

Features of Statistical Parsers

Mark Johnson

Brown University

April, 2005

Joint work with Eugene Charniak, Michael Collins and students in BLLIP

Supported by NSF grants LIS 9720368 and IIS0095940

Talk outline

- Who cares about parsing?
- PCFGs, maximum likelihood estimation and (lack of) independence
- Dependence, arbitrary features and exponential parsing models
- Joint vs. conditional maximum likelihood estimation of exponential models
- Features of an exponential parsing model

Why parsing?

- Applications
 - Recovering syntactic structure
 - * Information extraction
 - * Question answering
 - * Machine translation
 - Language modeling (distinguishing likely from unlikely sentences)
 - * Speech recognition
 - * Speech error detection and correction
- Science
 - First step on the path to meaning

Short history of parsing

- (1960s–1970s) basic dynamic programming algorithms, ATNs
 - inspired by linguistic theory (transformational grammar)
- (1980s–1990s) “unification” grammar (no deep structure)
 - parsers provably implement a specific linguistic grammar
 - grammar \sim axioms, UG \sim logic, parsing \sim deduction
 - emphasis on grammar formalisms
- (1980s–??) direct statistical models (e.g., n-gram models)
 - revolutionized speech recognition, practical applications
- (1990s–??) statistical models of linguistic structure
 - “curse of dimensionality” limits direct statistical models
 - emphasis on frequent constructions

N-gram models of strings

- Model dependencies between n adjacent words (*n-gram models*)
 - predict a word based on the $n - 1$ preceding words

$\$ \rightarrow the \rightarrow man \rightarrow in \rightarrow the \rightarrow hat \rightarrow drinks \rightarrow red \rightarrow wine \rightarrow \$$

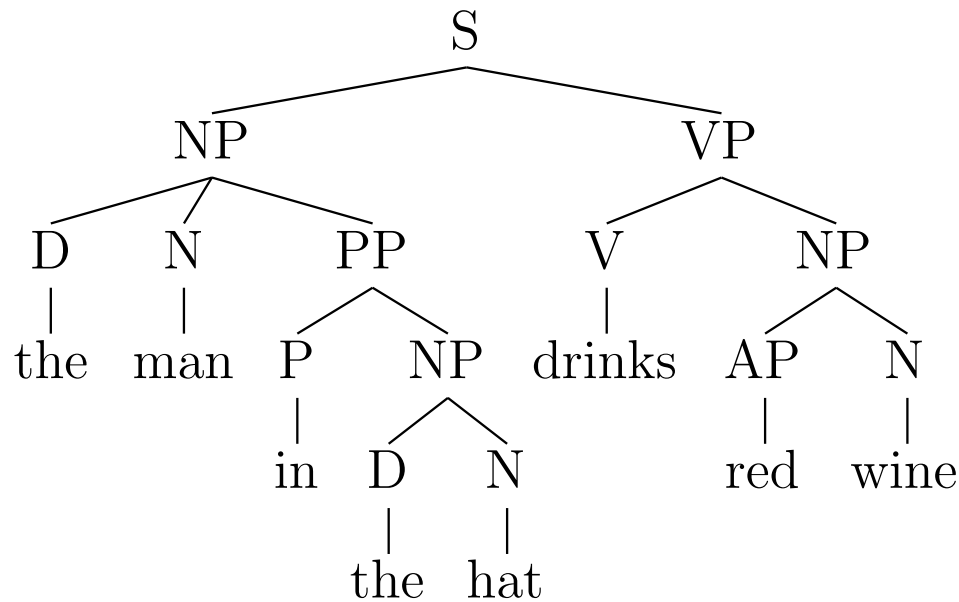
$n = 1 : P(hat); n = 2 : P(hat|the); n = 3 : P(hat|in, the)$

- N-gram models are *universal approximators* as $n \rightarrow \infty$
- Probabilities estimated from real *corpora*
- Probability distinguishes “good” from “bad” sentences
- These simple models work surprisingly well because *they are lexicalized* and *most dependencies are local*

Limitations of n-gram models

- *Curse of dimensionality*: the number of possible n-grams (sequences of n words) grows exponentially with n
 - $n = 1 \Rightarrow V^n \approx 1,000$
 - $n = 2 \Rightarrow V^n \approx 1,000,000$
 - $n = 3 \Rightarrow V^n \approx 1,000,000,000$
- *Bias-variance dilemma*: As n gets larger ...
 - the model becomes more accurate (lower bias)
 - more data is needed to learn the model (higher variance)
- *Structured models* collect statistics on linguistic structures
 - only model linguistic dependencies (strong bias)
 - ignore other dependencies (lower variance)

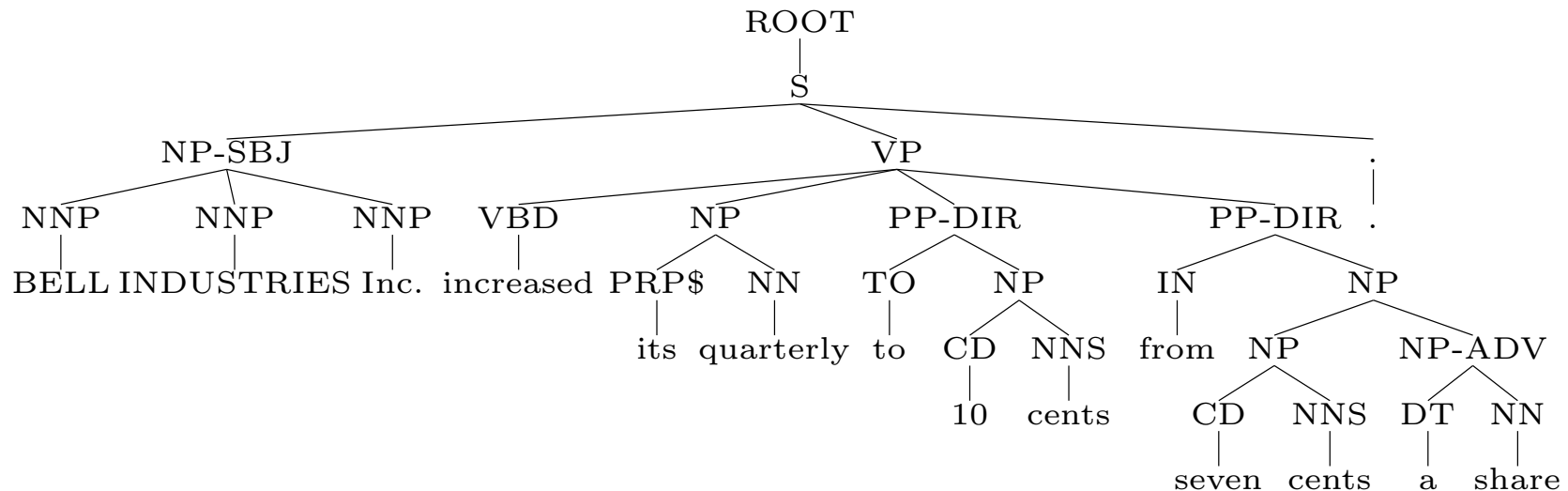
Probabilistic Context Free Grammars



$$\begin{aligned} P(\mathbf{t}) &= P(S \rightarrow NP VP) \times \\ &P(NP \rightarrow D N PP) \times \\ &P(D \rightarrow \text{the}) \times \\ &P(N \rightarrow \text{man}) \times \\ &\dots \end{aligned}$$

- Rules are associated with *probabilities*
- Tree probability is the *product of rule probabilities*
- *Most probable tree* is “best guess” at correct syntactic structure

Treebank corpora



- The Penn treebank contains hand-annotated parse trees for ~ 50,000 sentences
- Treebanks also exist for the Brown corpus, the Switchboard corpus (spontaneous telephone conversations) and Chinese and Arabic corpora

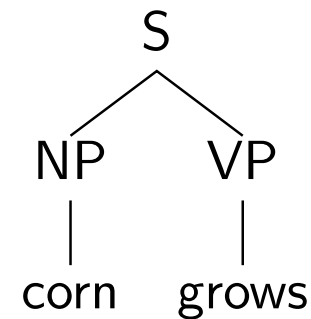
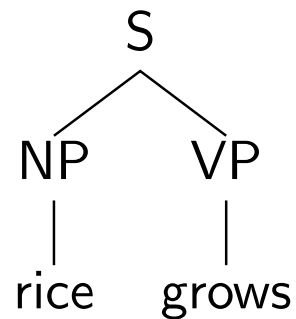
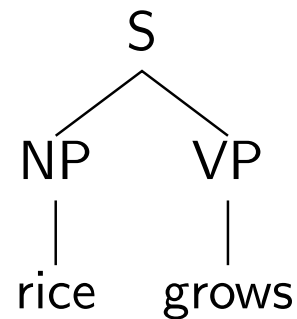
Estimating a grammar from a treebank

- *Maximum likelihood principle*: Choose the grammar and rule probabilities that make the trees in the corpus *as likely as possible*
 - read the rules off the trees
 - for PCFGs, set rule probabilities to the *relative frequency* of each rule in the treebank

$$P(\text{VP} \rightarrow \text{V NP}) = \frac{\text{Number of times VP} \rightarrow \text{V NP} \text{ occurs}}{\text{Number of times VP occurs}}$$

- *If the language is generated by a PCFG and the treebank trees are its derivation trees, the estimated grammar converges to the true grammar.*

Estimating PCFGs from visible data

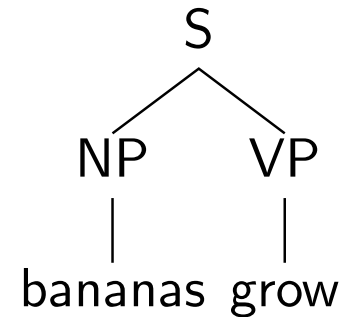
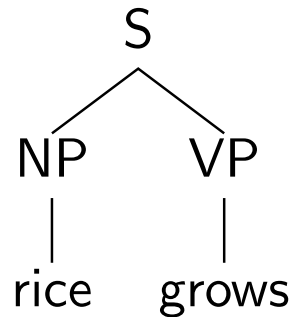
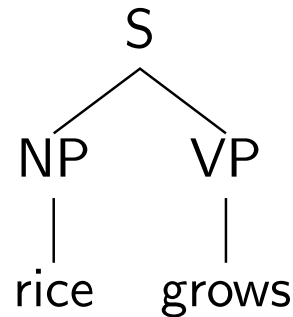


Rule	Count	Rel Freq
$S \rightarrow NP VP$	3	1
$NP \rightarrow \text{rice}$	2	$2/3$
$NP \rightarrow \text{corn}$	1	$1/3$
$VP \rightarrow \text{grows}$	3	1

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{rice} \quad \text{grows} \end{array} \right) = 2/3$$

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{corn} \quad \text{grows} \end{array} \right) = 1/3$$

Non-local constraints and PCFG MLE



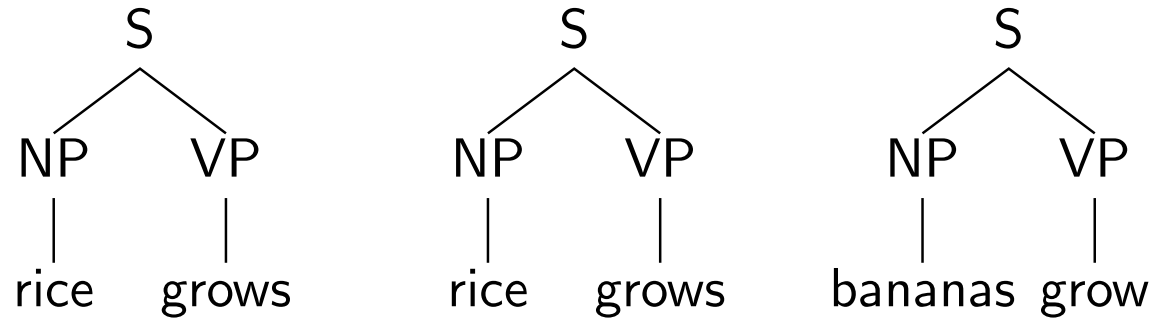
Rule	Count	Rel Freq
$S \rightarrow NP VP$	3	1
$NP \rightarrow \text{rice}$	2	$2/3$
$NP \rightarrow \text{bananas}$	1	$1/3$
$VP \rightarrow \text{grows}$	2	$2/3$
$VP \rightarrow \text{grow}$	1	$1/3$

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{rice} \quad \text{grows} \end{array} \right) = 4/9$$

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{bananas} \quad \text{grow} \end{array} \right) = 1/9$$

partition function $Z = 5/9$

Dividing by partition function Z



Rule

Count Rel Freq

$S \rightarrow NP VP$

3

1

$NP \rightarrow rice$

2

$2/3$

$NP \rightarrow bananas$

1

$1/3$

$VP \rightarrow grows$

2

$2/3$

$VP \rightarrow grow$

1

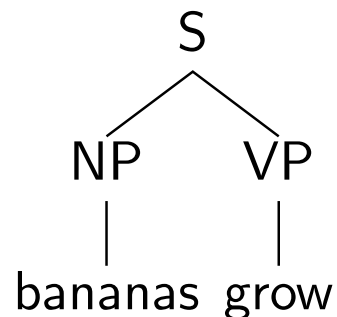
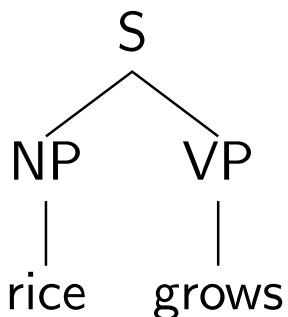
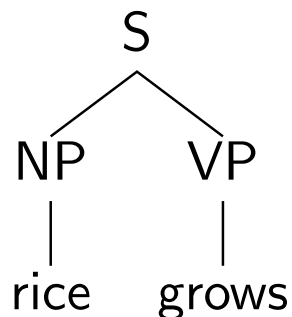
$1/3$

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ rice \quad grows \end{array} \right) = \frac{4}{9} \quad \frac{4}{5}$$

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ bananas \quad grow \end{array} \right) = \frac{1}{9} \quad \frac{1}{5}$$

$$Z = \frac{5}{9}$$

Other values do better!



Rule

Count

Rel Freq

S → NP VP

3

1

NP → rice

2

2/3

NP → bananas

1

1/3

VP → grows

2

1/2

VP → grow

1

1/2

(Abney 1997)

$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ rice \quad grows \end{array} \right) = \frac{2}{6} \quad \frac{2}{3}$$

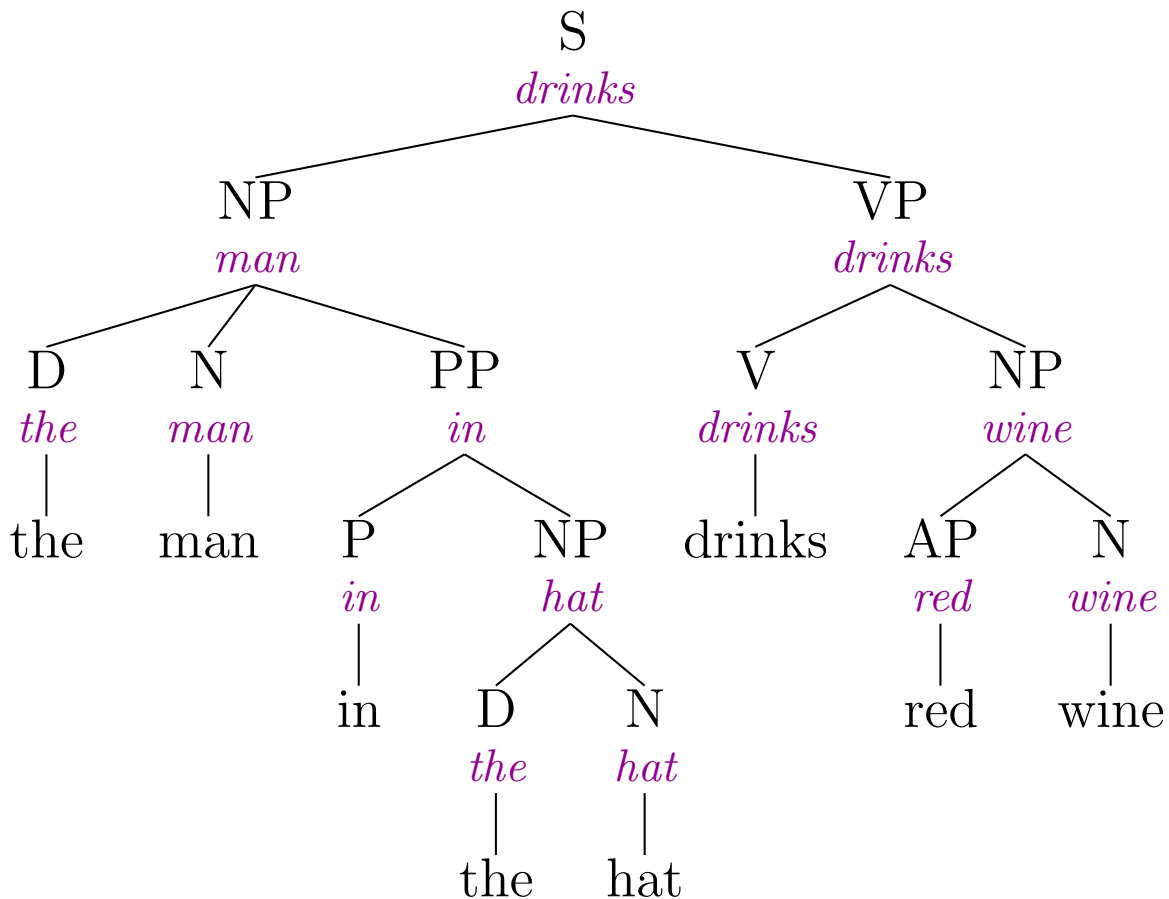
$$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ bananas \quad grow \end{array} \right) = \frac{1}{6} \quad \frac{1}{3}$$

$$Z = \frac{3}{6}$$

Make dependencies local – GPSG-style

rule	count	rel freq	
$S \rightarrow \begin{matrix} \text{NP} & \text{VP} \\ +\text{singular} & +\text{singular} \end{matrix}$	2	2/3	$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ \begin{matrix} \text{NP} & \text{VP} \\ +\text{singular} & +\text{singular} \end{matrix} \\ \downarrow \quad \downarrow \\ \text{rice} \quad \text{grows} \end{array} \right) = 2/3$
$S \rightarrow \begin{matrix} \text{NP} & \text{VP} \\ +\text{plural} & +\text{plural} \end{matrix}$	1	1/3	
$\begin{matrix} \text{NP} \\ +\text{singular} \end{matrix} \rightarrow \text{rice}$	2	1	
$\begin{matrix} \text{NP} \\ +\text{plural} \end{matrix} \rightarrow \text{bananas}$	1	1	$P \left(\begin{array}{c} S \\ \swarrow \quad \searrow \\ \begin{matrix} \text{NP} & \text{VP} \\ +\text{plural} & +\text{plural} \end{matrix} \\ \downarrow \quad \downarrow \\ \text{bananas} \quad \text{grow} \end{array} \right) = 1/3$
$\begin{matrix} \text{VP} \\ +\text{singular} \end{matrix} \rightarrow \text{grows}$	2	1	
$\begin{matrix} \text{VP} \\ +\text{plural} \end{matrix} \rightarrow \text{grow}$	1	1	

“Head to head” dependencies



Rules:

S → NP VP
drinks → man drinks

VP → V NP
drinks → drinks wine

NP → AP N
wine → red wine

...

- *Lexicalization* captures syntactic and semantic dependencies
- Lexicalized structural preferences may be most important

Summary so far

- Maximum likelihood is a good way of estimating a grammar
 - Maximum likelihood estimation of a PCFG from a treebank is easy *if the trees are accurate*
 - But real language has many more dependencies than treebank grammar describes
- ⇒ relative frequency estimator not MLE
- Make non-local dependencies local by splitting categories
- ⇒ Astronomical number of possible categories
- Or find some way of dealing with non-local dependencies ...

Exponential models

- Rules are not independent \Rightarrow in general $Z \neq 1$
- Models with dependencies between features are called *exponential models*
 - Universe \mathcal{T} (set of all possible parse trees)
 - Features $f = (f_1, \dots, f_m)$ ($f_j(\mathbf{t}) =$ value of j feature on $\mathbf{t} \in \mathcal{T}$)
 - Feature weights $w = (w_1, \dots, w_m)$

$$P(\mathbf{t}) = \frac{1}{Z} \exp w \cdot f(\mathbf{t})$$
$$Z = \sum_{\mathbf{t}' \in \mathcal{T}} \exp w \cdot f(\mathbf{t}')$$

Hint: Think of $\exp w \cdot f(\mathbf{t})$ as unnormalized probability of \mathbf{t}

PCFGs are exponential models

\mathcal{T} = set of all trees generated by PCFG G

$f_j(t)$ = number of times the j th rule is used in $t \in \mathcal{T}$

$p(r_j)$ = probability of j th rule in G

Set weight $w_j = \log p(r_j)$

$$f \left(\begin{array}{c} S \\ / \quad \backslash \\ NP \quad VP \\ | \quad | \\ rice \quad grows \end{array} \right) = [\underbrace{1}_{S \rightarrow NP \quad VP}, \underbrace{1}_{NP \rightarrow rice}, \underbrace{0}_{NP \rightarrow bananas}, \underbrace{1}_{VP \rightarrow grows}, \underbrace{0}_{VP \rightarrow grow}]$$

$$P_w(t) = \prod_{j=1}^m p(r_j)^{f_j(t)} = \prod_{j=1}^m (\exp w_j)^{f_j(t)} = \exp(w \cdot f(t))$$

So *a PCFG is just a special kind of exponential model* with $Z = 1$.

Advantages of exponential models

- Exponential models are very flexible ...
- Features f can be *any function of parses* ...
 - whether a particular structure occurs in a parse
 - conjunctions of prosodic and syntactic structure
- Parses t need not be trees, but *can be anything at all*
 - Feature structures (LFG, HPSG)
 - Minimalist derivations
- Exponential models are related to other popular models
 - Harmony theory and optimality theory
 - Maxent models

Modeling dependencies

- It's usually difficult to design a PCFG model that captures a particular set of dependencies
 - probability of the tree must be broken down into a product of *conditional probability distributions*
 - non-local dependencies must be expressed in terms of GPSG-style feature passing
- It's easy to make exponential models sensitive to new dependencies
 - add a new feature functions to existing feature functions
 - *figuring out what the right dependencies are is hard, but incorporating them into an exponential model is easy*

MLE of exponential models is hard

- An exponential model associates features $f(\mathbf{t}) = (f_1(\mathbf{t}), \dots, f_m(\mathbf{t}))$ with weights $\mathbf{w} = (w_1, \dots, w_m)$

$$P(\mathbf{t}) = \frac{1}{Z} \exp \mathbf{w} \cdot f(\mathbf{t})$$
$$Z = \sum_{\mathbf{t}' \in \mathcal{T}} \exp \mathbf{w} \cdot f(\mathbf{t}')$$

- Given treebank $(\mathbf{t}_1, \dots, \mathbf{t}_n)$, MLE chooses \mathbf{w} to maximize $P(\mathbf{t}_1) \times \dots \times P(\mathbf{t}_n)$, i.e., *make the treebank as likely as possible*
 - Computing $P(\mathbf{t})$ requires the partition function Z
 - Computing Z requires a sum over all parses \mathcal{T} for *all sentences*
- \Rightarrow *computing MLE of an exponential parsing model seems very hard*

Conditional ML estimation

- Conditional ML estimation chooses feature weights to maximize $P(t_1|s_1) \times \dots \times P(t_n|s_n)$, where s_i is string for t_i
 - choose feature weights to make t_i most likely relative to parses $\mathcal{T}(s_i)$ for s_i

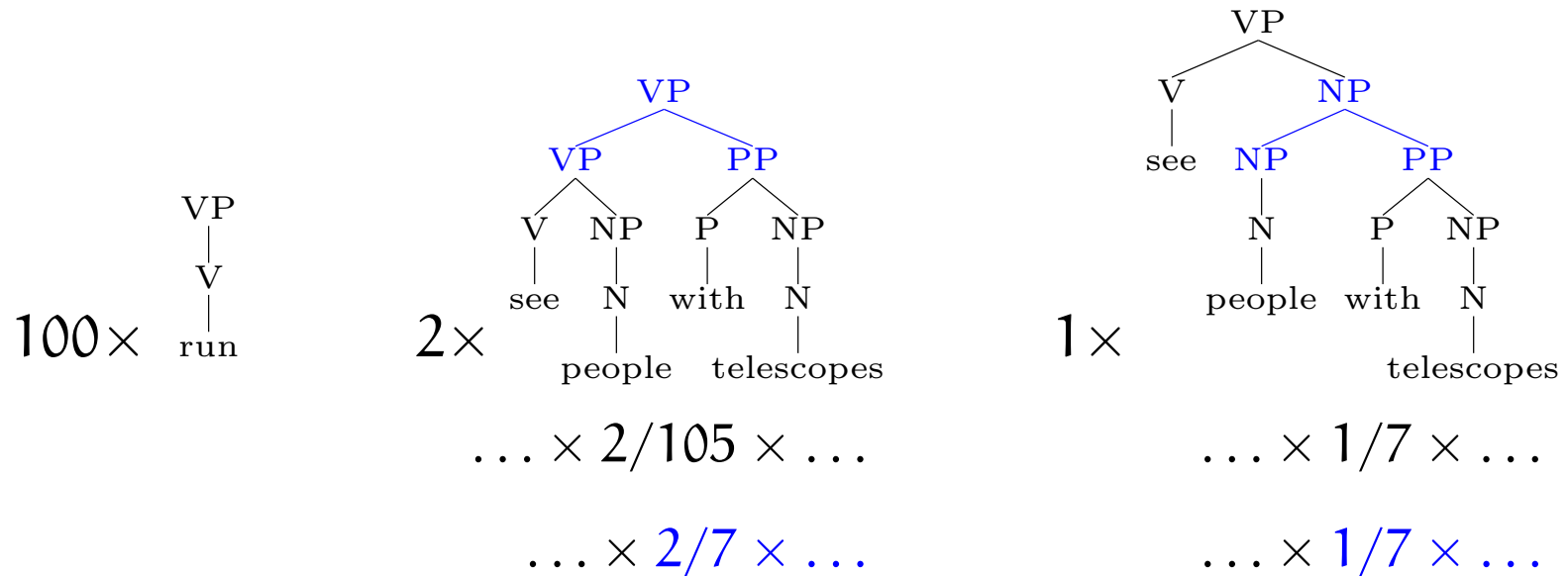
⇒ CMLE *doesn't involve other sentences*

$$P(t|s) = \frac{1}{Z(s)} \exp w \cdot f(t)$$

$$Z(s) = \sum_{t' \in \mathcal{T}(s)} \exp w \cdot f(t')$$

- CMLE “only” involves repeatedly parsing training data
- With “wrong” models, CMLE often produces a more accurate parser than joint MLE

Conditional vs joint MLE



Rule	count	rel freq	rel freq
VP \rightarrow V	100	100/105	4/7
VP \rightarrow V NP	3	3/105	1/7
VP \rightarrow VP PP	2	2/105	2/7
NP \rightarrow N	6	6/7	6/7
NP \rightarrow NP PP	1	1/7	1/7

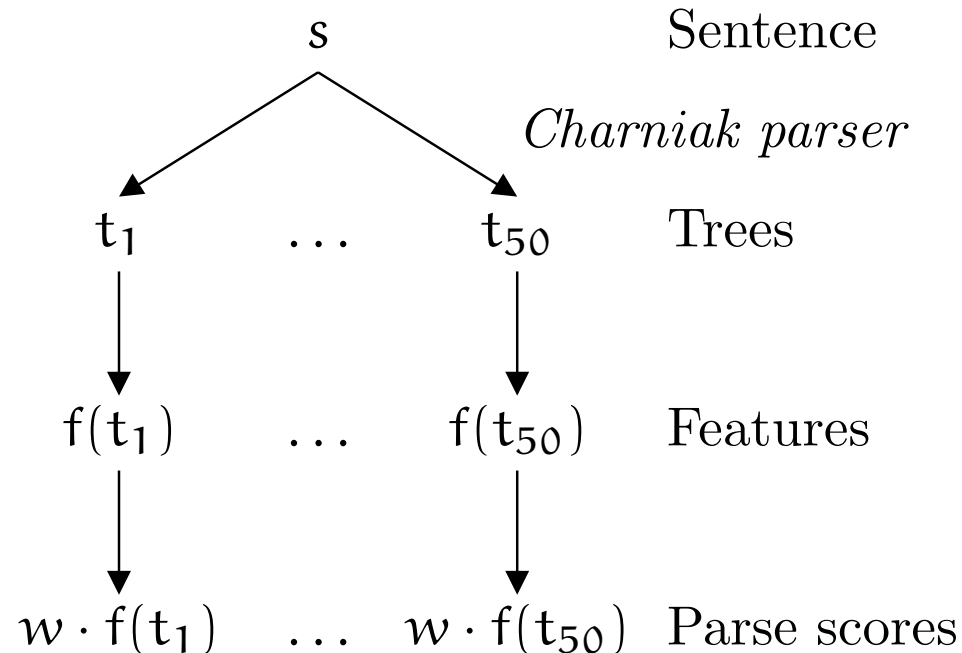
Conditional ML estimation

s	$f(\bar{t})$	$\{f(t) : t \in \mathcal{T}(s), t \neq \bar{t}(s)\}$
sentence 1	(1, 3, 2)	(2, 2, 3) (3, 1, 5) (2, 6, 3)
sentence 2	(7, 2, 1)	(2, 5, 5)
sentence 3	(2, 4, 2)	(1, 1, 7) (7, 2, 1)
...

- Parser designer specifies *feature functions* $f = (f_1, \dots, f_m)$
- A parser produces trees $\mathcal{T}(s)$ for each sentence s
- Treebank tells us correct tree $\bar{t}(s) \in \mathcal{T}(s)$ for sentence s
- Feature functions f apply to each tree $t \in \mathcal{T}(s)$, producing feature values $f(t) = (f_1(t), \dots, f_m(t))$
- MCLE estimates feature weights $w = (w_1, \dots, w_m)$

Coarse-to-fine parsing

- $Z(s)$ is still hard to compute \Rightarrow make $\mathcal{T}(s)$ even smaller!
- Restrict attention to 50-best parses produced by Charniak parser (a good PCFG-based parser)
- Exponential model is trained using CMLE to pick out best parse from Charniak's 50-best parses



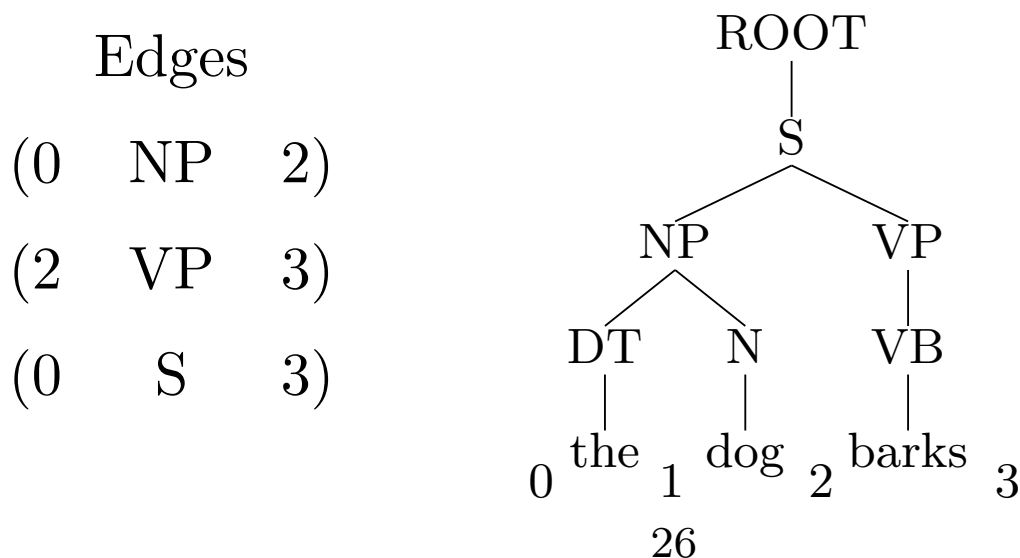
Parser evaluation

- A node's *edge* is its label and beginning and ending *string positions*
- $E(t)$ is the set of edges associated with a tree t
- If t is a parse tree and \bar{t} is the correct tree, then

precision $P_{\bar{t}}(t) = |E(t)| / |E(t) \cap E(\bar{t})|$

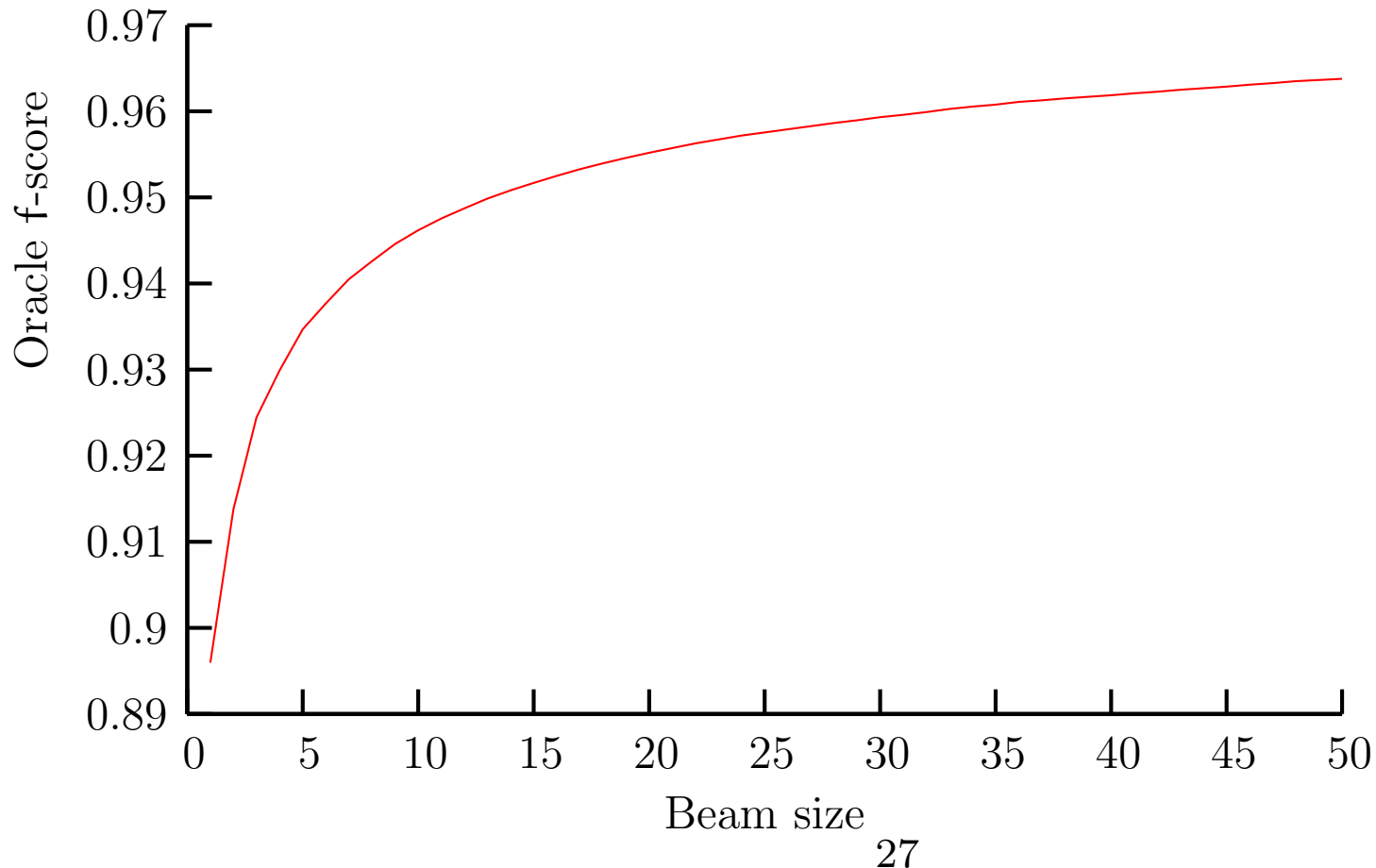
recall $R_{\bar{t}}(t) = |E(\bar{t})| / |E(t) \cap E(\bar{t})|$

f-score $F_{\bar{t}}(t) = 2 / (P_{\bar{t}}(t)^{-1} + R_{\bar{t}}(t)^{-1})$ (geometric mean of P and R)



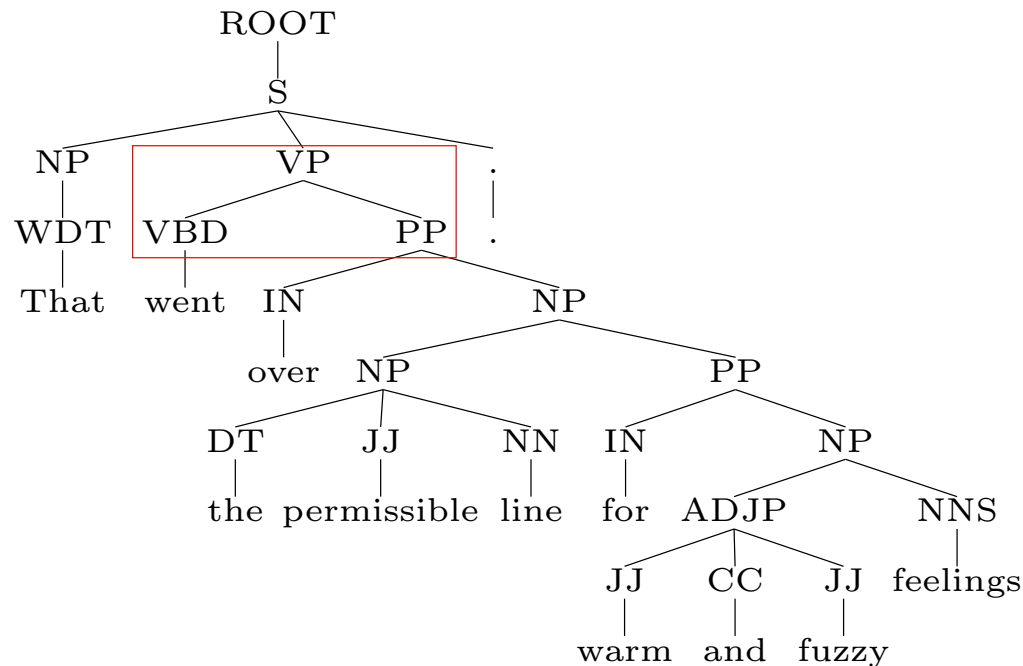
Performance of Charniak parser

- F-score of Charniak's most probable parse = 0.896 (cross-validated on PTB sections 2-19)
- Oracle f-score of Charniak's 50-best parses = 0.965 (66% redn)



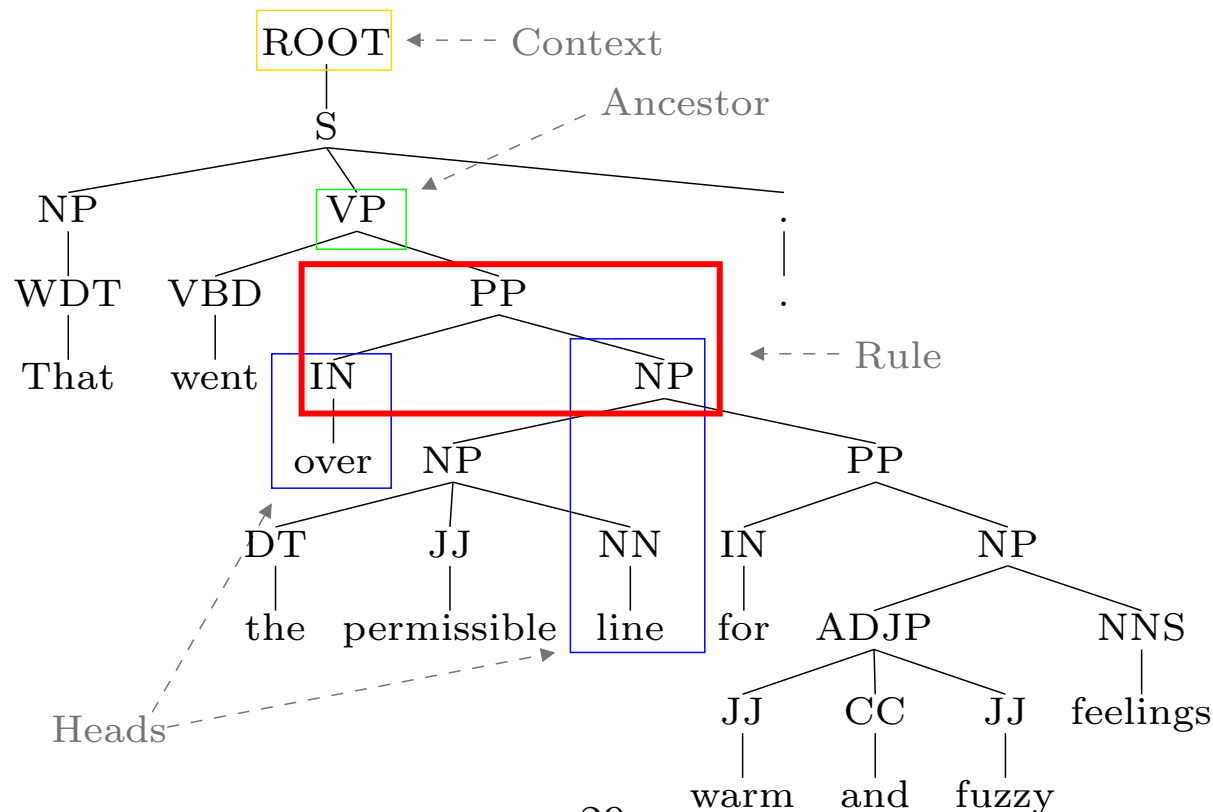
Expt 1: Only “old” features

- Features: 1 *log Charniak probability*, 10,124 Rule features
 - Charniak’s parser already conditions on local trees!
 - Feature selection: features must vary on 5 or more sentences
 - Results: f-score = 0.894; baseline = 0.890; $\approx 4\%$ error reduction
- ⇒ *discriminative training alone can improve accuracy*



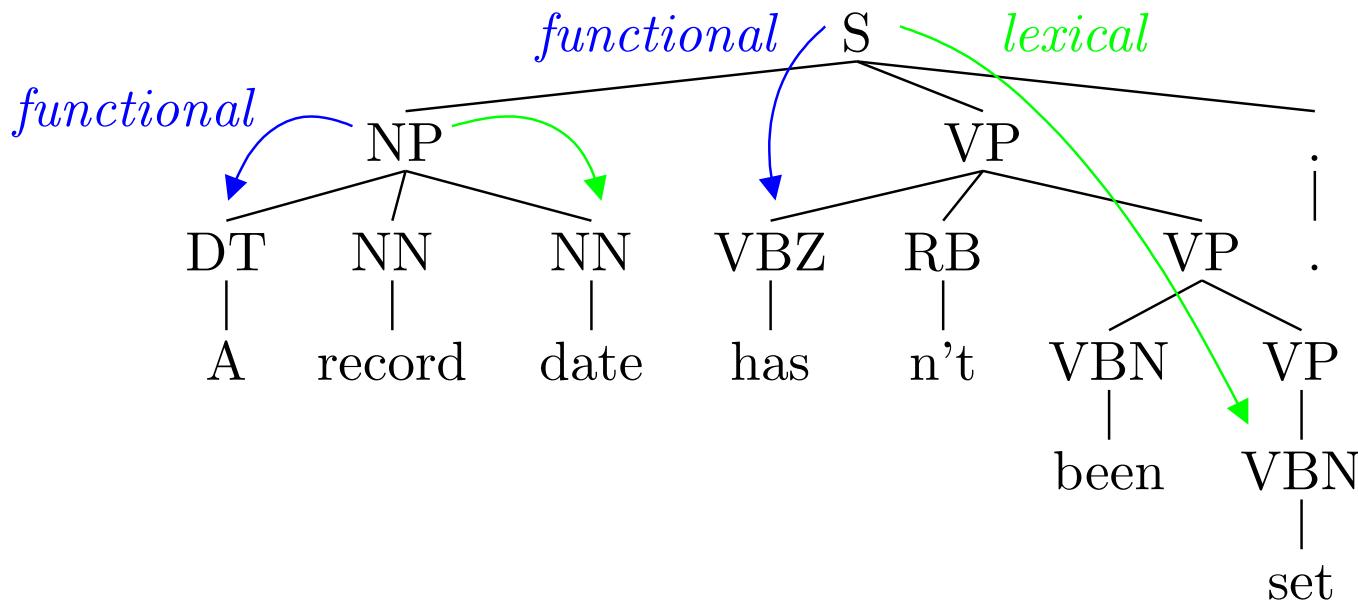
Lexicalized and parent-annotated rules

- *Lexicalization* associates each constituent with its head
- *Ancestor annotation* provides a little “vertical context”
- *Context annotation* indicates constructions that only occur in main clause (c.f., Emonds)



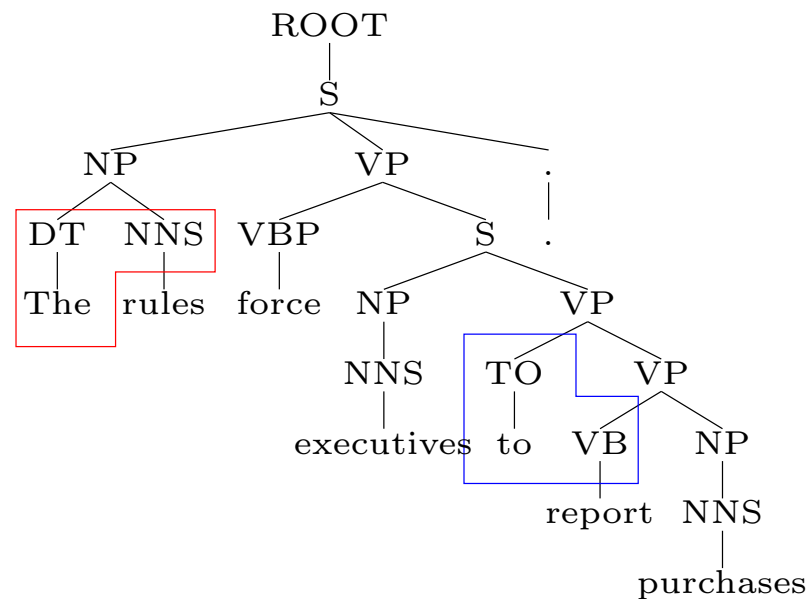
Functional and lexical heads

- There are at least two sensible notions of head (c.f., Grimshaw)
 - *Functional heads*: determiners of NPs, auxiliary verbs of VPs, etc.
 - *Lexical heads*: rightmost Ns of NPs, main verbs in VPs, etc.
- In a Maxent model, it is easy to use both!



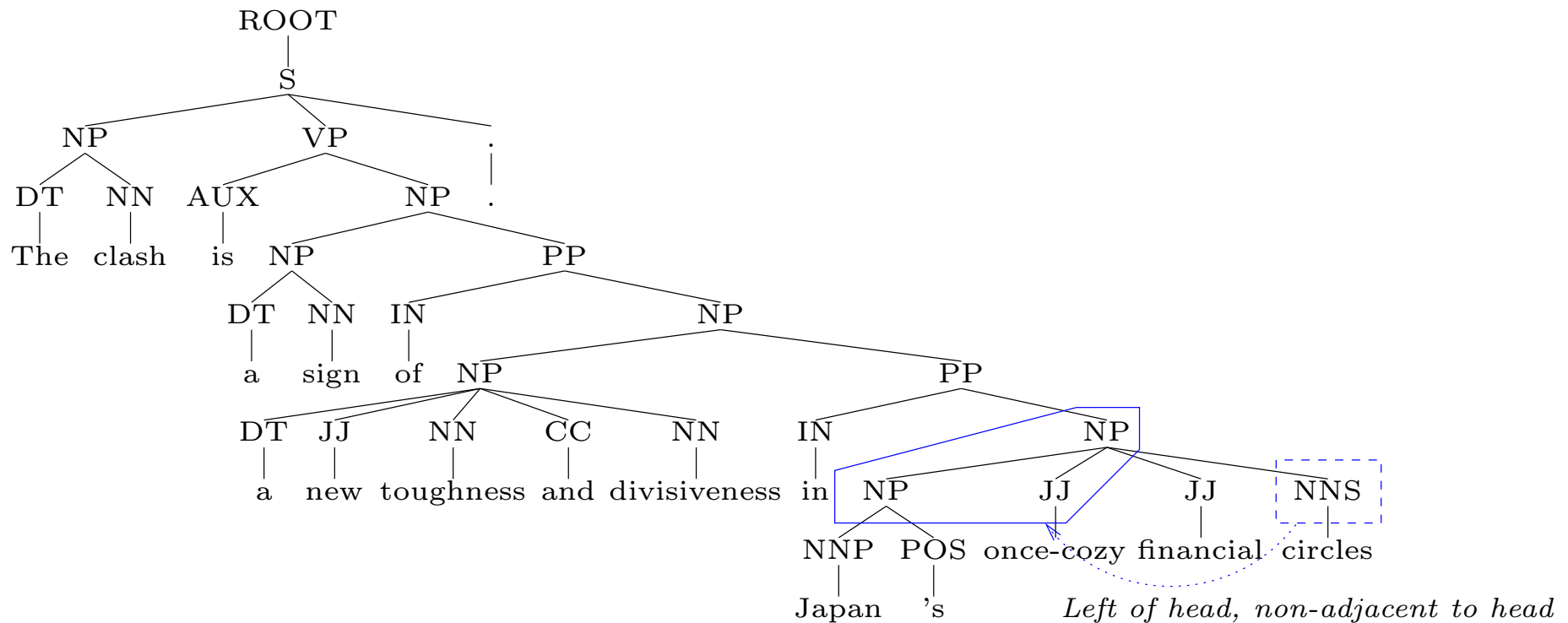
Functional-lexical head dependencies

- The SynSemHeads features collect pairs of functional and lexical heads of phrases
- This captures number agreement in NPs and aspects of other head-to-head dependencies



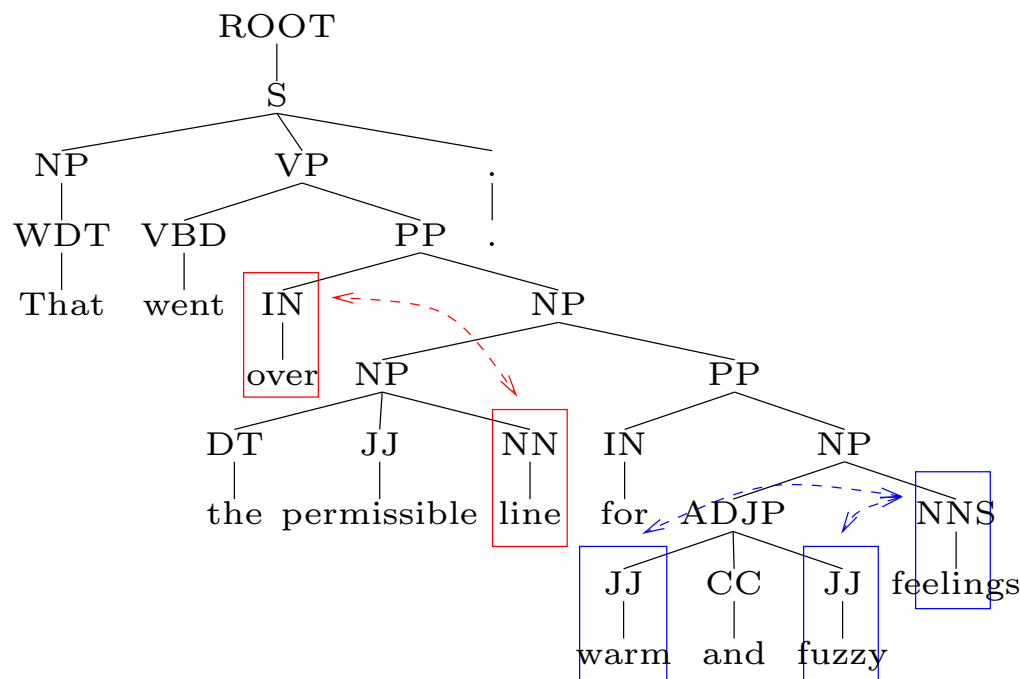
n-gram rule features generalize rules

- Collects *adjacent constituents* in a local tree
- Also includes *relationship to head*
- Constituents can be *ancestor-annotated* and *lexicalized*



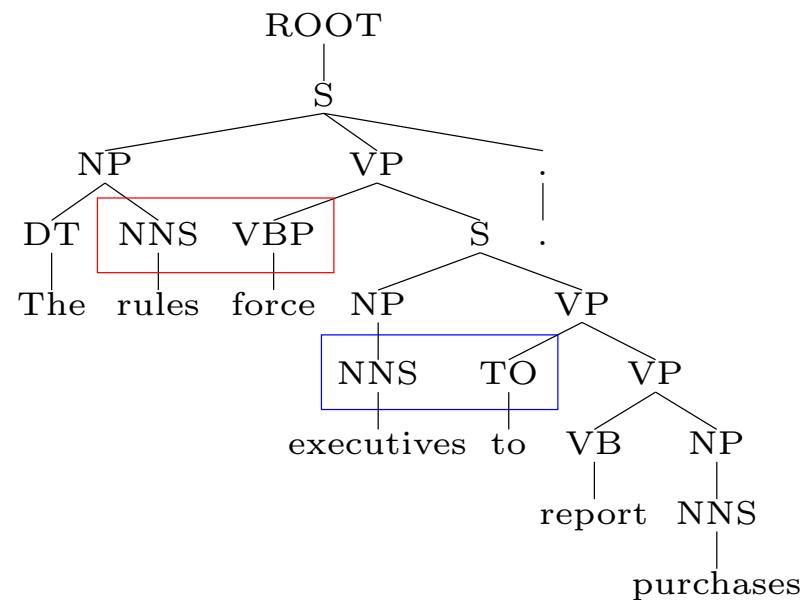
Head to head dependencies

- Head-to-head dependencies track the function-argument dependencies in a tree
- Co-ordination leads to phrases with multiple heads or functors
- Parameterized by head type and number of governors to include



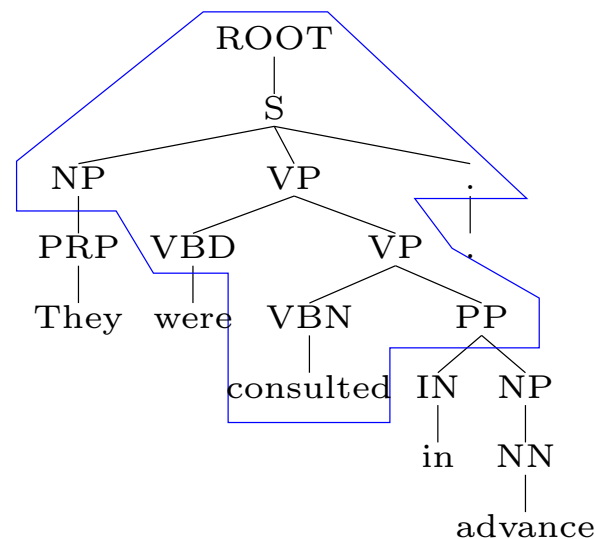
Subject-Verb Agreement

- The SubjVerbAgr features are the POS of the subject NP's lexical head and the VP's functional head



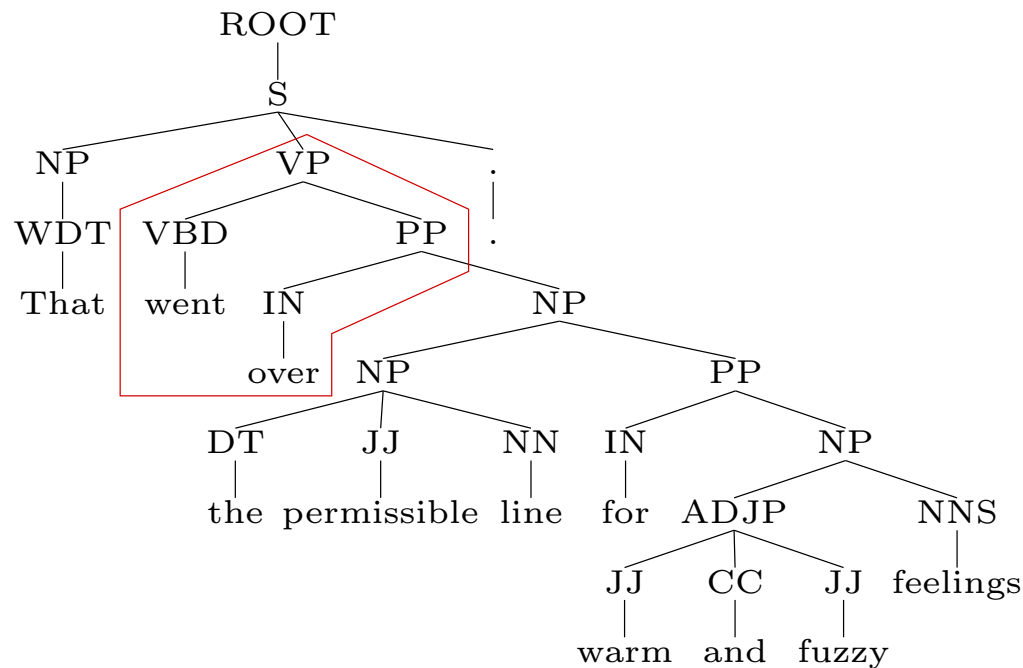
Head trees record all dependencies

- Head trees consist of a (lexical) head, all of its projections and (optionally) all of the siblings of these nodes
- These correspond roughly to TAG elementary trees



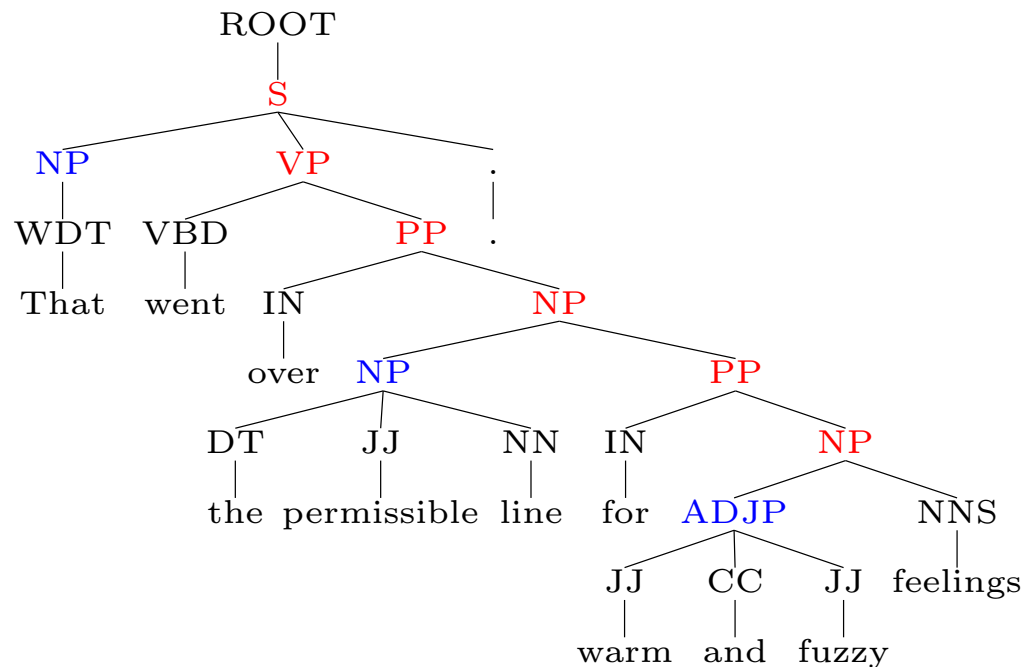
Tree n-gram

- A tree n-gram are tree fragments that connect sequences of adjacent n words



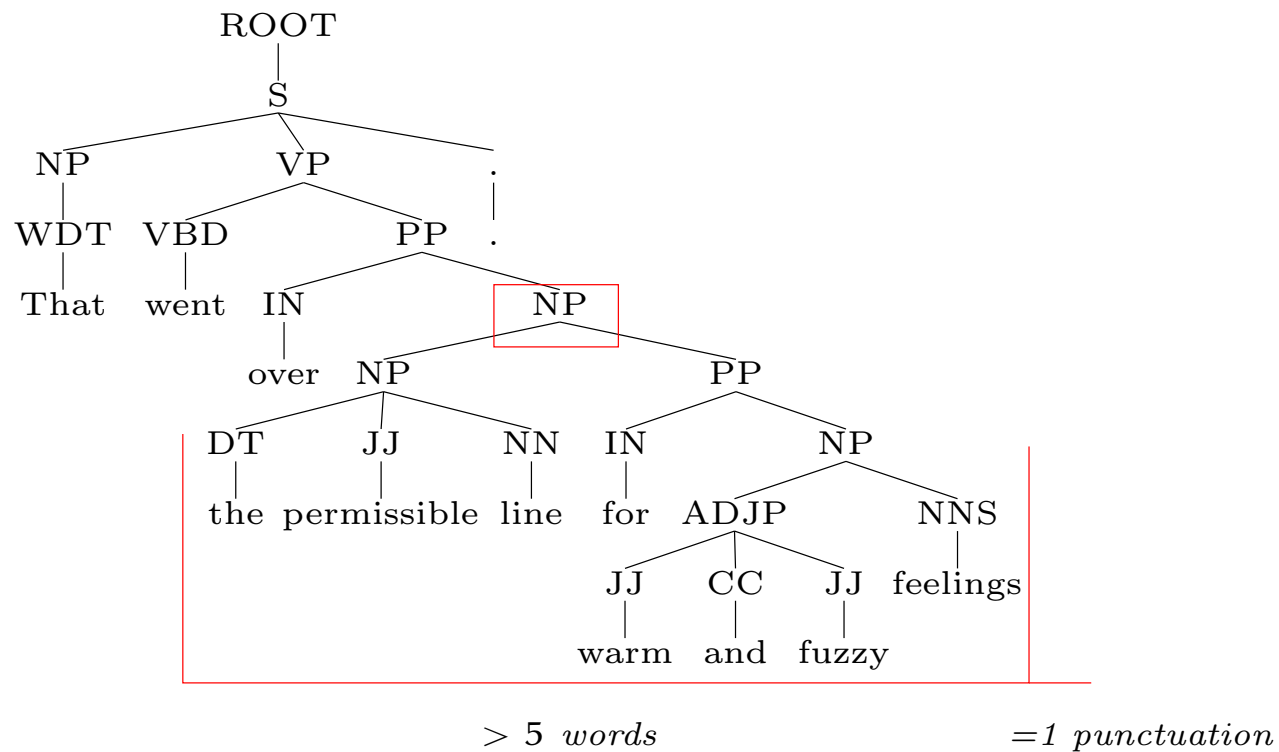
Rightmost branch bias

- The RightBranch feature's value is the number of nodes on the right-most branch (ignoring punctuation)
- Reflects the tendency toward right branching in English
- Only 2 different features, but very useful in final model!



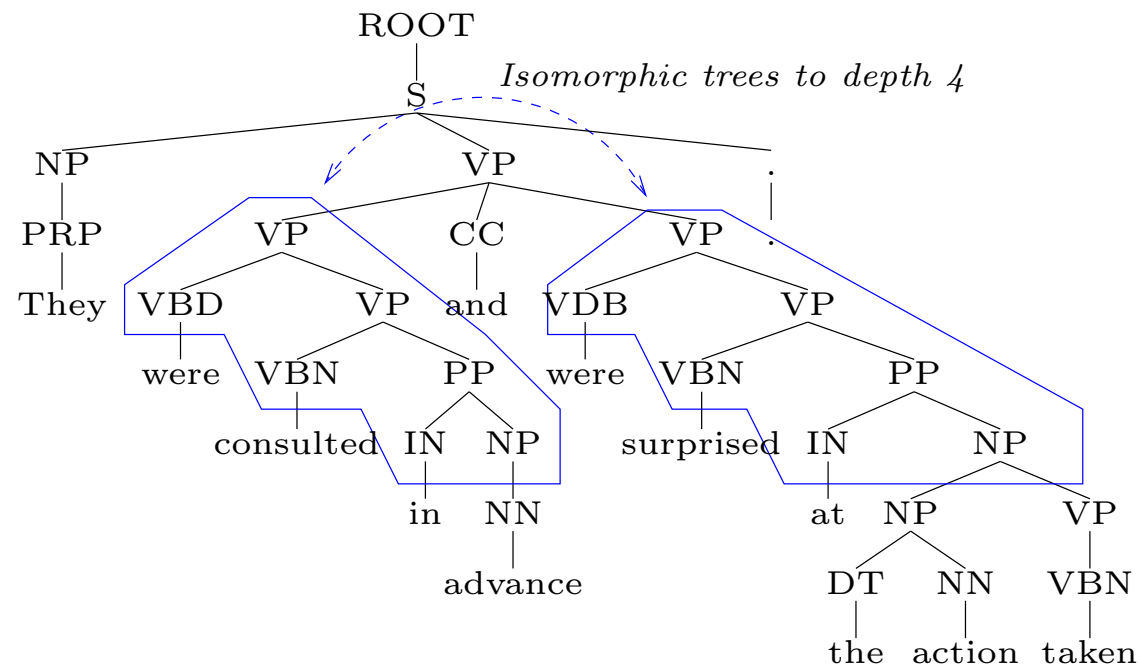
Constituent Heavyness and location

- Heavyness measures the constituent's category, its (binned) size and (binned) closeness to the end of the sentence



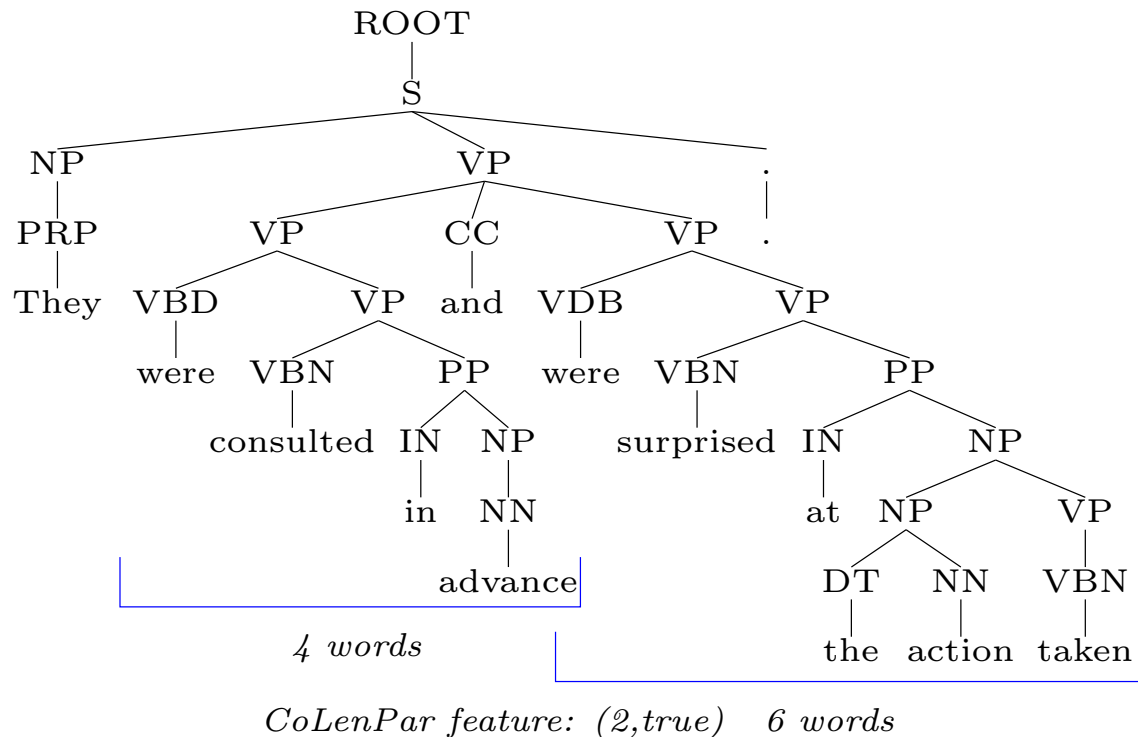
Coordination parallelism (1)

- The CoPar feature indicates the depth to which adjacent conjuncts are parallel



Coordination parallelism (2)

- The CoLenPar feature indicates the difference in length in adjacent conjuncts and whether this pair contains the last conjunct.

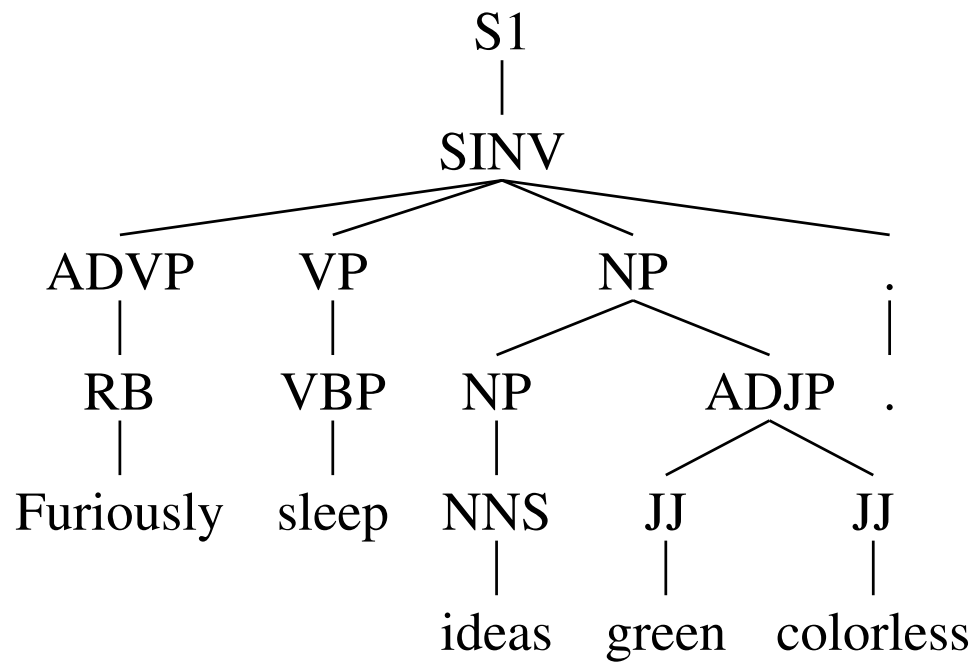
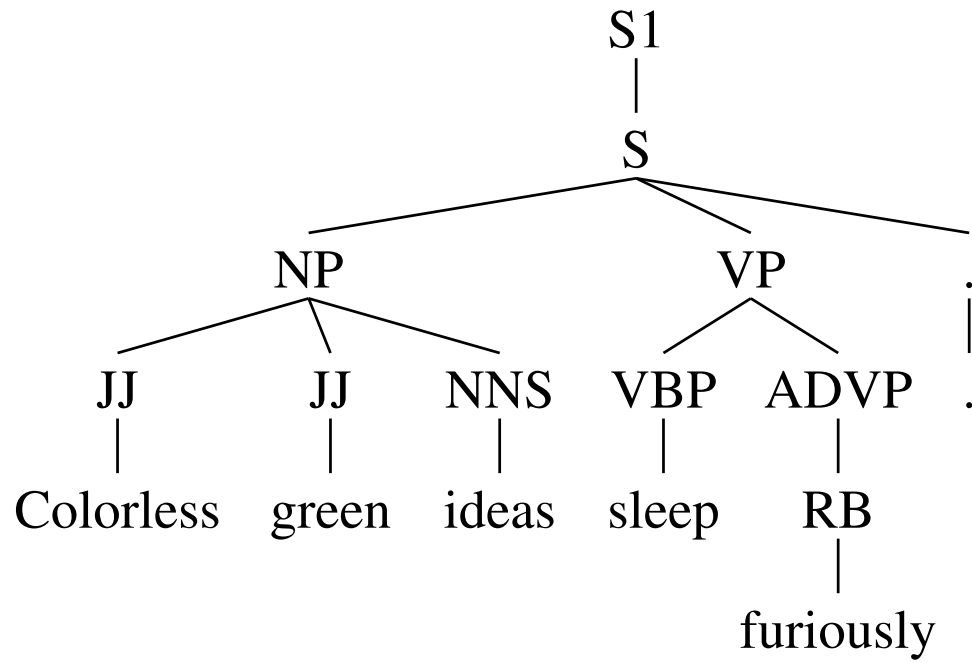


Experimental results with all features

- Features must vary on parses of at least 5 sentences in training data
- In this experiment, 724,550 features
- Gaussian regularization, adjusted via cross-validation on section 23
- *f-score on section 23* = 0.912 (15% error reduction over Charniak parser)

Conclusion

- It's possible to build (moderately) accurate, broad-coverage parsers
- Generative parsing models are easy to estimate, but make questionable independence assumptions
- Exponential models don't assume independence, so it's easy to add new features, but are difficult to estimate
- Coarse-to-fine conditional MLE for exponential models is a compromise
 - flexibility of exponential models
 - possible to estimate from treebank data
- Gives the currently best-reported parsing accuracy results



Sample parser errors

