

Natural Language Processing



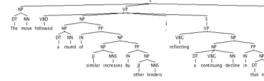
Berkeley

Syntax and Parsing

Dan Klein – UC Berkeley

Syntax

Parse Trees



The move followed a round of similar increases by other lenders, reflecting a continuing decline in that market

Phrase Structure Parsing

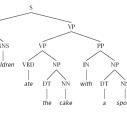
- Phrase structure parsing organizes syntax into constituents or brackets
- In general, this involves nested trees
- Linguists can, and do, argue about details
- Lots of ambiguity
- Not the only kind of syntax...

new art critics write reviews with computers



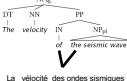
Constituency Tests

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

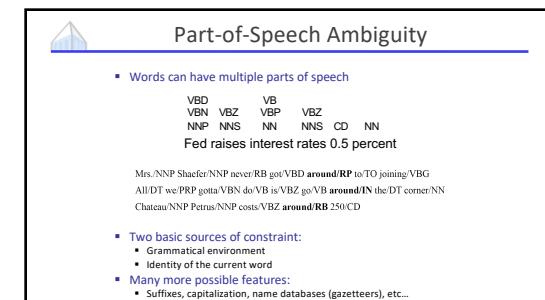
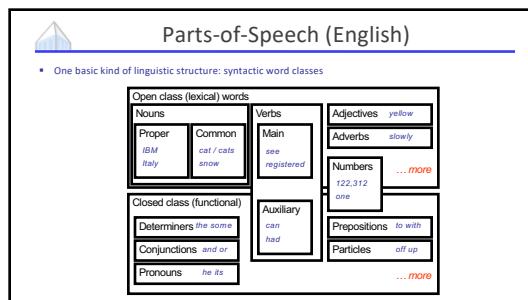
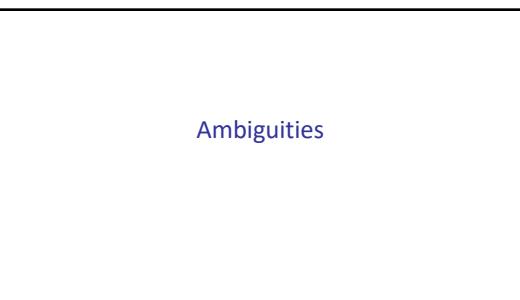
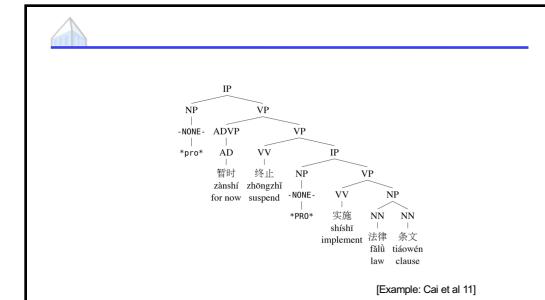
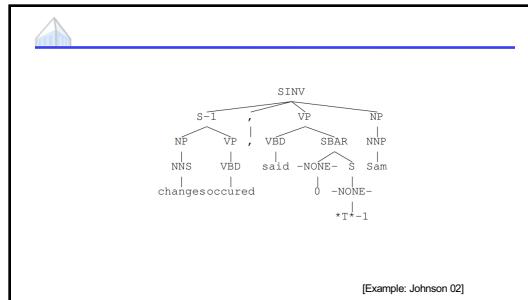
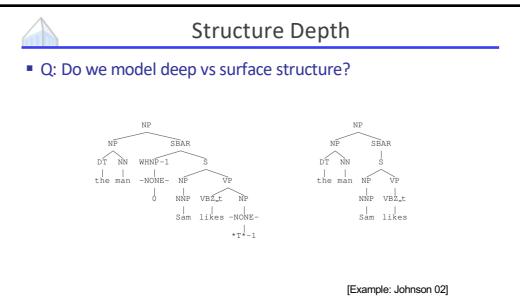


Conflicting Tests

- Constituency isn't always clear
 - Units of transfer:
 - think about ~ penser à
 - talk about ~ hablar de
- Phonological reduction:
 - I will go → I'll go
 - I want to go → I wanna go
 - a le centre → au centre
- Coordination
 - He went to and came from the store.



La vélodé des ondes sismiques





Why POS Tagging?

- Historically useful in and of itself (more than you'd think)
 - Text-to-speech: record, lead
 - Lemmatization: saw(v) → see, saw[n] → 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 ...

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



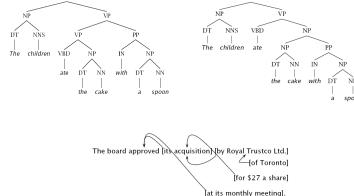
Classical NLP: Parsing

- Write symbolic or logical rules:

Grammar (CFG)	Lexicon
ROOT → S	NP → NP PP
S → NP VP	VP → VBP NP
NP → DT NN	VP → VBP NP PP
NP → NN NNS	PP → IN NP
- 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



Ambiguities: PP Attachment



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.

Inaccessible Ambiguities

- Inaccessible 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 $\langle N, T, S, R \rangle$
 - N : the set of non-terminals
 - Phrasal categories: S, NP, VP, ADJP, etc.
 - Parts-of-speech (pre-terminal): 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 $X \rightarrow Y_1 Y_2 \dots Y_n$, with $X, Y_i \in N$
 - Examples: $S \rightarrow NP VP$, $VP \rightarrow VP CO VP$
 - Also called rewrites, productions, or local trees
- A PCFG adds:
 - A top-down production probability per rule $P(Y_1 Y_2 \dots Y_n | X)$

Treebank Sentences

```
C (S (NP-SBJ The move)
     (VP followed
      (NP (NP a round)
       (PP of
        (NP (NP similar increases)
         (PP by
          (NP other lenders)))
         (PP against
          (NP Arizona real estate loans)))))

     (S-ADV (NP-SBJ *)
      (VP reflecting
       (NP (NP a continuing decline)
        (PP-LOC in
         (NP that market))))))

    .))
```

Treebank Grammars

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

ROOT → S	1
S → NP VP .	1
NP → PRP	1
VP → VBD ADJP	1
.....	

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

Treebank Grammar Scale

- Treebank grammars can be enormous
 - As FSAs, the raw grammar has $\sim 10^6$ states, excluding the lexicon
 - Better parsers usually make the grammars larger, not smaller

Chomsky Normal Form

- Chomsky normal form:
 - All rules of the form $X \rightarrow YZ$ 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 hard
 - The straightforward transformations don't preserve tree scores
- Makes parsing algorithms simpler!

CKY Parsing

A Recursive Parser

```

bestScore(X,i,j)
if (j == i+1)
  return tagScore(X,s[i])
else
  return max score(X->YZ) *
    bestScore(Y,i,k) *
    bestScore(Z,k,j)
  
```

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

A Memoized Parser

- One small change:

```

bestScore(X,i,j)
if (scores[X][i][j] == null)
  if (j == i+1)
    score = tagScore(X,s[i])
  else
    score = max score(X->YZ) *
      bestScore(Y,i,k) *
      bestScore(Z,k,j)
  scores[X][i][j] = score
return scores[X][i][j]
  
```

A Bottom-Up Parser (CKY)

- Can also organize things bottom-up

```

bestScore(s)
for (i : {0,n-1})
  for (X : tags[s[i]]) {
    scores[X][i][i+1] =
      tagScore(X,s[i])
    for (diff : {2,n})
      for (i : {0,n-diff})
        j = i + diff
        for (X->YZ : rule)
          for (k : {i+1, j-1})
            scores[X][i][j] = max scores[X][i][j],
              score(X->YZ) *
              score[Y][i][k] *
              score[Z][k][j]
  
```

Unary Rules

- Unary rules?

```

bestScore(X,i,j,s)
if (j == i+1)
  return tagScore(X,s[i])
else
  return max max score(X->YZ) *
    bestScore(Y,i,k) *
    bestScore(Z,k,j)
  max score(X->Y) *
    bestScore(Y,i,j)
  
```

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

```

bestScoreB(X,i,j,s)
    return max max score(X->Y) *
        bestScoreB(Y,i,k) *
        bestScoreU(Z,k,j)

bestScoreU(X,i,j,s)
    if (j == i+1)
        return tagScore(X,s[i])
    else
        return max max score(X->Y) *
            bestScoreB(Y,i,j)
  
```

Learning PCFGs

Treebank PCFGs [Charniak 96]

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

Model	F1
Baseline	72.0

Conditional Independence?

- 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

Non-Independence

- Independence assumptions are often too strong.

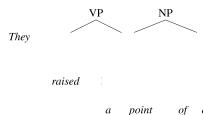
All NPs	NPs under S	NPs under VP
NP PP 11%	NP PP 9%	NP PP 23%
DT NN 9%	DT NN 9%	DT NN 7%
PRP 6%	PRP 21%	PRP 4%

- Example: the expansion of an NP is highly dependent on the parent of the NP (i.e., subjects vs. objects).
- Also: the subject and object expansions are correlated!



Grammar Refinement

- Example: PP attachment



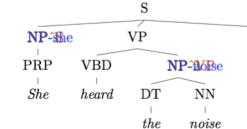
Natural Language Processing



Berkeley
Syntax and Parsing
Dan Klein – UC Berkeley



Grammar Refinement

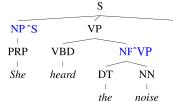


- Structure Annotation [Johnson '98, Klein&Manning '03]
- Lexicalization [Collins '99, Charniak '00]
- Latent Variables [Matsuzaki et al. 05, Petrov et al. '06]

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:
 - $VP \rightarrow VP\ PP$
 - $NP \rightarrow NP\ PP$
- 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?

Lexicalized Trees

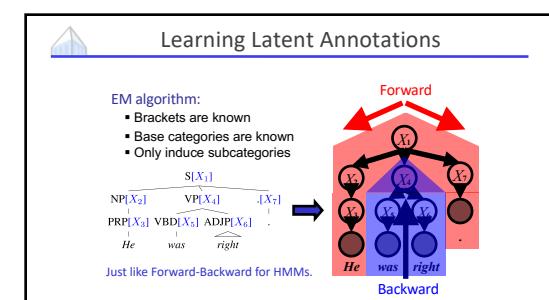
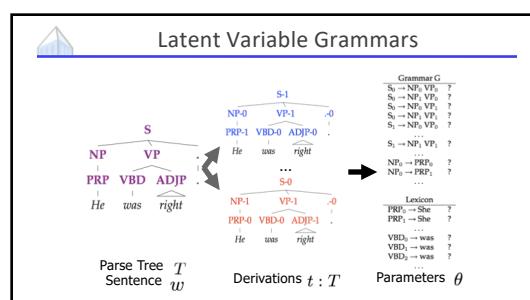
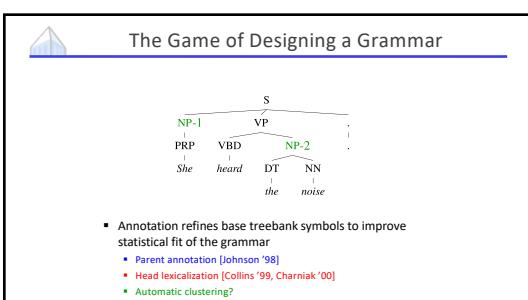
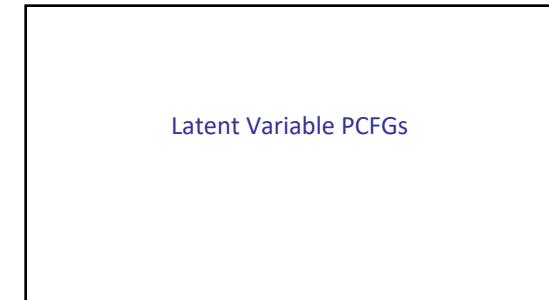
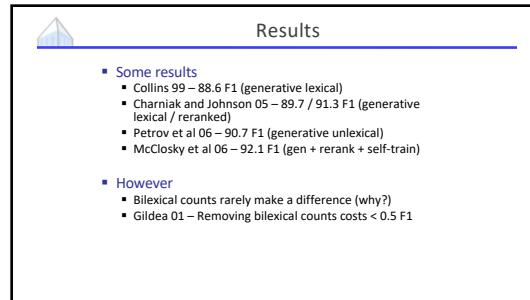
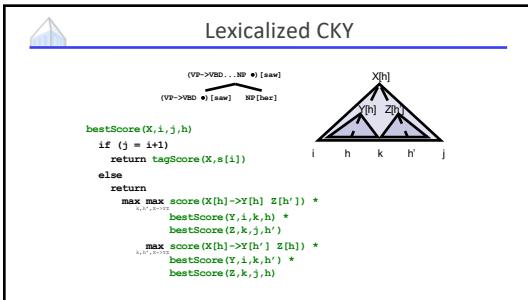
- Add "head words" to each phrasal node
 - Syntactic vs. semantic head
 - Headship not in (most) treebanks
- Usually use head rules, e.g.:
 - NP:
 - Take leftmost NP
 - Take rightmost N*
 - Take rightmost JJ
 - Take right child
 - VP:
 - Take leftmost VP*
 - Take leftmost VP
 - Take left child

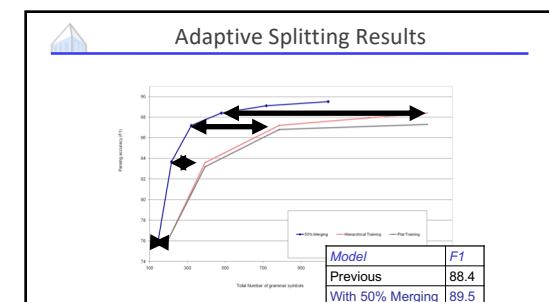
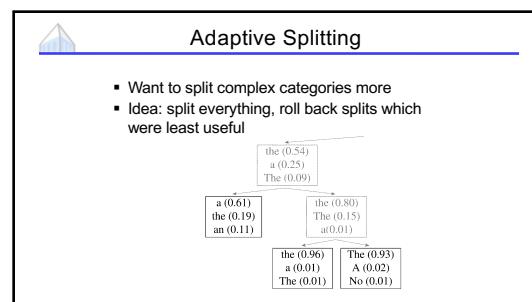
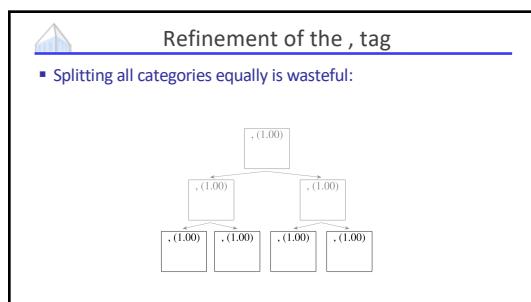
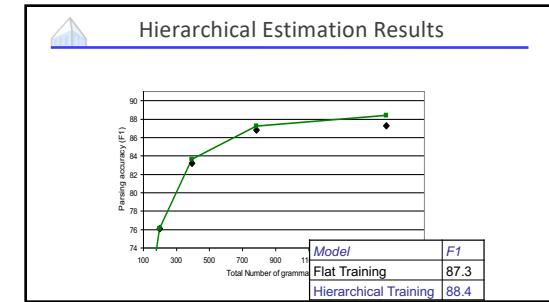
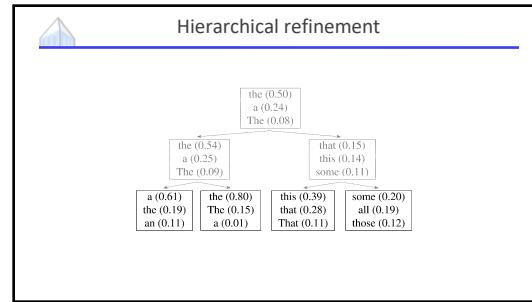
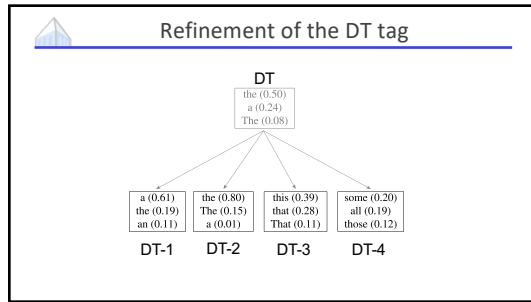
Lexicalized PCFGs?

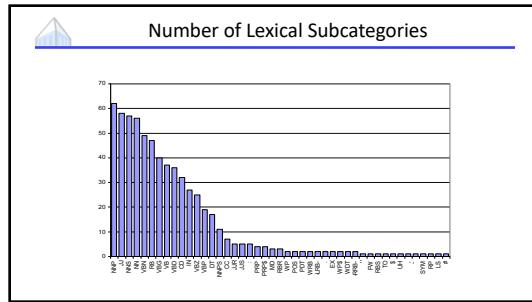
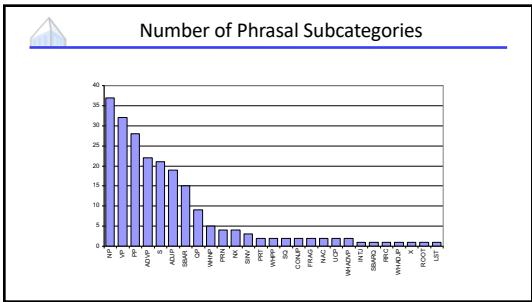
- Problem: we now have to estimate probabilities like
 $VP(\text{aw}) \rightarrow VBD(\text{aw})\ NP-C(\text{her})\ NP(\text{today})$
- Never going to get these atomically off of a treebank
- Solution: break up derivation into smaller steps

Lexical Derivation Steps

- A derivation of a local tree [Collins 99]







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	I	
PRP-1	it	he	they	
PRP-2	it	them	him	

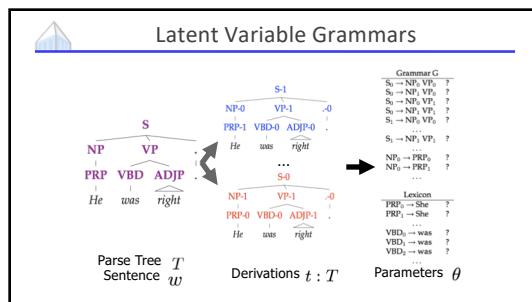


Learned Splits

- Relative adverbs (RBR):

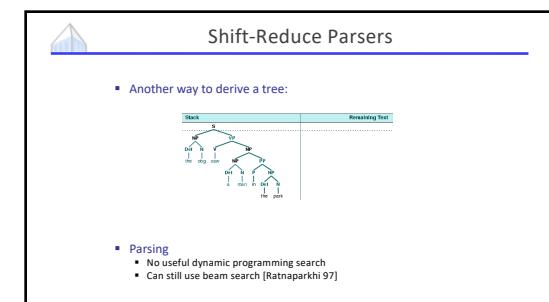
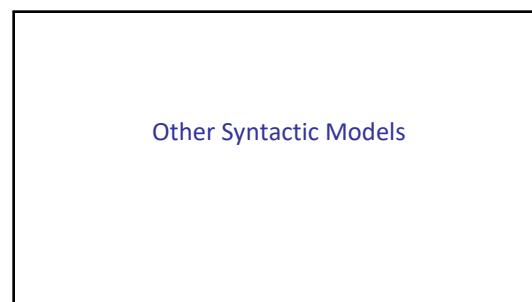
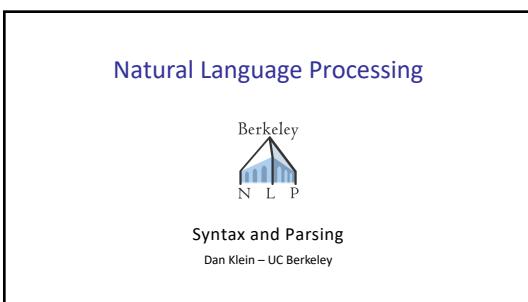
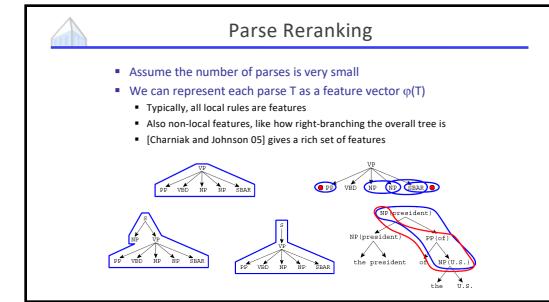
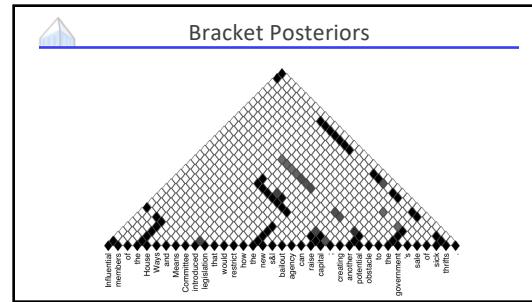
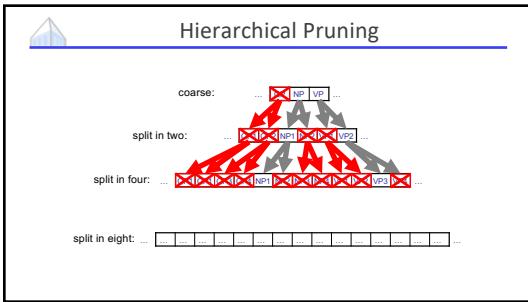
RBR-0	further	lower	higher
RBR-1	more	less	More
RBR-2	earlier	Earlier	later
- Cardinal Numbers (CD):

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



```

graph TD
    S --- NP
    S --- VP
    NP --- PRP
    PRP --- They
    VP --- "?"
    "??" --- V
    "??" --- NP
    "??" --- PP
    V --- raised
    NP --- DT
    NP --- NN
    DT --- a
    NN --- point
    PP --- IN
    PP --- NP
    IN --- of
    NP --- order
  
```



The diagram illustrates a tree structure representing transformations. The root node is 4, which has three children: 1, 2, and 3. Node 1 has four children: A, B, C, and D. Node 2 has one child: E. Node 3 has no children. Solid lines connect 4 to 1, 1 to 2, and 2 to 3. Dashed lines connect 4 to 3 and 1 to 3. Below the tree is a row of 10 symbols: / \ / \ / \ / \ / \ / \.

$\gamma = \gamma_0 + r$

Dependency Parsing

- Lexicalized parsers can be seen as producing dependency trees

```

graph TD
    S[questioned] --- NP_law[NP lawyer]
    S --- VP_q[VP questioned]
    NP_law --- DT_l[the]
    NP_law --- NN_l[NN lawyer]
    VP_q --- Vt_q[Vt questioned]
    VP_q --- NP_witness[NP witness]
    NN_l --- Vt_q
    DT_l --- Vt_q
    NN_witness --- Vt_q
    NN_witness --- NN_t[NN witness]
    Vt_q --- questioned[questioned]
    questioned --- lawyer[lawyer]
    questioned --- witness[witness]
    questioned --- the1[the]
    questioned --- the2[the]
  
```

Each local binary tree corresponds to an attachment in the dependency graph

 Dependency Parsing

- Pure dependency parsing is only cubic [Eisner 99]



- Some work on *non-projective* dependencies
 - Common in, e.g. Czech parsing
 - Can do with MST algorithms [McDonald and Pereira 05]



```

graph TD
    root --- John
    root --- saw
    root --- a
    root --- dog
    root --- yesterday
    root --- which
    root --- was
    root --- a
    root --- Yorkshire
    root --- Terrier
    John --- saw
    saw --- a
    a --- dog
    dog --- yesterday
    which --- was
    was --- a
    a --- Yorkshire
    a --- Terrier
  
```

The post office will hold out discounts and service concessions as incentives

- Rewrite large (possibly lexicalized) subtrees in a single step

- Formally, a *tree-insertion grammar*
- Derivational ambiguity whether subtrees were generated atomically or compositionally
- Most probable parse is NP-complete

TIG: Insertion

ϕ

 $A \downarrow$

ψ

 A

ϕ'

 A

NP_L S VP
 V
 saw

NP_D N
 man

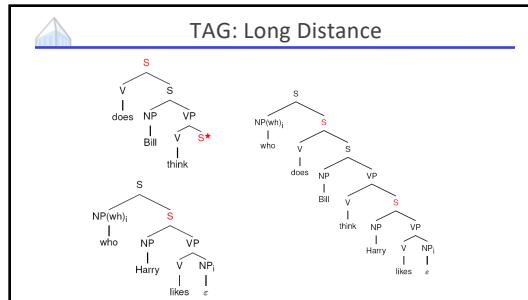
NP_D N S VP
 V
 NP_L
 saw

Tree-adjoining grammars

- Start with *local trees*
- Can insert structure with *adjunction* operation
- Mildly context-sensitive
- Models long-distance dependencies naturally
- ... as well as other weird stuff that CFGs don't capture well (e.g. cross-serial dependencies)

NP NNP NP S Qintex
NP NNP MD VP VP would
NP VB NP sell
NP NNS assets

NP NNP MD VP VP would
NP VB NP sell PRT NP RP NNS off assets



CCG Parsing

- Combinatory Categorial Grammar**
 - Fully (mono-) lexicalized grammar
 - Categories encode argument sequences
 - Very closely related to the lambda calculus (more later)
 - Can have spurious ambiguities (why?)

John ⊢ NP
shares ⊢ NP
buys ⊢ (S\NP)/NP
sleeps ⊢ S\NP
well ⊢ (S\NP)\(S\NP)

NP S
John (S\NP)/NP NP
buys shares