
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

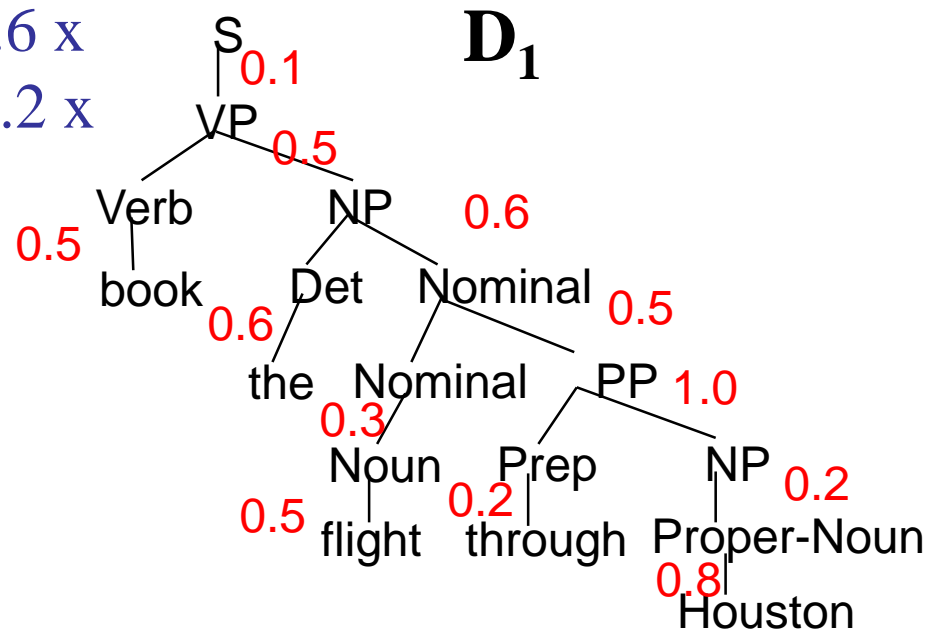
Grammar	Prob
$S \rightarrow NP VP$	0.8
$S \rightarrow Aux NP VP$	0.1
$S \rightarrow VP$	0.1
$NP \rightarrow Pronoun$	0.2
$NP \rightarrow Proper-Noun$	0.2
$NP \rightarrow Det Nominal$	0.6
$Nominal \rightarrow Noun$	0.3
$Nominal \rightarrow Nominal Noun$	0.2
$Nominal \rightarrow Nominal PP$	0.5
$VP \rightarrow Verb$	0.2
$VP \rightarrow Verb NP$	0.5
$VP \rightarrow VP PP$	0.3
$PP \rightarrow Prep NP$	1.0

Lexicon
$Det \rightarrow the \mid a \mid that \mid this$ 0.6 0.2 0.1 0.1
$Noun \rightarrow book \mid flight \mid meal \mid money$ 0.1 0.5 0.2 0.2
$Verb \rightarrow book \mid include \mid prefer$ 0.5 0.2 0.3
$Pronoun \rightarrow I \mid he \mid she \mid me$ 0.5 0.1 0.1 0.3
$Proper-Noun \rightarrow Houston \mid NWA$ 0.8 0.2
$Aux \rightarrow does$ 1.0
$Prep \rightarrow from \mid to \mid on \mid near \mid through$ 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.

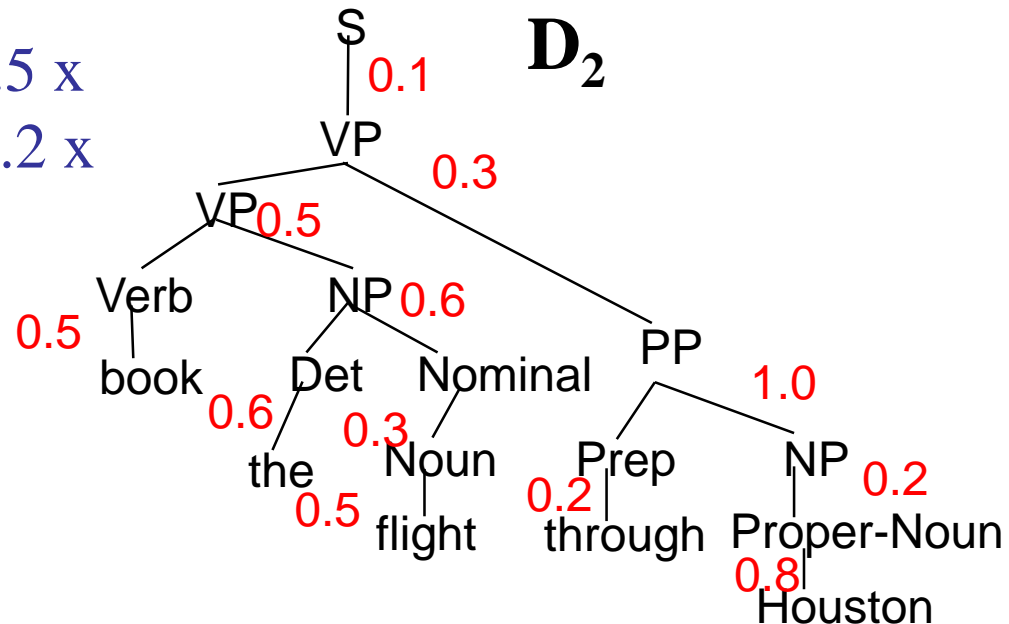
$$\begin{aligned} P(D_1) &= 0.1 \times 0.5 \times 0.5 \times 0.6 \times 0.6 \times \\ &\quad 0.5 \times 0.3 \times 1.0 \times 0.2 \times 0.2 \times \\ &\quad 0.5 \times 0.8 \\ &= 0.0000216 \end{aligned}$$



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.

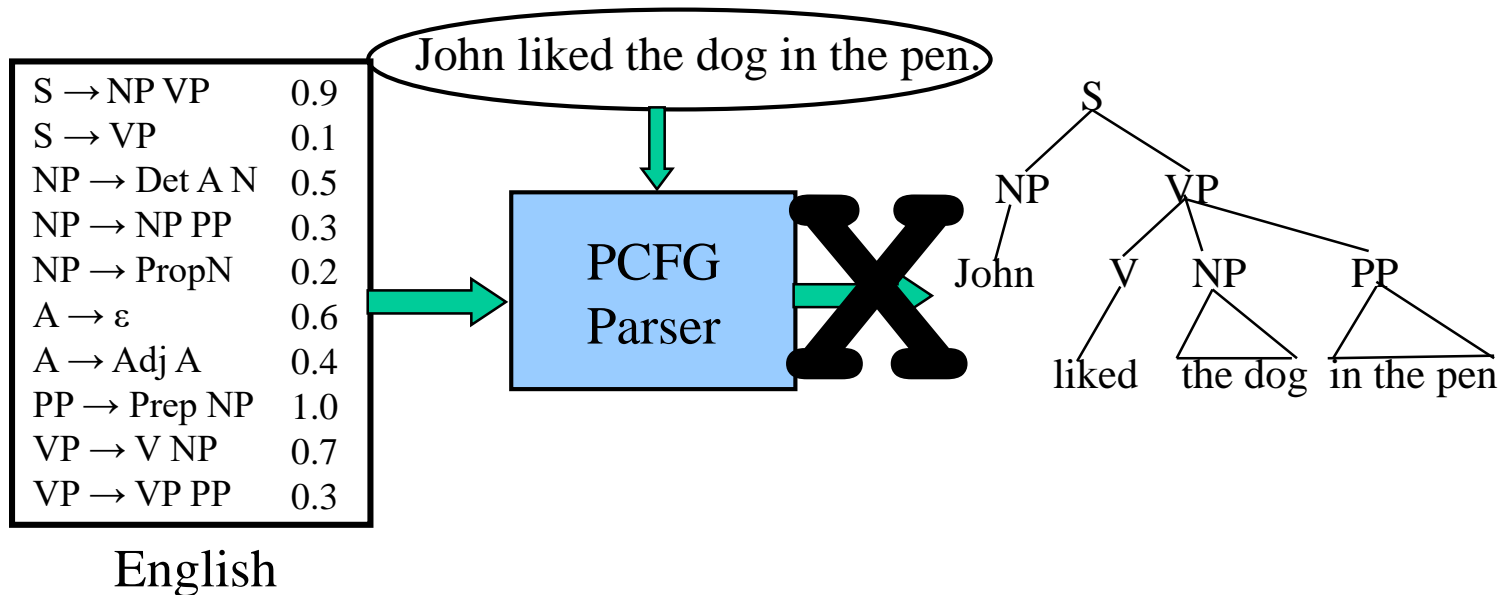
$$\begin{aligned} P(\text{"book the flight through Houston"}) &= \\ P(D_1) + P(D_2) &= 0.0000216 + 0.00001296 \\ &= 0.00003456 \end{aligned}$$

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.

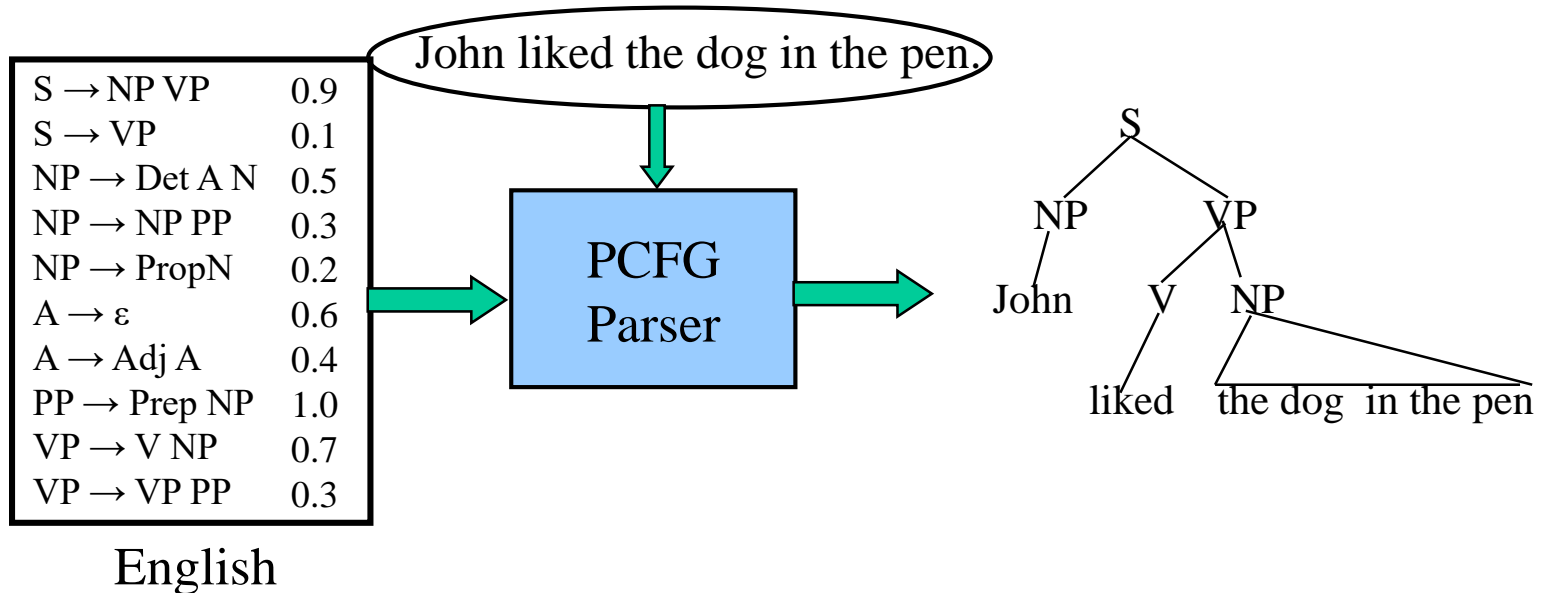
PCFG: Most Likely Derivation

- 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.
- $\text{Cell}[i,j]$ 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	
				NP:.6*.6* .0024 =.000864
	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	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

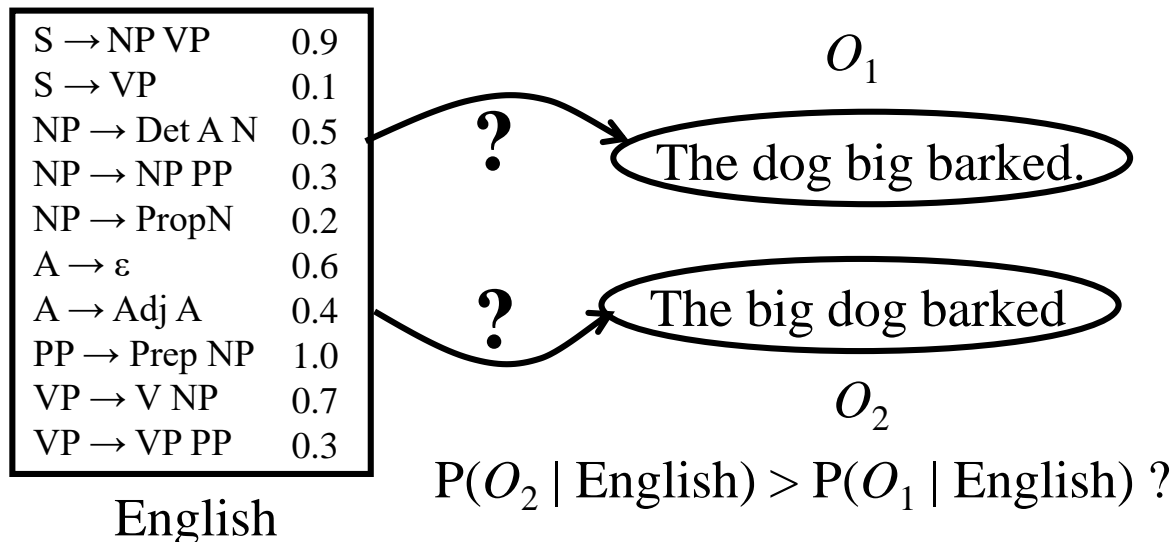
Book the flight through Houston

S :.01, VP:~.1, Verb:~.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:~.6	NP:~.6*.6*~.15 =.054	None	
				Nominal: ~.5*.15*.032 =.0024
		Nominal:~.15 Noun:~.5	None	
			Prep:~.2	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.



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

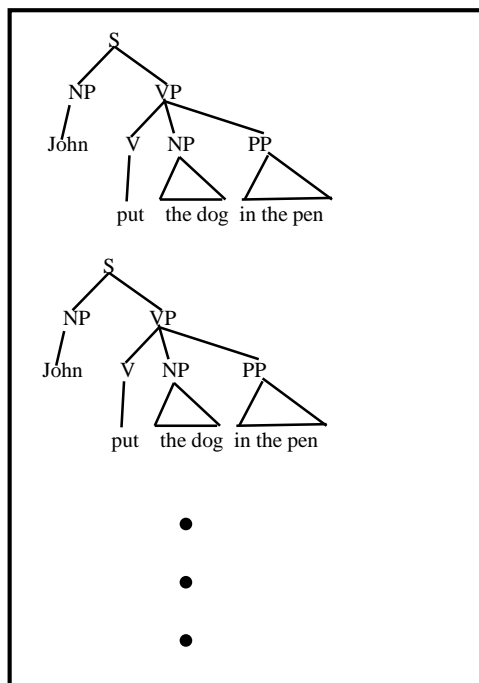
S :.01, VP:~.1, Verb:~.5 Nominal:~.03 Noun:~.1	None	S:~.05*~.5*~.054 =.00135 VP:~.5*~.5*~.054 =.0135	None	S: ~.00001296 +.0000216 =.00003456
	Det:~.6	NP:~.6*~.6*~.15 =.054	None	NP:~.6*~.6* .0024 =.000864
		Nominal:~.15 Noun:~.5	None	Nominal: ~.5*~.15*~.032 =.0024
			Prep:~.2	PP:~.1*~.2*~.16 =.032
		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

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

PCFG: Maximum Likelihood Training

- Given a set of sentences, induce a grammar that maximizes the probability that this data was generated from this grammar.
- Assume the number of non-terminals in the grammar is specified.
- Only need to have an unannotated set of sequences generated from the model. Does not need correct parse trees for these sentences. In this sense, it is **unsupervised**.

PCFG: Maximum Likelihood Training

Training Sentences

John ate the apple
A dog bit Mary
Mary hit the dog
John gave Mary the cat.
•
•
•

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

Inside-Outside

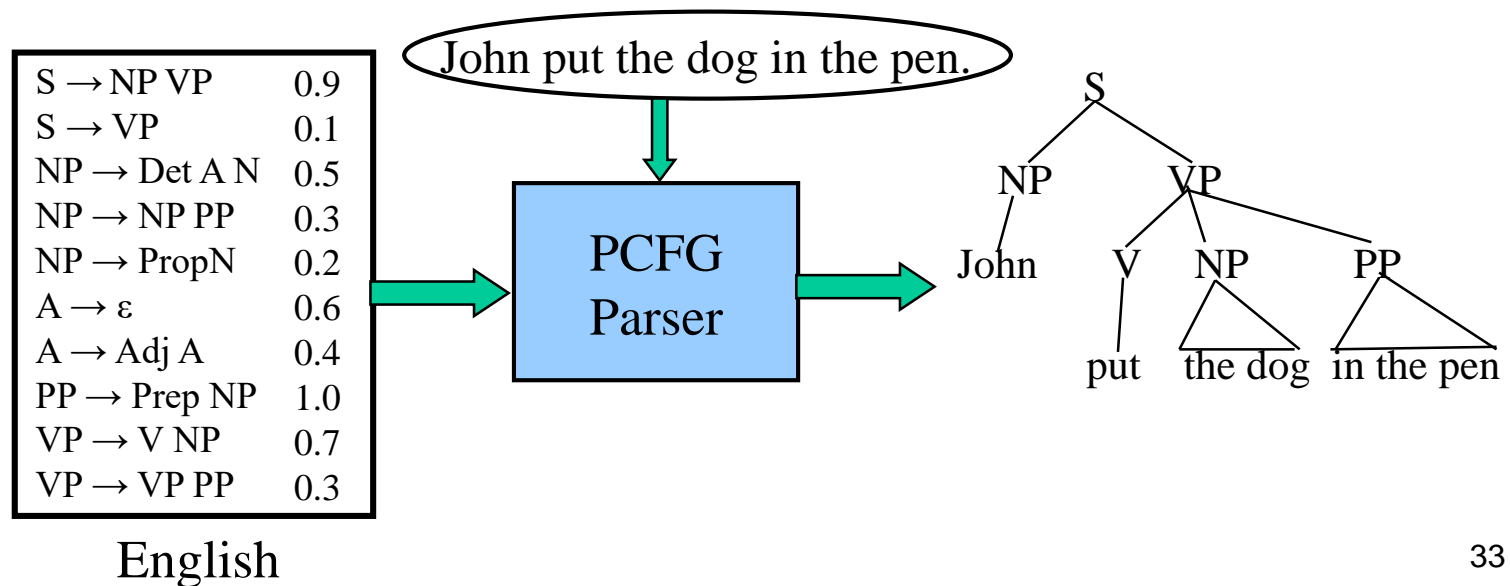
- The **Inside-Outside algorithm** is a version of EM for unsupervised learning of a PCFG.
 - Analogous to Baum-Welch (forward-backward) for HMMs
- Given the number of non-terminals, construct all possible CNF productions with these non-terminals and observed terminal symbols.
- Use EM to iteratively train the probabilities of these productions to locally maximize the likelihood of the data.
 - See Manning and Schütze text for details
- Experimental results are not impressive, but recent work imposes additional constraints to improve unsupervised grammar learning.

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).

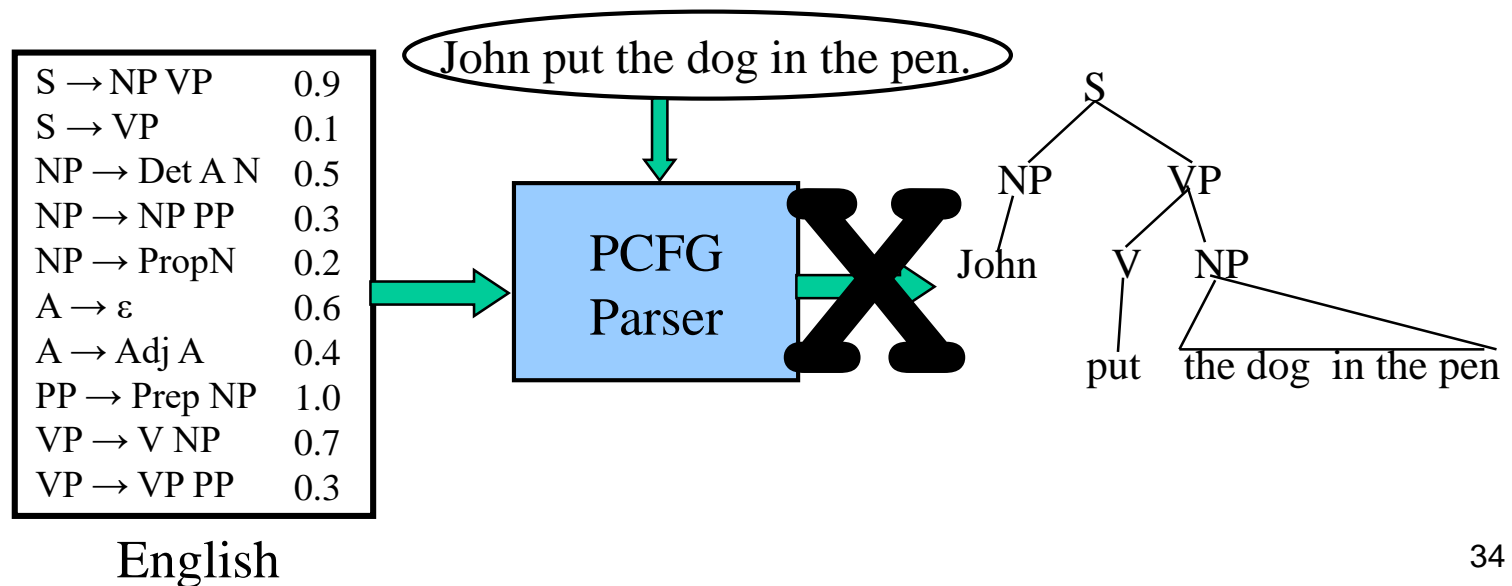
Example of Importance of Lexicalization

- 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.

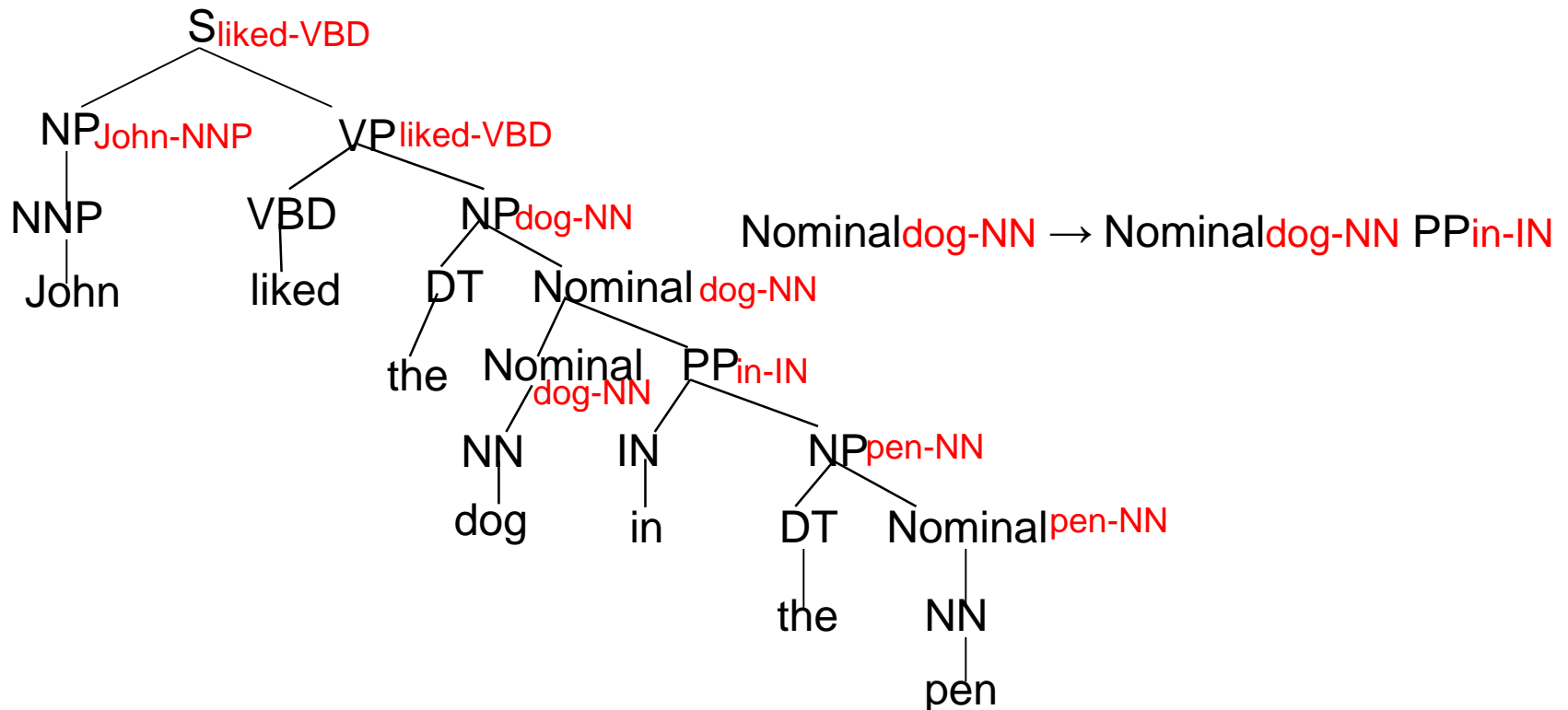


Head Words

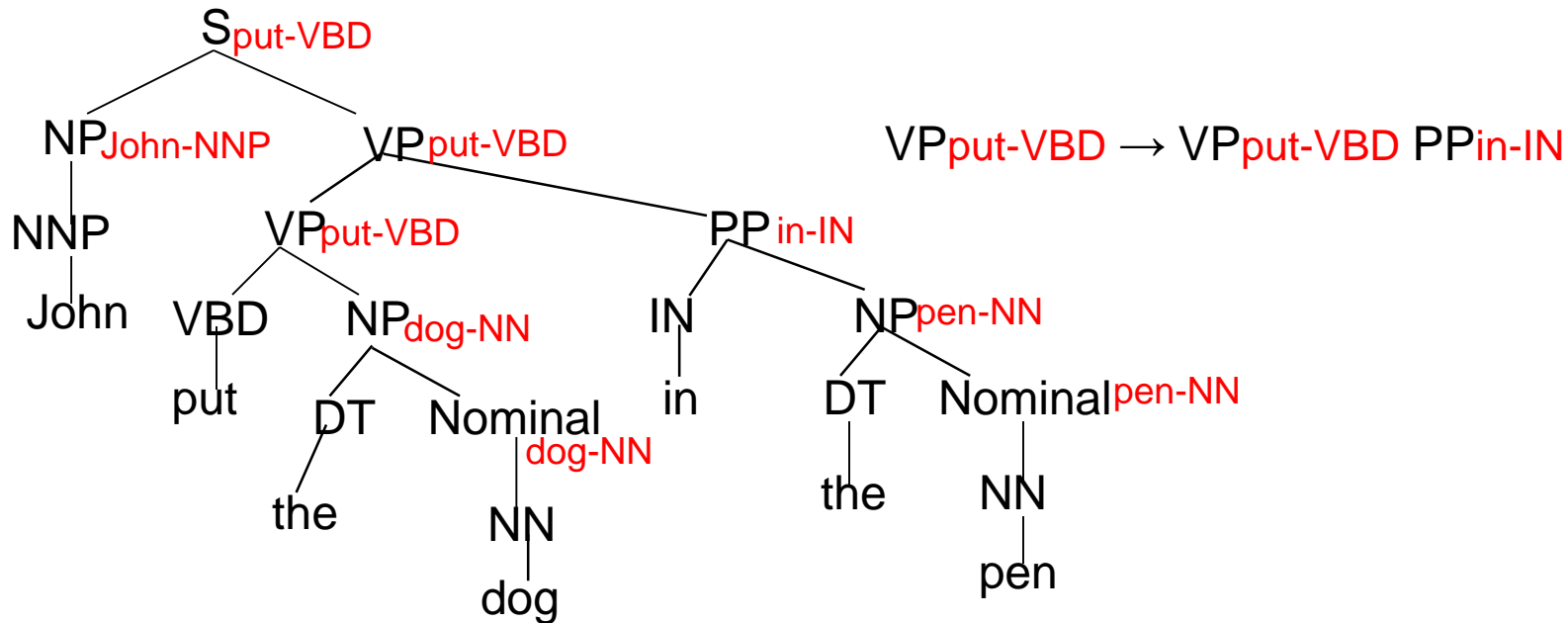
- Syntactic phrases usually have a word in them that is most “central” to the phrase.
- Linguists have defined the concept of a lexical **head** of a phrase.
- Simple rules can identify the head of any phrase by percolating head words up the parse tree.
 - Head of a VP is the main verb
 - Head of an 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.
- Collins (1999) introduced one approach to learning effective parameters for a lexicalized grammar.

Treebanks

- **English Penn Treebank**: Standard corpus for testing syntactic parsing consists of 1.2 M words of text from the Wall Street Journal (WSJ).
- Typical to train on about 40,000 parsed sentences and test on an additional standard disjoint test set of 2,416 sentences.
- **Chinese Penn Treebank**: 100K words from the Xinhua news service.
- Other corpora existing in many languages, see the Wikipedia article “Treebank”

First WSJ Sentence

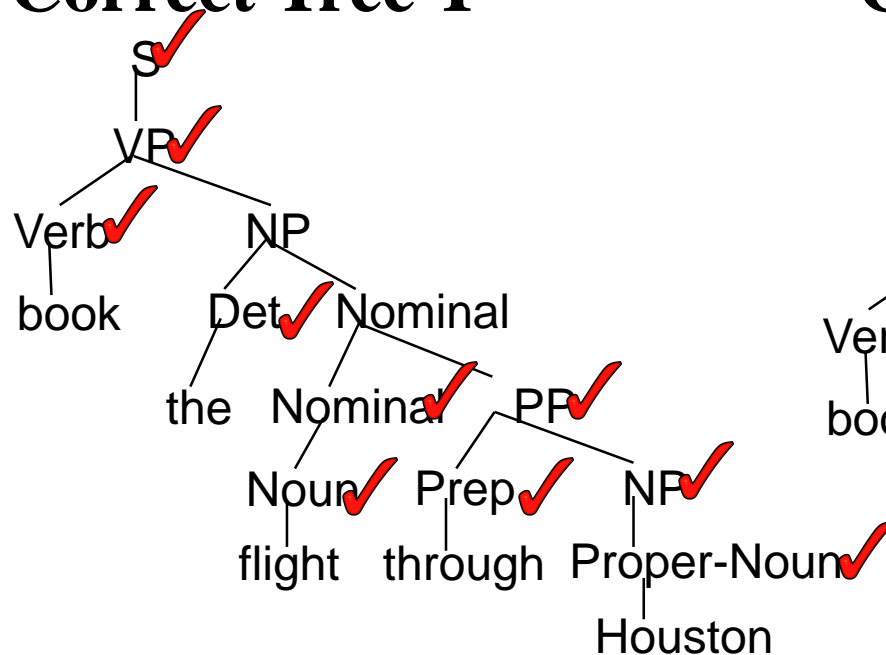
((S
 (NP-SBJ
 (NP (NNP Pierre) (NNP Vinken))
 (, ,)
 (ADJP
 (NP (CD 61) (NNS years))
 (JJ old))
 (, ,))
 (VP (MD will)
 (VP (VB join)
 (NP (DT the) (NN board))
 (PP-CLR (IN as)
 (NP (DT a) (JJ nonexecutive) (NN director)))
 (NP-TMP (NNP Nov.) (CD 29))))
 (. .)))

Parsing Evaluation Metrics

- PARSEVAL metrics measure the fraction of the constituents that match between the computed and human parse trees. If P is the system's parse tree and T is the human parse tree (the “gold standard”):
 - **Recall** = (# correct constituents in P) / (# constituents in T)
 - **Precision** = (# correct constituents in P) / (# constituents in P)
- **Labeled Precision** and **labeled recall** require getting the non-terminal label on the constituent node correct to count as correct.
- **F_1** is the harmonic mean of precision and recall.

Computing Evaluation Metrics

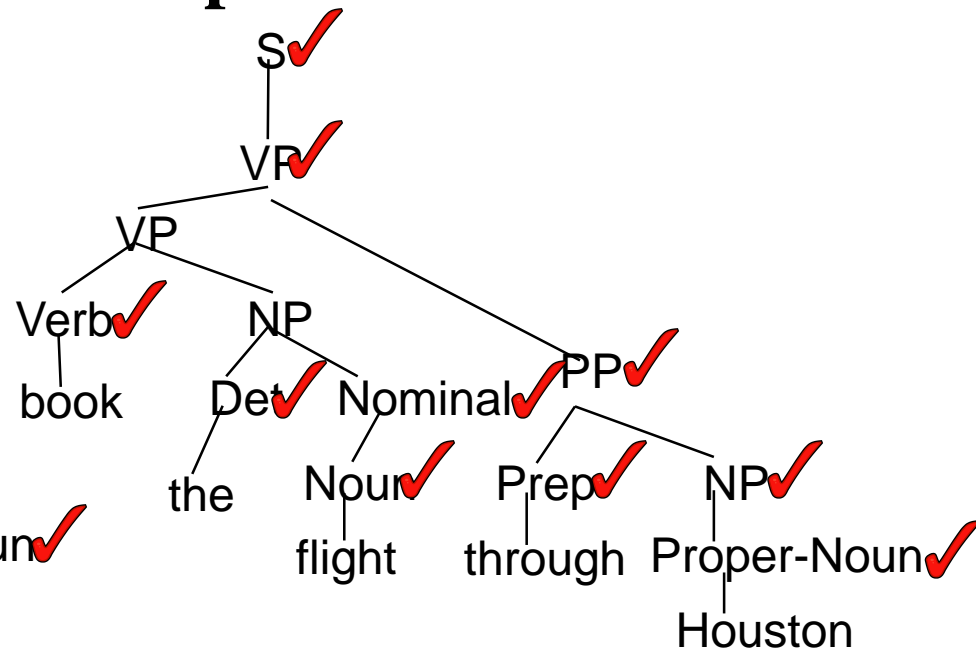
Correct Tree T



Constituents: 12

Correct Constituents: 10

Computed Tree P



Constituents: 12

Recall = $10/12 = 83.3\%$ Precision = $10/12 = 83.3\%$

$F_1 = 83.3\%$

Treebank Results

- Results of current state-of-the-art systems on the English Penn WSJ treebank are 91-92% labeled F_1 .

Human Parsing

- Computational parsers can be used to predict human reading time as measured by tracking the time taken to read each word in a sentence.
- Psycholinguistic studies show that words that are more probable given the preceding lexical and syntactic context are read faster.
 - John put the dog in the pen with a lock.
 - John put the dog in the pen with a bone in the car.
 - John liked the dog in the pen with a bone.
- Modeling these effects requires an *incremental* statistical parser that incorporates one word at a time into a continuously growing parse tree.

Garden Path Sentences

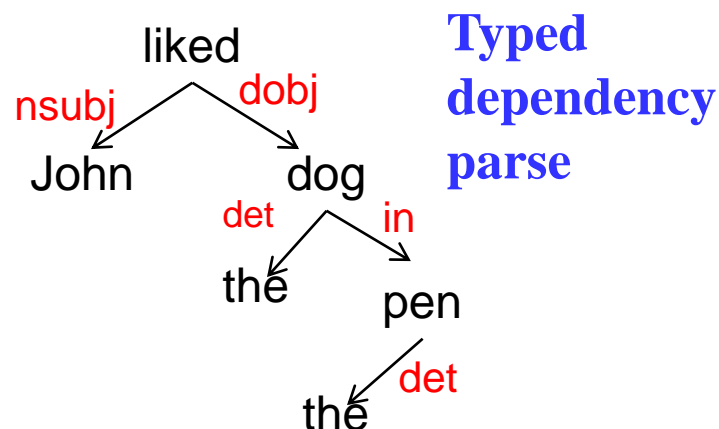
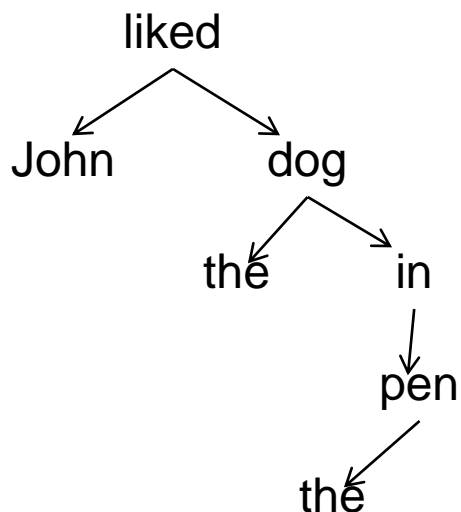
- People are confused by sentences that seem to have a particular syntactic structure but then suddenly violate this structure, so the listener is “lead down the garden path”.
 - The horse raced past the barn fell.
 - vs. The horse raced past the barn broke his leg.
 - The complex houses married students.
 - The old man the sea.
 - While Anna dressed the baby spit up on the bed.
- Incremental computational parsers can try to predict and explain the problems encountered parsing such sentences.

Center Embedding

- Nested expressions are hard for humans to process beyond 1 or 2 levels of nesting.
 - The rat the cat chased died.
 - The rat the cat the dog bit chased died.
 - The rat the cat the dog the boy owned bit chased died.
- Requires remembering and popping incomplete constituents from a stack and strains human short-term memory.
- Equivalent “tail embedded” (tail recursive) versions are easier to understand since no stack is required.
 - The boy owned a dog that bit a cat that chased a rat that died.

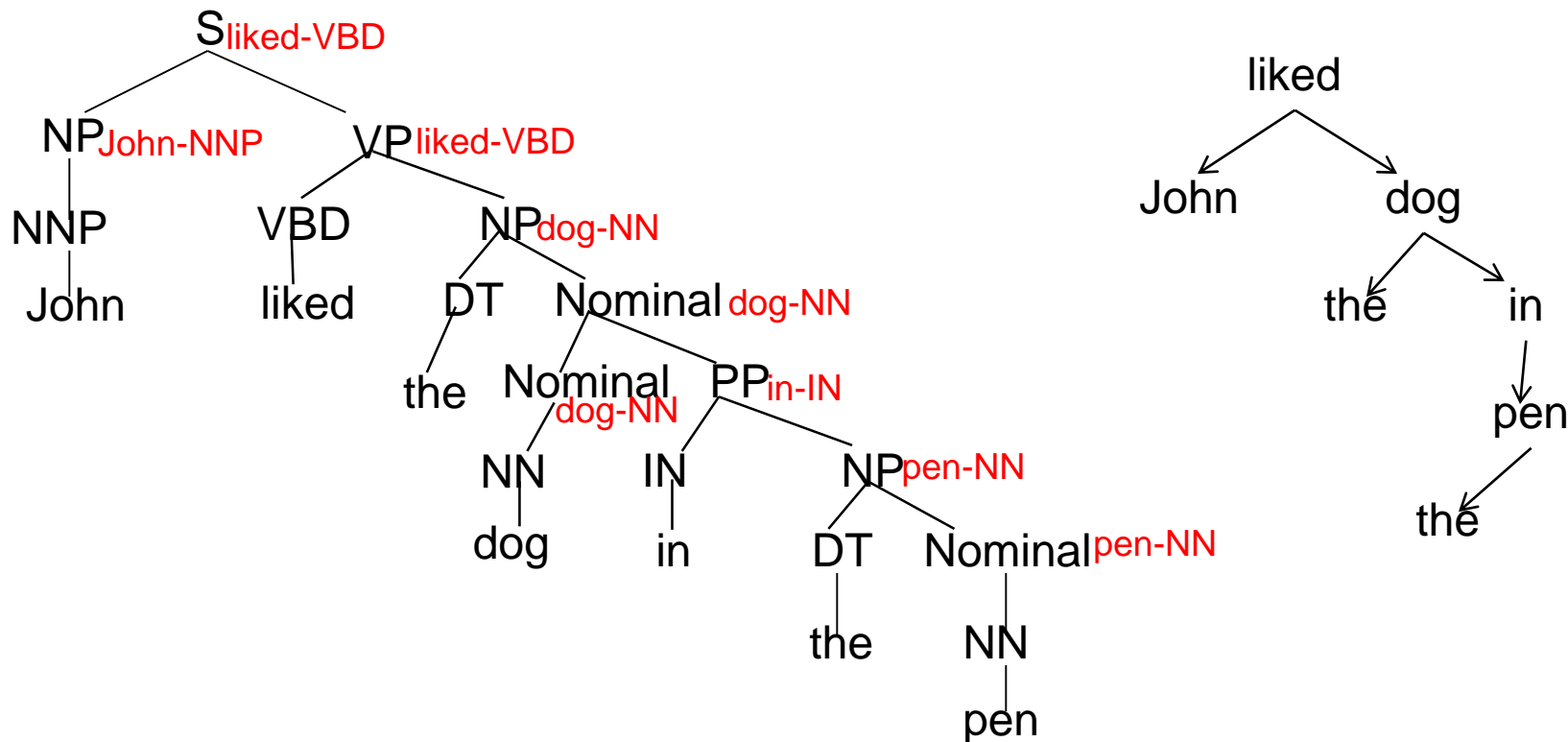
Dependency Grammars

- An alternative to phrase-structure grammar is to define a parse as a directed graph between the words of a sentence representing *dependencies* between the words.



Dependency Graph from Parse Tree

- Can convert a phrase structure parse to a dependency tree by making the head of each non-head child of a node depend on the head of the head child.



Unification Grammars

- In order to handle agreement issues more effectively, each constituent has a list of features such as number, person, gender, etc. which may or not be specified for a given constituent.
- In order for two constituents to combine to form a larger constituent, their features must *unify*, i.e. consistently combine into a merged set of features.
- Expressive grammars and parsers (e.g. HPSG) have been developed using this approach and have been partially integrated with modern statistical models of disambiguation.

Mildly Context-Sensitive Grammars

- Some grammatical formalisms provide a degree of context-sensitivity that helps capture aspects of NL syntax that are not easily handled by CFGs.
- Tree Adjoining Grammar (TAG) is based on combining tree fragments rather than individual phrases.
- Combinatory Categorical Grammar (CCG) consists of:
 - **Categorical Lexicon** that associates a syntactic and semantic category with each word.
 - **Combinatory Rules** that define how categories combine to form other categories.

Statistical Parsing Conclusions

- Statistical models such as PCFGs allow for probabilistic resolution of ambiguities.
- PCFGs can be easily learned from treebanks.
- Lexicalization and non-terminal splitting are required to effectively resolve many ambiguities.
- Current statistical parsers are quite accurate but not yet at the level of human-expert agreement.