

# Lecture 18

# Information Extraction

Three horizontal lines of different colors (orange, black, and green) stacked vertically, spanning most of the width of the slide.

**CS 6320**

# Overview

---

- Information extraction – turns unstructured information buried in texts into structured data
- Extract proper nouns – “named entity recognition”
- Reference resolution –
  - named entity mentions
  - Pronoun references
- Relation Detection and classification
- Event detection and classification
- Temporal analysis
- Template filling

# Sample text

---

- Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PERSON Tim Wagner] said. [ORG United Airlines] an unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].
- Identify named entities
- Identify relations

# Template Filling

---

Example template for “airfare raise”

FARE-RAISE ATTEMPT:	LEAD AIRLINE:	UNITED AIRLINES
	AMOUNT:	\$6
	EFFECTIVE DATE:	2006-10-26
	FOLLOWER:	AMERICAN AIRLINES

# Some Named Entity Types

---

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

# Examples of Named Entity Types

---

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>Mt. 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.

# Categorical Ambiguities

---

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

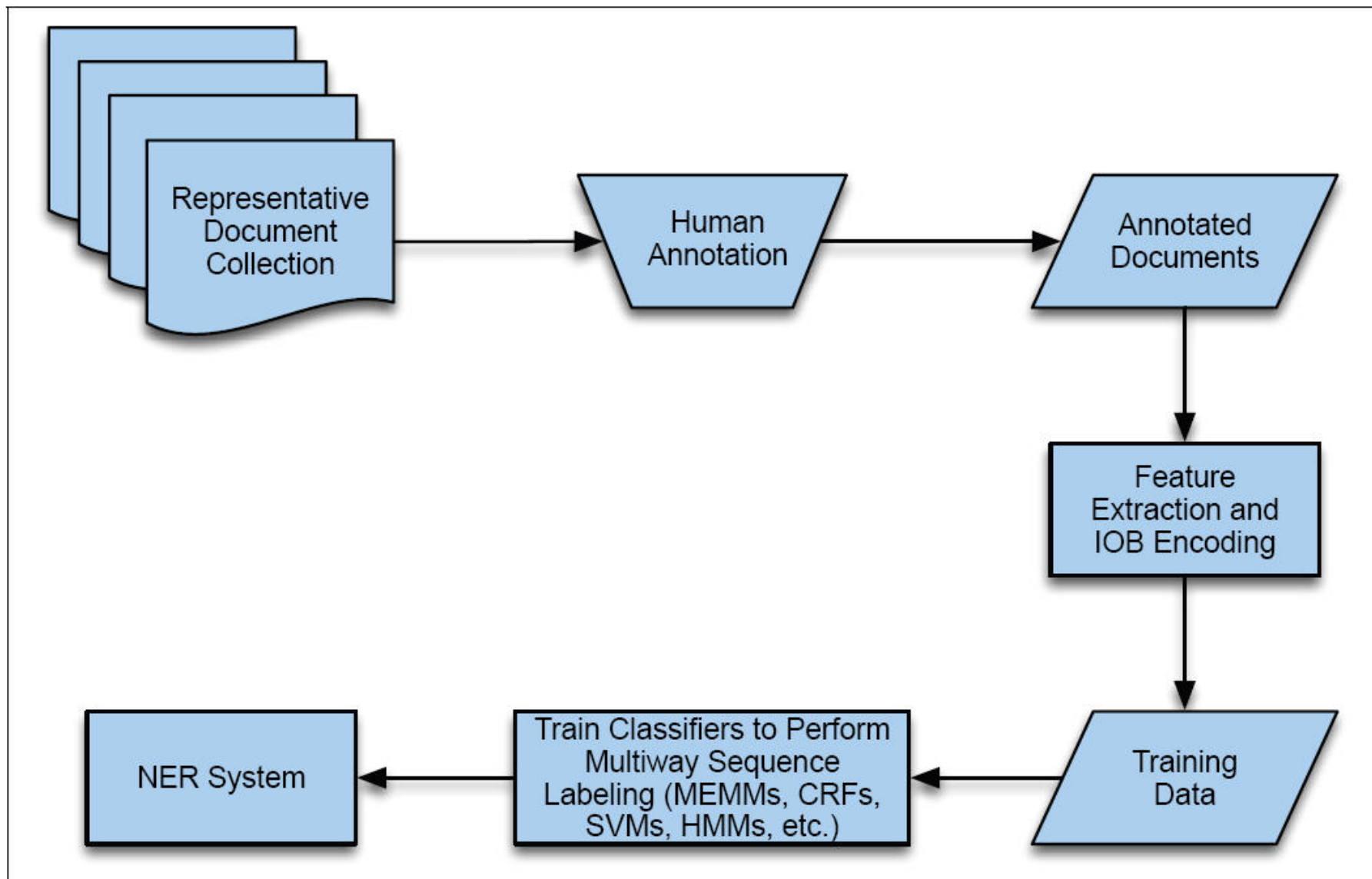
# Categorical Ambiguity

---

[*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...



# Statistical Seq. Labeling



# Features used in Training NER

Gazetteers – lists of place names

- [www.geonames.com](http://www.geonames.com)
- [www.census.gov](http://www.census.gov)

Feature	Explanation
Lexical items	The token to be labeled
Stemmed lexical items	Stemmed version of the target token
Shape	The orthographic pattern of the target word
Character affixes	Character-level affixes of the target and surrounding words
Part of speech	Part of speech of the word
Syntactic chunk labels	Base-phrase chunk label
Gazetteer or name list	Presence of the word in one or more named entity lists
Predictive token(s)	Presence of predictive words in surrounding text
Bag of words/Bag of N-grams	Words and/or <i>N</i> -grams occurring in the surrounding context

# Selected 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

# Boundary Detection

Words	IOB Label	IO Label
American	B-ORG	I-ORG
Airlines	I-ORG	I-ORG
,	O	O
a	O	O
unit	O	O
of	O	O
AMR	B-ORG	I-ORG
Corp.	I-ORG	I-ORG
,	O	O
immediately	O	O
matched	O	O
the	O	O
move	O	O
,	O	O
spokesman	O	O
Tim	B-PER	I-PER
Wagner	I-PER	I-PER
said	O	O
.	O	O

**Figure 17.4** Named entity tagging as a sequence model, showing IOB and IO encodings.

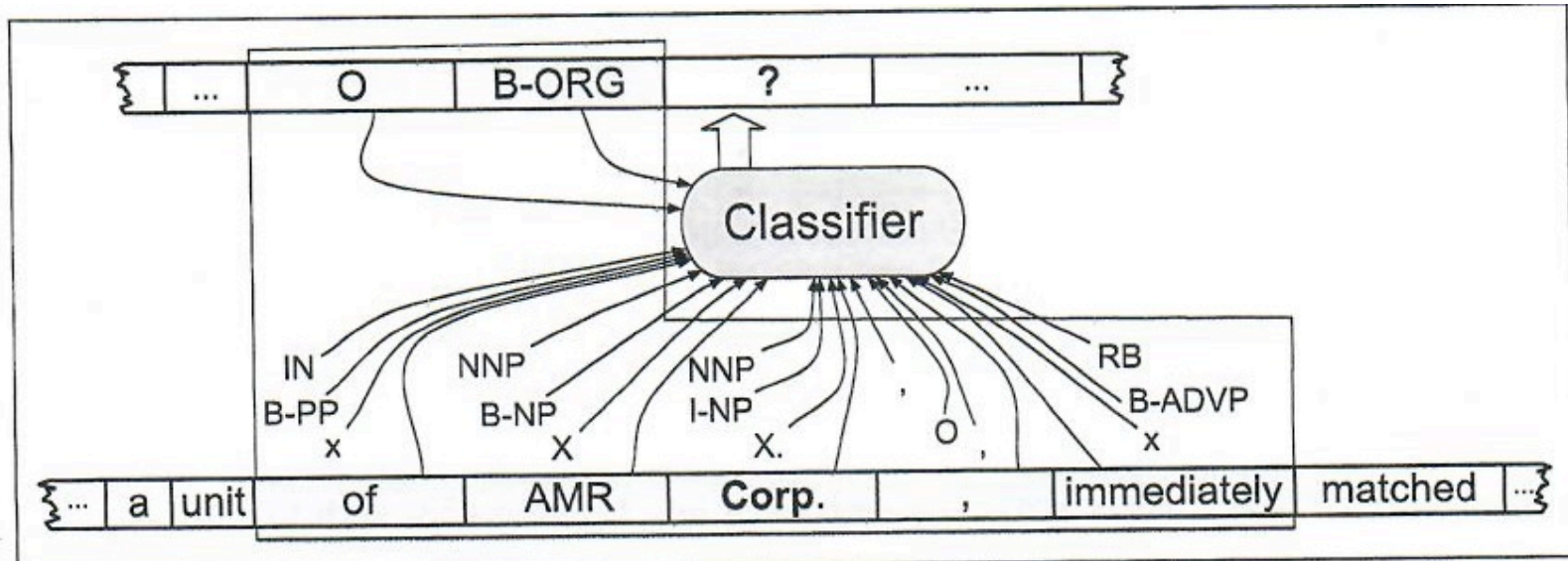
# BIO Encoding for NER

Word	POS	Chunk	Short shape	Label
American	NNP	B-NP	Xx	B-ORG
Airlines	NNPS	I-NP	Xx	I-ORG
,	,	O	,	O
a	DT	B-NP	x	O
unit	NN	I-NP	x	O
of	IN	B-PP	x	O
AMR	NNP	B-NP	X	B-ORG
Corp.	NNP	I-NP	Xx.	I-ORG
,	,	O	,	O
immediately	RB	B-ADVP	x	O
matched	VBD	B-VP	x	O
the	DT	B-NP	x	O
move	NN	I-NP	x	O
,	,	O	,	O
spokesman	NN	B-NP	x	O
Tim	NNP	I-NP	Xx	B-PER
Wagner	NNP	I-NP	Xx	I-PER
said	VBD	B-VP	x	O
.	.	O	.	O

**Figure 17.6** Word-by-word feature encoding for NER.

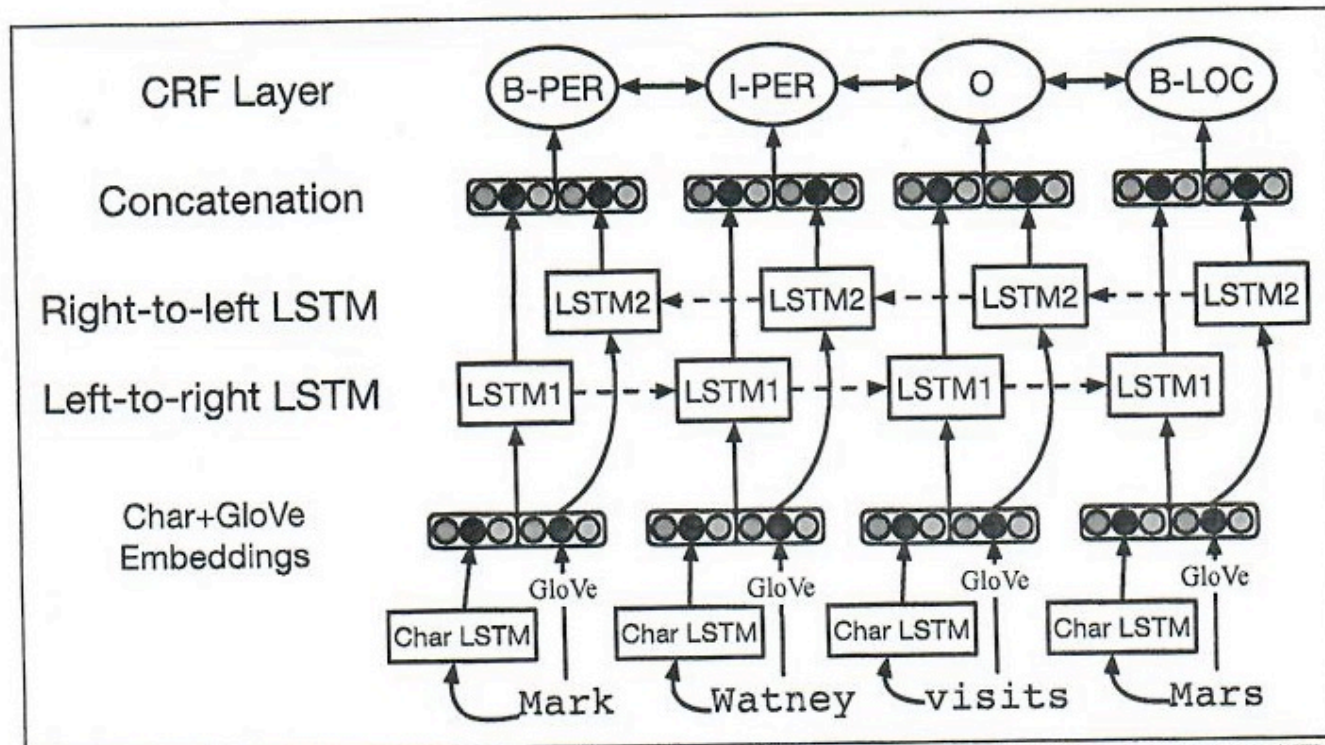


# NER Classifier with Input Features



**Figure 17.7** Named entity recognition as sequence labeling. The features available to the classifier during training and classification are those in the boxed area.

# A Neural Approach to NER



**Figure 17.8** Putting it all together: character embeddings and words together a bi-LSTM sequence model. After (Lample et al., 2016)

# Evaluation of N E R Systems

---

- Recall terms from Information retrieval
  - Recall = #correctly labeled / total # that should be labeled
  - Precision = # correctly labeled / total # labeled
- F- measure where  $\beta$  weights preferences
  - $\beta=1$  balanced
  - $\beta>1$  favors recall
  - $\beta<1$  favors precision

$$F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$



# NER Performance revisited

---

- NER performance revisited
  - Recall, Precision, F
  - High performance systems
    - »  $F \sim .92$  for PERSONS and LOCATIONS and  $\sim .84$  for ORG
- Practical Rule-based NER
  - Make several passes on text, allowing the results of one pass to influence the next

# Relation Detection and classification

---

- **Consider Sample text:**
  - Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PERSON Tim Wagner] said. [ORG United Airlines] an unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].
- **After identifying named entities what else can we extract?**
- **Relations**

# Example semantic relations

Relations		Examples	Types
Affiliations			
	Personal	<i>married to, mother of</i>	$\text{PER} \rightarrow \text{PER}$
	Organizational	<i>spokesman for, president of</i>	$\text{PER} \rightarrow \text{ORG}$
	Artifactual	<i>owns, invented, produces</i>	$(\text{PER} \mid \text{ORG}) \rightarrow \text{ART}$
Geospatial			
	Proximity	<i>near, on outskirts</i>	$\text{LOC} \rightarrow \text{LOC}$
	Directional	<i>southeast of</i>	$\text{LOC} \rightarrow \text{LOC}$
Part-Of			
	Organizational	<i>a unit of, parent of</i>	$\text{ORG} \rightarrow \text{ORG}$
	Political	<i>annexed, acquired</i>	$\text{GPE} \rightarrow \text{GPE}$

# Sub-Relations for Association

1. Communication – lexical constraint
  - a. COMMUNICATE
    - i. WRITE\_TO
      1. WRITTEN\_COMMUNICATION
      2. ELECTRONIC\_COMMUNICATION
    - ii. TELEPHONE\_TO
    - iii. SPEAK\_TO
    - iv. COMMANDS\_OR\_CONTROLS(person)
    - v. OTHER\_COMM
    - vi. RECRUITED
2. Meeting – lexical constraint
  - a. MEET
3. Joint work – lexical constraint
  - a. WORK\_WITH
  - b. SHARE\_TASK\_WITH
  - c. IS\_COWORKER\_OF
4. Economic/trade – usually lexical constraint
  - a. SEND\_TO
  - b. RECEIVE\_FROM
  - c. SELL\_TO
  - d. PURCHASE\_FROM
  - e. TRANSFER
  - f. IS\_EMPLOYER\_OF
5. Teacher/pupil – lexical constraint
  - a. TEACH
6. Directly described association – lexical constraint
  - a. IS\_AFFILIATED\_TO
7. Common group membership
  - a. MEMBER\_OF & MEMBER\_OF – no constraint
    - b. MEMBER\_OF & COMMANDS\_OR\_CONTROLS(org)
    - c. KINSHIP
8. Presence at shared location – no constraint
  - a. SHARE\_LOCATION
9. Common origin – no constraint
  - a. SHARE\_ORIGIN
10. Frequent travel to shared location – no constraint
  - a. TRAVEL\_TO

# Example Extraction

## 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

$$PartOf = \{\langle a, b \rangle, \langle c, d \rangle\}$$
$$OrgAff = \{\langle c, e \rangle\}$$
$$Serves = \{\langle a, f \rangle, \langle a, g \rangle, \langle a, h \rangle, \langle a, i \rangle\}$$

# Supervised Learning Approaches to Relation Analysis

---

Algorithm two step process

1. Identify whether pair of named entities are related
2. Classifier is trained to label relations

```
function FINDRELATIONS(words) returns relations
```

```
  relations  $\leftarrow$  nil
```

```
  entities  $\leftarrow$  FINDENTITIES(words)
```

```
  forall entity pairs  $\langle e1, e2 \rangle$  in entities do
```

```
    if RELATED?(e1, e2)
```

```
      relations  $\leftarrow$  relations + CLASSIFYRELATION(e1, e2)
```

# Factors used in Classifying

---

## Features of the named entities

- Named entity types of the two arguments
- Concatenation of the two entity types
- Headwords of the arguments
- Bag-of-words from each of the arguments
- Words in text
  - Bag-of-words and Bag-of-bigrams
  - Stemmed versions of the same
  - Distance between named entities (words / named entities)
- Syntactic structure
  - Parse related structures

# Sample features Extracted

## Entity-based features

Entity <sub>1</sub> type	ORG
Entity <sub>1</sub> head	<i>airlines</i>
Entity <sub>2</sub> type	PERS
Entity <sub>2</sub> 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 <sub>1</sub>	NONE
Word(s) after Entity <sub>2</sub>	<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	$Airlines \leftarrow_{subj} matched \leftarrow_{comp} said \rightarrow_{subj} Wagner$



# Semantic Role Labeling

---

- **SRL – is the task of finding semantic roles for each predicate.**

- **FrameNet**

**[You]            can't [blame]    [the program]    [for being unable to identify  
it]  
COGNIZER        TARGET    EVALUEE        REASON**

- **PropBank**

**[The San Francisco Examiner]    issued    [a special edition]  
[yesterday] ARG0                            TARGET    ARG1  
ARGM-TMP**

# Semantic Role Labeling Algorithm

---

- **Need syntactic parser.**
- **Extract features.**
- **Classify node.**

```
function SEMANTICROLELABEL(words) returns labeled tree
```

```
  parse  $\leftarrow$  PARSE(words)
```

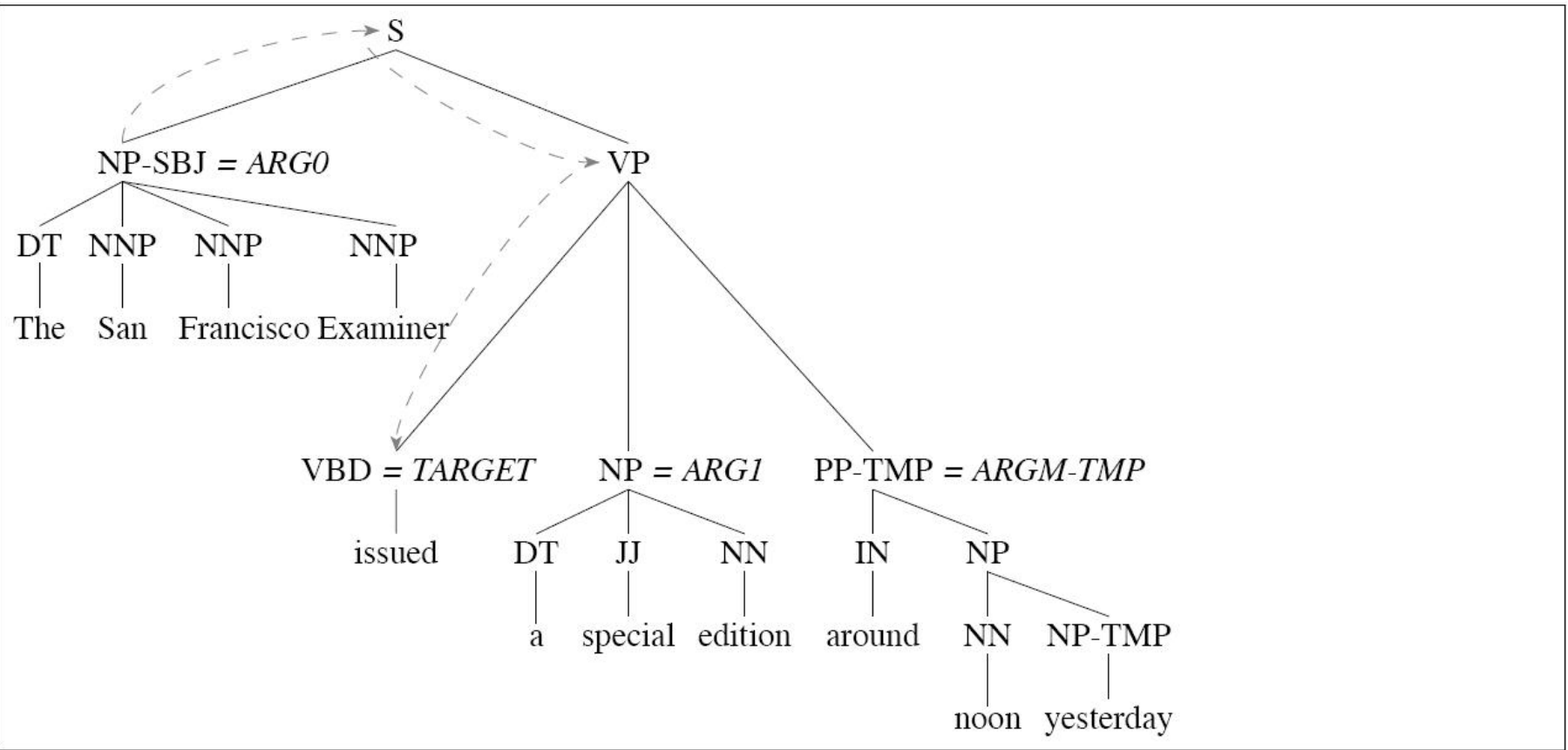
```
  for each predicate in parse do
```

```
    for each node in parse do
```

```
      featurevector  $\leftarrow$  EXTRACTFEATURES(node, predicate, parse)
```

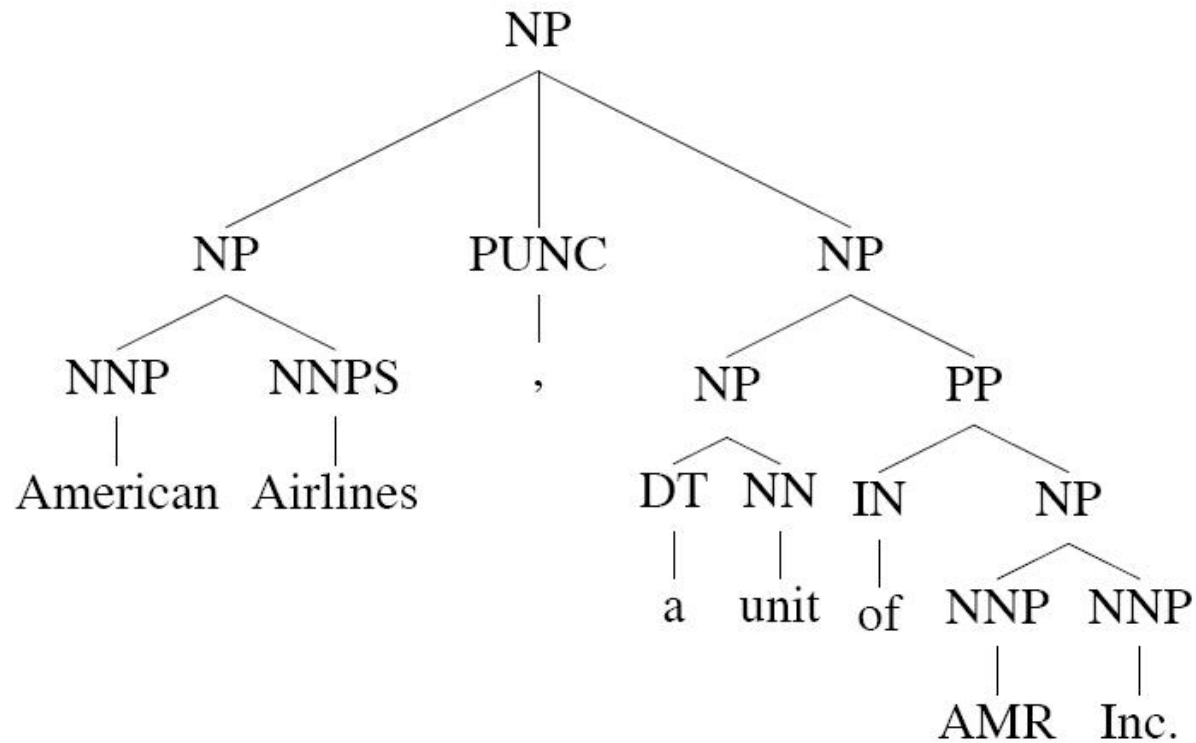
```
      CLASSIFYNODE(node, featurevector, parse)
```

# Semantic Role Labeling



# a-part-of relation

---



# Semantic Role Labeling-Features

---

- Governing Predicate.
- Phase type of constituent.
- Headword of constituent.
- Path in the parse tree from constituent to the predicate.
- Voice of the clause containing constituent.
- Binary respect to predicate (before or after).
- Sub categorization of predicate.

# Bootstrapping Example “Has a hub at”

---

Consider the pattern

/ \* has a hub at \* /

Google search

Milwaukee-based Midwest has a hub at KCI

Delta has a hub at LaGuardia

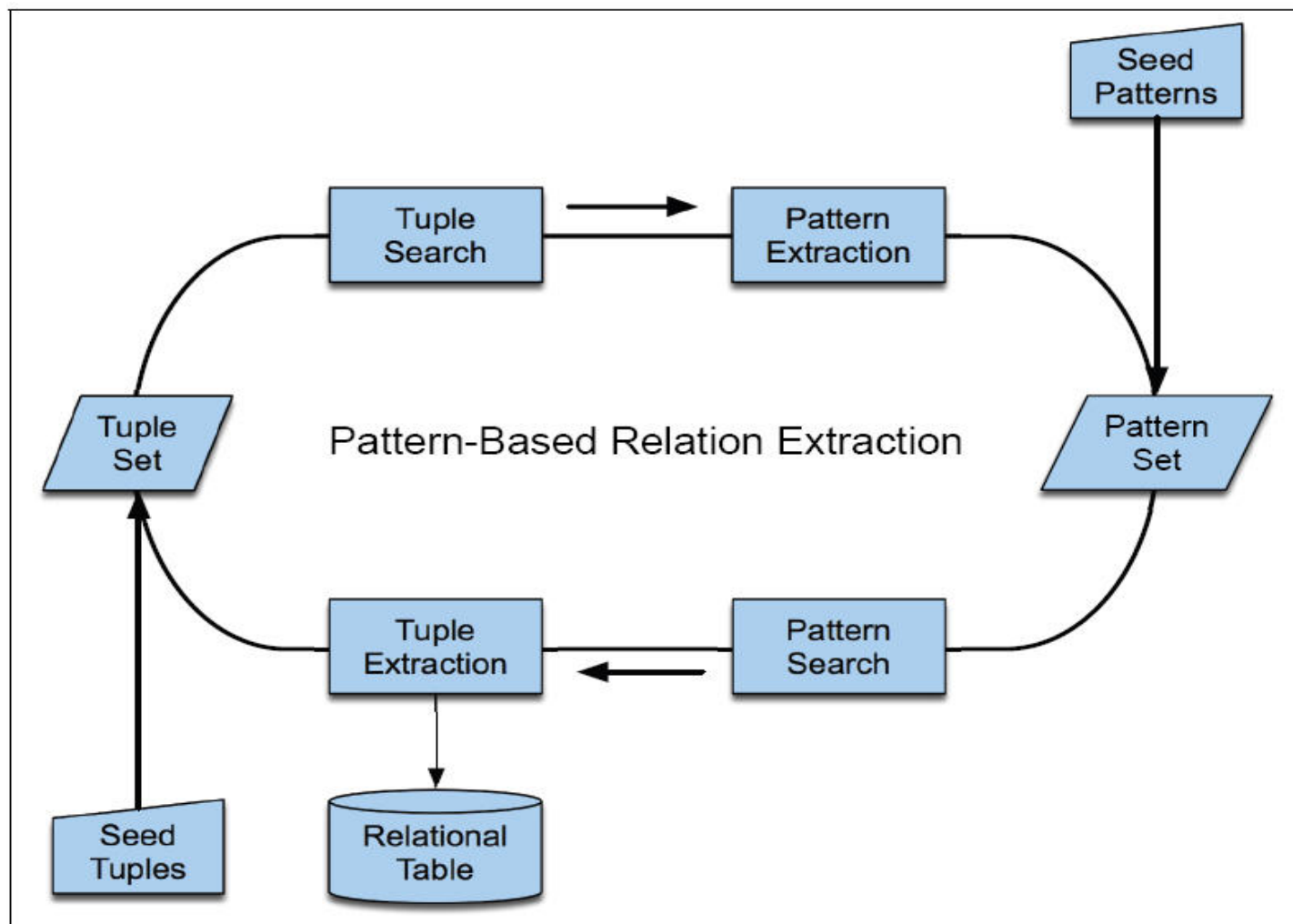
...

Two ways to fail

1. False positive: e.g. a star topology has a hub at its center
2. False negative? Just miss

No frill rival easyJet, which has established a hub at Liverpool

# Bootstrapping Relation Extraction



# Using Features to restrict patterns

---

Budget airline Ryanair, which uses Charleroi as a hub, scrapped all weekend flights ..

All flights in and out of Ryanair's Belgium hub at Charleroi airport were grounded on Friday ..

A spokesman at Charleroi, a main hub for Ryanair, estimated that ...

/ [ORG] , which uses a hub at [LOC] /

/ [ORG] ' s hub at [LOC] /

/ [LOC] a main hub for [ORG] /



# Temporal and Durational Expressions

---

Absolute temporal expressions

Relative 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

# Temporal lexical triggers

---

Category	Examples
Noun	<i>morning, noon, night, winter, dusk, dawn</i>
Proper Noun	<i>January, Monday, Ides, Easter, Rosh Hashana, Ramadan, Tet</i>
Adjective	<i>recent, past, annual, former</i>
Adverb	<i>hourly, daily, monthly, yearly</i>

# Temporal Normalization

---

iSO 8601 - standard for encoding temporal values

YYYY-MM-DD

# Sample ISO Patterns

---

Unit	Pattern	Sample Value
Fully specified dates	YYYY-MM-DD	1991-09-28
Weeks	YYYY-nnW	2007-27W
Weekends	PnWE	P1WE
24-hour clock times	HH:MM:SS	11:13:45
Dates and times	YYYY-MM-DDTHH:MM:SS	1991-09-28T11:00:00
Financial quarters	Qn	1999-Q3

# Event Detection and Analysis

---

## Event Detection and classification

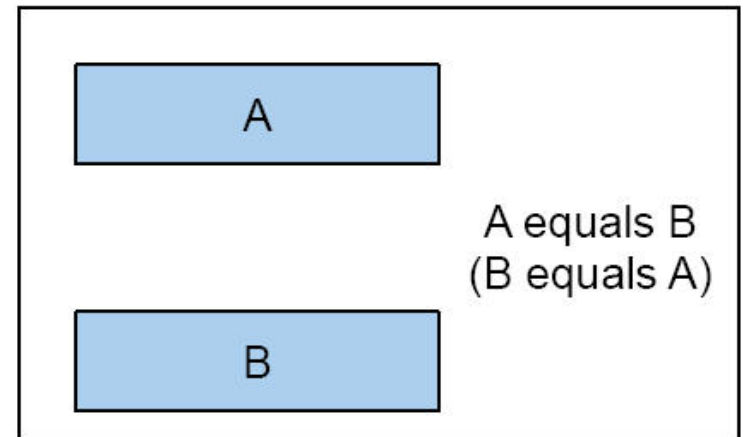
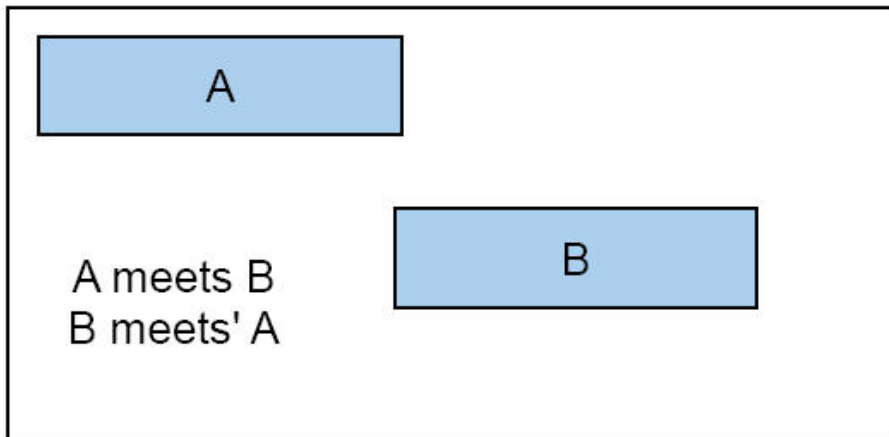
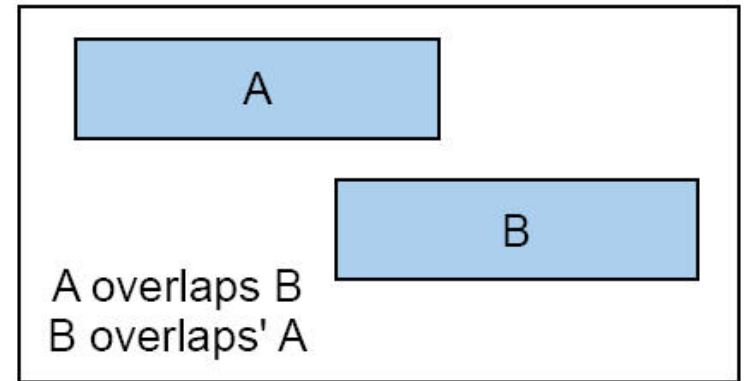
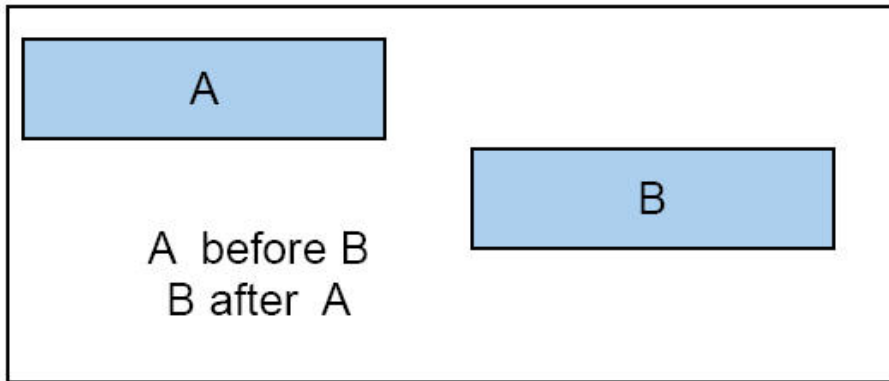
[*EVENT* Citing] high fuel prices, United Airlines [*EVENT* said] Friday it has [*EVENT* increased] fares by \$6 per round trip on flights to some cities also served by lower-cost carriers. American Airlines, a unit of AMR Corp., immediately [*EVENT* matched] [*EVENT* the move], spokesman Tim Wagner [*EVENT* said]. United, a unit of UAL Corp., [*EVENT* said] [*EVENT* the increase] took effect Thursday and [*EVENT* applies] to most routes where it [*EVENT* competes] against discount carriers, such as Chicago to Dallas and Denver to San Francisco.

# Features for Event Detection

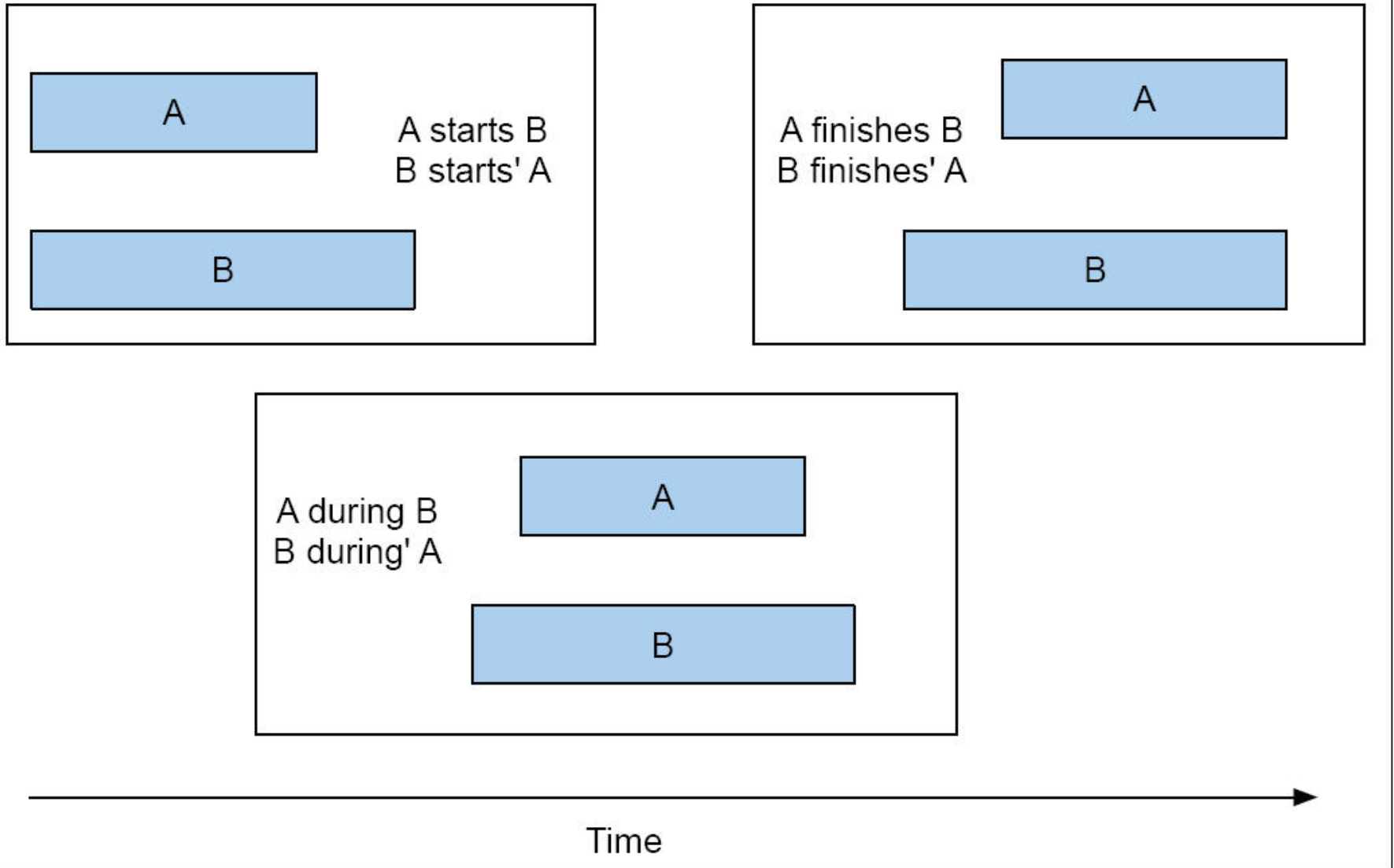
Features used in rule-based and statistical techniques

Feature	Explanation
Character affixes	Character-level prefixes and suffixes of target word
Nominalization suffix	Character level suffixes for nominalizations (e.g., <i>-tion</i> )
Part of speech	Part of speech of the target word
Light verb	Binary feature indicating that the target is governed by a light verb
Subject syntactic category	Syntactic category of the subject of the sentence
Morphological stem	Stemmed version of the target word
Verb root	Root form of the verb basis for a nominalization
WordNet hypernyms	Hypernym set for the target

# Allen's 13 temporal Relations



# continued





# Example from Timebank Corpus

---

```
<TIMEX3 tid="t57" type="DATE" value="1989-10-26" functionInDocument="CREATION_TIME">
10/26/89 </TIMEX3>
```

```
Delta Air Lines earnings <EVENT eid="e1" class="OCCURRENCE"> soared </EVENT> 33% to a
record in <TIMEX3 tid="t58" type="DATE" value="1989-Q1" anchorTimeID="t57"> the
fiscal first quarter </TIMEX3>, <EVENT eid="e3" class="OCCURRENCE"> bucking </EVENT>
the industry trend toward <EVENT eid="e4" class="OCCURRENCE"> declining </EVENT>
profits.
```