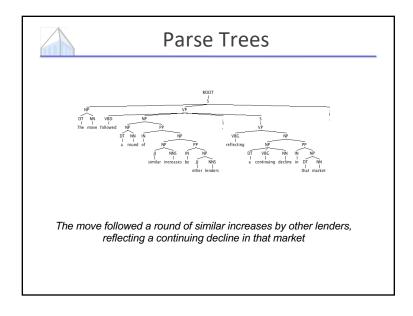
# **Natural Language Processing**

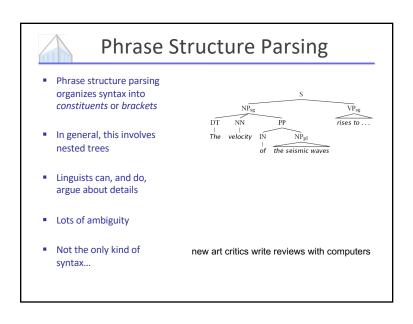


Syntax and Parsing

Dan Klein – UC Berkeley

# **Syntax**

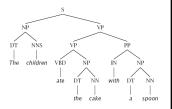






## **Constituency Tests**

- How do we know what nodes go in the tree?
- Classic constituency tests:
  - Substitution by proform
  - Question answers
  - Semantic gounds
    - Coherence
    - Reference
    - Idioms
  - Dislocation
  - Conjunction
- Cross-linguistic arguments, too



# **Conflicting Tests**

- Constituency isn't always clear
  - Units of transfer:
    - think about ~ penser à
    - talk about ~ hablar de
  - Phonological reduction:
    - I will go  $\rightarrow$  I'll go
    - I want to go → I wanna go
    - a le centre → au centre



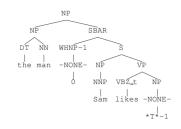
La vélocité des ondes sismiques

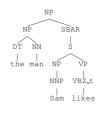
- Coordination
  - He went to and came from the store.



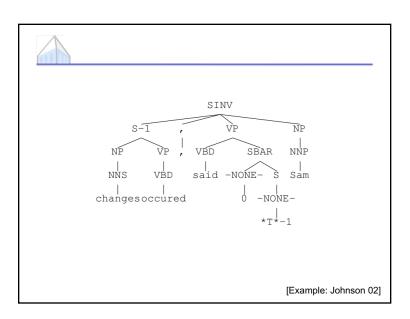
### Questions from Last Time

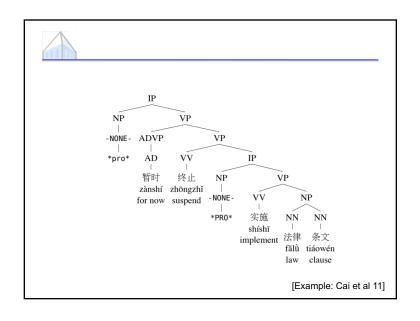
Q: Do we model deep vs surface structure?



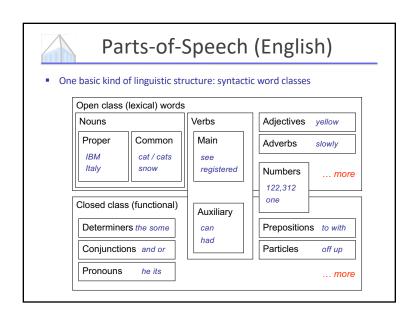


[Example: Johnson 02]





# **Ambiguities**





# Part-of-Speech Ambiguity

Words can have multiple parts of speech

VBD VB VBV VBZ VBP VBZ NNP NNS NN NNS CD NN Fed raises interest rates 0.5 percent

Mrs./NNP Shaefer/NNP never/RB got/VBD **around/RP** to/TO joining/VBG
All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB **around/IN** the/DT corner/NN
Chateau/NNP Petrus/NNP costs/VBZ **around/RB** 250/CD

- Two basic sources of constraint:
  - Grammatical environment
  - Identity of the current word
- Many more possible features:
  - Suffixes, capitalization, name databases (gazetteers), etc...



## Why POS Tagging?

- Useful in and of itself (more than you'd think)
  - Text-to-speech: record, lead
  - Lemmatization:  $saw[v] \rightarrow see$ ,  $saw[n] \rightarrow saw$
  - Quick-and-dirty NP-chunk detection: grep {JJ | NN}\* {NN | NNS}
- Useful as a pre-processing step for parsing
  - Less tag ambiguity means fewer parses
  - However, some tag choices are better decided by parsers

IN

DT NNP NN VBD VBN RP NN NNS
The Georgia branch had taken on loan commitments ...

VDN

DT NN IN NN VBD NNS VBD
The average of interbank offered rates plummeted ...



## Classical NLP: Parsing

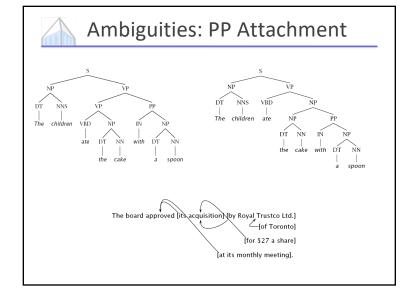
Write symbolic or logical rules:

Grammar (CFG)

Lexicon

...

- Use deduction systems to prove parses from words
  - Minimal grammar on "Fed raises" sentence: 36 parses
  - Simple 10-rule grammar: 592 parses
  - Real-size grammar: many millions of parses
- This scaled very badly, didn't yield broad-coverage tools





#### **Attachments**

- I cleaned the dishes from dinner
- I cleaned the dishes with detergent
- I cleaned the dishes in my pajamas
- I cleaned the dishes in the sink



## Syntactic Ambiguities I

- Prepositional phrases:
   They cooked the beans in the pot on the stove with handles.
- Particle vs. preposition: The puppy tore up the staircase.
- Complement structures
   The tourists objected to the guide that they couldn't hear.
   She knows you like the back of her hand.
- Gerund vs. participial adjective
   Visiting relatives can be boring.
   Changing schedules frequently confused passengers.



## Syntactic Ambiguities II

- Modifier scope within NPs impractical design requirements plastic cup holder
- Multiple gap constructions
   The chicken is ready to eat.

   The contractors are rich enough to sue.
- Coordination scope:
   Small rats and mice can squeeze into holes or cracks in the wall.

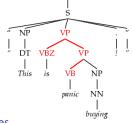


# **Dark Ambiguities**

 Dark ambiguities: most analyses are shockingly bad (meaning, they don't have an interpretation you can get your mind around)

This analysis corresponds to the correct parse of

"This will panic buyers!"



- Unknown words and new usages
- Solution: We need mechanisms to focus attention on the best ones, probabilistic techniques do this

**PCFGs** 



#### Probabilistic Context-Free Grammars

- A context-free grammar is a tuple <*N*, *T*, *S*, *R*>
  - N: the set of non-terminals
    - Phrasal categories: S, NP, VP, ADJP, etc.
    - Parts-of-speech (pre-terminals): NN, JJ, DT, VB
  - T: the set of terminals (the words)
  - S: the start symbol
    - Often written as ROOT or TOP
    - Not usually the sentence non-terminal S
  - R: the set of rules
    - Of the form  $X \rightarrow Y_1 Y_2 \dots Y_k$ , with  $X, Y_i \in N$
    - Examples: S → NP VP, VP → VP CC VP
    - Also called rewrites, productions, or local trees

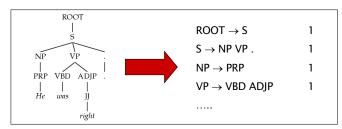
#### A PCFG adds:

■ A top-down production probability per rule P(Y<sub>1</sub> Y<sub>2</sub> ... Y<sub>k</sub> | X)



#### **Treebank Grammars**

- Need a PCFG for broad coverage parsing.
- Can take a grammar right off the trees (doesn't work well):



- Better results by enriching the grammar (e.g., lexicalization).
- Can also get state-of-the-art parsers without lexicalization.



#### **Treebank Sentences**



#### Treebank Grammar Scale

- Treebank grammars can be enormous
  - As FSAs, the raw grammar has ~10K states, excluding the lexicon
  - Better parsers usually make the grammars larger, not smaller

NN



# **Chomsky Normal Form**

- Chomsky normal form:
  - All rules of the form  $X \rightarrow Y Z$  or  $X \rightarrow w$
  - In principle, this is no limitation on the space of (P)CFGs
    - N-ary rules introduce new non-terminals



- Unaries / empties are "promoted"
- In practice it's kind of a pain:
  - Reconstructing n-aries is easy
  - · Reconstructing unaries is trickier
  - The straightforward transformations don't preserve tree scores
- Makes parsing algorithms simpler!

# **CKY Parsing**



#### A Recursive Parser

- Will this parser work?
- Why or why not?
- Memory requirements?



#### A Memoized Parser

One small change:



# A Bottom-Up Parser (CKY)

Can also organize things bottom-up



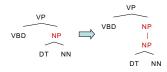
# **Unary Rules**

Unary rules?



## CNF + Unary Closure

- We need unaries to be non-cyclic
  - Can address by pre-calculating the unary closure
  - Rather than having zero or more unaries, always have exactly one



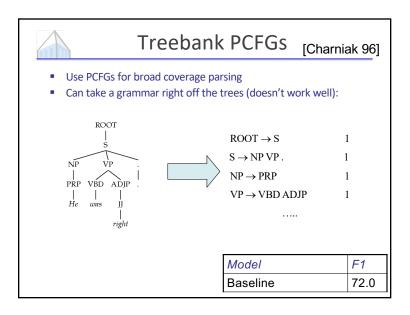


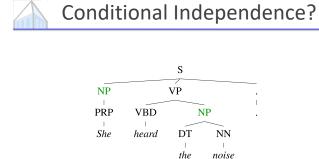
- Alternate unary and binary layers
- Reconstruct unary chains afterwards



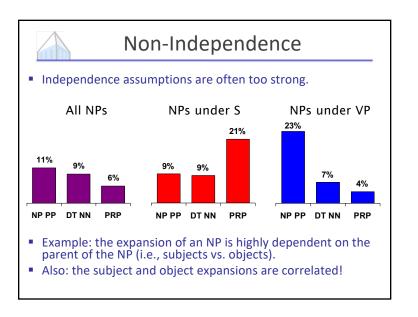
## Alternating Layers

# **Learning PCFGs**





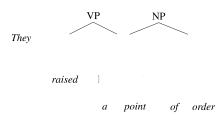
- Not every NP expansion can fill every NP slot
  - A grammar with symbols like "NP" won't be context-free
  - Statistically, conditional independence too strong





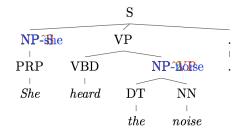
### **Grammar Refinement**

■ Example: PP attachment





### **Grammar Refinement**

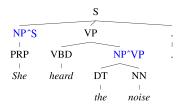


- Structure Annotation [Johnson '98, Klein&Manning '03]
- Lexicalization [Collins '99, Charniak '00]
- Latent Variables [Matsuzaki et al. 05, Petrov et al. '06]

### **Structural Annotation**



# The Game of Designing a Grammar

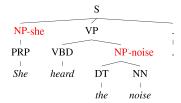


- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Structural annotation

## Lexicalization



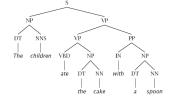
# The Game of Designing a Grammar

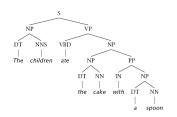


- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Structural annotation [Johnson '98, Klein and Manning 03]
  - Head lexicalization [Collins '99, Charniak '00]



#### Problems with PCFGs

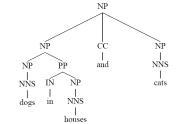


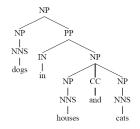


- If we do no annotation, these trees differ only in one rule:
  - VP → VP PP
  - NP → NP PP
- Parse will go one way or the other, regardless of words
- We addressed this in one way with unlexicalized grammars (how?)
- Lexicalization allows us to be sensitive to specific words

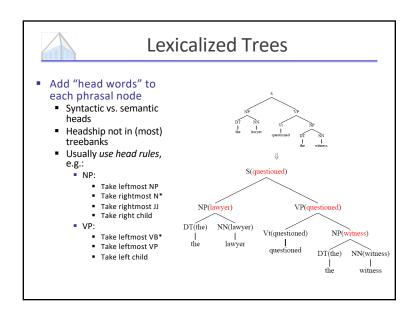


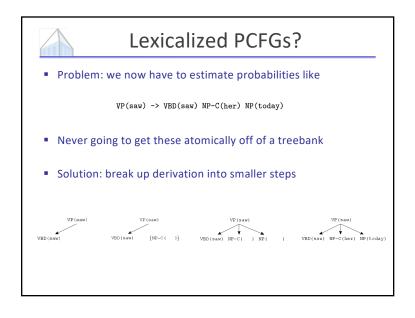
#### Problems with PCFGs

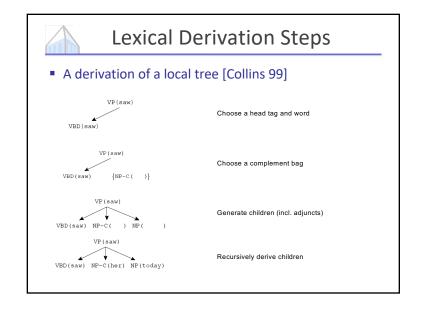


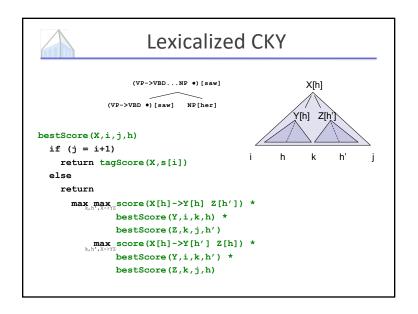


- What's different between basic PCFG scores here?
- What (lexical) correlations need to be scored?











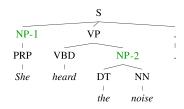
#### Results

- Some results
  - Collins 99 88.6 F1 (generative lexical)
  - Charniak and Johnson 05 89.7 / 91.3 F1 (generative lexical / reranked)
  - Petrov et al 06 90.7 F1 (generative unlexical)
  - McClosky et al 06 92.1 F1 (gen + rerank + self-train)
- However
  - Bilexical counts rarely make a difference (why?)
  - Gildea 01 Removing bilexical counts costs < 0.5 F1

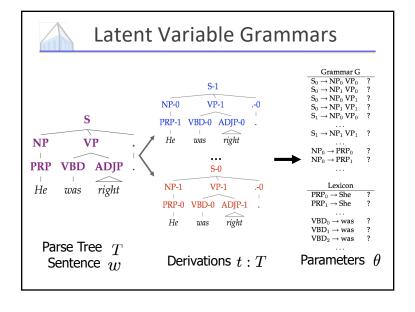
#### Latent Variable PCFGs



#### The Game of Designing a Grammar



- Annotation refines base treebank symbols to improve statistical fit of the grammar
  - Parent annotation [Johnson '98]
  - Head lexicalization [Collins '99, Charniak '00]
  - Automatic clustering?





# **Learned Splits**

Proper Nouns (NNP):

NNP-14	Oct.	Nov.	Sept.
NNP-12	John	Robert	James
NNP-2	J.	E.	L.
NNP-1	Bush	Noriega	Peters
NNP-15	New	San	Wall
NNP-3	York	Francisco	Street

Personal pronouns (PRP):

PRP-0	It	He	
PRP-1	it	he	they
PRP-2	it	them	him



# **Learned Splits**

• Relative adverbs (RBR):

RBR-0	further	lower	higher
RBR-1	more	less	More
RBR-2	earlier	Earlier	later

Cardinal Numbers (CD):

CD-7	one	two	Three
CD-4	1989	1990	1988
CD-11	million	billion	trillion
CD-0	1	50	100
CD-3	1	30	31
CD-9	78	58	34