

# CS674 Natural Language Processing

- **Partial parsing / Chunking**

- Task specification
- Error-driven pruning of Treebank grammars
- Comparison with TBL

## Partial parsing

When it's time for their biannual powwow, the nation's manufacturing titans typically jet off to the sunny confines of resort towns like Boca Raton and Hot Springs.

Partial Parser

When [<sub>S</sub> [<sub>NP</sub> it ] ] [<sub>V</sub> 's ] [<sub>Obj</sub> [<sub>NP</sub> time ] ] for [<sub>NP</sub> their biannual powwow ] , [<sub>NP</sub> the nation ] 's [<sub>S</sub> [<sub>NP</sub> manufacturing titans ] ] typically [<sub>V</sub> jet off ] to [<sub>NP</sub> the sunny confines ] of [<sub>NP</sub> resort towns ] like [<sub>NP</sub> Boca Raton ] and [<sub>NP</sub> Hot Springs ] .

## Why partial parsing?

- **Fast**
- **Supports a number of large-scale NLP tasks**
  - Information Extraction
  - Phrase identification for Information Retrieval
  - Question Answering

## Base noun phrases

### Non-recursive noun phrases (smallest NPs)

When [it] 's [time] for [their biannual powwow] , [the nation] 's [manufacturing titans] typically jet off to [the sunny confines] of [resort towns] like [Boca Raton] and [Hot Springs] .

## Inductive ML algorithm

- **Simple**

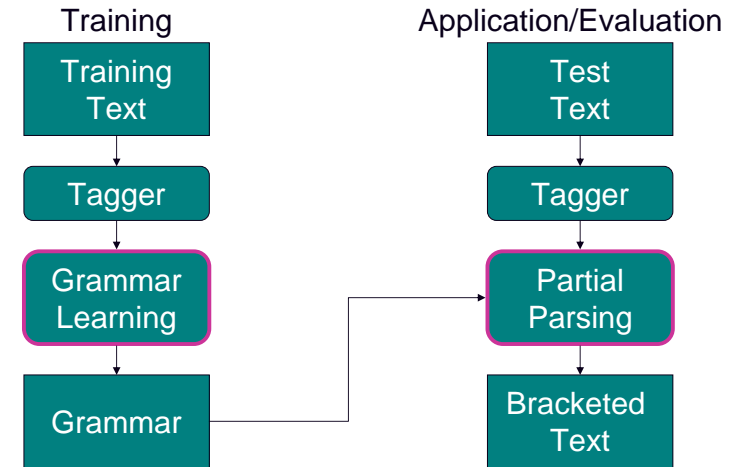
base NP = any string having the same part-of-speech tag sequence as a base NP from the training corpus

- **Combines components of existing techniques**

- Charniak (1996)
- Brill (1995)

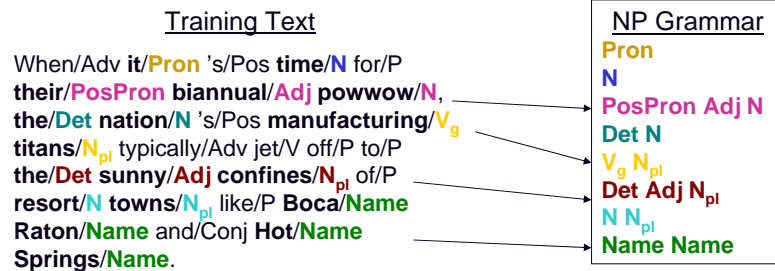
- **Achieves surprisingly high accuracies**

## Partial parsing framework



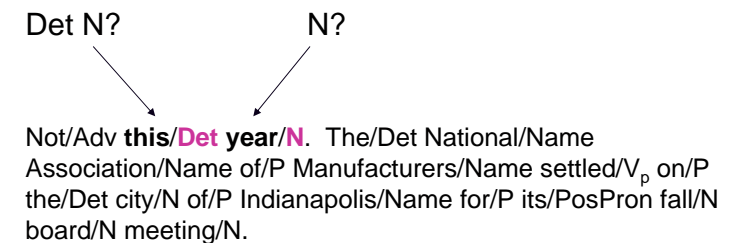
## Rule extraction

**rule = sequence of part-of-speech tags**



## Partial parsing bracketer

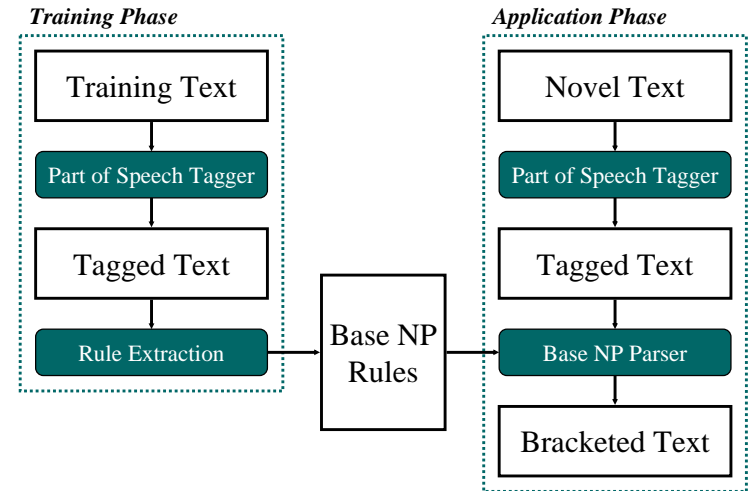
- Left-to-right
- Longest-match



## Parser (bracketer)

Bracket( $w_1, \dots, w_n$ ):  
assign p-o-s tags  $t_1, \dots, t_n$  to words  $w_1, \dots, w_n$   
 $i = 1$   
while  $i \leq n$  do  
   $\{r_1, \dots, r_k\} = \text{Matches}(w_i, \dots, w_n)$   
   $r = \text{longest}(r_1, \dots, r_k)$   
  make new NP from  $w_i, \dots, w_{i+|r|-1}$   
   $i = i + |r|$

## Overview of the method



## Poorly performing rules

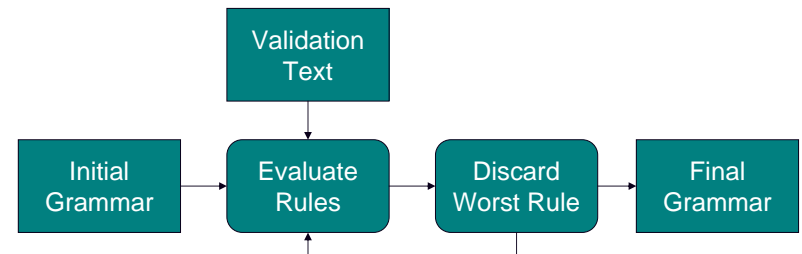
- **Sources of bad rules**

- errors in training data
- errors in part-of-speech tagging
- irregular & ambiguous constructs

...manufacturing/V<sub>g</sub> titans/N<sub>pl</sub>...

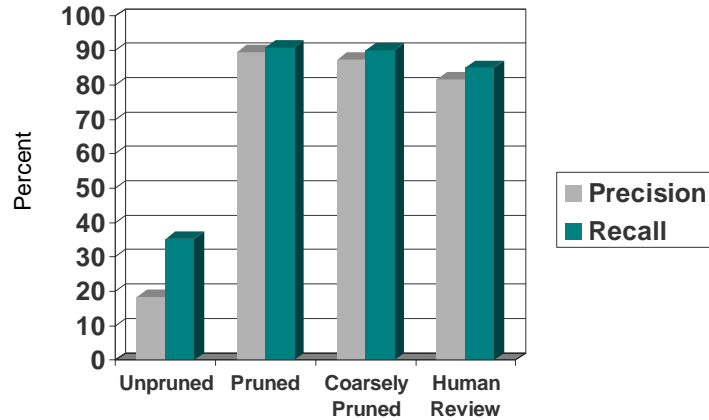
...the/Det executives/N<sub>pl</sub> began/V<sub>p</sub> boarding/V<sub>g</sub> buses/N<sub>pl</sub>...

## Grammar pruning



- $\text{score}(r) = \text{correct}(r) - \text{errors}(r)$
- stop when worst score is positive

## Results



## Results vs. TBL on R&M corpus

	TBL results	Pierce & Cardie [98]	Difference
w/lexical templates	93.1P/93.5R		-3.7P/-2.6R
w/o lexical templates	90.5P/90.7R	89.4P/90.9R	-0.9P/+0.2R

## Advantages of the approach

- **Good performance**
- **Simple**
  - Easy to understand, implement
  - Produces intelligible grammar rules
  - Easy to update for new text genre
- **Efficient**
  - Fastest bracketing procedure
- **State of the art circa 2003**
  - ~94% P/R for NP, VP, PP chunks
  - Using ensembles of SVM's (Kudo & Matsumoto, 2000) and Winnow as employed in Zhang et al. (2001)