



# INF5820

## Natural Language Processing - NLP



H2009

Jan Tore Lønning

[jtl@ifi.uio.no](mailto:jtl@ifi.uio.no)



# NER and Relation detection & classification

INF5830

Lecture 14

Nov 9, 2009

# [ Today ]

- Overview: Information extraction
- (NP-)chunking
- Named entity recognition
- The Constraint Grammar approach
- Relation detection and classification

# Information extraction

- Goal: Extract structured data from text
- Track business news, or
- Intelligence news (possible terrorist attacks)
- Similarities to the frames from last week

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

# [ Steps ]

- The bottom-up approach:
  1. Preprocessing: tokenization, segmentation
  2. Tagging
  3. Chunking
  4. Named entity recognition
  5. Reference resolution
  6. Relation detection and classification
  7. Temporal analysis
  8. Template filling

# [ Steps ]

- The bottom-up approach:
  1. Preprocessing: tokenization, segmentation
  2. Tagging
  3. **Chunking**
  4. **Named entity recognition**
  5. Reference resolution
  6. **Relation detection and classification**
  7. Temporal analysis
  8. Template filling

# [ Chunking ]

- Form of shallow parsing
- Flat structures
- Identify (some) phrases

[*NP* The morning flight] [*PP* from] [*NP* Denver] [*VP* has arrived.]

[*NP* a flight] [*PP* from] [*NP* Indianapolis][*PP* to][*NP* Houston][*PP* on][*NP* TWA].

- Non-overlapping phrases
- Compromises, cf. PP
- Sometimes only interested in NPs

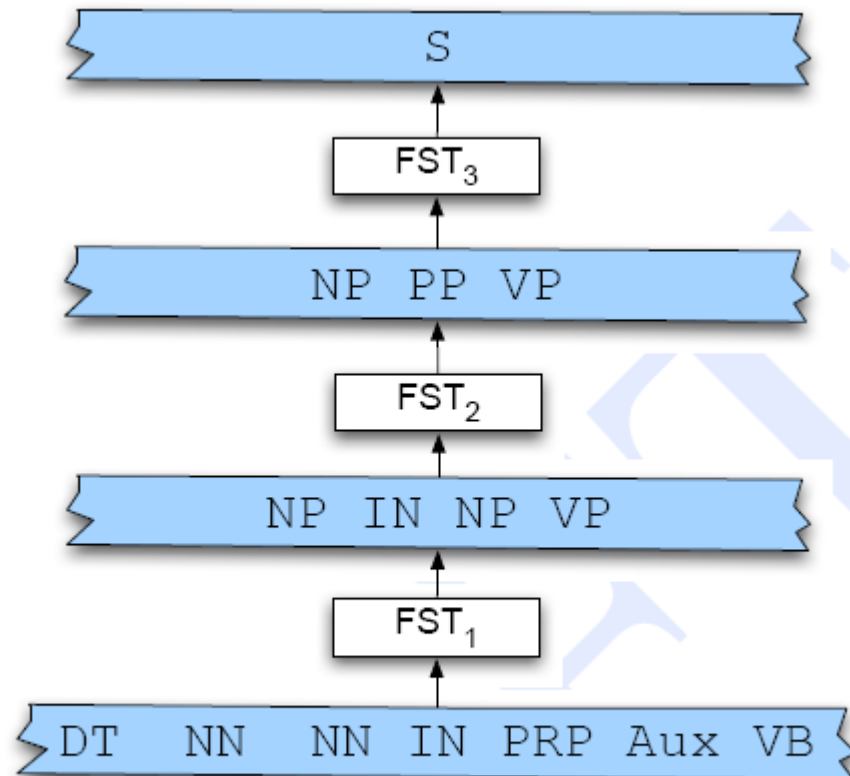
# Approach 1: Cascaded FSTs

*NP* → (*Det*) *Noun*\* *Noun*

*NP* → *Proper-Noun*

*VP* → *Verb*

*VP* → *Aux Verb*



The morning flight from Denver has arrived



# [ 2: ML-approaches ]

- Two tasks:
  - Identify the phrase (beginning-end)
  - Classify the phrase

# ML-approaches continued

- The two steps as a tagging task:

*The morning flight from Denver has arrived*  
B\_NP I\_NP I\_NP B\_PP B\_NP B\_VP I\_VP

The same sentence with only the base-NPs tagged

*The morning flight from Denver has arrived.*  
B\_NP I\_NP I\_NP O B\_NP O O

B\_NP : begin NP

I\_NP : inside NP

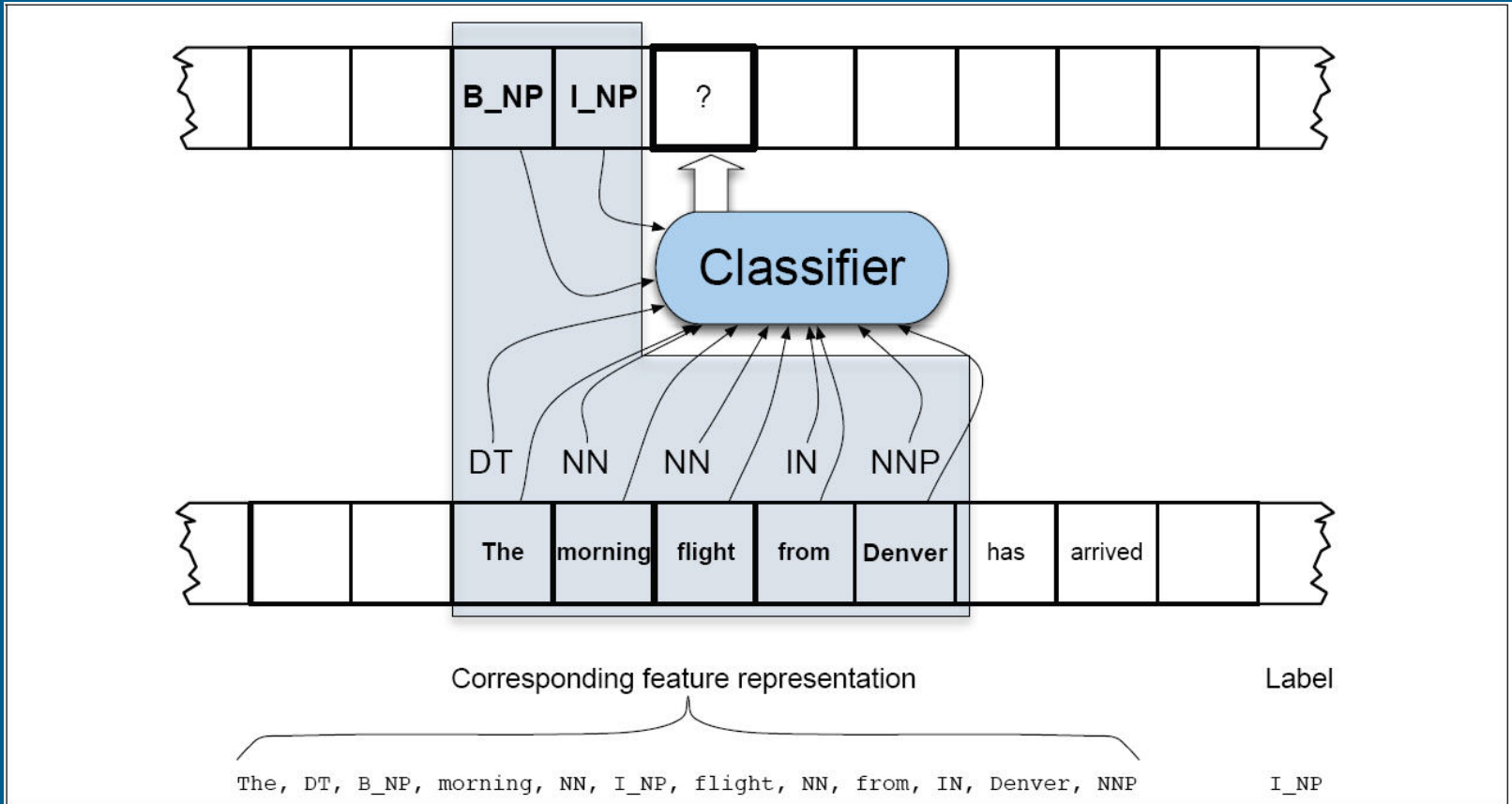
B\_VP : begin VP

I\_VP: inside VP

O: not part of a phrase

Etc.

# ML: Classifier



# [ ML: more ]

- Training data from Penn treebank
- Evaluation on found chunks
  - compared to test set
  - Recall and precision
    - (typo in hardcover book)

# [ Today ]

- Overview: Information extraction
- (NP-)chunking
- **Named entity recognition**
- The Constraint Grammar approach
- Relation detection and classification

# Named Entity Classes

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.

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

- Choice of types is/should be application specific

# [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

[*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...

# Named entity recognition

- Two tasks:
  - Identify the phrase (beginning-end)
  - Classify the phrase
- Similar to chunking but
  - Different/more fine-grained classification

Words	Label
American	<i>B<sub>ORG</sub></i>
Airlines	<i>I<sub>ORG</sub></i>
,	<i>O</i>
a	<i>O</i>
unit	<i>O</i>
of	<i>O</i>
AMR	<i>B<sub>ORG</sub></i>
Corp.	<i>I<sub>ORG</sub></i>
,	<i>O</i>
immediately	<i>O</i>
matched	<i>O</i>
the	<i>O</i>
move	<i>O</i>
,	<i>O</i>
spokesman	<i>O</i>
Tim	<i>B<sub>PERS</sub></i>
Wagner	<i>I<sub>PERS</sub></i>
said	<i>O</i>
.	<i>O</i>



# Features

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

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

# Training data

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

# [ Today ]

- Overview: Information extraction
- (NP-)chunking
- Named entity recognition
- **The Constraint Grammar approach**
- Relation detection and classification

# [ CG-approach to NER ]

- See OB-tagger
  - (text example: dn)
- Same approach as to tagging:
  - Start with all possible classes
  - Write rules which remove alternatives
  - May end with more than one answer

# [ Today ]

- Overview: Information extraction
- (NP-)chunking
- Named entity recognition
- The Constraint Grammar approach
- **Relation detection and classification**

# Relation detection and classification

- Two steps
  - Detection: **Is there a relation between two entities?**
  - Classification: **What kind of relation?**

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 [PERS Tim Wagner] said. [ORG United], a 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].

# As logical relations

## 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

- Corpus marked with NEs and relations

## 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

Features, examples

## 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	<i>NP ↑ NP ↑ S ↑ S ↓ NP</i>
Base syntactic chunk path	<i>NP → NP → PP → NP → VP → NP → NP</i>
Typed-dependency path	<i>Airlines ←<sub>subj</sub> matched ←<sub>comp</sub> said →<sub>subj</sub> Wagner</i>



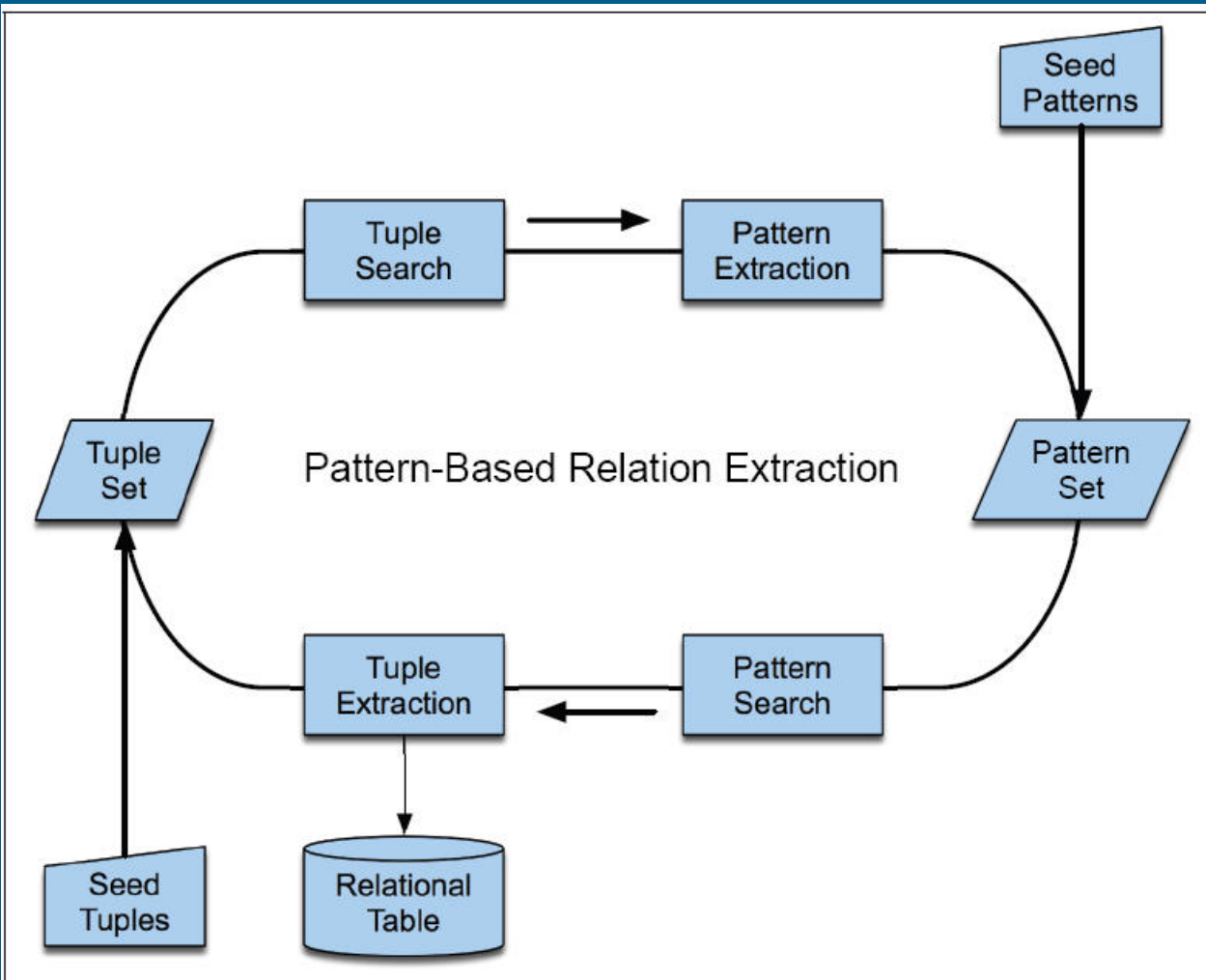
# [ Pattern-matching ]

1. Choose a pattern, e.g.
  - \* has a hub at +
2. Find pairs in the construction, e.g.
  - Milwaukee-based Midwest has a hub at KCI
  - Bulgaria Air has a hub at Sofia Airport
3. Extend/refine pattern, e.g.
  - [ORG] has a ADJ\* hub at [LOC]

# [ Bootstrapping ]

1. Choose a pattern
2. Find pairs in the construction
3. Find other occurrences of these pairs
4. Extract patterns from this, e.g.
  - [ORG] which uses [LOC] as hub
  - [ORG]'s hub at [LOC]
  - [LOC] a ADJ\* hum for [ORG]
5. Repeat from (2)

# Bootstrapping loop



# [ Beware ]

- Check that bootstrapping does not drift away
  - Control against original tuples
  - Skip details
- Evaluate, either
  - Count occurrences of relations in sentences in corpus and evaluate against them
  - Count whether relationship is entered into data base

# Information extraction

- The information extraction approach shaped by
  - MUC: message understanding conferences
  - Competition: make the best system
  - 1987-97
  - DARPA
- See
  - [Wikipedia](#)
  - Bibliographical and historical notes J&M, 22
- Overlap with other tasks/approaches
  - But comparison not always immediate