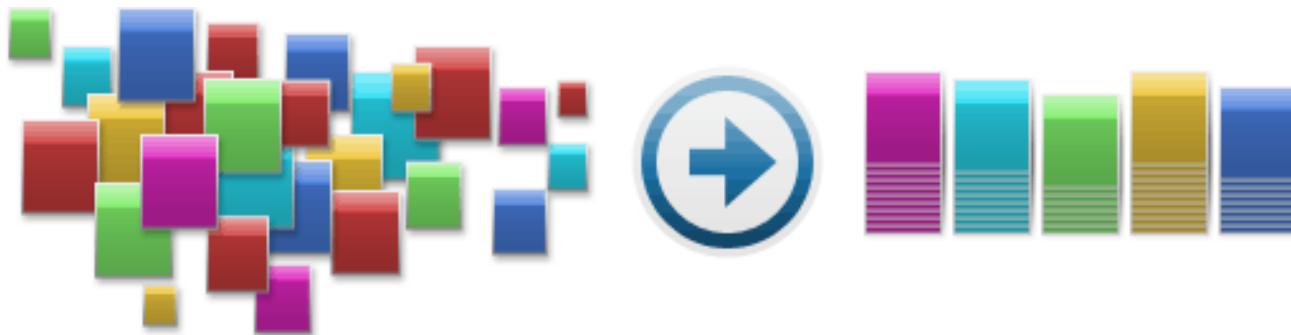


Information Extraction and Semantic Role Labeling

Information Extraction

- **Definition:** Information extraction is the process of turning unstructured information embedded in texts into structured information (example: relational databases)



Named Entity Recognition

- The first step is usually Named Entity Recognition (NER), which we talked about earlier in the course
 - Named entities often serve to define which entities appear in the text

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 lowercost carriers.

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

United, a unit of **UAL Corp.**, said the increase took effect Thursday and applies to most routes where it competes against discount carriers, such as **Chicago** to **Dallas** and **Denver** to San Francisco.

Legend:
Organization
Person
Location

Co-reference Resolution

- A common second step is finding equivalence classes of mentions that refer to the same entity

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 lowercost carriers.

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

United, a unit of **UAL Corp.**, said the increase took effect Thursday and applies to most routes where it competes against discount carriers, such as **Chicago** to **Dallas** and **Denver** to San Francisco.

Relation Extraction

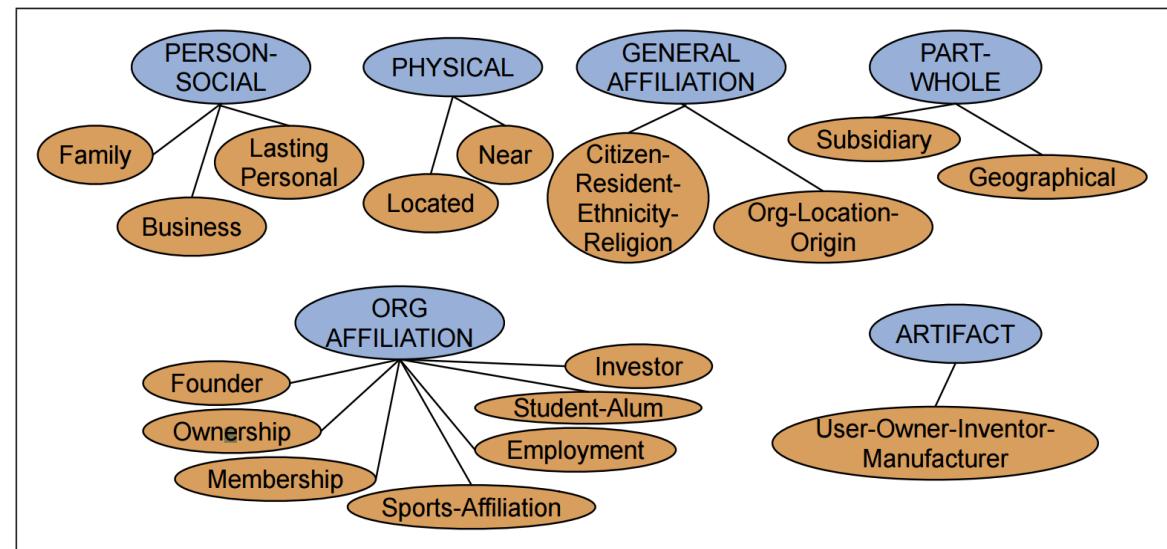
- Next on our list of tasks is to discern the relationships that exist among the detected entities

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 lowercost carriers. **American Airlines**, a unit of **AMR Corp.**, immediately matched the move, spokesman **Tim Wagner** said. **United**, a unit of **UAL Corp.**, said the increase took effect Thursday and applies to most routes where it competes against discount carriers, such as **Chicago** to **Dallas** and **Denver** to San Francisco.

- For instance, the text tells us that Tim Wagner is a spokesman for American Airlines, that United is a unit of UAL Corp., and that American Airlines is a unit of AMR

Relation Extraction

- Often the relations are mapped to a pre-defined ontology
 - This is an example from the ACE shared task:
- For instance the aforementioned relations are instances of the PART-WHOLE relation:
 - “United is a unit of UAL Corp.”
 - “American is a unit of AMR”



Ontologies can be Mapped to Model-Theoretic Semantics

Domain

United, UAL, American Airlines, AMR

Tim Wagner

Chicago, Dallas, Denver, and San Francisco

$$\mathcal{D} = \{a, b, c, d, e, f, g, h, i\}$$

$$a, b, c, d$$

$$e$$

$$f, g, h, i$$

Classes

United, UAL, American, and AMR are organizations

Tim Wagner is a person

Chicago, Dallas, Denver, and San Francisco are places

$$Org = \{a, b, c, d\}$$

$$Pers = \{e\}$$

$$Loc = \{f, g, h, i\}$$

Relations

United is a unit of UAL

American is a unit of AMR

Tim Wagner works for American Airlines

United serves Chicago, Dallas, Denver, and San Francisco

$$PartOf = \{\langle a, b \rangle, \langle c, d \rangle\}$$

$$OrgAff = \{\langle c, e \rangle\}$$

$$Serves = \{\langle a, f \rangle, \langle a, g \rangle, \langle a, h \rangle, \langle a, i \rangle\}$$

Relation Extraction

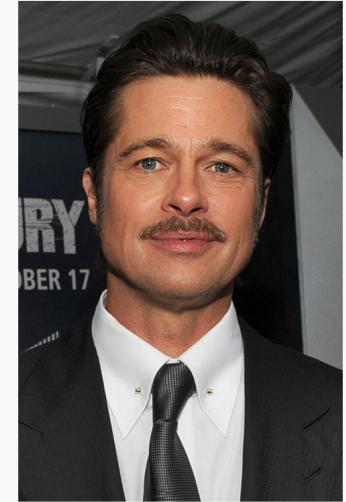
- There are many other ontologies of relations that have been defined
 - For example, UMLS, the Unified Medical Language System from the US National Library of Medicine has a network that defines 134 broad subject categories, such as:

Entity	Relation	Entity
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

Wikipedia Infoboxes

- Wikipedia also offers a range of relation types, drawn from the infoboxes
- For example, the Wikipedia infobox for Brad Pitt includes structured facts like *Occupation* = “Actor”/”Producer” or *Born* = “*Shawnee, Oklahoma*”
- These facts can be turned into relations like *born_in* or *occupied_as*
- These relations can complement relations extracted from text, or serve as features for them
- *Wikidata* (<https://www.wikidata.org>) is a knowledge base based on Wikipedia

Brad Pitt



Pitt at the premiere of *Fury* in Washington, D.C., October 2014

Born	William Bradley Pitt December 18, 1963 (age 53) <i>Shawnee, Oklahoma, U.S.</i>
Occupation	Actor • producer
Years active	1987–present
Works	Filmography
Home town	<i>Springfield, Missouri</i>
Spouse(s)	Jennifer Aniston (m. 2000; div. 2005) Angelina Jolie (m. 2014; separated 2016)
Children	6
Relatives	Douglas Pitt (brother)

Using Patterns to Discover Relations

- The earliest and still common approach for relation extraction is the use of lexico-syntactic patterns

NP {, NP}* {,} (and|or) other NP_H

temples, treasures, and other important **civic buildings**

NP_H such as {NP,}* {(or|and)} NP

red algae such as Gelidium

such NP_H as {NP,}* {(or|and)} NP

such **authors** as Herrick, Goldsmith, and Shakespeare

NP_H {,} including {NP,}* {(or|and)} NP

common-law countries, including Canada and England

NP_H {,} especially {NP,}* {(or|and)} NP

European countries, especially France, England, and Spain

Using Patterns to Discover Relations

- More modern approaches use additional features to define the patterns, such as named entity constraints
- For instance, if our goal is to answer questions about “who holds what office in which organization?”, we can use patterns like the following:

PER, POSITION of ORG:

George Marshall, Secretary of State of the United States

PER (named|appointed|chose|etc.) PER Prep? POSITION

Truman appointed Marshall Secretary of State

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

George Marshall was named US Secretary of State

Using Patterns to Discover Relations

- More accurate patterns still are ones that use syntactic information
- For instance, symmetric patterns can be used to discover synonymy or relatedness between entities

from X to Y

X and Y

X or Y

neither X nor Y

X as well as Y

Using Patterns to Discover Relations

- Syntactic information can filter a lot of noise
- What are the Xs and Ys in these cases:
 - “when they go to Austria, they like walking in the woods **as well as** skiing”
 - “apricots **and** other vegetables **and** fruit”
 - “Sandy is a Republican **and** proud of it”

from X to Y

X and Y

X or Y

neither X nor Y

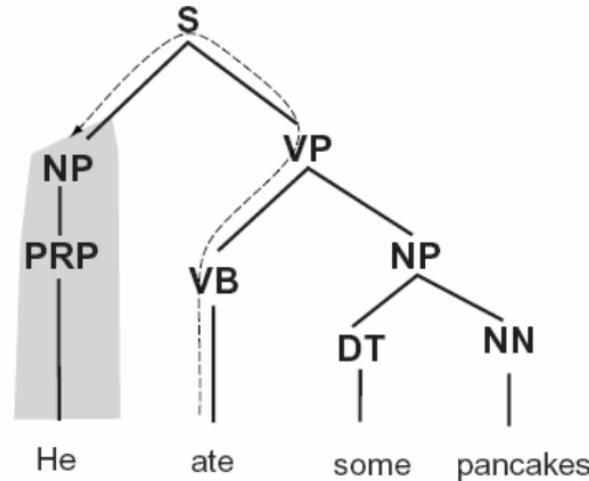
X as well as Y

Distant Supervision for Relation Extraction

- Hand-labeled data with relation labels is expensive to produce
- However, available resources such as Wikipedia infoboxes, have a great many relations that are in structured format
 - Wikipedia articles contain sentences that express these relations
 - This huge amount of examples allows us to define rich features

Path Features

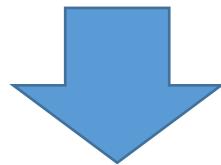
- Path features:
 - Encode the path between two nodes, through the direction of the edges and their labels
- Similar features exist for dependency parses
 - We will see a bit of this in Exercise 4



<i>Path</i>	<i>Description</i>
VB↑VP↓PP	PP argument/adjunct
VB↑VP↑S↓NP	subject
VB↑VP↓NP	object
VB↑VP↑VP↑S↓NP	subject (embedded VP)
VB↑VP↓ADVP	adverbial adjunct
NN↑NP↑NP↓PP	prepositional complement of noun

Examples of Rich Features

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



PER was born in LOC
PER, born *, LOC
PER's birthplace is LOC

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



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

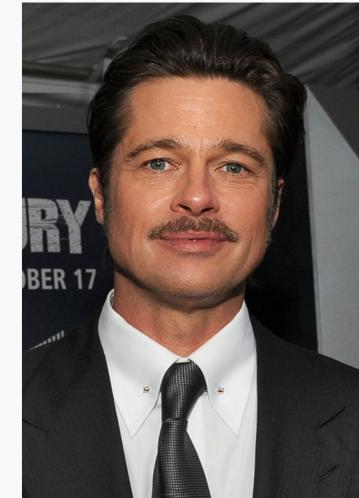
Schematized Algorithm for Distant Supervision

```
function DISTANT SUPERVISION(Database D, Text T) returns relation classifier C
    foreach relation R
        foreach tuple  $(e_1, e_2)$  of entities with relation R in D
            sentences  $\leftarrow$  Sentences in T that contain  $e_1$  and  $e_2$ 
            f  $\leftarrow$  Frequent features in sentences
            observations  $\leftarrow$  observations + new training tuple  $(e_1, e_2, f, R)$ 
        C  $\leftarrow$  Train supervised classifier on observations
    return C
```

Example: Distant Supervision

- From the infobox we learn that “Brad Pitt” and “Jennifer Aniston” were spouses
- The Wikipedia text states:
 - “In 2000, **he** married actress **Jennifer Aniston**; they divorced in 2005.”
 - “... playing a man with a grudge against Rachel Green, played by **Jennifer Aniston**, to whom **Pitt** was married at the time.”
 - “... Plan B Entertainment, a film production company **he** had founded two years earlier with **Jennifer Aniston** and Brad Grey”
- The first two examples are examples as to how spouse relation may be expressed in the text; the third example is noise

Brad Pitt



Pitt at the premiere of *Fury* in Washington, D.C., October 2014

Born	William Bradley Pitt December 18, 1963 (age 53) Shawnee, Oklahoma, U.S.
Occupation	Actor • producer
Years active	1987–present
Works	Filmography
Home town	Springfield, Missouri
Spouse(s)	Jennifer Aniston (m. 2000; div. 2005) Angelina Jolie (m. 2014; separated 2016)
Children	6
Relatives	Douglas Pitt (brother)

Evaluation of Relation Extraction

- Semi-supervised methods are much more difficult to evaluate than supervised methods
 - They extract **new** relations from the web or a large text
 - As methods use very large amounts of text, it is impossible to test them in a sand box with a small pre-annotated gold standard
- Evaluation is therefore done by:
 - Computing precision
 - If the system can rank its output, we can measure precision for different output sizes (e.g., precision for 100 relations, 1000 relations etc.)
 - Extrinsic evaluation by testing how well these relations help, say, question answering

Poon and Domingos, *Unsupervised Semantic Parsing*,
EMNLP 2009

Open Information Extraction

- Open Information Extraction (OpenIE) seeks to discover information from plain text
 - One way is to apply syntactic parsing and extract the relations it finds between named entities
 - This works well in some examples:
“**Pitt** was born in **Shawnee, Oklahoma**”
“**Jolie** gave birth to daughter **Shiloh Nouvel** in Swakopmund, Namibia, on May 27, 2006.”
- Extractions:
born_in(Pitt, Shawnee)
gave_birth_to_daughter(Jolie, Shiloh_Nouvel)

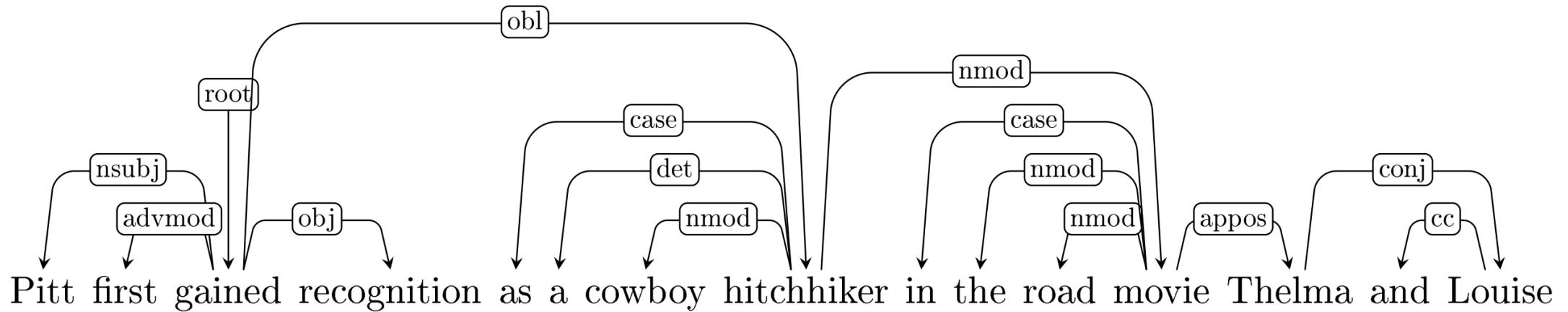
Open Information Extraction

- Open Information Extraction (OpenIE) seeks to discover information from plain text
 - One way is to apply syntactic parsing and extract the relations it finds between named entities
 - Consider this example:

Pitt first gained recognition as a cowboy hitchhiker in the road movie **Thelma and Louise** (1991). His first leading roles in big-budget productions came with the dramas **A River Runs Through It** (1992) and **Legends of the Fall** (1994), and **Interview with the Vampire** (1994).
 - What's the relation between "Pitt" and "Thelma and Louise"? Between "A River Runs Through It" and "Interview with the Vampire"?

Open Information Extraction

- The dependency path between *Pitt* and *Thelma and Louise* is $nsubj \uparrow obl \downarrow nmod \downarrow appos$



- This is an indirect syntactic relation
 - It would be useful to have a representation that directly encodes the **semantic** relations between the participants

Syntax doesn't always align w/ Semantics

Mary opened **the door**.

The door opened.

John slices **the bread** with **a knife**.

The bread slices easily.

The knife slices easily.

Mary loaded **the truck** with **hay**.

Mary loaded **hay** onto **the truck**.

The truck was loaded with **hay** (by **Mary**).

Hay was loaded onto **the truck** (by **Mary**).

John got **Mary** a present.

John got a present for **Mary**.

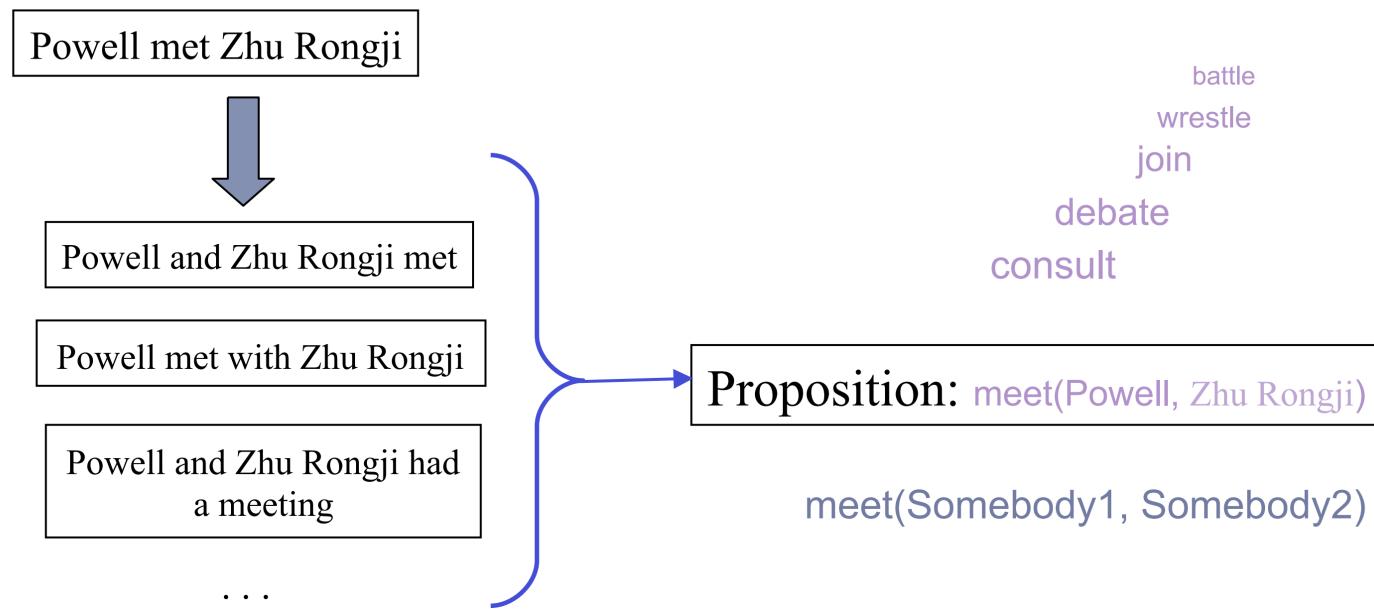
Mary got a present from **John**.

Even more variation
for indirect relations

Examples by Nathan Schneider

Semantic Role Labeling (SRL)

- The task of *Semantic Role Labeling* aims to represent a sentence in terms of its events or propositions
- Events consist of a main relation, participants and secondary relations



When Powell met Zhu Rongji on Thursday they discussed the return of the spy plane.

meet(Powell, Zhu) discuss([Powell, Zhu], return(X, plane))

Semantic Role Labeling

- Informally - determining **who** did **what** to **whom**, **where**, **when** and **how**
- More formally – **identifying**, for each predicate, its set of arguments and **establishing the semantic role** of each of them
 - Roles are usually taken from a pre-defined ontology
- Produces a flat (not hierarchical) structure for each of the predicates

SRL Example

Acceptor Modal Negation **Predicate** Thing Accepted

He would n't **accept** anything of value
from those he was writing about.

Accepted From

Semantic Roles

- Semantic role:
 - An underlying relationship an argument has with the main verb in a clause
 - This relation is (ideally) invariant to paraphrases
- Many semantic role lists (or ontologies) in the literature
- In this lecture – the prevailing approaches in NLP

Annotation Schemes of Semantic Roles

- We will briefly survey the three main representation approaches for SRL, which differ in:
 1. Organizing principles
 2. Scope (what is considered an event)
 3. Granularity
 4. Consistency within and across verbs

SRL Schemes: FrameNet

- Organized by *frames*: a schematized event type
- A frame is invoked by frame elements, generally words or morphemes
 - They constitute the anchor of the scene and determine what it is about
- Semantic roles are defined per frames
- Example:
 - The *Judgment* frame
 - Frame elements: *admire*, *appreciate*, *value*...

Judge

Evaluee

Reason

She **blames** the government for failing to do enough to help

FrameNet Judgement Frame Example

Judgment

[Lexical Unit Index](#)

Definition:

A **Cognizer** makes a judgment about an **Evaluee**. The judgment may be positive (e.g. respect) or negative (e.g. condemn), and this information is recorded in the semantic types Positive and Negative on the Lexical Units of this frame. There may be a specific **Reason** for the **Cognizer**'s judgment, or there may be a capacity or **Role** in which the **Evaluee** is **judged**.

This frame is distinct from the Judgment_communication frame in that this frame does not involve the Cognizer communicating his or her judgment to an Addressee.

JUDGMENT: She **ADMIRE**D Einstein for his character.

JUDGMENT_COMMUNICATION: She **ACCUSED** Einstein of collusion.

Currently, however, some lexical units and annotation for both remain in this frame.

FEs:

Core:

Cognizer [Cog]
Semantic Type: Sentient

The **Cognizer** makes the judgment. This role is typically expressed as the External Argument (or in a by-PP in passives).
The boss **APPRECIATES** you for your diligence.

The boss is very **APPRECIATIVE** of my work.

Evaluee [Eval]

Evaluee is the person or thing about whom/which a judgment is made. With verbs this FE is typically expressed as Object:
The boss **APPRECIATES** **you** for your diligence.

Expressor [Exr]

Expressor is the body part or action by a body part that conveys the judgment made by the **Cognizer**.
She viewed him with an **APPRECIATIVE** **gaze**.

Reason [Reas]
Semantic Type:
State_of_affairs

Typically, there is a constituent expressing the **REASON** for the **Judge**'s judgment. It is usually a for-PP, e.g.
I **ADMIRE** you **for your intellect**.

SRL Schemes: PropBank

- Predicate-specific core roles, with adjunct roles (such as temporal description, locations, manner adverbs, negation) are shared across predicates
 - So each predicate (e.g., *blame*) has core arguments which are indexed A0-A5
 - Non-core predicates are marked AM-*



PropBank includes Basic Predicate Sense Disambiguation

- **Decline.01** – “go down incrementally”

- *Arg1*: Entity going down
- *Arg2*: Amount gone down by
- *Arg3*: Start point
- *Arg4*: End point

[_{A1} Its net income] declining [_{A2} 42%] [_{A4} to \$121 million] [_{AM-TMP} in the first 9 months of 1989]

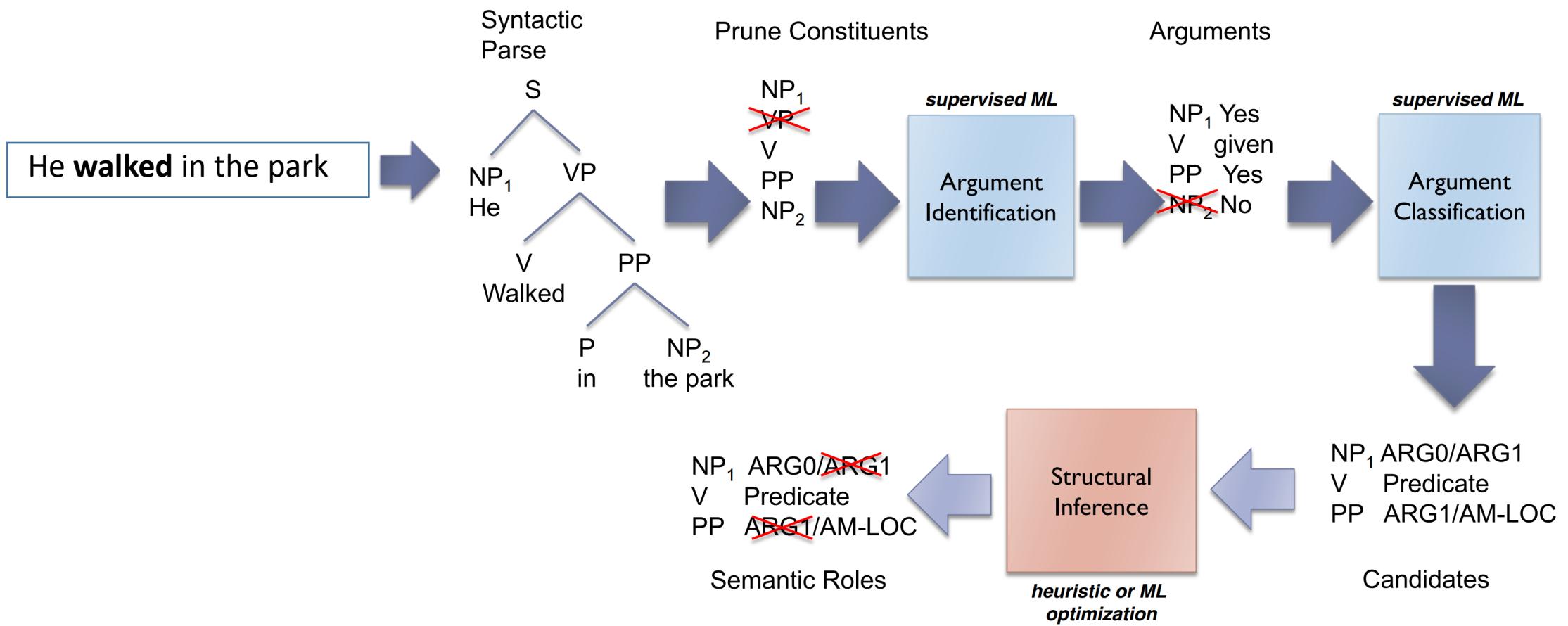


- **Decline.02** – “demure, reject”

- *Arg1*: Agent
- *Arg2*: Rejected Thing

[_{A1} A spokesman] declined [_{A2} to elaborate]

SRL Parsing: A Classic Architecture



SRL Evaluation

- Recall, Precision, F-score
- Variants when evaluating in SRL:
 - Arguments: Full span (CoNLL-2005), Headword only (CoNLL-2008)
- Predicates:
 - Given (CoNLL-2005)
 - System Identifies (CoNLL-2008)
 - Verb and nominal predicates (CoNLL-2008)

SRL Evaluation

Gold Standard Labels	SRL Output	Full	Head
Arg0: John	Arg0: John	+	+
Rel: mopped	Rel: mopped	+	+
Arg1: the floor	Arg1: the floor	+	+
Arg2: with the dress ... Thailand	Arg2: with the dress	-	+
Arg0: Mary	Arg0: Mary	+	+
Rel: bought	Rel: bought	+	+
Arg1: the dress	Arg1: the dress	+	+
Arg0: Mary		-	-
rel: studying		-	-
Argm-LOC: in Thailand		-	-
Arg0: Mary	Arg0: Mary	+	+
Rel: traveling	Rel: traveling	+	+
Argm-LOC: in Thailand		-	-

John mopped the floor with the
dress Mary bought while studying
and traveling in Thailand.

On Full Argument Span:

Precision = 8/9 = 88.9%

Recall = 8/13 = 61.5%

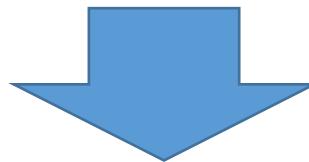
On Headword-only Evaluation:

Precision = 9/9 = 100%

Recall = 9/13 = 69.2%

Back to the Relation Extraction Example

Pitt first gained recognition as a cowboy hitchhiker in the road movie **Thelma & Louise** (1991)



gained recognition

Theme: Pitt

Role: as a cowboy hitchhiker

Means: Thelma & Louise

Extracting Times

- We have so far considered the analysis of events in terms of their predicates, participants and their roles
- Another important aspect of relation extraction is extracting the time when events happened
 - Absolute time
 - Relative time: events relative to one another

Extracting Times

- The most basic task of this type is extracting *temporal expressions*

Absolute	Relative	Durations
April 24, 1916	yesterday	four hours
The summer of '77	next semester	three weeks
10:15 AM	two weeks from yesterday	six days
The 3rd quarter of 2006	last quarter	the last three quarters

Figure 17.11 Examples of absolute, relational and durational temporal expressions.

- And temporal discourse connectives, such as *after*, *before*, *while* etc.

Extracting Times

- Even where the temporal expression is fairly explicit (e.g., ‘The summer of 77’), mapping it into an absolute time expression is a complex task
 - Ambiguity: what 77?
 - Underspecification: when exactly in the summer?
- Where the temporal expressions are relative (e.g., a week before), one has to identify the anchor, i.e., the reference point the relation refers to
 - The time of the anchor is usually not fully specified either

Extracting Event-Event Relations

- A related task is extracting the temporal relations of events depicted in the text
- The well-known TimeML format uses Allen’s interval algebra (Allen, 1983)
- Extracting inter-event relations and more generally timelines is still an open problem in NLP

a —	b —	a is <i>BEFORE</i> b b is <i>AFTER</i> a
a —	b —	a is <i>IBEFORE</i> b b is <i>IAFTER</i> a
a —		a <i>BEGINS</i> b
b —		b is <i>BEGUN_BY</i> a
	a —	a <i>ENDS</i> b
	b —	b is <i>ENDED_BY</i> a
a —		a is <i>DURING</i> b
b —		b is <i>DURING_INV</i> a
a —		a <i>INCLUDES</i> b
b —		b <i>IS_INCLUDED</i> in a
a —		a is <i>SIMULTANEOUS</i> with b
a —	b	a is <i>IDENTITY</i> with b