# Parsing
## CSCI 3130 Formal Languages and Automata Theory

### Siu On CHAN

Chinese University of Hong Kong

### Fall 2017

# Context-free versus regular

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

# Context-free versus regular

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

$$S \to U111$$
$$U \to 0\,U \mid 1\,U \mid \varepsilon$$

Can you do so for every regular language?

# Context-free versus regular

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

$$S \rightarrow U111$$
$$U \rightarrow 0\,U \mid 1\,U \mid \varepsilon$$

Can you do so for every regular language?

Every regular language is context-free

# From regular to context-free

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

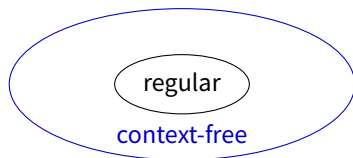$S$ becomes the new start variable

# Context-free versus regular

Is every context-free language regular?

# Context-free versus regular

Is every context-free language regular?

$$S \to 0S1 \qquad L = \{0^n 1^n \mid n \geqslant 0\}$$

Is context-free but not regular
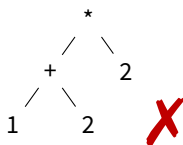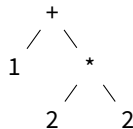
Ambiguity

# Ambiguity

$$E \rightarrow E\text{+}E \mid E^\star E \mid (E) \mid N$$
$$N \rightarrow 1 \mid 2$$

1+2*2



$= 6$          $= 5$

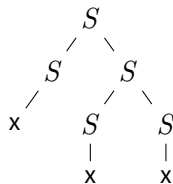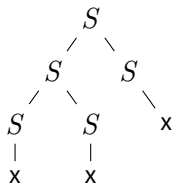A CFG is ambiguous if some string has more than one parse tree

# Example

Is $\boxed{S \rightarrow SS \mid \mathsf{x}}$ ambiguous?

# Example

Is $\boxed{S \rightarrow SS \mid \text{x}}$ ambiguous?

Yes, because



Two ways to derive xxx

# Disambiguation

$$S \rightarrow SS \mid \mathsf{x} \quad \Rightarrow \quad S \rightarrow S\mathsf{x} \mid \mathsf{x}$$



Sometimes we can rewrite the grammar to remove ambiguity

# Disambiguation

$$E \rightarrow E\texttt{+}E \mid E\texttt{*}E \mid (E) \mid N$$
$$N \rightarrow \texttt{1} \mid \texttt{2}$$

+ and * have the same precedence!
Divide expression into terms and factors

# Disambiguation

$$E \to E\text{+}E \mid E^{\star}E \mid (E) \mid N$$
$$N \to \text{1} \mid \text{2}$$

An expression is a sum of one or more terms
$$E \to T \mid E\text{+}T$$
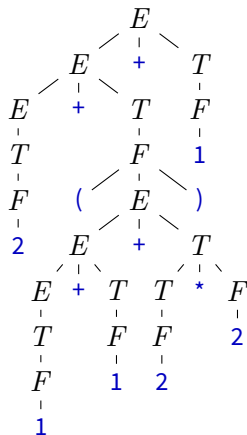Each term is a product of one or more factors
$$T \to F \mid T^{\star}F$$
Each factor is a parenthesized expression or a number
$$F \to (E) \mid \text{1} \mid \text{2}$$

# Parsing example

$$E \rightarrow T \mid E\texttt{+}T$$
$$T \rightarrow F \mid T\texttt{*}F$$
$$F \rightarrow (E) \mid \texttt{1} \mid \texttt{2}$$

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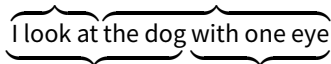
Disambiguation is not always possible because
There exists inherently ambiguous languages
There is no general procedure for disambiguation

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

$$S \to 0S1 \mid 1S0S \mid T$$
$$T \to S \mid \varepsilon$$

input: 0011

Is $0011 \in L$?
If so, how to build a parse tree with a program?

# Parsing

$$S \to 0S1 \mid 1S0S \mid T$$
$$T \to S \mid \varepsilon$$

input: 0011

Try all derivations?

# Parsing

$$S \to 0S1 \mid 1S0S \mid T$$
$$T \to S \mid \varepsilon$$

input: 0011

Try all derivations?

# Parsing

$$S \to 0S1 \mid 1S0S \mid T$$
$$T \to S \mid \varepsilon$$



input: 0011

Try all derivations?

# Parsing

$$S \to 0S1 \mid 1S0S \mid T$$
$$T \to S \mid \varepsilon$$

input: 0011



Try all derivations?

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

$$S \to \texttt{0}S\texttt{1} \mid \texttt{1}S\texttt{0}S \mid T$$
$$T \to \texttt{S} \mid \varepsilon$$

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

# When to stop

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

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

Problems:

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

Derived string may shrink
because of "$\varepsilon$-productions"

# When to stop

$$S \rightarrow 0S1 \mid 1S0S \mid T$$
$$T \rightarrow \mathsf{s} \mid \varepsilon$$

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

Problems:

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

Derived string may shrink
because of "$\varepsilon$-productions"

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

Derviation may loop because
of "unit productions"

Remove $\varepsilon$ and unit productions

# Removing $\varepsilon$-productions

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

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

Add a new start variable $T$
Add the rule $T \to S$

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

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

$$S \to ACD$$
$$A \to \mathsf{a}$$
$$B \to \varepsilon$$
$$C \to ED \mid \varepsilon$$
$$D \to BC \mid \mathsf{b}$$
$$E \to \mathsf{b}$$

# 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

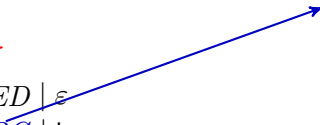Add a new start variable $T$
Add the rule $T \rightarrow S$

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 $B \rightarrow \alpha \beta$

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

$$D \rightarrow C$$

Removing $B \rightarrow \varepsilon$

# Removing $\varepsilon$-productions

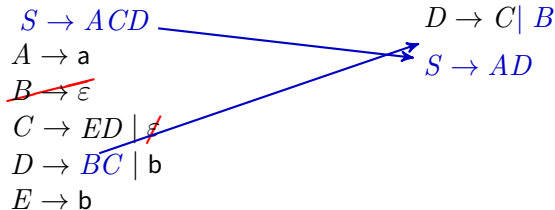Goal: remove all $A \to \varepsilon$ rules for every non-start variable $A$

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

Add a new start variable $T$
Add the rule $T \to S$

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

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

$S \to ACD$
$A \to \text{a}$
$B \to \varepsilon$
$C \to ED \mid \varepsilon$
$D \to BC \mid \text{b}$
$E \to \text{b}$

$D \to C \mid B$
$S \to AD$

Removing $C \to \varepsilon$

# Removing $\varepsilon$-productions

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

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

Add a new start variable $T$
Add the rule $T \to S$

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

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

$$S \to ACD$$
$$A \to \text{a}$$
$$B \to \varepsilon$$
$$C \to ED \mid \varepsilon$$
$$D \to BC \mid \text{b}$$
$$E \to \text{b}$$

$$D \to C \mid B$$
$$S \not\to AD$$
$$D \to \varepsilon$$

Removing $C \to \varepsilon$

# Removing $\varepsilon$-productions

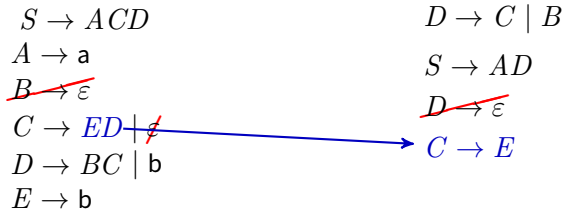Goal: remove all $A \to \varepsilon$ rules for every non-start variable $A$

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

Add a new start variable $T$
Add the rule $T \to S$

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

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

$S \to ACD$
$A \to \mathsf{a}$
$B \to \varepsilon$
$C \to ED \mid E$
$D \to BC \mid \mathsf{b}$
$E \to \mathsf{b}$

$D \to C \mid B$
$S \to AD$
$D \to \varepsilon$
$C \to E$

Removing $D \to \varepsilon$

# Removing $\varepsilon$-productions

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

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

Add a new start variable $T$
Add the rule $T \to S$

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

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

$$S \to ACD$$
$$A \to \text{a}$$
$$B \to \varepsilon$$
$$C \to ED \mid \not\varepsilon$$
$$D \to BC \mid \text{b}$$
$$E \to \text{b}$$

$$D \to C \mid B$$
$$S \to AD$$
$$D \to \varepsilon$$
$$C \to E$$
$$S \to A$$

Removing $D \to \varepsilon$

# Eliminating $\varepsilon$-productions

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

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

Do 2. every time $A$ appears

$B \to \alpha A \beta A \gamma$ yields
$$B \to \alpha \beta A \gamma \quad B \to \alpha A \beta \gamma$$
$$B \to \alpha \beta \gamma$$

# Eliminating $\varepsilon$-productions

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

1. Remove the rule $A \to \varepsilon$

2. If you see $B \to \alpha A \beta$
   Add a new rule $B \to \alpha \beta$

Do 2. every time $A$ appears

$B \to \alpha A \beta A \gamma$ yields
$B \to \alpha \beta A \gamma \quad B \to \alpha A \beta \gamma$
$B \to \alpha \beta \gamma$

$B \to A$ becomes $B \to \varepsilon$

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

# Eliminating unit productions

A unit production is a production of the form
$$A \to B$$

Grammar:
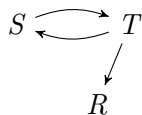
$S \to \mathtt{0}S\mathtt{1} \mid \mathtt{1}S\mathtt{0}S \mid T$
$T \to S \mid R \mid \varepsilon$
$R \to \mathtt{0}SR$

Unit production graph:

# Removing unit productions

① If there is a cycle of unit productions

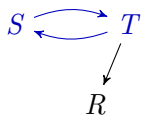$$A \rightarrow B \rightarrow \cdots \rightarrow C \rightarrow A$$

delete it and replace everything with $A$

$S \rightarrow \texttt{0}S\texttt{1} \mid \texttt{1}S\texttt{0}S \mid T$
$T \rightarrow S \mid R \mid \varepsilon$
$R \rightarrow \texttt{0}SR$

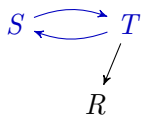# Removing unit productions

①If there is a cycle of unit productions

$$A \to B \to \cdots \to C \to A$$

delete it and replace everything with $A$

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

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

$R \to 0SR$

$S \rightleftarrows T$

$T$
$\downarrow$
$R$

$S \to 0S1 \mid 1S0S$

$S \to R \mid \varepsilon$

$R \to 0SR$

Replace $T$ by $S$

# Removal of unit productions

② replace any chain

$$A \to B \to \cdots \to C \to \alpha$$

$$\text{by} \quad A \to \alpha, \quad B \to \alpha, \quad \ldots, \quad C \to \alpha$$

$S \to \mathtt{0}S\mathtt{1} \mid \mathtt{1}S\mathtt{0}S$

$\quad \mid R \mid \varepsilon$

$R \to \mathtt{0}SR$

$S$

$\downarrow$

$R$

# Removal of unit productions

② replace any chain

$$A \to B \to \cdots \to C \to \alpha$$

$$\text{by} \quad A \to \alpha, \quad B \to \alpha, \quad \ldots, \quad C \to \alpha$$

$S \to 0S1 \mid 1S0S$

$\quad \mid R \mid \varepsilon$

$R \to 0SR$

$S$

$\downarrow$

$R$

$S \to 0S1 \mid 1S0S$

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

$R \to 0SR$

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

# 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
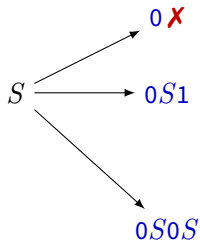$|\text{derived string}| > |\text{input}|$

# Example

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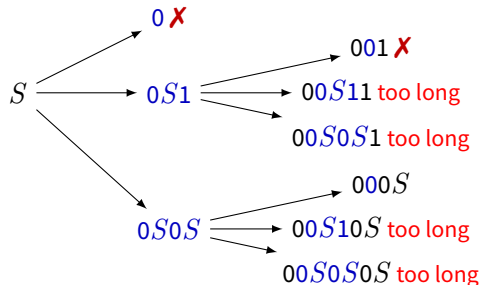$T \to S \mid 0$

$\implies$

$S \to 0S1 \mid 0S0S \mid 0$

input: 0011

# Example

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

$$\implies \qquad S \to 0S1 \mid 0S0S \mid 0$$

input: 0011

# Example

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

$$\implies \qquad S \to 0S1 \mid 0S0S \mid 0$$

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 perprocess the CFG:

Eliminate $\varepsilon$ productions
Eliminate unit productions
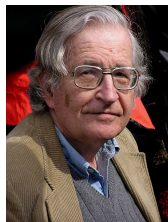Convert CFG to Chomsky Normal Form

# Chomsky Normal Form

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

$A \to BC$  or  $A \to$ a
where neither $B$ nor $C$ is the start variable

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



Noam Chomsky

Convert to Chomsky Normal Form:

$$A \to B\text{c}DE \quad \implies \quad A \to BCDE \quad \implies \quad A \to BX$$

| $A \to B\text{c}DE$ | $\implies$ | $A \to BCDE$ | $\implies$ | $A \to BX$ |
| | replace | $C \to$ c | break up | $X \to CY$ |
| | terminals | | sequences | $Y \to DE$ |
| | with new | | with new | $C \to$ c |
| | variables | | variables | |

# Cocke–Younger–Kasami algorithm

$$S \to AB \mid BC$$
$$A \to BA \mid \mathsf{a}$$
$$B \to CC \mid \mathsf{b}$$
$$C \to AB \mid \mathsf{a}$$

Input: $x = \mathsf{baaba}$

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



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

$$S \to AB \mid BC$$
$$A \to BA \mid \mathsf{a}$$
$$B \to CC \mid \mathsf{b}$$
$$C \to AB \mid \mathsf{a}$$

Input: $x = \mathsf{baaba}$

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

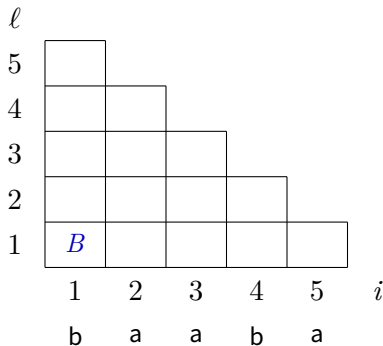$$S \rightarrow AB \mid BC$$
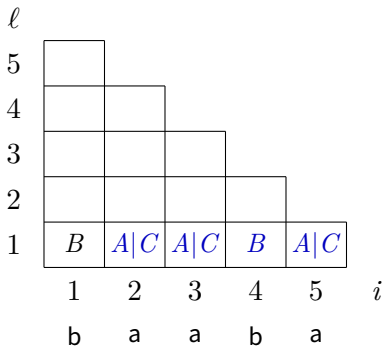$$A \rightarrow BA \mid \mathsf{a}$$
$$B \rightarrow CC \mid \mathsf{b}$$
$$C \rightarrow AB \mid \mathsf{a}$$

Input: $x = \mathsf{baaba}$

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



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

$$S \to AB \mid BC$$
$$A \to BA \mid \mathsf{a}$$
$$B \to CC \mid \mathsf{b}$$
$$C \to AB \mid \mathsf{a}$$

Input: $x = $ baaba

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

| $\ell$ | 1 | 2 | 3 | 4 | 5 | $i$ |
|---|---|---|---|---|---|---|
| 5 | | | | | | |
| 4 | | | | | | |
| 3 | | | | | | |
| 2 | $S\mid A$ | | | | | |
| 1 | $B$ | $A\mid C$ | $A\mid C$ | $B$ | $A\mid C$ | |
| | 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

$$S \to AB \mid BC$$
$$A \to BA \mid \mathsf{a}$$
$$B \to CC \mid \mathsf{b}$$
$$C \to AB \mid \mathsf{a}$$

Input: $x = \mathsf{baaba}$
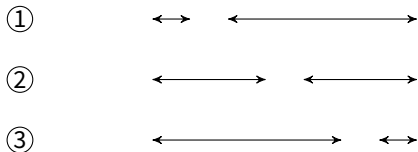
let
$$x[i, \ell] = x_i x_{i+1} \ldots x_{i+\ell-1}$$



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

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

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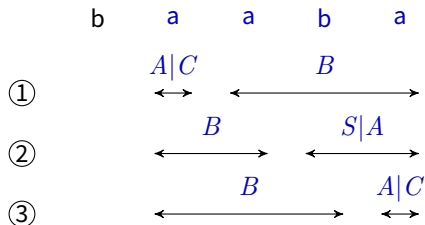
Try all possible ways to split $x[2, 4]$ into two substrings

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

<center>To compute $T[2, 4]$</center>

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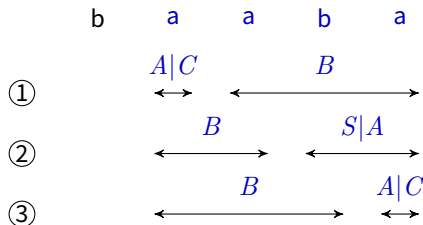


Look up entries regarding shorter substrings previously computed

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

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

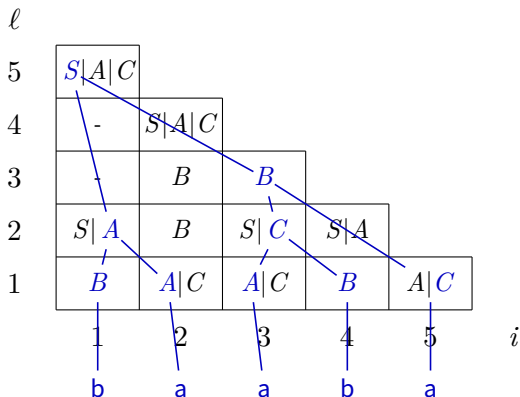$T[2, 4] = S \mid A \mid C$

# Cocke–Younger–Kasami algorithm

$S \rightarrow AB \mid BC$
$A \rightarrow BA \mid \text{a}$
$B \rightarrow CC \mid \text{b}$
$C \rightarrow AB \mid \text{a}$

Input: $x = \text{baaba}$



Get parse tree by tracing back derivations