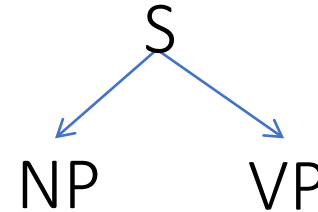# CS60075
# Natural Language Processing
# Autumn 2020

## Module 5: Part B
## L0 grammar and Phrase Structure Parsing

### Sep 30 2020

# Formal Grammars of English

# Context-free grammars (CFGs)

- Consist of
  - Rules
  - Terminals
  - Non-terminals
  - Start Symbol

- Specifies a set of tree structures that capture constituency and ordering in language

$N$    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** and a member of $N$
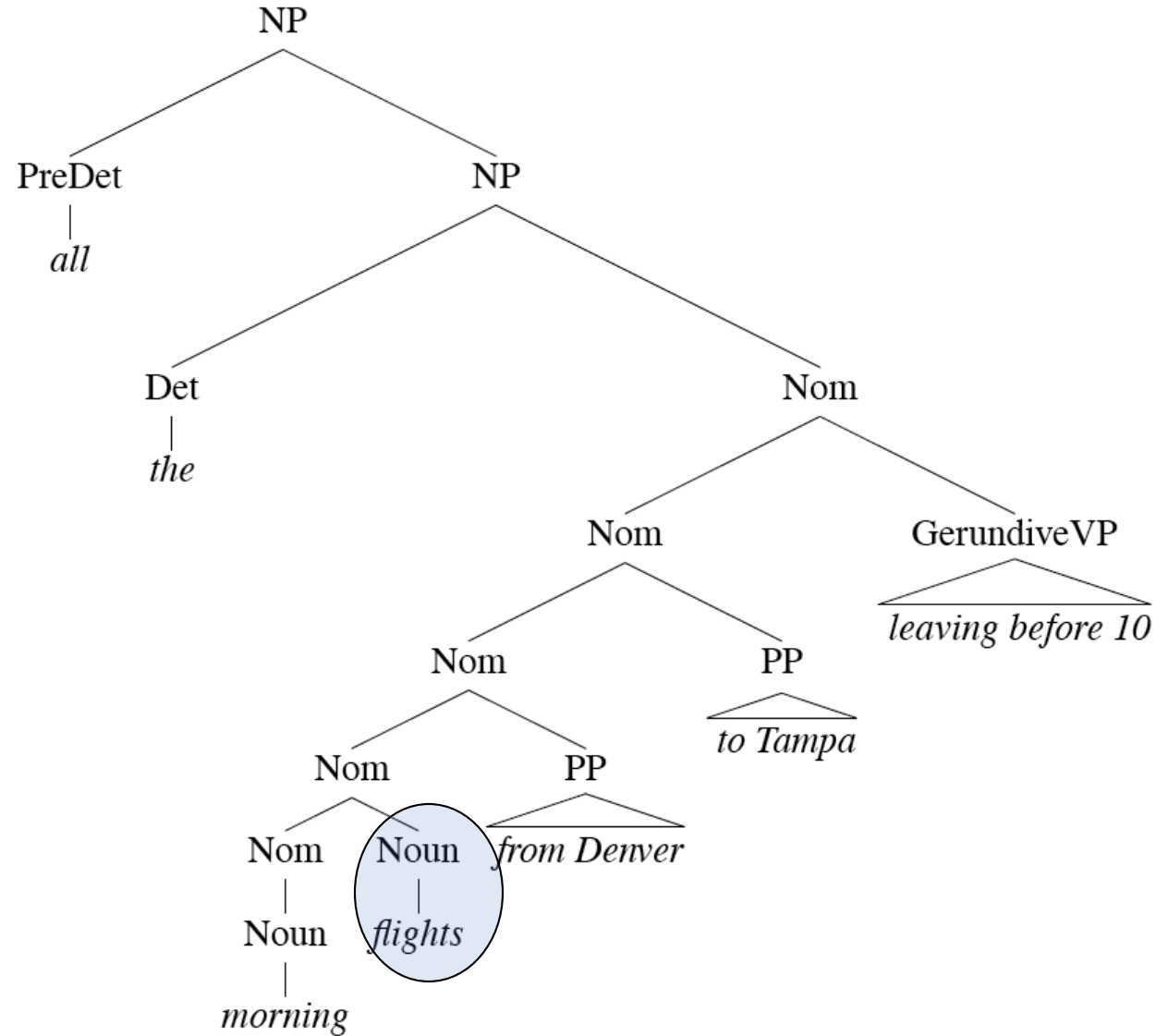
# Productions of CFG

- A CFG can be thought of in two ways:
  - a device for generating sentences
    (Derivation)
  - a device for assigning a structure to a given sentence.

- Some rules for noun phrases:

$$NP \rightarrow Det\ Nominal$$
$$NP \rightarrow ProperNoun$$
$$Nominal \rightarrow Noun\ |\ Nominal\ Noun$$

# Noun Phrases

# Nominals

- Contain the head and any pre- and post- modifiers of the head.
  - Pre-
    - Quantifiers, cardinals, ordinals...
      - *Three* cars
    - Adjectives
      - *large* cars

01-Oct-20

# Postmodifiers

- Three kinds
  - Prepositional phrases
    - *From Seattle*
  - Non-finite clauses
    - *Arriving before noon*
  - Relative clauses
    - *That serve breakfast*

- Same general (recursive) rules to handle these
  - *Nominal → Nominal PP*
  - *Nominal → Nominal GerundVP*
  - *Nominal → Nominal RelClause*

HINDI
*Nominal → PP Nominal*
*Nominal → Nominal GerundVP*
*Nominal → Nominal RelClause*

01-Oct-20

# Verb Phrases

- English *VP*s consist of a verb (the head) along with 0 or more *following* constituents which we'll call *arguments*.

| | | |
|---|---|---|
| VP | → Verb | disappear |
| VP | → Verb NP | prefer a morning flight |
| VP | → Verb NP PP | leave Boston in the morning |
| VP | → Verb PP | leaving on Thursday |

# Subcategorization

- 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 just an elaboration on the traditional notion of transitive/intransitive.

- Modern grammars have many such classes

01-Oct-20

# Subcategorization

- Sneeze:  John sneezed

- Find:  Please find [a flight to NY]$_{NP}$

- Give: Give [me]$_{NP}$[a cheaper fare]$_{NP}$

- Help: Can you help [me]$_{NP}$[with a flight]$_{PP}$

- Prefer: I prefer [to leave earlier]$_{TO-VP}$

- Told: I was told [United has a flight]$_S$

- …

01-Oct-20

# Generative Grammar

- The use of formal languages to model Generative natural languages is called **_generative grammar_** since the language is defined by the set of possible sentences "generated" by the grammar.

- You can view these rules as either analysis or synthesis engines
  - Generate strings in the language
  - Reject strings not in the language
  - Assign structures (trees) to strings in the language

# L0 Grammar

| Grammar Rules | | | Examples |
|---|---|---|---|
| *S* | → | *NP VP* | I + want a morning flight |
| *NP* | → | *Pronoun* | I |
| | \| | *Proper-Noun* | Los Angeles |
| | \| | *Det Nominal* | a + flight |
| *Nominal* | → | *Nominal Noun* | 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 |

# Sentence Types

- Declaratives: *A plane left.*

  *S ⟶ NP VP*

- Imperatives: *Leave!*
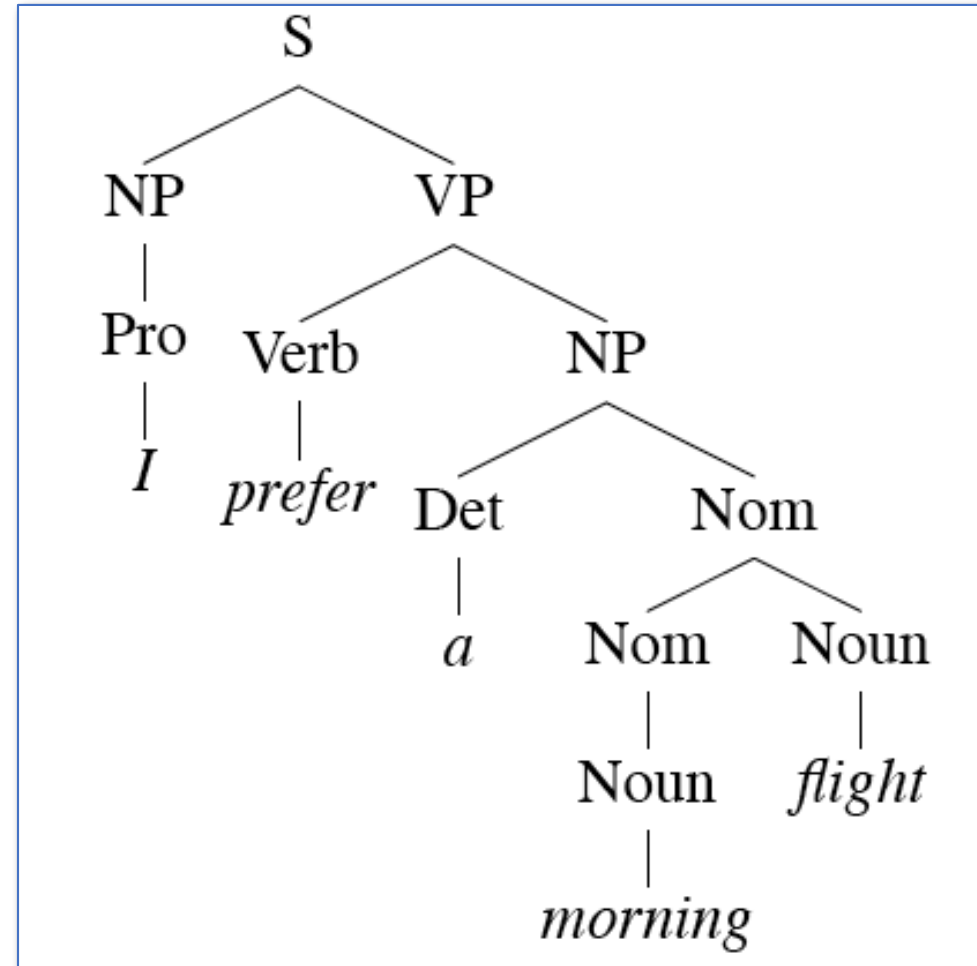
  *S ⟶ VP*

- Yes-No Questions: *Did the plane leave?*

  *S ⟶ Aux NP VP*

- WH Questions: *When did the plane leave?*

  *S ⟶ WH-NP Aux NP VP*

01-Oct-20

# 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

# Parsing

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

# Treebank

- A syntactically annotated corpus where every sentence is paired with a corresponding tree.

- The Penn Treebank project
  - treebanks from the Brown, Switchboard, ATIS, and Wall Street Journal corpora of English
  - treebanks in Arabic and Chinese.

- Others
  - the Prague Dependency Treebank for Czech,
  - the Negra treebank for German, and
  - the Susanne treebank for English
  - Universal Dependencies Treebank

# Penn Treebank

- Penn TreeBank is a widely used treebank.

Most well known part 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) ))
    (. .) ))
```

Speech and Language Processing - Jurafsky and Martin

```
((S
   (NP-SBJ (DT That)
     (JJ cold) (, ,)
     (JJ empty) (NN sky) )
   (VP (VBD was)
     (ADJP-PRD (JJ full)
       (PP (IN of)
         (NP (NN fire)
           (CC and)
           (NN light) ))))
   (. .) ))
                (a)
```

```
((S
   (NP-SBJ The/DT flight/NN )
   (VP should/MD
     (VP arrive/VB
       (PP-TMP at/IN
         (NP eleven/CD a.m/RB ))
       (NP-TMP tomorrow/NN )))))
                (b)
```

**Figure 11.7**   Parsed sentences from the LDC Treebank3 version of the Brown (a) and ATIS (b) corpora.
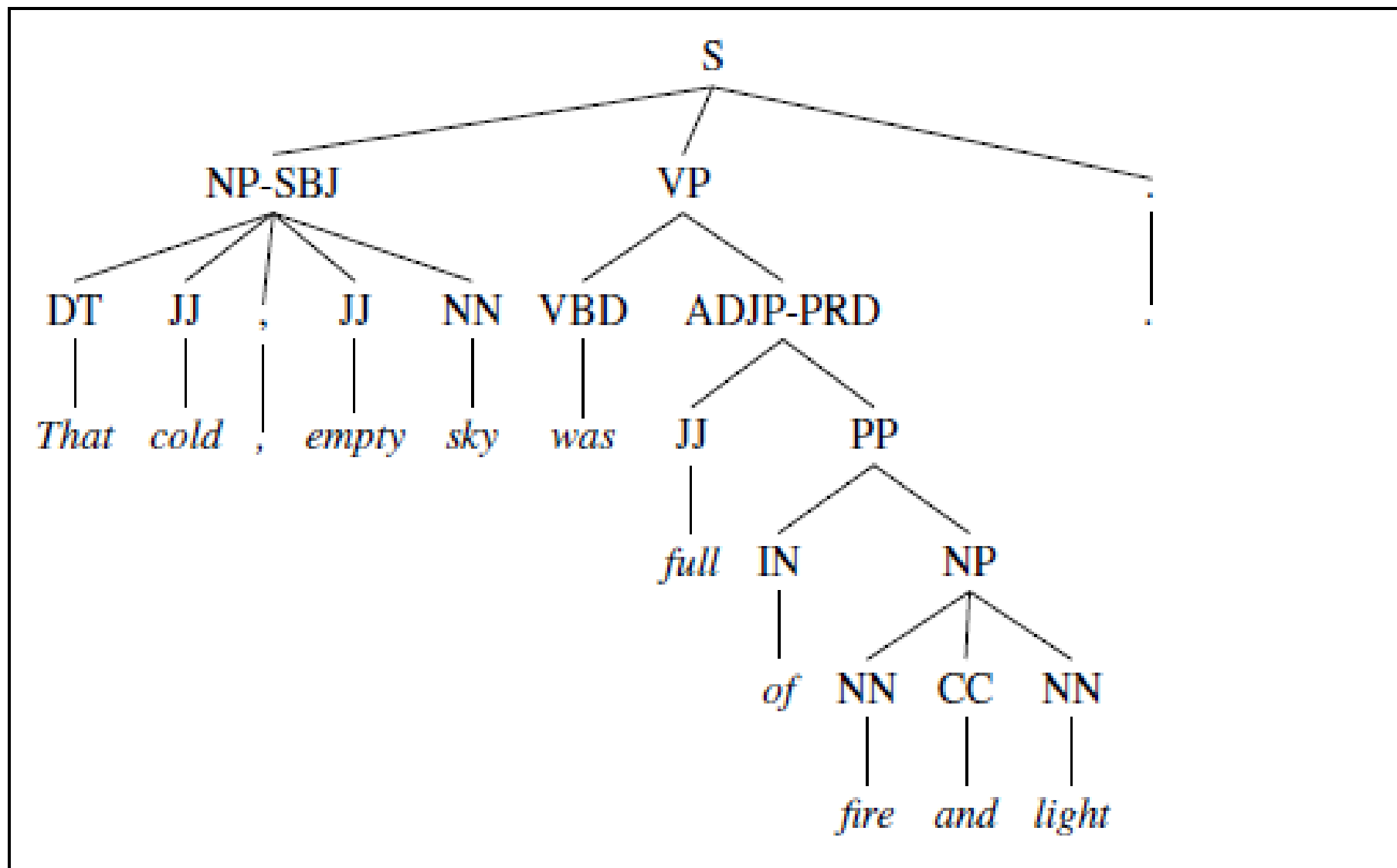
**Figure 11.8** The tree corresponding to the Brown corpus sentence in the previous figure.

# Treebanks as Grammars

- The sentences in a treebank implicitly constitute a grammar of the language represented by the corpus being annotated.

- Simply take the local rules that make up the sub-trees in all the trees in the collection and you have a grammar
  - The WSJ section gives us about 12k rules

# Parsing

- Parsing with CFGs refers to the task of assigning proper trees to input strings

- Proper here means a tree that covers <span style="color:darkred">all and only the elements of the input</span> and <span style="color:green">has an S at the top</span>

- It doesn't mean that the system can select the correct tree from among all the possible trees

# Syntactic Analysis (Parsing)

- Automatic methods of finding the syntactic structure for a sentence
  - Symbolic methods: a phrase grammar or another description of the structure of language is required.
  The chart parser.
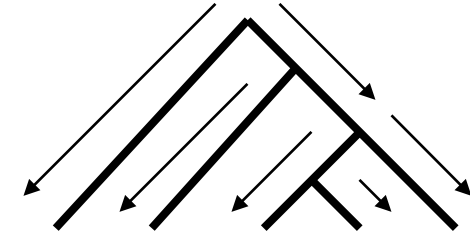  - Statistical methods: a text corpus with syntactic structures is needed (a treebank)

# Search Framework

- Think about parsing as a form of search…
  - A search through the space of possible trees given an input sentence and grammar

Speech and Language
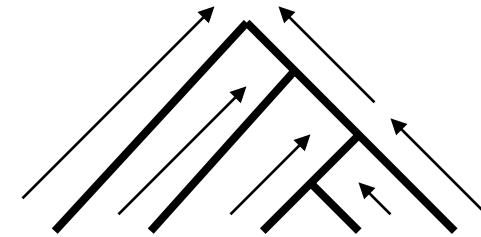Processing - Jurafsky and Martin

# How to parse

Top-down: Start at the top of the tree with an S node, and work your way down to the words.

Bottom-up: Look for small pieces that you know how to assemble, and work your way up to larger pieces.

# Summary

- CFGs appear to be just about what we need to account for a lot of basic syntactic structure in English.

- But there are problems
  - That can be dealt with adequately, although not elegantly, by staying within the CFG framework.

- There are simpler, more elegant, solutions that take us out of the CFG framework (beyond its formal power)
  - LFG, HPSG, Construction grammar, XTAG, etc.

# Top-Down Search

- Since we're trying to find trees rooted with an *S* (Sentences)
  - Start with the rules that give us an *S*.
  - Then we can work our way down from there to the words.

Speech and Language
Processing - Jurafsky and Martin

# 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.

Speech and Language
Processing - Jurafsky and Martin

# Bottom-Up Search

Book  that  flight

Speech and Language
Processing - Jurafsky and Martin

# Bottom-Up Search



$$\text{Verb} \quad \text{Det} \quad \text{Noun}$$
$$\text{Book} \quad \text{that} \quad \text{flight}$$

Speech and Language
Processing - Jurafsky and Martin

# Bottom-Up Search



```
                              Nominal
                                 |
         Verb     Det          Noun
          |        |             |
        Book     that         flight
```

Speech and Language
Processing - Jurafsky and Martin

# Bottom-Up Search

Speech and Language
Processing - Jurafsky and Martin

# Bottom-Up Search

Speech and Language
Processing - Jurafsky and Martin

# Issues

- Ambiguity
- Shared subproblems

# Ambiguity

Speech and Language
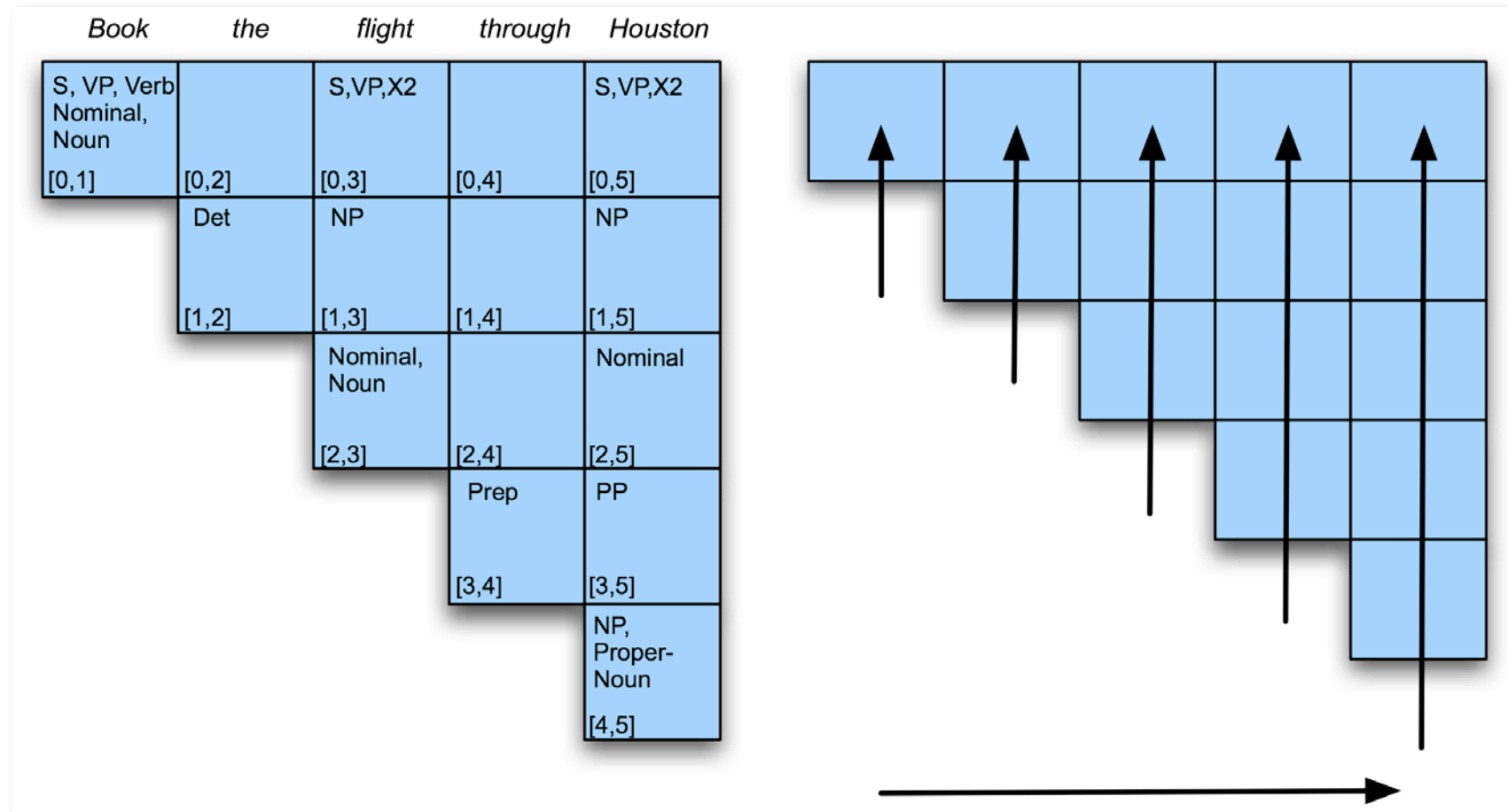Processing - Jurafsky and Martin

# Dynamic Programming

- DP search methods fill tables with partial results and thereby
  - Avoid doing avoidable repeated work
  - Solve exponential problems in polynomial time (ok, not really)
  - Efficiently store ambiguous structures with shared sub-parts.
- We'll cover one approach that corresponds to a bottom-up strategy
  - CKY

Speech and Language
Processing - Jurafsky and Martin

# CKY Algorithm



| | Book | the | flight | through | Houston |
|---|---|---|---|---|---|
| | S, VP, Verb Nominal, Noun [0,1] | [0,2] | S,VP,X2 [0,3] | [0,4] | S,VP,X2 [0,5] |
| | | Det [1,2] | NP [1,3] | [1,4] | NP [1,5] |
| | | | Nominal, Noun [2,3] | [2,4] | Nominal [2,5] |
| | | | | Prep [3,4] | PP [3,5] |
| | | | | | NP, Proper-Noun [4,5] |

Speech and Language Processing - Jurafsky and Martin

# CKY Algorithm

**function** CKY-PARSE(*words, grammar*) **returns** *table*

  **for** $j \leftarrow$ **from** $1$ **to** LENGTH(*words*) **do**

    $table[j-1, j] \leftarrow \{A \mid A \rightarrow words[j] \in gram$

    **for** $i \leftarrow$ **from** $j - 2$ **downto** $0$ **do**

      **for** $k \leftarrow i + 1$ **to** $j - 1$ **do**

        $table[i,j] \leftarrow table[i,j] \cup$

$$\{A \mid A \rightarrow BC \in grammar,$$
$$B \in table[i,k],$$
$$C \in table[k,j]\}$$

Looping over the columns

Filling the bottom cell

Filling row i in column j

Looping over the possible split locations between i and j.

Check the grammar for rules that link the constituents in [i,k] with those in [k,j]. For each rule found store the LHS of the rule in cell [i,j].

Speech and Language Processing - Jurafsky and Martin