

# Parsing

CSCI 3130 Formal Languages and Automata Theory

Siu On CHAN

Chinese University of Hong Kong

Fall 2016

## Context-free versus regular

Write a CFG for the language  $(0 + 1)^*111$

## Context-free versus regular

Write a CFG for the language  $(0 + 1)^*111$

$$S \rightarrow U111$$

$$U \rightarrow 0U \mid 1U \mid \varepsilon$$

Can you do so for **every** regular language?

## Context-free versus regular

Write a CFG for the language  $(0 + 1)^*111$

$$S \rightarrow U111$$

$$U \rightarrow 0U \mid 1U \mid \varepsilon$$

Can you do so for **every** regular language?

Every regular language is context-free



## From regular to context-free

regular expression	$\Rightarrow$ CFG
$\emptyset$	grammar with no rules
$\varepsilon$	$S \rightarrow \varepsilon$
a (alphabet symbol)	$S \rightarrow a$
$E_1 + E_2$	$S \rightarrow S_1 \mid S_2$
$E_1 E_2$	$S \rightarrow S_1 S_2$
$E_1^*$	$S \rightarrow S S_1 \mid \varepsilon$

$S$  becomes the new **start variable**

## Context-free versus regular

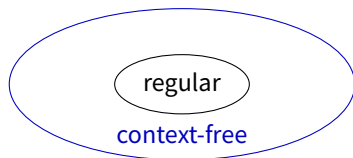
Is every context-free language regular?

## Context-free versus regular

Is every context-free language regular?

$$S \rightarrow 0S1 \quad L = \{0^n 1^n \mid n \geq 0\}$$

Is context-free but not regular



Ambiguity

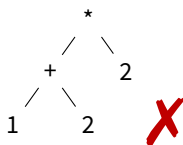


# Ambiguity

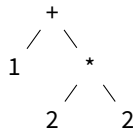
$$E \rightarrow E+E \mid E^*E \mid (E) \mid N$$

$$N \rightarrow 1N \mid 2N \mid 1 \mid 2$$

1+2\*2



= 6



= 5

A CFG is **ambiguous** if some string has more than one parse tree

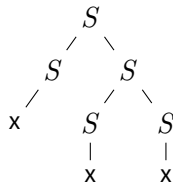
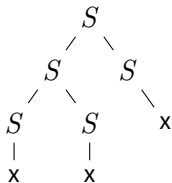
## Example

Is  $S \rightarrow SS|x$  ambiguous?

## Example

Is  $S \rightarrow SS|x$  ambiguous?

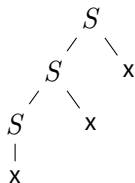
Yes, because



Two ways to derive  $xxx$

## Disambiguation

$$S \rightarrow SS|x \Rightarrow S \rightarrow Sx|x$$



Sometimes we can **rewrite the grammar** to remove ambiguity

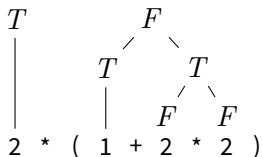
# Disambiguation

$$E \rightarrow E+E \mid E^*E \mid (E) \mid N$$

$$N \rightarrow 1N \mid 2N \mid 1 \mid 2$$

+ and \* have the same precedence!

Divide expression into **terms** and **factors**



# Disambiguation

$$E \rightarrow E+E \mid E^*E \mid (E) \mid N$$
$$N \rightarrow 1N \mid 2N \mid 1 \mid 2$$

An expression is a sum of one or more **terms**

$$E \rightarrow T \mid E+T$$

Each term is a product of one or more **factors**

$$T \rightarrow F \mid T^*F$$

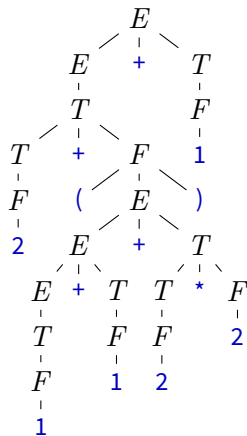
Each factor is a **parenthesized expression** or a **number**

$$F \rightarrow (E) \mid 1 \mid 2$$

## Parsing example

$$\begin{aligned} E &\rightarrow T \mid E+T \\ T &\rightarrow F \mid T*F \\ F &\rightarrow (E) \mid 1 \mid 2 \end{aligned}$$

Parse tree for  
 $2+(1+1+2*2)+1$



# Disambiguation

Disambiguation is **not always possible** because  
There exists **inherently ambiguous** languages  
There is **no general procedure** for disambiguation

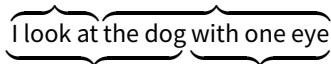


# Disambiguation

Disambiguation is **not always possible** because  
There exists **inherently ambiguous** languages  
There is **no general procedure** for disambiguation

In **programming languages**, ambiguity comes from the precedence rules,  
and we can resolve like in the example

In English, ambiguity is sometimes a problem:

  
I look at the dog with one eye

# Parsing

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

input: 0011

$$T \rightarrow S \mid \varepsilon$$

Is  $0011 \in L$ ?

If so, how to build a parse tree with a program?

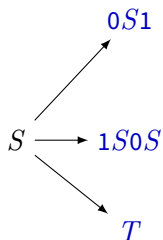
# Parsing

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

input: 0011

$$T \rightarrow S \mid \varepsilon$$

Try all derivations?



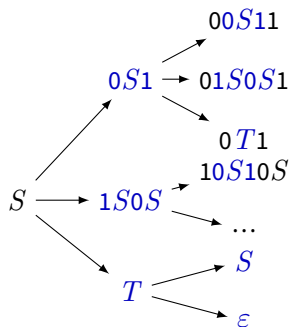
# Parsing

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

input: 0011

$$T \rightarrow S \mid \varepsilon$$

Try all derivations?

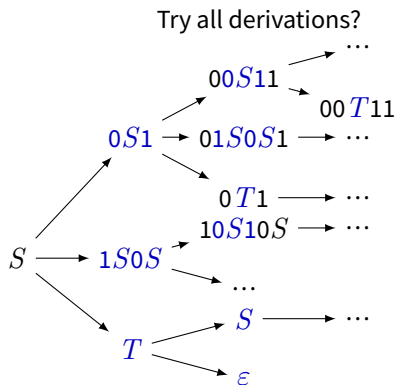


# Parsing

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

input: 0011

$$T \rightarrow S \mid \varepsilon$$

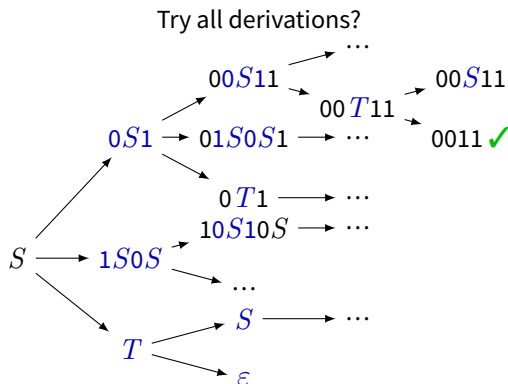


# Parsing

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

input: 0011

$$T \rightarrow S \mid \varepsilon$$



This is (part of) the **tree of all derivations**, not the parse tree

# Problems

1. Trying all derivations may take too long
2. If input is **not in the language**, parsing will never stop

Let's tackle the 2nd problem

## When to stop

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid \varepsilon$$

**Idea:** Stop when  
 $|\text{derived string}| > |\text{input}|$



## When to stop

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid \varepsilon$$

Idea: Stop when  
 $|\text{derived string}| > |\text{input}|$

Problems:

$$S \Rightarrow 0S1 \Rightarrow 0T1 \Rightarrow 01$$

Derived string may **shrink**  
because of “ $\varepsilon$ -productions”

## When to stop

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid \varepsilon$$

**Idea:** Stop when  
 $|\text{derived string}| > |\text{input}|$

Problems:

$$S \Rightarrow 0S1 \Rightarrow 0T1 \Rightarrow 01$$

Derived string may **shrink**  
because of “ $\varepsilon$ -productions”

$$S \Rightarrow T \Rightarrow S \Rightarrow T \Rightarrow \dots$$

Derivation may **loop** because  
of “unit productions”

Remove  $\varepsilon$  and unit productions

## Removing $\varepsilon$ -productions

Goal: remove all  $A \rightarrow \varepsilon$  rules for every non-start variable  $A$

If  $S$  is the start variable and the rule  $S \rightarrow \varepsilon$  exists

Add a new start variable  $T$

Add the rule  $T \rightarrow S$

$$S \rightarrow ACD$$

$$A \rightarrow a$$

$$B \rightarrow \varepsilon$$

$$C \rightarrow ED \mid \varepsilon$$

$$D \rightarrow BC \mid b$$

$$E \rightarrow b$$

For every rule  $A \rightarrow \varepsilon$  where  $A$  is not the (new) start variable

1. Remove the rule  $A \rightarrow \varepsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

## Removing $\varepsilon$ -productions

Goal: remove all  $A \rightarrow \varepsilon$  rules for every non-start variable  $A$

If  $S$  is the start variable and the rule  $S \rightarrow \varepsilon$  exists

Add a new start variable  $T$

Add the rule  $T \rightarrow S$

For every rule  $A \rightarrow \varepsilon$  where  $A$  is not the (new) start variable

1. Remove the rule  $A \rightarrow \varepsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

$S \rightarrow ACD$   
 $A \rightarrow a$   
 ~~$B \rightarrow \varepsilon$~~   
 $C \rightarrow ED \mid \varepsilon$   
 $D \rightarrow BC \mid b$   
 $E \rightarrow b$

$D \rightarrow C$

Removing  $B \rightarrow \varepsilon$

## Removing $\varepsilon$ -productions

Goal: remove all  $A \rightarrow \varepsilon$  rules for every non-start variable  $A$

If  $S$  is the start variable and the rule  $S \rightarrow \varepsilon$  exists

Add a new start variable  $T$

Add the rule  $T \rightarrow S$

For every rule  $A \rightarrow \varepsilon$  where  $A$  is not the (new) start variable

1. Remove the rule  $A \rightarrow \varepsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

$S \rightarrow ACD$        $D \rightarrow C$   
 $A \rightarrow a$        $S \rightarrow AD$   
 ~~$B \rightarrow \varepsilon$~~   
 $C \rightarrow ED \mid \varepsilon$   
 $D \rightarrow BC \mid b$   
 $E \rightarrow b$

Removing  $C \rightarrow \varepsilon$

## Removing $\epsilon$ -productions

Goal: remove all  $A \rightarrow \epsilon$  rules for every non-start variable  $A$

If  $S$  is the start variable and the rule  $S \rightarrow \epsilon$  exists

Add a new start variable  $T$

Add the rule  $T \rightarrow S$

$S \rightarrow ACD$   
 $A \rightarrow a$   
 ~~$B \rightarrow \epsilon$~~   
 $C \rightarrow ED \mid \epsilon$   
 $D \rightarrow BC \mid b$   
 $E \rightarrow b$

For every rule  $A \rightarrow \epsilon$  where  $A$  is not the (new) start variable

1. Remove the rule  $A \rightarrow \epsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

$D \rightarrow C$   
 $S \rightarrow AD$   
 $D \rightarrow \epsilon$

Removing  $C \rightarrow \epsilon$

## Removing $\epsilon$ -productions

Goal: remove all  $A \rightarrow \epsilon$  rules for every non-start variable  $A$

If  $S$  is the start variable and the rule  $S \rightarrow \epsilon$  exists

Add a new start variable  $T$

Add the rule  $T \rightarrow S$

For every rule  $A \rightarrow \epsilon$  where  $A$  is not the (new) start variable

1. Remove the rule  $A \rightarrow \epsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

$S \rightarrow ACD$

$A \rightarrow a$

~~$B \rightarrow \epsilon$~~

$C \rightarrow ED \mid \cancel{\epsilon}$

$D \rightarrow BC \mid b$

$E \rightarrow b$

$D \rightarrow C$

$S \rightarrow AD$

~~$D \rightarrow \epsilon$~~

$C \rightarrow E$

Removing  $D \rightarrow \epsilon$

## Removing $\varepsilon$ -productions

Goal: remove all  $A \rightarrow \varepsilon$  rules for every non-start variable  $A$

If  $S$  is the start variable and the rule  $S \rightarrow \varepsilon$  exists

Add a new start variable  $T$

Add the rule  $T \rightarrow S$

$S \rightarrow ACD$   
 $A \rightarrow a$   
 ~~$B \rightarrow \varepsilon$~~   
 $C \rightarrow ED \mid \varepsilon$   
 $D \rightarrow BC \mid b$   
 $E \rightarrow b$

For every rule  $A \rightarrow \varepsilon$  where  $A$  is not the (new) start variable

1. Remove the rule  $A \rightarrow \varepsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

$D \rightarrow C$   
 $S \rightarrow AD$   
 ~~$D \rightarrow \varepsilon$~~   
 $C \rightarrow E$   
 $S \rightarrow A$

Removing  $D \rightarrow \varepsilon$



## Eliminating $\varepsilon$ -productions

For every  $A \rightarrow \varepsilon$  rule where  $A$  is not the start variable

1. Remove the rule  $A \rightarrow \varepsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

Do 2. every time  $A$  appears

$$\begin{aligned} B \rightarrow \alpha A \beta A \gamma &\text{ yields} \\ B \rightarrow \alpha \beta A \gamma \quad B \rightarrow \alpha A \beta \gamma & \\ B \rightarrow \alpha \beta \gamma & \end{aligned}$$

## Eliminating $\varepsilon$ -productions

For every  $A \rightarrow \varepsilon$  rule where  $A$  is not the start variable

1. Remove the rule  $A \rightarrow \varepsilon$
2. If you see  $B \rightarrow \alpha A \beta$   
Add a new rule  $B \rightarrow \alpha \beta$

Do 2. every time  $A$  appears

$$\begin{array}{l} B \rightarrow \alpha A \beta A \gamma \text{ yields} \\ B \rightarrow \alpha \beta A \gamma \quad B \rightarrow \alpha A \beta \gamma \\ \quad B \rightarrow \alpha \beta \gamma \end{array}$$

$B \rightarrow A$  becomes  $B \rightarrow \varepsilon$

If  $B \rightarrow \varepsilon$  was removed earlier,  
don't add it back

## Eliminating unit productions

A **unit production** is a production of the form

$$A \rightarrow B$$

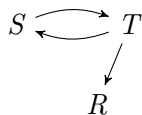
Grammar:

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid R \mid \varepsilon$$

$$R \rightarrow 0SR$$

Unit production graph:



## Removing unit productions

① If there is a **cycle** of unit productions

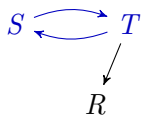
$$A \rightarrow B \rightarrow \dots \rightarrow C \rightarrow A$$

**delete it** and **replace** everything with  $A$

$$S \rightarrow 0S1 \mid 1S0S \mid T$$

$$T \rightarrow S \mid R \mid \varepsilon$$

$$R \rightarrow 0SR$$



## Removing unit productions

① If there is a **cycle** of unit productions

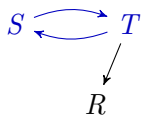
$$A \rightarrow B \rightarrow \dots \rightarrow C \rightarrow A$$

**delete it** and **replace** everything with  $A$

$$S \rightarrow 0S1 \mid 1S0S \mid \cancel{T}$$

$$\cancel{T} \rightarrow \cancel{S} \mid R \mid \varepsilon$$

$$R \rightarrow 0SR$$



$$S \rightarrow 0S1 \mid 1S0S$$

$$S \rightarrow R \mid \varepsilon$$

$$R \rightarrow 0SR$$

Replace  $T$  by  $S$

## Removal of unit productions

② replace any chain

$$A \rightarrow B \rightarrow \dots \rightarrow C \rightarrow \alpha$$

$$\text{by } A \rightarrow \alpha, \quad B \rightarrow \alpha, \quad \dots, \quad C \rightarrow \alpha$$

$$S \rightarrow 0S1 \mid 1S0S$$

$$\mid R \mid \varepsilon$$

$$R \rightarrow 0SR$$

$S$



$R$

## Removal of unit productions

② replace any chain

$$A \rightarrow B \rightarrow \dots \rightarrow C \rightarrow \alpha$$

$$\text{by } A \rightarrow \alpha, \quad B \rightarrow \alpha, \quad \dots, \quad C \rightarrow \alpha$$

$$S \rightarrow 0S1 \mid 1S0S$$

$$\mid R \mid \varepsilon$$

$$R \rightarrow 0SR$$

$$S$$

$$R$$

$$S \rightarrow 0S1 \mid 1S0S$$

$$\mid 0SR \mid \varepsilon$$

$$R \rightarrow 0SR$$

Replace  $S \rightarrow R \rightarrow 0SR$  by  $S \rightarrow 0SR, \quad R \rightarrow 0SR$

# Recap

## Problems:

1. Trying all derivations may take too long
2. If input is **not in the language**, parsing will never stop ✓

## Solution to problem 2:

1. Eliminate  $\varepsilon$  productions
2. Eliminate unit productions

Try all possible derivations but stop parsing when  
 $|\text{derived string}| > |\text{input}|$



## Example

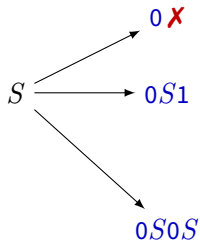
$$S \rightarrow 0S1 \mid 0S0S \mid T$$

$$T \rightarrow S \mid 0$$

$\implies$

$$S \rightarrow 0S1 \mid 0S0S \mid 0$$

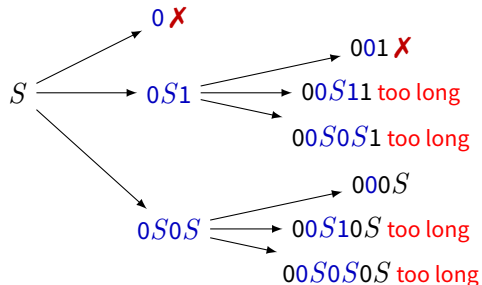
input: 0011



## Example

$$\begin{array}{l} S \rightarrow 0S1 \mid 0S0S \mid T \\ T \rightarrow S \mid 0 \end{array} \quad \Longrightarrow \quad S \rightarrow 0S1 \mid 0S0S \mid 0$$

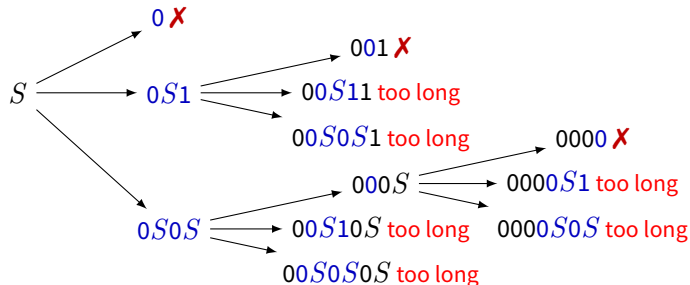
input: 0011



## Example

$$\begin{array}{l} S \rightarrow 0S1 \mid 0S0S \mid T \\ T \rightarrow S \mid 0 \end{array} \quad \Longrightarrow \quad S \rightarrow 0S1 \mid 0S0S \mid 0$$

input: 0011



Conclusion:  $0011 \notin L$

# Problems

1. Trying all derivations **may take too long**
2. If input is not in the language, parsing will never stop

# Preparations

A faster way to parse:

Cocke-Younger-Kasami algorithm

To use it we must preprocess the CFG:

Eliminate  $\epsilon$  productions

Eliminate unit productions

Convert CFG to **Chomsky Normal Form**

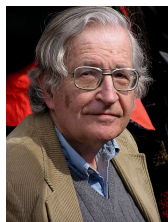
# Chomsky Normal Form

A CFG is in **Chomsky Normal Form** if every production has the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a$$

where neither  $B$  nor  $C$  is the start variable

but we also allow  $S \rightarrow \varepsilon$  for start variable  $S$



Noam Chomsky

Convert to Chomsky Normal Form:

$$\begin{array}{l} A \rightarrow BcDE \quad \Longrightarrow \quad A \rightarrow BCDE \quad \Longrightarrow \quad A \rightarrow BX \\ \text{replace } C \rightarrow c \quad \text{break up} \quad X \rightarrow CY \\ \text{terminals} \quad \text{sequences} \quad Y \rightarrow DE \\ \text{with new} \quad \text{with new} \quad C \rightarrow c \\ \text{variables} \quad \text{variables} \end{array}$$

# Cocke-Younger-Kasami algorithm

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

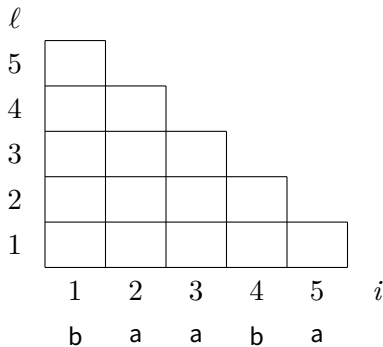
$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$

**Input:**  $x = \text{baaba}$

let

$$x[i, \ell] = x_i x_{i+1} \dots x_{i+\ell-1}$$



For every substring  $x[i, \ell]$ , remember all variables  $R$  that derive  $x[i, \ell]$   
Store in a table  $T[i, \ell]$

# Cocke-Younger-Kasami algorithm

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

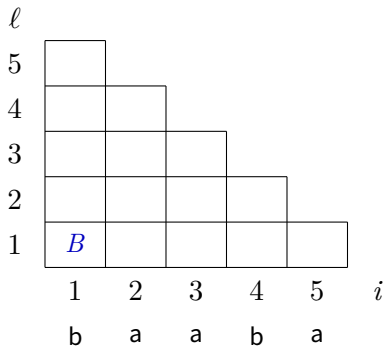
$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$

Input:  $x = \text{baaba}$

let

$$x[i, \ell] = x_i x_{i+1} \dots x_{i+\ell-1}$$



For every substring  $x[i, \ell]$ , remember all variables  $R$  that derive  $x[i, \ell]$   
Store in a table  $T[i, \ell]$



# Cocke-Younger-Kasami algorithm

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

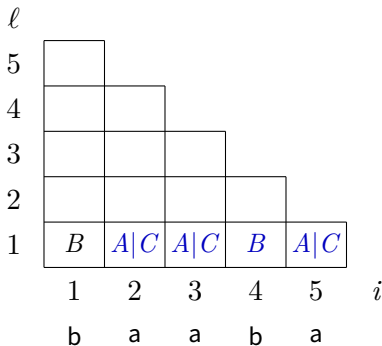
$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$

Input:  $x = \text{baaba}$

let

$$x[i, \ell] = x_i x_{i+1} \dots x_{i+\ell-1}$$



For every substring  $x[i, \ell]$ , remember all variables  $R$  that derive  $x[i, \ell]$   
Store in a table  $T[i, \ell]$

# Cocke-Younger-Kasami algorithm

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$

Input:  $x = \text{baaba}$

let

$$x[i, \ell] = x_i x_{i+1} \dots x_{i+\ell-1}$$

$\ell$						
5						
4						
3						
2	$S A$					
1	$B$	$A C$	$A C$	$B$	$A C$	
	1	2	3	4	5	$i$
	b	a	a	b	a	

For every substring  $x[i, \ell]$ , remember all variables  $R$  that derive  $x[i, \ell]$   
Store in a table  $T[i, \ell]$

# Cocke-Younger-Kasami algorithm

$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

**Input:**  $x = \text{baaba}$

let

$x[i, \ell] = x_i x_{i+1} \dots x_{i+\ell-1}$

$\ell$						
5						
4						
3						
2	S A	B	S C	S A		
1	B	A C	A C	B	A C	
	1	2	3	4	5	$i$
	b	a	a	b	a	

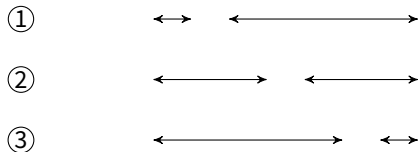
For every substring  $x[i, \ell]$ , remember all variables  $R$  that derive  $x[i, \ell]$   
 Store in a table  $T[i, \ell]$

## Computing $T[i, \ell]$ for $\ell \geq 2$

To compute  $T[2, 4]$

Try all possible ways to split  $x[2, 4]$  into two substrings

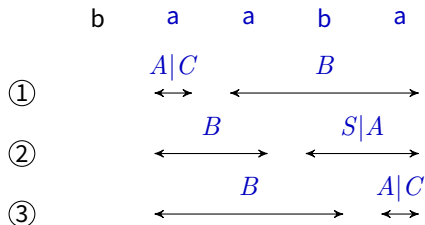
b    a    a    b    a



## Computing $T[i, \ell]$ for $\ell \geq 2$

To compute  $T[2, 4]$

Try all possible ways to split  $x[2, 4]$  into two substrings

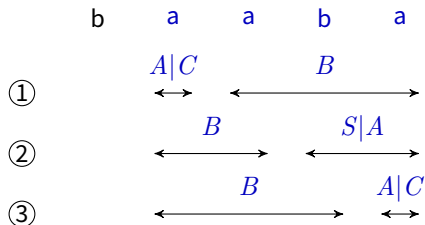


Look up entries regarding **shorter substrings** previously computed

## Computing $T[i, \ell]$ for $\ell \geq 2$

To compute  $T[2, 4]$

Try all possible ways to split  $x[2, 4]$  into two substrings



Look up entries regarding **shorter substrings** previously computed

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$

$$T[2, 4] = S|A|C$$

# Cocke-Younger-Kasami algorithm

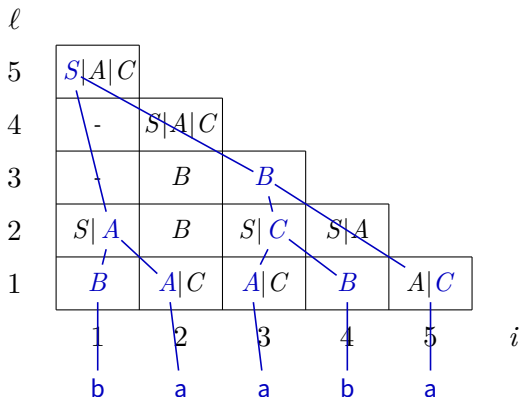
$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

Input:  $x = \text{baaba}$



Get **parse tree** by tracing back derivations