

Speech and Language Processing

Parsing

Today

- Parsing with CFGs
 - Bottom-up, top-down
 - Ambiguity
 - Early & CKY parsing

Parsing

- Parsing with CFGs refers to the task of assigning proper trees to input strings
- Proper here means a tree that covers **all and only the elements of the input** and **has an S at the top**
- It doesn't actually mean that the system can select the correct tree from among all the possible trees

Parsing

- As with everything of interest, parsing involves a **search** which involves the making of choices
- We'll start with some basic (meaning bad) methods before moving on to the one or two that you need to know

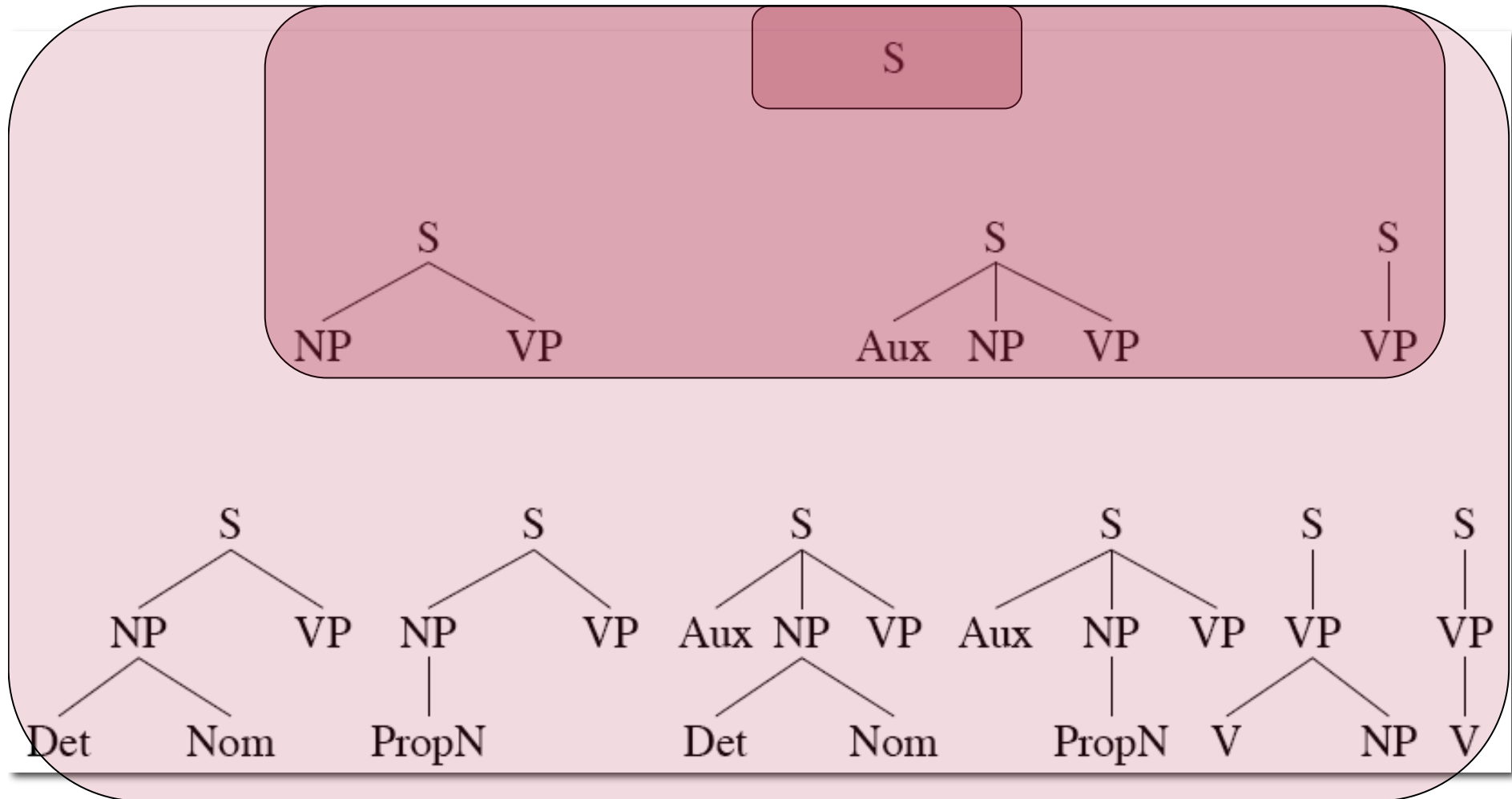
For Now

- Assume...
 - You have all the words already in some buffer
 - The input isn't POS tagged
 - We won't worry about morphological analysis
 - All the words are known
- These are all problematic in various ways, and would have to be addressed in real applications.

Top-Down Search

- Since we're trying to find trees rooted with an S (Sentences), why not start with the rules that give us an S .
- Then we can work our way down from there to the words.

Top Down Space



Bottom-Up Parsing

- Of course, we also want trees that cover the input words. So we might also start with trees that link up with the words in the right way.
- Then work your way up from there to larger and larger trees.

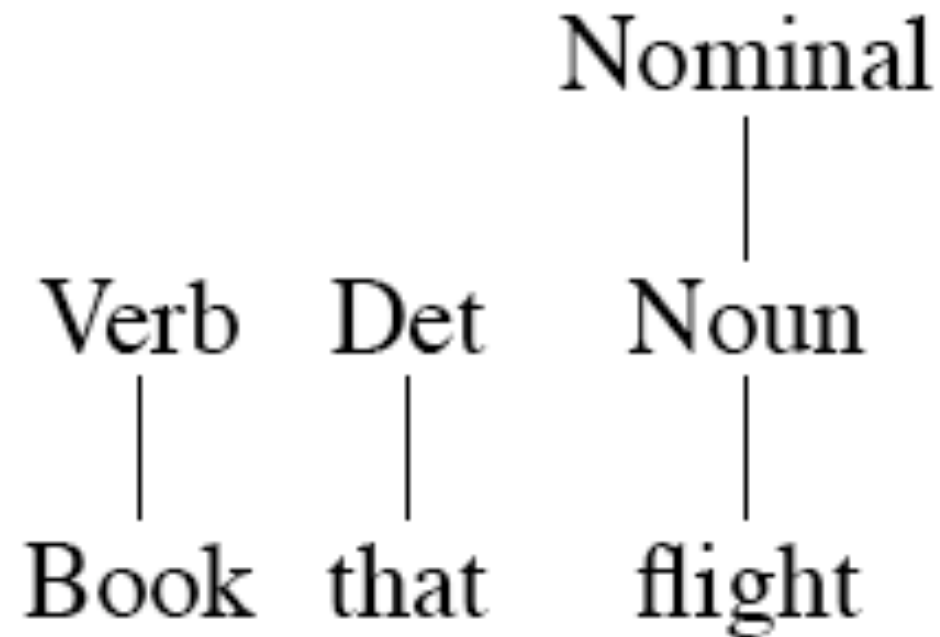
Bottom-Up Search

Book that flight

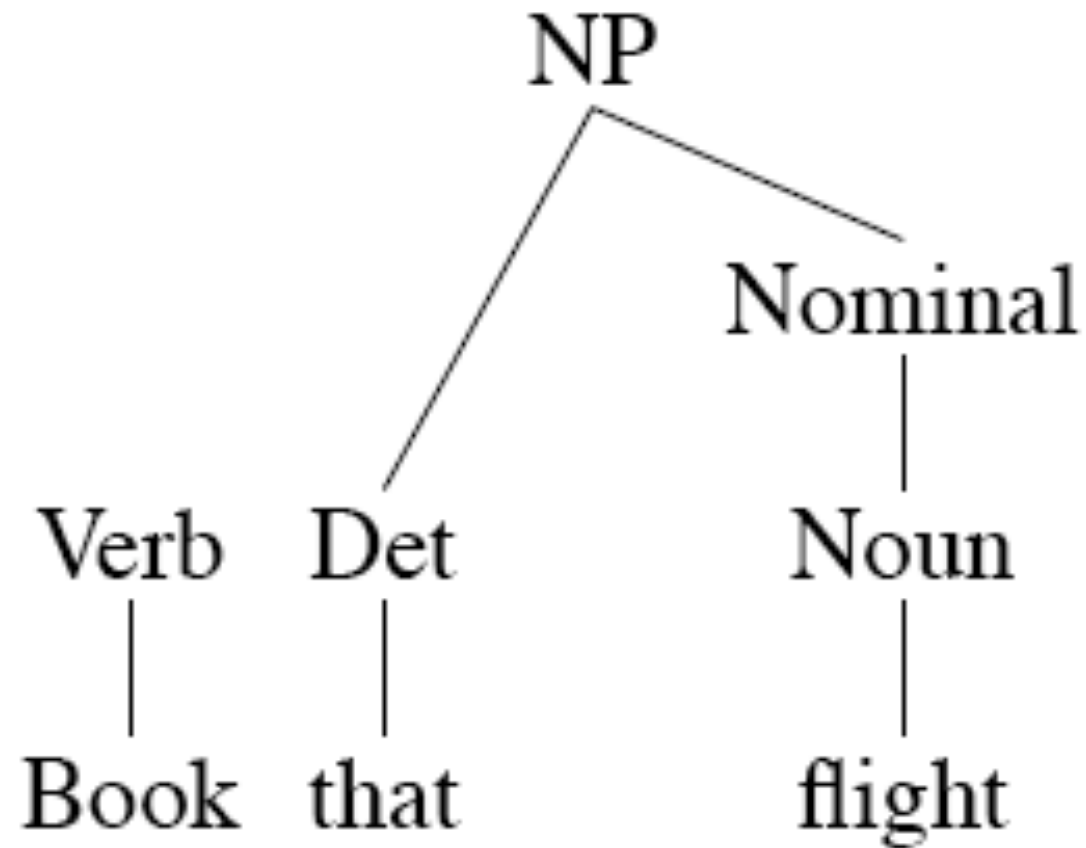
Bottom-Up Search

Verb Det Noun
| | |
Book that flight

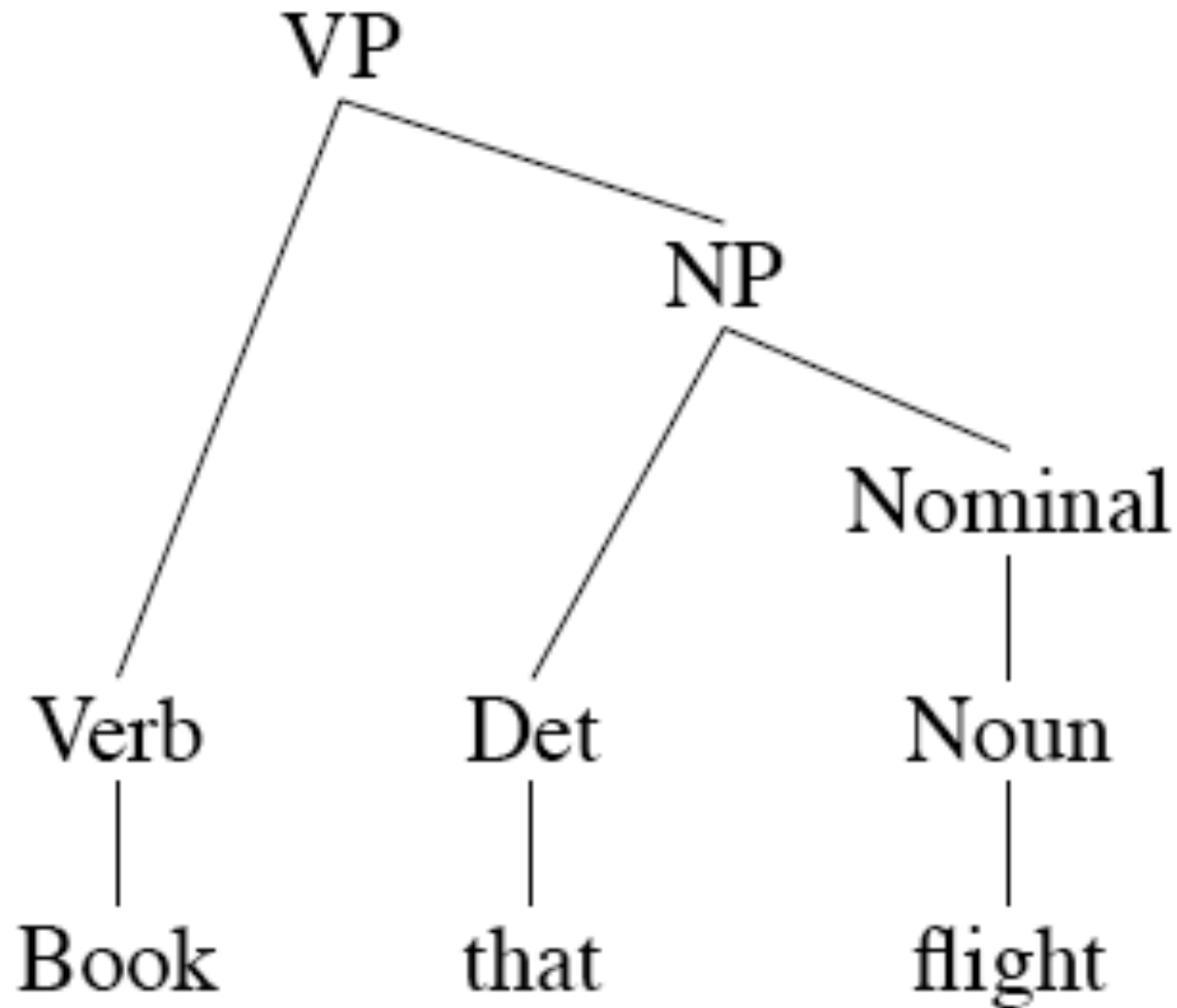
Bottom-Up Search



Bottom-Up Search



Bottom-Up Search



Top-Down and Bottom-Up

- **Top-down**
 - Only searches for trees that can be answers (i.e. S' 's)
 - But also suggests trees that are not consistent with any of the words
- **Bottom-up**
 - Only forms trees consistent with the words
 - But suggests trees that make no sense globally

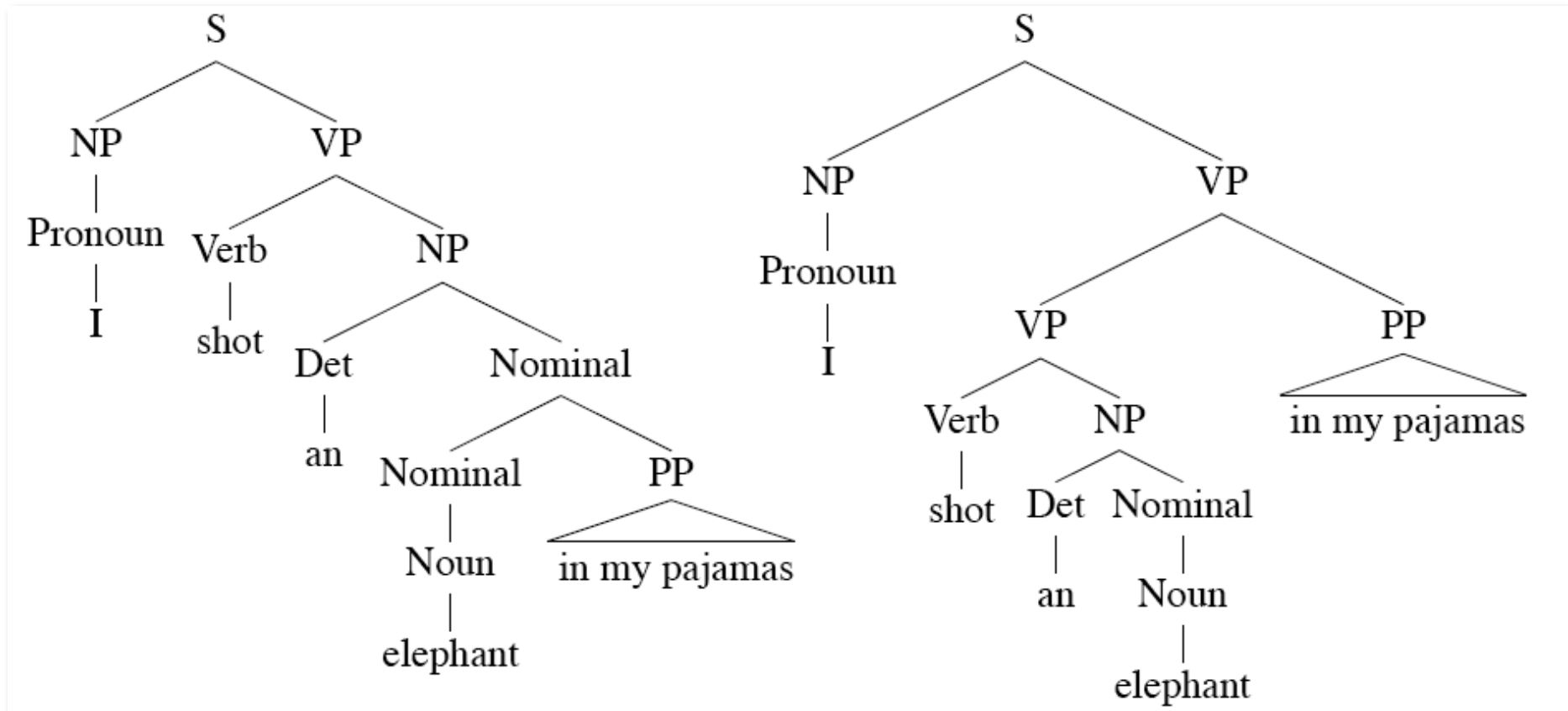
Control

- Of course, in both cases we left out how to keep track of the search space and how to make choices
 - Which node to try to expand next
 - Which grammar rule to use to expand a node
- One approach is called backtracking.
 - Make a choice, if it works out then fine
 - If not then back up and make a different choice

Problems

- Even with the best filtering, backtracking methods are doomed because of two inter-related problems
 - Ambiguity
 - Shared subproblems

Ambiguity

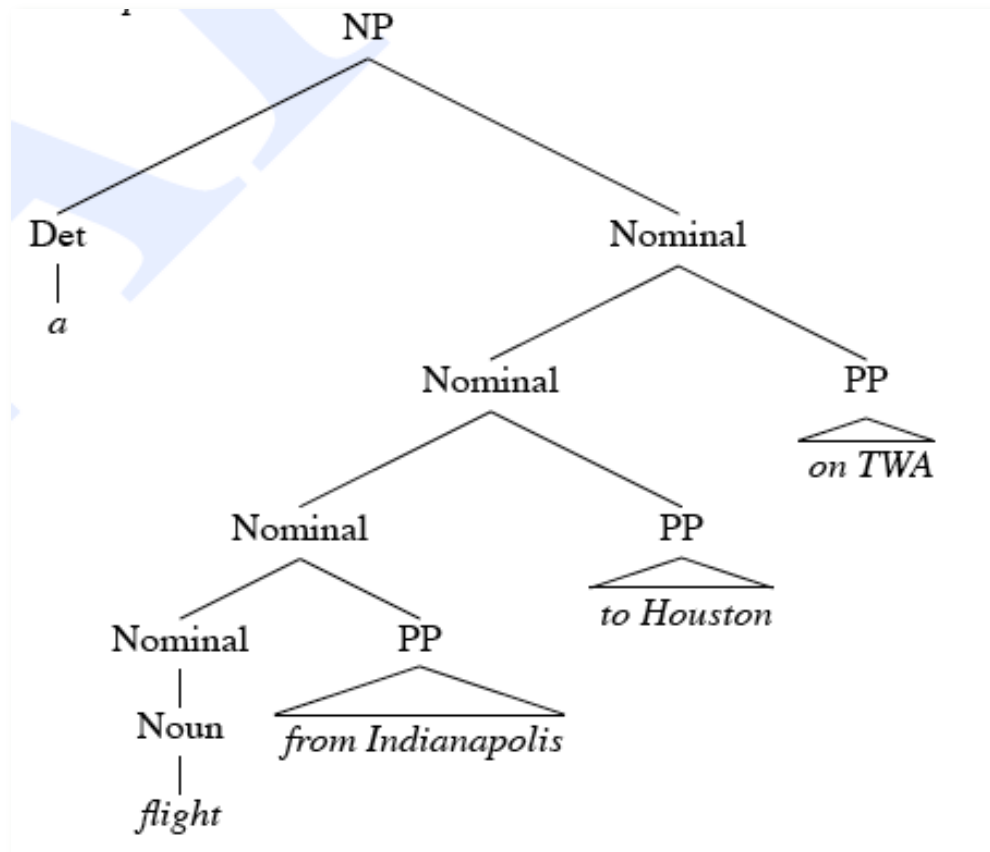


Shared Sub-Problems

- No matter what kind of search (top-down or bottom-up or mixed) that we choose.
 - We don't want to redo work we've already done.
 - Unfortunately, naïve backtracking will lead to duplicated work.

Shared Sub-Problems

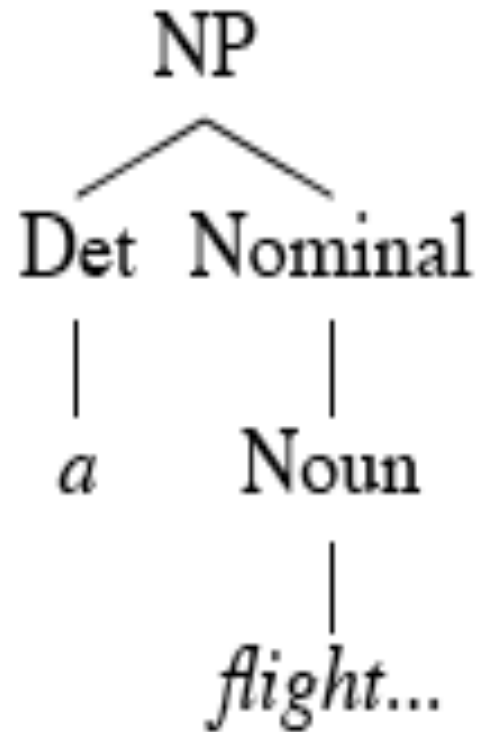
- Consider
 - A flight from Indianapolis to Houston on TWA



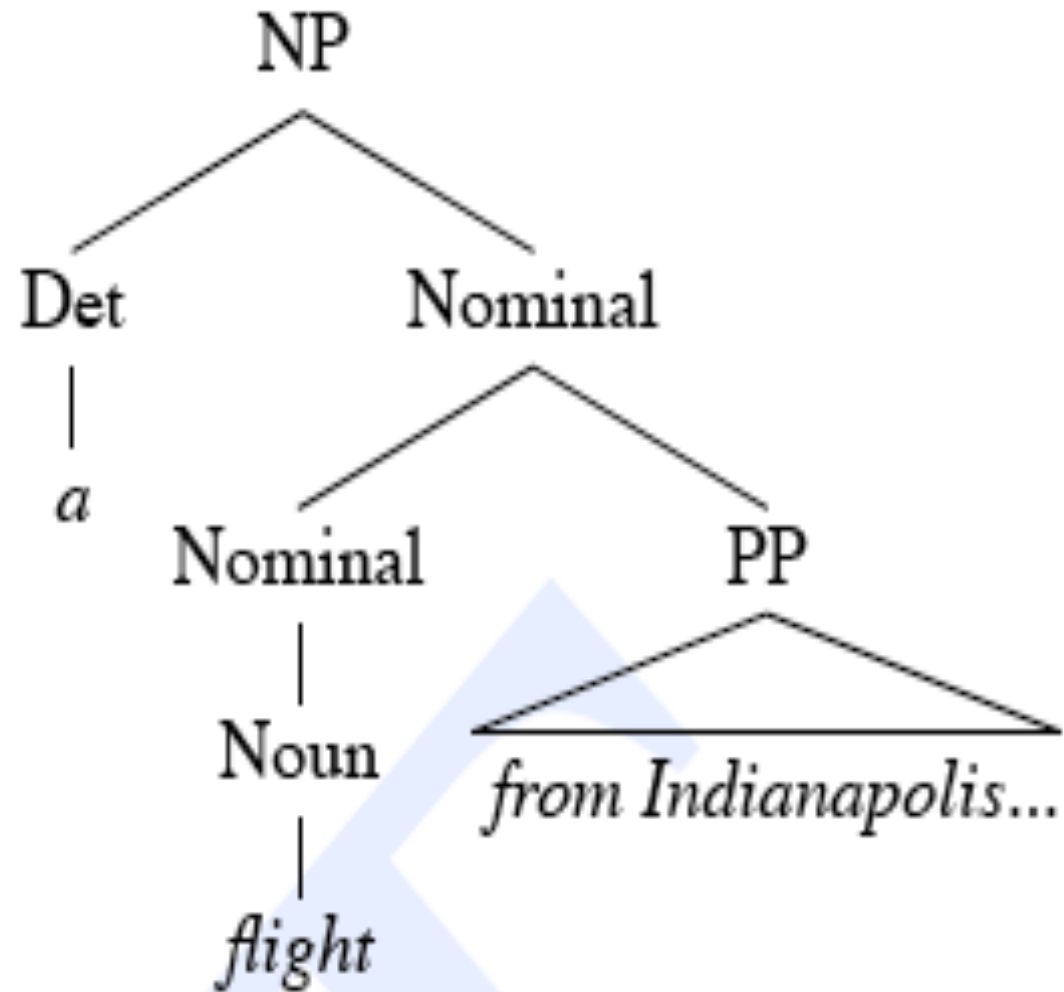
Shared Sub-Problems

- Assume a top-down parse making choices among the various Nominal rules.
- In particular, between these two
 - Nominal -> Noun
 - Nominal -> Nominal PP
- Statically choosing the rules in this order leads to the following bad results...

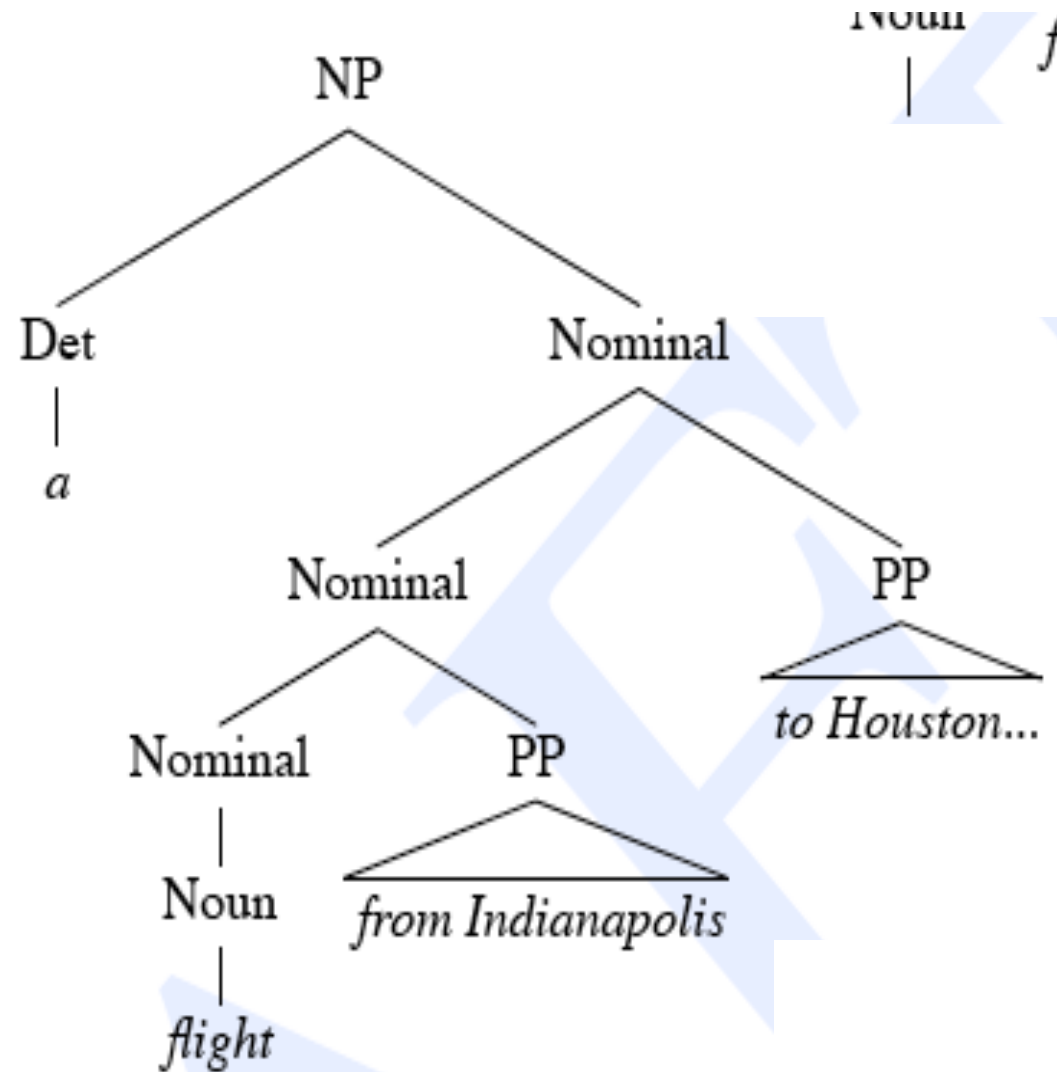
Shared Sub-Problems



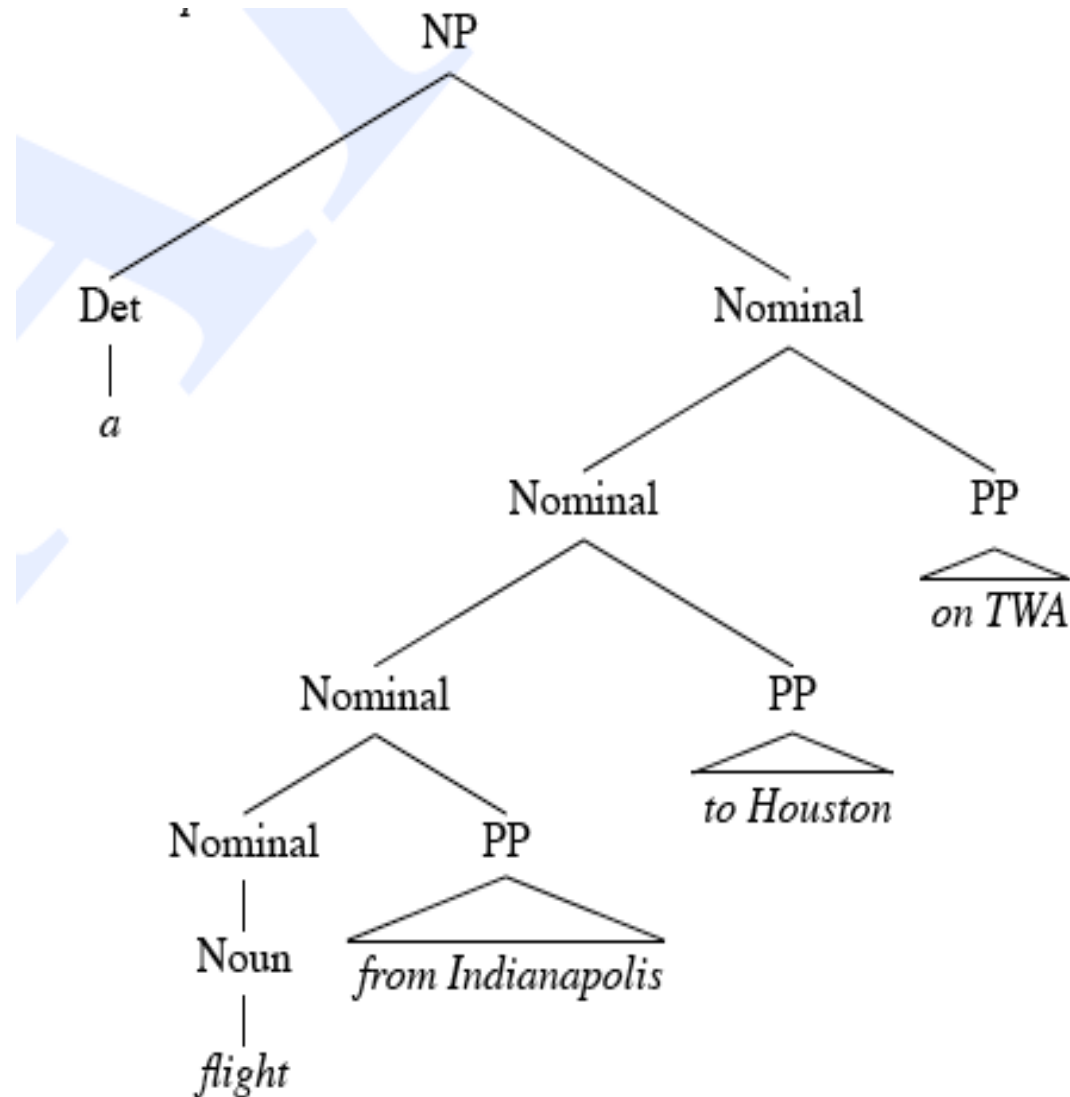
Shared Sub-Problems



Shared Sub-Problems



Shared Sub-Problems



Dynamic Programming

- DP search methods fill tables with partial results and thereby
 - Avoid doing avoidable repeated work
 - Solve exponential problems in polynomial time (well, no not really)
 - Efficiently store ambiguous structures with shared sub-parts.
- Two important approaches that roughly correspond to top-down and bottom-up approaches.
 - CKY
 - Earley

Ambiguity

- Both CKY and Earley will result in multiple **S** structures.
- They both efficiently store the sub-parts that are shared between multiple parses.
- And they obviously avoid re-deriving those sub-parts.
- But neither can tell us which one is right.

Ambiguity

- In most cases, humans don't notice incidental ambiguity (lexical or syntactic). It is resolved on the fly and never noticed.
- We'll try to model that with probabilities.