# CSC 594 Topics in AI – Applied Natural Language Processing

Fall 2009/2010

4. Grammar and Parsing

## Lexical Categories: Parts-of-Speech (POS)

•	<ul> <li>8 (ish) traditional parts of speech</li> <li>Noun, verb, adjective, preposition, adverb, article, interjection pronoun, conjunction, etc.</li> </ul>							
•	Ν	noun	chair, bandwidth, pacing					
•	V	verb	study, debate, munch					
•	ADJ	adjective	purple, tall, ridiculous					
•	ADV	adverb	unfortunately, slowly					
•	Р	preposition	of, by, to					
•	PRO	pronoun	I, me, mine					
•	DET	determiner	the, a, that, those					

2

## Language Structure and Meaning We want to know how meaning is mapped onto what language structures. Commonly in English in ways like this: (THING The dog) is (PROPERTY fierce) (ACTION (THING The dog) is chasing (THING the cat]) (STATE (THING The dog) was sitting (PLACE in the garden) (TIME yesterday]) (ACTION (THING We) ran (PATH OUT into the water]) (ACTION (THING The dog) barked (PROPERTY/MAINER loudly]] (ACTION (THING The dog) barked (PROPERTY/MAINER loudly]] (ACTION (THING The dog) barked (PROPERTY/MAINER loudly]]

## Constituency

• Sentences have parts, some of which appear to have subparts. These groupings of words that go together we will call **constituents**.

I hit the man with a cleaver I hit [the man with a cleaver] I hit [the man] with a cleaver

You could not go to her party You [could not] go to her party You could [not go] to her party

Source: Andrew McCallum, UMass Amherst

## **Constituent Phrases**

• For constituents, we usually name them as phrases based on the word that **heads** the constituent:

- Noun phrase e.g. [the man], [the man [with a cleaver]]
- Verb phrase e.g. [hit the man], [hit the man with a cleaver]
   Adjective phrase e.g. [extremely significant]
- Prepositional phrase e.g. [with a cleaver], [in my room]

## Grammar

- A formal specification of how a constituent should be made of and combined with other constituents (to form a larger constituent).
- Unlike programming languages (which are formal languages), natural languages are born without formal specifications.
- So, grammars for natural languages are empirically derived by observation, which is also subjective.
   → There are many different grammars for a given language.



## **Context-Free Grammar (CFG)**

- The most common way of specifying natural language grammars (because most natural languages are context-free in the generative capacity).
- Formally, a CFG consists of:
  - N a set of non-terminal symbols (or variables)
  - $\Sigma$  a set of **terminal symbols** (disjoint from *N*)
  - *R* a set of **rules** or productions, each of the form  $A \rightarrow \beta$ , where *A* is a non-terminal,
  - $\beta$  is a string of symbols from the infinite set of strings  $(\Sigma \cup N)*$
  - S a designated start symbol

8

## Grammaticality

- A CFG defines a formal language = the set of all sentences (strings of words) that can be derived by the grammar.
- · Sentences in this set said to be grammatical.
- Sentences outside this set said to be ungrammatical.

Source: Andrew McCallum, UMass Amherst













## Agreement

- Agreement is a constraint that holds between constituents to make the sentence grammatical.
- For example, in English, determiners and the head nouns in NPs have to agree in their number.
  - This flight
  - (\*) This flights
  - Those flights
  - (\*) Those flight
- Another example: the head nouns in the subject NPs must agree with the person of the main verb (if 3<sup>rd</sup> person singular, present tense).
  - He thinks ..
  - (\*) He think

14

15

## Problem

- Our earlier NP rules are clearly deficient since they don't capture the agreement constraints
  - NP → Det Nom
    - Accepts, and assigns correct structures, to grammatical examples
       (*this flight*)
  - But its also happy with incorrect examples (\*these flight)
  - Such a rule is said to overgenerate → causes ambiguity

## Verb Subcategorization

- · English VPs consist of a head verb along with 0 or more following constituents which we'll call arguments.
  - Find: Please find [a flight to NY]<sub>NP</sub>
  - Give: Give [me]<sub>NP</sub>[a cheaper fare]<sub>NP</sub>
  - Help: Can you help [me]<sub>NP</sub>[with a flight]<sub>PP</sub>
  - Prefer: I prefer [to leave earlier]<sub>TO-VP</sub>
  - Told: I was told [United has a flight]<sub>s</sub>
  - (\*) John sneezed the book
  - (\*) I prefer United has a flight
  - (\*) Give with a flight

Source: Jurafsky & Martin "Speech and Language Processing"

## **Possible CFG Solution** SgS -> SgNP SgVP Possible solution for

agreement. Can use the same trick for all the verb/VP

Source: Jurafsky & Martin "Speech and Language Processing"

classes.

PIS -> PINp PIVP

16

17

18

- SgNP -> SgDet SgNom PINP -> PIDet PINom
- PIVP -> PIV NP
- SgVP ->SgV Np
- ...

# **CFG Solution for Agreement**

- · It works and stays within the power of CFGs
- · But its ugly
- And it doesn't scale all that well because of the interaction among the various constraints explodes the number of rules in our grammar.

Source: Jurafsky & Martin "Speech and Language Processing"

## **The Point**

- CFGs appear to be just about what we need to account for a lot of basic syntactic structure in English.
- But there are problems

   That can be dealt with adequately, although not elegantly, by staying within the CFG framework.
- There are simpler, more elegant, solutions that take us

19

20

- out of the CFG framework (beyond its formal power) - Lexical Functional Grammar (LFG)
- Head-driven Phrase Structure Grammar (HPSG)
- Construction grammar
- X-Tree Adjoining Grammar (XTAG)
- Unification Grammar
- etc.

Source: Jurafsky & Martin "Speech and Language Processing"

## **Treebanks**

• Treebanks are corpora in which each sentence has been paired with a parse tree.

- · These are generally created
  - By first parsing the collection with an automatic parser
     And then having human annotators correct each parse as
  - And then having human annotators correct each parse as necessary.
- This generally requires detailed annotation guidelines that provide a POS tagset, a grammar and instructions for how to deal with particular grammatical constructions.

Source: Jurafsky & Martin "Speech and Language Processing"







- Treebanks implicitly define a grammar for the language covered in the treebank.
- Not complete, but if you have decent size corpus, you'll have a grammar with decent coverage.
- Grammar rules tend to be 'flat' to avoid recursive rules.
- For example, the Penn Treebank has 4500 different rules for VPs. Among them...

$$\begin{array}{rcl} VP & \rightarrow & VBD & PP \\ VP & \rightarrow & VBD & PP & PP \\ VP & \rightarrow & VBD & PP & PP & PP \\ VP & \rightarrow & VBD & PP & PP & PP & PP \end{array}$$

Source: Jurafsky & Martin "Speech and Language Processing"





Source: Jurafsky & Martin "Speech and Language Processing"

23

22









## **Dependency Grammars**

- In CFG-style phrase-structure grammars the main focus is on *constituents*.
- But it turns out you can get a lot done with just binary relations among the words in an utterance.
- In a dependency grammar framework, a parse is a tree where
  - the nodes stand for the words in an utterance
  - The links between the words represent dependency relations between pairs of words.
     Relations may be typed (labeled), or not.

Source: Jurafsky & Martin "Speech and Language Processing"

Source: Jurafsky & Martin "Speech and Language Processing"

# **Dependency Relations**

Argument Dependencies	Description
nsubj	nominal subject
csubj	clausal subject
dobj	direct object
iobj	indirect object
pobj	object of preposition
Modifier Dependencies	Description
tmod	temporal modifier
appos	appositional modifier
det	determiner
prep	prepositional modifier

27

26





# **Dependency Parsing**

• The dependency approach has a number of advantages over full phrase-structure parsing.

- Deals well with free word order languages where the constituent structure is quite fluid
- Parsing is much faster than CFG-bases parsers
- Dependency structure often captures the syntactic relations needed by later applications
  - CFG-based approaches often extract this same information from trees anyway.

29

### Source: Jurafsky & Martin "Speech and Language Processing"



## **Parsing Algorithms**

- Top-down Parsing -- (top-down) derivation
- Bottom-up Parsing
- Chart Parsing
- Earley's Algorithm most efficient, O(n<sup>3</sup>)
- Left-corner Parsing optimization of Earley's

31

· and lots more...

	upottom)،تر"	ohn ate the		ing
Grammar $S \rightarrow NP VP$ $NP \rightarrow Det N$	(11) reduce $(0,4,S \rightarrow NP VP \bullet)$	± 4		
$VP \rightarrow VG NP$ $VG \rightarrow V$ $NP \rightarrow "John"$		(10) reduce $\langle l, 4, VP \rightarrow VG NP \bullet \rangle$		
$V \rightarrow$ "ate" Det $\rightarrow$ "the" $N \rightarrow$ "cake"		(5) shift 2 $\langle l, 2, VP \rightarrow VG \bullet NP \rangle$	(9) reduce $\langle 2, 4, NP \rightarrow Det N \bullet \rangle$	
	(2) shift 2 $(0,1,S \rightarrow NP \bullet VP)$	(4) shift 2 $\langle l, 2, VG \rightarrow V \bullet \rangle$	(7) shift 2 $\langle 2,3, NP \rightarrow Det \bullet N \rangle$	
	(1) shift 1 "John" $(0,1, NP \rightarrow "John" \bullet)$	(3) shift 1 "ate" $\langle l,2, V \rightarrow$ "ate" • $\rangle$	(6) shift 1 "the" $\langle 2,3, \text{Det} \rightarrow$ " the" • $\rangle$	(8) shift 1 "cake" $\langle 3,4, N \rightarrow$ "cake" • $\rangle$
	0	1 2	2	3







# **Probabilistic Parsing**

- For ambiguous sentences, we'd like to know which parse tree is more likely than others.
- So we must assign probability to each parse tree ... but how?
- A probability of a parse tree *t* is

 $p(t) = \sum p(r)$  where *r* is a rule used in *t*.

and p(r) is obtained from a (annotated) corpus.

35

# **Partial Parsing**

- Parsing fails when the coverage of the grammar is not complete – but it's almost impossible to write out all legal syntax (without accepting ungrammatical sentences).
- We'd like to at least get pieces even when full parsing fails.
- Why not abandon full parsing and aim for partial parsing from the start...

Semantic Analysis (1)	
<ul> <li>Derive the meaning of a sentence.</li> <li>Often applied on the result of syntactic analysis.         <ul> <li><u>"John ate the cake."</u> NP V NP                 ((action INGEST) ; syntactic verb                 (actorJOHN-01) ; syntactic subj                 (object FOOD)) ; syntactic obj</li> </ul> </li> <li>To do semantic analysis, we need a (semantic)         dictionary (e.g. WordNet,         <ul> <li><u>http://www.cogsci.princeton.edu/~wn/</u>).</li> </ul> </li> </ul>	
	37

# Semantic Analysis (2)

Semantics is a double-edged sword...

- Can resolve syntactic ambiguity
  - "I saw a man on the hill with a telescope"
  - "I saw a man on the hill with a hat"
- But introduces semantic ambiguity
  - "She walked towards the bank"
- But in human sentence processing, we seem to resolve both types of ambiguities simultaneously (and in linear time)...

38

