
CS 388:
Natural Language Processing:
Statistical Parsing

Raymond J. Mooney
University of Texas at Austin

Statistical Parsing

- Statistical parsing uses a probabilistic model of syntax in order to assign probabilities to each parse tree.
- Provides principled approach to resolving syntactic ambiguity.
- Allows supervised learning of parsers from tree-banks of parse trees provided by human linguists.
- Also allows unsupervised learning of parsers from unannotated text, but the accuracy of such parsers has been limited.

Probabilistic Context Free Grammar (PCFG)

- A PCFG is a probabilistic version of a CFG where each production has a probability.
- Probabilities of all productions rewriting a given non-terminal must add to 1, defining a distribution for each non-terminal.
- String generation is now probabilistic where production probabilities are used to non-deterministically select a production for rewriting a given non-terminal.

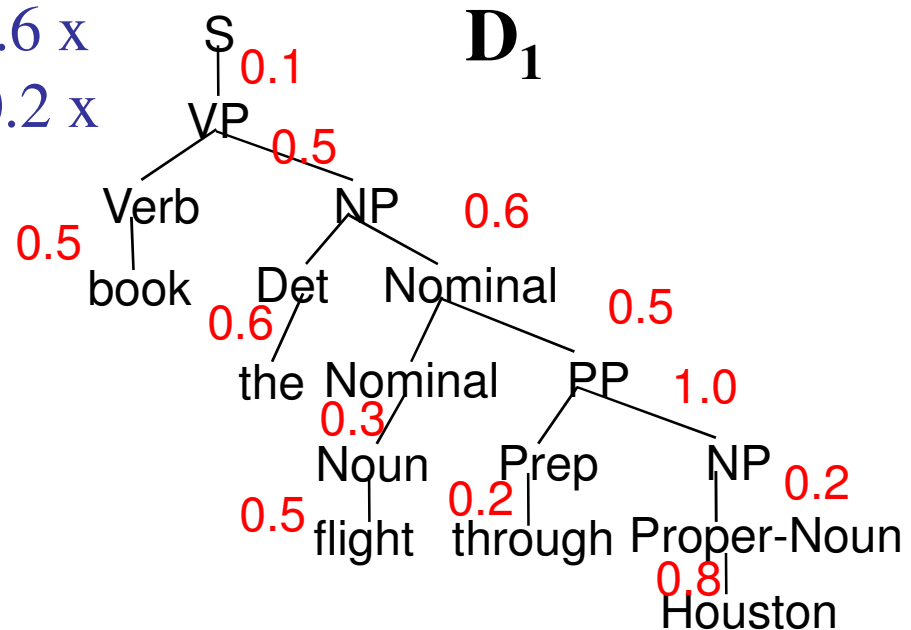
Simple PCFG for ATIS English

Grammar	Prob	Lexicon
S → NP VP	0.8	Det → the a that this
S → Aux NP VP	0.1	0.6 0.2 0.1 0.1
S → VP	0.1	Noun → book flight meal money
NP → Pronoun	0.2	0.1 0.5 0.2 0.2
NP → Proper-Noun	0.2	Verb → book include prefer
NP → Det Nominal	0.6	0.5 0.2 0.3
Nominal → Noun	0.3	Pronoun → I he she me
Nominal → Nominal Noun	0.2	0.5 0.1 0.1 0.3
Nominal → Nominal PP	0.5	Proper-Noun → Houston NWA
VP → Verb	0.2	0.8 0.2
VP → Verb NP	0.5	Aux → does
VP → VP PP	0.3	1.0
PP → Prep NP	1.0	Prep → from to on near through
		0.25 0.25 0.1 0.2 0.2

Sentence Probability

- Assume productions for each node are chosen independently.
- Probability of derivation is the product of the probabilities of its productions.

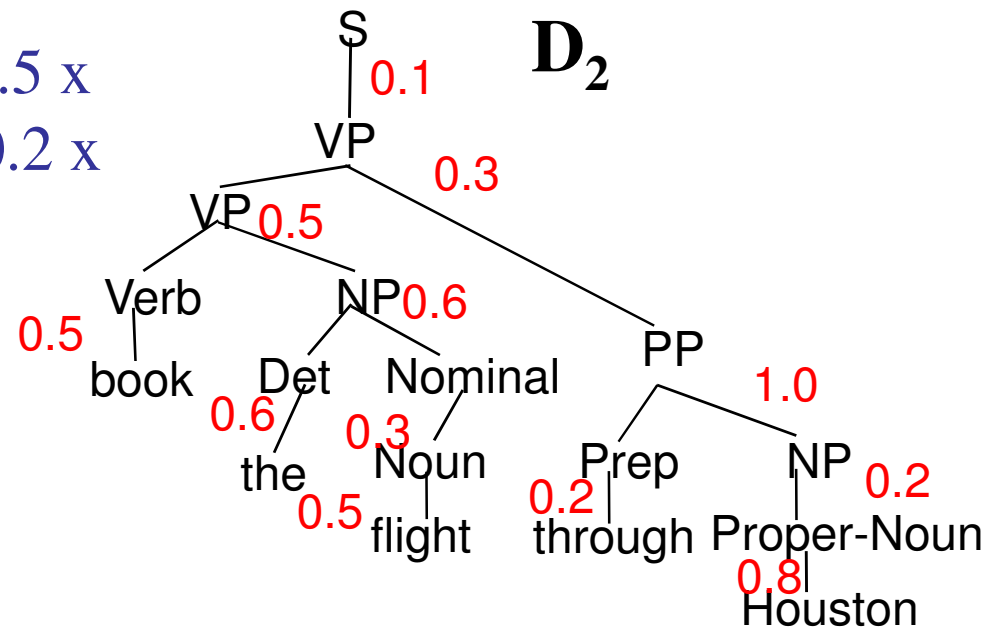
$$\begin{aligned} P(D_1) &= 0.1 \times 0.5 \times 0.5 \times 0.6 \times 0.6 \times \\ &\quad 0.5 \times 0.3 \times 1.0 \times 0.2 \times 0.2 \times \\ &\quad 0.5 \times 0.8 \\ &= 0.0000216 \end{aligned}$$



Syntactic Disambiguation

- Resolve ambiguity by picking most probable parse tree.

$$\begin{aligned} P(D_2) &= 0.1 \times 0.3 \times 0.5 \times 0.6 \times 0.5 \times \\ &\quad 0.6 \times 0.3 \times 1.0 \times 0.5 \times 0.2 \times \\ &\quad 0.2 \times 0.8 \\ &= 0.00001296 \end{aligned}$$



Sentence Probability

- Probability of a sentence is the sum of the probabilities of all of its derivations.

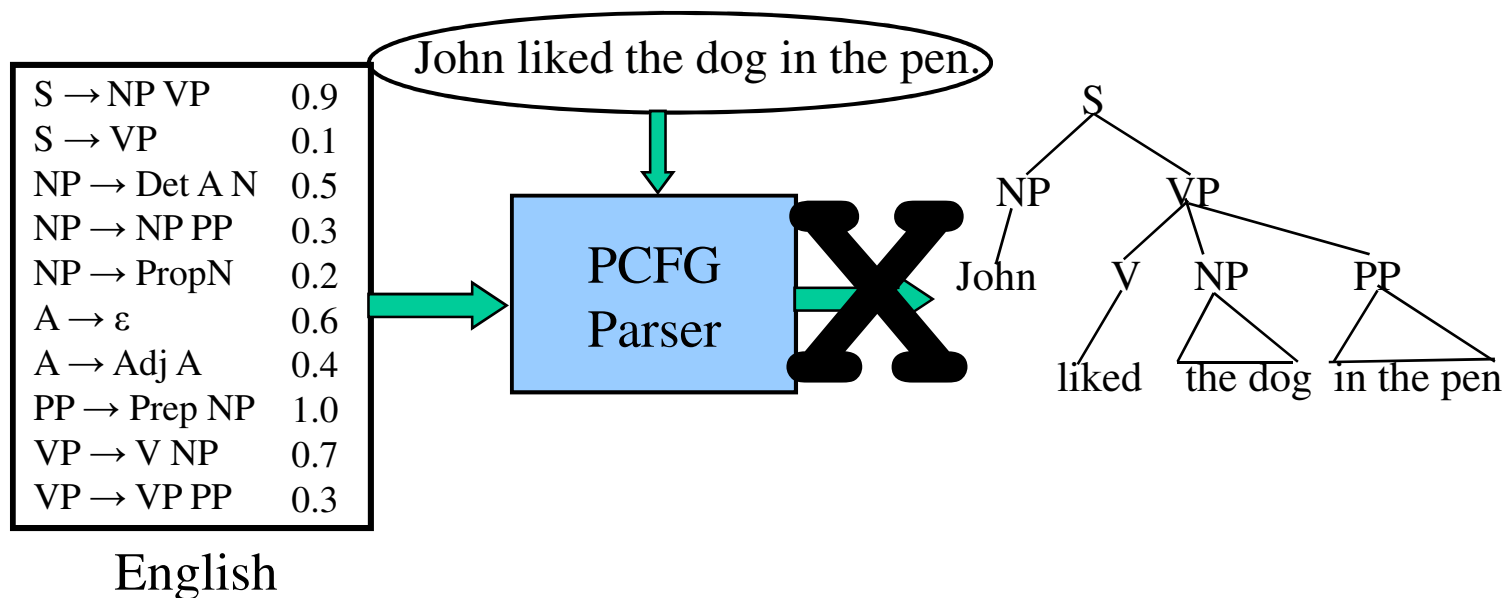
$$\begin{aligned} P(\text{"book the flight through Houston"}) &= \\ P(D_1) + P(D_2) &= 0.0000216 + 0.00001296 \\ &= 0.00003456 \end{aligned}$$

Three Useful PCFG Tasks

- **Observation likelihood:** To classify and order sentences.
- **Most likely derivation:** To determine the most likely parse tree for a sentence.
- **Maximum likelihood training:** To train a PCFG to fit empirical training data
 - We will not discuss this

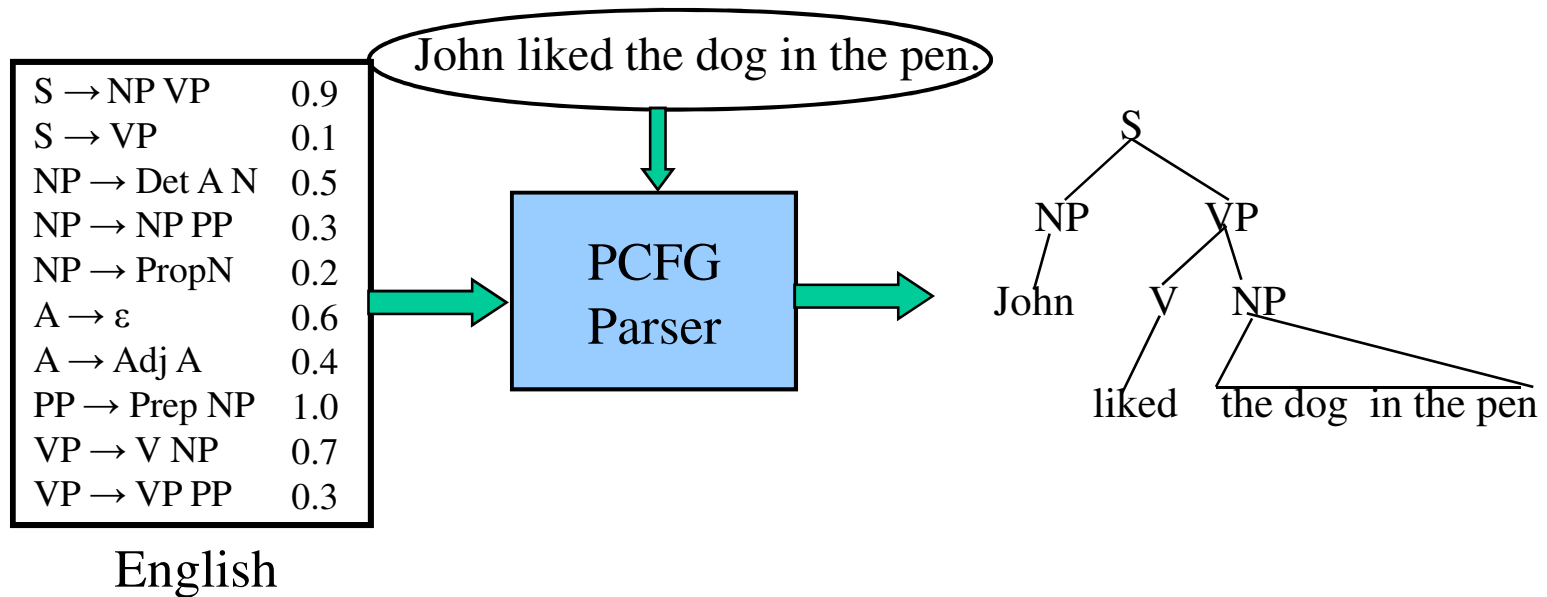
PCFG: Most Likely Derivation

- There is an analog to the Viterbi algorithm to efficiently determine the most probable derivation (parse tree) for a sentence.



PCFG: Most Likely Derivation

- There is an analog to the Viterbi algorithm to efficiently determine the most probable derivation (parse tree) for a sentence.



Probabilistic CKY

- CKY can be modified for PCFG parsing by including in each cell a probability for each non-terminal.
- Cell[i,j] must retain the *most probable* derivation of each constituent (non-terminal) covering words $i + 1$ through j together with its associated probability.
- When transforming the grammar to CNF, must set production probabilities to preserve the probability of derivations.

Probabilistic Grammar Conversion

Original Grammar

Chomsky Normal Form

S → NP VP	0.8	S → NP VP	0.8
S → Aux NP VP	0.1	S → X1 VP	0.1
		X1 → Aux NP	1.0
S → VP	0.1	S → book include prefer	
		0.01 0.004 0.006	
		S → Verb NP	0.05
		S → VP PP	0.03
NP → Pronoun	0.2	NP → I he she me	
		0.1 0.02 0.02 0.06	
NP → Proper-Noun	0.2	NP → Houston NWA	
		0.16 .04	
NP → Det Nominal	0.6	NP → Det Nominal	0.6
Nominal → Noun	0.3	Nominal → book flight meal money	
		0.03 0.15 0.06 0.06	
Nominal → Nominal Noun	0.2	Nominal → Nominal Noun	0.2
Nominal → Nominal PP	0.5	Nominal → Nominal PP	0.5
VP → Verb	0.2	VP → book include prefer	
		0.1 0.04 0.06	
VP → Verb NP	0.5	VP → Verb NP	0.5
VP → VP PP	0.3	VP → VP PP	0.3
PP → Prep NP	1.0	PP → Prep NP	1.0

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None			
	Det:.6 ←	NP:.6*.6*.15 =.054		
		Nominal:.15 Noun:.5		

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 ← Nominal:.03 Noun:.1	None	VP:.5*.5*.054 =.0135		
	Det:.6	NP:.6*.6*.15 =.054		
		Nominal:.15 Noun:.5		

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 ← Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135		
	Det:.6	↓ NP:.6*.6*.15 =.054		
		Nominal:.15 Noun:.5		

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6	NP:.6*.6*.15 =.054	None	
		Nominal:.15 Noun:.5	None	
			Prep: .2	

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6	NP:.6*.6*.15 =.054	None	
		Nominal:.15 Noun:.5	None	
			Prep: .2 ←	PP:1.0*.2*.16 =.032
				NP: .16 PropNoun: .8

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6	NP:.6*.6*.15 =.054	None	
		Nominal:.15 Noun:.5	← None	Nominal: .5*.15*.032 =.0024
			Prep: .2	↓ PP:1.0*.2*.16 =.032
				NP: .16 PropNoun: .8

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	
	Det:.6 ←	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep: .2	PP:1.0*.2*.16 =.032
				NP: .16 PropNoun: .8

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:. 5 Nominal:.03 Noun:.1		S:.05*.5*.054 =.00135		S:.05*.5* .000864 =.0000216
	None	VP:.5*.5*.054 =.0135	None	
	Det:.6	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep: .2	PP:1.0*.2*.16 =.032
				NP: .16 PropNoun: .8

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	S:.03*.0135* .032 =.00001296 S:.0000216
	Det:.6	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep: .2	PP:1.0*.2*.16 =.032
				NP: .16 PropNoun: .8

Probabilistic CKY Parser

Book the flight through Houston

S :.01, VP:.1, Verb:. \leftarrow .5 Nominal:.03 Noun:.1		S:. $.05*.5*.054$ =.00135		S:.0000216
	None	VP:. $.5*.5*.054$ =.0135	None	
				NP:. $.6*.6*$.0024 =.000864
	Det:. \leftarrow .6	NP:. $.6*.6*.15$ =.054	None	
				Nominal: .5*.15*.032 =.0024
		Nominal:.15 Noun:.5	None	
			Prep: .2 \leftarrow	PP:1.0*.2*.16 =.032
				NP: .16 PropNoun: .8

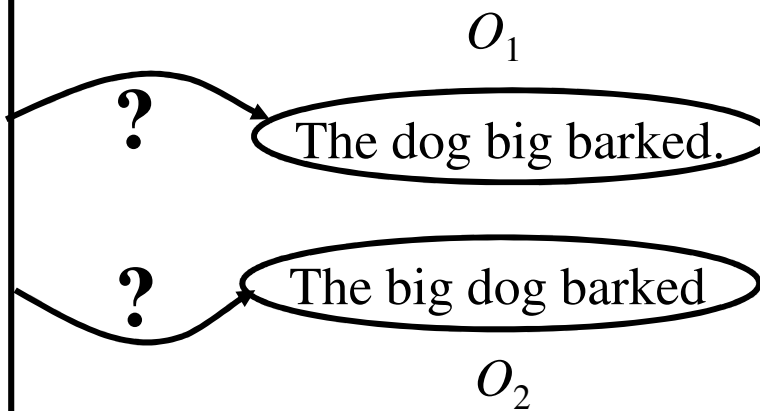
Pick most probable parse, i.e. take max to combine probabilities of multiple derivations of each constituent in each cell.

PCFG: Observation Likelihood

- There is an analog to Forward algorithm for HMMs called the **Inside algorithm** for efficiently determining how likely a string is to be produced by a PCFG.
- Can use a PCFG as a language model to choose between alternative sentences for speech recognition or machine translation.

S → NP VP	0.9
S → VP	0.1
NP → Det A N	0.5
NP → NP PP	0.3
NP → PropN	0.2
A → ε	0.6
A → Adj A	0.4
PP → Prep NP	1.0
VP → V NP	0.7
VP → VP PP	0.3

English



$P(O_2 | \text{English}) > P(O_1 | \text{English}) ?$

Inside Algorithm

- Use CKY probabilistic parsing algorithm but combine probabilities of multiple derivations of any constituent using **addition** instead of **max**.

Probabilistic CKY Parser for Inside Computation

Book the flight through Houston

S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	S:..00001296 S:.0000216
	Det:.6	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep:.2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

Probabilistic CKY Parser for Inside Computation

Book the flight through Houston

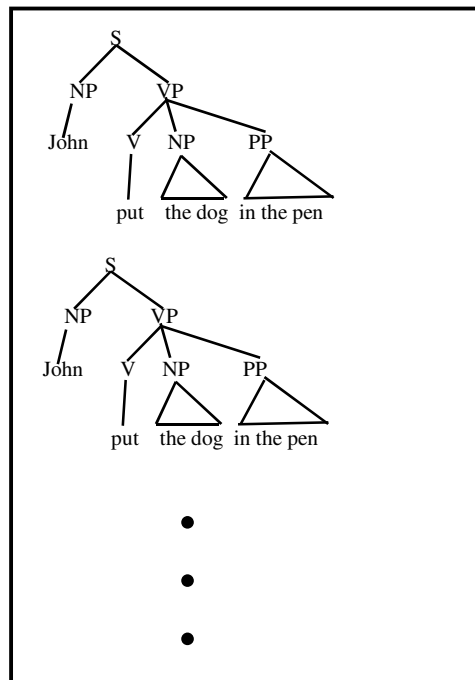
S :.01, VP:.1, Verb:.5 Nominal:.03 Noun:.1	None	S:.05*.5*.054 =.00135 VP:.5*.5*.054 =.0135	None	S: .00001296 +.0000216 =.00003456
	Det:.6	NP:.6*.6*.15 =.054	None	NP:.6*.6* .0024 =.000864
		Nominal:.15 Noun:.5	None	Nominal: .5*.15*.032 =.0024
			Prep:.2	PP:1.0*.2*.16 =.032
				NP:.16 PropNoun:.8

**Sum probabilities
of multiple derivations
of each constituent in
each cell.**

PCFG: Supervised Training

- If parse trees are provided for training sentences, a grammar and its parameters can all be estimated directly from counts accumulated from the **tree-bank** (with appropriate smoothing).

Tree Bank



Supervised
PCFG
Training

$S \rightarrow NP VP$	0.9
$S \rightarrow VP$	0.1
$NP \rightarrow Det A N$	0.5
$NP \rightarrow NP PP$	0.3
$NP \rightarrow PropN$	0.2
$A \rightarrow \epsilon$	0.6
$A \rightarrow Adj A$	0.4
$PP \rightarrow Prep NP$	1.0
$VP \rightarrow V NP$	0.7
$VP \rightarrow VP PP$	0.3

English

Estimating Production Probabilities

- Set of production rules can be taken directly from the set of rewrites in the treebank.
- Parameters can be directly estimated from frequency counts in the treebank.

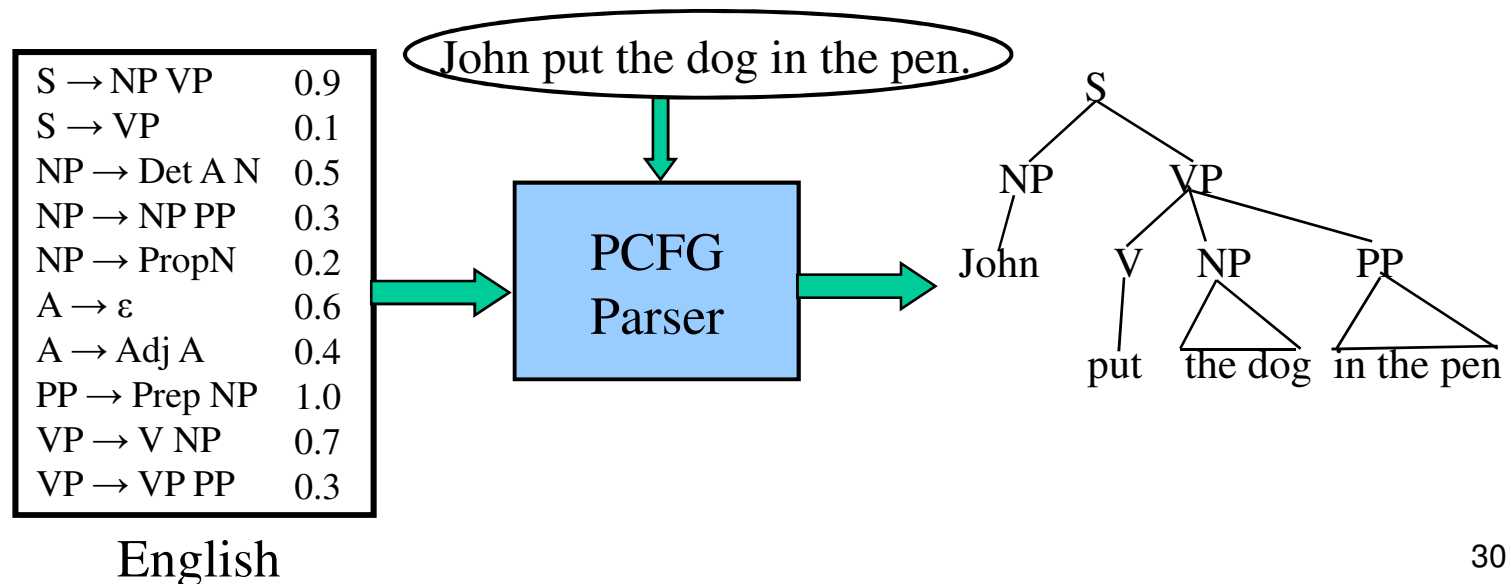
$$P(\alpha \rightarrow \beta | \alpha) = \frac{\text{count}(\alpha \rightarrow \beta)}{\sum_{\gamma} \text{count}(\alpha \rightarrow \gamma)} = \frac{\text{count}(\alpha \rightarrow \beta)}{\text{count}(\alpha)}$$

Vanilla PCFG Limitations

- Since probabilities of productions do not rely on specific words or concepts, only general structural disambiguation is possible (e.g. prefer to attach PPs to Nominals).
- Consequently, vanilla PCFGs cannot resolve syntactic ambiguities that require semantics to resolve, e.g. ate with fork vs. meatballs.
- In order to work well, PCFGs must be **lexicalized**, i.e. productions must be specialized to specific words by including their head-word in their LHS non-terminals (e.g. VP-ate).

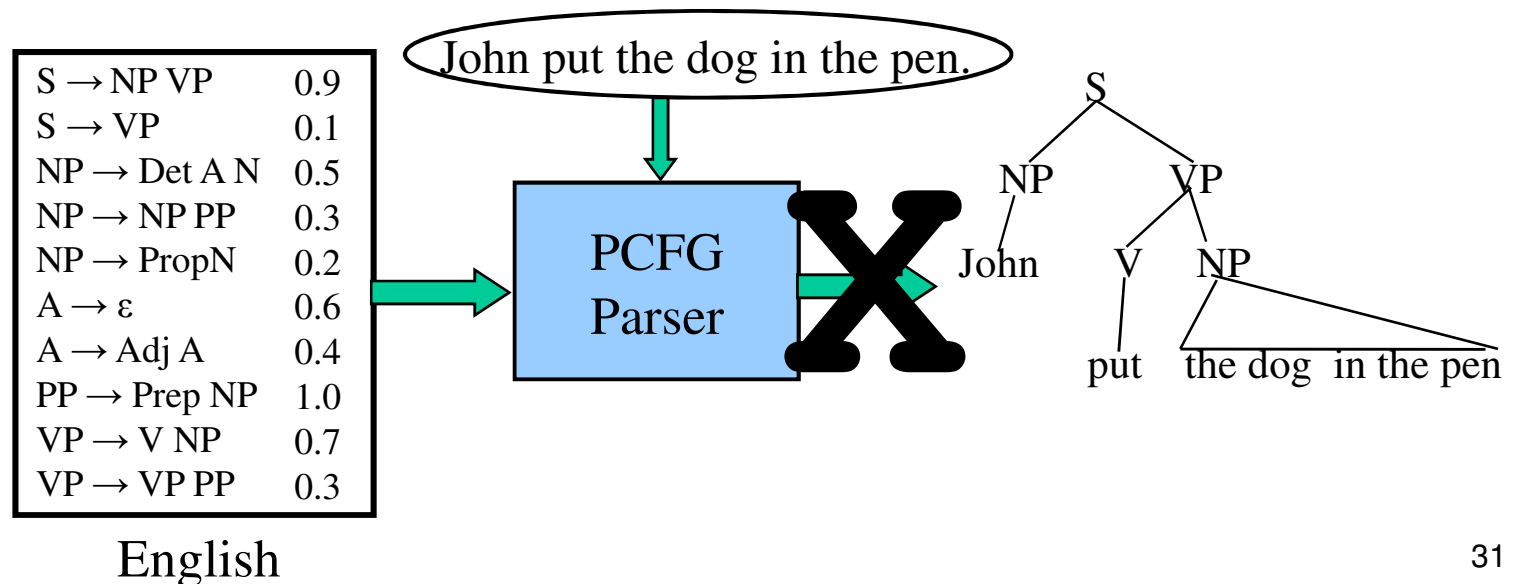
Example of Importance of Lexicalization

- A general preference for attaching PPs to NPs rather than VPs can be learned by a vanilla PCFG.
- But the desired preference can depend on specific words.



Example of Importance of Lexicalization

- A general preference for attaching PPs to NPs rather than VPs can be learned by a vanilla PCFG.
- But the desired preference can depend on specific words.

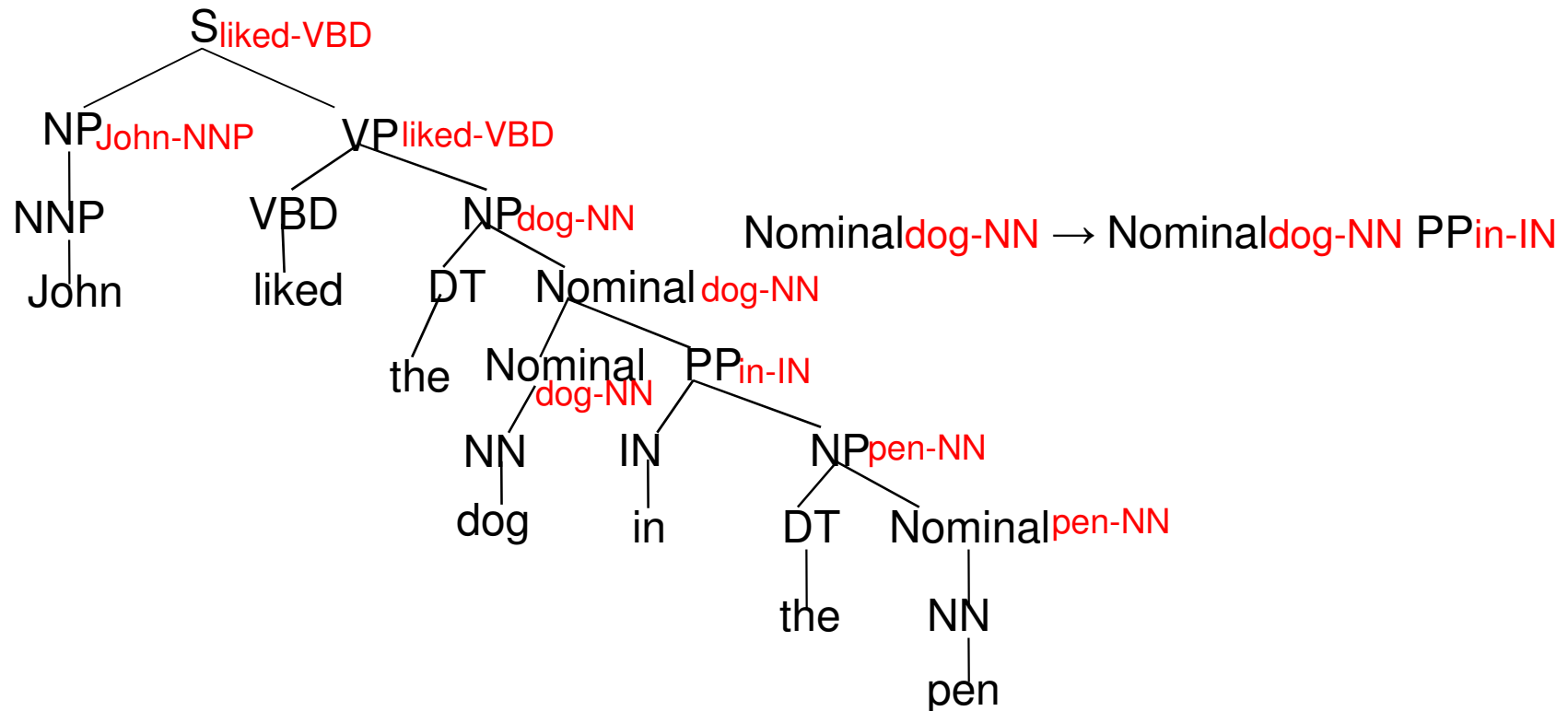


Head Words

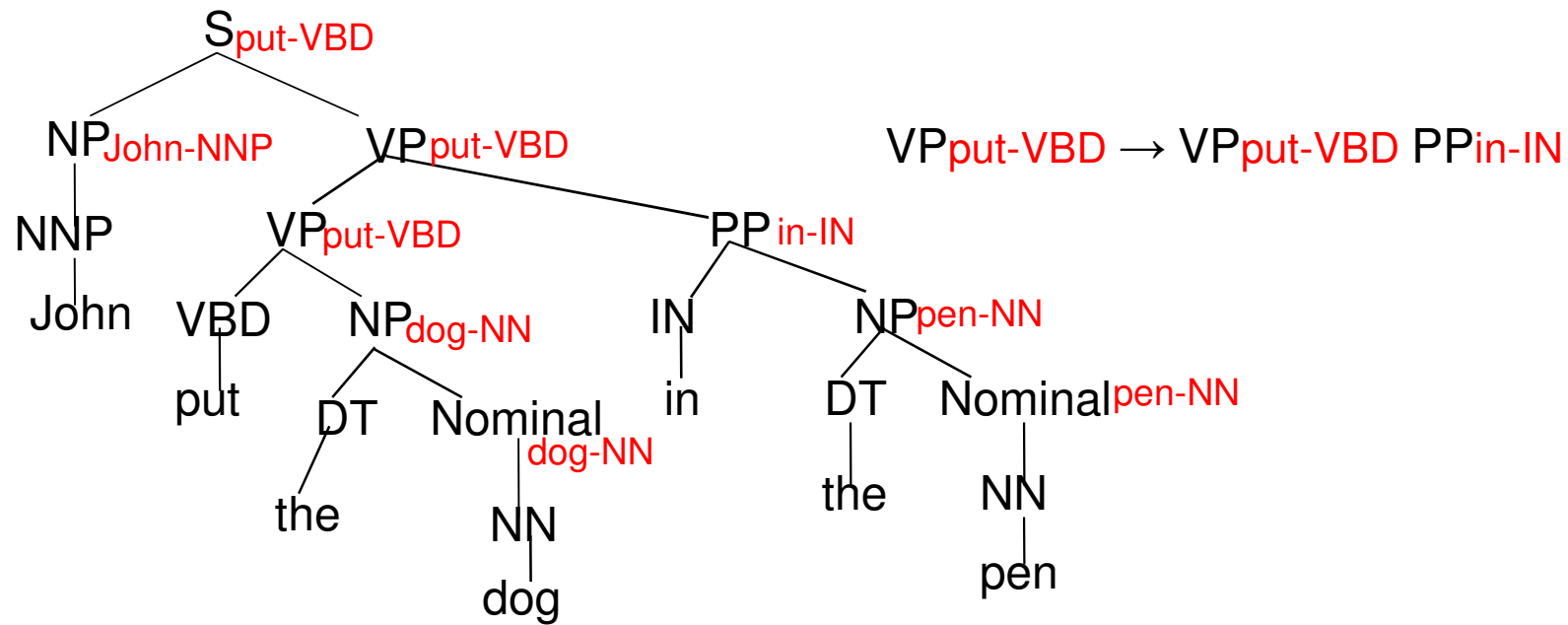
- Syntactic phrases usually have a word in them that is most “central” to the phrase.
- Linguists have defined the concept of a lexical **head** of a phrase.
- Simple rules can identify the head of any phrase by percolating head words up the parse tree.
 - Head of a VP is the main verb
 - Head of an NP is the main noun
 - Head of a PP is the preposition
 - Head of a sentence is the head of its VP

Lexicalized Productions

- Specialized productions can be generated by including the head word and its POS of each non-terminal as part of that non-terminal's symbol.



Lexicalized Productions



Parameterizing Lexicalized Productions

- Accurately estimating parameters on such a large number of very specialized productions could require enormous amounts of treebank data.
- Need some way of estimating parameters for lexicalized productions that makes reasonable independence assumptions so that accurate probabilities for very specific rules can be learned.
- See Collins Parser: Chapter 14.6
 - (we do not cover it in this course)