# Intelligent Systems (AI-2)

## Computer Science cpsc422, Lecture 26

# Learning Goals for today's class

## You can:

- Provide a formal definition of a PCFG

- Apply a PCFG to compute the probability of a parse tree of a sentence as well as the probability of a sentence

- Describe the content of a treebank

- Describe the process to identify a head of a syntactic constituent

- Compute the probability distribution of a PCFG from a treebank

# Lecture Overview

- Recap English Syntax and Parsing
- Key Problem with parsing: Ambiguity
- Probabilistic Context Free Grammars (PCFG)
- Treebanks and Grammar Learning

# Key Constituents: Examples

*Head*

**NP → N**
**NP → Det N**

**(Specifier) X (Complement)**

- **Noun phrases (NP)**
  - (Det)    N    (PP)
    - the    cat   on the table

- **Verb phrases (VP)**
  - (Qual)   V    (NP)
    - never   eat    a cat

- **Prepositional phrases (PP)**
  - (Deg)    P     (NP)
    - almost   in    the net

- **Adjective phrases(AP)**
  - (Deg)    A     (PP)
    - very    happy   about it

- **Sentences (S)**
  - (NP)      (-)  (VP)
    - a mouse    --    ate it

CPSC 422, Lecture 26

# Context Free Grammar (CFG)

- **4-tuple** (non-term., term., productions, start)

- **(** $N, \Sigma, P, S$ **)**

- P is a set of rules $A \rightarrow \alpha$; $A \in \mathbf{N}$, $\alpha \in (\Sigma \cup N)^*$

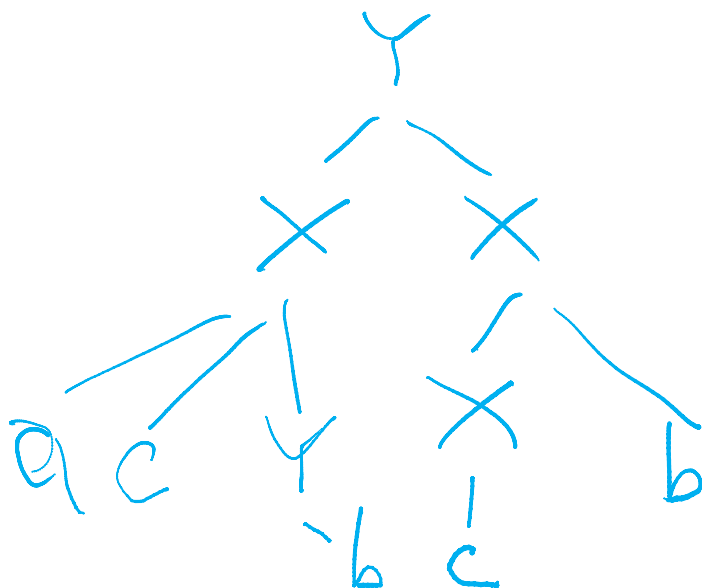$N = \{X \, Y\}$  $\Sigma = \{a \, b \, c\}$  $P = \quad X \rightarrow X b$

$Y \rightarrow X X$

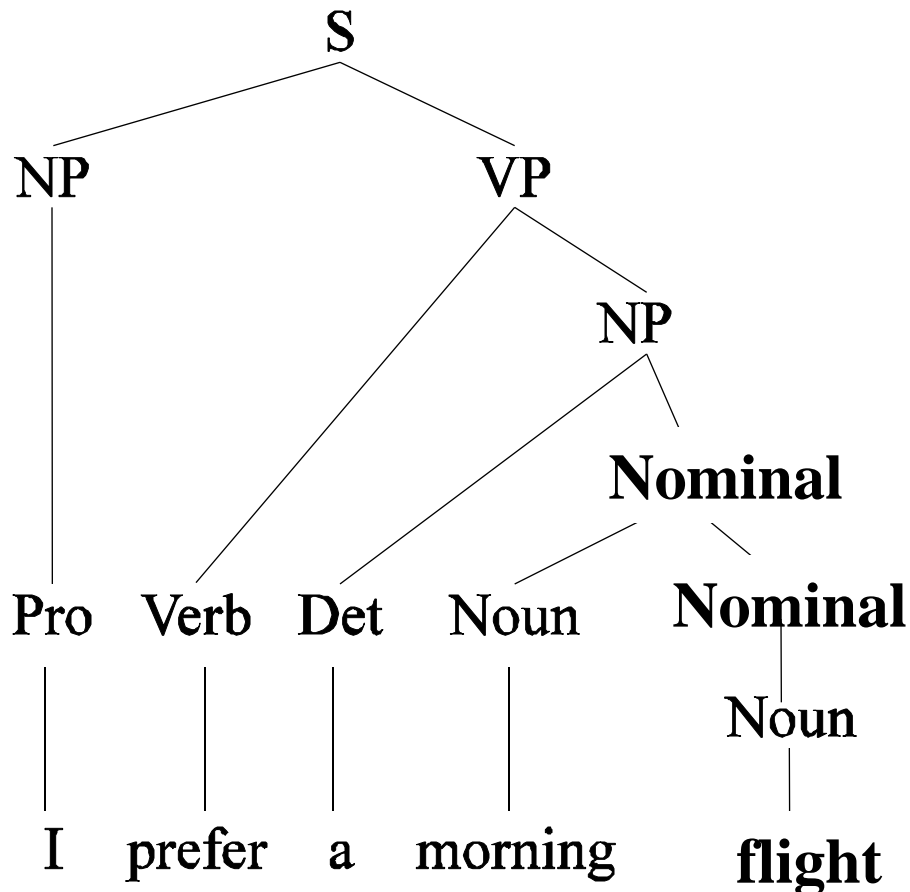$X \rightarrow ac \, Y$

$X \rightarrow c$

$Y \rightarrow b$

# CFG Example

## Grammar with example phrases

| | | |
|---|---|---|
| $S$ → | $NP\ VP$ | I + want a morning flight |
| $NP$ → | $Pronoun$ | I |
| | $Proper$-$Noun$ | Los Angeles |
| | $Det\ Nominal$ | a + flight |
| $Nominal$ → | $Noun\ Nominal$ | morning + flight |
| | $Noun$ | flights |
| $VP$ → | $Verb$ | do |
| | $Verb\ NP$ | want + a flight |
| | $Verb\ NP\ PP$ | leave + Boston + in the morning |
| | $Verb\ PP$ | leaving + on Thursday |
| $PP$ → | $Preposition\ NP$ | from + Los Angeles |

## Lexicon

| | |
|---|---|
| $Noun$ → | $flights \mid breeze \mid trip \mid morning \mid \ldots$ |
| $Verb$ → | $is \mid prefer \mid like \mid need \mid want \mid fly$ |
| $Adjective$ → | $cheapest \mid non-stop \mid first \mid latest$ $\mid other \mid direct \mid \ldots$ |
| $Pronoun$ → | $me \mid I \mid you \mid it \mid \ldots$ |
| $Proper$-$Noun$ → | $Alaska \mid Baltimore \mid Los\ Angeles$ $\mid Chicago \mid United \mid American \mid \ldots$ |
| $Determiner$ → | $the \mid a \mid an \mid this \mid these \mid that \mid \ldots$ |
| $Preposition$ → | $from \mid to \mid on \mid near \mid \ldots$ |
| $Conjunction$ → | $and \mid or \mid but \mid \ldots$ |

# Derivations as Trees



$$S \rightarrow NP\ VP$$

$$NP \rightarrow Pronoun$$
$$\mid\ Proper\text{-}Noun$$
$$\mid\ Det\ Nominal$$
$$Nominal \rightarrow Noun\ Nominal$$
$$\mid\ Noun$$

$$VP \rightarrow Verb$$
$$\mid\ Verb\ NP$$
$$\mid\ Verb\ NP\ PP$$
$$\mid\ Verb\ PP$$

$$PP \rightarrow Preposition\ NP$$

# Example of relatively complex parse tree



NP→JJ JJ NN
NP → NNS
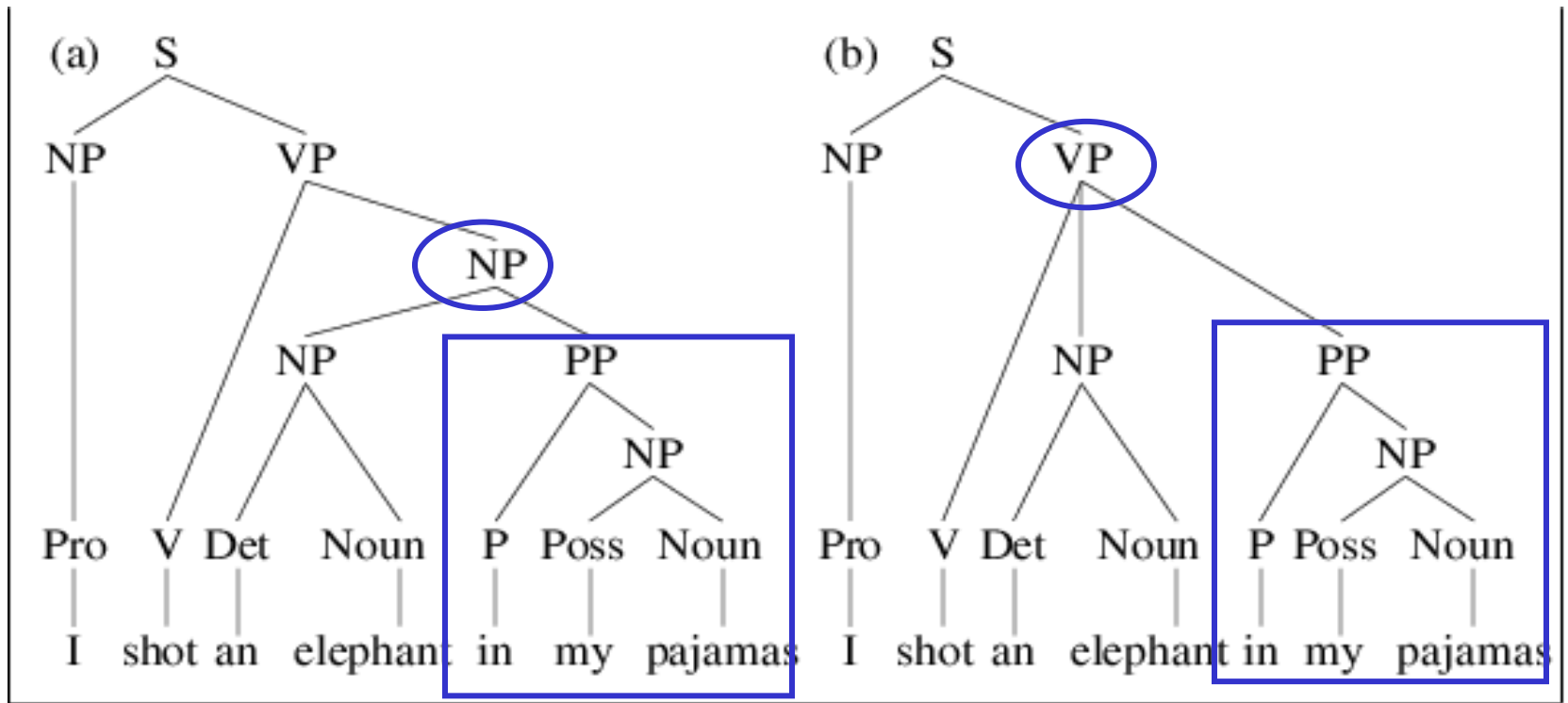NP → VBN NN

# Lecture Overview

- Recap English Syntax and Parsing
- **Key Problem with parsing: Ambiguity**
- Probabilistic Context Free Grammars (PCFG)
- Treebanks and Grammar Learning

# Structural Ambiguity (Ex. 1)

VP -> V NP ; NP -> NP PP
VP -> V NP PP

*"I shot an elephant in my pajamas"*

# Structural Ambiguity (Ex.2)

"I saw Mary passing by cs2"

"I saw Mary passing by cs2"

```
(ROOT                              (ROOT
  (S                                 (S
    (NP (PRP I))                       (NP (PRP I))
    (VP (VBD saw)                      (VP (VBD saw)
      (S                                 (NP (NNP Mary))
        (NP (NNP Mary))                  (S
        (VP (VBG passing)                  (VP (VBG passing)
          (PP (IN by)                        (PP (IN by)
            (NP (NNP cs2))))))                  (NP (NNP cs2)))))))
```

# Structural Ambiguity (Ex. 3)

- **Coordination** JJ N N **"new students and profs"**

$$NP \rightarrow NP \text{ and } NP$$

$$NP \rightarrow JJ \; NP$$

$$NP \rightarrow N$$

# Structural Ambiguity (Ex. 4)

- **NP-bracketing "French language teacher"**

NP → JJ NP

NP → N

NP → NP NP

# Lecture Overview

- Recap English Syntax and Parsing
- Key Problem with parsing: Ambiguity
- **Probabilistic Context Free Grammars (PCFG)**
- Treebanks and Grammar Learning (acquiring the probabilities)
- Intro to Parsing PCFG

# Probabilistic CFGs (PCFGs)

- GOAL: assign a probability to parse trees and to sentences
- Each grammar rule is augmented with a conditional probability

- If these are <u>all the rules for VP</u> and .55 is P(VP->Verb | VP)

| | |
|---|---|
| VP -> Verb | .55 |
| VP -> Verb NP | .40 |
| VP -> Verb NP NP | **??** |

A. 1
B. 0
C. 0.05
D. 0.42
E. None of the above

i>clicker.

- **What should ?? be ?**

# Probabilistic CFGs (PCFGs)

- GOAL: assign a probability to parse trees and to sentences

- Each grammar rule is augmented with a conditional probability

  - **The expansions for a given non-terminal sum to 1**

    VP -> Verb                    .55

    VP  -> Verb NP                .40

    VP  -> Verb NP NP             .05

$P(VP \rightarrow Verb \mid VP)$

$P(VP \rightarrow Verb\, NP \mid VP)$

$P(VP \rightarrow Verb\, NP\, NP \mid VP)$

**Formal Def**: 5-tuple $(N, \Sigma, P, S, D)$

# Sample PCFG

| | | | |
|---|---|---|---|
| $S \rightarrow NP\ VP$ | [.80] | $Det \rightarrow that\ [.05] \mid the\ [.80] \mid a\ [.15]$ | |
| $S \rightarrow Aux\ NP\ VP$ | [.15] | $Noun \rightarrow book$ | [.10] |
| $S \rightarrow VP$ | [.05] | $Noun \rightarrow flights$ | [.50] |
| $NP \rightarrow Det\ Nom$ | [.20] | $Noun \rightarrow meal$ | [.40] |
| $NP \rightarrow Proper\text{-}Noun$ | [.35] | $Verb \rightarrow book$ | [.30] |
| $NP \rightarrow Nom$ | [.05] | $Verb \rightarrow include$ | [.30] |
| $NP \rightarrow Pronoun$ | [.40] | $Verb \rightarrow want$ | [.40] |
| $Nom \rightarrow Noun$ | [.75] | $Aux \rightarrow can$ | [.40] |
| $Nom \rightarrow Noun\ Nom$ | [.20] | $Aux \rightarrow does$ | [.30] |
| $Nom \rightarrow Proper\text{-}Noun\ Nom$ | [.05] | $Aux \rightarrow do$ | [.30] |
| $VP \rightarrow Verb$ | [.55] | $Proper\text{-}Noun \rightarrow TWA$ | [.40] |
| $VP \rightarrow Verb\ NP$ | [.40] | $Proper\text{-}Noun \rightarrow Denver$ | [.40] |
| $VP \rightarrow Verb\ NP\ NP$ | [.05] | $Pronoun \rightarrow you\ [.40] \mid I\ [.60]$ | |

# PCFGs are used to....
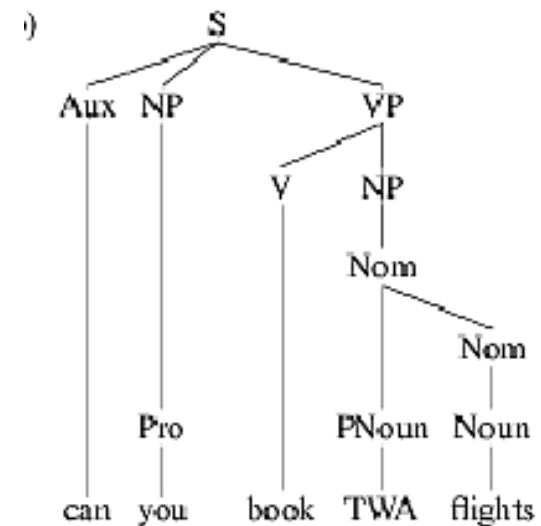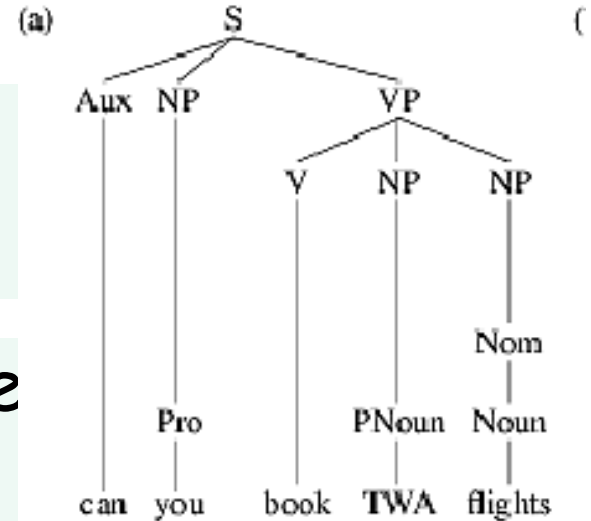
• **Estimate Prob. of parse tree**

> A. Sum of the probs of all the rules applied

> B. Product of the probs of all the rules applied

• **Estimate Prob. of a sentence**

A. Sum of the probs of all the parse trees

B. Product of the probs of all the parse trees

# PCFGs are used to….

- **Estimate Prob. of parse tree**

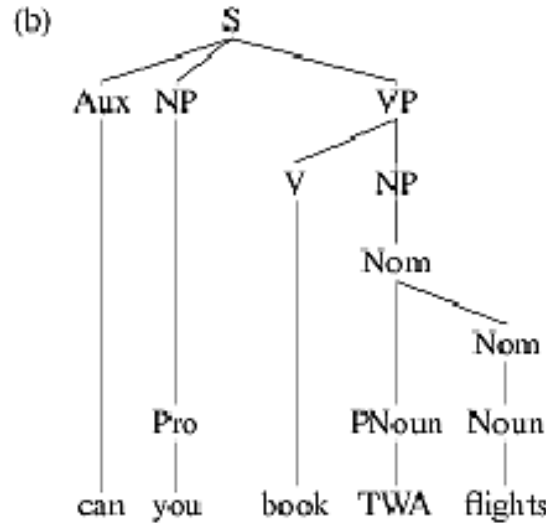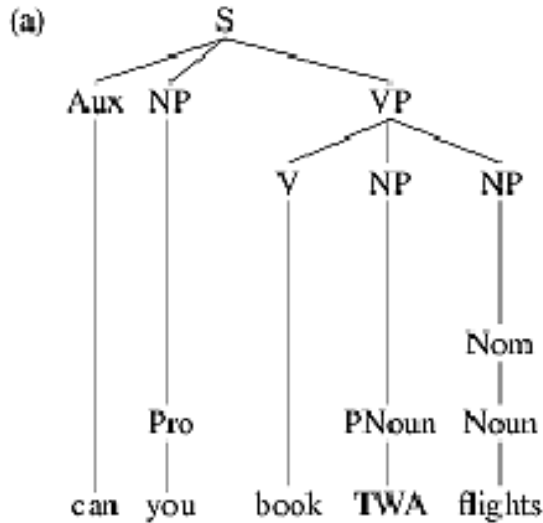$$P(Tree) = \prod_{node \in Tree} P(expansion\ for\ node)$$

- **Estimate Prob. to sentences**

$$P(Sentence) = \sum_{Tree \in Sentence\ parses} P(Tree)$$

# Example

$$P(Tree^a) = .15 * .4 * \ldots = \boxed{3.8 \cdot 10^{-7}} \qquad P(Tree^b) = .15 * .4 * \ldots = \boxed{4.3 \cdot 10^{-7}}$$



P("can you book TWA flights")

$$\boxed{8.1 \cdot 10^{-7}}$$

(a)

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

(a) →

(b)

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

← (b)

21

# Lecture Overview

- Recap English Syntax and Parsing
- Key Problem with parsing: Ambiguity
- Probabilistic Context Free Grammars (PCFG)
- **Treebanks and Grammar Learning (acquiring the probabilities)**

# Treebanks

- **Definition: corpora in which each sentence has been paired with a parse tree**
- **These are generally created**
  - Parse collection with parser
  - human annotators revise each parse
- **Requires detailed annotation guidelines**
  - POS tagset
  - Grammar
  - instructions for how to deal with particular grammatical constructions.

# Penn Treebank

- **Penn TreeBank is a widely used treebank.**

▪Most well known is the Wall Street Journal section of the Penn TreeBank.

▪1 M words from the 1987-1989 Wall Street Journal.

```
(  (S ('' '')
     (S-TPC-2
       (NP-SBJ-1 (PRP We) )
       (VP (MD would)
         (VP (VB have)
           (S
             (NP-SBJ (-NONE- *-1) )
             (VP (TO to)
               (VP (VB wait)
                 (SBAR-TMP (IN until)
                   (S
                     (NP-SBJ (PRP we) )
                     (VP (VBP have)
                       (VP (VBN collected)
                         (PP-CLR (IN on)
                           (NP (DT those)(NNS assets))))))))))))))))
   (,  ,) ('' '')
   (NP-SBJ (PRP he) )
   (VP (VBD said)
     (S (-NONE- *T*-2) ))
   (. .) ))
```

# Treebank Grammars

- **Such grammars tend to contain lots of rules….**

- **For example, the Penn Treebank has 4500 different rules for VPs! Among them...**
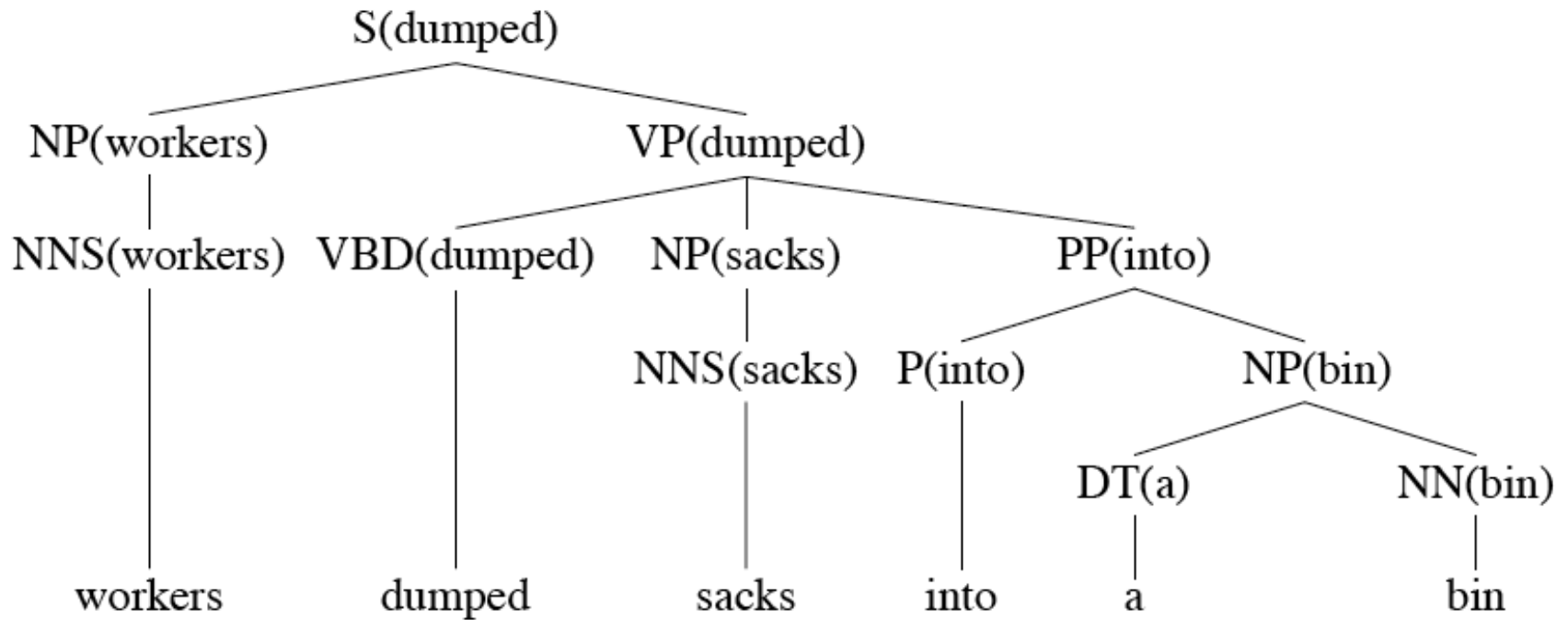
```
VP  →  VBD  PP
VP  →  VBD  PP  PP
VP  →  VBD  PP  PP  PP
VP  →  VBD  PP  PP  PP  PP
```

# Heads in Trees

- **Finding heads in treebank trees is a task that arises frequently in many applications.**

  - **Particularly important in statistical parsing**

- **We can visualize this task by annotating the nodes of a parse tree with the heads of each corresponding node.**
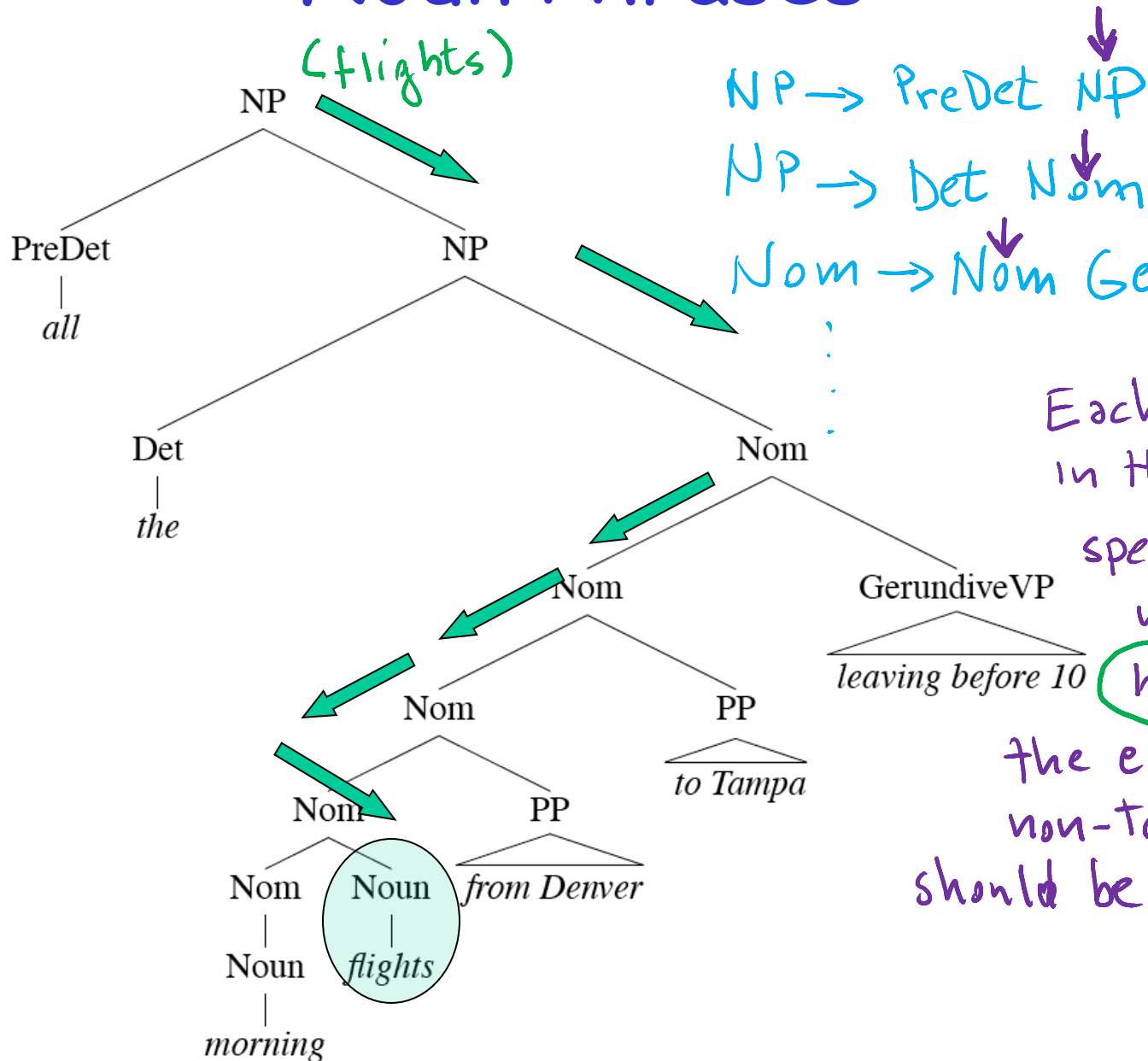
# Lexically Decorated Tree

# Head Finding

- **The standard way to do head finding is to use a simple set of tree traversal rules specific to each non-terminal in the grammar.**

- **Each rule in the PCFG specifies where the head of the expanded non-terminal should be found**

# Noun Phrases

(flights)

NP

PreDet
|
*all*

NP

Det
|
*the*

Nom

Nom

GerundiveVP

*leaving before 10*

Nom

PP

*to Tampa*

Nom

PP

*from Denver*

Nom

Noun
|
*flights*

Noun
|
*morning*

NP → PreDet NP

NP → Det Nom

Nom → Nom GerundiveVP

Each rule in the PCFG specifies where the head of the expanded non-terminal should be found

29

# Acquiring Grammars and Probabilities

**Manually parsed text corpora (e.g., PennTreebank)**

- **Grammar**: read it off the parse trees

 Ex: if an NP contains an ART, ADJ, and NOUN then we create the rule  NP -> ART ADJ NOUN.

- **Probabilities**:

$$P(A \to \alpha | A) = \frac{count(A \to \alpha)}{\sum_{\forall \gamma} count(A \to \gamma)} = \frac{count(A \to \alpha)}{count(A)}$$

 Ex: if the *NP -> ART ADJ NOUN* rule is used 50 times and all NP rules are used 5000 times, then the rule's probability is … .01

# Example

If you look at all the parse trees in the
bank you find three rules for NP

How many times

① NP → ART ADJ NOUN                     50

② NP → NOUN                            4000

③ NP → PRONOUN                   $\dfrac{950}{5000}$ total #
                                            of NP
                                            expansions

$P(①|NP) = 50/5000 = .01$

$P(②|NP) = 4000/5000 = .8$

$P(③|NP) = 950/5000 = .19$     $also = 1 - (.01 + .8)$

# Next class (Nov 13)

- **Parsing Probabilistic CFG: CKY parsing**
- **PCFG in practice: Modeling Structural and Lexical Dependencies**