

# CSCI 5832

## Natural Language Processing

**Lecture 23**  
**Jim Martin**

4/24/07

CSCI 5832 Spring 2006

1

## Today: 4/17

- **Finish Lexical Semantics**
- **Wrap up Information Extraction**

4/24/07

CSCI 5832 Spring 2006

2

## Inside Words

- **Thematic roles: more on the stuff that goes on inside verbs.**

4/24/07

CSCI 5832 Spring 2006

3

## Inside Verbs

- **Semantic generalizations over the specific roles that occur with specific verbs.**
- **I.e. Takers, givers, eaters, makers, doers, killers, all have something in common**
  - -er
  - They're all the **agents** of the actions
- **We can generalize (or try to) across other roles as well**

4/24/07

CSCI 5832 Spring 2006

4

## Thematic Roles

Thematic Role	Definition
AGENT	The volitional causer of an event
EXPERIENCER	The experiencer of an event
FORCE	The non-volitional causer of the event
THEME	The participant most directly affected by an event
RESULT	The end product of an event
CONTENT	The proposition or content of a propositional event
INSTRUMENT	An instrument used in an event
BENEFICIARY	The beneficiary of an event
SOURCE	The origin of the object of a transfer event
GOAL	The destination of an object of a transfer event

4/24/07

CSCI 5832 Spring 2006

5

## Thematic Role Examples

Thematic Role	Example
AGENT	<i>The waiter</i> spilled the soup.
EXPERIENCER	<i>John</i> has a headache.
FORCE	<i>The wind</i> blows debris from the mall into our yards.
THEME	Only after Benjamin Franklin broke <i>the ice...</i>
RESULT	The French government has built a <i>regulation-size baseball diamond...</i>
CONTENT	Mona asked " <i>You met Mary Ann at a supermarket</i> "?
INSTRUMENT	He turned to poaching catfish, stunning them <i>with a shocking device...</i>
BENEFICIARY	Whenever Ann Callahan makes hotel reservations <i>for her boss...</i>
SOURCE	I flew in <i>from Boston.</i>
GOAL	I drove <i>to Portland.</i>

4/24/07

CSCI 5832 Spring 2006

6

## Why Thematic Roles?

- It's not the case that every verb is unique and has to introduce unique labels for all of its roles; thematic roles let us specify a fixed set of roles.
- More importantly it permits us to distinguish surface level shallow semantics from deeper semantics

4/24/07

CSCI 5832 Spring 2006

7

## Example

- From the WSJ...
  - *He melted her reserve with a husky-voiced paeon to her eyes.*
  - If we label the constituents **He** and **reserve** as the **Melter** and **Melted**, then those labels lose any meaning they might have had literally.
  - If we make them **Agent** and **Theme** then we don't have the same problems

4/24/07

CSCI 5832 Spring 2006

8

## Tasks

- **Shallow semantic analysis is defined as**
  - Assigning the right labels to the arguments of verb in a sentence. Aka
    - Case role assignment
    - Thematic role assignment

Thematic Role	De nition
AGENT	The volitional causer of an event
EXPERIENCER	The experiencer of an event
FORCE	The non-volitional causer of the event
THEME	The participant most directly affected by an event
RESULT	The end product of an event
CONTENT	The proposition or content of a propositional event
INSTRUMENT	An instrument used in an event
BENEFICIARY	The bene ciary of an event
SOURCE	The origin of the object of a transfer event
GOAL	The destination of an object of a transfer event

4/24/07

CSCI 5832 Spring 2006

9

## Example

- Newswire text
  - [*British forces* *agent*] [*believe* *target*] that [*Ali* *theme*] *was killed in a recent air raid*
  - British forces believe that [*Ali* *theme*] was [*killed* *target*] [*in a recent air raid* *temporal*]

4/24/07

CSCI 5832 Spring 2006

10

## Resources

- **PropBank**
  - Annotate every verb in the Penn Treebank with its semantic arguments.
  - Use a fixed (25 or so) set of role labels (Arg0, Arg1...)
  - Every verb has a set of frames associated with it that indicate what its roles are.
    - So for *Give* we're told that Arg0 -> *Giver*

4/24/07

CSCI 5832 Spring 2006

11

## Resources

- **Propbank**
  - Since it's built on the treebank we have the trees and the parts of speech for all the words in each sentence.
  - Since it's a corpus we have the statistical coverage information we need for training machine learning systems.

4/24/07

CSCI 5832 Spring 2006

12

## Resources

- **Propbank**
  - Since it's the WSJ it contains some fairly odd (domain specific) word uses that don't match our intuitions of the normal use of the words
  - Similarly, the word distribution is skewed by the genre from "normal" English (whatever that means).
  - There's no unifying semantic theory behind the various frame files (*buy* and *sell* are essentially unrelated).

4/24/07

CSCI 5832 Spring 2006

13

## Resources

- **FrameNet**
  - Instead of annotating a corpus, annotate domains of human knowledge a domain at a time (called frames)
    - Then within a domain annotate lexical items from within that domain.
    - Develop a set of semantic roles (called frame elements) that are based on the domain and **shared across** the lexical items in the frame.

4/24/07

CSCI 5832 Spring 2006

14

## Cause\_Harm Frame

### Cause\_harm

#### Definition:

The words in this frame describe situations in which an **Agent** or a **Cause** injures a **Victim**. The **Body\_part** of the **Victim** which is most directly affected may also be mentioned in the place of the **Victim**. In such cases, the **Victim** is often indicated as a genitive modifier of the **Body\_part**, in which case the **Victim** FE is indicated on a second FE layer.

4/24/07

CSCI 5832 Spring 2006

15

## Lexical Units

### Lexical Units

*bash.v, batter.v, bayonet.v, beat up.v, beat.v, belt.v, biff.v, bludgeon.v, boil.v, break.v, bruise.v, buffet.v, burn.v, butt.v, cane.v, chop.v, claw.v, clout.v, club.v, crack.v, crush.v, cudgel.v, cuff.v, cut.v, elbow.v, electrocute.v, electrocution.n, flagellate.v, flog.v, fracture.v, gash.v, hammer.v, hit.v, horsewhip.v, hurt.v, impale.v, injure.v, jab.v, kick.v, knee.v, knife.v, knock.v, lash.v, maim.v, maul.v, mutilate.v, pelt.v, poison.v, poisoning.n, pummel.v, punch.v, slap.v, slice.v, smack.v, smash.v, spear.v, squash.v, stab.v, stone.v, strike.v, thwack.v, torture.v, transfix.v, welt.v, whip.v, wound.v*

Created by hcb on Wed May 23 15:26:58 PDT 2001

4/24/07

CSCI 5832 Spring 2006

16

## FrameNet

- Frames and frame elements are entities in a hierarchy.
  - **Cause\_Harm** inherits from **Transitive\_Action**
  - **Corporal\_Punishment** inherits from **Cause\_Harm**
  - The **victim** FE in **Cause\_Harm** inherits from the **patient** FE of **Transitive\_Action**
  - And the **evaluatee** of the **Corporal\_Punishment** frame inherits from the **victim** of the **Cause\_Harm** frame.

4/24/07

CSCI 5832 Spring 2006

17

## FrameNet

- [Framenet.icsi.berkeley.edu](http://Framenet.icsi.berkeley.edu)

4/24/07

CSCI 5832 Spring 2006

18

## Break

**Thursday we'll turn to discourse (Chapter 20).**

**Next week Stat MT**

**Final quiz will be on May 1.**

4/24/07

CSCI 5832 Spring 2006

19

## HLT Certificate

**You may be on your way to the...**

***Human Language Technology Certificate***

**For typical CS students**

**5 courses**

**CS: NLP, UI design, AI**

**Ling: Syntax and Morphology, Phonetics**

4/24/07

CSCI 5832 Spring 2006

20

## Information Extraction

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 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

4/24/07

CSCI 5832 Spring 2006

21

## Information Extraction

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 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.

4/24/07

CSCI 5832 Spring 2006

22

## Named Entity Recognition

- Find the named entities and classify them by type.
- Typical approach
  - Acquire training data
  - Encode using IOB labeling
  - Train a sequential supervised classifier
  - Augment with pre- and post-processing using available list resources (census data, gazeteers, etc.)

4/24/07

CSCI 5832 Spring 2006

23

## Information Extraction

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 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

4/24/07

CSCI 5832 Spring 2006

24

## Temporal and Numerical Expressions

- **Temporals**
  - Find all the temporal expressions
  - Normalize them based on some reference point
- **Numerical Expressions**
  - Find all the expressions
  - Classify by type
  - Normalize

4/24/07

CSCI 5832 Spring 2006

25

## Information Extraction

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 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

4/24/07

CSCI 5832 Spring 2006

26

## Event Detection

- Find and classify all the events in a text.

4/24/07

CSCI 5832 Spring 2006

27

## Information Extraction

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 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

4/24/07

CSCI 5832 Spring 2006

28

## Relation Extraction

- **Basic task: find all the classifiable relations among the named entities in a text (populate a database)...**
  - **Employs**
    - { <American, Tim Wagner> }
  - **Part-Of**
    - { <United, UAL>, {American, AMR} >

4/24/07

CSCI 5832 Spring 2006

29

## Relation Extraction

- **Typical approach:**
  - For all pairs of entities in a text
    - **Extract features from the text span that just covers both of the entities**
      - Use a binary classifier to decide if there is likely to be a relation
      - If yes: then apply each of the known classifiers to the pair to decide which one it is
  - **Use supervised ML to train the required classifiers from an annotated corpus**

4/24/07

CSCI 5832 Spring 2006

30

## Information Extraction

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 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

4/24/07

CSCI 5832 Spring 2006

31

## Template Analysis

- **Many news stories have a script-like flavor to them. They have fixed sets of expected events, entities, relations, etc.**
- **Template, schemas or script processing involves:**
  - **Recognizing that a story matches a known script**
  - **Extracting the parts of that script**

4/24/07

CSCI 5832 Spring 2006

32

## Template Analysis

- So airlines often try to raise fares. Sometimes it sticks, sometimes it doesn't; it depends on how the other airlines react to the increase.
  - Airline that starts it off: **United**
  - Effective date of the increase: **Thursday**
  - Amount of the increase: **\$6**
  - Followers: **American**
  - Routes: ...

4/24/07

CSCI 5832 Spring 2006

33

## Template Processing

- Builds on earlier steps; obviously helps to know the entity types of the things that can fill the slots in the script.
- One approach...
  - Use supervised ML (with IOB labeling) to label all the candidate segments with their roles.
  - Collect all the candidate slots and resolve
    - If there's only one candidate take it
    - If not then vote or take the candidate with highest confidence score

4/24/07

CSCI 5832 Spring 2006

34

## Information Extraction

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 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

4/24/07

CSCI 5832 Spring 2006

35

## Information Extraction Summary

- **Named entity recognition and classification**
- **Coreference analysis**
- **Temporal and numerical expression analysis**
- **Event detection and classification**
- **Relation extraction**
- **Template analysis**

4/24/07

CSCI 5832 Spring 2006

36

## Information Extraction

- **Ordinary newswire text is often used in typical examples.**
  - **And there's an argument that there are useful applications there**
- **The real interest/money is in specialized domains**
  - **Bioinformatics**
  - **Patent analysis**
  - **Specific market segments for stock analysis**
  - **Intelligence analysis**
  - **Etc.**