

## Named Entity Recognition

Natural Language Processing  
CSCI 6350  
Dr. Sal Barbosa

## Named Entity Recognition

- Specific type of information extraction in which the goal is to extract formal names of particular types of entities such as people, places, organizations, etc.
- Usually a preprocessing step for subsequent task-specific IE, or other tasks such as question answering.

2

from slides by Raymond Mooney, University of Austin

## Named Entity Recognition Example

### U.S. Supreme Court quashes 'illegal' Guantanamo trials

Military trials arranged by the Bush administration for detainees at Guantanamo Bay are illegal, the United States Supreme Court ruled Thursday. The court found that the trials — known as military commissions — for people detained on suspicion of terrorist activity abroad do not conform to any act of Congress. The justices also rejected the government's argument that the Geneva Conventions regarding prisoners of war do not apply to those held at Guantanamo Bay. Writing for the 5-3 majority, Justice Stephen Breyer said the White House had overstepped its powers under the U.S. Constitution. "Congress has not issued the executive a blank cheque," Breyer wrote.

President George W. Bush said he takes the ruling very seriously and would find a way to both respect the court's findings and protect the American people.

3

from slides by Raymond Mooney, University of Austin

## Named Entity Recognition Example

**people**                      **places**                      **organizations**  
**U.S. Supreme Court** quashes 'illegal' Guantanamo trials

Military trials arranged by the Bush administration for detainees at Guantanamo Bay are illegal, the United States Supreme Court ruled Thursday. The court found that the trials — known as military commissions — for people detained on suspicion of terrorist activity abroad do not conform to any act of Congress. The justices also rejected the government's argument that the Geneva Conventions regarding prisoners of war do not apply to those held at Guantanamo Bay. Writing for the 5-3 majority, Justice Stephen Breyer said the White House had overstepped its powers under the U.S. Constitution. "Congress has not issued the executive a blank cheque," Breyer wrote.

President George W. Bush said he takes the ruling very seriously and would find a way to both respect the court's findings and protect the American people.

4

from slides by Raymond Mooney, University of Austin

## Named Entity Types and Examples

Type	Tag	Sample Categories
People	PER	Individuals, fictional characters, small groups
Organization	ORG	Companies, agencies, political parties, religious groups, sports teams
Location	LOC	Physical extents, mountains, lakes, seas
Geo-Political Entity	GPE	Countries, states, provinces, counties
Facility	FAC	Bridges, buildings, airports
Vehicles	VEH	Planes, trains, and automobiles

Type	Example
People	<i>Turing</i> is often considered to be the father of modern computer science.
Organization	The <i>IPCC</i> said it is likely that future tropical cyclones will become more intense.
Location	The <i>Mr. Sanitas</i> loop hike begins at the base of <i>Sunshine Canyon</i> .
Geo-Political Entity	<i>Palo Alto</i> is looking at raising the fees for parking in the University Avenue district.
Facility	Drivers were advised to consider either the <i>Tappan Zee Bridge</i> or the <i>Lincoln Tunnel</i> .
Vehicles	The updated <i>Mini Cooper</i> retains its charm and agility.

5

from Speech and Language Processing (2nd ed) by Jurafsky and Martin

## Ambiguous Named Entity Types

Name	Possible Categories
<i>Washington</i>	Person, Location, Political Entity, Organization, Facility
<i>Downing St.</i>	Location, Organization
<i>IRA</i>	Person, Organization, Monetary Instrument
<i>Louis Vuitton</i>	Person, Organization, Commercial Product

[*PERS* Washington] was born into slavery on the farm of James Burroughs.  
 [*ORG* Washington] went up 2 games to 1 in the four-game series.  
 Blair arrived in [*LOC* Washington] for what may well be his last state visit.  
 In June, [*GPE* Washington] passed a primary seatbelt law.  
 The [*FAC* Washington] had proved to be a leaky ship, every passage I made...

6

adapted from Speech and Language Processing (2nd ed) by Jurafsky and Martin

## Features Used in Training NER

identity of  $w_i$ , identity of neighboring words  
 embeddings for  $w_i$ , embeddings for neighboring words  
 part of speech of  $w_i$ , part of speech of neighboring words  
 base-phrase syntactic chunk label of  $w_i$  and neighboring words  
 presence of  $w_i$  in a **gazetteer**  
 $w_i$  contains a particular prefix (from all prefixes of length  $\leq 4$ )  
 $w_i$  contains a particular suffix (from all suffixes of length  $\leq 4$ )  
 $w_i$  is all upper case  
 word shape of  $w_i$ , word shape of neighboring words  
 short word shape of  $w_i$ , short word shape of neighboring words  
 presence of hyphen

7

from Speech and Language Processing (2nd ed) by Jurafsky and Martin

## Shape Features

Shape	Example
Lower	cummings
Capitalized	Washington
All caps	IRA
Mixed case	eBay
Capitalized character with period	H.
Ends in digit	A9
Contains hyphen	H-P

8

adapted from Speech and Language Processing by Jurafsky and Martin

# Feature Encoding for NER

Features			Label
American	NNP	<i>B<sub>NP</sub></i>	cap
Airlines	NNPS	<i>I<sub>NP</sub></i>	cap
,	PUNC	O	punc
a	DT	<i>B<sub>NP</sub></i>	lower
unit	NN	<i>I<sub>NP</sub></i>	lower
of	IN	<i>B<sub>PP</sub></i>	lower
AMR	NNP	<i>B<sub>NP</sub></i>	upper
Corp.	NNP	<i>I<sub>NP</sub></i>	cap.punc
,	PUNC	O	punc
immediately	RB	<i>B<sub>ADJP</sub></i>	lower
matched	VBD	<i>B<sub>VP</sub></i>	lower
the	DT	<i>B<sub>NP</sub></i>	lower
move	NN	<i>I<sub>VP</sub></i>	lower
,	PUNC	O	punc
spokesman	NN	<i>B<sub>NP</sub></i>	lower
Tim	NNP	<i>I<sub>NP</sub></i>	cap
Wagner	NNP	<i>I<sub>NP</sub></i>	cap
said	VBD	<i>B<sub>VP</sub></i>	lower
,	PUNC	O	punc

9

from Speech and Language Processing by Jurafsky and Martin

# NER as Sequence Labeling Task

- The standard algorithm for named entity recognition is as a word-by-word sequence labeling task
- Assigned tags capture both the boundary and the NE type
- Sequence classifiers like MEMM/CRF, bi-LSTM, or transformers are trained to label the tokens with tags that indicate the presence of particular NE types

The diagram illustrates the NER as a sequence labeling task. It shows a sequence of tokens: "a", "unit", "of", "AMR", "Corp.", ",", "immediately", "matched". Above the tokens are labels: "O", "B\_ORG", "?", and "...". A "Classifier" box is shown with arrows pointing to the labels. The classifier also outputs various NE types: IN, B\_PP, lower, NNP, B\_NP, upper, NNP, I\_NP, cap\_punc, PUNC, O, punc, RB, B\_ADVP, lower, VBD, B\_VP, lower.

# Statistical Sequence Labeling

```
graph TD; A[Representative Document Collection] --> B[/Human Annotation/]; B --> C[/Annotated Documents/]; C --> D[Feature Extraction and IOB Encoding]; D --> E[/Training Data/]; E --> F[Train Classifiers to Perform Multiway Sequence Labeling (MEMMs, CRFs, SVMs, HMMs, etc.)]; F --> G[NER System];
```

The diagram illustrates the workflow for Statistical Sequence Labeling. It begins with a 'Representative Document Collection' (represented by a stack of blue rectangles), which leads to 'Human Annotation' (a blue trapezoid). This step produces 'Annotated Documents' (a blue parallelogram). These documents then undergo 'Feature Extraction and IOB Encoding' (a blue rectangle), resulting in 'Training Data' (a blue parallelogram). The training data is used to 'Train Classifiers to Perform Multiway Sequence Labeling (MEMMs, CRFs, SVMs, HMMs, etc.)' (a blue rectangle). Finally, the trained classifiers are deployed as the 'NER System' (a blue rectangle).

11

from Speech and Language Processing by Jurafsky and Martin

# Code Demo