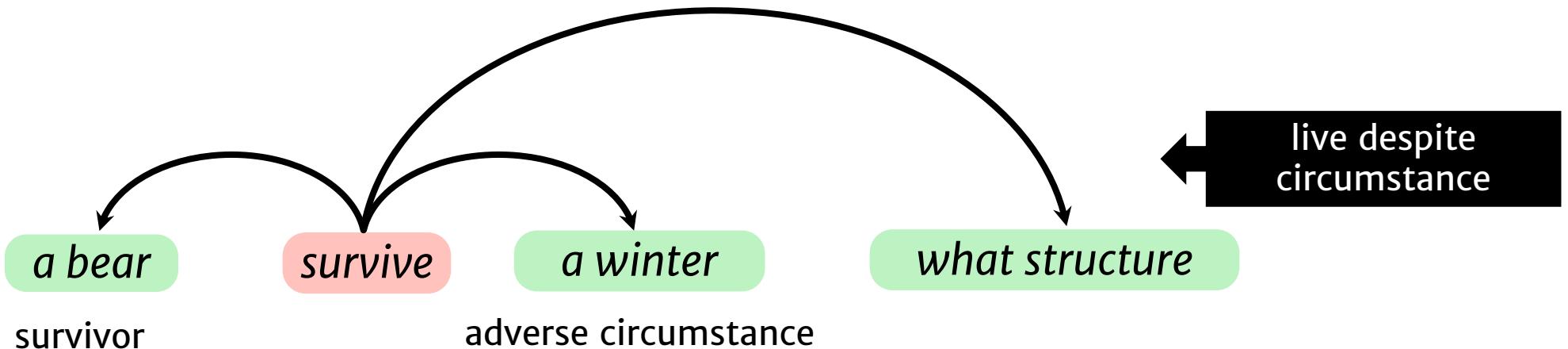
# “Lifting” Meaning as Semantic Graphs



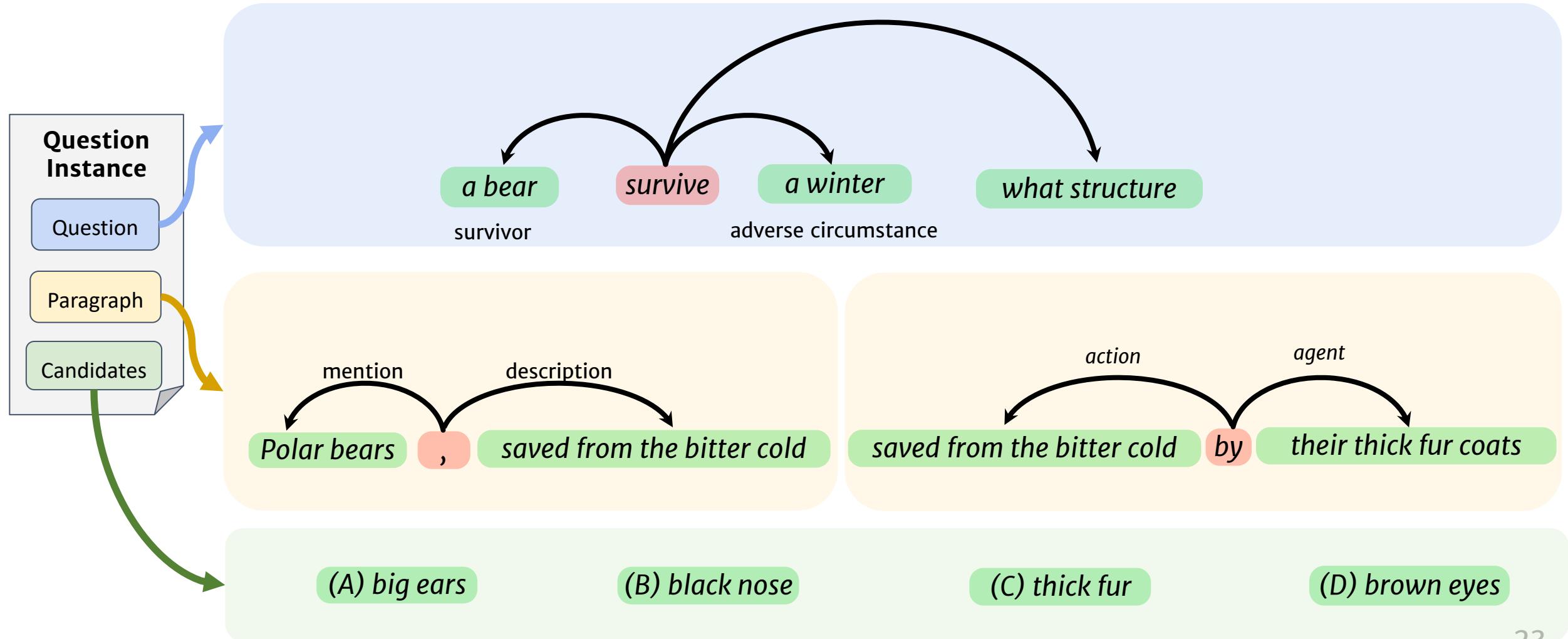
Question: A bear **survives** winters with what structure?

+

- (A) big ears
- (B) black nose
- (C) **thick fur**
- (D) brown eyes



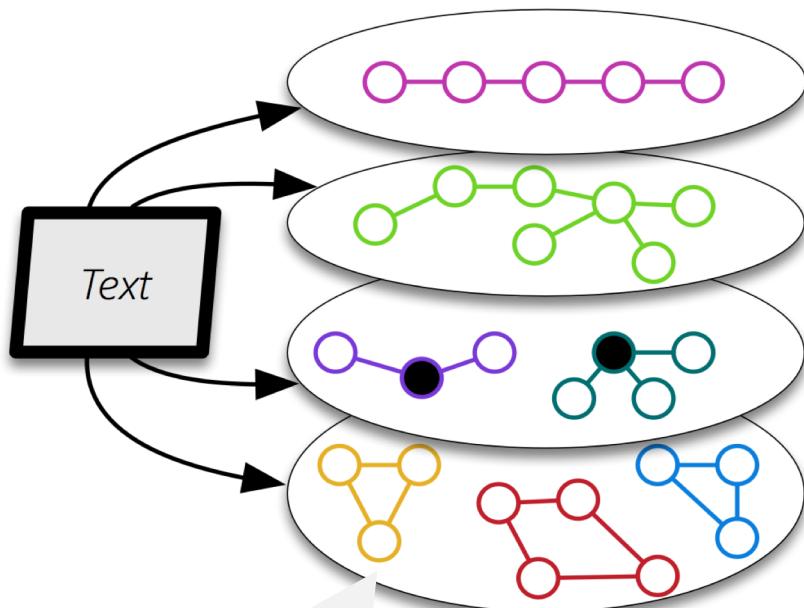
# Semantic Representations Altogether



# Collections of Semantic Graphs

- Create a **unified representation** of **families of graphs**

available in our  
software pipeline.  
K et al. LREC'18



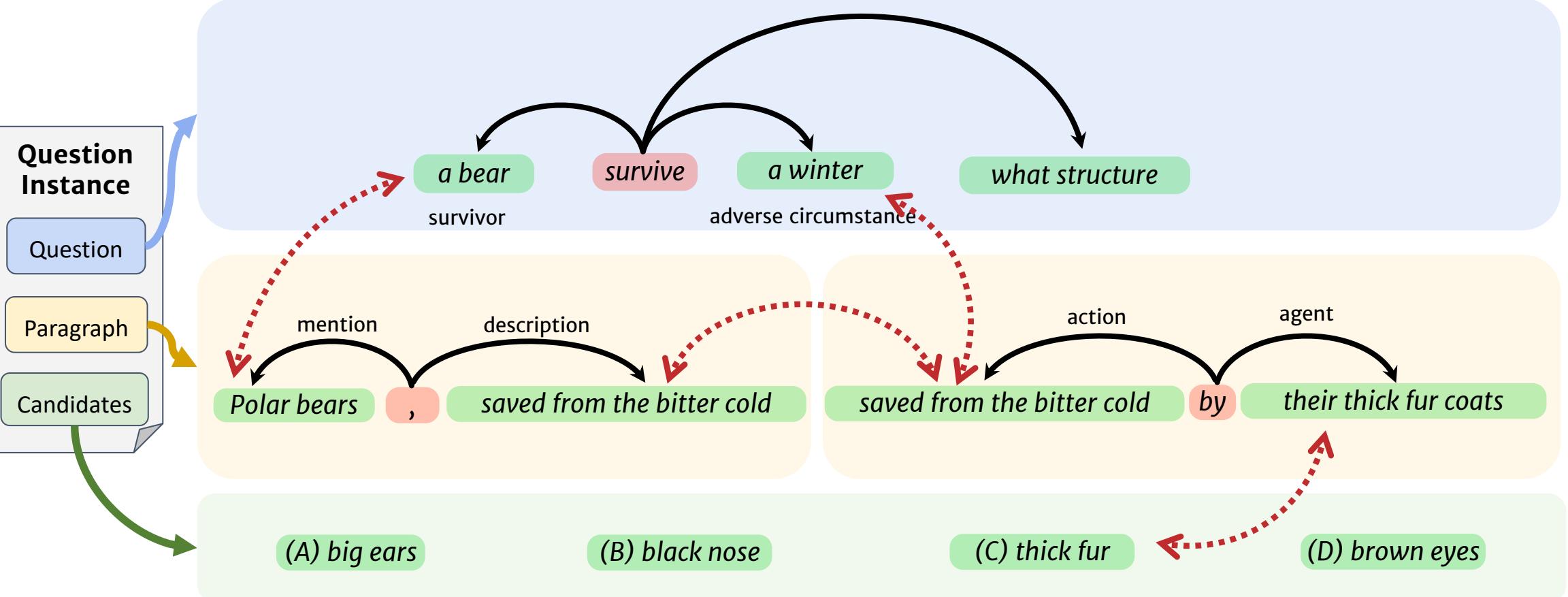
- Verb Semantic Roles [Punyakanok et al. 2008]
- Preposition Semantic Roles [Srikumar & Roth 2013]
- Comma Semantic Roles [Arivazhagan et al. 2016]
- Coreference [Chang et al. 2012]
- ...

Our representation is **not** QA-specific.  
It reflects our understanding of the language

- Surface word
- Semantic labels
- Other representation
- ...

Consequently, we expect these representations to  
be useful for a range of tasks

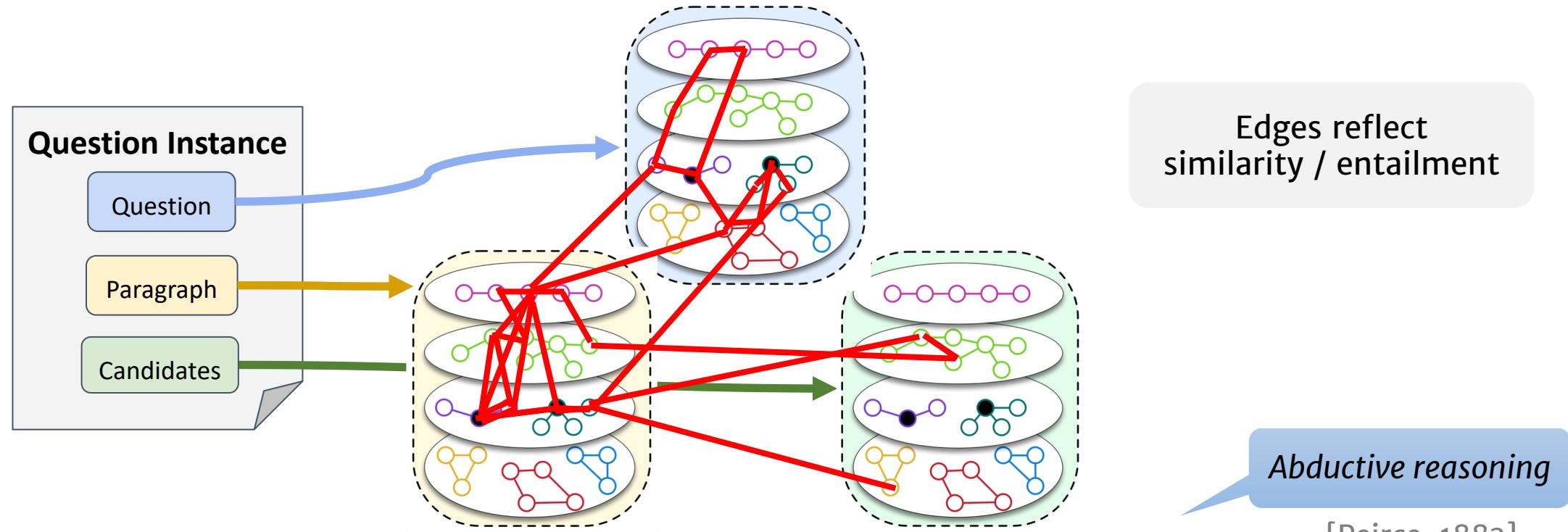
# Support Graph



Search for the best **Support Graph** connecting the Question to an Answer through the knowledge graph.

# Reasoning With a Meaning Representation

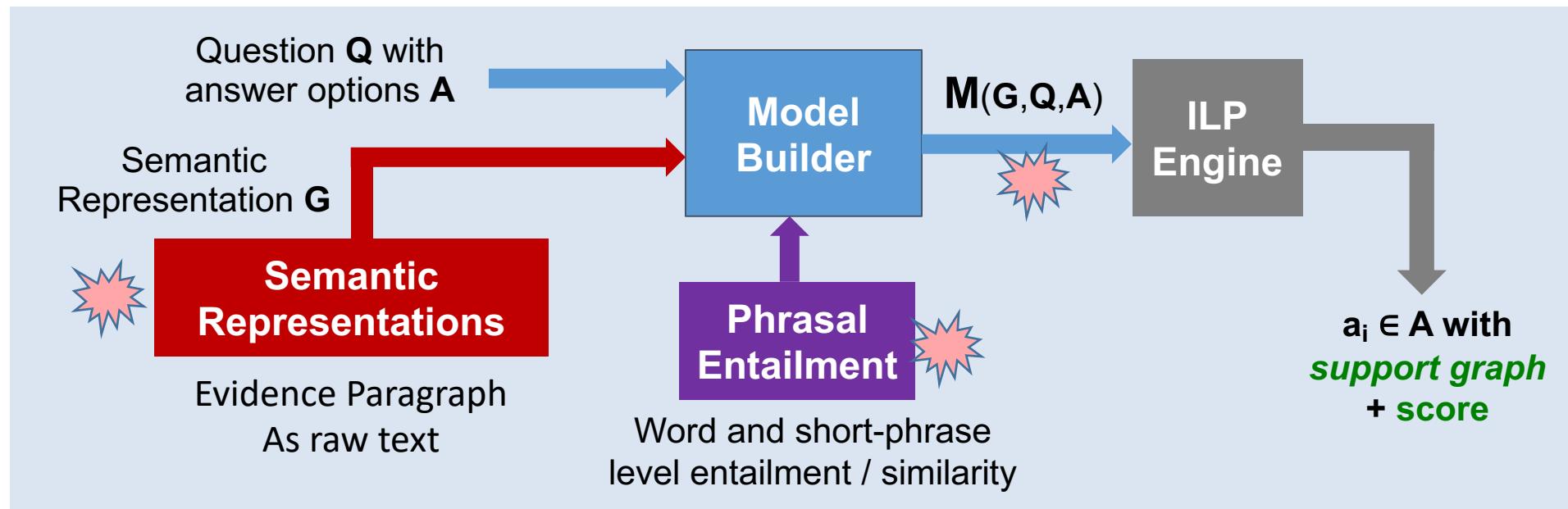
- **Support Graph** creates potential alignments between various semantic abstractions.



QA Reasoning formulated as finding “best” explanation – subgraph connecting the Question to the Answers via the Knowledge

# Framework Overview

- A discrete **optimization** approach to QA for multiple-choice questions



$M(G, Q, A)$

$$\max \sum_i c_i x_i$$
$$\forall x_i \in \mathbb{N} \cup \{0\}$$
$$\begin{cases} \sum_i a_{1i} x_i \leq b_1 \\ \dots \\ \sum_i a_{ki} x_i \leq b_k \end{cases}$$

Optimization using Integer Linear Program (**ILP**) formalism

# ILP Model: Design Challenges

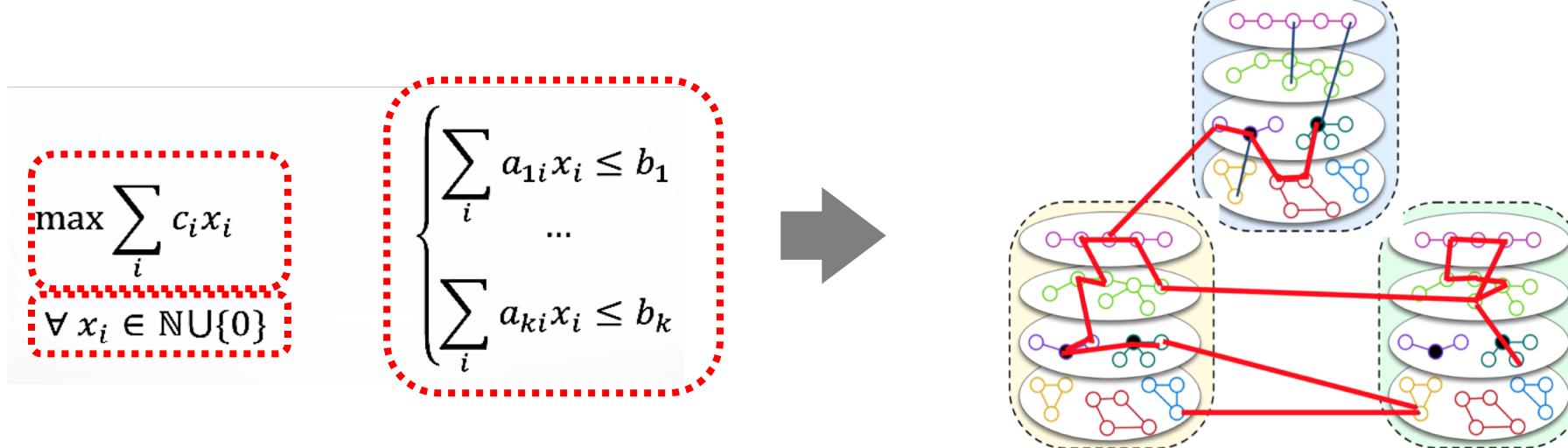
***Goal:*** Design ILP objective function, s.t. maximizing it subject to the constraints yields a “desirable” support graph

Not so straightforward!

$$\begin{aligned} \max \sum_i c_i x_i \\ \forall x_i \in \mathbb{N} \cup \{0\} \end{aligned} \quad \left\{ \begin{array}{l} \sum_i a_{1i} x_i \leq b_1 \\ \dots \\ \sum_i a_{ki} x_i \leq b_k \end{array} \right.$$

- Many possible “proof structures”
- Imperfect lexical “similarity” blackbox
- Partial or missing knowledge
- Question logic (negation, conjunction, comparison)
- Scalability of ILP solvers
- ...

# ILP Model: Some Details

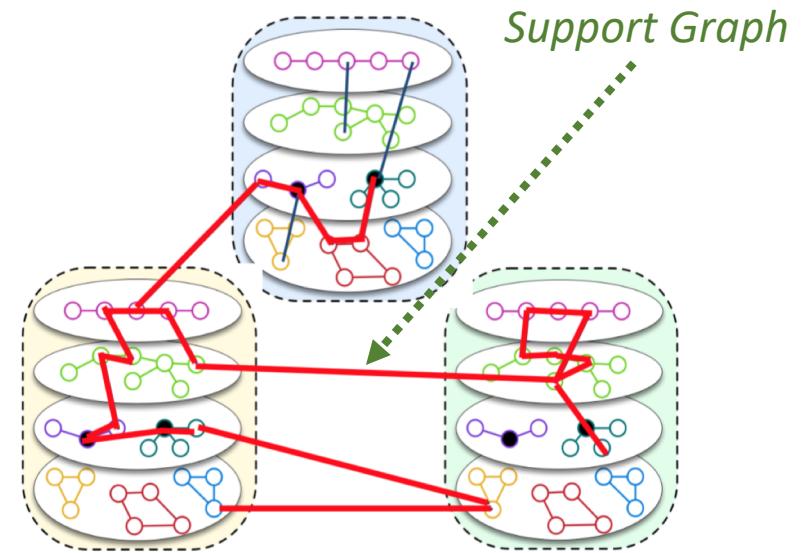
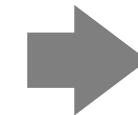


# ILP Model: Some Details

Variables define the space of “support graphs”:

- Each variable corresponds to a node or edge.
- $x=1$  iff nodes / edges are part of the semantic graph.

$$\max \sum_i c_i x_i$$
$$\left\{ \begin{array}{l} \sum_i a_{1i} x_i \leq b_1 \\ \dots \\ \sum_i a_{ki} x_i \leq b_k \end{array} \right.$$

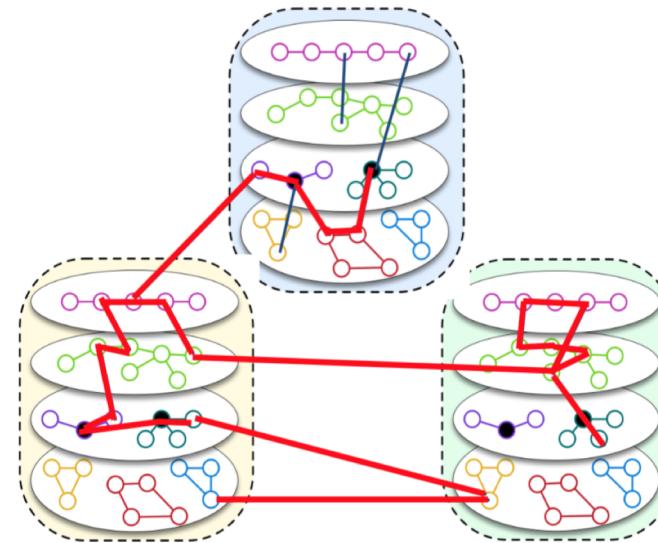


# ILP Model: Some Details

**Objective Function:** “better” support graphs = higher objective value

- Reward good behavior:
  - High lexical match links, nearby alignments, using the subject if using a predicate-argument structure, WH-terms (“which of energy ...”), etc.
- Penalize spurious overuse of frequently occurring terms

$$\max \sum_i c_i x_i$$
$$\forall x_i \in \mathbb{N} \cup \{0\}$$
$$\left\{ \begin{array}{l} \sum_i a_{1i} x_i \leq b_1 \\ \dots \\ \sum_i a_{ki} x_i \leq b_k \end{array} \right.$$



# ILP Model: Some Details

**Dual goal:** scalability, consider only meaningful support graphs

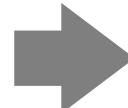
Incorporate global and local structure.

- **Structural Constraints**

- Meaningful proof structures
  - connectedness, question coverage, etc.
  - single/multi-graph, etc.

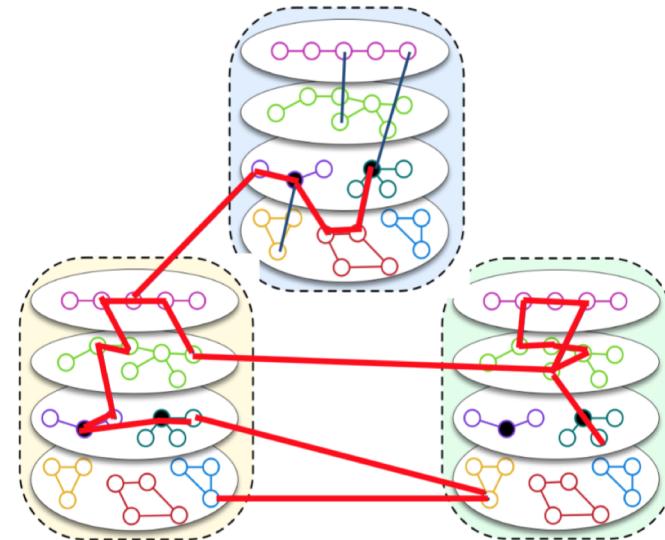
$$\begin{aligned} \max \sum_i c_i x_i \\ \forall x_i \in \mathbb{N} \cup \{0\} \end{aligned}$$

$$\left\{ \begin{array}{l} \sum_i a_{1i} x_i \leq b_1 \\ \dots \\ \sum_i a_{ki} x_i \leq b_k \end{array} \right.$$



- **Semantic Constraints**

- If using a predicate-argument graphs,
  - use at least predicate and argument



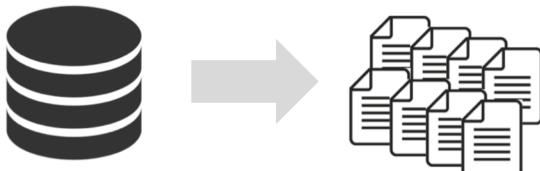
# Evaluation: Notable Baselines

## Information Retrieval (IR)

[Clark et al. AAAI'15]

Information retrieval baseline (Lucene)

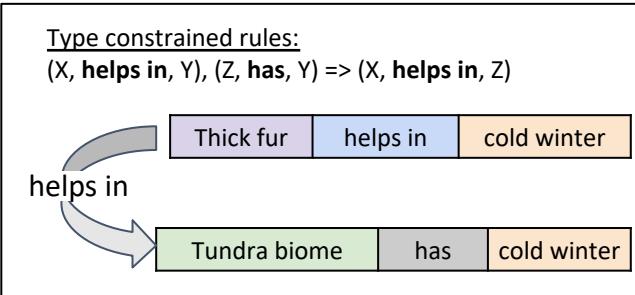
Using 280 GB of plain text



## Inference over structure (TupleInf)

[Khot et al. ACL'17]

Inference over  
auto-generated short triples  
And type-constrained rules

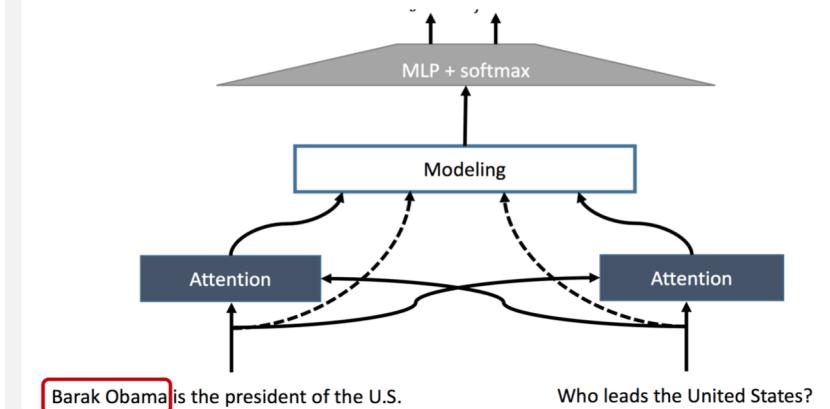


## Neural Network (BiDAF)

[Seo et al. ICLR'16]

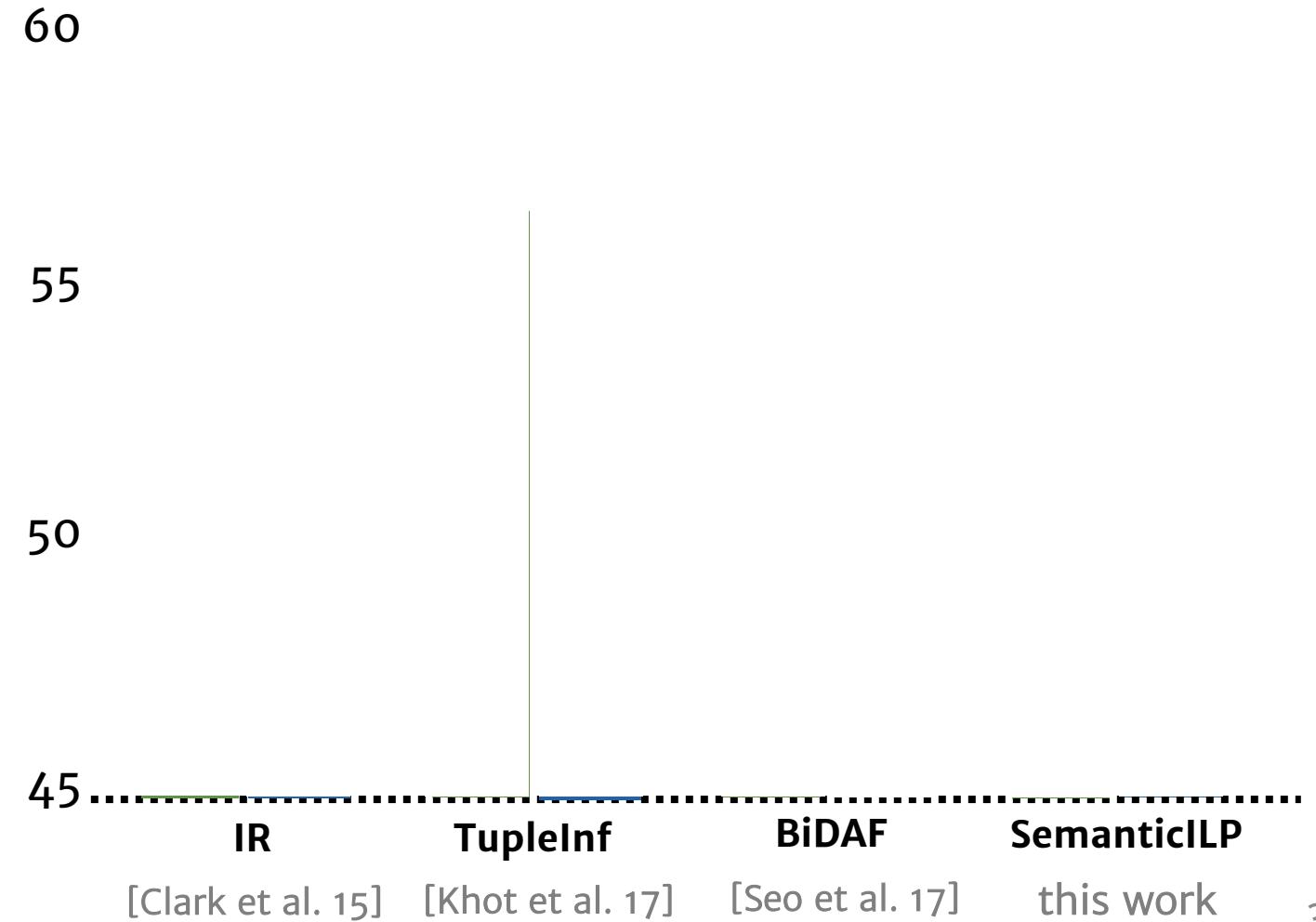
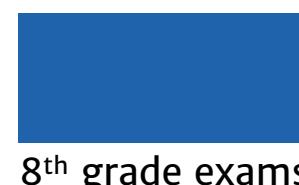
Attention & LSTM

Extractive, i.e select a contiguous phrase in a given paragraph



# Empirical results: Science Domain [ZKTR'18]

SemanticILP consistently outperforms the best baselines in each case by 3%-5%.

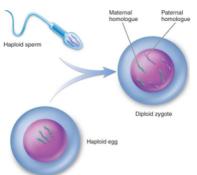


# Answering Questions: Biology Exams

## ▪ **Biology exams** [Berant et al, 2014]

- Technical terms and answer not easy to find.
- Requires understanding complex relations.

We use **the same** version of our systems across our datasets.



**Question:** *What does meiosis directly produce?*



- (A) Gametes  
(B) Haploid cells

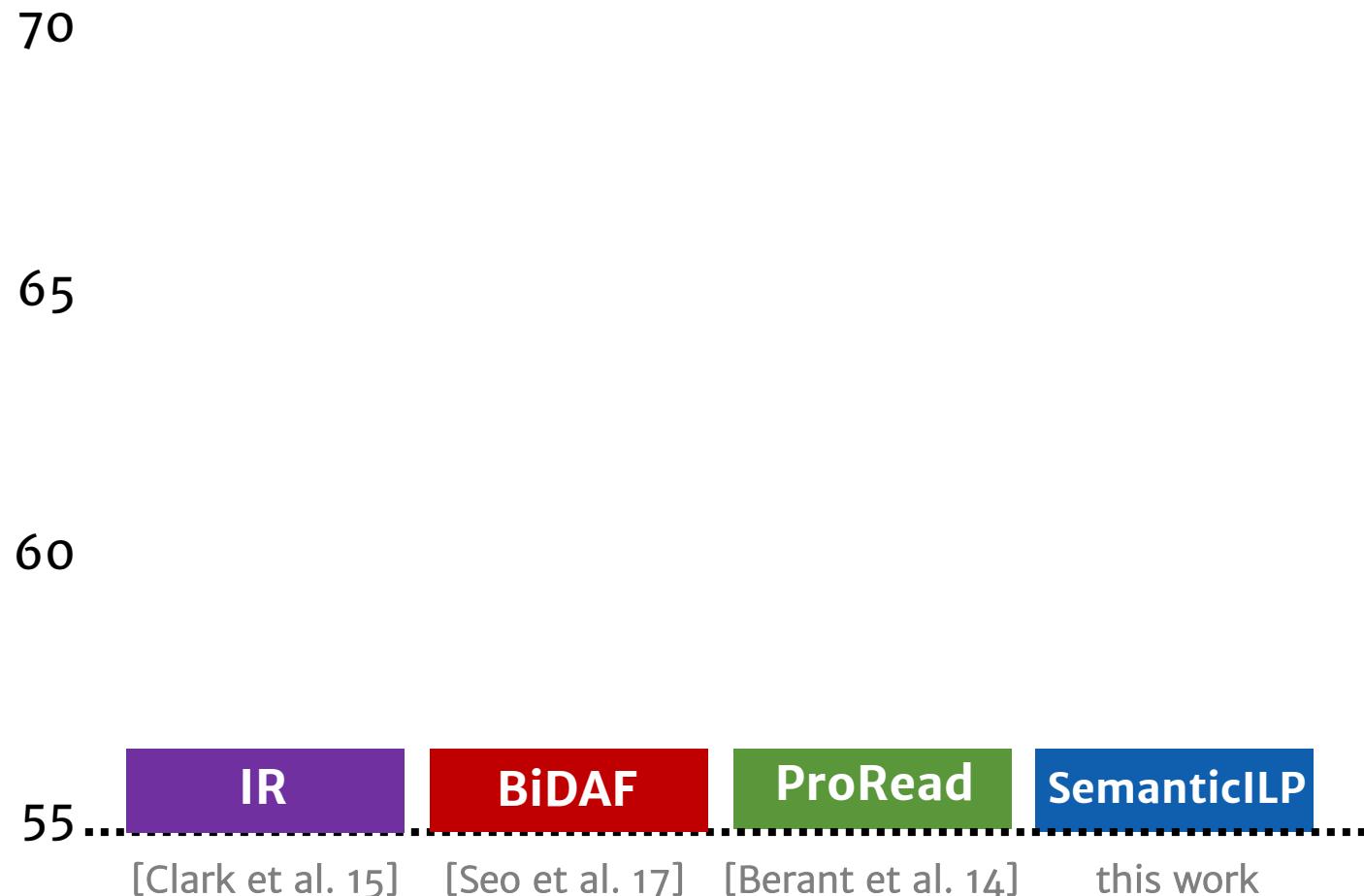
## Evidence paragraph



*... Meiosis produces not gametes but haploid cells that then divide by mitosis and give rise to either unicellular descendants or a haploid multicellular adult organism. Subsequently, the haploid organism carries out further mitoses, producing the cells that develop into gametes....*

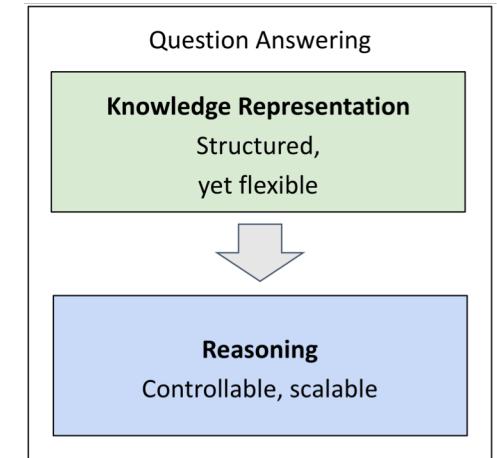
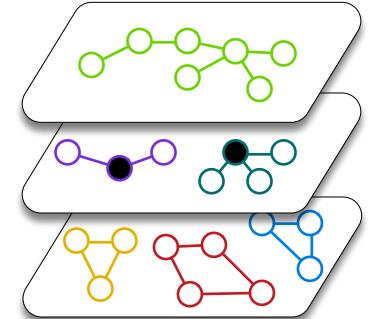
# Empirical results: Biology Domain [ZKTR'18]

SemanticILP generalizes to a **different** domain and achieves on-par score with the best domain-specific system.



# Lessons

- Reasoning over language requires dealing with a diverse set of semantic phenomena.
- Collection of semantic representations of language, independent of the task (**indirect supervision**).
- Better generalization across two different domains.



# **ENTITY TYPING** *with minimal supervision*

Zhou, K et al. Zero-Shot Open Entity Typing as Type-Compatible Grounding. EMNLP 18.

Fei, K et al. Illinois-Profiler: Knowledge Schemas at Scale. IJCAI (Cognitum) 15.

# SEMANTIC TYPING OF ENTITIES

Label mentions with their semantic **types**.

*A handful of professors in the  
**CMU Department of Chemistry**  
are being recognized for their  
efforts and contributions to the  
scientific community.*



**CMU:**

/organization

/organization/education\_institution

**Department of Chemistry:**

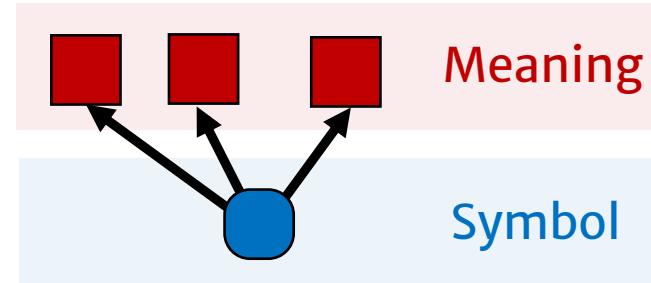
/organization

/education

/education/department

# Semantic Entity Typing: The Necessity (1)

- Dealing with ambiguity



*Our break in **Paris** was quite memorable.*

city

*I met a girl named **Paris**.*

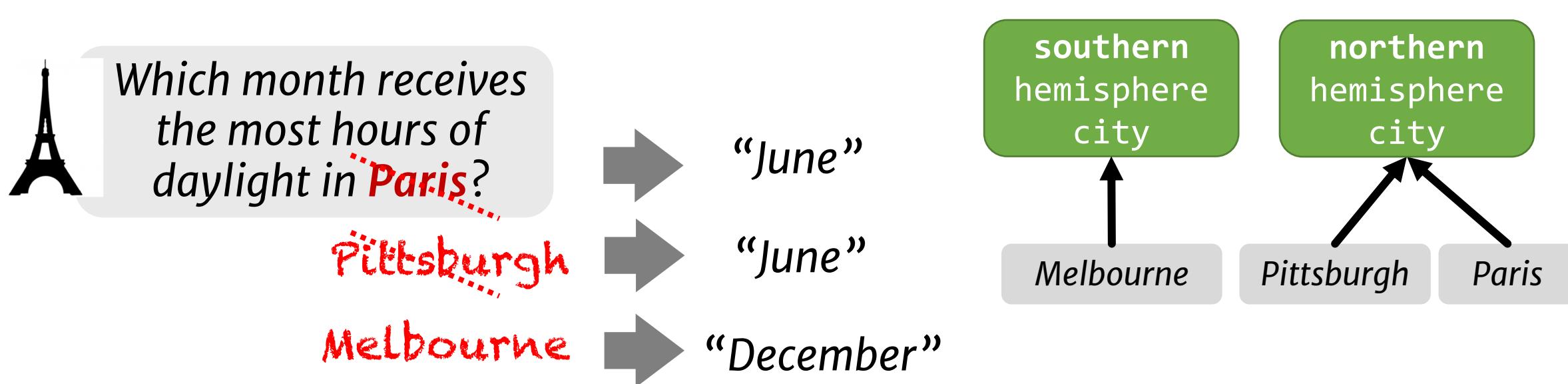
person

***Paris** issued a statement condemning the proposal.*

government

# Semantic Entity Typing: The Necessity (2)

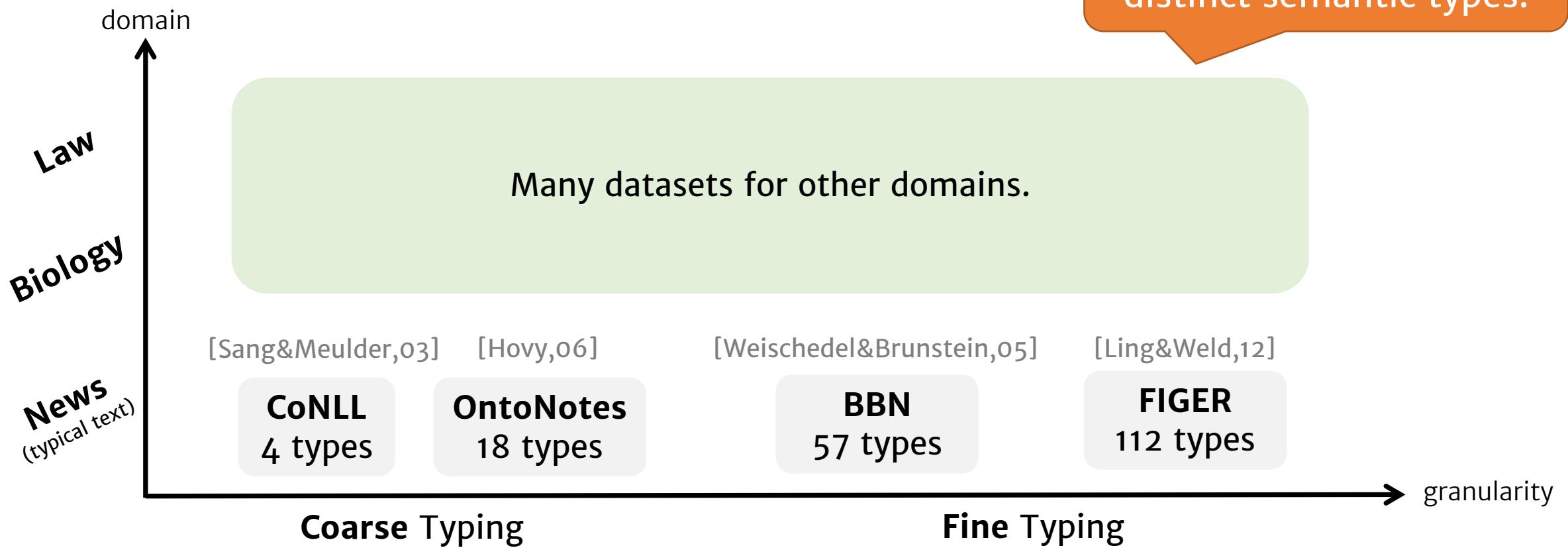
- Linguistic **generalization** requires “abstractions”.
  - Ex: Question answering [Yavuz,16]; Information Extraction [Ling,12].



# Entity Typing: Existing Work

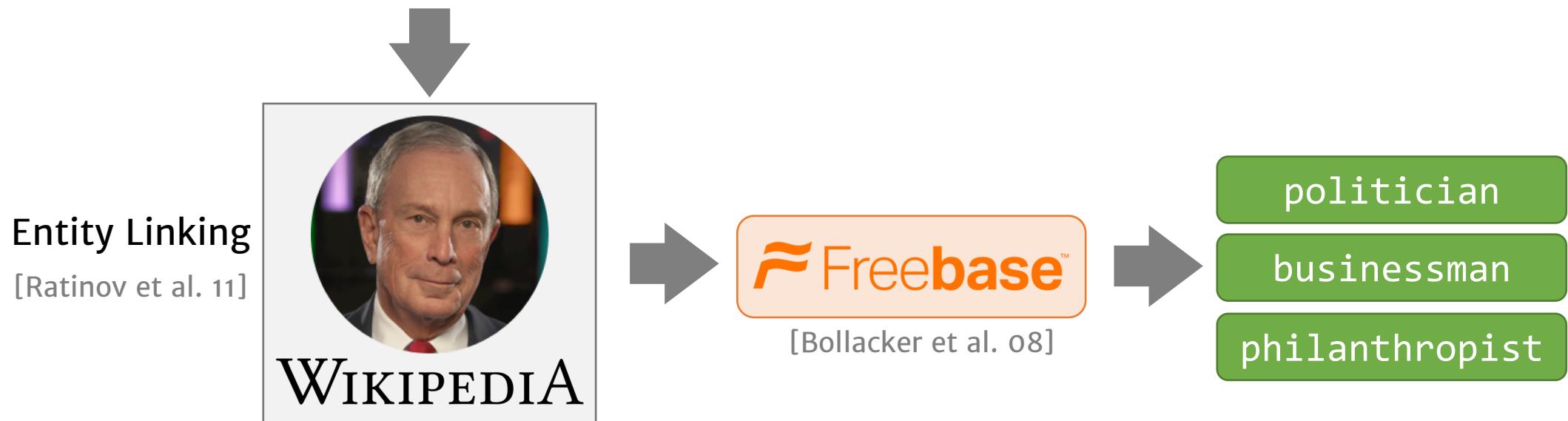
- Multiple datasets for semantic typing

Many datasets, each with distinct semantic types.



# “Cheap” Typing with Wikipedia

A former Democrat, **Bloomberg** switched his party registration in 2001.

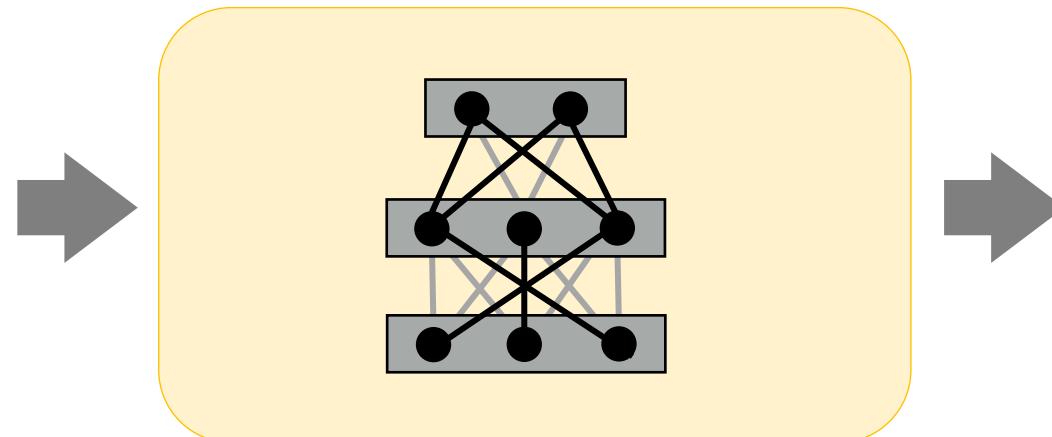


Not consistent with the context

# A Common Approach: Supervised Learning

- **Input:** sentence, mention.
- **Output:** a set of types.

*A former Democrat,  
**Bloomberg** switched his  
party registration in 2001.*



person  
politician

Taxonomy is [indirectly] defined  
during the **training** time.

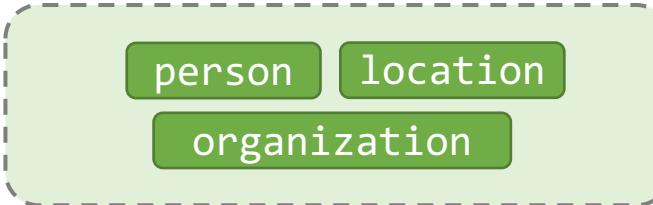
# Zero-Shot Open Entity Typing

- **Input:** sentence, mention, **target taxonomy**.
- **Output:** a set of types (according to the target type taxonomy).

A former Democrat,  
**Bloomberg** switched his  
party registration in 2001.



Target  
Taxonomy



Classifies into a  
given taxonomy

# Zero-Shot Open Entity Typing

- **Input:** sentence, mention, **target taxonomy**.
- **Output:** a set of types (according to the target type taxonomy).

A former Democrat,  
**Bloomberg** switched his  
party registration in 2001.

**Open:**  
Classifies into a  
given taxonomy

Target  
Taxonomy



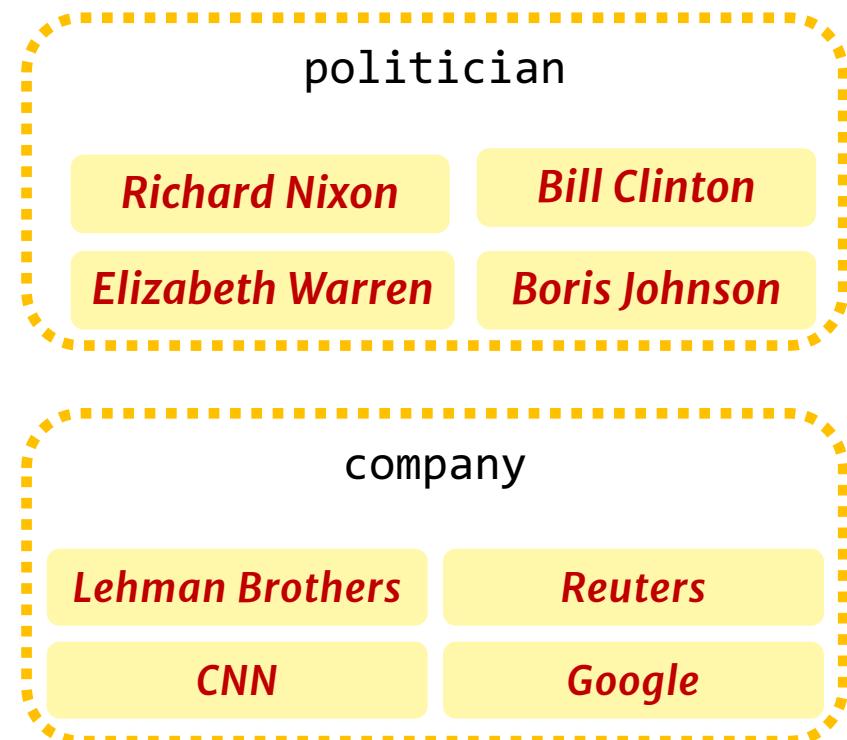
**Zero-shot:**  
no taxonomy-specific  
supervision

# ZOE: Type-Compatible Grounding

- “Type” as conceptual container binding entities together.

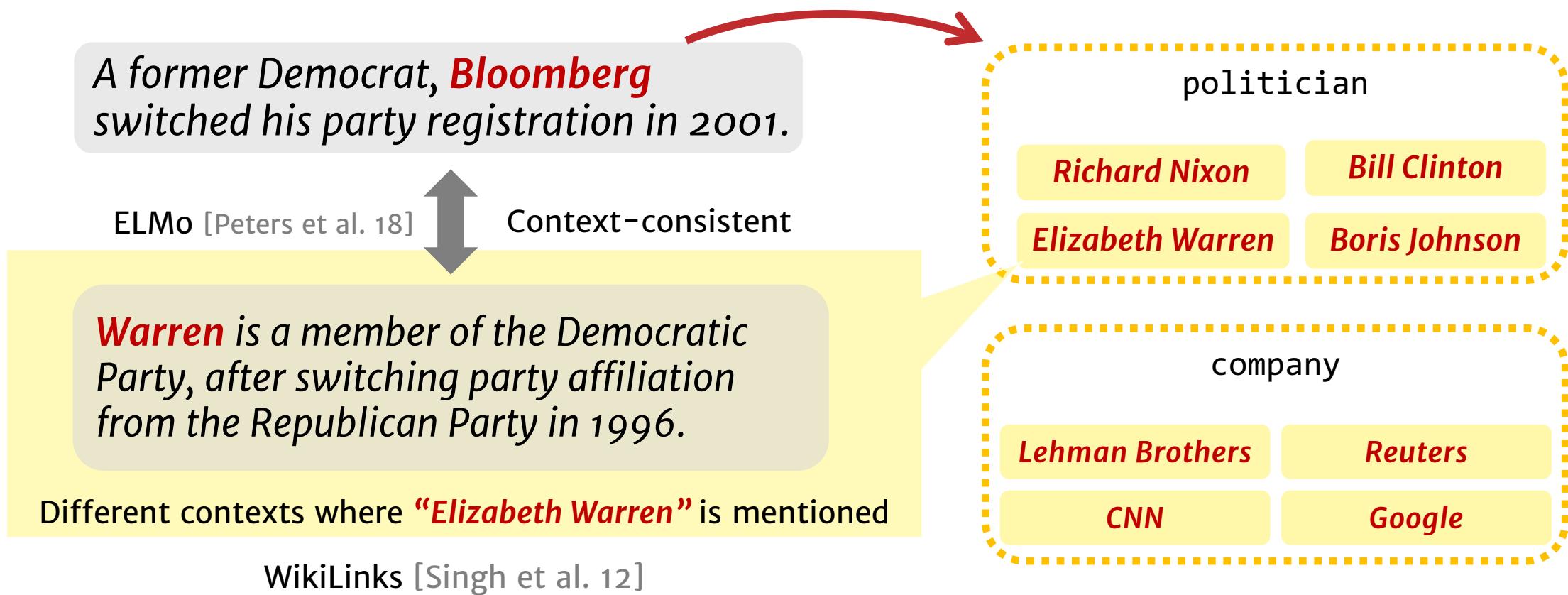
A former Democrat, **Bloomberg** switched his party registration in 2001.

**Key idea:** Determine the **type** of an input mention by finding entities in the **type defining set** that share a similar context



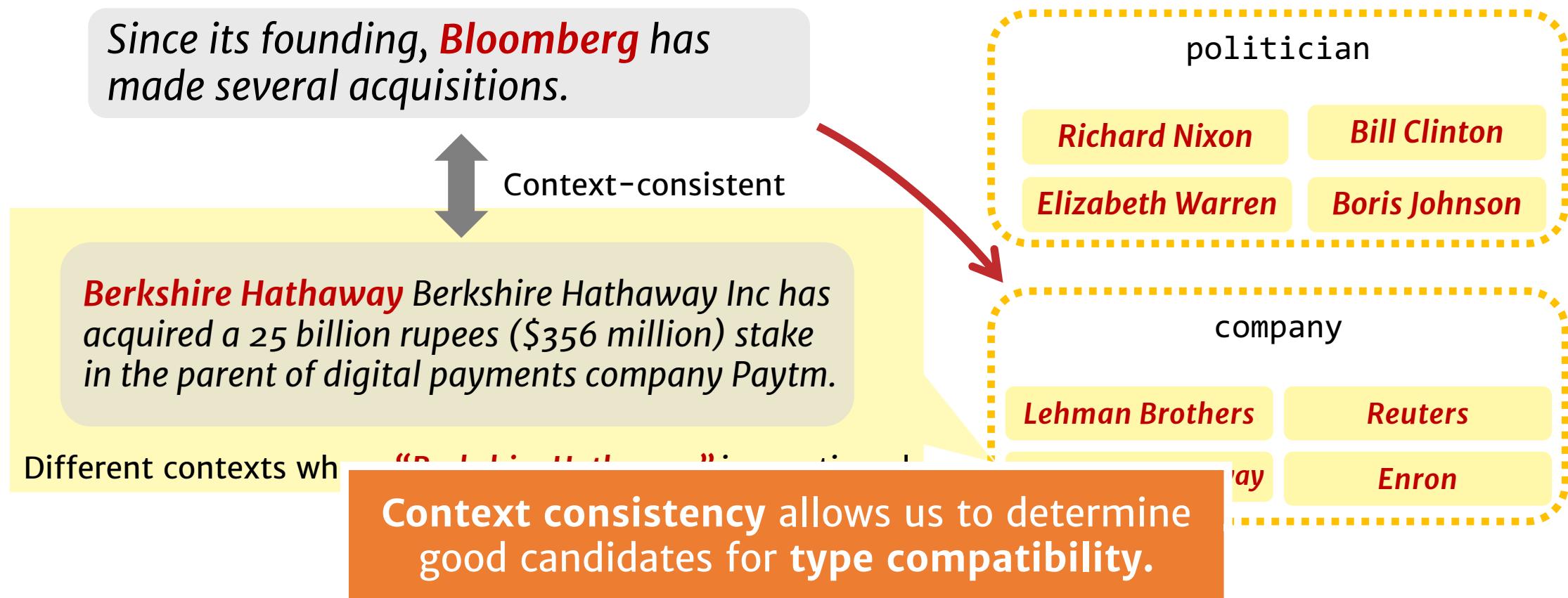
# ZOE: Type-Compatible Grounding

- “Type” as conceptual container binding entities together.



# ZOE: Type-Compatible Grounding

- “Type” as conceptual container binding entities together.



# Zero-Shot Open Typing: Big Picture

A mention & its context

*A former Democrat,  
**Bloomberg** switched his  
party registration in 2001.*



Mapping type-compatible Wikipedia entities

**Richard Nixon**

person  
politician  
president

**Bill de Blasio**

mayor  
politician  
person

**Elizabeth Warren**

person  
politician  
scholar



Inference: aggregate and rank the consistency scores.



person

politician

official

## High-level Algorithm:

1. Map the mention to context-consistent Wikipedia concepts
2. Rank candidate titles by context-consistency and infer the types according to the type taxonomy.

# Zero-Shot Open Typing: Big Picture

A mention & its context

*A former Democrat,  
**Bloomberg** switched his  
party registration in 2001.*



Mapping type-compatible Wikipedia entities

**Richard Nixon**

person  
politician  
president

**Bill de Blasio**

mayor  
politician  
person

**Elizabeth Warren**

person  
politician  
scholar



Inference: aggregate and rank the consistency scores.



person

politician

official

Resources



WIKIPEDIA

 Freebase™

[Bollacker et al. 08]

WikiLinks  
[Singh et al. 12]

Contextualized  
Representations



[Peters et al. 18]

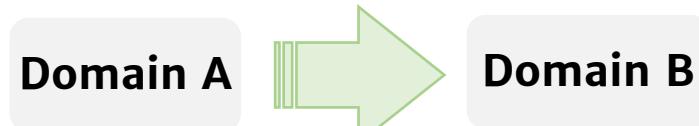
# Empirical Results: Fine-Typing [zKTR18]

- Outperforms supervised system in cross-domain.

System	Trained on	Evaluated on		
		FIGER	BBN	Ontonotes

- Comparable results with supervised systems.

# Lessons



- Reformulating the task and using weak signals helps us reduce our dependence on direct “supervision”.
- This type-aware approach leads to the ability to transfer across domains & taxonomies.

# Beyond Supervision-rich “tasks”

- We will never have enough annotated data to train all the models for all the tasks.
  - Annotation for complex tasks is difficult, costly and sometimes impossible.
- We don't even know what are “all the tasks”.

# Beyond Supervision-rich “tasks”

- Two samples of research projects in an attempt to utilize hints in data to infer supervision signals:
  - Representation
  - Structure
- Not just two systems:
  - Initial steps towards a broader theory of using “incidental” signals.

[Roth, AAAI'17]

# **BIG PICTURE + LOOK AHEAD**

## Machine Learning, Optimization & applications

KSKCSSR. StartAI'18  
KKCMSR. COLING'16  
QK. NourIPS'15  
KNJF. TIP'14  
NKTNJ. SMC'11

## Natural Language Processing

### Semantics

*Semantic Role Labeling, Name Entities, Semantic Language models, Coreference, etc.*

ZKCR. EMNLP'18

KCRUR. NAACL'18

FKPWR. Cognitum'15

PKR. NAACL'15

### Learning & Inference

*Question Answering, Textual Entailment, etc.*

KKSR. AAAI'18

KKSR. CoNLL'17

CEKSTT<sup>K</sup>. AAAI'16

KKSR. IJCAI'16

### NLP tools/software

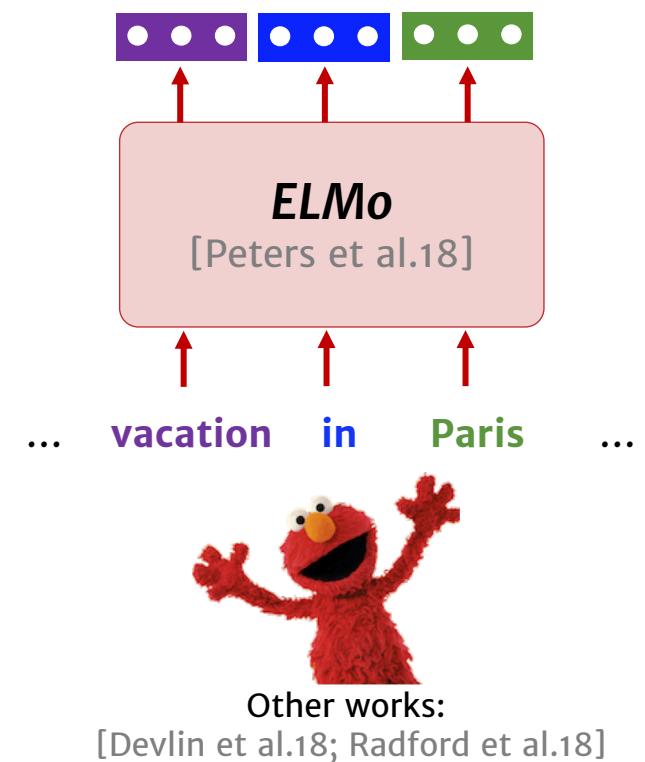
K et al. LREC'18

SCKKSVBWR. LREC'16

# Beyond Supervision-rich “tasks”

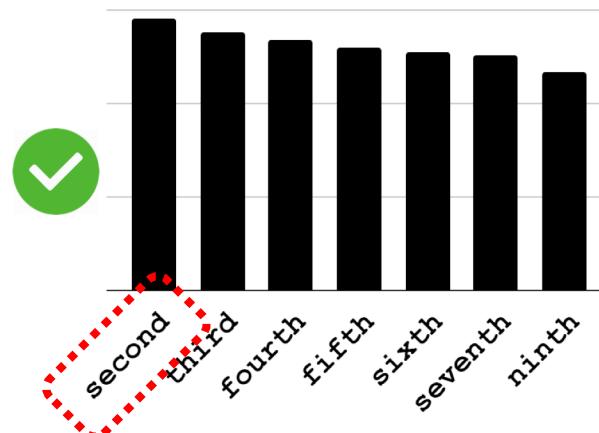
- A major shift in the field:

- Being able to make use of massive loads of **unlabeled** data in the form of language models.
  - Compatible with the philosophy I advocated for here.

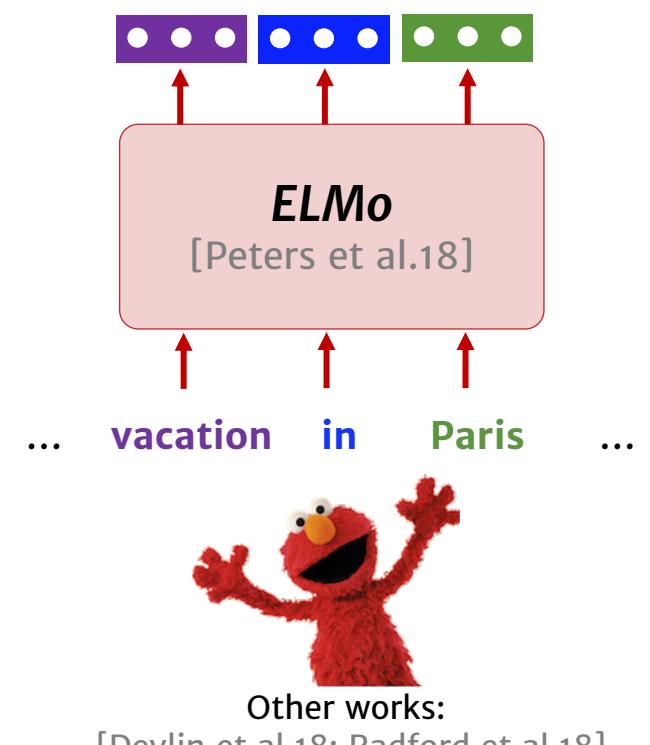


# Language Models: Means to Access Knowledge

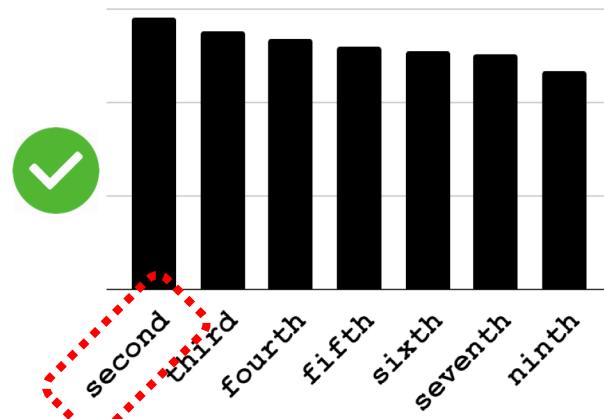
- They let you “query” for knowledge:



Pittsburgh is the \_\_\_\_ -largest populated city in Pennsylvania.



# Language Models: Means to Access Knowledge



Pittsburgh is the \_\_\_\_ -largest populated city in Pennsylvania.

- **What is known:**
  - What is the nature the knowledge that they have internalized?
- **Know what you know:**
  - Is there a mechanism to decide whether something is [not]?
- **Inference with knowledge:**
  - Access what is known and be able to solve bigger problems.

# Language Models: Biases

- What does this mean for the NLP systems built out of such systems?
- **Discovery:**
  - How can we automate the discovery of issues?
- **Mitigation:**
  - How can we resolve the such biases?

## Machine Learning, Optimization & applications

KSKCSSR. StartAI'18  
KKCMSR. COLING'16  
QK. NourIPS'15  
KNJF. TIP'14  
NKTNJ. SMC'11

## Natural Language Processing

### Semantics

*Semantic Role Labeling, Name Entities, Semantic Language models, Coreference, etc.*

ZKCR. EMNLP'18

KCRUR. NAACL'18

FKPWR. Cognitum'15

PKR. NAACL'15

### Learning & Inference

*Question Answering, Textual Entailment, etc.*

KKSR. AAAI'18

KKSR. CoNLL'17

CEKSTT<sup>K</sup>. AAAI'16

KKSR. IJCAI'16

### NLP tools/software

K et al. LREC'18

SCKKSVBWR. LREC'16

### Information Pollution

CKWCR. NAACL'19



# Information Pollution

- Information Technology started with much optimism:
  - Democratizing information and greater liberties.
- Few foresaw the huge radical impact of the information revolution.
  - Massive amount of Information pollution:



“The contamination of the information supply with irrelevant, redundant, unsolicited, incorrect, and otherwise low-value information.”

[Levent Orman '15]

# Information Pollution: Not Just Politics

- Medical Domain, Education, Public Policy, etc.

- “Best treatment for X;” “Side effects of X.”

- Are they consistent?
  - Are they trustworthy?
  - Are they written by someone with an agenda?



# Information Pollution: Not Just Fact-Checking

- Many issues don't have a single “answer.”
  - “Should X be legalized?”
    - Possible answers are subject to situations, world views or background.
    - Moral, utilitarian, libertarian, philosophy, etc.



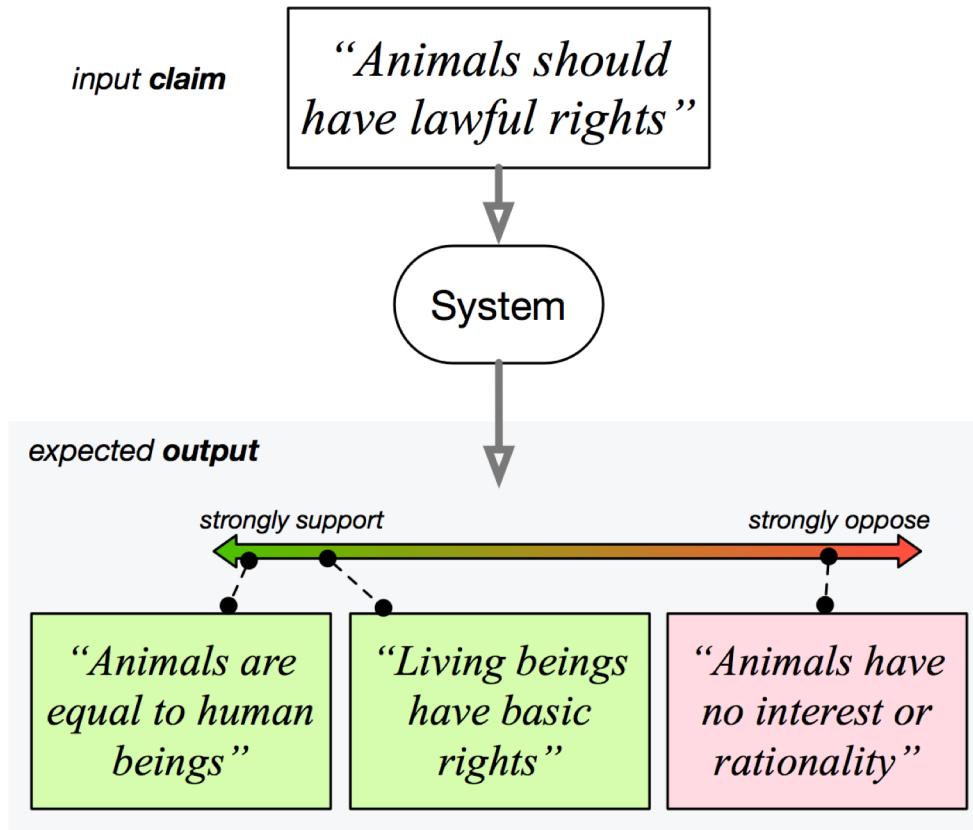
Factual information (or lack of) is  
**not** really the core of the problem.

# Information Pollution, as NLU Problems

- **Understanding Sources**
- But **what should we believe**, and who should we trust?
- Sources may
  - Have their own, often hidden, motivations
  - Make different or even contradictory claims
- **Understanding the evidence**
- Sources may present different, but legitimate, perspectives
- Most interesting issues/questions have multiple “right” answers
  - Perspectives must be supported by evidence

*Not only applications for NLP, but also drive the research in important directions.*

# Discovering Diverse “Perspectives”



- Our recent work: provide users with the understanding that each “story” has more than one “perspective.”
- Goal:
  - Perspectives could give a fuller understanding of an issue.
  - Make us more open-minded, less afraid & more likely to consider other views.

# Information Pollution: an NLU Challenge

- Suffering from this pollution is not a forgone conclusion.
- A computational model that will help us navigate the polluted world.
  - Natural Language Processing/Understanding + Algorithmic Components
  - Collaborative efforts involving experts from the social sciences, policy, and others.
- Overreliance on fully annotated data, unlikely to solve the problem.
- Interesting challenge, important, and will have societal impact.



Tushar Khot  
(AI2)



Dan Roth  
(UPenn)



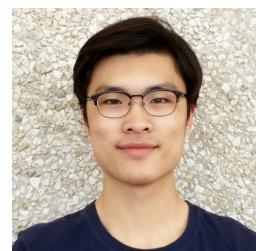
Ashish Sabharwal  
(AI2)



Peter Clark  
(AI2)



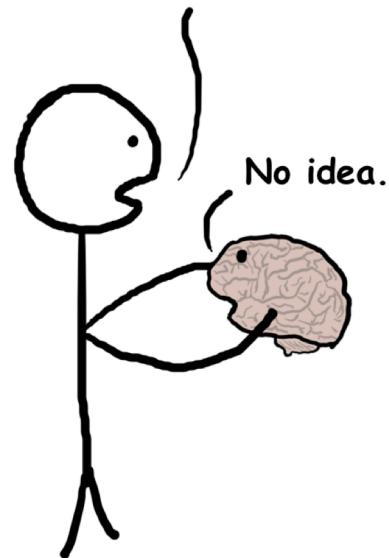
Chen-Tse Tsai  
(Bloomberg)



Ben Zhou  
(UIUC → UPenn)

# That's it, folks!

How do you work?



# Natural Language Processing

## Semantics

*Semantic Role Labeling, Name Entities, Semantic Language models, Coreference, etc.*

ZKCR. EMNLP'18

KCRUR. NAACL'18

FKPWR. Cognitum'15

PKR. NAACL'15

## Learning & Inference

*Question Answering, Textual Entailment, etc.*

KKSR. AAAI'18

KKSR. CoNLL'17

CEKSTT<sup>K</sup>. AAAI'16

KKSR. IJCAI'16

## NLP tools/software

K et al. LREC'18

SCKKSVBWR. LREC'16

## Machine Learning, Optimization & applications

KS<sup>K</sup>CSSR. StartAI'18

KKCMSR. COLING'16

Q<sup>K</sup>. NourIPS'15

KNJF. TIP'14

NKTNJ. SMC'11

Theoretical

Empirical

