

Natural Language Processing *Session 5*

Nick Kadochnikov

University of Chicago Professional Education



Session 5 Agenda

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



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 of these dinners three years 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

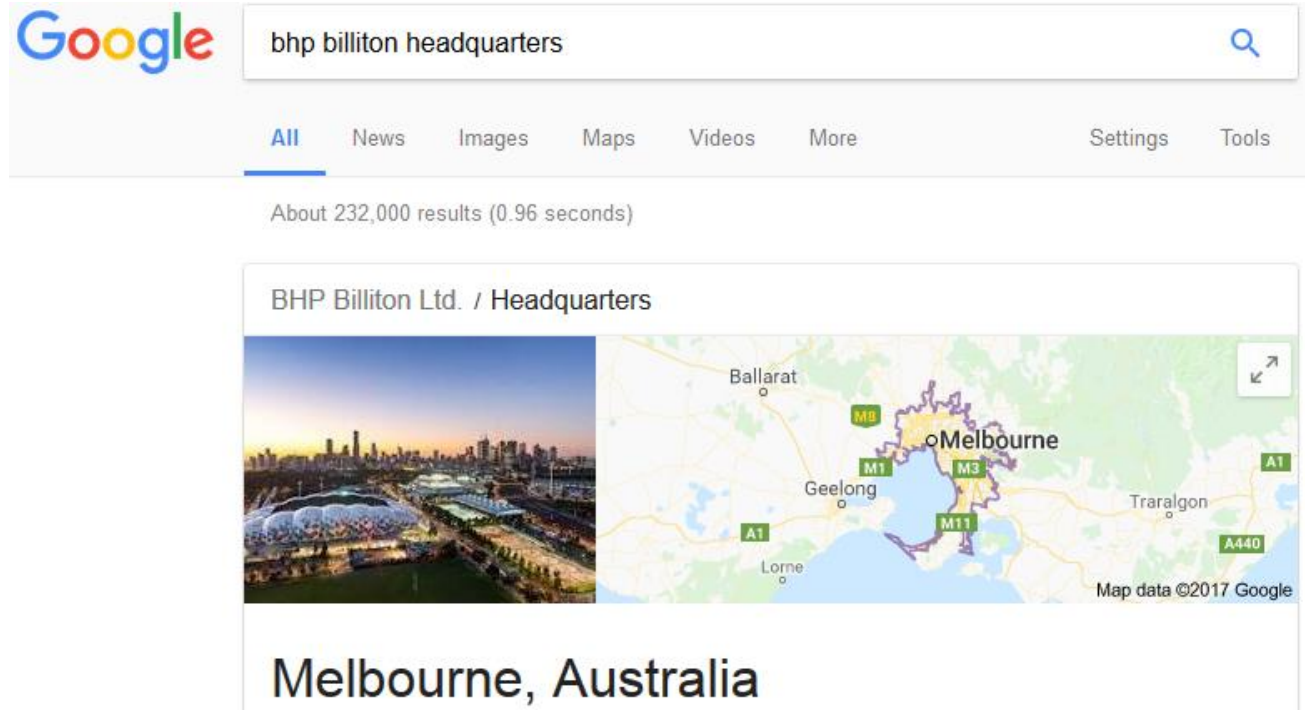
Merger of BHP & Billiton 2001 (creation of a DLC). **Headquarters, Melbourne, Australia** (BHP Billiton Limited and BHP Billiton Group) **London, United Kingdom ...**

Shopping

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



As Google matured from rules to ML / AI



Exercise

What are the top 20 main **characters** (people) and **locations** in the “Book_CMC_AD” book?

And how frequently were they mentioned?



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, IBM NLU, 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



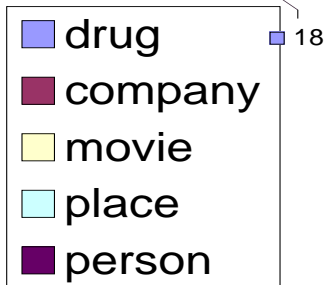
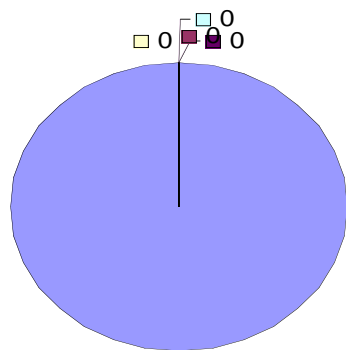
Features for sequence labeling

- 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

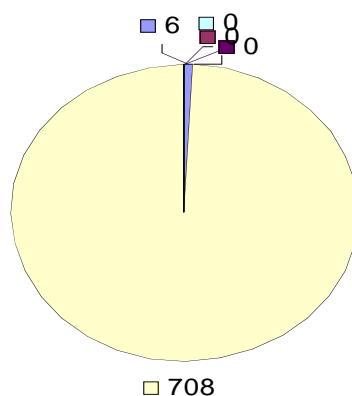


Features: Word substrings

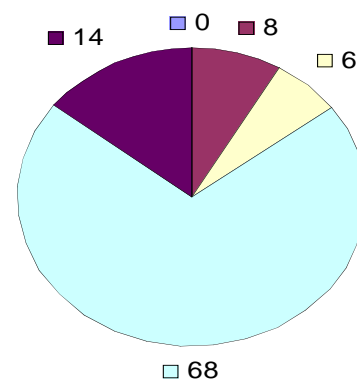
oxa



:



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

NER in Python





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](#)

Stanford University
Leland Stanford Junior University

Seal of Stanford University

Motto *Die Luft der Freiheit weht* (German)^[1]

Motto in English The wind of free thought blows!^[1]

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



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



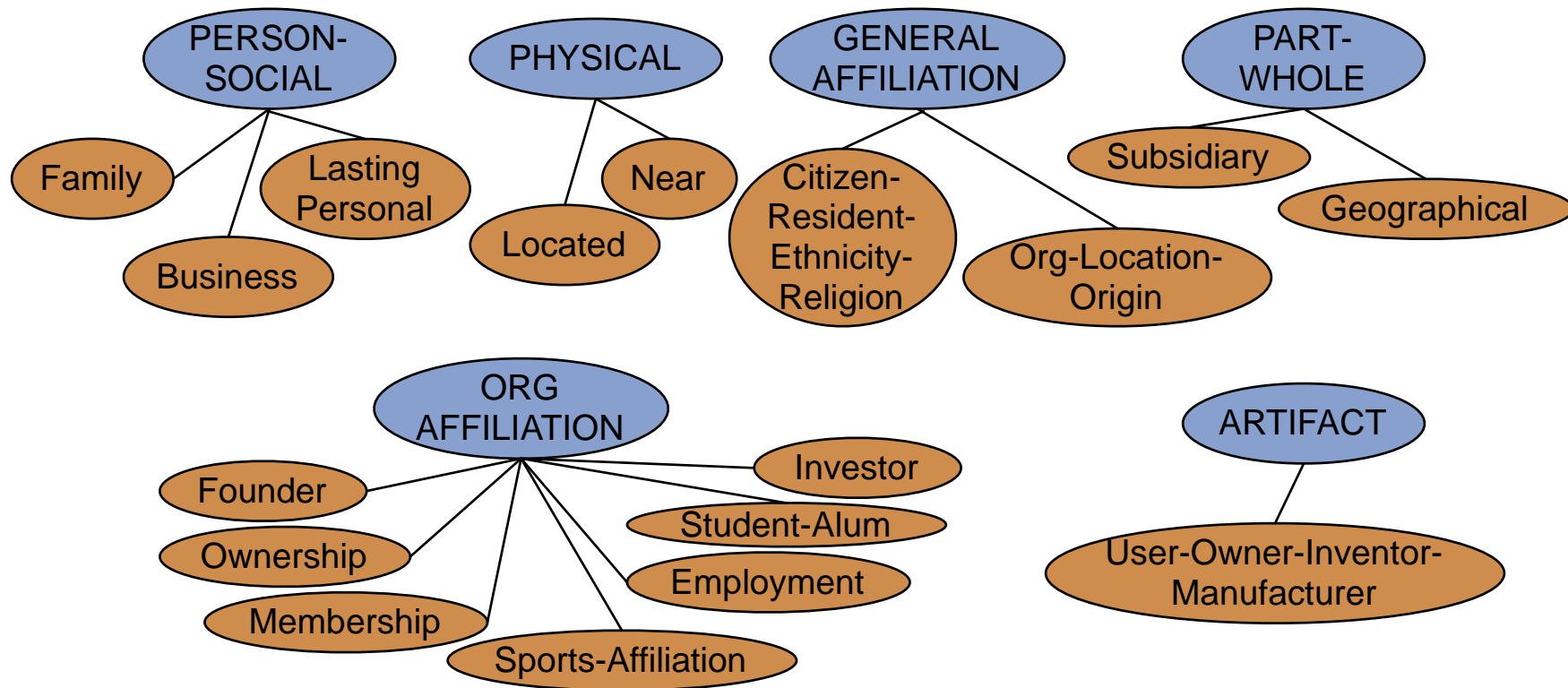
Why Relation Extraction?

- 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:

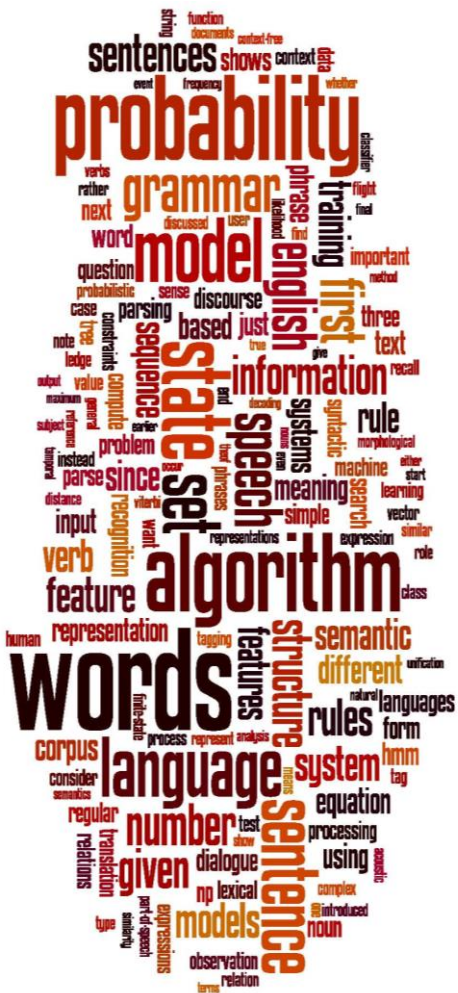
people/person/nationality,	location/location/contains
people/person/profession,	people/person/place-of-birth
biology/organism_higher_classification	film/film/genre



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



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



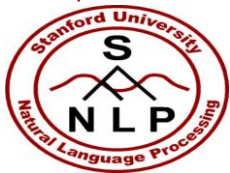
Hand-built patterns for relations

- 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 in Python

[illegible]

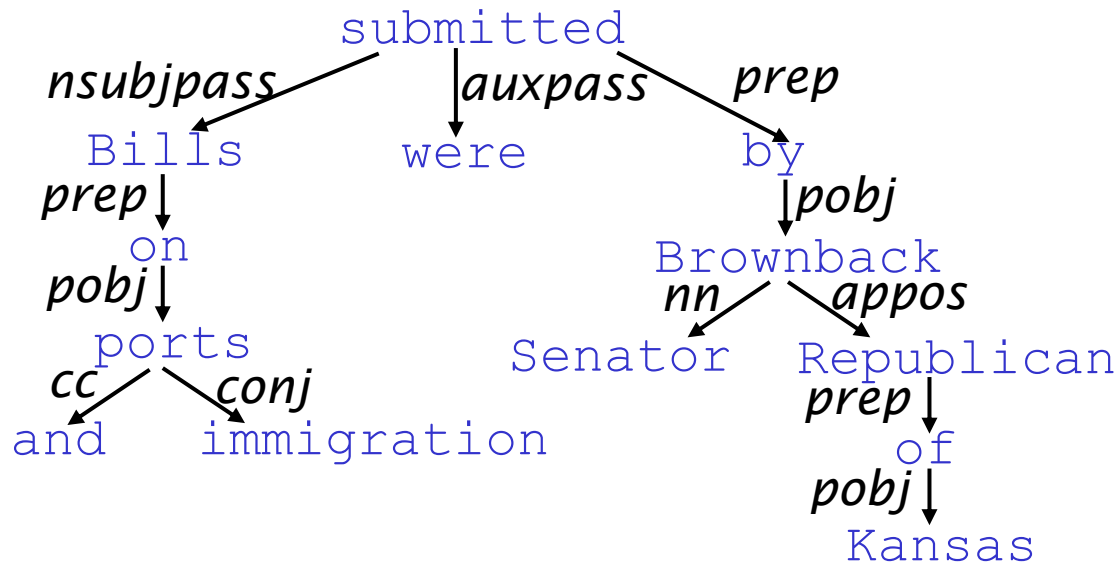
Introduction

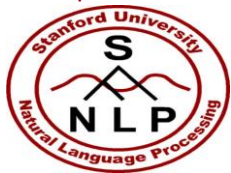


Dependency Grammar and Dependency Structure

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.)



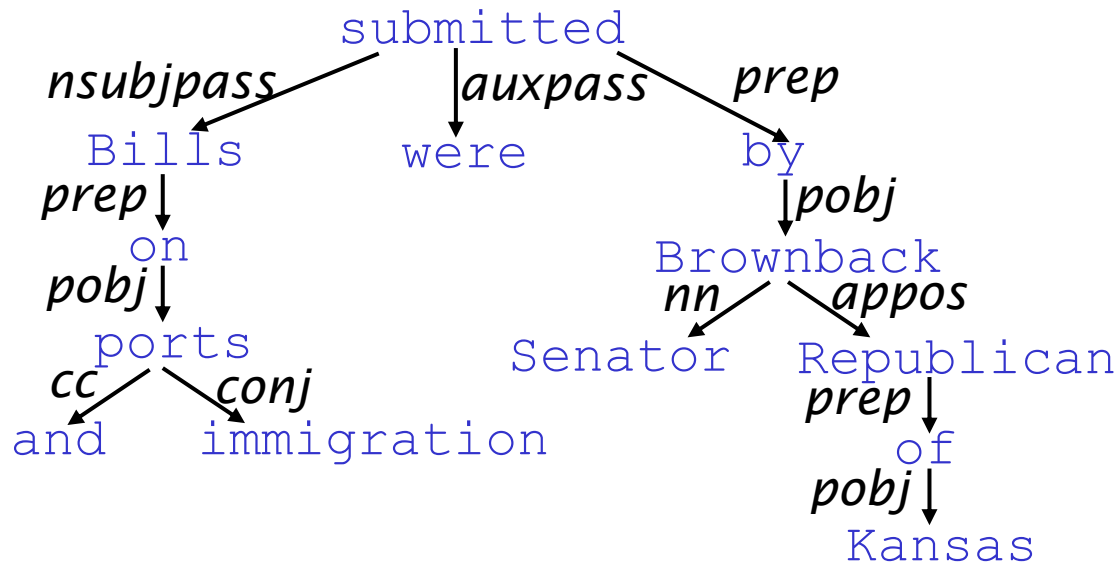


Dependency Grammar and Dependency Structure

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)



Dependency Parsing in Python

Watson NLU Demo



Thank You!

