

Probabilistic Context-Free Grammar

By:

B.Senthil Kumar

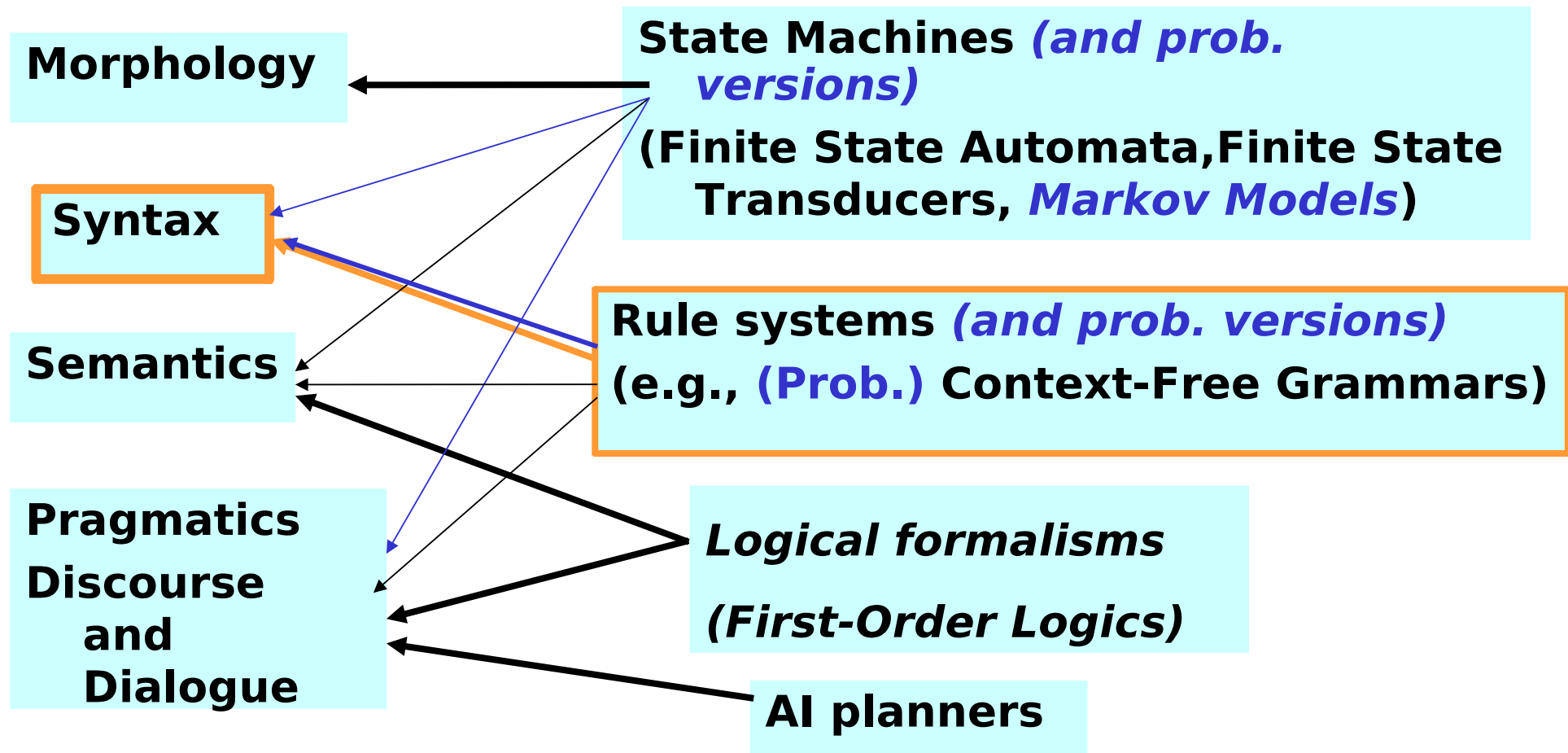
Asst. Prof / CSE

Natural Language Processing

Agenda

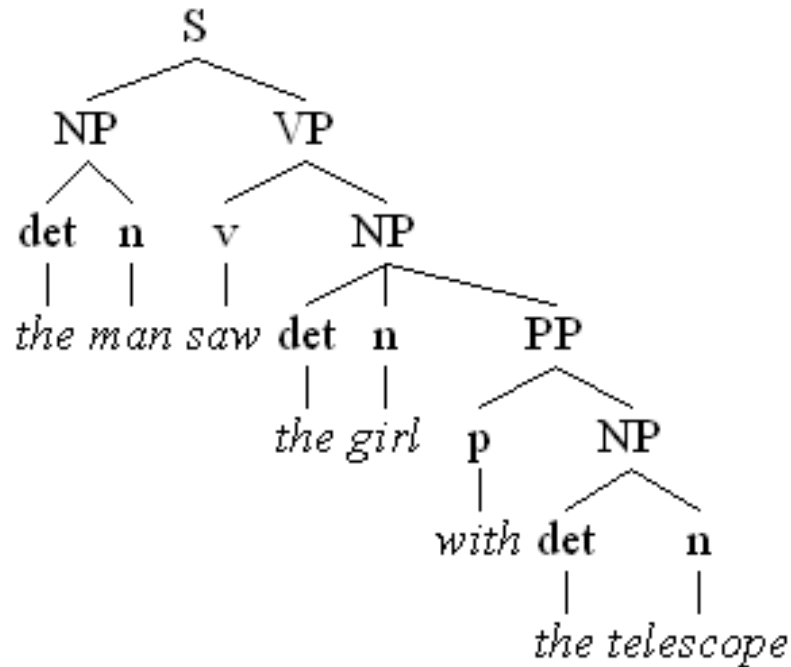
- Probabilistic CFG
- CYK Algorithm
- Lexicalized CFG

Knowledge-Formalism Map

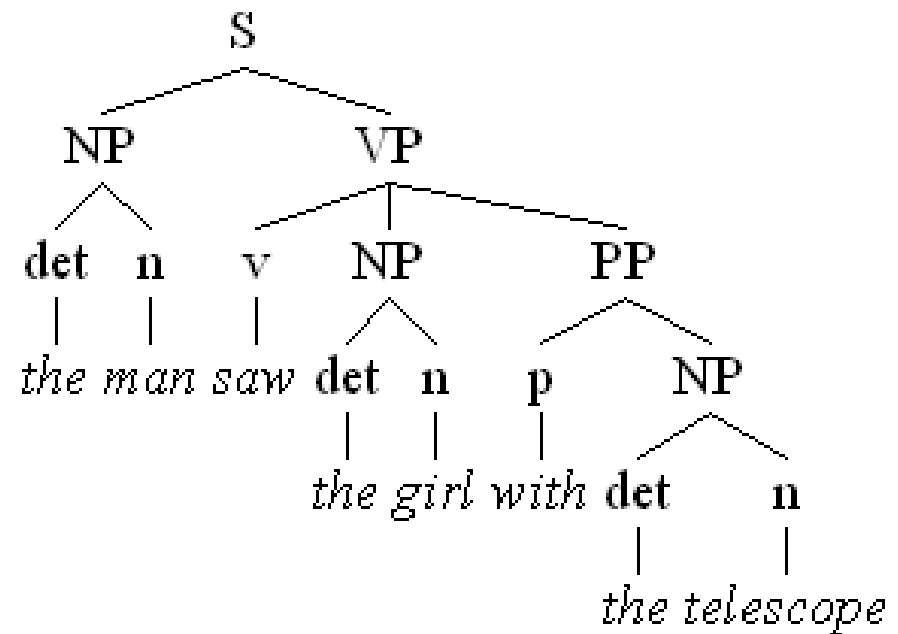


Ambiguity

the man saw the girl with the telescope



The girl has the telescope



The man has the telescope

Probabilistic Context-Free Grammar

- The simplest augmentation of the context-free grammar is the Probabilistic Context-Free Grammar (PCFG).
- Context-free grammar G is defined by four parameters:

N a set of **non-terminal symbols** (or **variables**)
 Σ a set of **terminal symbols** (disjoint from N)
 R a set of **rules** or productions, each of the form $A \rightarrow \beta$,
where A is a non-terminal,
 β is a string of symbols from the infinite set of strings $(\Sigma \cup N)^*$
 S a designated **start symbol**

Probabilistic Context-Free Grammar

- Definition of a CFG:
 - Set of non-terminals (N)
 - Set of terminals (T)
 - Set of rules/productions (P), of the form $A \rightarrow \alpha$
 - Designated start symbol (S)
- Definition of a PCFG:
 - Same as a CFG, but with one more function, D
 - D assigns probabilities to each rule in P

Probabilistic Context-Free Grammar

- A PCFG augments each rule in P with conditional probability:

$$A \rightarrow \beta [p]$$

- A PCFG is a 5-tuple $G=(N, \Sigma, P, S, D)$ where D is a function.
- This function expresses the probability p that the given non-terminal A will be expanded to the sequence β as:

$$\Pr(A \rightarrow \beta) \quad \text{or} \quad \Pr(A \rightarrow \beta | A)$$

- If we consider all the possible expansions of a non-terminal, the sum of their probabilities must be 1.

Probabilistic Context-Free Grammar

- Attach probabilities to each grammar rule:
- $VP \rightarrow \text{Verb}$.55
- $VP \rightarrow \text{Verb NP}$.40
- $VP \rightarrow \text{Verb NP NP}$.05

A PCFG

$S \rightarrow NP VP$	[.80]	$Det \rightarrow that [.10] \mid a [.30] \mid the [.60]$
$S \rightarrow Aux NP VP$	[.15]	$Noun \rightarrow book [.10] \mid flight [.30]$
$S \rightarrow VP$	[.05]	$\mid meal [.15] \mid money [.05]$
$NP \rightarrow Pronoun$	[.35]	$\mid flights [.40] \mid dinner [.10]$
$NP \rightarrow Proper-Noun$	[.30]	$Verb \rightarrow book [.30] \mid include [.30]$
$NP \rightarrow Det Nominal$	[.20]	$\mid prefer; [.40]$
$NP \rightarrow Nominal$	[.15]	$Pronoun \rightarrow I [.40] \mid she [.05]$
$Nominal \rightarrow Noun$	[.75]	$\mid me [.15] \mid you [.40]$
$Nominal \rightarrow Nominal Noun$	[.20]	$Proper-Noun \rightarrow Houston [.60]$
$Nominal \rightarrow Nominal PP$	[.05]	$\mid TWA [.40]$
$VP \rightarrow Verb$	[.35]	$Aux \rightarrow does [.60] \mid can [.40]$
$VP \rightarrow Verb NP$	[.20]	$Preposition \rightarrow from [.30] \mid to [.30]$
$VP \rightarrow Verb NP PP$	[.10]	$\mid on [.20] \mid near [.15]$
$VP \rightarrow Verb PP$	[.15]	$\mid through [.05]$
$VP \rightarrow Verb NP NP$	[.05]	
$VP \rightarrow VP PP$	[.15]	
$PP \rightarrow Preposition NP$	[1.0]	

A probabilistic augmentation of the miniature English grammar and lexicon

Using Probabilities

- A PCFG assigns a probability to each parse-tree T of a sentence S .
- This attribute is useful in disambiguation.
- The probability of a parse T is the product of the probabilities of all rules r used to expand each node n in the parse tree:

$$P(T,S) = \prod_{n \in T} p(r(n))$$

- Probability $P(T,S)$ is joint probability of the parse and the sentence

$$P(T,S) = P(T) P(S|T)$$

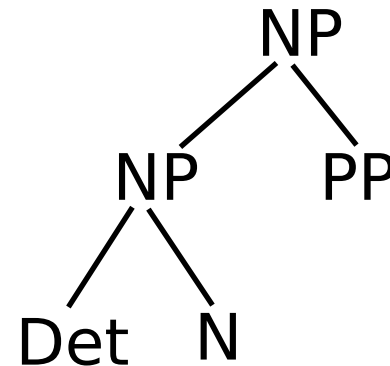
since a parse tree includes all the words of the sentence, $P(S|T) = 1$

Thus: $P(T,S) = P(T)$

$$P(T) = \prod_{n \in T} p(r(n))$$

Using Probabilities : A Sample

NP → Det N	:	0.4
NP → NP _{poss} N	:	0.1
NP → Pronoun	:	0.2
NP → NP PP	:	0.1
NP → N	:	0.2

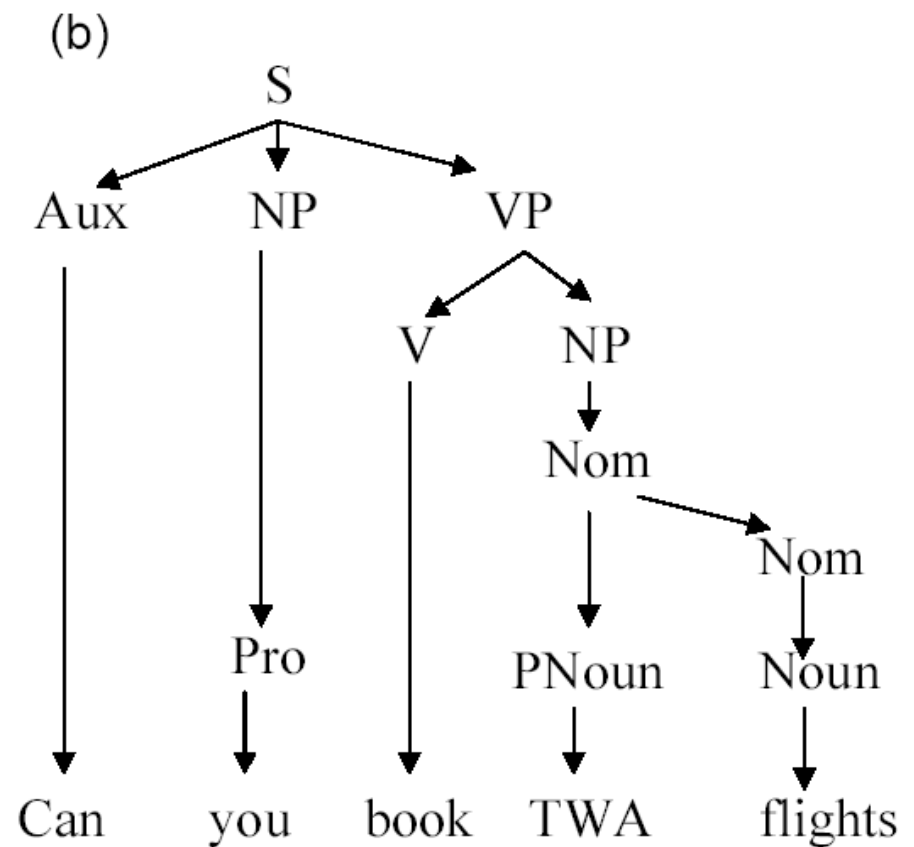
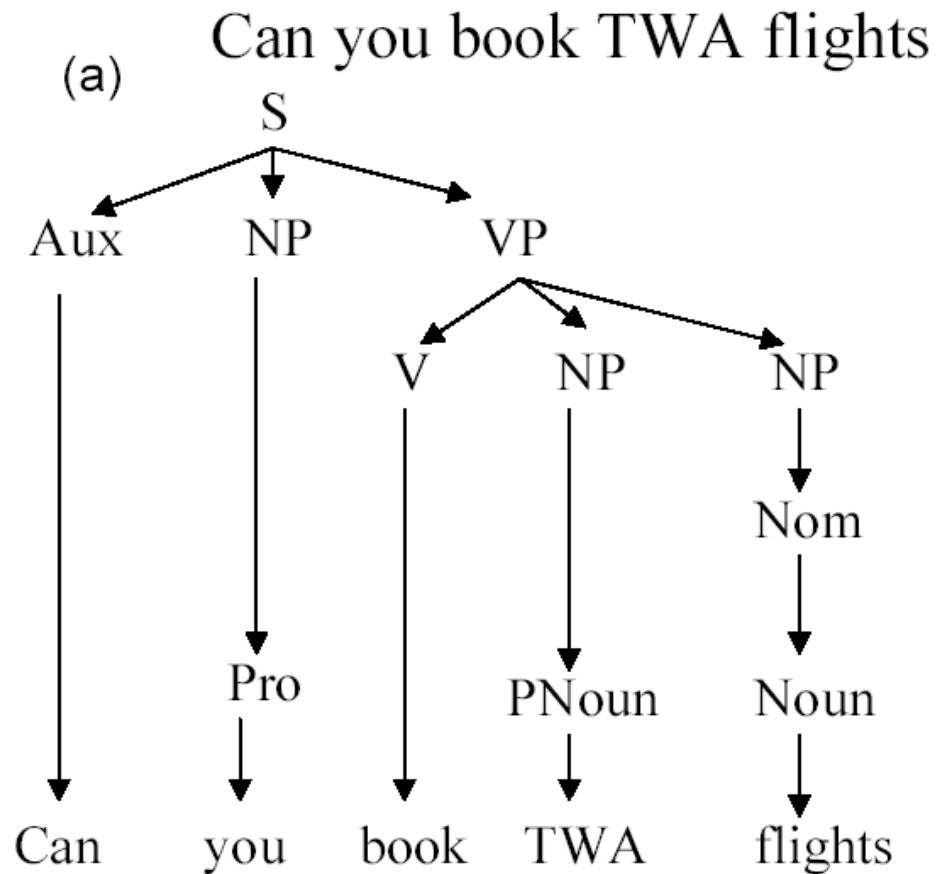


$$P(\text{subtree above}) = 0.1 \times 0.4 = 0.04$$

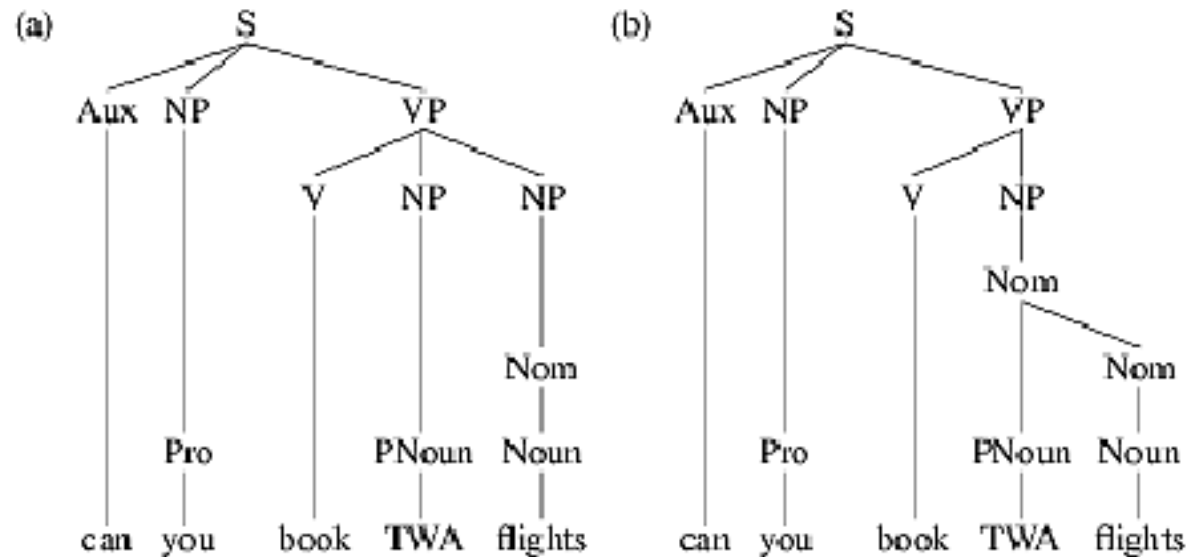
English Practice

- What do you understand from the sentence:
“Can you book TWA flights?”
- Can you book flights on behalf of TWA?
 - [TWA] [flights]
- Can you book flights run by TWA?
 - [TWA flights]

A Sample Parse



A Sample Parse with PCFG

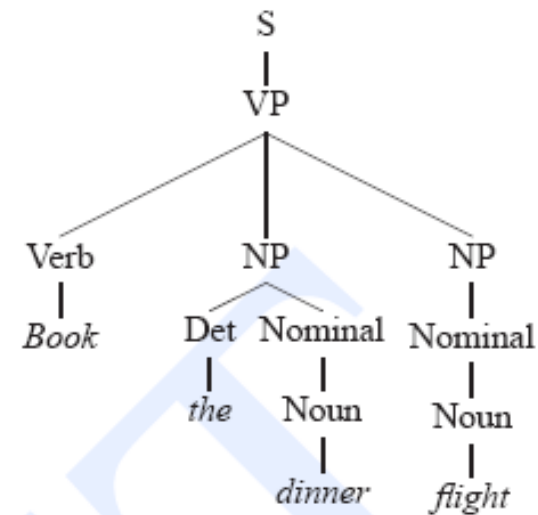
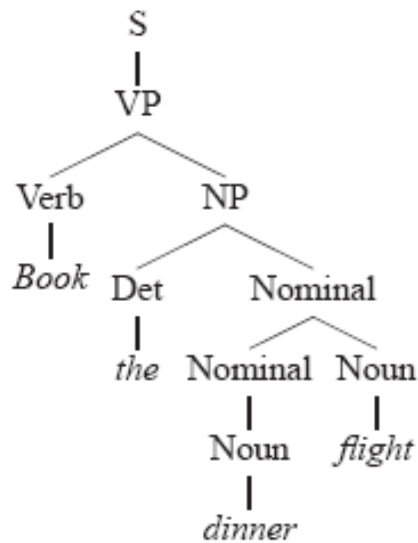


Rules	P	Rules	P
S → Aux NP VP	.15	S → Aux NP VP	.15
NP → Pro	.40	NP → Pro	.40
VP → V NP NP	.05	VP → V NP	.40
NP → Nom	.05	NP → Nom	.05
NP → PNoun	.35	Nom → PNoun Nom	.05
Nom → Noun	.75	Nom → Noun	.75
Aux → Can	.40	Aux → Can	.40
NP → Pro	.40	NP → Pro	.40
Pro → you	.40	Pro → you	.40
Verb → book	.30	Verb → book	.30
PNoun → TWA	.40	PNoun → TWA	.40
Noun → flights	.50	Noun → flights	.50

A Sample Parse with PCFG

- $P(T_1) = .15 \times .40 \times .05 \times .05 \times .35 \times .75 \times .40 \times .40 \times .40 \times .30 \times .40 \times .50$
 $= 1.5 \times 10^{-6}$
- $P(T_r) = .15 \times .40 \times .40 \times .05 \times .05 \times .75 \times .40 \times .40 \times .40 \times .30 \times .40 \times .50$
 $= 1.7 \times 10^{-6}$
- The right tree T_r has a higher probability and would be chosen by a disambiguation algorithm.
- $P(\text{Sentence}) = 1.5 \times 10^{-6} + 1.7 \times 10^{-6}$
 $= 3.2 \times 10^{-6}$

A Sam



Rules			P	Rules			P
S	→	VP	.05	S	→	VP	.05
VP	→	Verb NP	.20	VP	→	Verb NP NP	.10
NP	→	Det Nominal	.20	NP	→	Det Nominal	.20
Nominal	→	Nominal Noun	.20	NP	→	Nominal	.15
Nominal	→	Noun	.75	Nominal	→	Noun	.75
Verb	→	book	.30	Nominal	→	Noun	.75
Det	→	the	.60	Verb	→	book	.30
Noun	→	dinner	.10	Det	→	the	.60
Noun	→	flights	.40	Noun	→	dinner	.10
				Noun	→	flights	.40

Figure 14.2 Two parse trees for an ambiguous sentence, The transitive parse (a) corresponds to the sensible meaning “Book flights that serve dinner”, while the ditransitive parse (b) to the nonsensical meaning “Book flights on behalf of ‘the dinner’”.

$$P(T_{left}) = .05 * .20 * .20 * .20 * .75 * .30 * .60 * .10 * .40 = 2.2 \times 10^{-6}$$

$$P(T_{right}) = .05 * .10 * .20 * .15 * .75 * .75 * .30 * .60 * .10 * .40 = 6.1 \times 10^{-7}$$

Picking the best parse

- Picking the parse with the highest probability is the correct way to do disambiguation.
- Pick the best tree for a sentence S out of the set of parse trees for S .

$$\hat{T}(S) = \operatorname{argmax}_{T \in \tau(S)} P(T|S)$$

$$= \operatorname{argmax} \frac{P(T, S)}{P(S)}$$

$$= \operatorname{argmax} P(T, S)$$

$$= \operatorname{argmax} P(T)$$

Parse tree T which is most likely given the sentence S .

$P(T|S)$ can be rewritten as $P(T, S)/P(S)$.

Since we are maximizing over all parse for the same sentence $P(S)$ will be a constant.

Since $P(T, S) = P(T)$.

The most likely parse is choosing the **parse with the highest probability**

Getting Probabilities

- From an annotated database (a treebank)
- Learned from a corpus

Getting Probabilities – from treebank

- Get a large collection of parsed sentences
- Compute probability for each non-terminal rule expansion in the collection
- Normalize
- Done

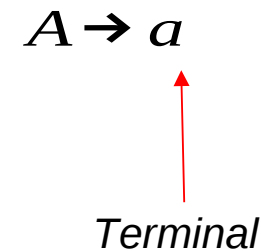
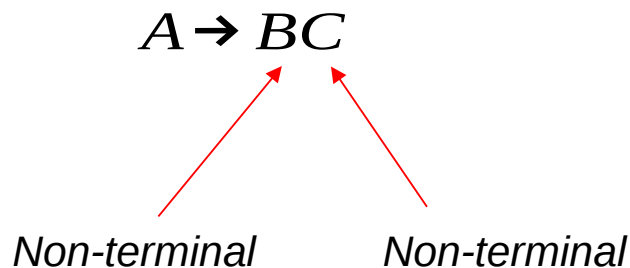
Getting Probabilities – learn from corpus

- What if you don't have a treebank (and can't get one)
- Take a large collection of text and parse it.
- In the case of syntactically ambiguous sentences collect all the possible parses
- Prorate the rule statistics gathered for rules in the ambiguous case by their probability
- Proceed as you did with a treebank.
- **Inside-Outside** algorithm

Parsing of PCFGs

Probabilistic CYK

- Probabilistic CYK (Cocke-Younger-Kasami) algorithm for parsing Probabilistic CFG.
- Bottom-up dynamic parsing algorithm.
- Assume PCFG is in Chomsky Normal Form (CNF):
production is either of the form $A \rightarrow B C$ or $A \rightarrow a$



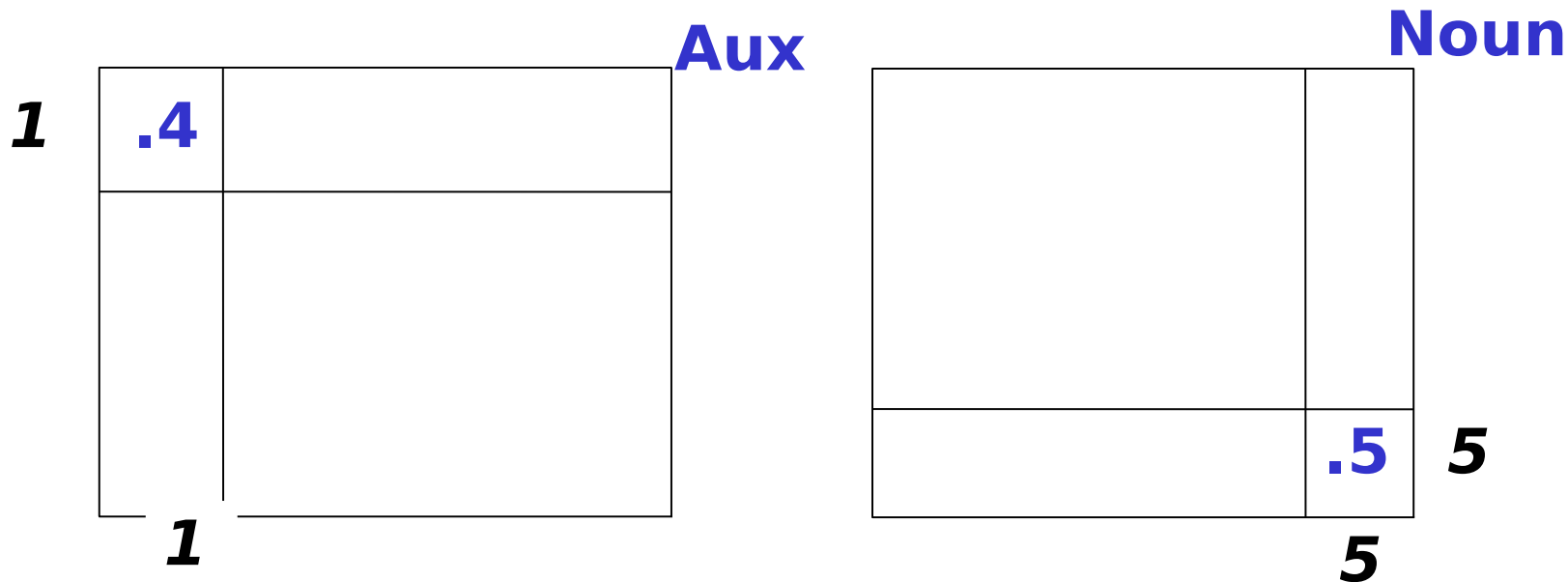
Probabilistic CYK

- CYK Algorithm: bottom-up parser
- Input:
 - A **Chomsky normal form** PCFG, $G = (N, \Sigma, P, S, D)$. Assume that the $|N|$ non-terminals have indices $1, 2, \dots, |N|$, and the start symbol S has index 1.
 - n words w_1, \dots, w_n
- Data Structure:
 - A dynamic programming **array** $\Pi\pi[i, j, a]$ holds the maximum probability for a constituent with non-terminal index a spanning words $i..j$.
- Output:
 - The **maximum** probability parse $\Pi\pi[1, n, 1]$: the parse tree whose root is S and which spans the entire string of words w_1, \dots, w_n

CYK Algorithm : Base Case

- Consider the input *strings of length one* (individual words w_i)
- In CNF, the probability of a given non-terminal A expanding to a single word w_i must come only from the rule $A \rightarrow w_i$ i.e., $P(A \rightarrow w_i)$
- $\Pi\pi[i,j,a] = P(A \rightarrow w_i)$

“Can₁ you₂ book₃ TWA₄ flights₅ ?”



CYK Algorithm : Recursive Case

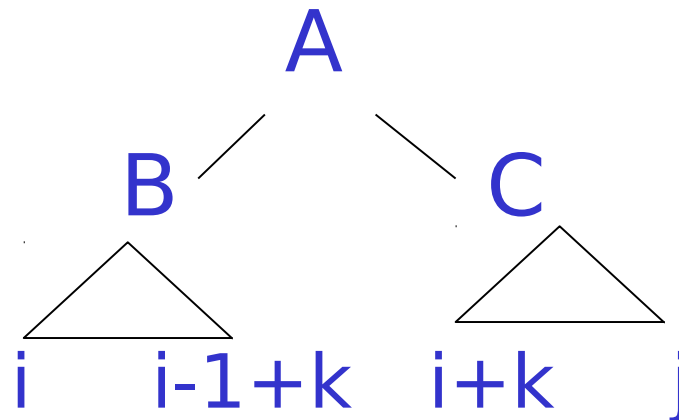
- For strings of words of length > 1 ,

$A \Rightarrow w_{ij}$ iff there is at least one rule $A \rightarrow BC$

where B derives the first k words (between i and $i-1+k$) and

C derives the remaining ones (between $i+k$ and j)

- (for each non-terminal) Choose the max among all possibilities
- $\Pi\pi[i,j,A] = \Pi\pi[i, i-1+k, B] \times \Pi\pi[i+k, j, C] \times P(A \rightarrow BC)$



CYK Algorithm :

Function CYK(*words*, *grammar*)

return the most probable parse and its probability

For $i = \leftarrow 1$ **to** *num_words* #base case

for $a = \leftarrow 1$ **to** *num_nonterminals*

If ($A \rightarrow w_i$) is in grammar **then** $\pi\Pi[i, i, a] = \leftarrow P(A \rightarrow w_i)$

For $j = \leftarrow 2$ **to** *num_words* #recursive

For $i = \leftarrow 1$ **to** *num_words* - *j* + 1

For $k = \leftarrow 1$ **to** *j* - 1

for $A = \leftarrow 1$ **to** *num_nonterminals*

for $B = \leftarrow 1$ **to** *num_nonterminals*

for $C = \leftarrow 1$ **to** *num_nonterminals*

$prob = \leftarrow \Pi\Pi[i, k, B] \times \Pi\Pi[i+k, j-k, C] \times P(A \rightarrow BC)$

If ($prob > \pi\Pi[i, j, A]$) **then**

$\pi\Pi[i, j, A] = prob$

$back[i, j, A] = \{k, A, B\}$

Return *build_tree(back[1, num_words, 1]), $\pi[1, num_words, 1]$*

CYK Algorithm :

Det: .40 [0,1]	NP: $.30 * .40 * .02$ = .0024 [0,2]	[0,3]	[0,4]	[0,5]
	N: .02 [1,2]	[1,3]	[1,4]	[1,5]
		V: .05 [2,3]	[2,4]	[3,5]
			[3,4]	[3,5]
				[4,5]

$S \rightarrow NP VP$.80	$Det \rightarrow the$.50
$NP \rightarrow Det N$.30	$Det \rightarrow a$.40
$VP \rightarrow V NP$.20	$N \rightarrow meal$.01
$V \rightarrow includes$.05	$N \rightarrow flight$.02

The flight includes a meal

Problems with PCFG

- Poor independence assumptions
 - CFG impose an independence assumption on probabilities
 - Results in poor modeling of structural dependencies
- Lack of lexical conditioning
 - CFG rules don't model syntactic facts about words leading to problems with:
 - subcategorization ambiguities
 - preposition attachment
 - coordinate structure ambiguities

Problems with PCFG – 1

- The expansion of any one non-terminal is independent of the expansion of any other non-terminal – assumption
- Statistics of English syntax shows that the choice of how a node expands is dependent on the location of the node in the parse tree.
- NP → Det NN .28
NP → ProNoun .25
- PCFG don't allow a rule probability to be conditioned on surrounding context – hence equal probability

Problems with PCFG – 1

- The probability of expanding an *NP* as a pronoun versus a lexical *NP* were *conditioned* on whether the *NP* was a subject or an object.
- $P[(NP \rightarrow \text{Pronoun}) \mid NP = \text{subj}] \gg P[(NP \rightarrow \text{Pronoun}) \mid NP = \text{obj}]$

	Pronoun	Non-Pronoun
Subject	91%	9%
Object	34%	66%

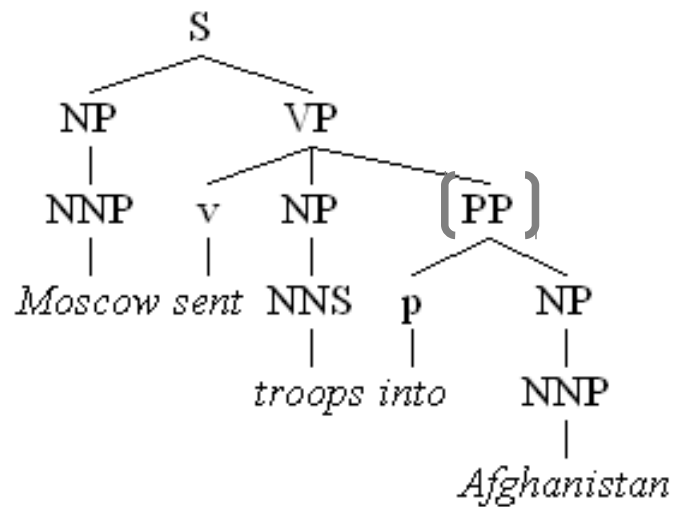
Problems with PCFG – 2

- Lexical information plays an important role in selecting the correct parsing of an ambiguous *PP* attachment
- Moscow *sent* more than 100,000 soldiers *into Afghanistan*.
- Here the *PP* [*into Afghanistan*] can be attached either to the *NP* [*more than 100,000 soldiers*] or to the *VP* headed by *sent*
- In PCFG, the choice between two rules:

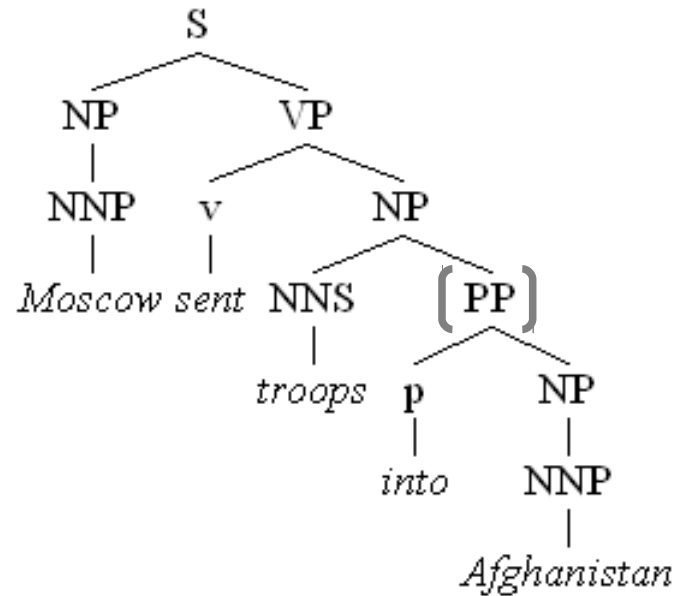
	$NP \rightarrow NP\ PP$	$VP \rightarrow NP\ PP$
<i>Hindle and Rooth (1991)</i>	67%	33%
<i>Collins (1999)</i>	52%	48%

- Verb *send* subcategorizes for destination, which can be expressed with *into*
- Keep separate **lexical dependency** statistics for different verbs

Problems with PCFG – 2



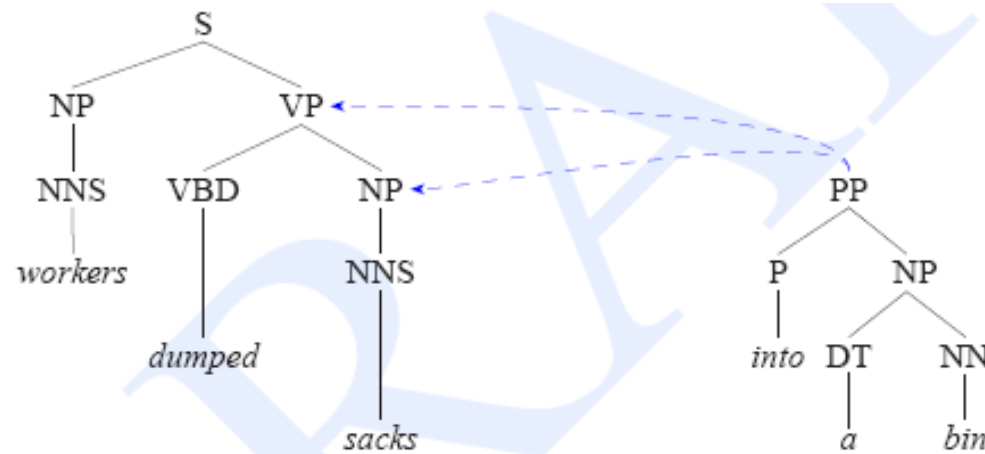
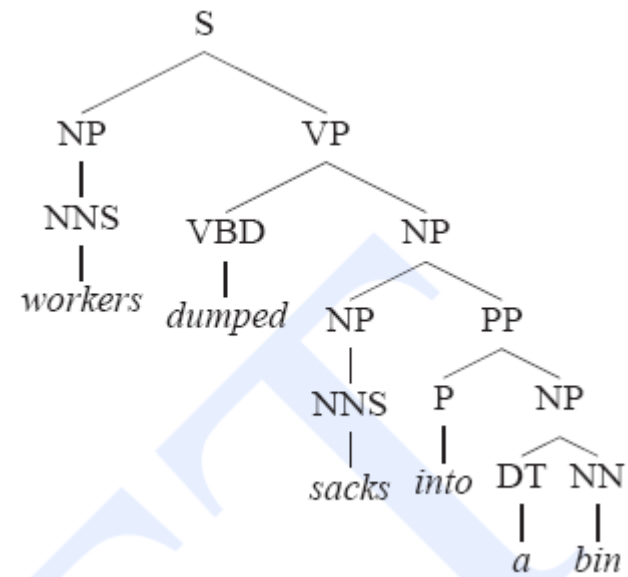
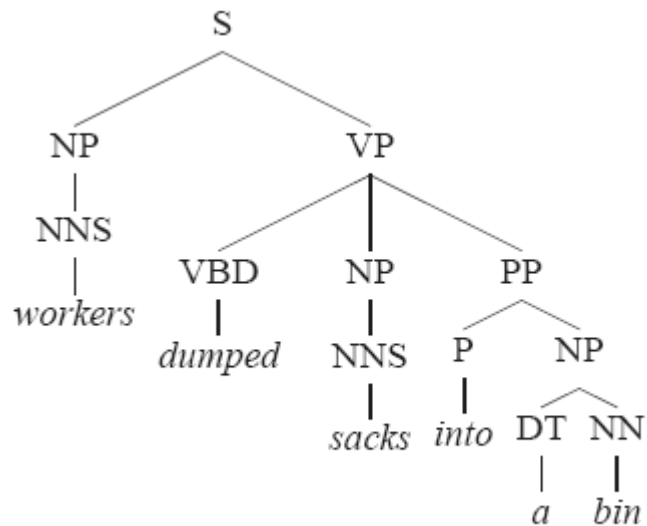
VP-attachment



NP-attachment

Typically NP-attachment more frequent than VP-attachment

Problems with PCFG – 2

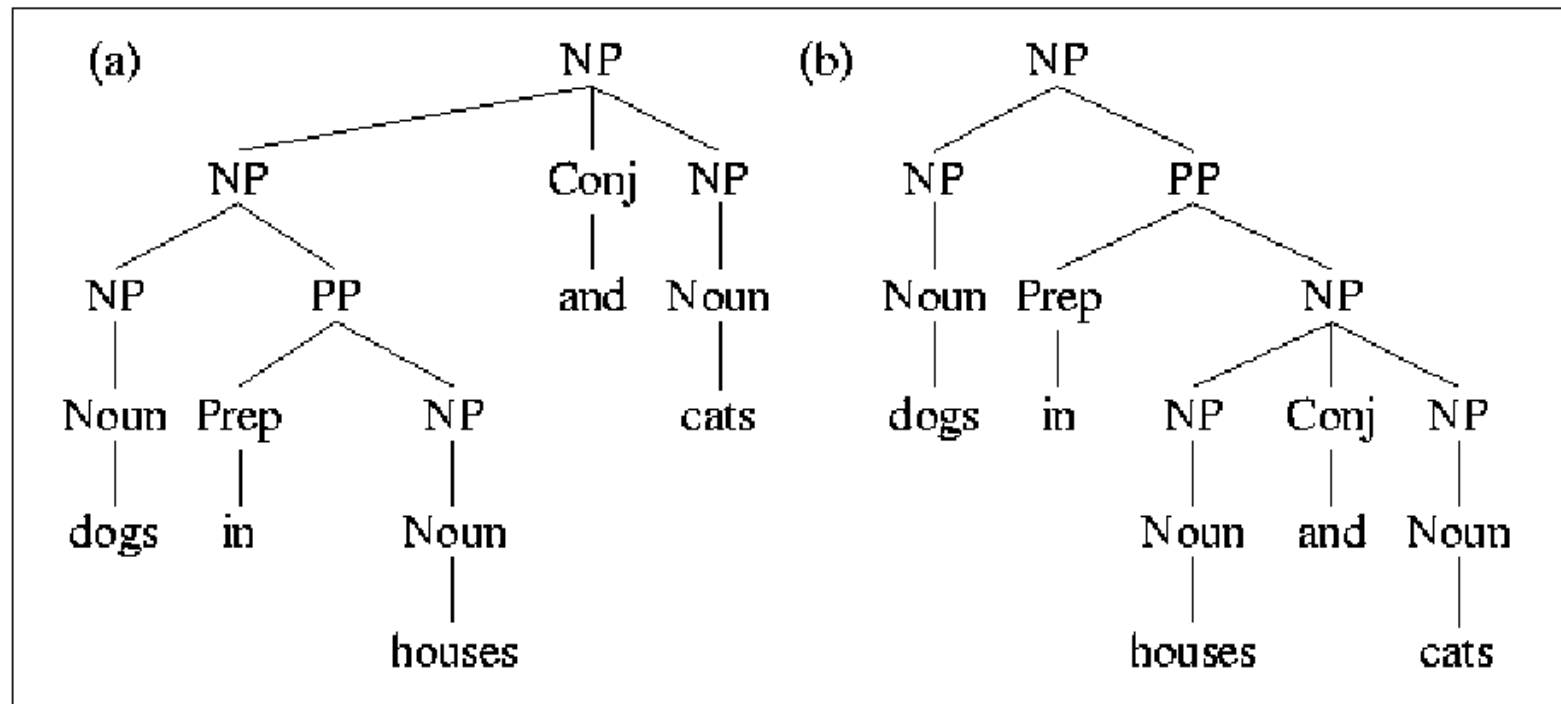


$VP \rightarrow VBD \ NP \ PP$

$NP \rightarrow NP \ PP$

Problems with PCFG – 3

- Coordination ambiguity:
- [[*dogs in houses*] *and* [*cats*]]
dogs in [[*houses*] *and* [*cats*]]



- A PCFG assigns identical probabilities, since both structure use the exact same rules

Probabilistic Lexicalized CFG

Probabilistic Lexicalized CFG

- Syntactic constituents could be associated with a lexical head.
- Each non-terminal in a parse tree is annotated with its lexical head.
- Each PCFG rule must be augmented to identify one RHS constituent to be the head daughter.
- The headword for a node is set to the headword of its head daughter.

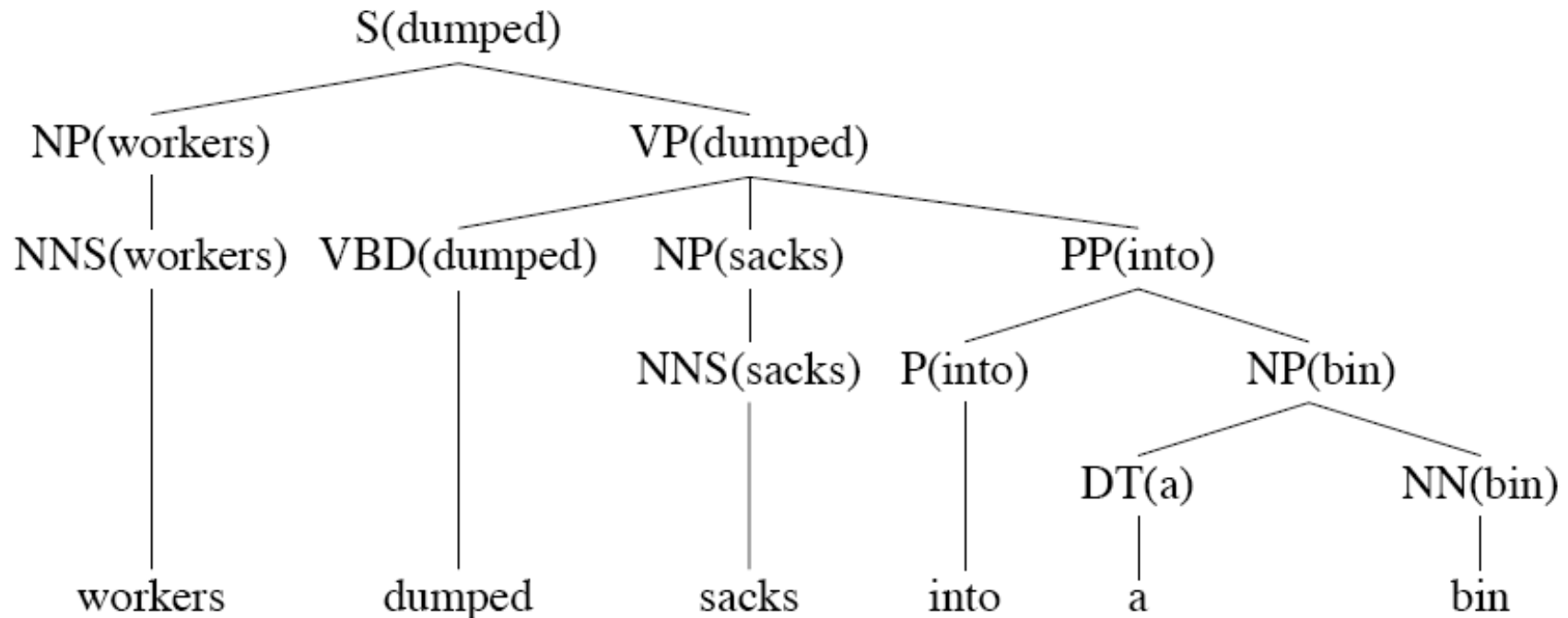
VP -> VBD NP PP

VP(dumped) -> VBD(dumped) NP(sacks) PP(into)

- In a standard lexicalized grammar, associate the head tag with the POS tags of the headwords:

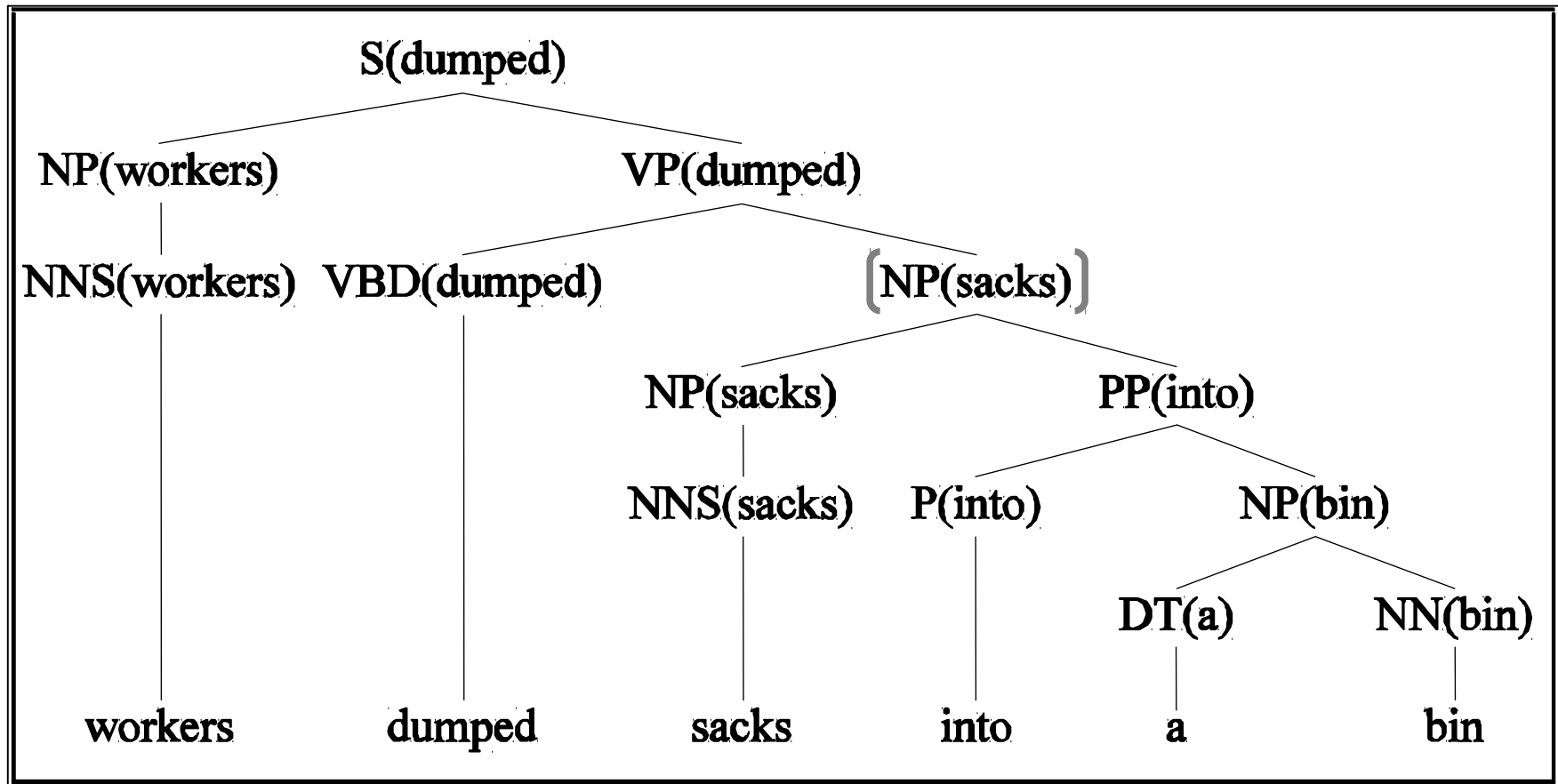
VP(dumped,VBD) -> VBD(dumped,VBD) NP(sacks,NNS) PP(into,IN)

Probabilistic Lexicalized CFG



A lexicalized tree from Collins (1999)

Probabilistic Lexicalized CFG – Incorrect parse



An incorrect parse of the sentence from Collins (1999)

Probabilistic Lexicalized CFG

- Add the headword of the node $h(n)$ for computing the probability of a node being expanded via rule r .

$$p(r(n)|n, h(n)) \quad [p(r(n)) - \text{syntactic category of } n] - (1)$$

$$r = VP \rightarrow VBD NP PP$$

$$p(r|VP, \textit{dumped})$$

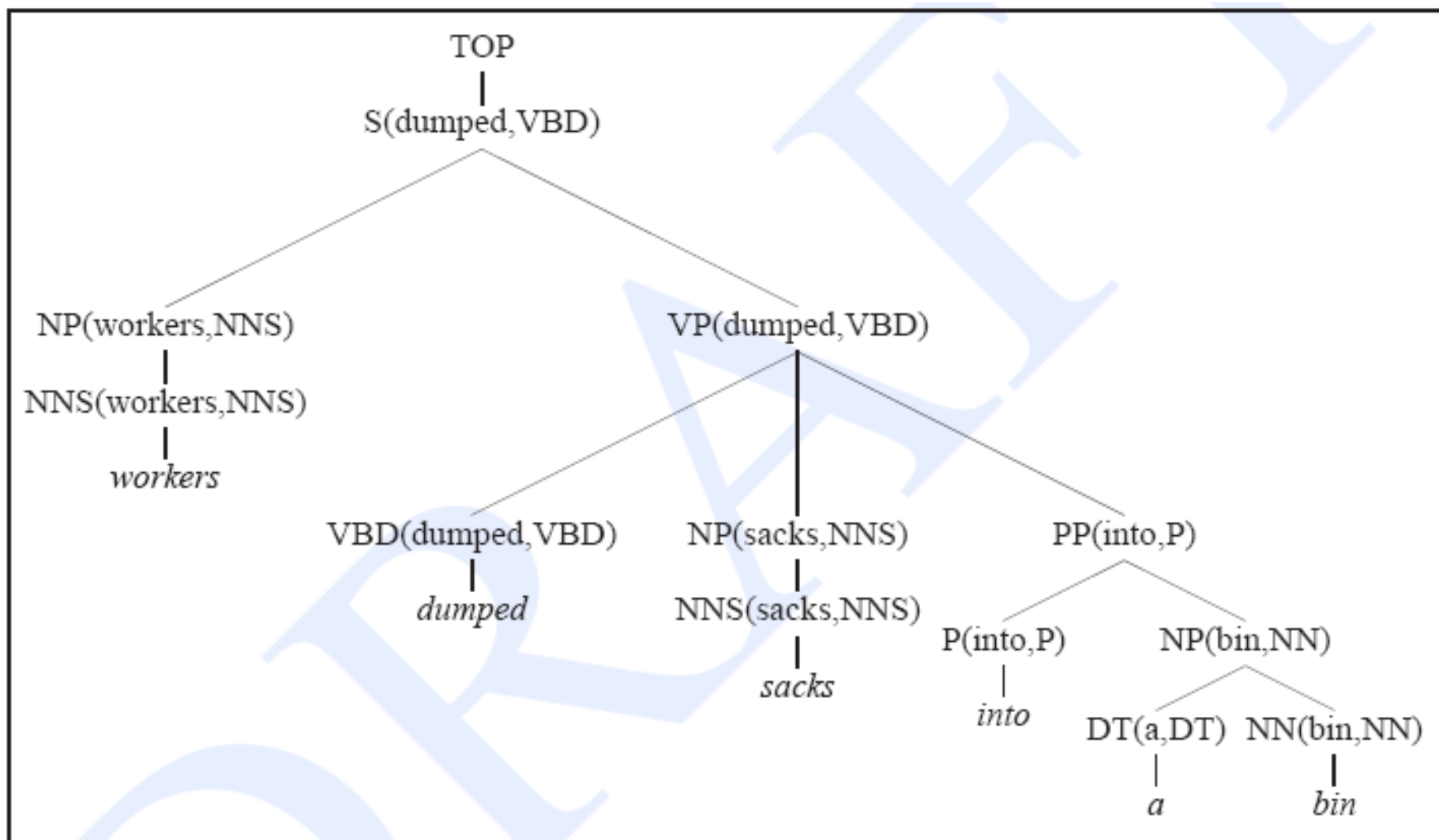
- What is the probability that a VP headed by *dumped* will be expanded as $VBD NP PP$?
- How to compute the probability of a head?
- The probability of a node n having a head h on two factors:
 - The syntactic category of the node n
 - The head of the node's mother $h(m(n))$

Probabilistic Lexicalized CFG

- $P(h(n) = word_i \mid n, h(m(n)))$ – (2)
- $p(head(n)=sacks \mid n=NP, h(m(n))=dumped)$
- What is the probability that an *NP* whose mother's head is *dumped* has the head *sacks*?

$X(dumped)$
|
 $NP(?sacks?)$

- Head-probability is capturing *dependency information* between the words *dumped* and *sacks*



Internal Rules

TOP	→	S(dumped, VBD)	
S(dumped, VBD)	→	NP(workers, NNS)	VP(dumped, VBD)
NP(workers, NNS)	→	NNS(workers, NNS)	
VP(dumped, VBD)	→	VBD(dumped, VBD)	NP(sacks, NNS) PP(into, P)
PP(into, P)	→	P(into, P)	NP(bin, NN)
NP(bin, NN)	→	DT(a, DT)	NN(bin, NN)

Lexical Rules

NNS(workers, NNS)	→	workers
VBD(dumped, VBD)	→	dumped
NNS(sacks, NNS)	→	sacks
P(into, P)	→	into
DT(a, DT)	→	a
NN(bin, NN)	→	bin

Probabilistic Lexicalized CFG

- Probability of a parse is given by:

$$P(T,S) = \prod_{n \in T} p(r(n)|n, h(n)) \times p(h(n)|n, h(m(n)))$$

- The head-rule and head-head probabilities will correctly choose the *VP* attachment over the *NP* attachment.
- Head-rule probabilities calculated as:

$$P(VP(dumped, VBD) \rightarrow VBD(dumped, VBD)NP(sacks, NNS)PP(into, P)) \\ = \frac{Count(VP(dumped, VBD) \rightarrow VBD(dumped, VBD)NP(sacks, NNS)PP(into, P))}{Count(VP(dumped, VBD))}$$

how many times the rule occurs with $h(n)$ as the headword versus how many times the mother/headword combination appear in total.

Probabilistic Lexicalized CFG

- From the Brown corpus: head-rule probabilities

$$p(\text{VP} \rightarrow \text{VBD NP PP} \mid \text{VP}, \textit{dumped}) = 6 / 9 = .67$$

$$p(\text{VP} \rightarrow \text{VBD NP} \mid \text{VP}, \textit{dumped}) = 0 / 9 = 0$$

- What about the head-head probabilities?

$$p(\textit{into} \mid \text{PP}, \textit{dumped}) = 2 / 9 = .22$$

$$p(\textit{into} \mid \text{PP}, \textit{sacks}) = 0 / 0 = 0$$

- So, the contribution of this part of the parse to the total scores for the two candidates is:

$$[\textit{dumped into}] .67 \times .22 = .147$$

$$[\textit{sacks into}] \quad 0 \times 0 = 0$$

References

- Slides were adapted from:
Speech and Language Processing, *Jurfsky and Martin*