



**THE OHIO STATE
UNIVERSITY**

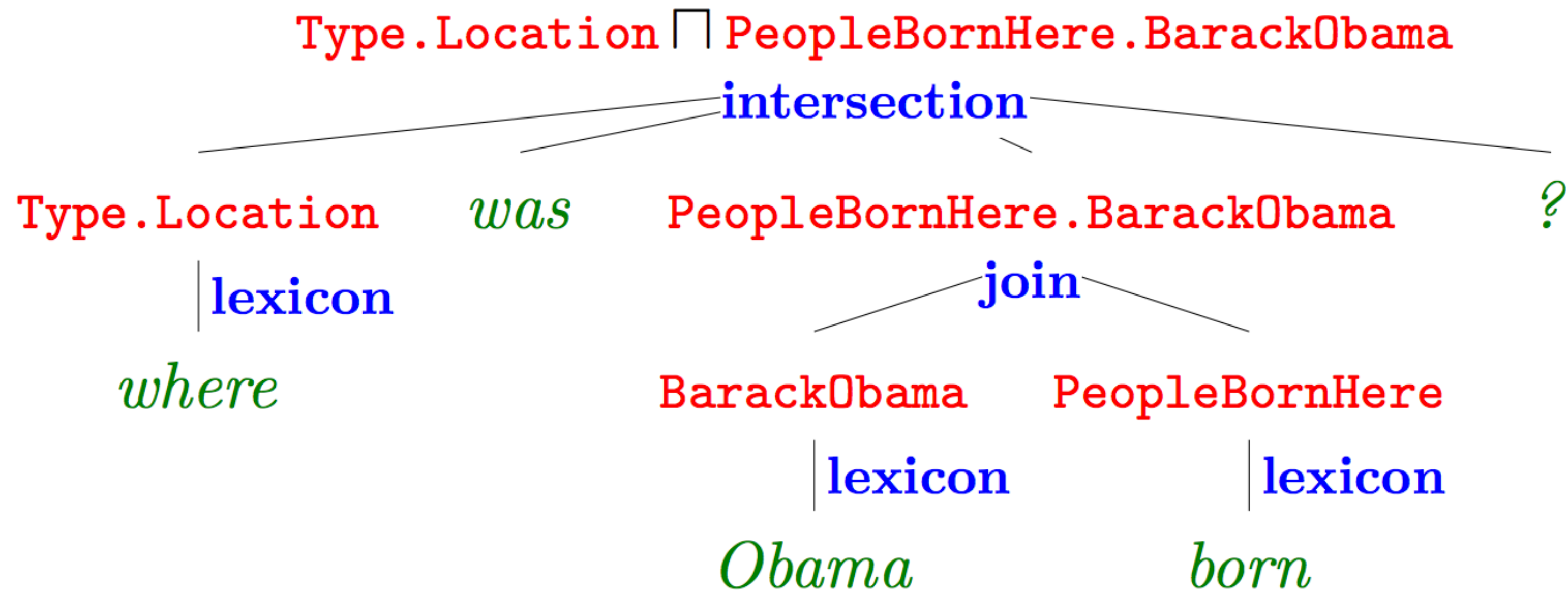
CSE 5525: Foundations of Speech and Language Processing

Information Extraction (IE)

Huan Sun (CSE@OSU)

Slides were largely adapted from Prof. Greg Durrett @ UT Austin.

Why Information Extraction?



Connecting with QA previously discussed

This Lecture

- ▶ How do we represent information for information extraction?
- ▶ Relation extraction
- ▶ Slot filling
- ▶ Open Information Extraction

Representing Information

Semantic Representations

- ▶ “World” is a set of entities and predicates

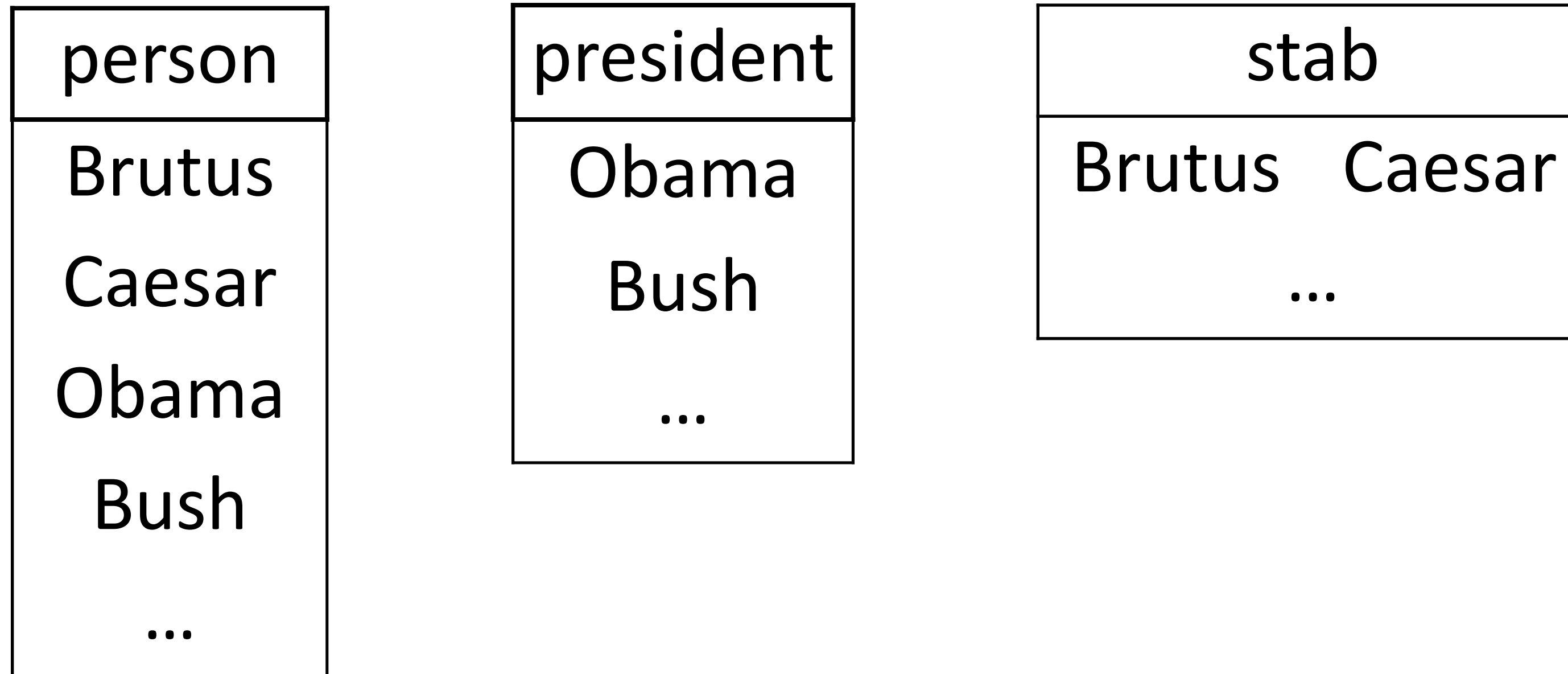
person
Brutus
Caesar
Obama
Bush
...

president
Obama
Bush
...

stab
Brutus Caesar
...

Semantic Representations

- ▶ “World” is a set of entities and predicates



- ▶ Statements are logical expressions that evaluate to true or false

Brutus stabs Caesar

$\text{stab}(\text{Brutus}, \text{Caesar}) \Rightarrow \text{true}$

Caesar was stabbed

$\exists x \text{stab}(x, \text{Caesar}) \Rightarrow \text{true}$

Example credit: Asad Sayeed

Neo-Davidsonian Event Semantics

Brutus stabbed Caesar with a knife at the theater on the Ides of March

$\exists e \text{ stabs}(e, \text{Brutus}, \text{Caesar}) \wedge \text{with}(e, \text{knife}) \wedge \text{location}(e, \text{theater})$
 $\wedge \text{time}(e, \text{Ides of March})$

- ▶ Lets describe events as having properties
- ▶ Unified representation of events and entities:

some clever driver in America

$\exists x \text{ driver}(x) \wedge \text{clever}(x) \wedge \text{location}(x, \text{America})$

Real Text

Barack Obama signed the Affordable Care act on Tuesday. He gave a speech later that afternoon on how the act would help the American people. Several prominent Republicans were quick to denounce the new law.

$\exists e \text{ sign}(e, \text{Barack Obama}) \wedge \text{patient}(e, \text{ACA}) \wedge \text{time}(e, \text{Tuesday})$

Real Text

which afternoon?

who?

Barack Obama signed the Affordable Care act on Tuesday. He gave a speech later that afternoon on how the act would help the American people. Several prominent Republicans were quick to denounce the new law.

which Tuesday?

$\exists e \text{ sign}(e, \text{Barack Obama}) \wedge \text{patient}(e, \text{ACA}) \wedge \text{time}(e, \text{Tuesday})$

- ▶ Need to impute missing information, resolve coreference, etc.

Real Text

which afternoon?

who?

Barack Obama signed the Affordable Care act on Tuesday. He gave a speech later that afternoon on how the act would help the American people. Several prominent Republicans were quick to denounce the new law.

???

which Tuesday?

$\exists e \text{ sign}(e, \text{Barack Obama}) \wedge \text{patient}(e, \text{ACA}) \wedge \text{time}(e, \text{Tuesday})$

- ▶ Need to impute missing information, resolve coreference, etc.
- ▶ Still unclear how to represent some things precisely or how that information could be leveraged (several prominent Republicans)

Real Text

which afternoon?

who?

Barack Obama signed the Affordable Care act on Tuesday. He gave a speech later that afternoon on how the act would help the American people. Several prominent Republicans were quick to denounce the new law.

???

which Tuesday?

$\exists e \text{ sign}(e, \text{Barack Obama}) \wedge \text{patient}(e, \text{ACA}) \wedge \text{time}(e, \text{Tuesday})$

- ▶ Need to impute missing information, resolve coreference, etc.
- ▶ Still unclear how to represent some things precisely or how that information could be leveraged (several prominent Republicans)

Other Challenges

Bob and Alice were friends until he moved away to attend college

$\exists e1 \exists e2 \text{ friends}(e1, \text{Bob}, \text{Alice}) \wedge \text{moved}(e2, \text{Bob}) \wedge \text{end_of}(e1, e2)$

► How to represent temporal information?

*Bob and Alice were friends until **around the time** he moved away to attend college*

Other Challenges

Bob and Alice were friends until he moved away to attend college

$\exists e1 \exists e2 \text{ friends}(e1, \text{Bob}, \text{Alice}) \wedge \text{moved}(e2, \text{Bob}) \wedge \text{end_of}(e1, e2)$

- ▶ How to represent temporal information?

*Bob and Alice were friends until **around the time** he moved away to attend college*

- ▶ Representing truly open-domain information is very complicated! We don't have a formal representation that can capture everything

(At least) Three Solutions

- ▶ Crafted annotations to capture some subset of phenomena: predicate-argument structures (semantic role labeling), time (temporal relations), ...
- ▶ Slot filling: specific ontology, populate information in a predefined way

(Earthquake: magnitude=8.0, epicenter=central Italy, ...)

- ▶ Entity-relation-entity triples: focus on entities and their relations (note that entities are pretty broad: can include events like *World War II*, etc.)

(Lady Gaga, singerOf, Bad Romance)

Open IE

- ▶ Entity-relation-entity triples **aren't necessarily grounded in an ontology**
- ▶ Extract strings and let a downstream system figure it out

Barack Obama signed the Affordable Care act on Tuesday. He gave a speech later that afternoon on how the act would help the American people. Several prominent Republicans were quick to denounce the new law.

(Barack Obama, signed, the Affordable Care act)

(Several prominent Republicans, denounce, the new law)

IE: The Big Picture

- ▶ How do we represent information? What do we extract?
 - ▶ Semantic roles
 - ▶ Abstract meaning representation
 - ▶ Slot fillers
 - ▶ Entity-relation-entity triples (fixed ontology or open)

Semantic Role Labeling/ Abstract Meaning Representation

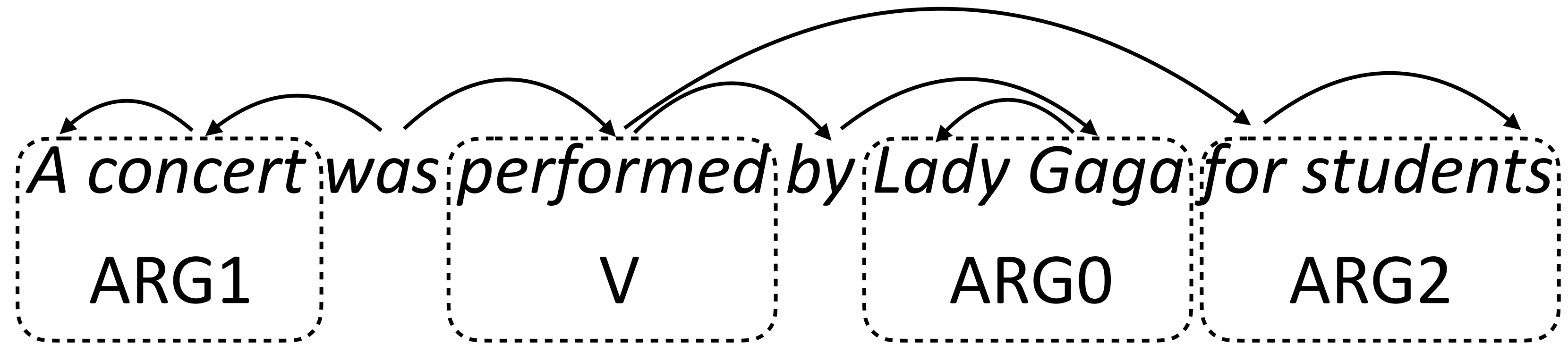
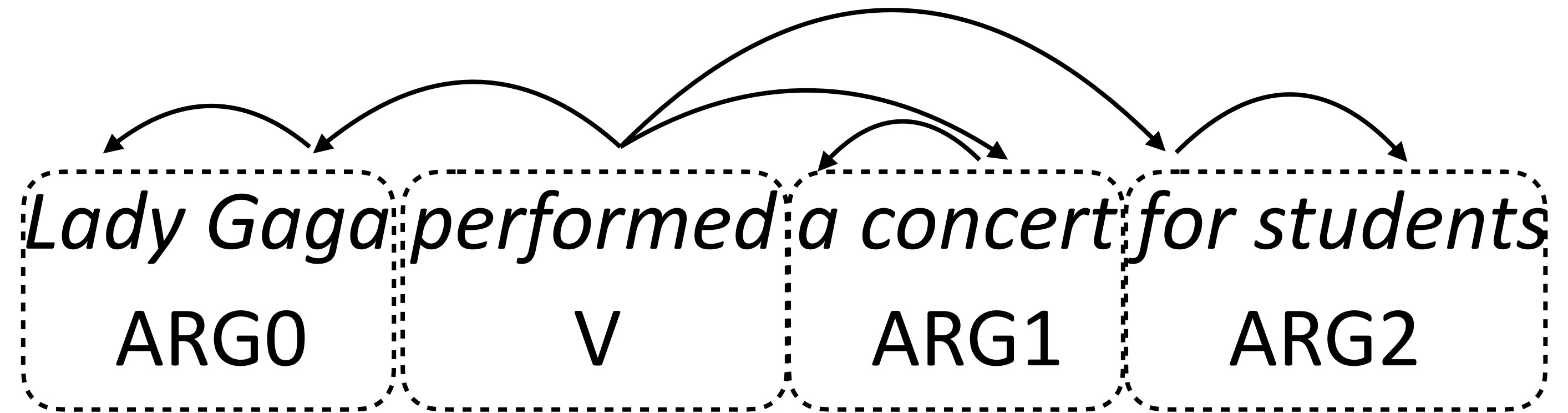
Semantic Role Labeling

- ▶ Performing event

- ▶ Subject: Lady Gaga

- ▶ Object: a concert

- ▶ Audience: students



- ▶ Same event described but the representation looks different

- ▶ Verb (predicate) associated with several arguments (roles): “Agent”, “Theme”, and “Beneficiary”

VerbNet

percentage.n	(GROUPING)
percentile.n	(GROUPING)
perception.n	(GROUPING)
perch	(IN PLACING), (GROUPING)
percolate	
perfect	
perforate	
perform	
performance.n	
perfume	
perfuse	
peril.n	

FRAMES	
NP V NP	
EXAMPLE	"Sandy sang a song."
SYNTAX	<u>AGENT</u> V <u>THEME</u>
SEMANTICS	PERFORM (DURING(E), AGENT, THEME)
NP V	
EXAMPLE	"Sandy sang."
SYNTAX	<u>AGENT</u> V
SEMANTICS	PERFORM (DURING(E), AGENT, ?THEME)
NP V NP PP.BENEFICIARY	
EXAMPLE	"Sandy sang a song for me."
SYNTAX	<u>AGENT</u> V <u>THEME</u> {FOR} <u>BENEFICIARY</u>
SEMANTICS	PERFORM (DURING(E), AGENT, THEME) BENEFIT (E, BENEFICIARY)

- ▶ Defines the semantics of verbs, arguments for every verb in English

Semantic Roles

- ▶ “Postprocessing” layer on top of dependency parsing that exposes useful information, canonicalizes across grammatical constructions
- ▶ Related to theta roles in linguistics
- ▶ Agent (~ subject), patient/theme (~ object), goal (~ indirect object)
ARG0 ARG1 ARG2+ (semantics vary)

Semantic Role Labeling

- ▶ Identify predicate, disambiguate it, identify that predicate's arguments
- ▶ Verb roles from Propbank (Palmer et al., 2005)

Gold

ARG1

V

ARG2

ARG3

Housing starts are expected to quicken a bit from August's pace

quicken:

Arg0-PAG: *causer of speed-up*

Arg1-PPT: *thing becoming faster* (vnrole: 45.4-patient)

Arg2-EXT: *EXT*

Arg3-DIR: *old speed*

Arg4-PRD: *new speed*

Figure from He et al. (2017)

Semantic Role Labeling

- ▶ Identify predicates (*love*) using a classifier (not shown)
- ▶ Identify ARG0, ARG1, etc. as a tagging task with a BiLSTM conditioned on *love*
- ▶ Other systems incorporate syntax, joint predicate-argument finding

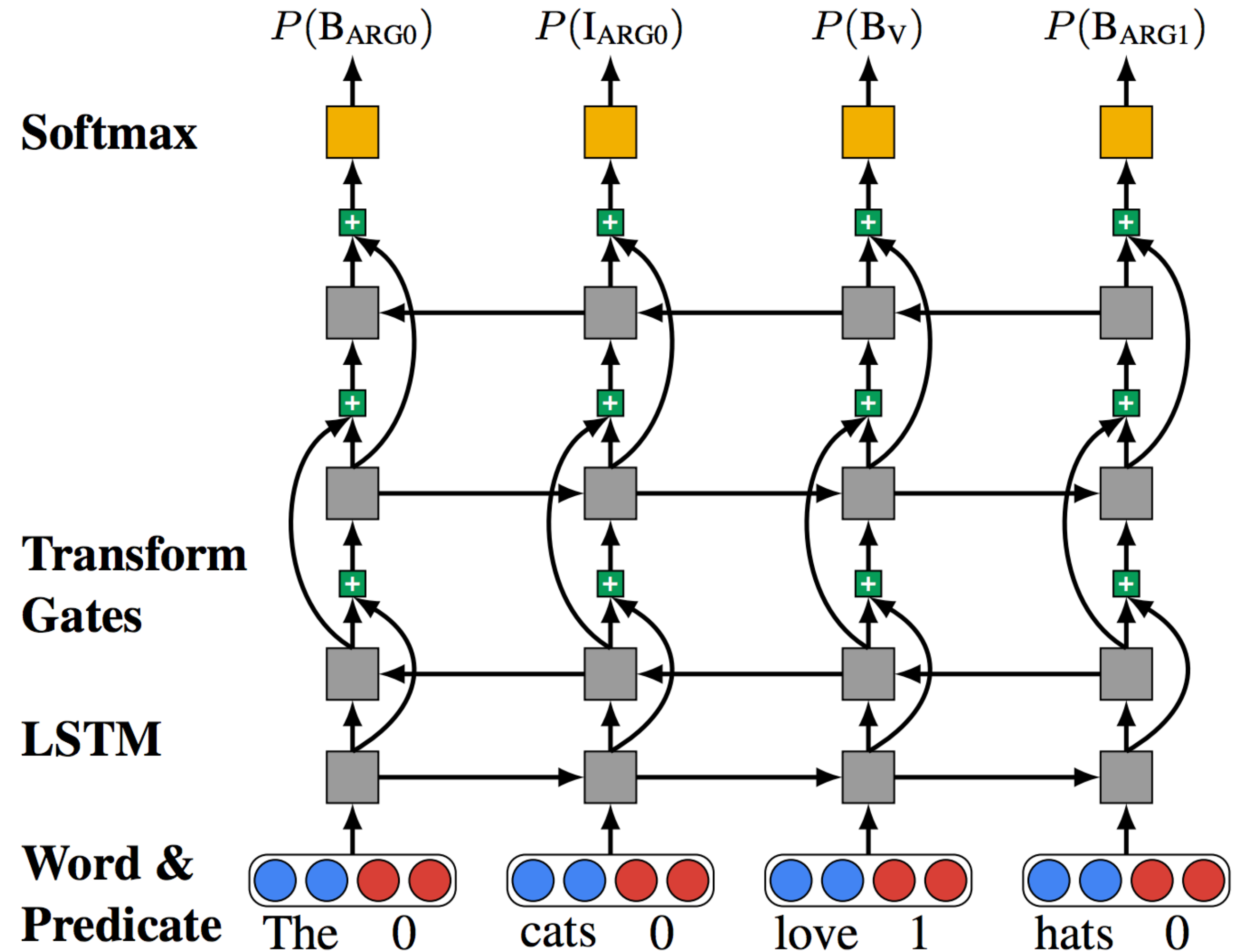


Figure from He et al. (2017)

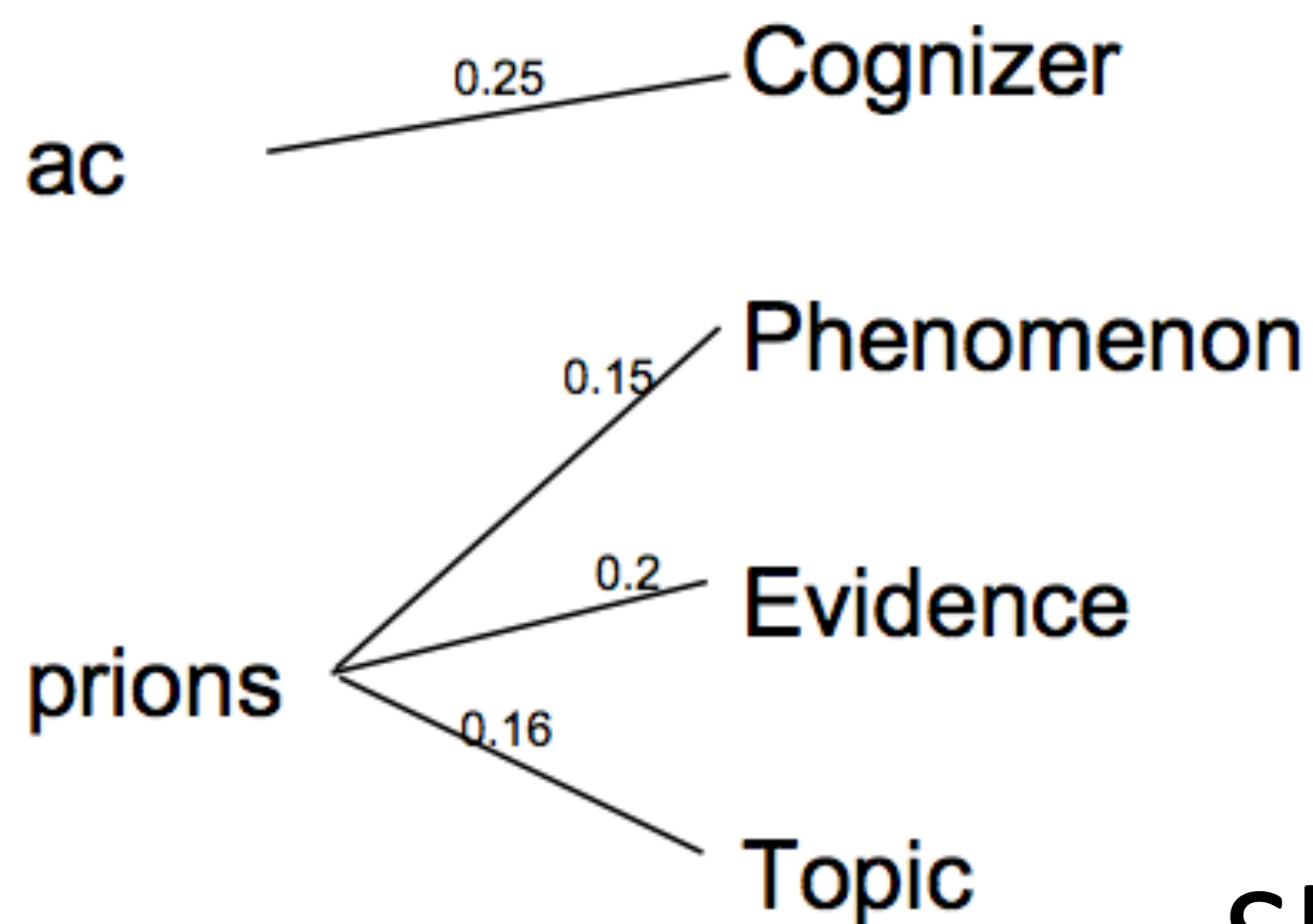
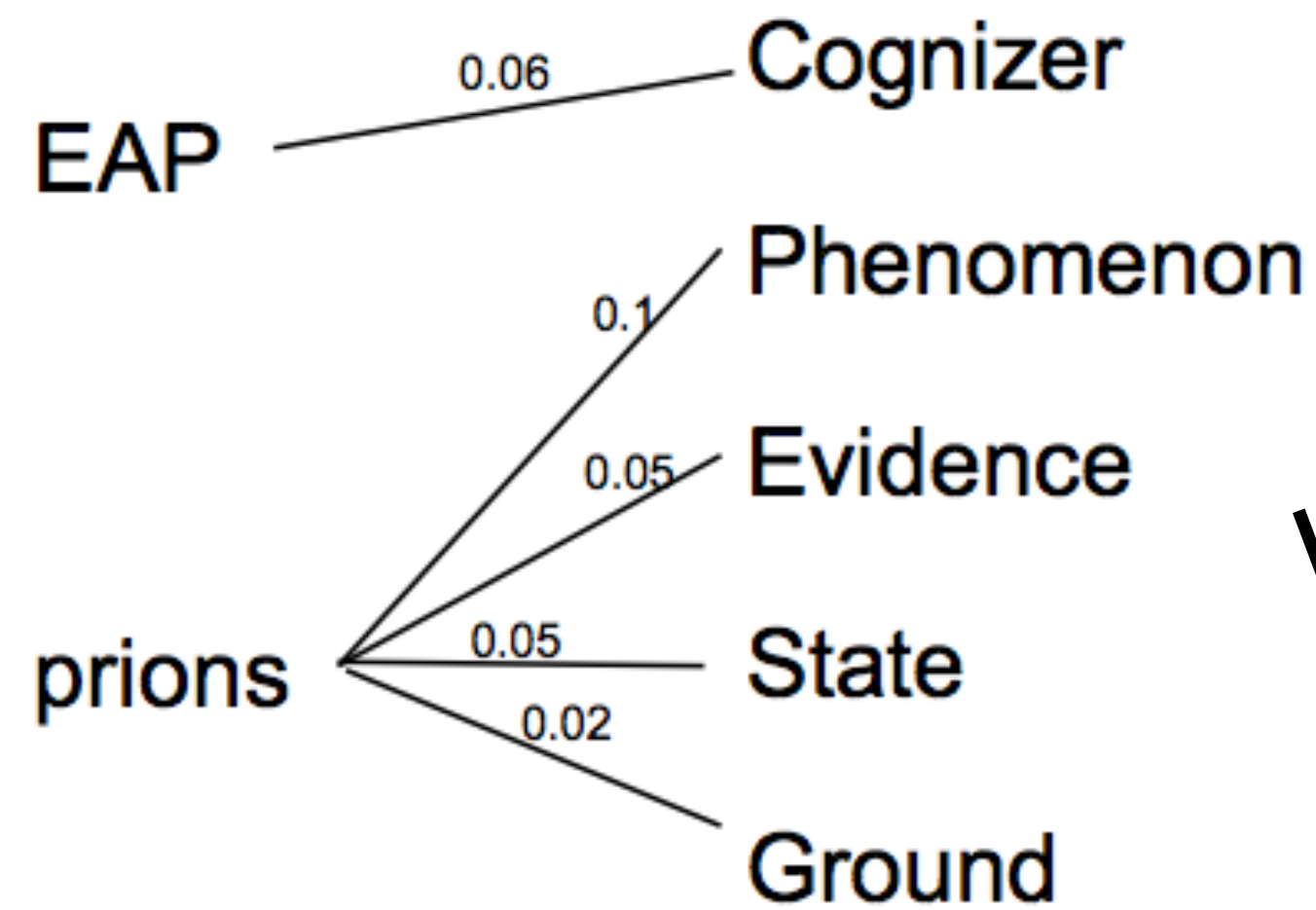
SRL for QA

- ▶ Question and several answer candidates

Q: *Who discovered prions?*

AC1: *In 1997, Stanley B. Prusiner, a scientist in the United States, discovered prions...*

AC2: *Prions were researched by...*



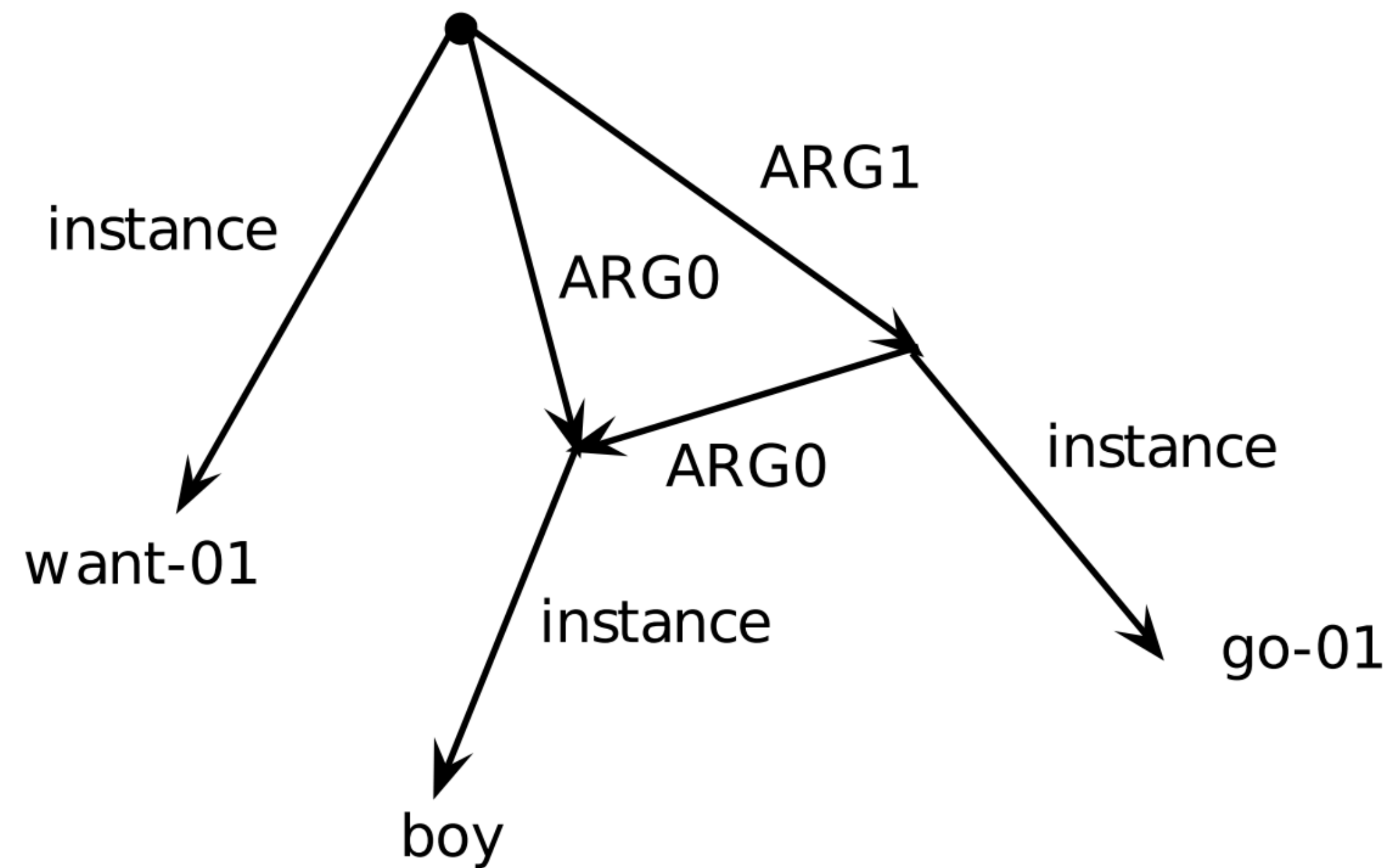
Score by matching expected answer phrase (EAP) against answer candidate (AC)

Shen and Lapata (2007)

Abstract Meaning Representation

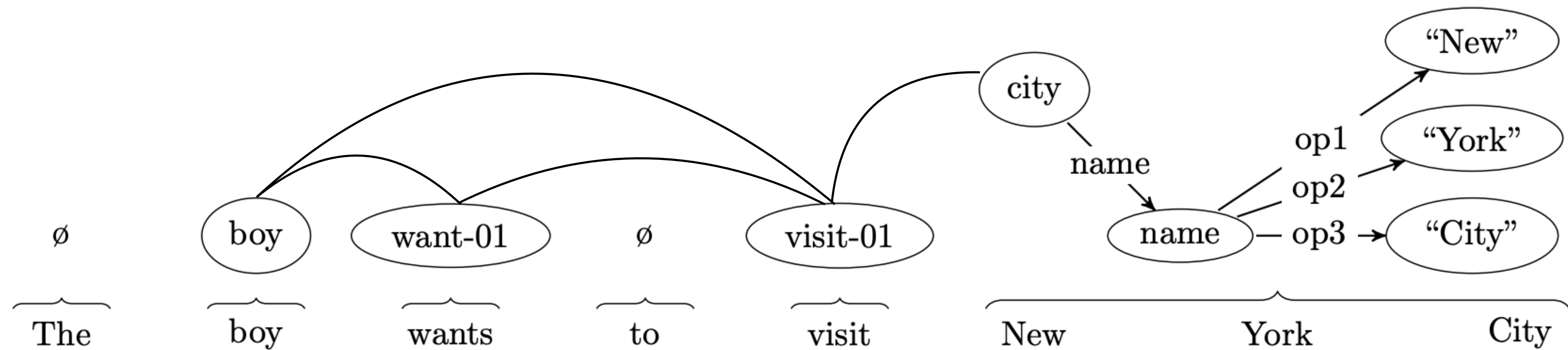
Banarescu et al. (2014)

- ▶ Graph-structured annotation
- ▶ Superset of SRL: full sentence analyses, contains coreference and multi-word expressions as well
- ▶ F1 scores in the 60s: hard!
- ▶ So comprehensive that it's hard to predict, but still doesn't handle tense or some other things...



The boy wants to go

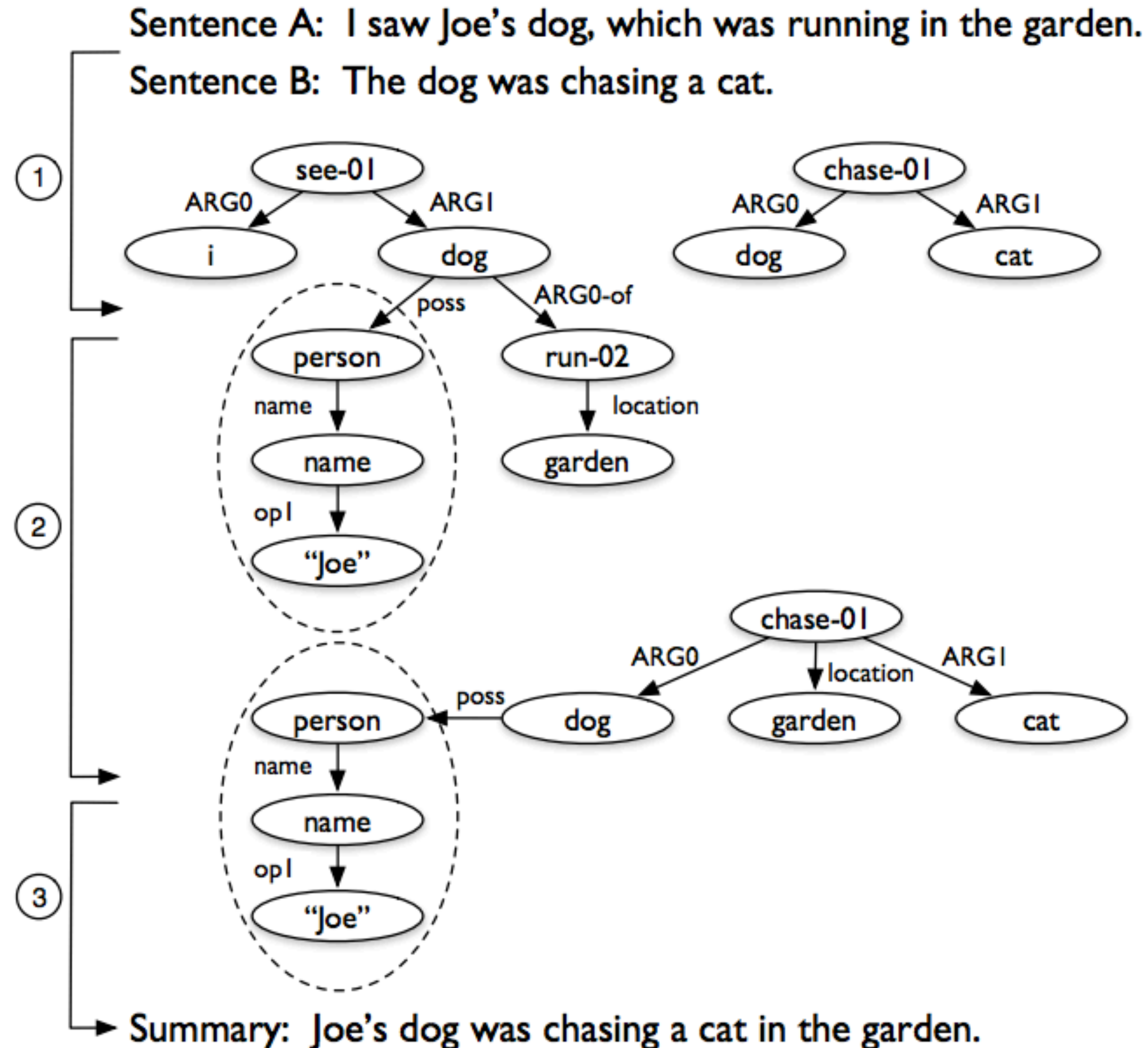
Abstract Meaning Representation



- ▶ First predict mapping from concepts to graph nodes (many-to-many)
- ▶ Then use an edge scoring module similar to dependency parsers to predict edges
- ▶ Predicting a coherent graph is *hard*, lots of constraints on it and no dynamic program

Flanigan et al. (2016), Lyu et al. (2018)

Summarization with AMR



- ▶ Merge AMRs across multiple sentences
- ▶ Summarization = subgraph extraction
- ▶ No real systems actually work this way (more when we talk about summarization)

Slot Filling

Slot Filling

- ▶ Most conservative, narrow form of IE

magnitude

time

Indian Express — A massive earthquake of magnitude 7.3 struck Iraq on Sunday, 103 kms (64 miles) southeast of the city of As-Sulaymaniyah, the US Geological Survey said, reports Reuters. US Geological Survey initially said the quake was of a magnitude 7.2, before revising it to 7.3.

epicenter

Slot Filling

- ▶ Most conservative, narrow form of IE

magnitude

time

Indian Express — A massive earthquake of **magnitude 7.3** struck Iraq on **Sunday**, 103 kms **(64 miles) southeast of the city of As-Sulaymaniyah**, the US Geological Survey said, reports Reuters. US Geological Survey initially said the quake was of a magnitude 7.2, before revising it to 7.3.

epicenter

Speaker: **Alan Clark** speaker

“Gender Roles in the Holy Roman Empire” title

Allagher Center Main Auditorium location

This talk will discuss...

- ▶ Old work: HMMs, later CRFs trained per role

Slot Filling: MUC

Template

(a)

SELLER	BUSINESS	ACQUIRED	PURCHASER
CSR Limited	Oil and Gas	Delhi Fund	Esso Inc.

Document

(b) [S CSR] has said that [S it] has sold [S its] [B oil interests] held in [A Delhi Fund]. [P Esso Inc.] did not disclose how much [P they] paid for [A Dehli].

- ▶ Key aspect: need to combine information across multiple mentions of an entity using coreference

Relation Extraction

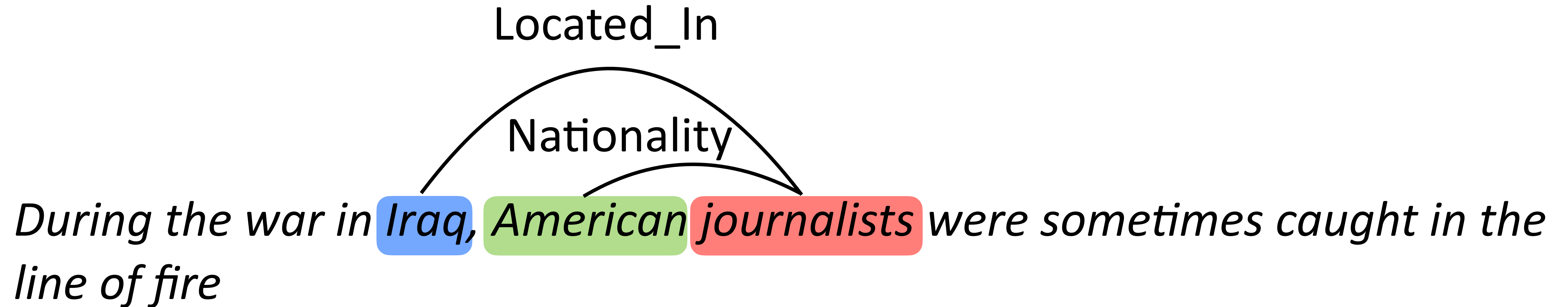
Relation Extraction

- ▶ Extract entity-relation-entity triples from a fixed inventory



Relation Extraction

- ▶ Extract entity-relation-entity triples from a fixed inventory



- ▶ Use NER-like system to identify entity spans, classify relations between entity pairs with a classifier
- ▶ Systems can be feature-based or neural, look at surface words, syntactic features (dependency paths), semantic roles

Relation Extraction

- ▶ Extract entity-relation-entity triples from a fixed inventory



- ▶ Use NER-like system to identify entity spans, classify relations between entity pairs with a classifier
- ▶ Systems can be feature-based or neural, look at surface words, syntactic features (dependency paths), semantic roles
- ▶ **Problem: limited data for scaling to big ontologies** ACE (2003-2005)

Hearst Patterns

- ▶ Syntactic patterns especially for finding hypernym-hyponym pairs (“is a” relations)

Y is a X

Berlin is a city

X such as [list]

cities such as Berlin, Paris, and London.

other X including Y

other cities including Berlin

- ▶ Totally unsupervised way of harvesting world knowledge for tasks like parsing and coreference (Bansal and Klein, 2011-2012)

Distant Supervision

- ▶ Lots of relations in our knowledge base already (e.g., 23,000 film-director relations); use these to bootstrap more training data
- ▶ If two entities in a relation appear in the same sentence, assume the sentence expresses the relation

Distant Supervision

- ▶ Lots of relations in our knowledge base already (e.g., 23,000 film-director relations); use these to bootstrap more training data
- ▶ If two entities in a relation appear in the same sentence, assume the sentence expresses the relation

Director

[Steven Spielberg]'s film [Saving Private Ryan] is loosely based on the brothers' story

Allison co-produced the Academy Award-winning [Saving Private Ryan], directed by [Steven Spielberg]

Director

Distant Supervision

- ▶ Learn decently accurate classifiers for ~100 Freebase relations
- ▶ Could be used to crawl the web and expand our knowledge base

Relation name	100 instances			1000 instances		
	Syn	Lex	Both	Syn	Lex	Both
/film/director/film	0.49	0.43	0.44	0.49	0.41	0.46
/film/writer/film	0.70	0.60	0.65	0.71	0.61	0.69
/geography/river/basin_countries	0.65	0.64	0.67	0.73	0.71	0.64
/location/country/administrative_divisions	0.68	0.59	0.70	0.72	0.68	0.72
/location/location/contains	0.81	0.89	0.84	0.85	0.83	0.84
/location/us_county/county_seat	0.51	0.51	0.53	0.47	0.57	0.42
/music/artist/origin	0.64	0.66	0.71	0.61	0.63	0.60
/people/deceased_person/place_of_death	0.80	0.79	0.81	0.80	0.81	0.78
/people/person/nationality	0.61	0.70	0.72	0.56	0.61	0.63
/people/person/place_of_birth	0.78	0.77	0.78	0.88	0.85	0.91
Average	0.67	0.66	0.69	0.68	0.67	0.67

FewRel: more challenging setting

- ▶ Treats relation classification as a few-shot classification problem
- ▶ 100 classes x 700 instances, goal is to generalize to each class **with just a few instances**
- ▶ BERT can handle this fairly well (Soares et al., 2019)
- ▶ “FewRel 2.0”: new dataset with “none of the above” type, which makes things much harder

Supporting Set	
(A) capital_of	(1) <i>London</i> is the capital of <i>the U.K.</i> (2) <i>Washington</i> is the capital of <i>the U.S.A.</i>
(B) member_of	(1) <i>Newton</i> served as the president of <i>the Royal Society.</i> (2) <i>Leibniz</i> was a member of <i>the Prussian Academy of Sciences.</i>
(C) birth_name	(1) <i>Samuel Langhorne Clemens</i> , better known by his pen name <i>Mark Twain</i> , was an American writer. (2) <i>Alexei Maximovich Peshkov</i> , primarily known as <i>Maxim Gorky</i> , was a Russian and Soviet writer.
Test Instance	
(A) or (B) or (C)	<i>Euler</i> was elected a foreign member of <i>the Royal Swedish Academy of Sciences.</i>

Han et al. (2018), Gao et al. (2019)

Open IE

Open Information Extraction

- ▶ “Open”ness — want to be able to extract all kinds of information from open-domain text
- ▶ Acquire commonsense knowledge just from “reading” about it, but need to process lots of text (“machine reading”)
- ▶ Typically no fixed relation inventory

TextRunner

- ▶ Extract positive examples of (e, r, e) triples via parsing and heuristics
- ▶ Train a Naive Bayes classifier to filter triples from raw text: uses features on POS tags, lexical features, stopwords, etc.

TextRunner

- ▶ Extract positive examples of (e, r, e) triples via parsing and heuristics
- ▶ Train a Naive Bayes classifier to filter triples from raw text: uses features on POS tags, lexical features, stopwords, etc.

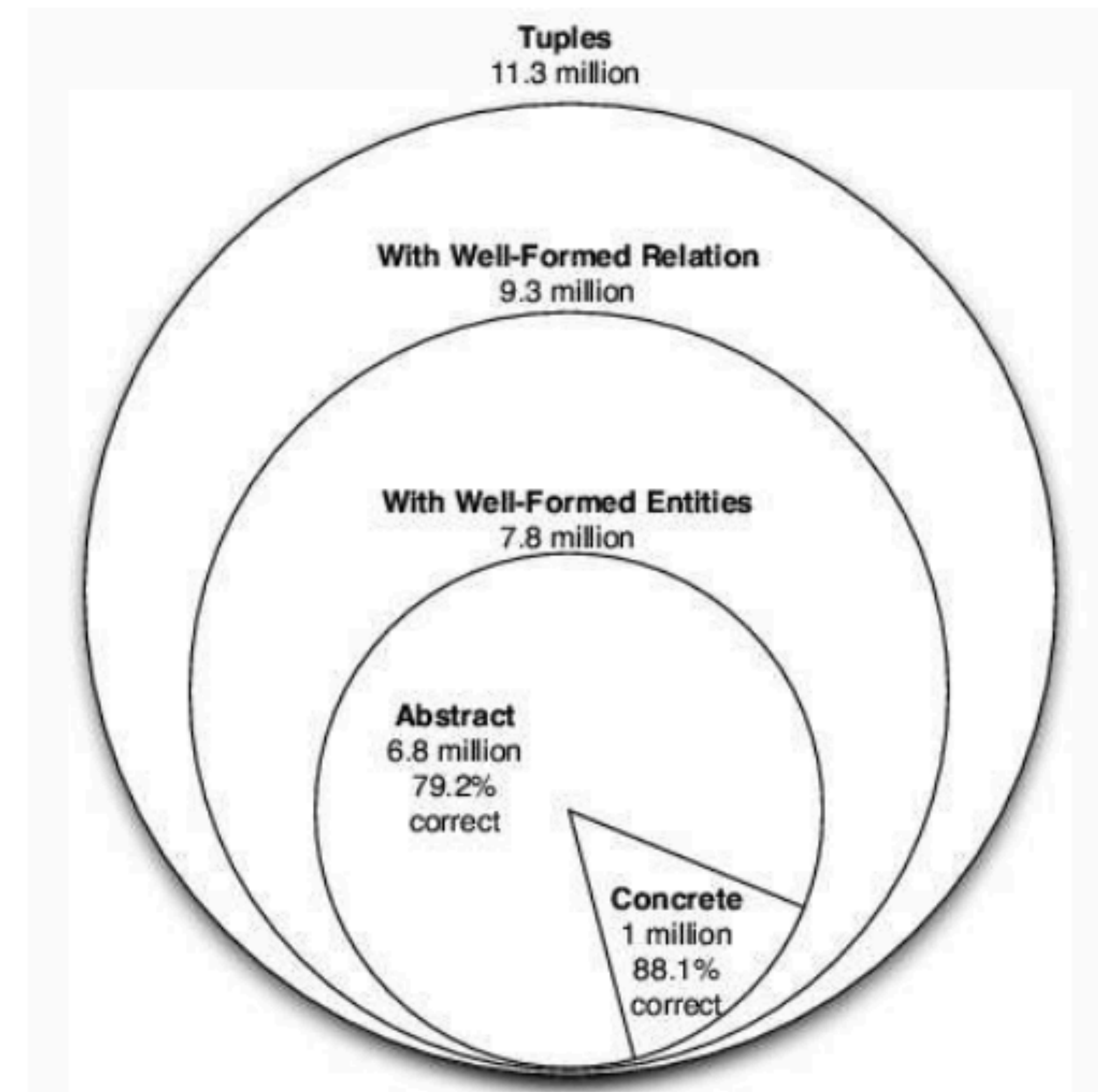
Barack Obama, 44th president of the United States, was born on August 4, 1961 in Honolulu

=> Barack_Obama, was born in, Honolulu

- ▶ 80x faster than running a parser (which was slow in 2007...)
- ▶ Use multiple instances of extractions to assign probability to a relation

Exploiting Redundancy

- ▶ 9M web pages / 133M sentences
- ▶ 2.2 tuples extracted per sentence, filter based on probabilities
- ▶ **Concrete:** e.g., (Tesla, invented, coil transformer)
Abstract: underspecified, e.g., (Einstein, derived, theory)
- ▶ Hard to evaluate: can assess precision of extracted facts, but how do we know recall?



Banko et al. (2007)

ReVerb

- ▶ More constraints: open relations have to begin with verb, end with preposition, be contiguous (e.g., *was born on*)

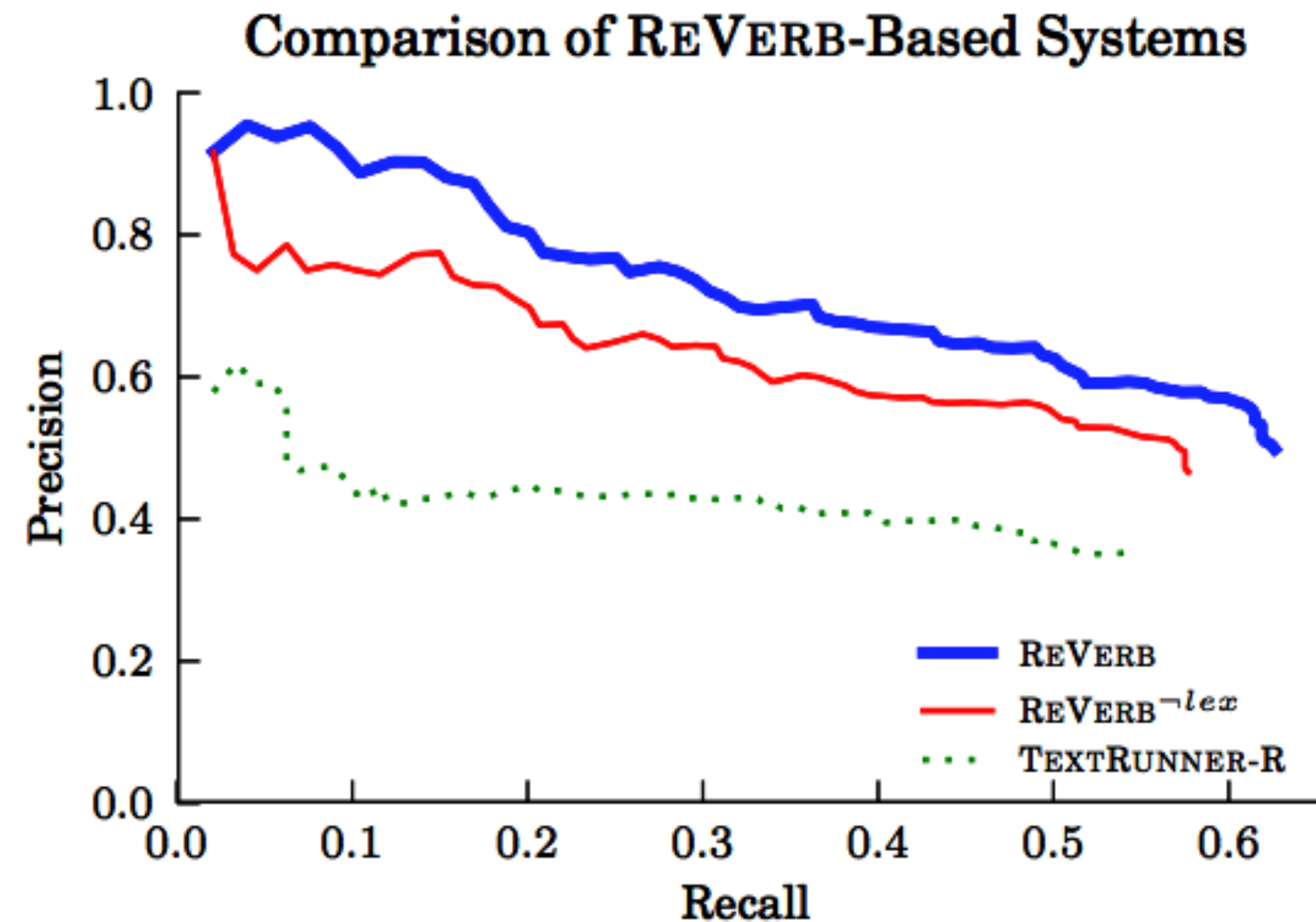
ReVerb

- ▶ More constraints: open relations have to begin with verb, end with preposition, be contiguous (e.g., *was born on*)
- ▶ Extract more meaningful relations, particularly with light verbs

is	is an album by, is the author of, is a city in
has	has a population of, has a Ph.D. in, has a cameo in
made	made a deal with, made a promise to
took	took place in, took control over, took advantage of
gave	gave birth to, gave a talk at, gave new meaning to
got	got tickets to, got a deal on, got funding from

ReVerb

- ▶ For each verb, identify the longest sequence of words following the verb that satisfy a POS regex ($V \cdot^* P$) and which satisfy heuristic lexical constraints on specificity
- ▶ Find the nearest arguments on either side of the relation
- ▶ Annotators labeled relations in 500 documents to assess recall



QA from Open IE

(a) **CCG parse** builds an underspecified semantic representation of the sentence.

Former	municipalities	in	Brandenburgh
N/N	N	$N \setminus N/NP$	NP
$\lambda f \lambda x. f(x) \wedge former(x)$	$\lambda x. municipalities(x)$	$\lambda f \lambda x \lambda y. f(y) \wedge in(y, x)$	Brandenburg
N	$N \setminus N$	$N \setminus N$	$N \setminus N$
$\lambda x. former(x) \wedge municipalities(x)$	$\lambda f \lambda y. f(y) \wedge in(y, Brandenburg)$	$\lambda f \lambda y. f(y) \wedge in(y, Brandenburg)$	$\lambda f \lambda y. f(y) \wedge in(y, Brandenburg)$
N	N	N	N
$l_0 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$			

(b) **Constant matches** replace underspecified constants with Freebase concepts

- $l_0 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$
- $l_1 = \lambda x. former(x) \wedge municipalities(x) \wedge in(x, Brandenburg)$
- $l_2 = \lambda x. former(x) \wedge municipalities(x) \wedge location.containedby(x, Brandenburg)$
- $l_3 = \lambda x. former(x) \wedge OpenRel(x, Municipality) \wedge location.containedby(x, Brandenburg)$
- $l_4 = \lambda x. OpenType(x) \wedge OpenRel(x, Municipality) \wedge location.containedby(x, Brandenburg)$

Takeaways

- ▶ Relation extraction: can collect data with distant supervision, use this to expand knowledge bases
- ▶ Slot filling: tied to a specific ontology, but gives fine-grained information
- ▶ Open IE: extracts lots of things, but hard to know how good or useful they are
 - ▶ Can combine with standard question answering
 - ▶ Add new facts to knowledge bases
- ▶ Many, many applications and techniques