

# FORMAL GRAMMARS

---

Ing. R. Tedesco. PhD, AA 20-21

(mostly from: Speech and Language Processing - Jurafsky and Martin)

# Today

- Context-free grammar
- Grammars for English
- Treebanks

# Syntax

- Why should you care?
- Grammars (and parsing) are key components in many applications
  - Grammar checkers
  - Dialogue management
  - Question answering
  - Information extraction
  - Machine translation

# Syntax

- Key notions that we'll cover
  - Constituency
  - Grammatical relations and Dependency
    - Heads
- Key formalism
  - Context-free grammars
- Resources
  - Treebanks

# Constituency

- The basic idea here is that **groups of words** within utterances can be shown to act as **single units** → *constituent*
- And, in a given language, these units form **coherent classes** that can be shown to behave in similar ways
  - With respect to their internal structure
  - And with respect to other units in the language

# Constituency

- Internal structure
  - Constituents have internal structures
  - We can describe such internal structures
- External behavior
  - For example, we can say that noun phrases can come before verbs

# Constituency

- For example, it makes sense that the following are all *noun phrases* in English...

Harry the Horse  
the Broadway coppers  
they

a high-class spot such as Mindy's  
the reason he comes into the Hot Box  
three parties from Brooklyn

- Why? One piece of evidence is that they can all precede verbs.
  - This is external evidence

# Grammars and Constituency

- Of course, there's nothing easy or obvious about how we come up with right set of constituents and the rules that govern how they combine...
- That's why there are so many different theories of grammar and competing analyses of the same data.
- The approach to grammar, and the analyses, adopted here are very generic...
- ...and don't correspond to any modern linguistic theory of grammar



# Context-Free Grammars

- Context-free grammars (CFGs)
  - Also known as
    - Phrase structure grammars
    - Backus-Naur form
- Consist of
  - Rules
  - Terminals
  - Non-terminals

# Context-Free Grammars

- **Terminals**
  - We'll take these to be words
- **Non-terminals**
  - The *constituents*
    - Like noun phrase, verb phrase, etc.
  - The *pre-terminals* (e.g., POS tags)
  - A special symbol  $S$  to start from
- **Rules**
  - a single non-terminal  $\rightarrow$  any number of terminals and non-terminals

# Some NP Rules

- Here are some rules for our noun phrases

$NP \rightarrow Det\ Nominal$

$NP \rightarrow ProperNoun$

$Nominal \rightarrow Noun \mid Nominal\ Noun$

- Together, these describe two kinds of NPs.
  - One that consists of a determiner followed by a nominal
  - And another that says that proper names are NPs.
  - The third rule illustrates two things
    - An explicit disjunction
      - Two kinds of nominals
    - A recursive definition
      - Same non-terminal on the right and left-side of the rule

# The L<sub>0</sub> Grammar

GRAMMAR RULES	EXAMPLES
$S \rightarrow NP VP$	I + want a morning flight
$NP \rightarrow$ <ul style="list-style-type: none"> <li><math>Pronoun</math></li> <li><math>Proper-Noun</math></li> <li><math>Det Nominal</math></li> </ul>	I Los Angeles a + flight
$Nominal \rightarrow$ <ul style="list-style-type: none"> <li><math>Nominal Noun</math></li> <li><math>Noun</math></li> </ul>	morning + flight flights
$VP \rightarrow$ <ul style="list-style-type: none"> <li><math>Verb</math></li> <li><math>Verb NP</math></li> <li><math>Verb NP PP</math></li> <li><math>Verb PP</math></li> </ul>	do want + a flight leave + Boston + in the morning leaving + on Thursday
$PP \rightarrow Preposition NP$	from + Los Angeles

# The $L_0$ Grammar

## GRAMMAR LEXICON

*Noun* → *flights* | *breeze* | *trip* | *morning* | ...

*Verb* → *is* | *prefer* | *like* | *need* | *want* | *fly*

*Adjective* → *cheapest* | *non-stop* | *first* | *latest*  
| *other* | *direct* | ...

*Pronoun* → *me* | *I* | *you* | *it* | ...

*Proper-Noun* → *Alaska* | *Baltimore* | *Los Angeles*  
| *Chicago* | *United* | *American* | ...

*Determiner* → *the* | *a* | *an* | *this* | *these* | *that* | ...

*Preposition* → *from* | *to* | *on* | *near* | ...

*Conjunction* → *and* | *or* | *but* | ...

Pre-terminals (POS tags)

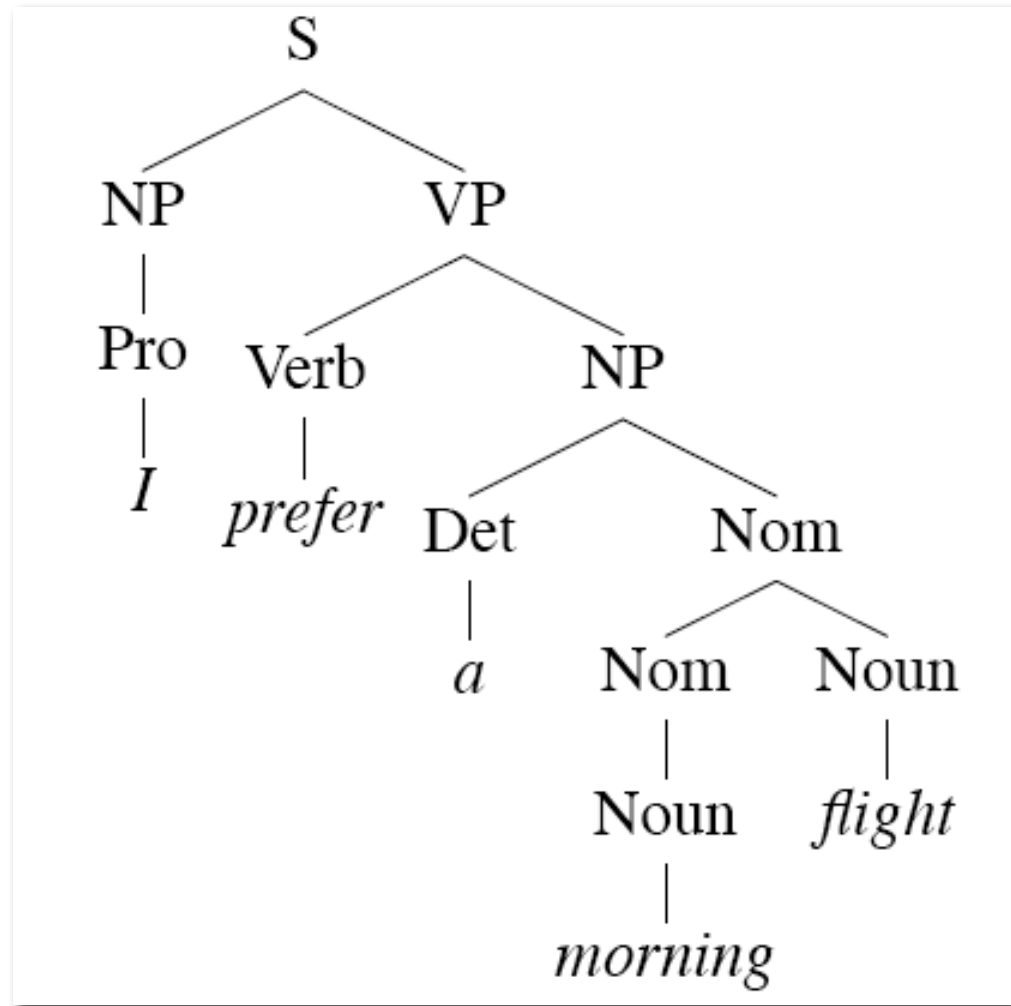
Terminals (words)

# Generativity

- You can view these rules as either analysis or synthesis machines
  - Generate strings in the language
  - Reject strings not in the language
  - Impose structures (trees) on strings in the language

# Derivations

- A **derivation** is a sequence of rules applied to a string that *accounts* for that string
  - Covers **all** the elements in the string
  - Covers **only** the elements in the string



# Definition

- More formally, a CFG consists of

$N$  a set of **non-terminal symbols** (or **variables**)

$\Sigma$  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,

$\beta$  is a string of symbols from the infinite set of strings  $(\Sigma \cup N)^*$

$S$  a designated **start symbol**

- **Parsing** is the process of taking a string and a grammar and returning a (multiple?) parse tree(s) for that string



# Noun Phrases

- Let's consider the following rule in more detail...

*NP* → *Det Nominal*

# Agreement

- By ***agreement***, we have in mind constraints that hold among various constituents that take part in a rule or set of rules
- For example, in English, determiners and the head nouns in NPs have to agree in their number.

This flight

Those flights

\*This flights

\*Those flight

\*=wrong sentence

# Problem

- Our NP rules are clearly deficient since they don't capture this constraint
  - *NP* → *Det Nominal*
    - Accepts, and assigns correct structures, to grammatical examples (*this flight*)
    - But its also happy with incorrect examples (\*these flight)
  - Such a rule is said to *overgenerate*.
  - We'll come back to this in a bit

# Verb Phrases

- English VPs consist of a head verb along with zero or more following constituents which we'll call *arguments*.

$VP \rightarrow Verb$  disappear

$VP \rightarrow Verb$   $NP$  prefer a morning flight

$VP \rightarrow Verb$   $NP$   $PP$  leave Boston in the morning

$VP \rightarrow Verb$   $PP$  leaving on Thursday

# Subcategorization

- But, even though there are many valid VP rules in English, not all verbs are allowed to participate in all those VP rules.
- We can subcategorize the verbs in a language according to the sets of VP rules that they participate in.
- This is a modern take on the traditional notion of transitive/intransitive and complements.
- Modern grammars may have 100s or such classes.

# Subcategorization

- Sneeze: John sneezed
- Find: Please find [a flight to NY]<sub>NP</sub>
- Give: Give [me]<sub>NP</sub>[a cheaper fare]<sub>NP</sub>
- Help: Can you help [me]<sub>NP</sub>[with a flight]<sub>PP</sub>
- Prefer: I prefer [to leave earlier]<sub>TO-VP</sub>
- ...

# Subcategorization

- \*John sneezed the book
  - \*I prefer United has a flight
  - \*Give with a flight
- 
- As with agreement phenomena, we need a way to formally express the constraints

# Possible CFG Solution

- Possible solution for agreement and subcategorization: create ad-hoc sub-classes
- $S \rightarrow NP VP$ 
  - $SG\_S \rightarrow SG\_NP SG\_VP$
  - $PL\_S \rightarrow PL\_NP PL\_VP$
- $NP \rightarrow Det Nom$ 
  - $SG\_NP \rightarrow SG\_Det SG\_Nom$
  - $PL\_NP \rightarrow PL\_Det PL\_Nom$
- $VP \rightarrow V NP$ 
  - $PL\_VP \rightarrow PL\_V NP$
  - $SL\_VP \rightarrow SG\_V NP$
- ...
- ...



# CFG Solution for Agreement

- It works and stays within the power of CFGs
- But it's ugly and it doesn't scale because of the interaction among the various constraints explodes the number of rules in our grammar
- CFGs appear to be just about what we need to account for a lot of basic syntactic structure in

- Other solutions:

- Feature-based grammars

$S \rightarrow NP VP$

$\langle NP \text{ AGREEMENT} \rangle = \langle VP \text{ AGREEMENT} \rangle$

A feature associated  
to NP and VP  
Words carry features

- LFG, HPSG, Construction grammar, XTAG, etc.
  - Do not try to find a solution!

# Treebanks

- Treebanks are corpora in which each sentence has been paired with a parse tree (presumably the right one).
- These are generally created
  - By first parsing the collection with an automatic parser
  - And then having human annotators correct each parse as necessary.
- This generally requires detailed annotation guidelines that provide a POS tagset, a grammar and 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

- Treebanks implicitly define a grammar for the language covered in the treebank.
- Simply take the local rules that make up the sub-trees in all the trees in the collection and you have a grammar.
- Not complete, but if you have decent size corpus, you'll have a grammar with decent coverage.

# Treebank Grammars

- Such grammars tend to be very flat due to the fact that they tend to avoid recursion.
  - To ease the annotators burden
- For example, the Penn Treebank has 4500 different rules for VPs. Among them...

VP → VBD PP

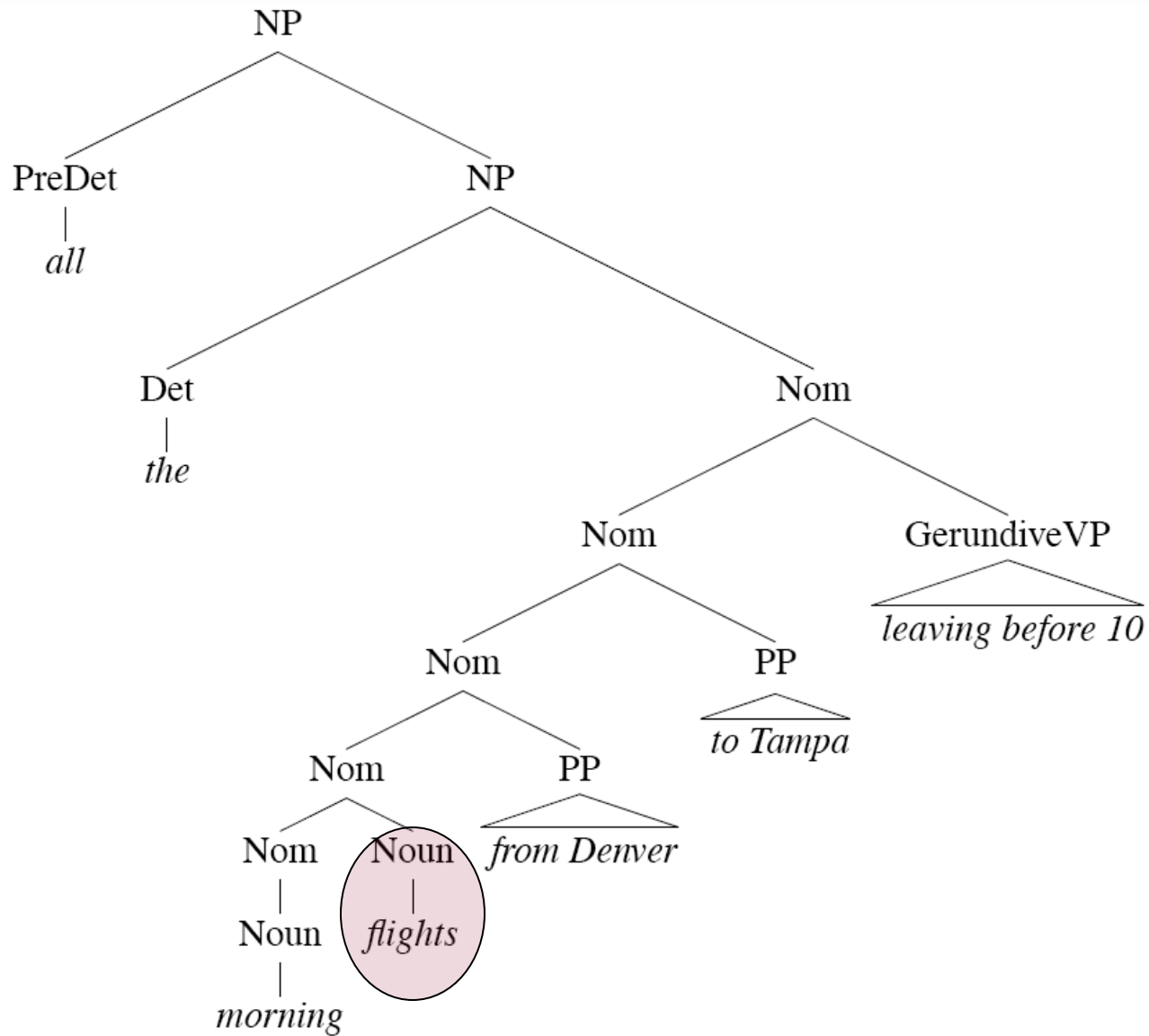
VP → VBD PP PP

VP → VBD PP PP PP

VP → VBD PP PP PP PP

- Instead of  $VP \rightarrow VBD PP_2$   
 $PP_2 \rightarrow PP_2 PP$

# Head



# Head

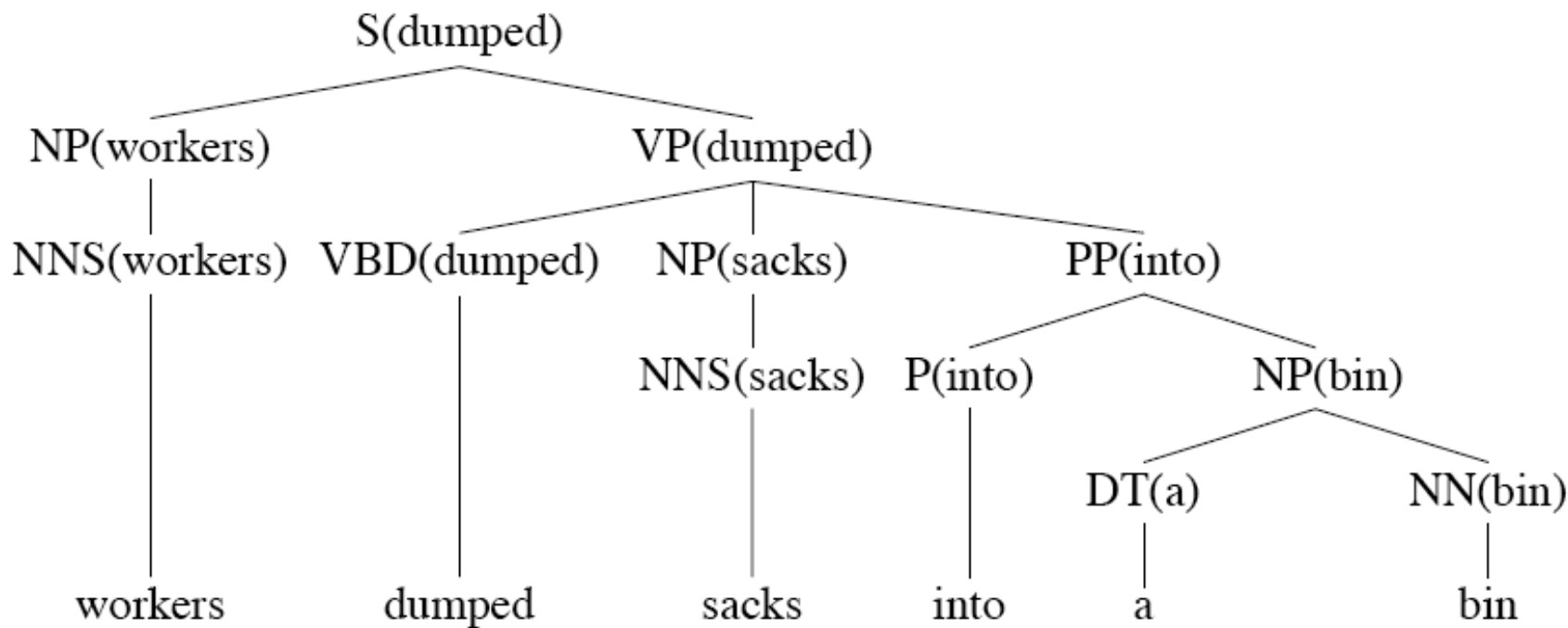
- Clearly this NP is really about *flights*.  
That's the central critical noun in this NP.  
Let's call that the *head*.
- We can dissect this kind of NP into the stuff that can come before the head, and the stuff that can come after it.
- Not only for NP...

# Heads in Trees

- Finding heads in treebank trees is a task that arises frequently in many applications.
  - Particularly important in statistical parsing
  - More on that later...
- We can visualize this task by annotating the nodes of a parse tree with the heads of each corresponding node.



# Lexically Decorated Tree



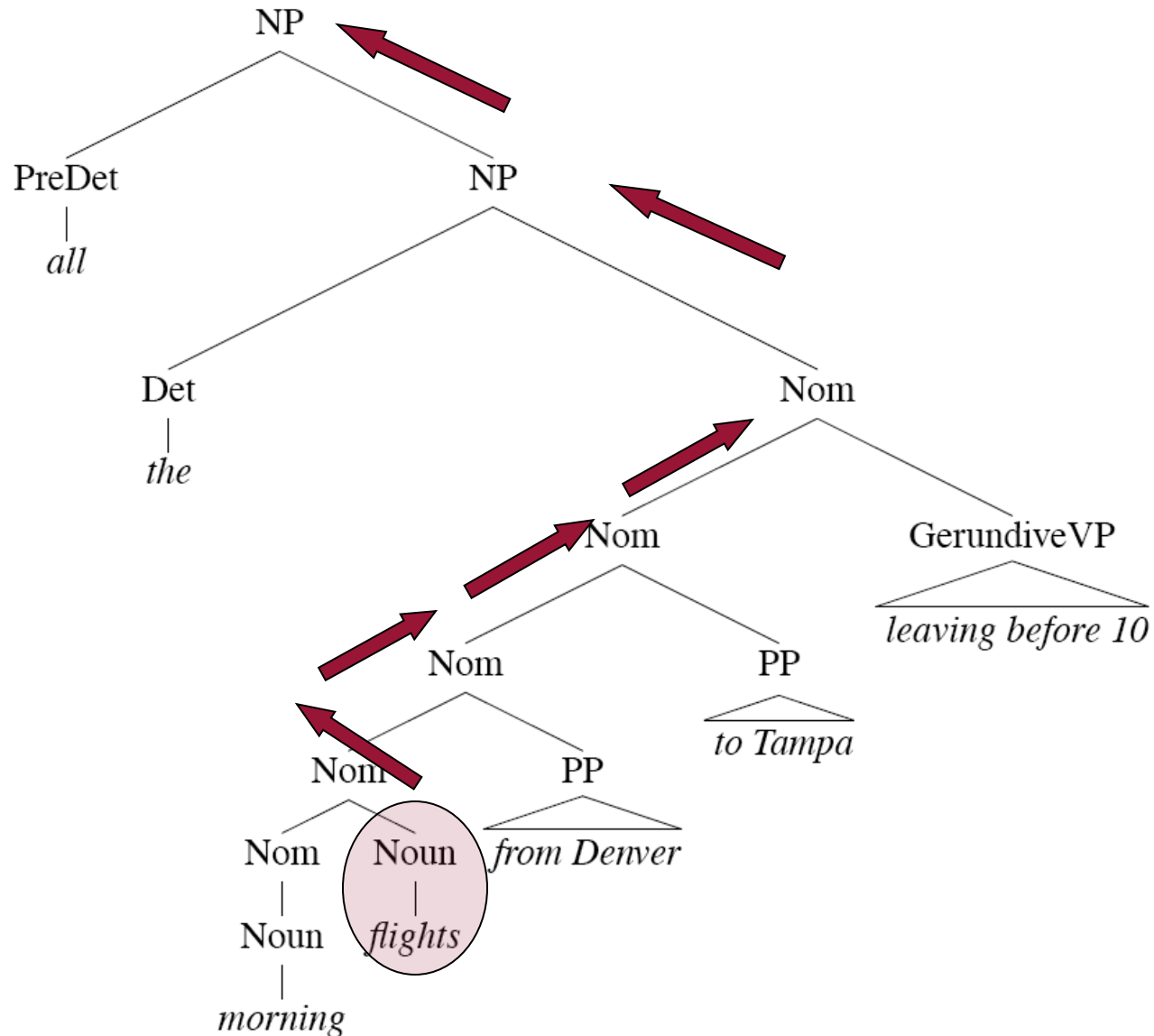
# Head Finding

- 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.
  - such traversal rules are “encoded” into *lexicalized grammar rules* (more on that later...)
- Derived from lexicalized treebanks

*VP* (dumped)  $\rightarrow$  *VDB* (dumped) *NP* (sacks) *PP* (into)

Means that in  $VP \rightarrow VDB\ NP\ PP$ , the head of *VDB* goes to *VP*

# Noun Phrases



# Summary

- Context-free grammars can be used to model various facts about the syntax of a language.
- When paired with parsers, such grammars constitute a critical component in many applications.
- Constituency is a key phenomena easily captured with CFG rules.
  - But agreement and subcategorization do pose significant problems
- Treebanks pair sentences in corpus with their corresponding trees.