# CS 3313 Foundations of Computing:

## **CYK Parsing Algorithm**

http://gw-cs3313.github.io

1

#### **Simplification and Parsing**

- 1. Simplification rules: transform a grammar such that:
  - Resulting grammar generates the same language
  - and has "more efficient" production rules in a specific format
- 2. Normal Forms: express all CFGs using a standard "format" for how the production rules are specified
  - Definition of CFGs places no restrictions on RHS of production
  - It is convenient (for parsing algorithms) to restrict to a standard form
    - Chomsky Normal Form (CNF) or Greiback Normal Form (GNF)
- 3. Parsing Algorithm: Design a parsing algorithm that takes a grammar in a standard form (CNF) to check if string w is generated by grammar G.

# Procedure to transform any CFG to Chomsky Normal Form

- A CFG is said to be in *Chomsky Normal Form* if every production is of one of these two forms:
  - 1.  $A \rightarrow BC$  (right hand side is two variables).
  - 2. A  $\rightarrow a$  (right hand side is a single terminal).
- Theorem: If L is a CFL, then L { λ} has a CFG in CNF.
  - Note: Theorem 2.4 implies every string on RHS of prodution is either a single terminal or has length > 2.
    - This is our starting point when converting to CNF form
- Question: property of parse trees for CNF grammars ?

3

3

#### **CNF**

- $G_1$  with production rules:
  - $\circ S \rightarrow AS \mid a$
  - $\circ A \rightarrow SA \mid b$
- Is  $G_1$  in CNF?
- *G*<sub>2</sub> with production rules:
- $P: S \rightarrow ABa \ A \rightarrow aab \ B \rightarrow Ac$
- Is  $G_2$  in CNF?

#### Testing for Membership – a Parsing Algorithm

- Simple algorithm: Convert CFG to a Greibach Normal Form (all productions are of the form  $A \rightarrow a\alpha$ )
  - For string w of length n, we have n derivation steps.
  - At each step, explore all productions.
  - Time:  $O(|P|^n)$  this is exponential (in length of input string w)
- Example:  $S \rightarrow aSB \mid bSA \mid aB \quad A \rightarrow a \quad B \rightarrow b$
- w = baba

5

### **Testing Membership (Parsing)**

- Determine if w is in L(G).
- Can we do better than the simple method with GNF with  $O(|P|^n)$  for string of length n?
- Assume G is in CNF.
  - Or convert the given grammar to CNF.
  - $w = \lambda$  is a special case, solved by testing if the start symbol is nullable.
- Cocke Younger Kashimi Algorithm (*CYK*) is a good example of dynamic programming and runs in time  $O(n^3)$ , where n = |w|.

6

#### **CYK Algorithm notations**

- Important: these notations are a bit different from notations in the book, but the end algorithm works in the same manner
- Input string w has length n i.e, consists of n terminal symbols:

$$w = a_1 a_2 ... a_n$$
 where each  $a_i \in T$ 

- Ex: w = abcaab  $a_1 = a a_2 = b a_3 = c,...$
- Define a substring  $x_{ii}$  (of w) as the substring starting at position i and having length j
  - $x_{13} = abc$   $x_{22} = bc$   $x_{33} = caa$
- For a substring  $x_{ij}$ , define  $V_{ij}$  to be set of variables that derive  $x_{ij}$ 
  - $V_{ij} = \{ A \mid A =>^* x_{ij} \}$
  - Ex:  $V_{33} = \{ A \mid A => *x_{33} = caa \}$

#### **Setting up our solution/algorithm: Notations**

- Input string w = abcaab
- Define a substring  $x_{ii}$  (of w) as the substring starting at position i and having length *j* 
  - Ex:
- $x_{13} = abc x_{22} = bc$
- $x_{33} = caa$
- $x_{15}$ = abcaa
- $w=x_{16}=abcaab$
- For a substring  $x_{ij}$ , define  $V_{ij}$  to be set of variables that derive  $x_{ij}$ 
  - $V_{ii} = \{ A \mid A = >^* x_{ii} \}$

#### **Algorithm**

- Claim is that we can construct  $V_{ij}$  iteratively
- Basis:  $V_{il} = \{ A \mid A \rightarrow x_{il} \text{ is a production } \}$
- Ind.  $A = >^* x_{ij}$  iff  $A \to BC$  and for some k, 1 <= k <= j,  $B = >^* x_{ik}$  and  $C = >^* x_{i+k,j-k}$
- Since k, j-k are  $\leq j$  the IH holds.
- w is in L(G) iff  $S \in V_{In}$  (since  $w = x_{In}$ )

```
V_{ij} = \{A \mid A \rightarrow BC, \text{ and} \} for some k, B is in V_{ik} and C is in V_{i+k,j-k}
```

9

#### **CYK Algorithm**

Input: CFG G=(V,T,P,S) in CNF, Input string w of length n

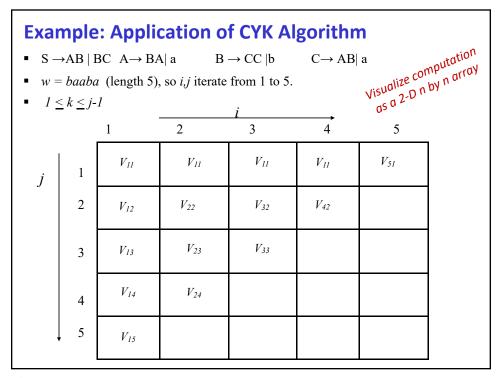
#### **Time Complexity**

- Step 1: takes O(n) to examine each of the n symbols
  - Assume P is a constant.
- Step 2:  $O(n^3)$ 
  - Outer j loop iterates O(n)
  - The *i* loop iterates O(n)
  - For each of the  $n^2$  iterations, the k loop iterates O(n)
- Dynamic programming formulation
  - Construct solution for size *n* in terms of sizes *n-1* 
    - Principle of optimality needs to hold

11

#### **Example: Application of CYK Algorithm**

- $S \rightarrow AB \mid BC$   $A \rightarrow BA \mid a$
- $B \rightarrow CC | b \quad C \rightarrow AB | a$
- w = baaba (length 5), so i,j iterate from 1 to 5.
- Some sample  $V_{ij}$
- To compute  $V_{3l}$ ,  $x_{3l} = a$ .  $V_{3l} = \{ X \mid X \rightarrow a \text{ is in } P \}$ 
  - $V_{31} = \{A, C\}$
- To compute  $V_{12}$ :  $X \rightarrow YZ$  in P and
  - check if  $Y \in V_{11}$  and  $Z \in V_{21}$
- To compute  $V_{23}: X \rightarrow YZ$  in P and
  - Check for Y in  $V_{21}$  and Z in  $V_{32}$
  - Check for Y in  $V_{22}$  and Z in  $V_{41}$



| Example: Application of CYK Algorithm  |   |                |  |   |  |  |  |
|--|---|----------------|--|---|--|--|--|
| ■ $S \rightarrow AB \mid BC \mid A \rightarrow BA \mid a$ $B \rightarrow CC \mid b$ $C \rightarrow AB \mid a$  |   |                |  |   |  |  |  |
| <b>Example: Application of CYK Algorithm</b> S $\rightarrow$ AB   BC A $\rightarrow$ BA  a B $\rightarrow$ CC  b C $\rightarrow$ AB  a $w = baaba$ (length 5), so $i,j$ iterate from 1 to 5.  Visualize computation of CYK Algorithm  Visualize computation of CYK Algorithm  Visualize computation of CYK Algorithm |   |                |  |   |  |  |  |
|  |   |                |  | i |  |  |  |
|  |   |                |  |   |  |  |  |
| j  |   |                |  |   |  |  |  |
|  |   |                |  |   |  |  |  |
|  |   |                |  |   |  |  |  |
|  | , | <i>i=1 j=5</i> |  |   |  |  |  |

#### Visualize computation **Example: Application of CYK Algorithm** visuumee computurion as a 2-D n by n array • $S \rightarrow AB \mid BC \mid A \rightarrow BA \mid a$ $B \rightarrow CC \mid b$ $C \rightarrow AB | a$ w = baaba (length 5), so *i,j* iterate from 1 to 5 • $V_{12} = \{ X \mid X \rightarrow YZ, Y \text{ in } V_{11} = (B) \text{ and } Z \text{ in } V_{21} = (A, C) \}$ 3 $x_{31}=a$ $\sum_{x_{2I}=a}^{2}$ $x_{4I} = b$ $x_{51}=a$ $1 \atop x_{II} = b$ $V_{31} =$ $V_{41} =$ $V_{5I} =$ $V_{2I} =$ $V_{II} =$ { *A*, *C*} { B} { *A*,*C*} { *A*,*C*} { B} 1 j $V_{12} = {S,A}$ 2 3 $V_{51} = {B}$ 4 5 $V_{15}$

15

| Example: CYK Algorithm  |     |   |  |   |              |              |                |   |
|---|-----|---|--|---|--------------|--------------|----------------|---|
| ■ $S \rightarrow AB \mid BC  A \rightarrow BA \mid a$ $B \rightarrow CC \mid b$ $C \rightarrow AB \mid a$   |     |   |  |   |              |              |                |   |
| • $V_{24} = \{X \mid X \rightarrow YZ, Y \text{ in } V_{21} \text{ and } Z \text{ in } V_{33}\} \cup \{X \mid X \rightarrow YZ, Y \text{ in } V_{22} \text{ and } Z \text{ in } V_{42}\}$ |     |   |  |   |              |              |                |   |
| $U\{X \mid X \rightarrow YZ, Y \text{ in } V_{23} \text{ and } Z \text{ in } V_{51}\}$  |     |   |  |   |              |              |                |   |
| i   |     |   |  |   |              |              |                |   |
|   | _   | $\begin{matrix} 1 \\ x_{II} = b \end{matrix}$ | $     \begin{array}{c}       2 \\       x_{2I} = a     \end{array} $ |   | $x_{3I} = a$ | $x_{41} = b$ | $x_{51} = a$ 5 | _ |
| l ,   |     | В   | A,C  |   | A,C          | В            | A,C            |   |
| j   | 1   |   |  |   |              |              |                |   |
|   | 2   | S,A   | В  |   | S,C          | S,A          |                |   |
|   | _   |   |  |   |              |              |                |   |
|   | 3   | Ø   | В  |   | B //         |              |                |   |
|   | 3   | -   | Į  | ļ |              |              |                |   |
|   | 4   | Ø   |  |   |              |              |                |   |
|   | 4   |   |  |   |              |              |                |   |
|   | , 5 |   |  |   |              |              |                |   |
|   |     | $V_{I5}$                                      |  |   |              |              |                | ] |

#### **Application of CYK Algorithm**

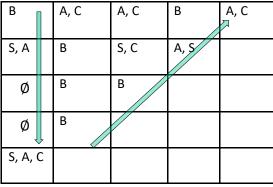
•  $S \rightarrow AB \mid BC$ 

$$A \rightarrow BA \mid a$$

$$A \rightarrow BA|a \qquad B \rightarrow CC|b \qquad C \rightarrow AB|a$$

$$C \rightarrow AB|a$$

- $w_1 = baaba$
- $w_2 = aab$ ?



S is in V<sub>15</sub> therefore w is in L(G)

17

### **Application of CYK Algorithm**

 $\bullet \quad S \rightarrow AB \mid BC \qquad \quad A \rightarrow BA \mid a \qquad B \rightarrow CC \mid b \qquad C \rightarrow AB \mid a$ 

$$A \rightarrow BA \mid a$$

$$B \rightarrow CC \mid b \mid$$

$$C \rightarrow AB \mid a$$

• Parse tree for  $w_1 = baaba$ 

| В       | A, C | A, C | В    | A, C |
|---------|------|------|------|------|
| S, A    | В    | S, C | A, S |      |
| Ø       | В    | В    |      |      |
| Ø       | В    |      |      |      |
| S, A, C |      |      |      |      |

S is in V<sub>15</sub> therefore w is in L(G)