

# Artificial Intelligence

Subject Code: 20A05502T

**UNIT IV - Natural Language for Communication**

**SYNTACTIC ANALYSIS - PARSING**

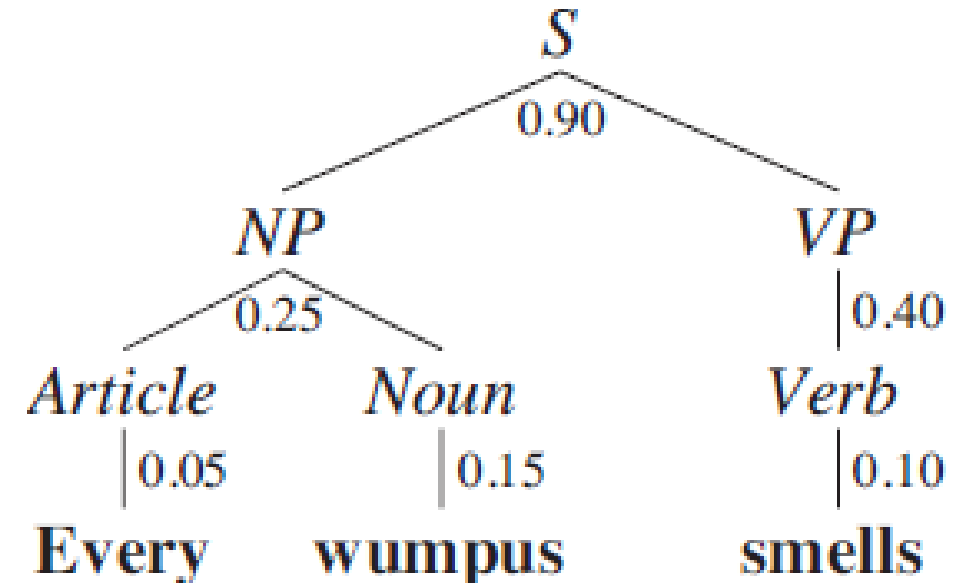
# SYNTACTIC ANALYSIS (PARSING)

- **Parsing** is the process of analyzing a string of words to uncover its **phrase structure**, according to the rules of a grammar.
- **Top Down Parsing**
- The **S**, starting symbol and search top down for a tree that has the words as its leaves, or
- **Bottom up Parsing**
- Start with the words and search bottom up for a tree that culminates in an **S**.
- Both top-down and bottom-up parsing can be inefficient, because they can end up repeating effort in areas of the search space that lead to dead ends.

<i>List of items</i>	<i>Rule</i>
<i>S</i>	
<i>NP VP</i>	$S \rightarrow NP VP$
<i>NP VP Adjective</i>	$VP \rightarrow VP Adjective$
<i>NP Verb Adjective</i>	$VP \rightarrow Verb$
<i>NP Verb <b>dead</b></i>	$Adjective \rightarrow \mathbf{dead}$
<i>NP <b>is dead</b></i>	$Verb \rightarrow \mathbf{is}$
<i>Article Noun <b>is dead</b></i>	$NP \rightarrow Article Noun$
<i>Article <b>wumpus is dead</b></i>	$Noun \rightarrow \mathbf{wumpus}$
<i>the wumpus is dead</i>	$Article \rightarrow \mathbf{the}$

# Parse Tree

- Parse tree for the sentence “**Every wumpus smells**”, according to the grammar  $E_0$ .
- Each interior node of the tree is labeled with its probability.
- The probability of the tree as a whole is  $0.9 \times 0.25 \times 0.05 \times 0.15 \times 0.40 \times 0.10 = 0.0000675$ .
- Since this tree is the only parse of the sentence, that number is also the probability of the sentence.
- The tree can also be written in linear form as
- **[S [NP [Article every] [Noun wumpus]][VP [Verb smells]]]**.



- The E0 grammar generates a wide range of English sentences such as the following:
- John is in the pit
- The wumpus that stinks is in 2 2
- Mary is in Boston and the wumpus is near 3 2

# Drawbacks

- Consider the following two sentences:
  - Have the students in section A of Computer Science III take the exam.
  - Have the students in section A of Computer Science III taken the exam?
- Even though they share the first 10 words, these sentences have very different parses, because the first is a command and the second is a question.
- A **left-to-right parsing** algorithm would have to guess whether
- the **first word** is part of a **command or a question** and
- will not be able to tell if the guess is correct until at least the eleventh word, *take or taken*.
- If the algorithm guesses wrong, it will have to backtrack all the way to the first word and reanalyze the whole sentence under the other interpretation.
- To avoid this source of inefficiency we can use **dynamic programming**.

# Dynamic Programming:

- *In dynamic programming, “every time we analyze a substring, store the results so we won’t have to reanalyze it later”.*
- For example,
- “the students in section B of Computer Science III” is an NP,
- record that result in a data structure as a **chart**.
- Algorithms that do this are called **chart parsers**.
- In **context-free grammars**, any phrase that was found in the context of one branch of the search space, can work as well in any other branch of the search space.
- There are many types of chart parsers;
- a bottom-up version called the **CYK algorithm**, after its inventors, **John Cocke, Daniel Younger, and Tadeo Kasami**



# CYK algorithm

- it requires a grammar with all rules in one of two very specific formats:
  - Lexical rules of the form  $X \rightarrow \mathbf{word}$ , and
  - Syntactic rules of the form  $X \rightarrow Y Z \mid \mathbf{word}$
- This grammar format, called **Chomsky Normal Form**,
- any context-free grammar can be automatically transformed into Chomsky Normal Form.
- The CYK algorithm uses
  - space of  $O(n^2m)$  for the P table,
  - time  $O(n^3m)$ .
- where  $n$  is the number of words in the sentence, and
- $m$  is the number of nonterminal symbols in the grammar.



Thank You