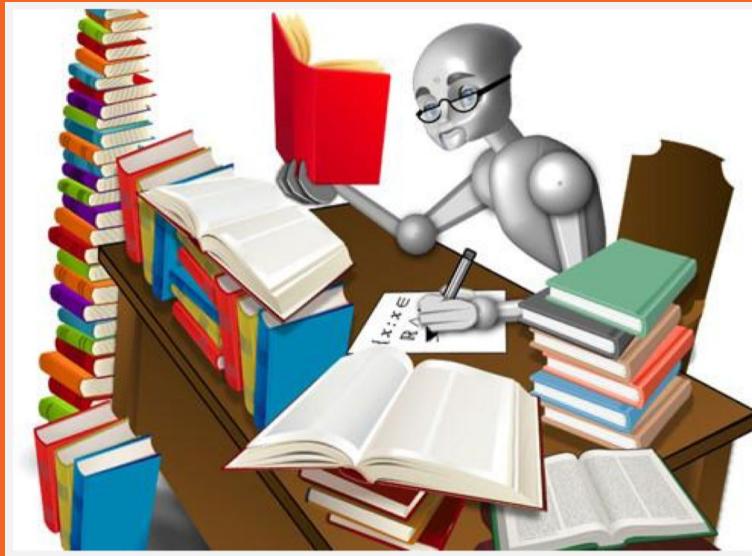


# Relation Extraction & Knowledge Bases



DARPA

Sam Bowman

*including slides & figures from Bill MacCartney, Dan Jurafsky, Rion Snow, Jim Martin, Chris Manning, William Cohen, Michele Banko, Mike Mintz, Steven Bills, and others.*

---

# Announcements

- HW2 grades out
- Proposal due today, expect feedback no later than the end of break
- HPC access available

---

# The Big Question

What kind of a thing is the meaning of a sentence?

---

# The Big Question

~~What kind of a thing is the meaning of a sentence?~~

What can you do with a sentence if you know its meaning?

---

---

# The Big Question

~~What kind of a thing is the meaning of a sentence?~~

What can you do with a sentence if you know its meaning?

One answer:

- Judge its sentiment.

Another answer:

- Extract the literal information it contains.
-

---

# The Big Question

How do we represent this information?

This week:  $(a, R, b)$

---

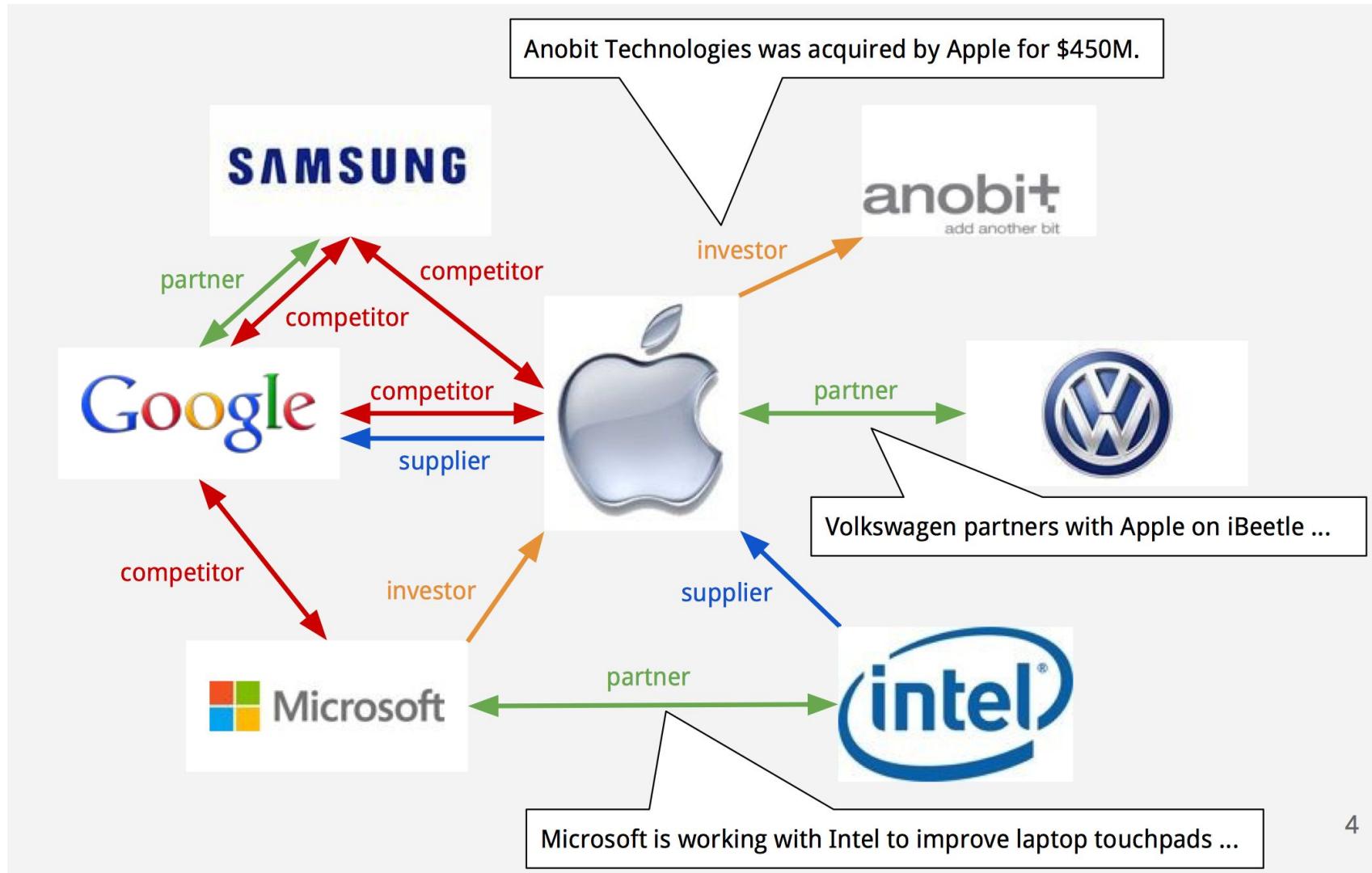
# Relations, Knowledge Bases, & Relation Extraction

# Relations

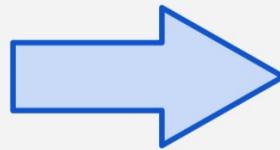
CHICAGO (AP) — Citing high fuel prices, United Airlines said Friday it has increased fares by \$6 per round trip on flights to some cities also served by lower-cost carriers. **American Airlines**, a **unit of AMR**, immediately matched the move, spokesman **Tim Wagner** said. **United**, a **unit of UAL**, said the increase took effect Thursday night and applies to most routes where it competes against discount carriers, such as Chicago to Dallas and Atlanta and Denver to San Francisco, Los Angeles and New York.

Subject	Relation	Object
American Airlines	subsidiary	AMR
Tim Wagner	employee	American Airlines
United Airlines	subsidiary	UAL

# A Knowledge Base



# Biomedical Information Extraction: Gene Regulation



Subject	Relation	Object
p53	is_a	protein
Bax	is_a	protein
p53	has_function	apoptosis
Bax	has_function	induction
apoptosis	involved_in	cell_death
Bax	is_in	mitochondrial outer membrane
Bax	is_in	cytoplasm
apoptosis	related_to	caspase activation
...	...	...

textual abstract:  
summary for human

structured knowledge extraction:  
summary for machine

# WordNet as Knowledge Base

```
vehicle
  craft
    aircraft
      airplane
      dirigible
      helicopter
    spacecraft
    watercraft
      boat
      ship
      yacht
  rocket
    missile
    multistage rocket
wheeled vehicle
  automobile
  bicycle
  locomotive
  wagon
```

(airplane, hyponym\_of, vehicle)  
(rocket, hypernym\_of, missile)  
...

# WordNet as *Incomplete* Knowledge Base

In WordNet 3.1	Not in WordNet 3.1
insulin	leptin
progesterone	pregnenolone
combustibility	affordability
navigability	reusability
HTML	XML
Google, Yahoo	Microsoft, IBM

Esp. for specific domains: restaurants, auto parts, finance

Esp. neologisms: iPad, selfie, bitcoin, twerking, Hadoop, dubstep

# Expanding WordNet

Mirror ran a headline questioning whether the killer's actions were a result of playing **Call of Duty, a first-person shooter game ...**



Melee, in video game terms, is a style of elbow-drop hand-to-hand combat popular in **first-person shooters and other shooters.**



Tower defense is a kind of real-time strategy game in which the goal is to protect an area/place/locality and prevent enemies from reaching ...



video game  
action game  
ball and paddle game  
Breakout  
platform game  
Donkey Kong  
shooter  
arcade shooter  
Space Invaders  
first-person shooter  
Call of Duty  
third-person shooter  
Tomb Raider  
adventure game  
text adventure  
graphic adventure  
strategy game  
4X game  
Civilization  
tower defense  
Plants vs. Zombies

# Freebase (now WikiData)

Freebase: 20K relations, 40M entities, 600M assertions

Curation is an ongoing challenge — things change!

Relies heavily on relation extraction from the web

## /film/film/starring

Bad Words	Jason Bateman
Divergent	Shailene Woodley
Non-Stop	Liam Neeson

## /organization/organization/parent

WhatsApp	Facebook
Nest Labs	Google
Nokia	Microsoft

## /music/artist/track

Macklemore	White Privilege
Phantogram	Mouthful of Diamonds
Lorde	Royals

## /people/person/date\_of\_death

Nelson Mandela	2013-12-05
Paul Walker	2013-11-30
Lou Reed	2013-10-27

# Information Extraction Methods

# Information Extraction Methods

## Hand-Built Patterns

# Hand-Built Rules

- Intuition from Hearst (1992)

*Agar is a substance prepared from a mixture of red algae, such as **Gelidium**, for laboratory or industrial use.*

- What does *Gelidium* mean?
- How do you know?



# Some Rules for Hyponymy

Ys such as X ((, X)\* (, and/or) X)

such Ys as X...

X... or other Ys

X... and other Ys

Ys including X...

Ys, especially X...

## Some Results: X, especially Y

The best part of the night was seeing all of the tweets of the performers, especially Miley Cyrus and Drake. ✓

Those child stars, especially Miley Cyrus, I feel like you have to put the fault on the media. ✓

Kelly wasn't shy about sharing her feelings about some of the musical acts, especially Miley Cyrus. ✓

Rihanna was bored with everything at the MTV VMAs, especially Miley Cyrus. ✗

The celebrities enjoyed themselves while sipping on delicious cocktails, especially Miley Cyrus who landed the coveted #1 spot. ✗

None of these girls are good idols or role models, especially Miley Cyrus. ✗

# This isn't ideal.

- Requires hand-building patterns for each relation!
  - and every language!
  - hard to write; hard to maintain
  - there are zillions of them
  - domain-dependent
- Don't want to do this for all possible relations!
- Plus, we'd like better accuracy
  - Hearst: 66% accuracy on hyponym extraction
  - Berland & Charniak: 55% accuracy on meronyms

# Information Extraction Methods

Bootstrapping

# Bootstrapping

- If you don't have enough annotated text to train on ...
- But you do have:
  - some **seed instances** of the relation
  - (or some patterns that work pretty well)
  - and lots & lots of **unannotated text** (e.g., the web)
- ... can you use those seeds to do something useful?
- Bootstrapping can be considered *semi-supervised*

# Bootstrapping: Collecting Patterns

- Target relation: *burial place*
- Seed tuple: [Mark Twain, Elmira]
- Grep/Google for “Mark Twain” and “Elmira”
  - “Mark Twain is buried in Elmira, NY.”
    - X is buried in Y
  - “The grave of Mark Twain is in Elmira”
    - The grave of X is in Y
  - “Elmira is Mark Twain’s final resting place”
    - Y is X’s final resting place
- Use those patterns to search for new tuples

# Bootstrapping: Collecting Relations

**Lincoln is buried in Springfield, Illinois - May 04, 1865 - HISTORY.com**

[www.history.com>this-day-in-history/lincoln-is-buried-in-springfield-illinois](http://www.history.com>this-day-in-history/lincoln-is-buried-in-springfield-illinois) ▾

On this day in 1865, Abraham Lincoln is laid to rest in his hometown of Springfield, Illinois. His funeral train had traveled through 180 cities and seven states before reaching Springfield. ... Lincoln's son Willie, who died at age 11 from typhoid fever in 1862 and had originally ...

**Tsvangirai is buried in rural home | Free & Fair Zimbabwe Election**

[zimbabwelection.com/2018/02/21/tsvangirai-buried-rural-home/](http://zimbabweelection.com/2018/02/21/tsvangirai-buried-rural-home/) ▾

Feb 21, 2018 - **Tsvangirai is buried in rural home** Thousands of people in Zimbabwe have gathered for the burial of opposition leader Morgan Tsvangirai, whose death from cancer exacerbated divisions within a movement preparing for elections this year. A hearse carrying Tsvangirai's body in a white casket on Tuesday ...

**Zimbabwean opposition leader is buried in rural home | Fox News**

[www.foxnews.com/world/.../zimbabwean-opposition-leader-is-buried-in-rural-home.htm...](http://www.foxnews.com/world/.../zimbabwean-opposition-leader-is-buried-in-rural-home.htm...)

Feb 20, 2018 - Thousands of people in Zimbabwe have gathered for the burial of opposition leader Morgan Tsvangirai, whose death from cancer exacerbated divisions within a movement preparing for elections this year.

**Who is Buried in the Famous Shrine of St James in Santiago de ...**

[www.ancient-origins.net/.../who-buried-famous-shrine-st-james-santiago-de-compostel...](http://www.ancient-origins.net/.../who-buried-famous-shrine-st-james-santiago-de-compostel...) ▾

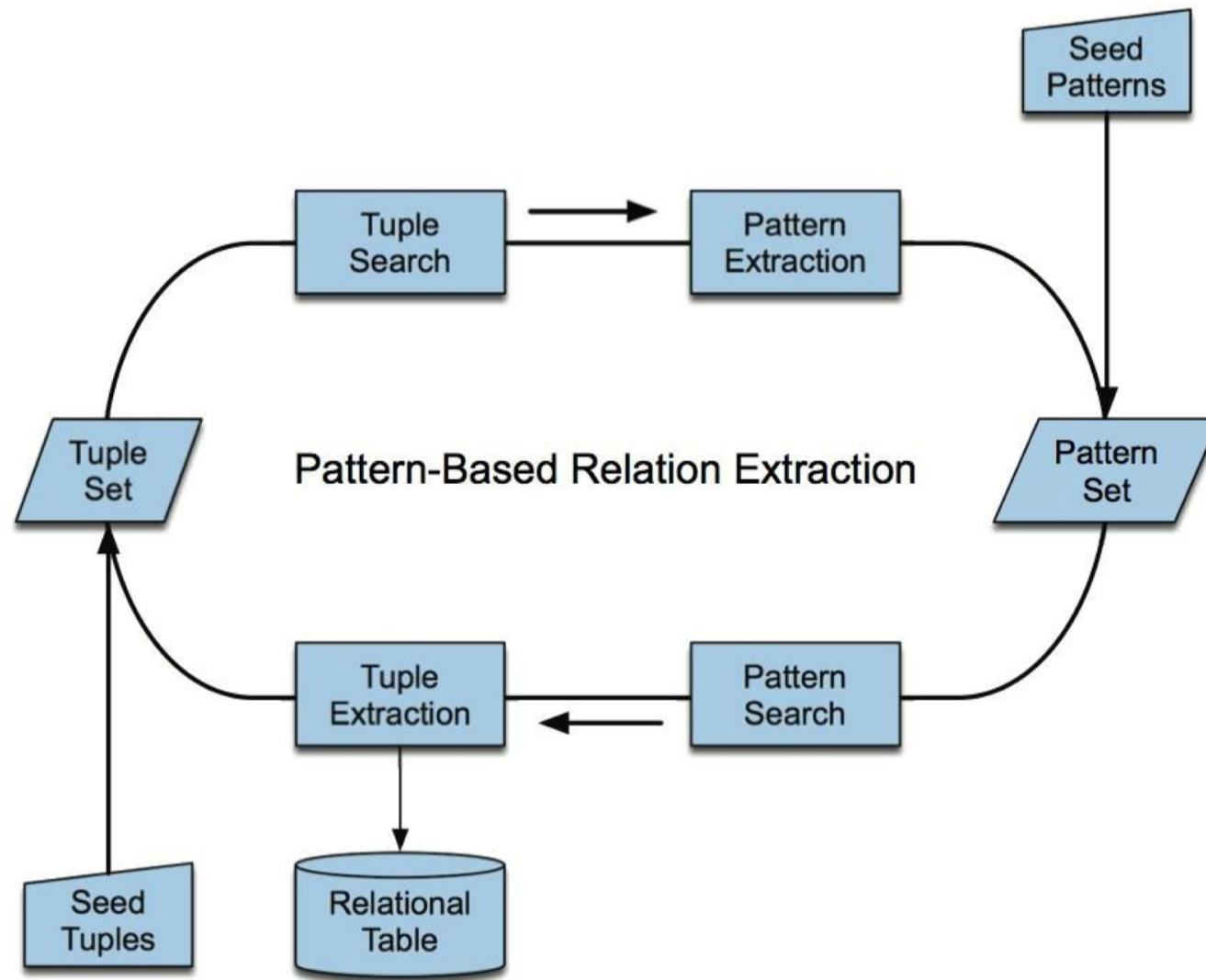
Apr 18, 2016 - The legend of the Apostle John's brother is one of the most important stories in Spanish Christianity. According to legend, a man who was a friend of Jesus **is buried in** the cathedral in Santiago de Compostela. Legends say that the remains of St James were transported from Jerusalem to Galicia, where he ...

**Who is buried in the Hoover Dam? - io9 - Gizmodo**

<https://io9.gizmodo.com/5893183/who-is-buried-in-the-hoover-dam> ▾

Mar 16, 2012 - The Hoover Dam is one of the most phenomenal structures in modern history. This 1244 feet long, 660 feet thick, and 726 feet high concrete behemoth holds back so much water that it deformed the earth's crust and caused 600 small earthquakes in the decade after its construction.

# Bootstrapping: The Loop



# DIPRE (Brin et al. '99)

Extract (author, book) pairs

Start with these 5 seeds:

Author	Book
Isaac Asimov	The Robots of Dawn
David Brin	Startide Rising
James Gleick	Chaos: Making a New Science
Charles Dickens	Great Expectations
William Shakespeare	The Comedy of Errors

Learn these patterns:

URL Prefix	Text Pattern
<code>www.sff.net/locus/c.*</code>	<code>&lt;LI&gt;&lt;B&gt;title&lt;/B&gt; by author (</code>
<code>dns.city-net.com/~lmann/awards/hugos/1984.html</code>	<code>&lt;i&gt;title&lt;/i&gt; by author (</code>
<code>dolphin.upenn.edu/~dcummins/texts/sf-award.htm</code>	<code>author    title    (</code>

Iterate: use these patterns to get more instances & patterns...

# DIPRE (Brin et al. '99)

Extract (author, book) pairs  
Start with these 5 seeds:

Author	Book
Isaac Asimov	The Robots of Dawn
David Brin	Startide Rising
James Gleick	Chaos: Making a New Science
Charles Dickens	Great Expectations
William Shakespeare	The Comedy of Errors



Learn these patterns:

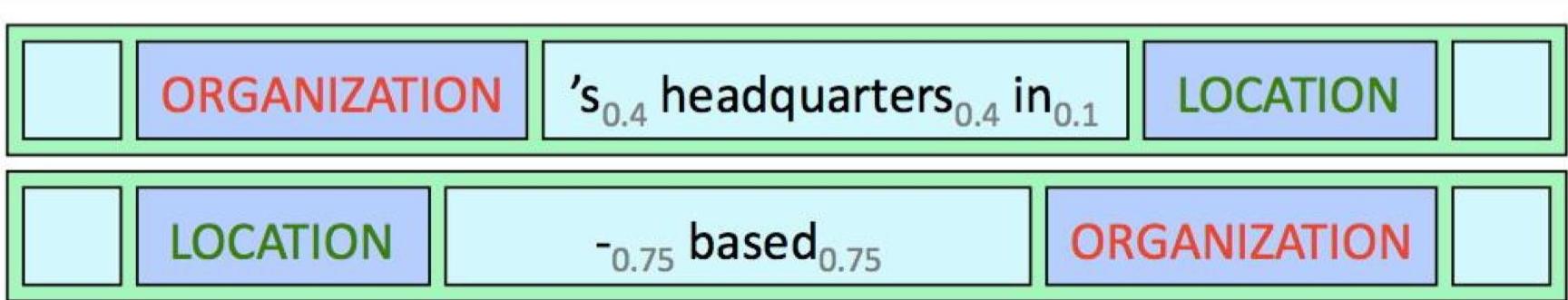
URL Prefix	Text Pattern
<code>www.sff.net/locus/c.*</code>	<code>&lt;LI&gt;&lt;B&gt;title&lt;/B&gt; by author (</code>
<code>dns.city-net.com/~lmann/awards/hugos/1984.html</code>	<code>&lt;i&gt;title&lt;/i&gt; by author (</code>
<code>dolphin.upenn.edu/~dcummins/texts/sf-award.htm</code>	<code>author    title    (</code>

Iterate: use these patterns to get more instances & patterns...

# SNOWBALL (Agichtein & Gravano '00)

New idea: require that X and Y be named entities of particular types

Organization	Location of Headquarters
Microsoft	Redmond
Exxon	Irving
IBM	Armonk
Boeing	Seattle
Intel	Santa Clara



-

# Brief Detour: Named Entity Recognition

# Named Entity Recognition

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

# Entity Types

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	The Mt. Sanitas loop is in Sunshine Canyon.
Geo-Political Entity	GPE	countries, states, provinces	Palo Alto is raising the fees for parking.
Facility	FAC	bridges, buildings, airports	Consider the Tappan Zee Bridge.
Vehicles	VEH	planes, trains, automobiles	It was a classic Ford Falcon.

**Figure 21.1** A list of generic named entity types with the kinds of entities they refer to.

# Entity Types & Ambiguity

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	The Mt. Sanitas loop is in Sunshine Canyon.
Geo-Political	GPE	countries, states, provinces	Palo Alto is raising the fees for parking.
Entity			
Facility	FAC	bridges, buildings, airports	Consider the Tappan Zee Bridge.
Vehicles	VEH	planes, trains, automobiles	It was a classic Ford Falcon.

**Figure 21.1** A list of generic named entity types with the kinds of entities they refer to.

Name	Possible Categories
Washington	Person, Location, Political Entity, Organization, Vehicle
Downing St.	Location, Organization
IRA	Person, Organization, Monetary Instrument
Louis Vuitton	Person, Organization, Commercial Product

**Figure 21.2** Common categorical ambiguities associated with various proper names.

# Span Tagging as Classification

Words	IOB Label
American	B-ORG
Airlines	I-ORG
,	O
a	O
unit	O
of	O
AMR	B-ORG
Corp.	I-ORG
,	O
immediately	O
matched	O
the	O
move	O
,	O
spokesman	O
Tim	B-PER
Wagner	I-PER
said	O
.	O



# NER in Practice

- We have reasonably mature systems for most languages
- Standard packages
  - Stanford CoreNLP (legacy, Java)
  - spaCy (Python)
- If targeting a new domain, it can help to do some feature engineering.
  - ...especially *gazetteers* and other lists of proper names
- Closely related to the problem of *coreference resolution* (also useful for information extraction):
  - Working out whether two named entity or pronoun mentions refer to the same thing:

[Apple]<sup>1</sup> ... [Microsoft]<sup>2</sup> ... it<sup>2</sup> ... [the company]<sup>2</sup>

-

# Poll:

# Named Entity Recognition

-

# Nishant: Exercise on Bootstrapping and OpenIE

Intermission

# Bootstrapping is still not ideal.

- Requires that we have seeds for each relation
  - Sensitive to original set of seeds
- Big problem of semantic drift at each iteration
- Precision tends to be not that high
- Generally have lots of parameters to be tuned
- No probabilistic interpretation
  - Hard to know how confident to be in each result

# Information Extraction Methods

## Supervised Learning

# Supervised Learning

For each pair of entities in a sentence, predict the *relation type* (if any) that holds between them.

The supervised approach requires:

- Defining an inventory of relation types
- Collecting labeled training data (the hard part!)
- Designing a feature representation
- Choosing a classifier: Naïve Bayes, MaxEnt, SVM, ...
- Evaluating the results

# The Automatic Content Extraction (ACE) Data

Type	Subtype
ART (artifact)	User-Owner-Inventor-Manufacturer
GEN-AFF (General affiliation)	Citizen-Resident-Religion-Ethnicity, Org-Location
METONYMY*	<i>None</i>
ORG-AFF (Org-affiliation)	Employment, Founder, Ownership, Student-Alum, Sports-Affiliation, Investor-Shareholder, Membership
PART-WHOLE (part-to-whole)	Artifact, Geographical, Subsidiary
PER-SOC* (person-social)	Business, Family, Lasting-Personal
PHYS* (physical)	Located, Near

Relation types used in the ACE 2008 evaluation

# Feature Engineering!

- Lightweight features — require little pre-processing
  - Bags of words & bigrams between, before, and after the entities
  - Stemmed versions of the same
  - The types of the entities
  - The distance (number of words) between the entities
- Medium-weight features — require base phrase chunking
  - Base-phrase chunk paths
  - Bags of chunk heads
- Heavyweight features — require full syntactic parsing
  - Dependency-tree paths between the entities
  - Constituent-tree paths between the entities
  - Tree distance between the entities
  - Presence of particular constructions in a constituent structure

# Some Results

Methods	F
Sentence level in Ji and Grishman (2008)	59.7
MaxEnt with local features in Li et al. (2013b)	64.7
Joint beam search with local features in Li et al. (2013b)	63.7
Joint beam search with local and global features in Li et al. (2013b)	65.6
CNN1: CNN without any external features	<b>67.6</b>

Nguyen & Grishman '15

# A Caveat!

Type	Subtype	#Testing Instances	#Correct	#Error	P	R	F
AT		<b>392</b>	<b>224</b>	<b>105</b>	<b>68.1</b>	<b>57.1</b>	<b>62.1</b>
	Based-In	85	39	10	79.6	45.9	58.2
	Located	241	132	120	52.4	54.8	53.5
	Residence	66	19	9	67.9	28.8	40.4
NEAR		<b>35</b>	<b>8</b>	<b>1</b>	<b>88.9</b>	<b>22.9</b>	<b>36.4</b>
	Relative-Location	35	8	1	<b>88.9</b>	22.9	36.4
PART		<b>164</b>	<b>106</b>	<b>39</b>	<b>75.1</b>	<b>64.6</b>	<b>68.6</b>
	Part-Of	136	76	32	70.4	55.9	62.3
	Subsidiary	27	14	23	37.8	51.9	43.8
ROLE		<b>699</b>	<b>443</b>	<b>82</b>	<b>84.4</b>	<b>63.4</b>	<b>72.4</b>
	Citizen-Of	36	25	8	75.8	69.4	72.6
	General-Staff	201	108	46	71.1	53.7	62.3
	Management	165	106	72	59.6	64.2	61.8
	Member	224	104	36	74.3	46.4	57.1
SOCIAL		<b>95</b>	<b>60</b>	<b>21</b>	<b>74.1</b>	<b>63.2</b>	<b>68.5</b>
	Other-Professional	29	16	32	<b>33.3</b>	55.2	41.6
	Parent	25	17	0	<b>100</b>	68.0	81.0

Zhou et al. '05

# Supervised Learning

- Supervised approach can achieve high accuracy
  - At least, for *some* relations
  - If we have lots of hand-labeled training data
- But has significant limitations!
  - Labeling 5,000 relations (+ named entities) is expensive
  - Doesn't generalize to different relations, languages
- Next: beyond supervised relation extraction
  - Distantly supervised relation extraction
  - Unsupervised relation extraction

# Information Extraction Methods

## Distant Supervision

# Distant Supervision

Snow, Jurafsky, Ng. 2005. Learning syntactic patterns for automatic hypernym discovery. NIPS 17

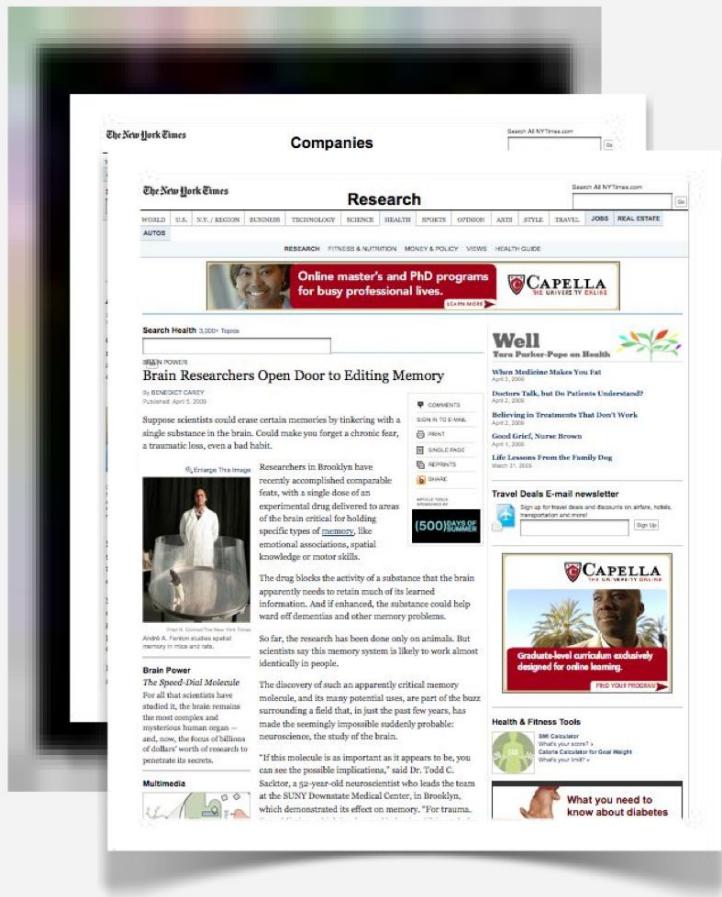
Mintz, Bills, Snow, Jurafsky. 2009. Distant supervision for relation extraction without labeled data. ACL-2009.



- Hypothesis: If two entities belong to a certain relation, any sentence containing those two entities is likely to express that relation
- Key idea: use a *database* of relations to get lots of training examples
  - instead of hand-creating a few seed tuples (bootstrapping)
  - instead of using hand-labeled corpus (supervised)

# Distant Supervision

We construct a noisy training set consisting of occurrences from our corpus that contain a hyponym-hypernym pair from WordNet.



This yields high-signal examples like:

“...consider **authors** like **Shakespeare**...”

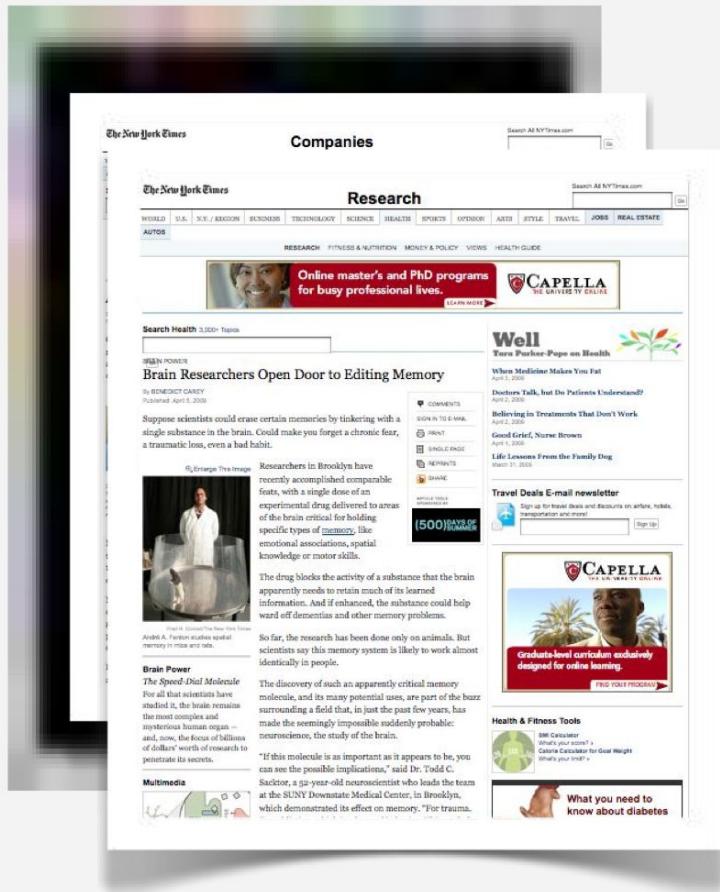
“Some **authors** (including **Shakespeare**)...”

“**Shakespeare** was the **author** of several...”

“**Shakespeare**, **author** of *The Tempest*...”

# Distant Supervision

We construct a noisy training set consisting of occurrences from our corpus that contain a hyponym-hypernym pair from WordNet.



This yields high-signal examples like:

“...consider **authors** like **Shakespeare**...”

“Some **authors** (including **Shakespeare**)...”

“**Shakespeare** was the **author** of several...”

“**Shakespeare**, **author** of *The Tempest*...”

But also noisy examples like:

“The **author** of *Shakespeare in Love*...”

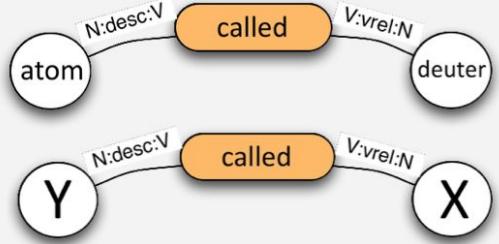
“...**authors** at the **Shakespeare** Festival...”

# The Procedure

1. Take 6M newswire sentences

*... doubly heavy hydrogen atom called deuterium ...*
2. Collect noun pairs

e.g. (atom, deuterium)  
752,311 pairs from 6M sentences of newswire
3. Is pair a hypernym in WordNet?

14,387 yes; 737,924 no
4. Parse the sentences
5. Extract patterns

69,592 dependency paths with >5 pairs
6. Train classifier on patterns

logistic regression with 70K features  
(converted to 974,288 bucketed binary features)

# One of the 70,000 patterns

Pattern: <superordinate> called <subordinate>  
or: <Y> called <X>

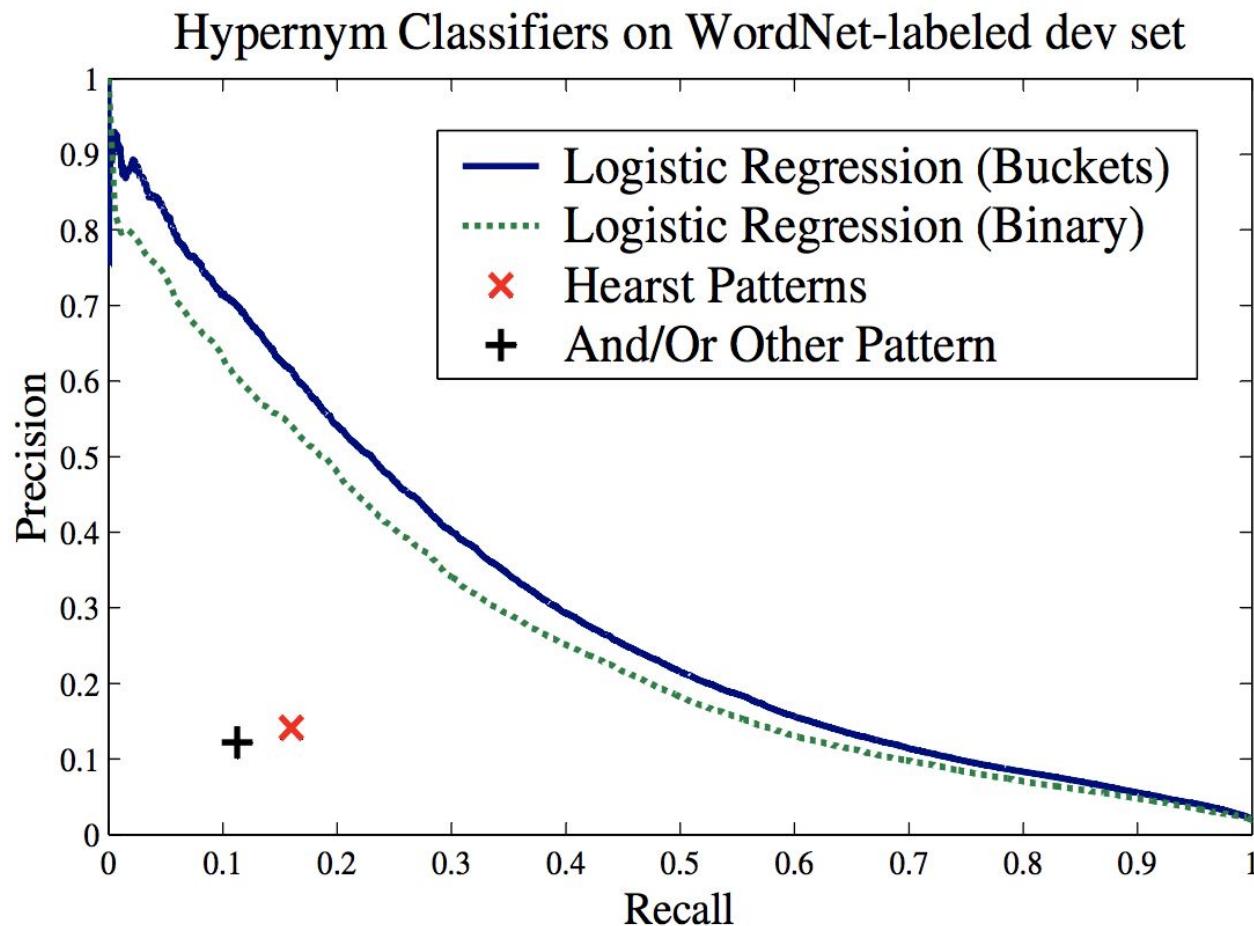
Learned from cases such as:

(**sarcoma**, cancer) ...an uncommon bone **cancer** called osteogenic **sarcoma** and to...  
(**deuterium**, atom) ...heavy water rich in the doubly heavy hydrogen **atom** called **deuterium**.

New pairs discovered:

(**efflorescence**, condition) ...and a **condition** called **efflorescence** are other reasons for...  
(**O'Neal\_inc**, company) ...The **company**, now called **O'Neal Inc.**, was sole distributor of...  
(**hat\_creek\_outfit**, ranch) ...run a small **ranch** called the **Hat Creek Outfit**.  
(**HIV-1**, AIDS virus) ...infected by the **AIDS virus**, called **HIV-1**.  
(**bateau\_mouche**, attraction) ...local sightseeing **attraction** called the **Bateau Mouche**...

# Some Results: Snow '05



-

# Poll: Precision & Recall

# Beyond Hyponymy

Mintz, Bills, Snow, Jurafsky (2009).

Distant supervision for relation extraction without labeled data.



**Training set**



102 relations  
940,000 entities  
1.8 million instances

**Corpus**



1.8 million articles  
25.7 million sentences

# Distant Supervision on Freebase

## Corpus text

Bill Gates founded Microsoft in 1975.  
Bill Gates, founder of Microsoft, ...  
Bill Gates attended Harvard from...  
Google was founded by Larry Page ...

## Freebase

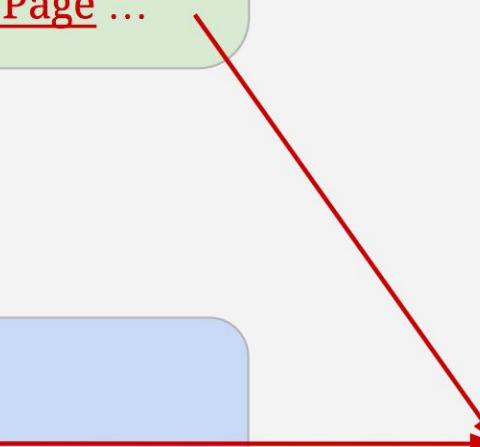
Founder: (Bill Gates, Microsoft)  
Founder: (Larry Page, Google)  
CollegeAttended: (Bill Gates, Harvard)

## Training data

(Bill Gates, Microsoft)  
Label: Founder  
Feature: X founded Y  
Feature: X, founder of Y

(Bill Gates, Harvard)  
Label: CollegeAttended  
Feature: X attended Y

(Larry Page, Google)  
Label: Founder  
Feature: Y was founded by X



# Negative Examples

Can't train a classifier with only positive data! Need negative training data too!

Solution?

Sample 1% of unrelated pairs of entities.

Result: roughly balanced data.

## Corpus text

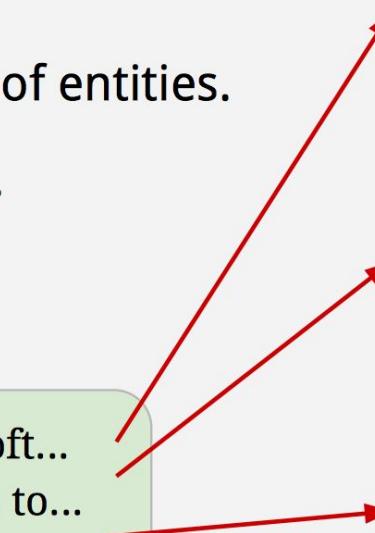
Larry Page took a swipe at Microsoft...  
...after Harvard invited Larry Page to...  
Google is Bill Gates' worst fear ...

## Training data

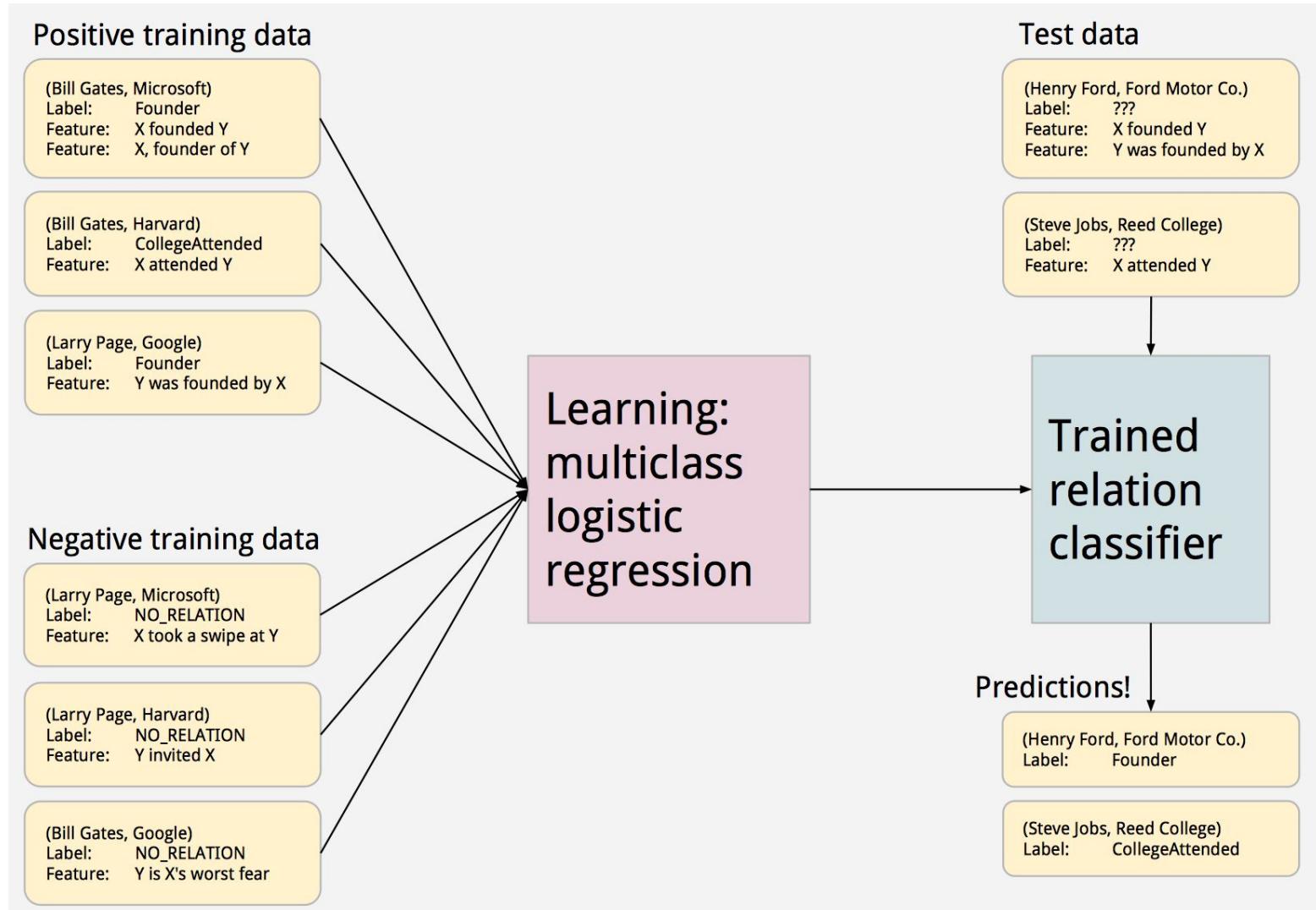
(Larry Page, Microsoft)  
Label: NO\_RELATION  
Feature: X took a swipe at Y

(Larry Page, Harvard)  
Label: NO\_RELATION  
Feature: Y invited X

(Bill Gates, Google)  
Label: NO\_RELATION  
Feature: Y is X's worst fear



# Training



# Distant Supervision

- Has advantages of supervised approach
  - leverage rich, reliable hand-created knowledge
  - relations have canonical names
  - can use rich features (e.g. syntactic features)
- Has advantages of unsupervised approach
  - leverage unlimited amounts of text data
  - allows for very large number of weak features
  - not sensitive to training corpus: genre-independent

# Distant Supervision

- 1.8 million relation instances used for training
  - Compared to 17,000 relation instances in ACE
- 800,000 Wikipedia articles used for training,  
400,000 different articles used for testing
- Only extract relation instances not already in  
Freebase

# Success!

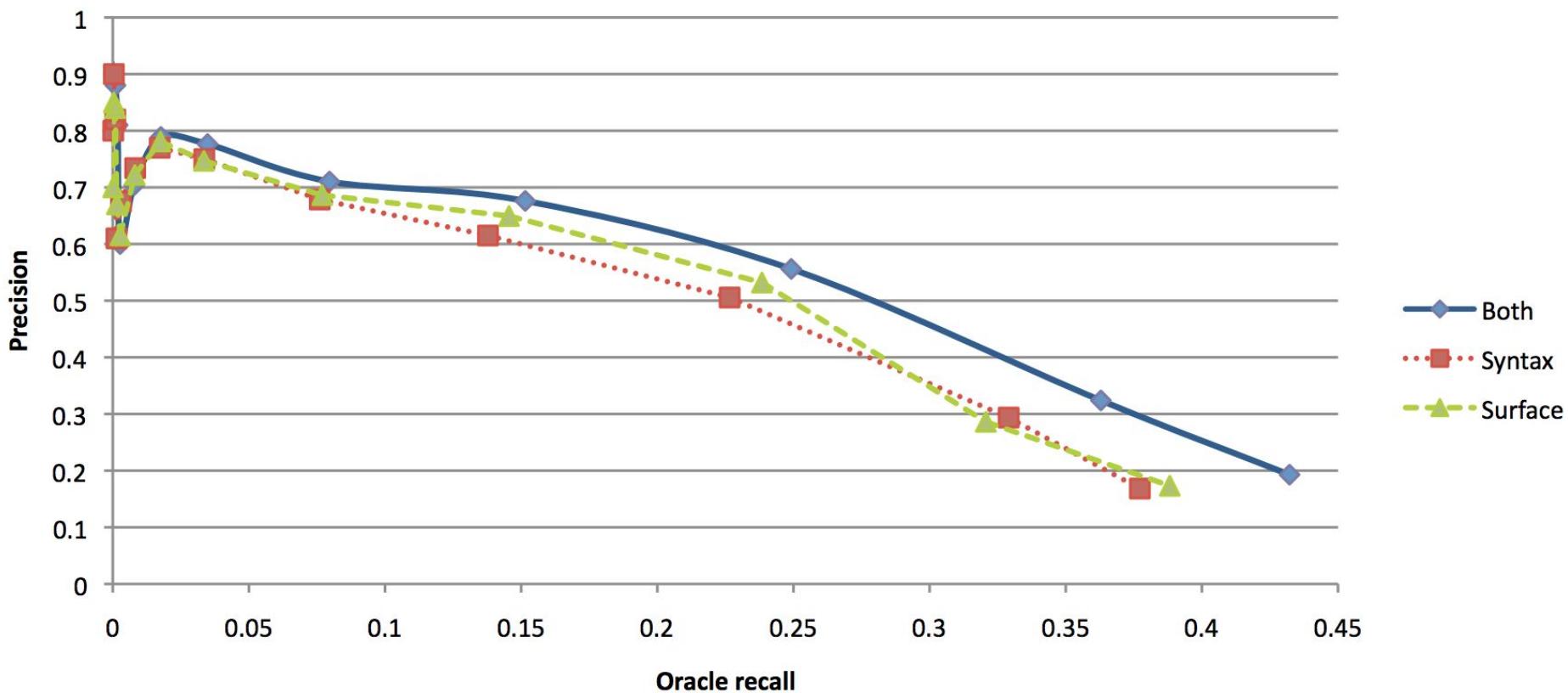
Ten relation instances extracted by the system that weren't in Freebase

Relation name	New instance
/location/location/contains	Paris, Montmartre
/location/location/contains	Ontario, Fort Erie
/music/artist/origin	Mighty Wagon, Cincinnati
/people/deceased_person/place_of_death	Fyodor Kamensky, Clearwater
/people/person/nationality	Marianne Yvonne Heemskerk, Netherlands
/people/person/place_of_birth	Wavell Wayne Hinds, Kingston
/book/author/works_written	Upton Sinclair, Lanny Budd
/business/company/founders	WWE, Vince McMahon
/people/person/profession	Thomas Mellon, judge

# Evaluation

- Held-out evaluation
  - Train on 50% of gold-standard Freebase relation instances, test on other 50%
  - Used to tune parameters quickly without having to wait for human evaluation
- Human evaluation
  - Performed by evaluators on Amazon Mechanical Turk
  - Calculated precision at 100 and 1000 recall levels for the ten most common relations

# Some Results (Snow et al. '09)



# Distant Supervision

- The distant supervision approach uses a database of known relation instances as a source of supervision
- We're classifying pairs of entities, not pairs of entity mentions
- The features for a pair of entities describe the patterns in which the two entities have co-occurred across many sentences in a large corpus
- Can make use of 100x or even 1000x more data than in the supervised paradigm

Wrapping Up

---

# What else is out there?

Important related work:

- Work on *open information extraction* (OpenIE) attempts to extract *all* relations from some text, without any fixed set of relations.

(OpenIE, extracts\_all\_relations\_from, text)

- Work on *knowledge base completion* attempts to use information in a knowledge base to fill in missing entries.

(AB, country\_of\_birth, Iceland)

=> (AB, speaks\_language, Icelandic)

---

---

# Coming Up

Tomorrow at NLP/Text as Data (4p, 7th Fl.):

- Justine Zhang (Cornell)  
*Unsupervised Models of Conversational Dynamics*

Next week:

- *Spring break*

March 21:

- HW3 due
- Guest lecture by Paloma  
*Foundations: Formal approaches to sentence meaning*



# THE END

