



THE UNIVERSITY OF CHICAGO
GRAHAM SCHOOL
CONTINUING LIBERAL AND PROFESSIONAL STUDIES



Natural Language Processing

Session 5

Nick Kadochnikov



Session 5 Agenda

- Information extraction
- Named entity recognition
- Relation extraction
 - Automatic content extraction annotation guidelines for entities
- Natural language parsing
- Dependency parsing





Information Extraction and Named Entity Recognition

Introducing the tasks:
Getting simple structured
information out of text



Information Extraction

- Information extraction (IE) systems
 - Find and understand limited relevant parts of texts
 - Gather information from many pieces of text
 - Produce a structured representation of relevant information:
 - *relations* (in the database sense), a.k.a.,
 - *a knowledge base*
 - Goals:
 1. Organize information so that it is useful to people
 2. Put information in a semantically precise form that allows further inferences to be made by computer algorithms



Information Extraction (IE)

- IE systems extract clear, factual information
 - Roughly: *Who did what to whom when?*
- E.g.,
 - Gathering earnings, profits, board members, headquarters, etc. from company reports
 - The headquarters of BHP Billiton Limited, and the global headquarters of the combined BHP Billiton Group, are located in Melbourne, Australia.
 - *headquarters("BHP Biliton Limited", "Melbourne, Australia")*
 - Learn drug-gene product interactions from medical research literature



Low-level information extraction

- Is now available – and I think popular – in applications like Apple or Google mail, and web indexing

The Los Altos Robotics Board of Directors is having a potluck dinner Friday January 6, 2012 and the upcoming [Botball](#) and FRC ([MVHS](#) [Eagle Strike Robotics](#)) seasons. You are back and it was a

Create New iCal Event...
Show This Date in iCal...
Copy

- Often seems to be based on regular expressions and name lists



Low-level information extraction



bhp billiton headquarters

Search

About 123,000 results (0.23 seconds)

Everything

Best guess for BHP Billiton Ltd. Headquarters is **Melbourne, London**

Images

Mentioned on at least 9 websites including wikipedia.org, bhpbilliton.com and bhpbilliton.com - [Feedback](#)

Maps

[BHP Billiton - Wikipedia, the free encyclopedia](http://en.wikipedia.org/wiki/BHP_Billiton)

Videos

en.wikipedia.org/wiki/BHP_Billiton

News

Merger of BHP & Billiton 2001 (creation of a DLC). **Headquarters, Melbourne, Australia** (BHP Billiton Limited and BHP Billiton Group) **London, United Kingdom ...**

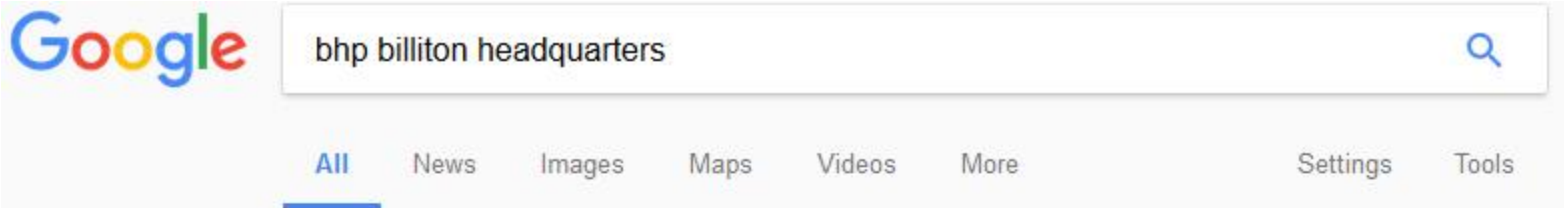
Shopping

[History](#) - [Corporate affairs](#) - [Operations](#) - [Accidents](#)





As Google matured from rules to ML / AI



About 232,000 results (0.96 seconds)

BHP Billiton Ltd. / Headquarters

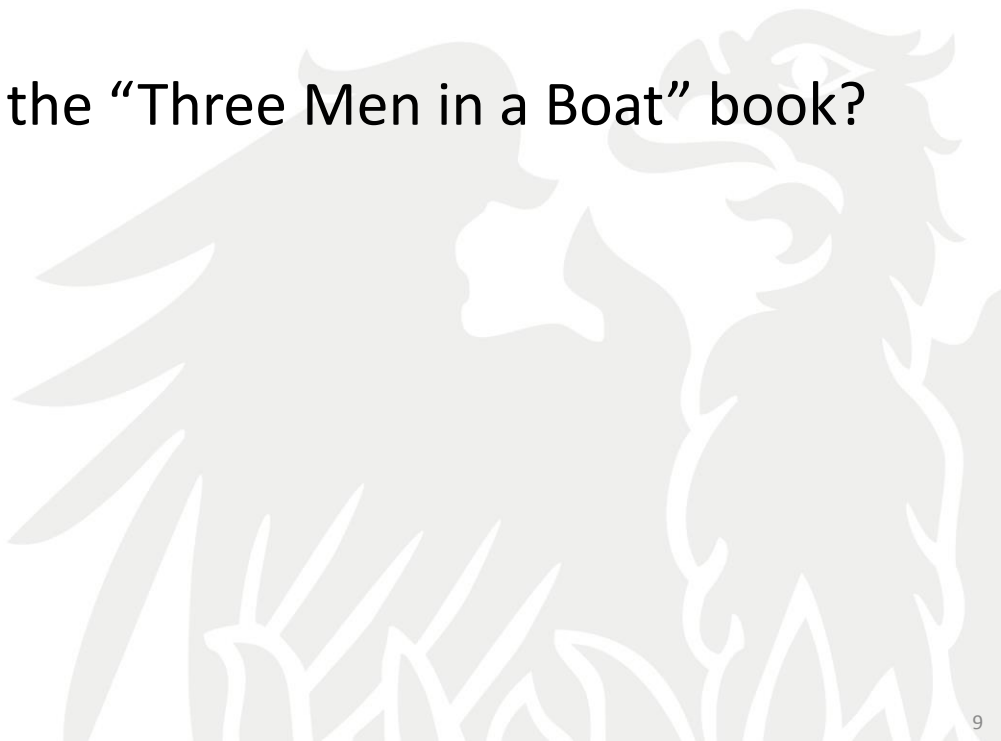


Melbourne, Australia



Exercise

- Who are the main characters in the “Three Men in a Boat” book?





Named Entity Recognition (NER)

- A very important sub-task: **find** and **classify** names in text, for example:
 - The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.



Named Entity Recognition (NER)

- A very important sub-task: **find** and **classify** names in text, for example:
 - The decision by the independent MP **Andrew Wilkie** to withdraw his support for the minority **Labor** government sounded dramatic but it should not further threaten its stability. When, after the **2010** election, **Wilkie, Rob Oakeshott, Tony Windsor** and the **Greens** agreed to support **Labor**, they gave just two guarantees: confidence and supply.



Named Entity Recognition (NER)

- A very important sub-task: **find** and **classify** names in text, for example:

- The decision by the independent MP **Andrew Wilkie** to withdraw his support for the minority **Labor** government sounded dramatic but it should not further threaten its stability. When, after the **2010** election, **Wilkie**, **Rob Oakeshott**, **Tony Windsor** and the **Greens** agreed to support **Labor**, they gave just two guarantees: confidence and supply.

Person
Date
Location
**Organi-
zation**



Named Entity Recognition (NER)

- The uses:
 - Named entities can be indexed, linked off, etc.
 - Sentiment can be attributed to companies or products
 - A lot of IE relations are associations between named entities
 - For question answering, answers are often named entities.
- Concretely:
 - Many web pages tag various entities, with links to bio or topic pages, etc.
 - Reuters' OpenCalais, Evri, AlchemyAPI, Yahoo's Term Extraction, ...
 - Apple/Google/Microsoft/... smart recognizers for document content



The Named Entity Recognition Task

Task: Predict entities in a text

Foreign **ORG**

Ministry **ORG**

spokesman **O**

Shen **PER**

Guofang **PER**

told **O**

Reuters **ORG**

:



Standard
evaluation
is per entity,
not per token



Precision/Recall/F1 for IE/NER

- Recall and precision are straightforward for tasks like IR and text categorization, where there is only one grain size (documents)
- The measure behaves a bit funnily for IE/NER when there are *boundary errors* (which are *common*):
 - First Bank of Chicago announced earnings ...
- This counts as both a fp and a fn
- Selecting *nothing* would have been better
- Some other metrics (e.g., MUC scorer) give partial credit (according to complex rules)



The ML sequence model approach to NER

Training

1. Collect a set of representative training documents
2. Label each token for its entity class or other (O)
3. Design feature extractors appropriate to the text and classes
4. Train a sequence classifier to predict the labels from the data

Testing

1. Receive a set of testing documents
2. Run sequence model inference to label each token
3. Appropriately output the recognized entities



Encoding classes for sequence labeling

IO encoding

IOB encoding

Fred

PER

B-PER

showed

O

O

Sue

PER

B-PER

Mengqiu

PER

B-PER

Huang

PER

I-PER

's

O

O

new

O

O

painting

O

O

B-PER indicates the beginning of a person name, *I-PER* indicates inside a person name, and so forth



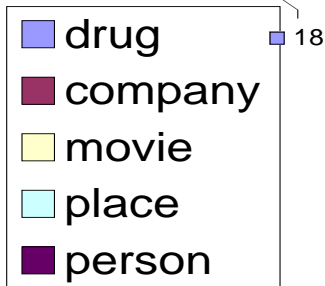
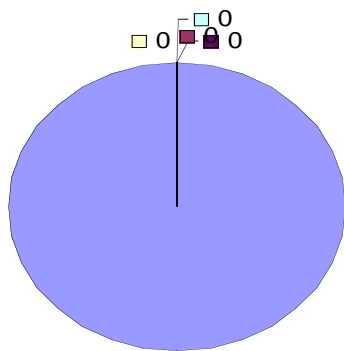
Features for sequence labeling

- Words
 - Current word (essentially like a learned dictionary)
 - Previous/next word (context)
- Other kinds of inferred linguistic classification
 - Part-of-speech tags
- Label context
 - Previous (and perhaps next) label

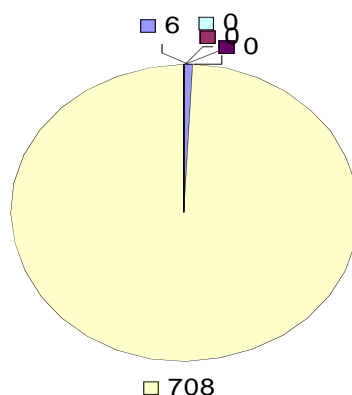


Features: Word substrings

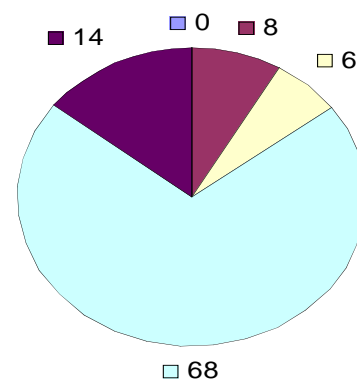
oxa



:



field



Cotrimoxazole

Wethersfield

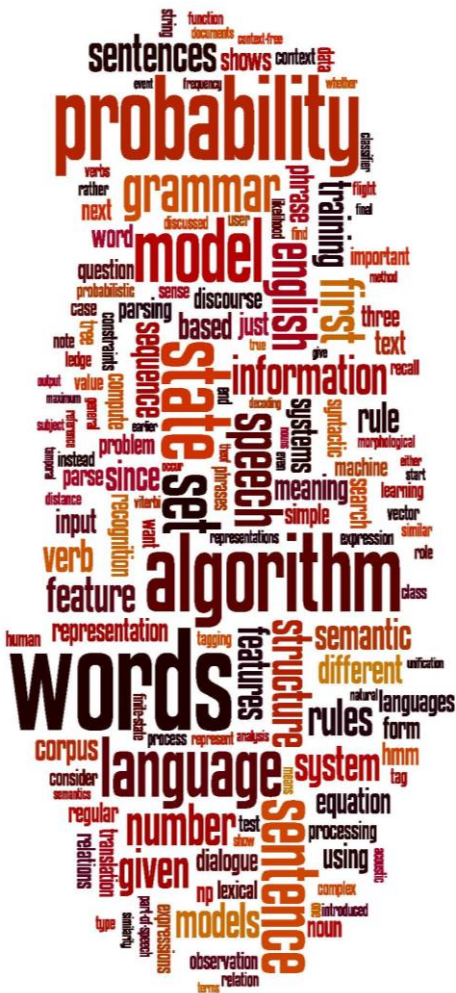
Alien Fury: Countdown to Invasion



Features: Word shapes

- Word Shapes
 - Map words to simplified representation that encodes attributes such as length, capitalization, numerals, Greek letters, internal punctuation, etc.

Varicella-zoster	Xx-xxx
mRNA	xXXX
CPA1	XXXd



Relation Extraction

What is relation extraction?



Extracting relations from text

- Company report: “International Business Machines Corporation (IBM or the company) was incorporated in the State of New York on June 16, 1911, as the Computing-Tabulating-Recording Co. (C-T-R)...”

- Extracted Complex Relation:

Company-Founding

Company	IBM
Location	New York
Date	June 16, 1911
Original-Name	Computing-Tabulating-Recording Co.

- But we will focus on the simpler task of extracting relation **triples**

Founding-year(IBM,1911)

Founding-location(IBM,New York)



Extracting Relation Triples from Text



Article Talk Read Edit View history Search

Stanford University

From Wikipedia, the free encyclopedia

Coordinates: 37°43′N 122°17′W﻿ / ﻿37.43°N 122.17°W﻿ / 37.43; -122.17

"Stanford" redirects here. For other uses, see [Stanford \(disambiguation\)](#).

Not to be confused with [Stanford University \(disambiguation\)](#).

The **Leland Stanford Junior University**, commonly referred to as **Stanford University** or **Stanford**, is an American [private research university](#) located in [Stanford, California](#) on an 8,180-acre (3,310 ha) campus near [Palo Alto, California, United States](#). It is situated in the northwestern [Santa Clara Valley](#) on the [San Francisco Peninsula](#), approximately 20 miles (32 km) northwest of [San Jose](#) and 37 miles (60 km) southeast of [San Francisco](#).^[6]

[Leland Stanford](#), a Californian railroad tycoon and politician, founded the university in 1891 in honor of his son, [Leland Stanford, Jr.](#), who died of [typhoid](#) two months before his 16th birthday. The university was established as a coeducational and nondenominational institution, but struggled financially after the senior Stanford's 1893 death and after much of the campus was damaged by the 1906 [San Francisco earthquake](#). Following [World War II](#), Provost [Frederick Terman](#) supported faculty and graduates' entrepreneurialism to build a self-sufficient local industry in what would become known as [Silicon Valley](#). By 1970, Stanford was home to a [linear accelerator](#), was one of the original four [ARPANET](#) nodes, and had transformed itself into a major research university in [computer science](#), [mathematics](#), [natural sciences](#), and [social sciences](#). More than 50 Stanford faculty, staff, and alumni have won the [Nobel Prize](#) and Stanford has the largest number of [Turing award](#) winners for a single institution. Stanford faculty and alumni have founded many prominent technology companies including [Cisco Systems](#), [Google](#), [Hewlett-Packard](#), [LinkedIn](#), [Rambus](#), [Silicon Graphics](#), [Sun Microsystems](#), [Varian Associates](#), and [Yahoo!](#)^[7]

The university is organized into seven schools including academic schools of [Humanities](#), [Life Sciences](#), [Physical Sciences](#), [Engineering](#), [Business](#), [Education](#), and [Law](#).

Stanford University
Leland Stanford Junior University

Seal of Stanford University

Motto *Die Luft der Freiheit weht* (German)^[1]
The wind of free- dom blows!^[1]

Motto in English

Leland Stanford Junior University,
founded in 1891, is an American
private research university located in
Stanford, California, near Palo Alto,
California.



Stanford is a research university
located in California
near Palo Alto,
founded in 1891
by Leland Stanford.



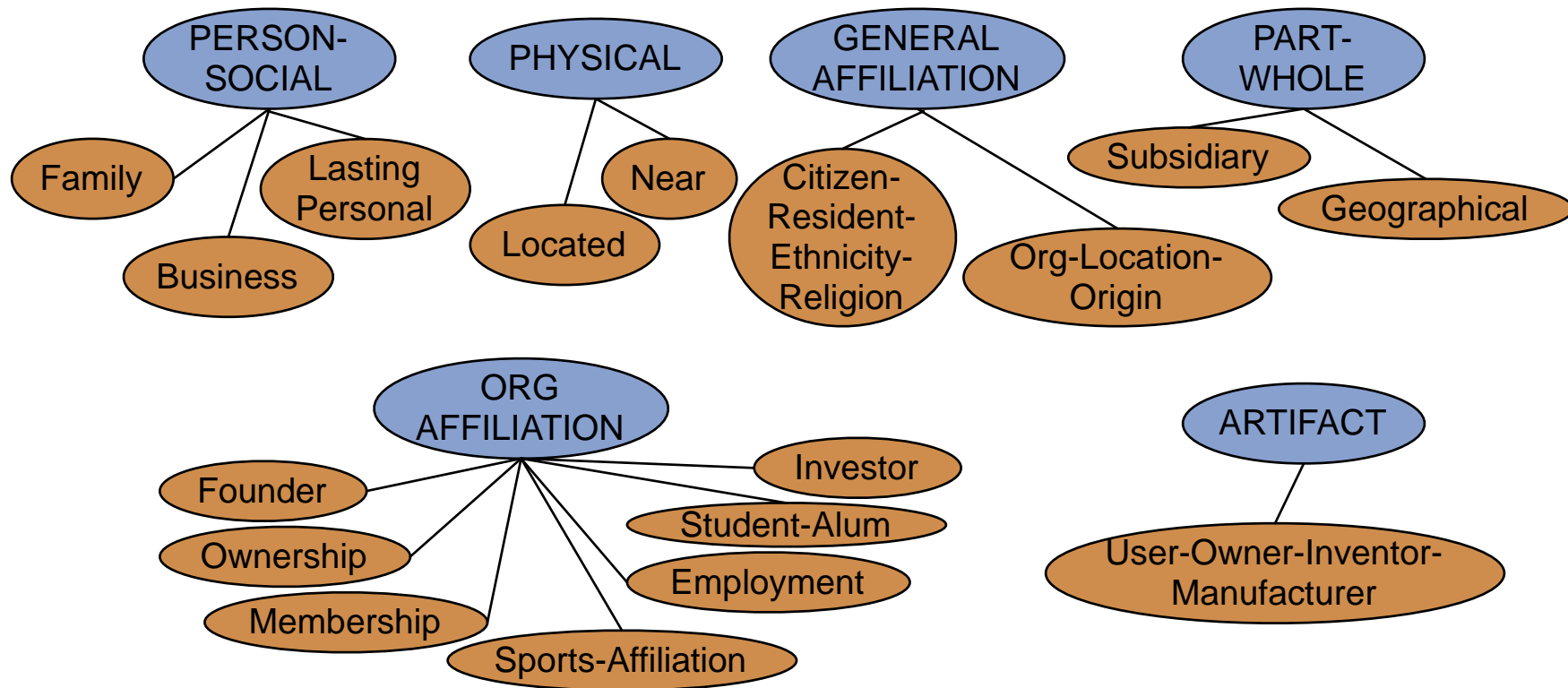
Why Relation Extraction?

- Create new structured knowledge bases, useful for any app
- Augment current knowledge bases
 - Adding words to WordNet thesaurus, facts to FreeBase or DBPedia
- Support question answering
 - The granddaughter of which actor starred in the movie “E.T.”?
 (acted-in ?x “E.T.”) (is-a ?y actor) (granddaughter-of ?x ?y)
- But which relations should we extract?



Automated Content Extraction (ACE)

17 relations from 2008 “Relation Extraction Task”





Automated Content Extraction (ACE)

- Physical-Located **PER-GPE**

He was in Tennessee

- Part-Whole-Subsidiary **ORG-ORG**

XYZ, the parent company of ABC

- Person-Social-Family **PER-PER**

John's wife Yoko

- Org-AFF-Founder **PER-ORG**

Steve Jobs, co-founder of Apple...

Persons (PER)
Geographical (GPE)
Organizations (ORG)



UMLS: Unified Medical Language System

- 134 entity types, 54 relations

Injury	disrupts	Physiological Function
Bodily Location	location-of	Biologic Function
Anatomical Structure	part-of	Organism
Pharmacologic Substance	causes	Pathological Function
Pharmacologic Substance	treats	Pathologic Function

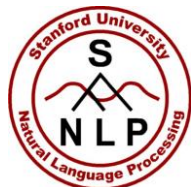


Extracting UMLS relations from a sentence

Doppler echocardiography can be used to
diagnose left anterior descending artery
stenosis in patients with type 2 diabetes



Echocardiography, Doppler **DIAGNOSES** Acquired stenosis



Databases of Wikipedia Relations

Wikipedia Infobox

Relations extracted from Infobox

Stanford [state](#) California

Stanford [motto](#) “Die Luft der Freiheit weht”

{{Infobox university

|image_name= Stanford University seal.svg

|image_size= 210px

|caption = Seal of Stanford University

|name =Stanford University

|native_name =Leland Stanford Junior Uni

|motto = {{lang|de|"Die Luft der Freiheit v

name="casper">{{cite speech|title=Die Lu

Casper|first=Gerhard|last=Casper|author

05|url=http://www.stanford.edu/dept/pr

|mottoeng = The wind of freedom blows<

|established = 1891<ref>{{cite web |

url=http://www.stanford.edu/home/stan

publisher = Stanford University | accessd:

|type = [[private university|Private]]

|calendar= Quarter

|president = [[John L. Hennessy]]

|provost = [[John Etchemendy]]

|city = [[Stanford, California|Stanford]]

|state = California

|country = U.S.

Type

[Private](#)

Endowment

US\$ 16.5 [billion](#) (2011)^[3]

President

[John L. Hennessy](#)

Provost

[John Etchemendy](#)

Academic staff

1,910^[4]

Students

15,319

Undergraduates

6,878^[5]

Postgraduates

8,441^[5]

Location

[Stanford](#), California, U.S.

Campus

[Suburban](#), 8,180 acres (3,310 ha)^[6]

Colors

Cardinal red and white



1

tml}}</ref>

ty History |



Relation databases that draw from Wikipedia

- Resource Description Framework (RDF) triples
subject predicate object
Golden Gate Park `location` San Francisco
`dbpedia:Golden_Gate_Park` `dbpedia-owl:location` `dbpedia:San_Francisco`
- DBPedia: 1 billion RDF triples, 385 from English Wikipedia
- Frequent Freebase relations:

people/person/nationality,	location/location/contains
people/person/profession,	people/person/place-of-birth
biology/organism_higher_classification	film/film/genre



Ontological relations

Examples from the WordNet Thesaurus

- IS-A (hypernym): subsumption between classes
 - Giraffe IS-A ruminant IS-A ungulate IS-A mammal IS-A vertebrate IS-A animal...
- Instance-of: relation between individual and class
 - San Francisco instance-of city



How to build relation extractors

1. Hand-written patterns
2. Supervised machine learning
3. Semi-supervised and unsupervised
 - Bootstrapping (using seeds)
 - Distant supervision
 - Unsupervised learning from the web



Relation Extraction

Using patterns to extract relations



Rules for extracting IS-A relation

Early 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?



Rules for extracting IS-A relation

Early 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?



Hearst's Patterns for extracting IS-A relations

(Hearst, 1992): Automatic Acquisition of Hyponyms

"Y such as X ((, X) * (, and|or) X) "

"such Y as X"

"X or other Y"

"X and other Y"

"Y including X"

"Y, especially X"



Hearst's Patterns for extracting IS-A relations

Hearst pattern	Example occurrences
X and other Y	...temples, treasures, and other important civic buildings.
X or other Y	Bruises, wounds, broken bones or other injuries...
Y such as X	The bow lute, such as the Bambara ndang...
Such Y as X	... such authors as Herrick, Goldsmith, and Shakespeare.
Y including X	...common-law countries, including Canada and England...
Y , especially X	European countries, especially France, England, and Spain...



Extracting Richer Relations Using Rules

- Intuition: relations often hold between specific entities
 - **located-in** (ORGANIZATION, LOCATION)
 - **founded** (PERSON, ORGANIZATION)
 - **cures** (DRUG, DISEASE)
- Start with Named Entity tags to help extract relation!



Named Entities aren't quite enough. Which relations hold between 2 entities?



Drug

Cure?
Prevent?
Cause?



Disease



What relations hold between 2 entities?



PERSON

Founder?

Investor?

Member?

Employee?

President?



ORGANIZATION



Extracting Richer Relations Using Rules and Named Entities

Who holds what office in what organization?

PERSON, POSITION of ORG

- George Marshall, Secretary of State of the United States

PERSON (named | appointed | chose | *etc.*) PERSON Prep? POSITION

- Truman appointed Marshall Secretary of State

PERSON [be]? (named | appointed | *etc.*) Prep? ORG POSITION

- George Marshall was named US Secretary of State



Hand-built patterns for relations

- Plus:
 - Human patterns tend to be high-precision
 - Can be tailored to specific domains
- Minus
 - Human patterns are often low-recall
 - A lot of work to think of all possible patterns!
 - Don't want to have to do this for every relation!
 - We'd like better accuracy



Relation Extraction

Supervised relation extraction



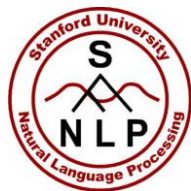
Supervised machine learning for relations

- Choose a set of relations we'd like to extract
- Choose a set of relevant named entities
- Find and label data
 - Choose a representative corpus
 - Label the named entities in the corpus
 - Hand-label the relations between these entities
 - Break into training, development, and test
- Train a classifier on the training set



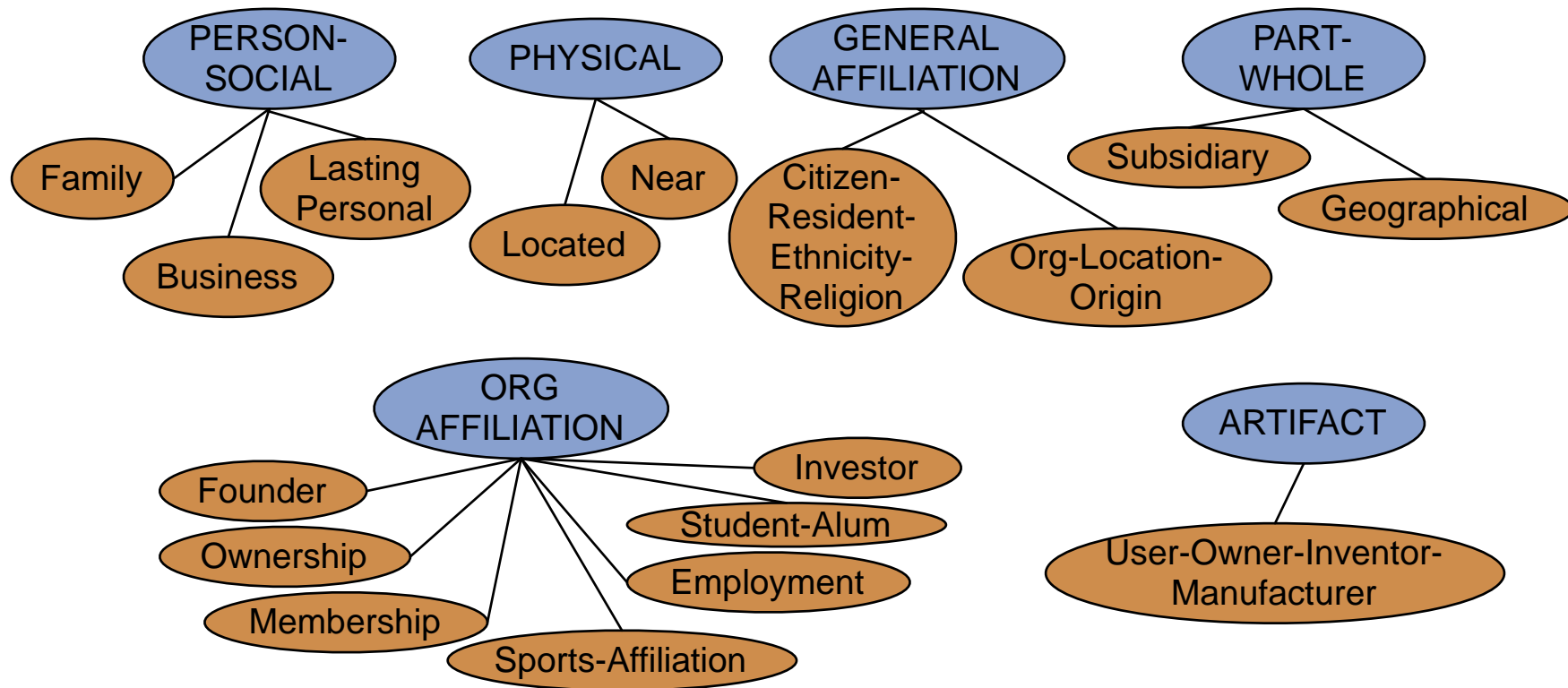
How to do classification in supervised relation extraction

1. Find all pairs of named entities (usually in same sentence)
2. Decide if 2 entities are related
3. If yes, classify the relation
 - Why the extra step?
 - Faster classification training by eliminating most pairs
 - Can use distinct feature-sets appropriate for each task.



Automated Content Extraction (ACE)

17 sub-relations of 6 relations from 2008 "Relation Extraction Task"





Relation Extraction

Classify the relation between two entities in a sentence

***American Airlines**, a unit of AMR, immediately matched the move, spokesman **Tim Wagner** said.*

FAMILY

CITIZEN

SUBSIDIARY

FOUNDER



NIL

EMPLOYMENT

INVENTOR

...



Word Features for Relation Extraction

***American Airlines**, a unit of AMR, immediately matched the move, spokesman **Tim Wagner** said*

Mention 1

Mention 2

- Headwords of M1 and M2, and combination
 Airlines Wagner Airlines-Wagner
- Bag of words and bigrams in M1 and M2
 {American, Airlines, Tim, Wagner, American Airlines, Tim Wagner}
- Words or bigrams in particular positions left and right of M1/M2
 M2: -1 *spokesman*
 M2: +1 *said*
- Bag of words or bigrams between the two entities
 {a, AMR, of, immediately, matched, move, spokesman, the, unit}



Named Entity Type and Mention Level Features for Relation Extraction

American Airlines, a unit of AMR, immediately matched the move, spokesman *Tim Wagner* said

Mention 1 Mention 2

- Named-entity types
 - M1: **ORG**
 - M2: **PERSON**
- Concatenation of the two named-entity types
 - **ORG-PERSON**
- Entity Level of M1 and M2 (NAME, NOMINAL, PRONOUN)
 - M1: **NAME** [it or he would be **PRONOUN**]
 - M2: **NAME** [the company would be **NOMINAL**]



Parse Features for Relation Extraction

American Airlines, a unit of AMR, immediately matched the move, spokesman *Tim Wagner* said
 Mention 1 Mention 2

- Base syntactic chunk sequence from one to the other

NP NP PP VP NP NP

- Constituent path through the tree from one to the other

NP ↑ NP ↑ S ↑ S ↓ NP

- Dependency path

Airlines matched Wagner said



Gazeteer and trigger word features for relation extraction

- Trigger list for family: kinship terms
 - [parent](#), [wife](#), [husband](#), [grandparent](#), [etc.](#) [from WordNet]
- Gazetteer:
 - Lists of useful geo or geopolitical words
 - Country name list
 - Other sub-entities



***American Airlines**, a unit of AMR, immediately matched the move, spokesman **Tim Wagner** said.*

Entity-based features

Entity ₁ type	ORG
Entity ₁ head	<i>airlines</i>
Entity ₂ type	PERS
Entity ₂ head	<i>Wagner</i>
Concatenated types	ORGPERS

Word-based features

Between-entity bag of words	{ <i>a, unit, of, AMR, Inc., immediately, matched, the, move, spokesman</i> }
Word(s) before Entity ₁	NONE
Word(s) after Entity ₂	<i>said</i>

Syntactic features

Constituent path	$NP \uparrow NP \uparrow S \uparrow S \downarrow NP$
Base syntactic chunk path	$NP \rightarrow NP \rightarrow PP \rightarrow NP \rightarrow VP \rightarrow NP \rightarrow NP$
Typed-dependency path	$Airlines \leftarrow_{subj} matched \leftarrow_{comp} said \rightarrow_{subj} Wagner$



Classifiers for supervised methods

- Now you can use any classifier you like
 - MaxEnt
 - Naïve Bayes
 - SVM
 - ...
- Train it on the training set, tune on the dev set, test on the test set



Evaluation of Supervised Relation Extraction

- Compute P/R/ F_1 for each relation

$$P = \frac{\text{\# of correctly extracted relations}}{\text{Total \# of extracted relations}}$$

$$R = \frac{\text{\# of correctly extracted relations}}{\text{Total \# of gold relations}}$$

$$F_1 = \frac{2PR}{P + R}$$



Summary: Supervised Relation Extraction

- + Can get high accuracies with enough hand-labeled training data, if test similar enough to training
- Labeling a large training set is expensive
- Supervised models are brittle, don't generalize well to different genres



Semi-supervised and unsupervised relation extraction



Seed-based or bootstrapping approaches to relation extraction

- No training set? Maybe you have:
 - A few seed tuples or
 - A few high-precision patterns
- Can you use those seeds to do something useful?
 - Bootstrapping: use the seeds to directly learn to populate a relation



Relation Bootstrapping (Hearst 1992)

- Gather a set of seed pairs that have relation R
- Iterate:
 1. Find sentences with these pairs
 2. Look at the context between or around the pair and generalize the context to create patterns
 3. Use the patterns for grep for more pairs



Bootstrapping

- <Mark Twain, Elmira> **Seed tuple**
 - Grep (google) for the environments of the seed tuple

“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 grep for new tuples
- Iterate



Dipre: Extract <author,book> pairs

Brin, Sergei. 1998. Extracting Patterns and Relations from the World Wide Web.

Start with 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

- Find Instances:

The Comedy of Errors, by William Shakespeare, was

The Comedy of Errors, by William Shakespeare, is

The Comedy of Errors, one of William Shakespeare's earliest attempts

The Comedy of Errors, one of William Shakespeare's most

- Extract patterns (group by middle, take longest common prefix/suffix)

?x , by ?y , ?x , one of ?y 's

- Now iterate, finding new seeds that match the pattern



Snowball

E. Agichtein and L. Gravano 2000. Snowball: Extracting Relations from Large Plain-Text Collections. ICDL

- Similar iterative algorithm

Organization	Location of Headquarters
Microsoft	Redmond
Exxon	Irving
IBM	Armonk

- Group instances w/similar prefix, middle, suffix, extract patterns
 - But require that X and Y be named entities
 - And compute a confidence for each pattern

.69

ORGANIZATION

{ 's, in, headquarters }

LOCATION

.75

LOCATION

{ in, based }

ORGANIZATION



Distant Supervision

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

Fei Wu and Daniel S. Weld. 2007. Autonomously Semantifying Wikipeida. CIKM 2007

Mintz, Bills, Snow, Jurafsky. 2009. Distant supervision for relation extraction without labeled data. ACL09

- Combine bootstrapping with supervised learning
 - Instead of 5 seeds,
 - Use a large database to get huge # of seed examples
 - Create lots of features from all these examples
 - Combine in a supervised classifier



Distant supervision paradigm

- Like supervised classification:
 - Uses a classifier with lots of features
 - Supervised by detailed hand-created knowledge
 - Doesn't require iteratively expanding patterns
- Like unsupervised classification:
 - Uses very large amounts of unlabeled data
 - Not sensitive to genre issues in training corpus



Distantly supervised learning of relation extraction patterns

- 1 For each relation
- 2 For each tuple in big database
- 3 Find sentences in large corpus with both entities
- 4 Extract frequent features (parse, words, etc)
- 5 Train supervised classifier using thousands of patterns

Born-In

<Edwin Hubble, Marshfield>
<Albert Einstein, Ulm>

Hubble was born in Marshfield
Einstein, born (1879), Ulm
Hubble's birthplace in Marshfield

PER was born in LOC
PER, born (XXXX), LOC
PER's birthplace in LOC

$P(\text{born-in} \mid f_1, f_2, f_3, \dots, f_{70000})$



Unsupervised relation extraction

M. Banko, M. Cararella, S. Soderland, M. Broadhead, and O. Etzioni.
2007. Open information extraction from the web. IJCAI

- Open Information Extraction:
 - extract relations from the web with no training data, no list of relations
- 1. Use parsed data to train a “trustworthy tuple” classifier
- 2. Single-pass extract all relations between NPs, keep if trustworthy
- 3. Assessor ranks relations based on text redundancy
 - (FCI, specializes in, software development)
 - (Tesla, invented, coil transformer)



Evaluation of Semi-supervised and Unsupervised Relation Extraction

- Since it extracts totally new relations from the web
 - There is no gold set of correct instances of relations!
 - Can't compute precision (don't know which ones are correct)
 - Can't compute recall (don't know which ones were missed)
- Instead, we can approximate precision (only)
 - Draw a random sample of relations from output, check precision manually

$$\hat{P} = \frac{\text{\# of correctly extracted relations in the sample}}{\text{Total \# of extracted relations in the sample}}$$

- Can also compute precision at different levels of recall.
 - Precision for top 1000 new relations, top 10,000 new relations, top 100,000
 - In each case taking a random sample of that set

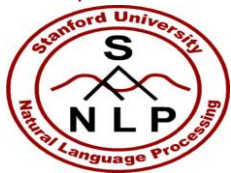


NER in Python



Natural Language Parsing

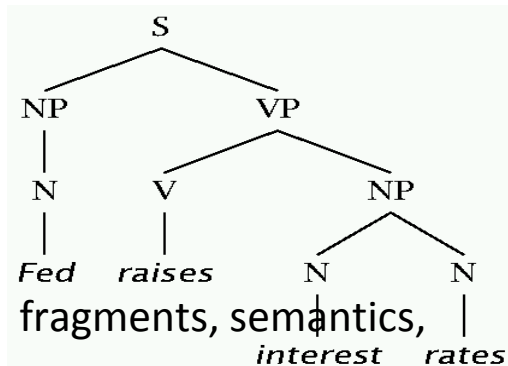
Two views of syntactic structure

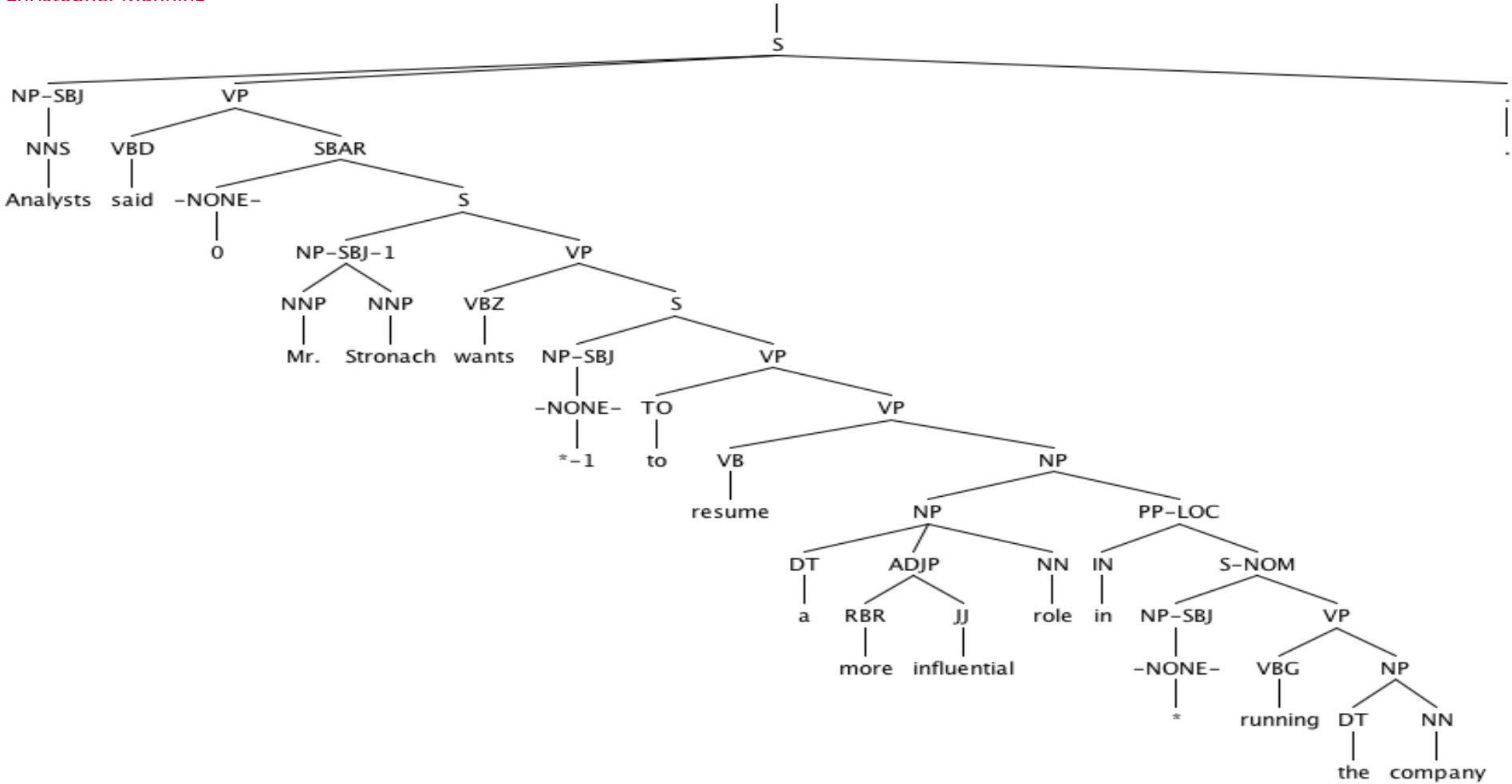


Two views of linguistic structure:

1. Constituency (phrase structure)

- Phrase structure organizes words into nested constituents.
- How do we know what is a **constituent**? (Not that linguists don't argue about some cases.)
 - Distribution: a constituent behaves as a unit that can appear in different places:
 - John talked [to the children] [about drugs].
 - John talked [about drugs] [to the children].
 - *John talked drugs to the children about
 - Substitution/expansion/pro-forms:
 - I sat [on the box/right on top of the box/there].
 - Coordination, regular internal structure, no intrusion, ...





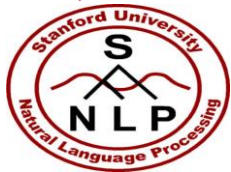


Headed phrase structure

- $VP \rightarrow \dots VB^* \dots$
- $NP \rightarrow \dots NN^* \dots$
- $ADJP \rightarrow \dots JJ^* \dots$
- $ADVP \rightarrow \dots RB^* \dots$

- $SBAR(Q) \rightarrow S | SINV | SQ \rightarrow \dots NP VP \dots$

- Plus minor phrase types:
 - QP (quantifier phrase in NP), CONJP (multi word constructions: *as well as*), INTJ (interjections), etc.



Two views of linguistic structure:

2. Dependency structure

- Dependency structure shows which words depend on (modify or are arguments of) which other words.

Statistical Natural Language Parsing

Parsing: The rise of data and statistics



Pre 1990 (“Classical”) NLP Parsing

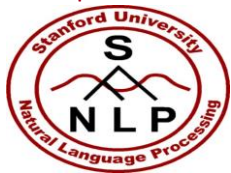
- Wrote symbolic grammar (CFG or often richer) and lexicon
 - $S \rightarrow NP VP$
 - $NP \rightarrow (DT) NN$
 - $NP \rightarrow NN NNS$
 - $NP \rightarrow NNP$
 - $VP \rightarrow V NP$
 - $NN \rightarrow interest$
 - $NNS \rightarrow rates$
 - $NNS \rightarrow raises$
 - $VBP \rightarrow interest$
 - $VBZ \rightarrow rates$
- Used grammar/proof systems to prove parses from words
- This scaled very badly and didn't give coverage. For sentence:
Fed raises interest rates 0.5% in effort to control inflation
 - Minimal grammar: 36 parses
 - Simple 10 rule grammar: 592 parses
 - Real-size broad-coverage grammar: millions of parses



Classical NLP Parsing:

The problem and its solution

- Categorical constraints can be added to grammars to limit unlikely/weird parses for sentences
 - But the attempt make the grammars not robust
 - In traditional systems, commonly 30% of sentences in even an edited text would have *no* parse.
- A less constrained grammar can parse more sentences
 - But simple sentences end up with ever more parses with no way to choose between them
- We need mechanisms that allow us to find the most likely parse(s) for a sentence
 - Statistical parsing lets us work with very loose grammars that admit millions of parses for sentences but still quickly find the best parse(s)



The rise of annotated data: The Penn Treebank

[Marcus et al. 1993, *Computational Linguistics*]

```
( (S
  (NP-SBJ (DT The) (NN move))
  (VP (VBD followed)
    (NP
      (NP (DT a) (NN round))
      (PP (IN of)
        (NP
          (NP (JJ similar) (NNS increases))
          (PP (IN by)
            (NP (JJ other) (NNS lenders))))
          (PP (IN against)
            (NP (NNP Arizona) (JJ real) (NN estate) (NNS loans))))))
    (, ,)
    (S-ADV
      (NP-SBJ (-NONE- *))
      (VP (VBG reflecting)
        (NP
          (NP (DT a) (VBG continuing) (NN decline))
          (PP-LOC (IN in)
            (NP (DT that) (NN market))))))
      (. .)))
```

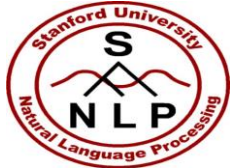


The rise of annotated data

- Starting off, building a treebank seems a lot slower and less useful than building a grammar
- But a treebank gives us many things
 - Reusability of the labor
 - Many parsers, POS taggers, etc.
 - Valuable resource for linguistics
 - Broad coverage
 - Frequencies and distributional information
 - A way to evaluate systems

Statistical Natural Language Parsing

An exponential number of attachments



Attachment ambiguities

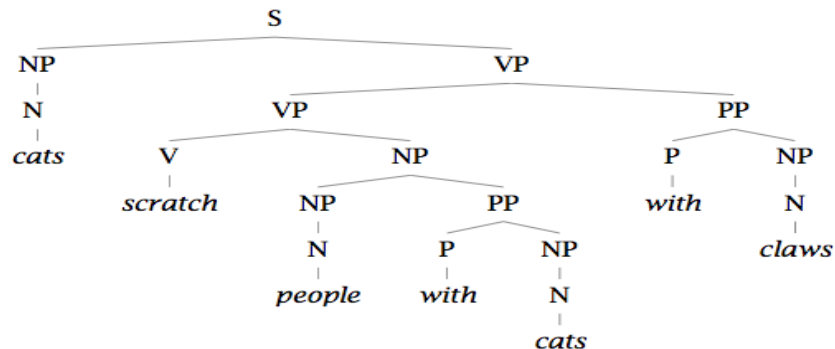
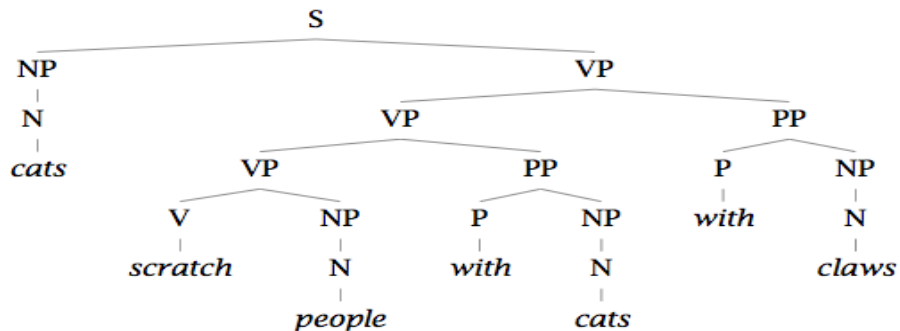
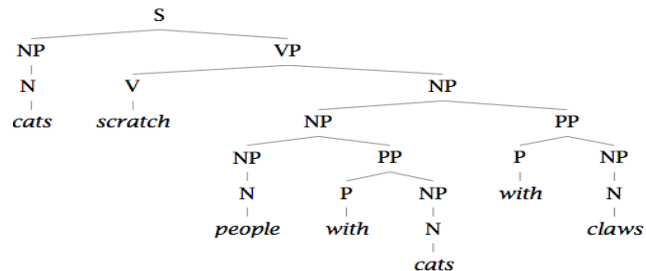
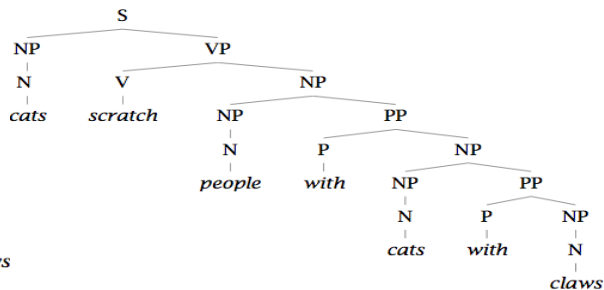
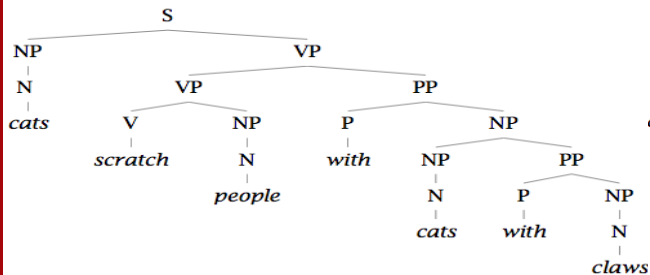
- A key parsing decision is how we ‘attach’ various constituents
 - PPs, adverbial or participial phrases, infinitives, coordinations, etc.

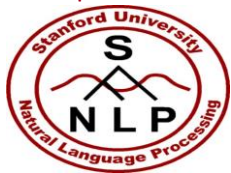
The board approved [its acquisition] [by Royal Trustco Ltd.]
[of Toronto]
[for \$27 a share]
[at its monthly meeting].

- Catalan numbers: $C_n = (2n)! / [(n+1)!n!]$
- An exponentially growing series, which arises in many tree-like contexts:
 - E.g., the number of possible triangulations of a polygon with $n+2$ sides
 - Turns up in triangulation of probabilistic graphical models....

Two problems to solve:

1. Repeated words...





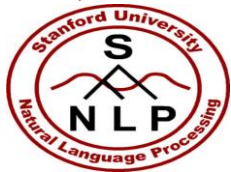
Two problems to solve:

2. Choosing the correct parse

- How do we work out the correct attachment:
 - She saw the man with a telescope
- Is the problem 'AI complete'? Yes, but ...
- Words are good predictors of attachment
 - Even absent full understanding
 - Moscow sent more than 100,000 soldiers into Afghanistan ...
 - Sydney Water breached an agreement with NSW Health ...
- Our statistical parsers will try to exploit such statistics.

[illegible]

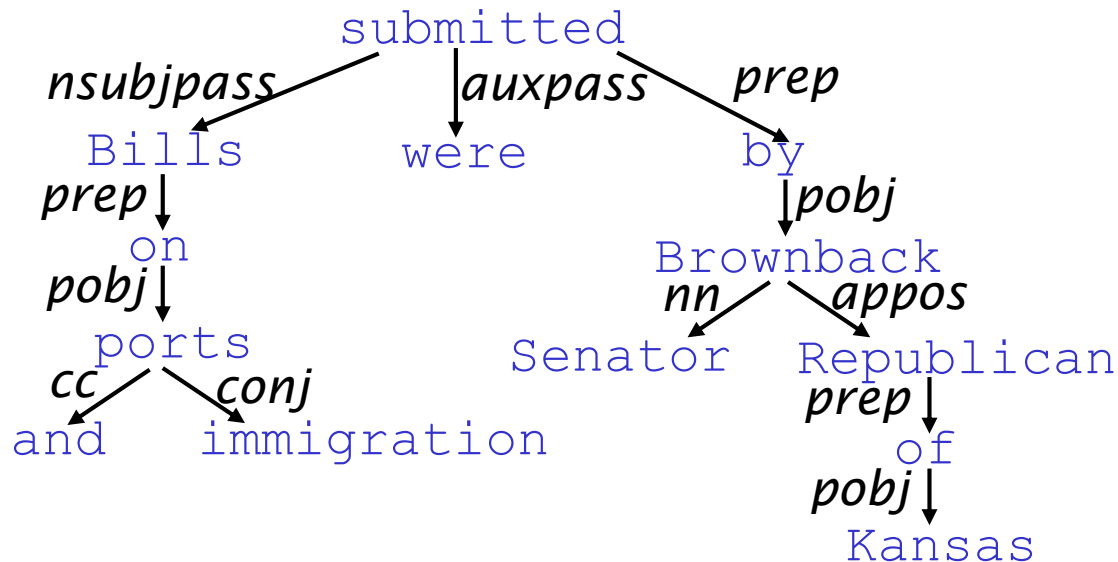
Introduction

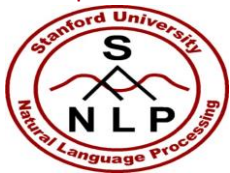


Dependency Grammar and Dependency Structure

Dependency syntax postulates that syntactic structure consists of lexical items linked by binary asymmetric relations (“arrows”) called dependencies

The arrows are commonly **typed** with the name of grammatical relations (subject, prepositional object, apposition, etc.)



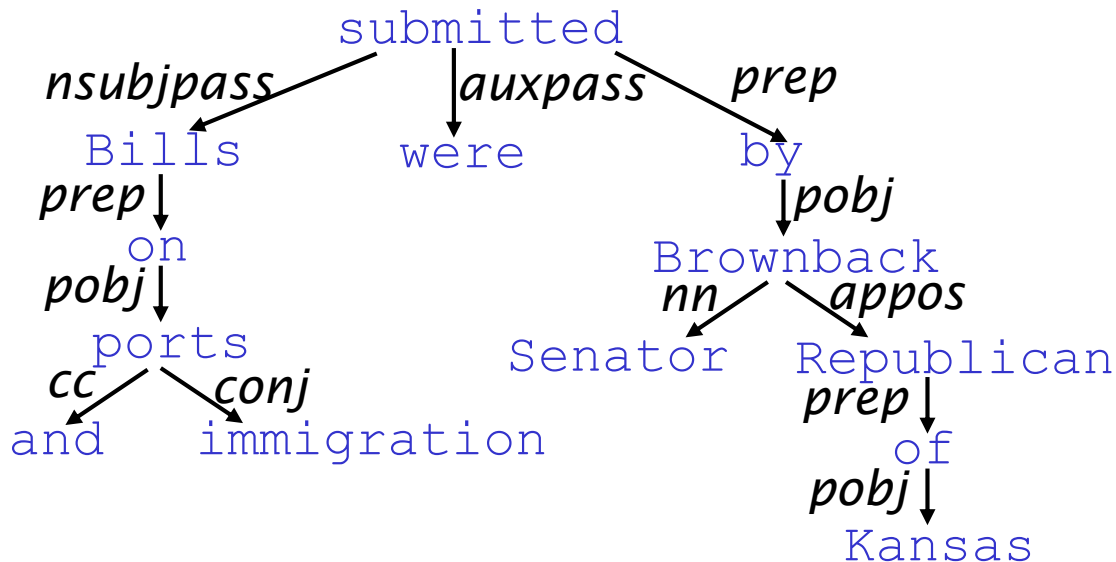


Dependency Grammar and Dependency Structure

Dependency syntax postulates that syntactic structure consists of lexical items linked by binary asymmetric relations (“arrows”) called dependencies

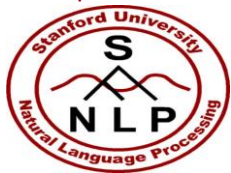
The arrow connects a **head** (governor, superior, regent) with a **dependent** (modifier, inferior, subordinate)

Usually, dependencies form a tree (connected, acyclic, single-head)



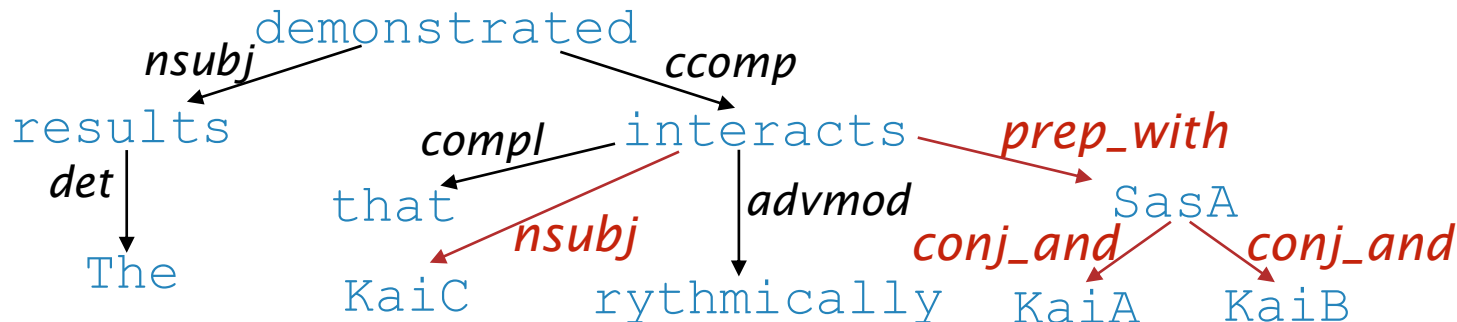
Dependencies encode relational structure

Relation Extraction with Stanford Dependencies



Dependency paths identify relations like protein interaction

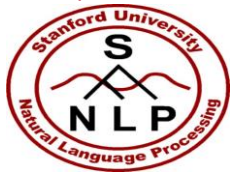
[Erkan et al. EMNLP 07, Fundel et al. 2007]



KaiC \leftarrow nsubj interacts prep_with \rightarrow SasA

KaiC \leftarrow nsubj interacts prep_with \rightarrow SasA conj_and \rightarrow KaiA

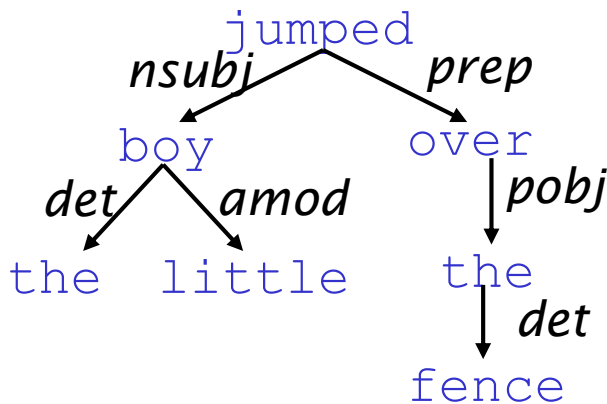
KaiC \leftarrow nsubj interacts prep_with \rightarrow SasA conj_and \rightarrow KaiB

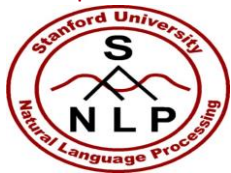


Stanford Dependencies

[de Marneffe et al. LREC 2006]

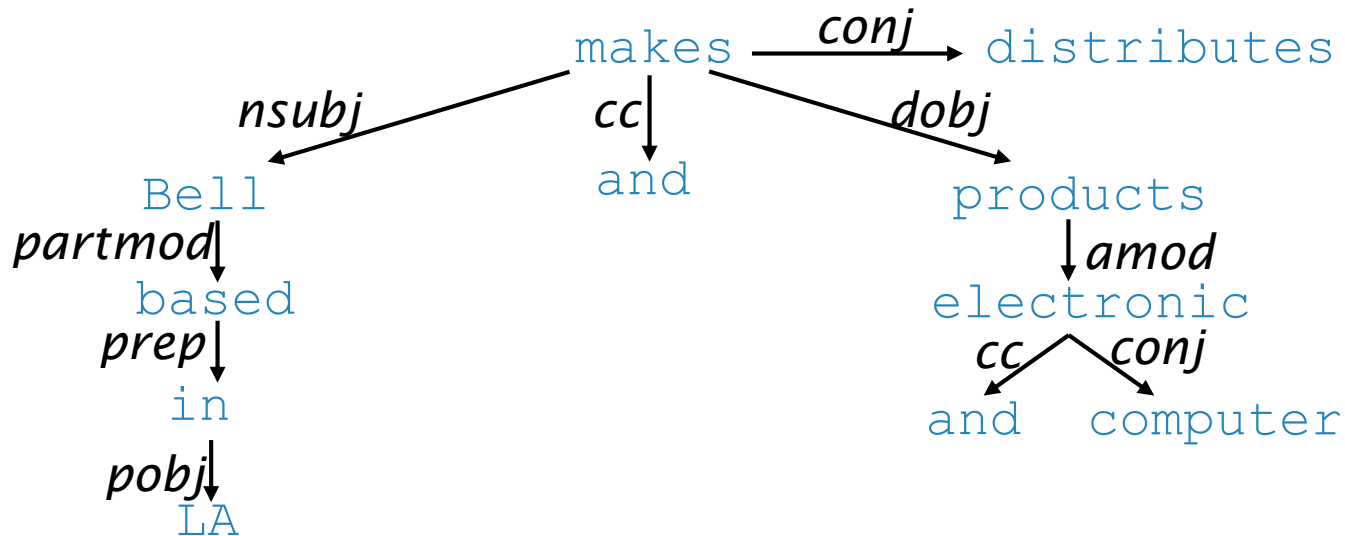
- The basic dependency representation is projective
- It can be generated by postprocessing headed phrase structure parses (Penn Treebank syntax)
- It can also be generated directly by dependency parsers, such as MaltParser, or the Easy-First Parser

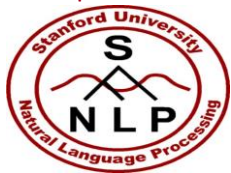




Graph modification to facilitate semantic analysis

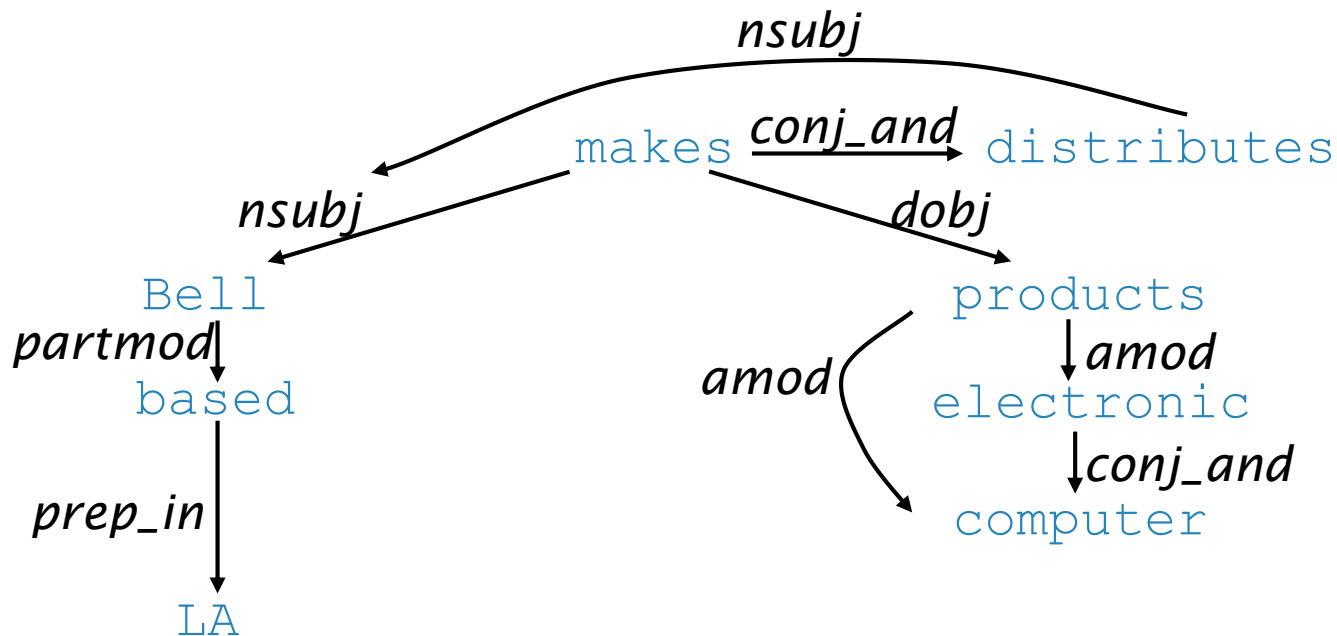
Bell, based in LA, makes and distributes electronic and computer products.





Graph modification to facilitate semantic analysis

Bell, based in LA, makes and distributes electronic and computer products.





THE UNIVERSITY OF CHICAGO
GRAHAM SCHOOL
CONTINUING LIBERAL AND PROFESSIONAL STUDIES

Dependency Parsing in Python





THE UNIVERSITY OF CHICAGO
GRAHAM SCHOOL
CONTINUING LIBERAL AND PROFESSIONAL STUDIES

Watson NLU Demo





THE UNIVERSITY OF CHICAGO
GRAHAM SCHOOL
CONTINUING LIBERAL AND PROFESSIONAL STUDIES

Thank You!

