

Chomsky Normal Form

CSE 211 (Theory of Computation)

Atif Hasan Rahman

Assistant Professor
Department of Computer Science and Engineering
Bangladesh University of Engineering & Technology

Adapted from slides by
Dr. Muhammad Masroor Ali

Chomsky Normal Form

CFG দেওয়া
 String \Rightarrow } একটি specific form এ থাকা CFG (CNF)
 grammar এ করে বিনা।

- Every nonempty CFL without ϵ has a grammar G in which all productions are in one of two simple forms, either:
 - $A \rightarrow BC$, where A , B , and C , are each variables, or
 - $A \rightarrow a$, where A is a variable and a is a terminal.
- Further, G has no useless symbols.
- Such a grammar is said to be in **Chomsky Normal Form**, or **CNF**.

→ production rule গুলো কিছু নিয়ম ফল করবে:

variable

(i) $A \rightarrow e$ ✓ থারেন্সি: $s \rightarrow e$ OK ?

(ii) $A \rightarrow BC$ ✓

$A \cancel{\rightarrow} B$, $A \cancel{\rightarrow} BCD$ → একটি থারেন্সি।

variable থারে ২টি

variable থাইটি

terminal থারে ১টি

$A \rightarrow a$ ✓

$A \rightarrow aa$ ✗

$A \rightarrow nAX$

Noam Chomsky

- “the father of modern linguistics” - wiki
- Linguist, philosopher, cognitive scientist, historian, social critic, and political activist
- Developed the theory of transformational grammar
- Author of many books and articles
 - Anti-war essay “The Responsibility of Intellectuals”
 - Criticism of media in “Manufacturing Consent”



Testing Membership in a CFL

CNF, $w \leftarrow$ string

- There is an efficient technique based on the idea of “dynamic programming”
- The algorithm is known as the **CYK Algorithm**.
 - Cocke-Younger-Kasami algorithm
- It starts with a CNF grammar $G = (V, \Sigma, R, S)$ for a language L .
- The input to the algorithm is a string $w = a_1 a_2 \dots a_n$ in Σ^* . In $O(n^3)$ time, the algorithm constructs a table that tells whether w is in L .

Testing Membership in a CFL

					X_{15}
				X_{14}	X_{25}
			X_{13}	X_{24}	X_{35}
		X_{12}	X_{23}	X_{34}	X_{45}
	X_{11}	X_{22}	X_{33}	X_{44}	X_{55}
<hr/>					
	a_1	a_2	a_3	a_4	a_5

Figure 7.12: The table constructed by the CYK algorithm

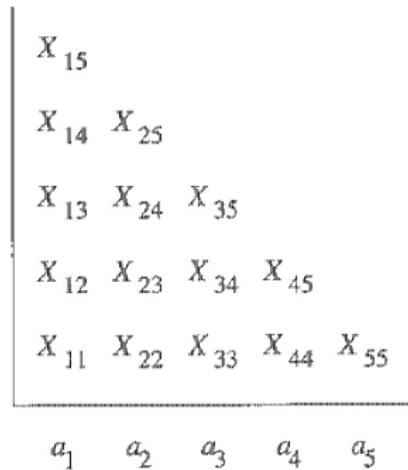
Testing Membership in a CFL

X_{15}
$X_{14} \quad X_{25}$
$X_{13} \quad X_{24} \quad X_{35}$
$X_{12} \quad X_{23} \quad X_{34} \quad X_{45}$
$X_{11} \quad X_{22} \quad X_{33} \quad X_{44} \quad X_{55}$

$a_1 \quad a_2 \quad a_3 \quad a_4 \quad a_5$

- We construct a triangular table.

Testing Membership in a CFL



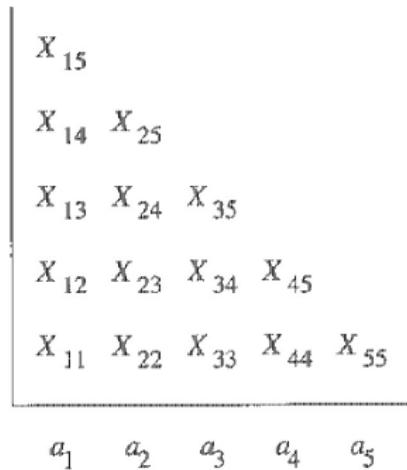
- The horizontal axis corresponds to the positions of the string $w = a_1 a_2 \dots a_n$.

Testing Membership in a CFL

	X_{15}				
	X_{14}	X_{25}			
	X_{13}	X_{24}	X_{35}		
	X_{12}	X_{23}	X_{34}	X_{45}	
	X_{11}	X_{22}	X_{33}	X_{44}	X_{55}
	a_1	a_2	a_3	a_4	a_5

- The table entry X_{ij} is the set of variables A such that $A \xrightarrow{*} a_i a_{i+1} \dots a_j$.

Testing Membership in a CFL



- We are interested in whether S is in the set X_{1n} , because that is the same as saying $S \xrightarrow{*} w$, i.e., w is in L .

Testing Membership in a CFL

- X_{ij} is the set of variables A such that $A \rightarrow a_i$ is a production of G .
- In order for A to be in X_{ij} , we must find variables B and C , and integer k such that:
 - ✓ $i \leq k < j$.
 - B is in X_{ik} .
 - C is in $X_{k+1,j}$.
 - $A \rightarrow BC$ is a production of G .

Example

The following are the productions of a CNF grammar G

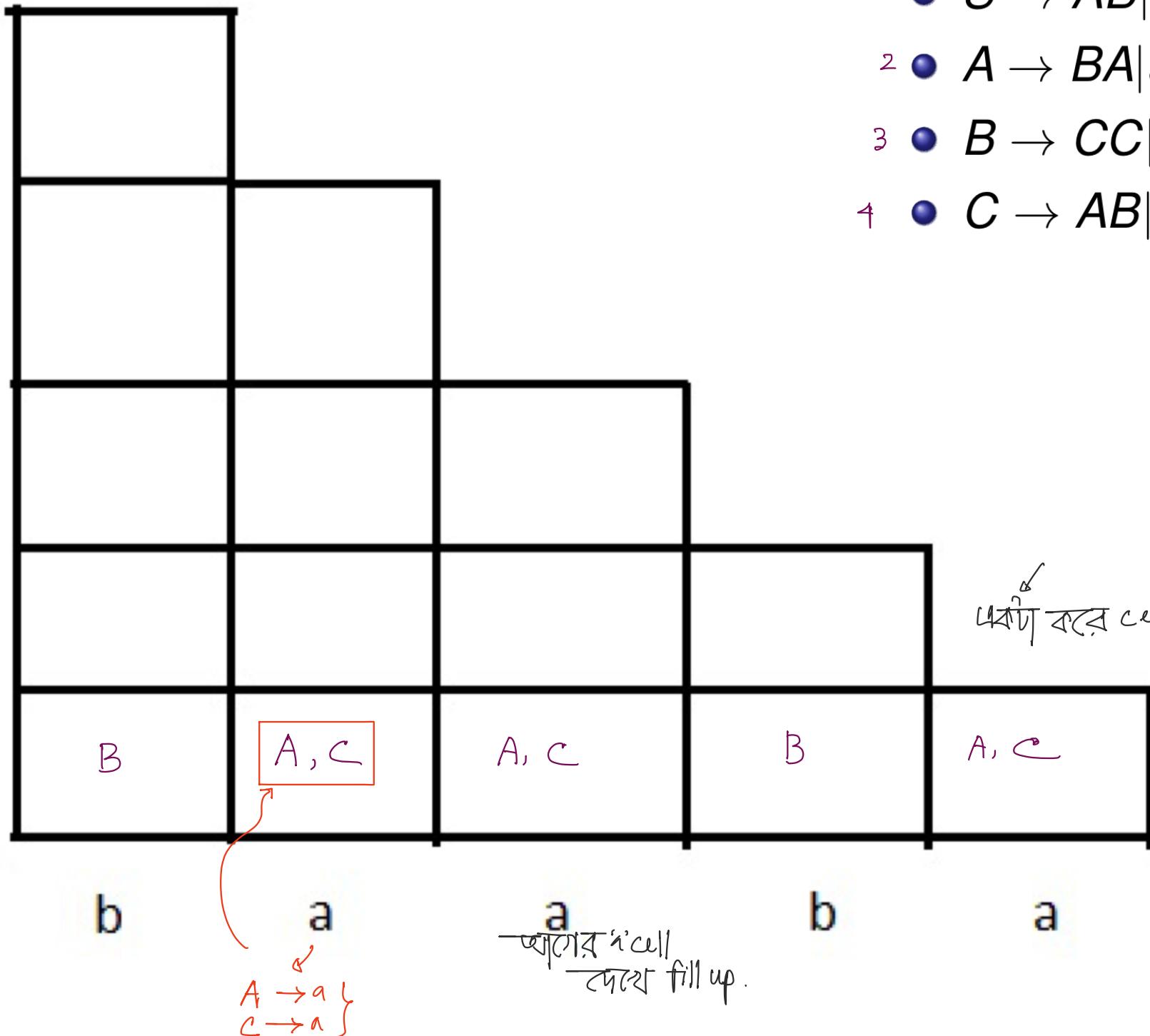
- $S \rightarrow AB|BC$
 - $A \rightarrow BA|a$
 - $B \rightarrow CC|b$
 - $C \rightarrow AB|a$
- এগুলো rule এ রয়েছে কিন্তু variable বা একটা terminal আছে।
CNF-ত আছে।

We shall test for membership in $L(G)$ the string $baaba$.

all possible derivation try কর
ক্ষমতা।

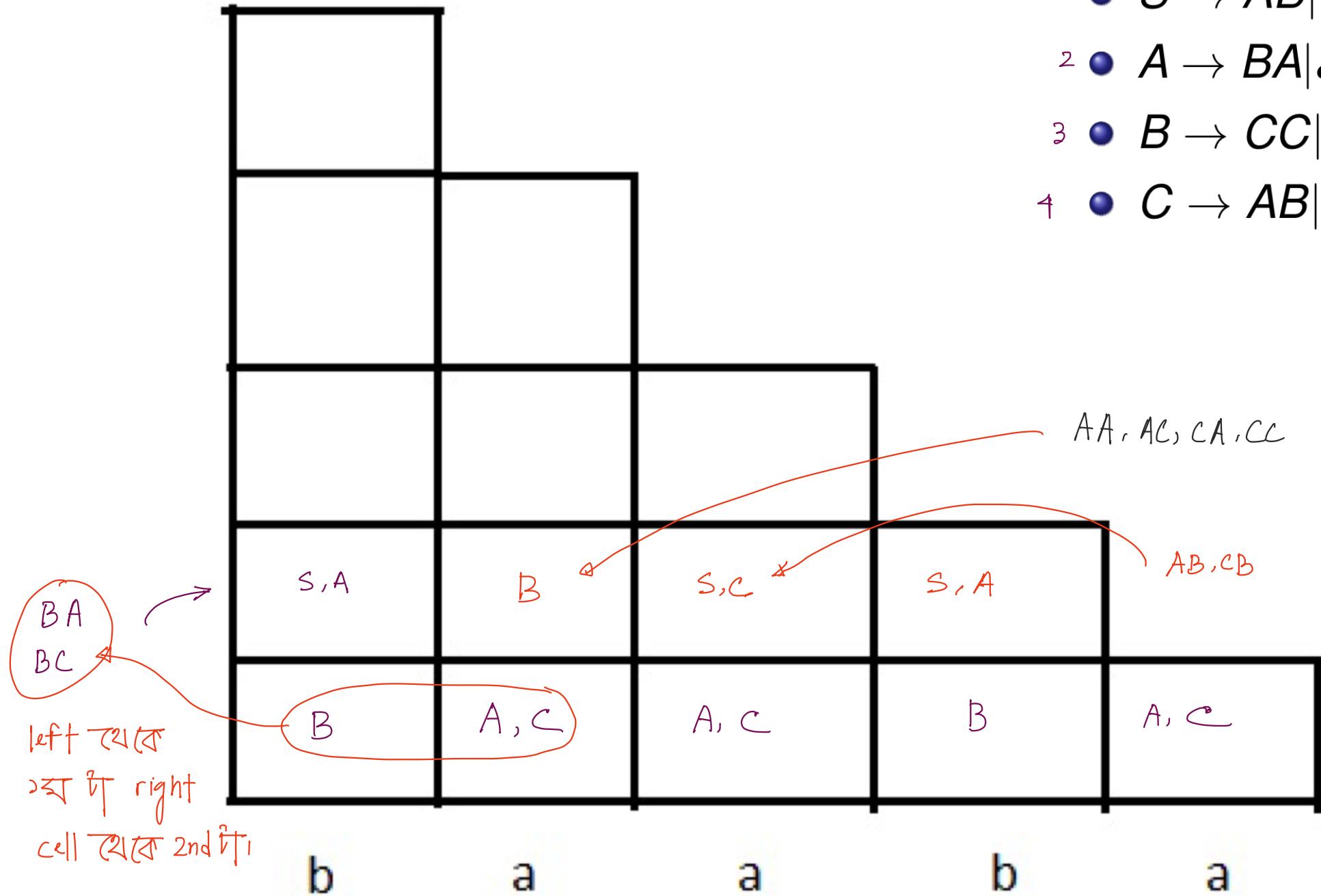
row by row fill up করবে।

* pencil use

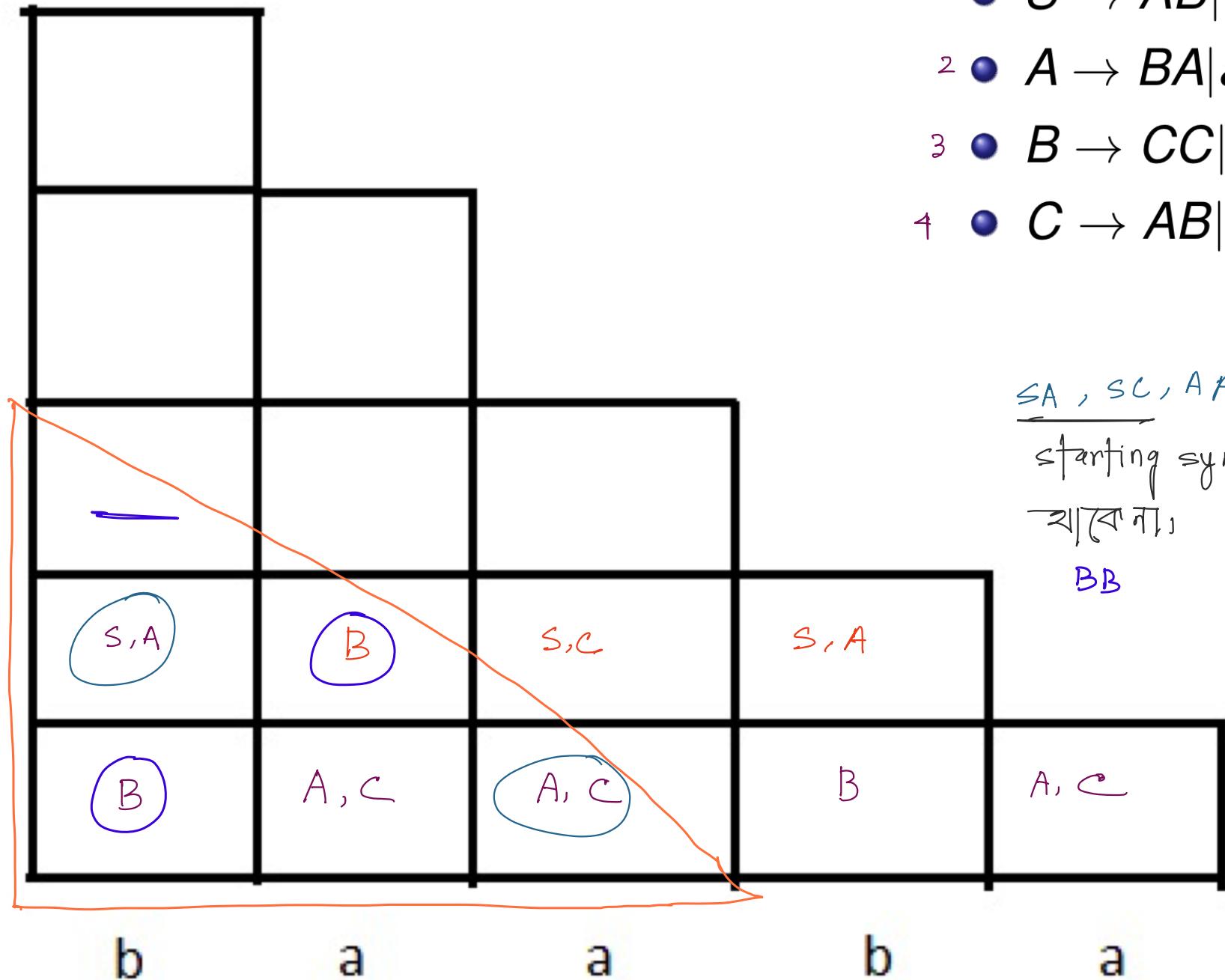


- 1 ● $S \rightarrow AB|BC$
- 2 ● $A \rightarrow BA|a$
- 3 ● $B \rightarrow CC|b$
- 4 ● $C \rightarrow AB|a$

- 1 • $S \rightarrow AB|BC$
- 2 • $A \rightarrow BA|a$
- 3 • $B \rightarrow CC|b$
- 4 • $C \rightarrow AB|a$

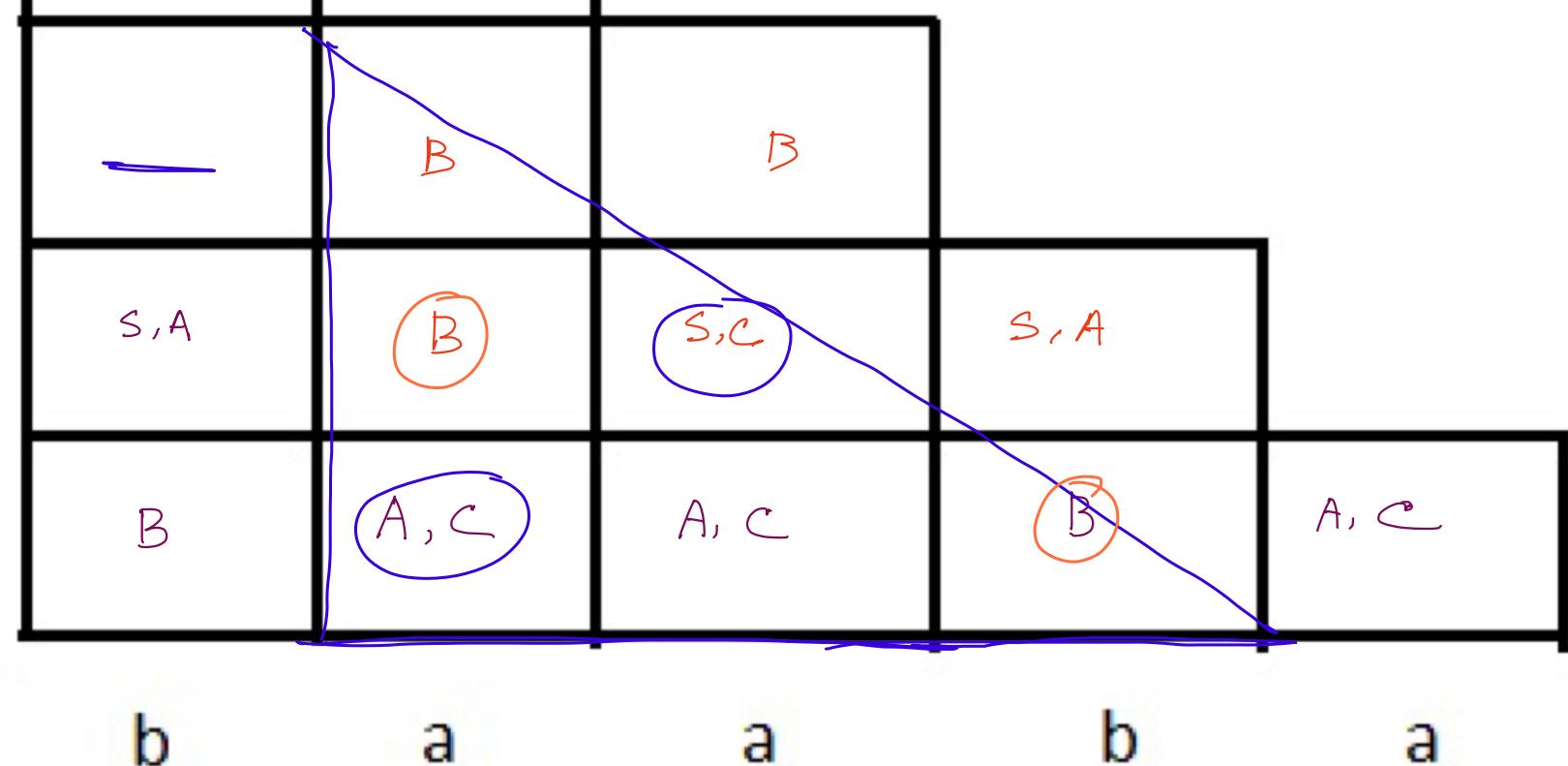


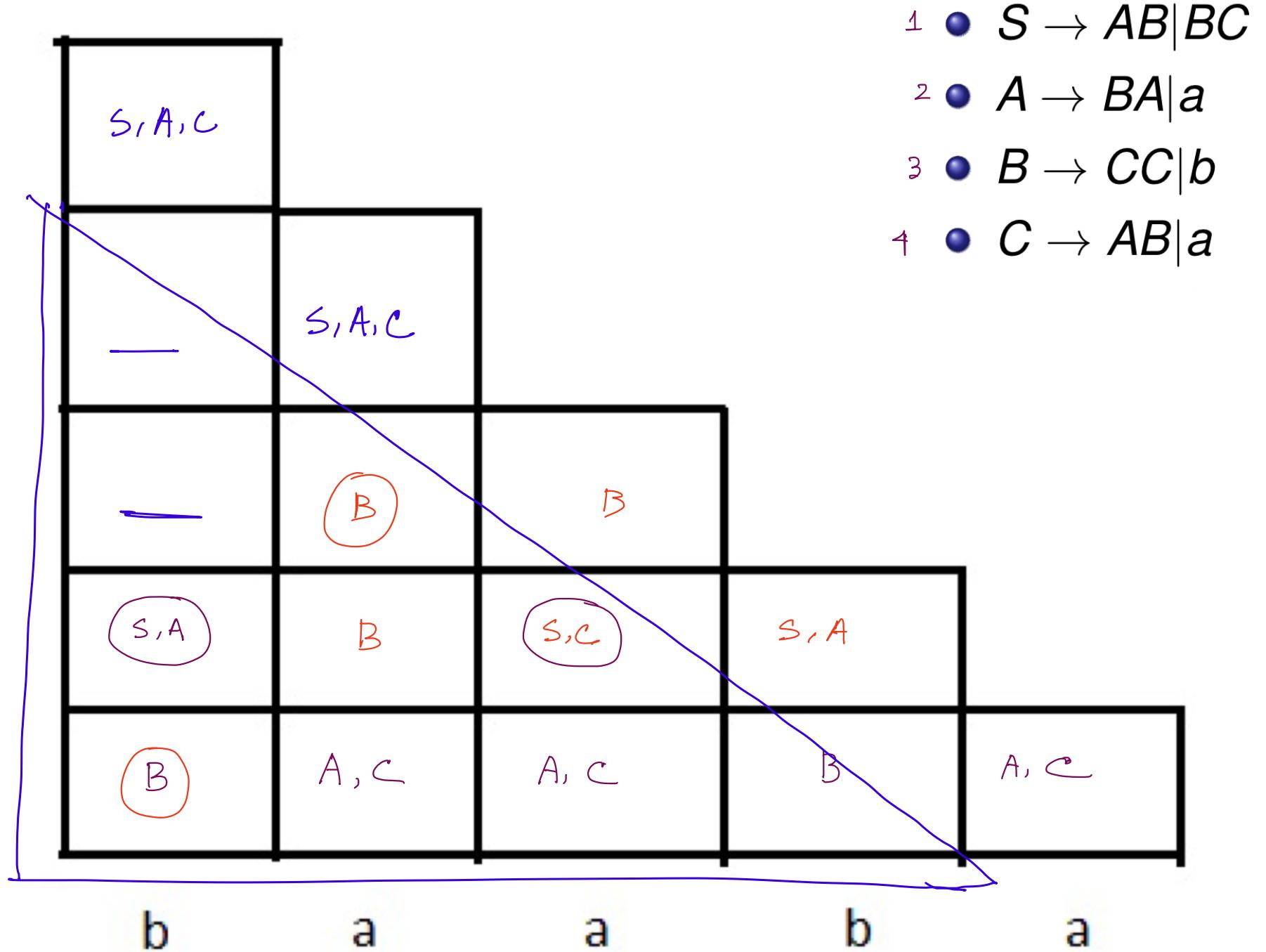
- 1 • $S \rightarrow AB|BC$
- 2 • $A \rightarrow BA|a$
- 3 • $B \rightarrow CC|b$
- 4 • $C \rightarrow AB|a$



- $\bullet S \rightarrow AB|BC$
- $\bullet A \rightarrow BA|a$
- $\bullet B \rightarrow CC|b$
- $\bullet C \rightarrow AB|a$

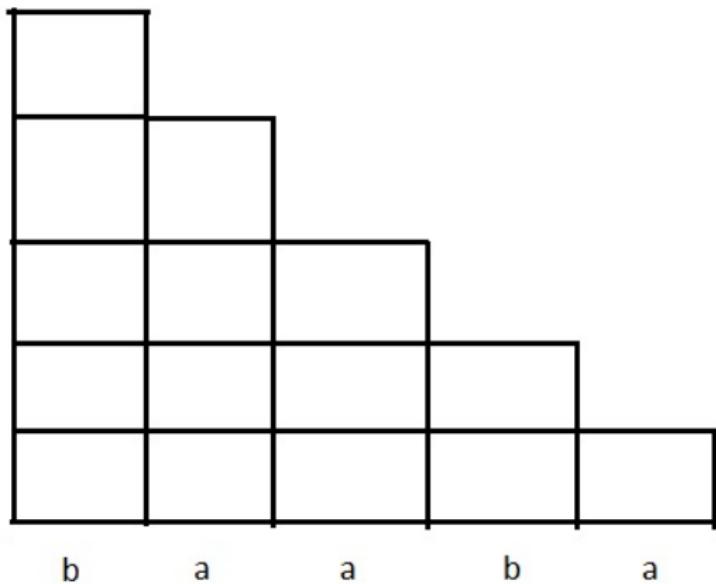
BB AS, AC, CS, CC





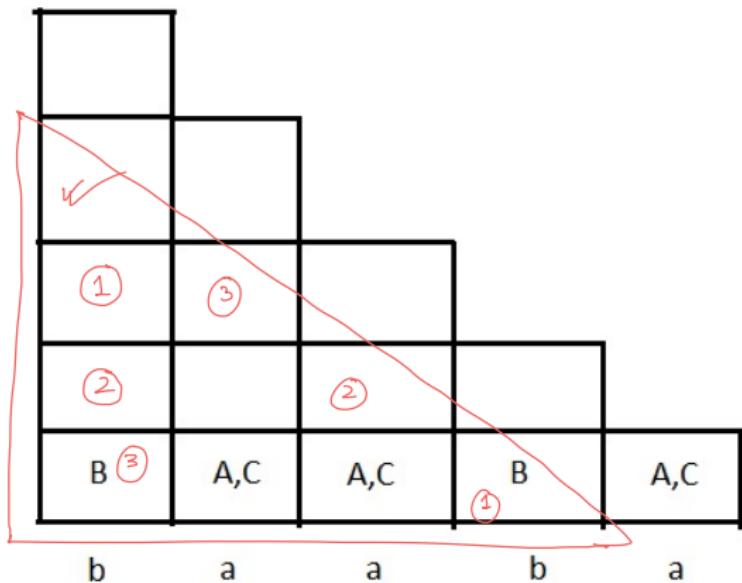
Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$



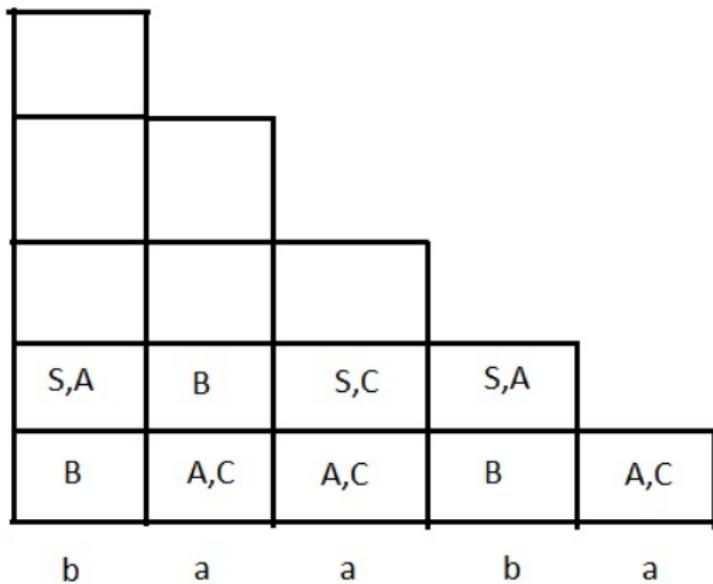
Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$



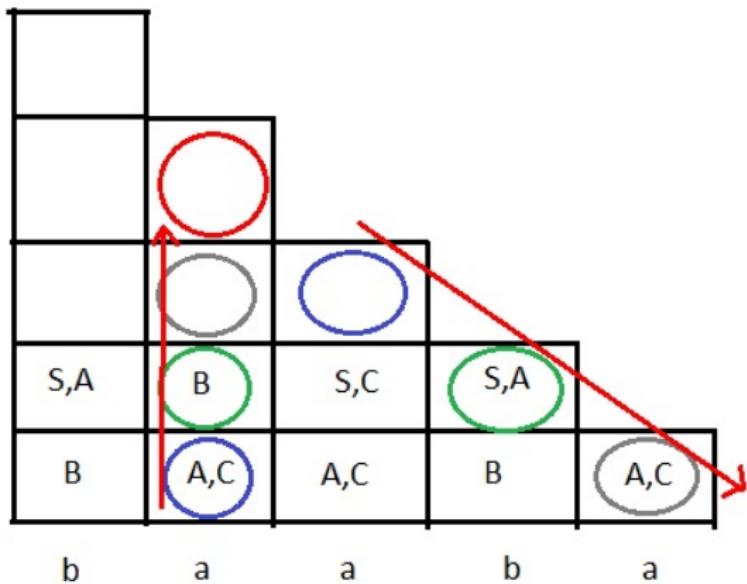
Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$



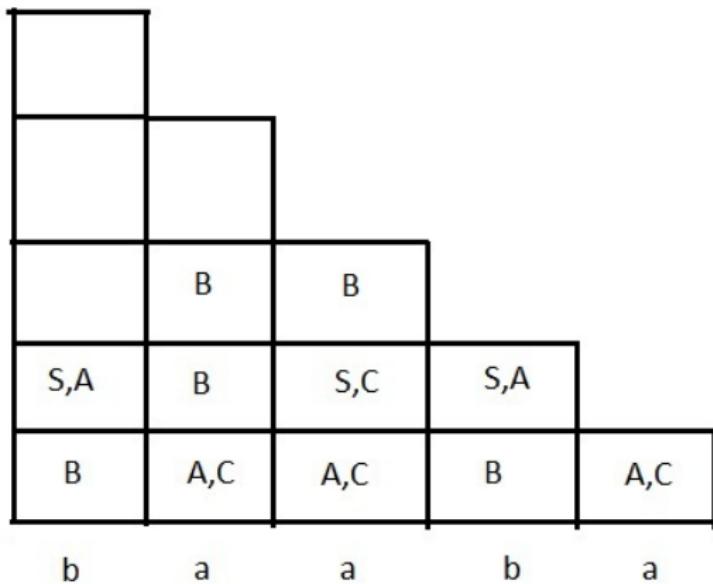
Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$



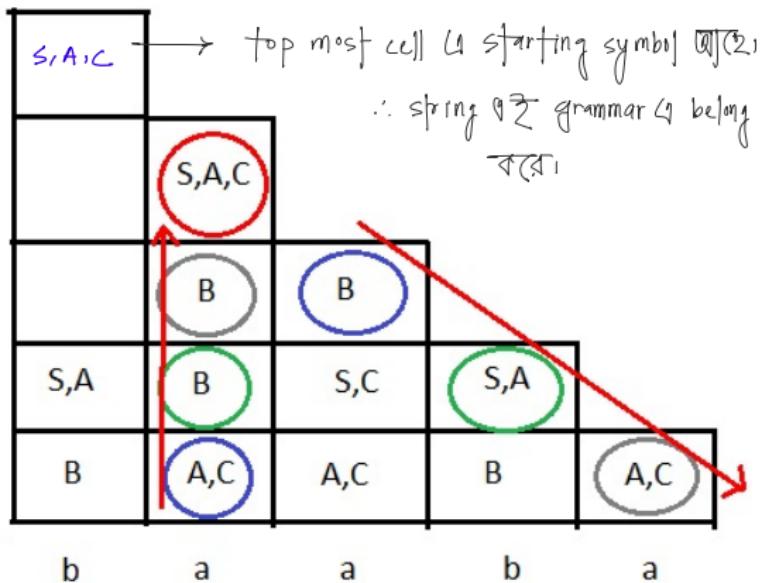
Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$



Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$



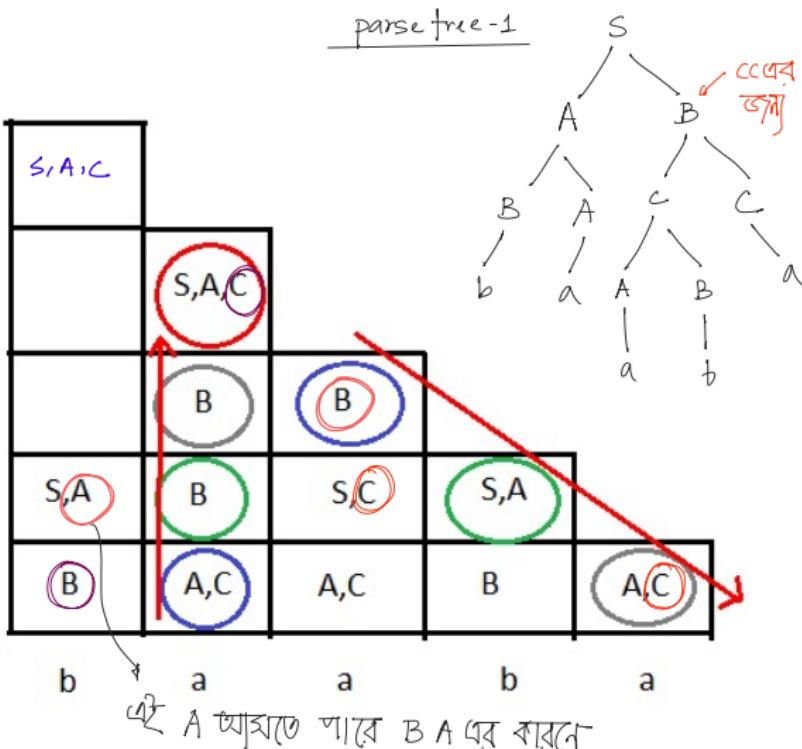
Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$

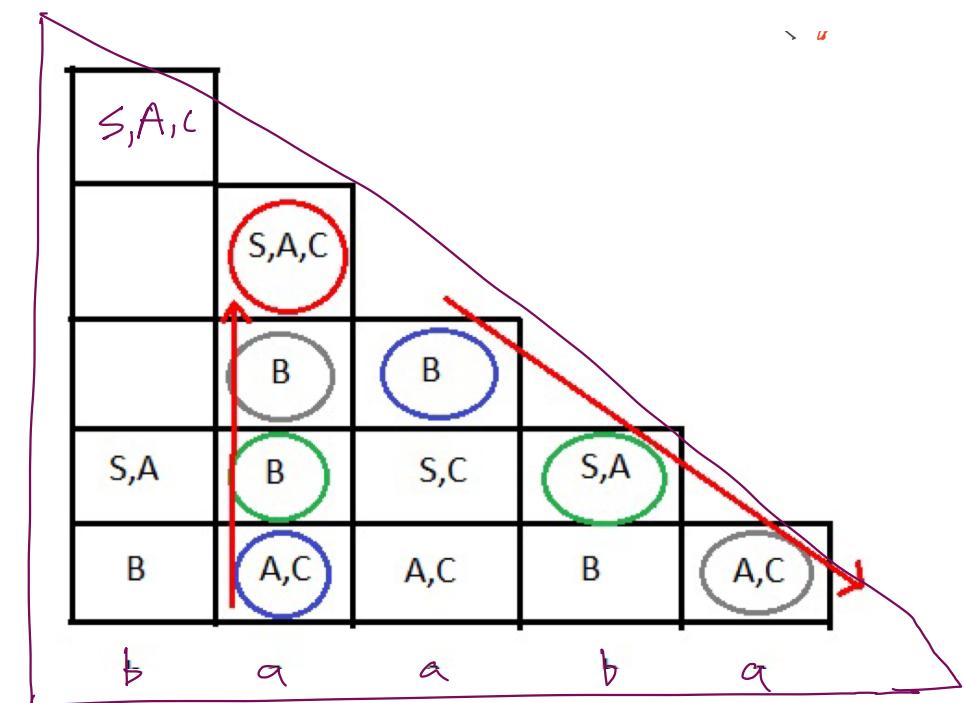
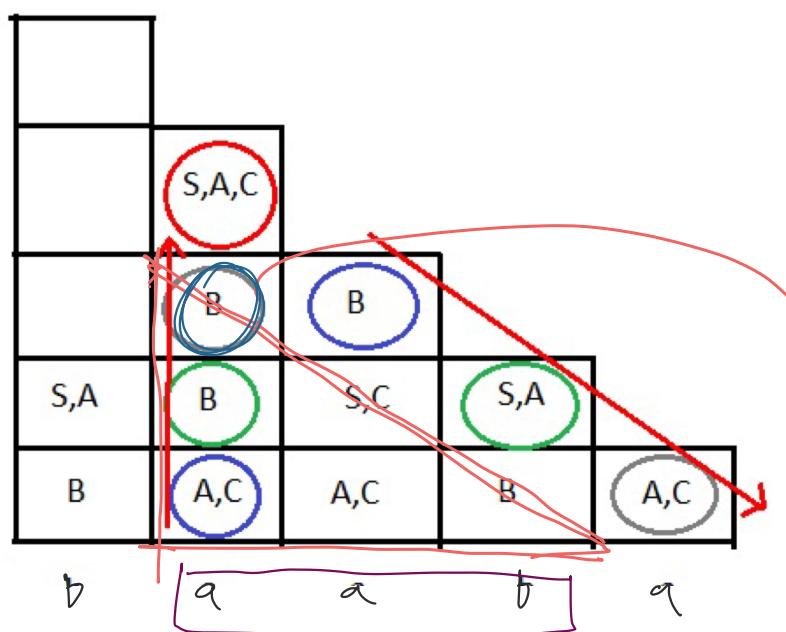
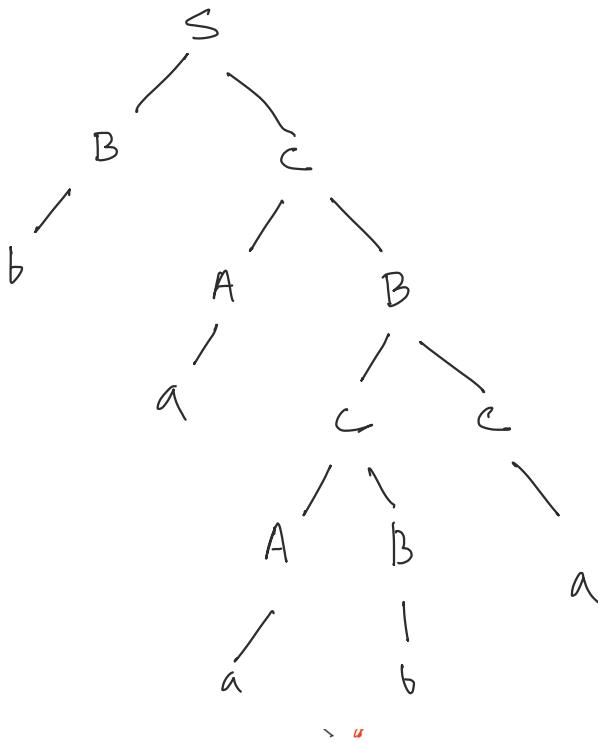
S থেকে শুরু করে দেখবো যান

মুঠ এর বারনে এটা এমেছু।

S আব্দতে পারে AB বা BC থেকে।



parse tree - 2



triangle এর topmost cell এর variable

মূলে থেকে derive করে করে

base string টা পাই।

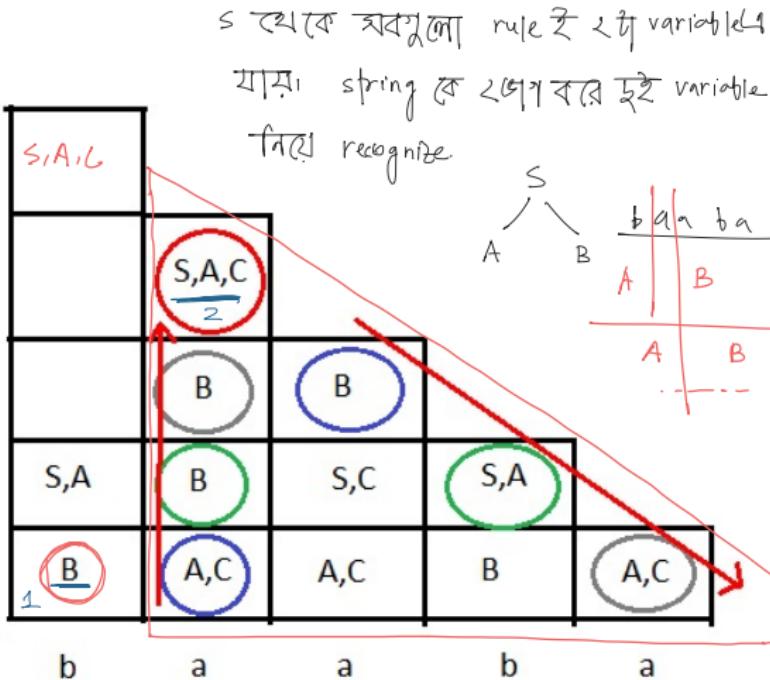
top most cell টা S নাই। so 1 9 6 এই CFG টা নাই।

Example

- $S \rightarrow AB|BC$
- $A \rightarrow BA|a$
- $B \rightarrow CC|b$
- $C \rightarrow AB|a$

1, 2 അ null ലൈംഗിക്ക്. So

belong കരാം!



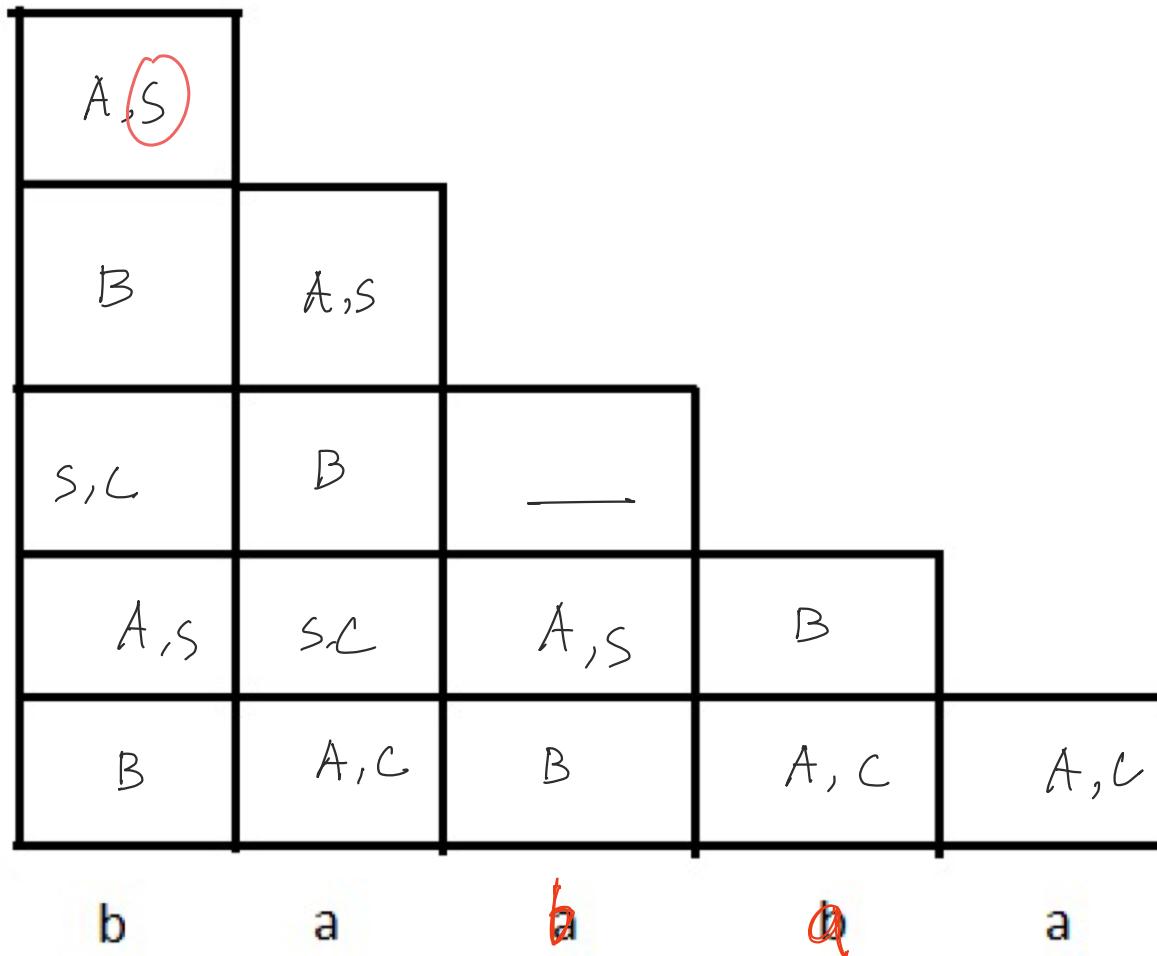
Example

					$\{S, A, C\}$
-					$\{S, A, C\}$
-					$\{B\} \quad \{B\}$
	$\{S, A\}$	$\{B\}$	$\{S, C\}$	$\{S, A\}$	
	$\{B\}$	$\{A, C\}$	$\{A, C\}$	$\{B\}$	$\{A, C\}$
	<i>b</i>	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>

Figure 7.14: The table for string *baaba* constructed by the CYK algorithm

practice baba

String Grammar \hookrightarrow belong করে
বিনা — practice



— belong করে

Q: CFG to CNF

then একটি string CNFG

আছে বিনা।

Chomsky Normal Form

- To convert a grammar to its Chomsky Normal Form, we need to make a number of preliminary simplifications:
 - We must eliminate *useless symbols*, those variables or terminals that do not appear in any derivation of a terminal string from the start symbol.
 - We must eliminate ϵ -*productions*, those of the form $A \rightarrow \epsilon$ for some variable A .
 - We must eliminate unit productions, those of the form $A \rightarrow B$ for variables A and B .

Chomsky Normal Form

Theorem

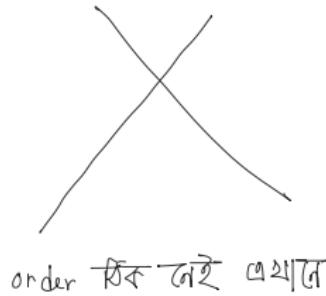
Any context-free language is generated by a context-free grammar in Chomsky normal form.

CT - 03

Converting to CNF

Perform the following steps in this order:

- i Eliminate ϵ productions
- ii Eliminate unit productions
- iii Eliminate useless symbols
- iv Convert to CNF
 - a Arrange that all bodies of length 2 or more consist only of variables
 - b Break bodies of length 3 or more into a cascade of productions, each with a body consisting of two variables



Steps:

1. useless symbol eliminate
2. new start symbol introduce
3. nullable variable handle
4. unit production handle
5. single terminal handle
6. CNF বানানো

①

e.g. $S \rightarrow AB | \alpha$

$$A \rightarrow b$$

- useless variable } start symbol থেকে start করে কিছু step পর terminal এ লীচাত চাই
1. start থেকে না লীচাত দারলে
 2. এই variable থেকে রানো terminal এ না লীচাত দারলে

$$S \rightarrow AB \mid a$$

$$A \rightarrow b$$

B থেকে ফোনো rule দিয়ে terminal produce করে গা,
— useless

not generating
useless symbol.



$$\begin{aligned} S &\rightarrow a \\ A &\rightarrow b \end{aligned}$$

— B related ঘরে rule বাদ দিয়ে দেলবো,



$$S \rightarrow a$$



CNF এ আববে।

(2)

grammer এর start symbol right hand side এ আছবে না,

$$S \rightarrow AB \mid ASB$$

$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b$$

start symbol S

→ 1st rule এর left

symbol

start symbol right এ আছবে,

∴ একটা start symbol introduce.

$$S' \rightarrow S$$

start symbol $\rightarrow S'$

$$S \rightarrow AB \mid ASB$$

$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b$$

③ nullable \rightarrow এমন বোন variable থেকে null generate করার rule আছে।

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aAA \mid \epsilon \\ B &\rightarrow bBB \mid \epsilon \end{aligned}$$

$\rightarrow S$ থেকে recursively null generate করা যাব।

$\therefore S, A, B \rightarrow \text{nullable variable}$

$$S \rightarrow ABC$$

$$A \rightarrow aAA \mid \epsilon$$

$$B \rightarrow bBB \mid \epsilon$$

nullable নঁ

$$S \rightarrow ABC \mid A B . \quad \rightarrow \text{nullable}$$

$$A \rightarrow aAA \mid \epsilon$$

$$B \rightarrow bBB \mid \epsilon$$

$$S \rightarrow AB$$

$$A \rightarrow aAA | \epsilon$$

$$B \rightarrow bBB | \epsilon$$

rule 1,2 \rightarrow done

3 \rightarrow এটি start symbol nullable হ্যাঁ

e.g. $S \rightarrow AB\epsilon$ \rightarrow extra rule add এম।
এবাটো null থাকবে। start থেকে,

$A \rightarrow aAA$ \rightarrow right side'এ যতগুলো null করা যাবে তা'রে সব consider কৰা
পারবে,

$$A \rightarrow \underline{aAA} \mid \underline{aA} \mid \underline{a}$$

প্রথমটি
null না

$\xrightarrow{\text{2nd/3rd } A}$

যেখনটি null

$\xrightarrow{A \rightarrow \text{null}}$

প্রতিটি rule-এ গায়ে check করবে nullable variable আছে কিনা।

$$B \rightarrow bBB \mid bB \mid b$$

$$S \rightarrow AB \mid A \mid B \mid \underline{\epsilon}$$

AB ≥ 1 এ null হলে

start থেকে null consider করছি,

④

unit production:

Left \sqsubseteq exactly variable

right \sqsubseteq exactly exactly variable

একটি variable-থেকে আরেকটি variable এ transform.

$$S \rightarrow AB \mid A \mid B \mid \epsilon$$

unit production ≥ 1

$$A \rightarrow aAA \mid aA \mid a$$

$$S \rightarrow A$$

$$B \rightarrow b \mid bBB \mid bB$$

$$S \rightarrow B$$

যে rule ও unit production থাবলৈ যেটা বাদ দিয়ে A থেকে যা আসছে যেগুলো
এখন S থেকে আসব।

$$S \rightarrow AB | \frac{1AA|aA|a}{A} | b | bB | bBB | \in$$

$S \rightarrow A$
handled

$S \rightarrow B$
handled

⑤ একের বেশি terminal বা terminal and variable এবং যাখে আছে right ও
যেগুলো দেখ কো। CNF ও right side ও একটীই terminal বা টু বু variable থাবলৈ
পারতো,

valid so change কো যাবেন।

$$S \rightarrow AB | 1AA | aA | a | b | bB | bBB | \in$$

$$\begin{array}{l} x \rightarrow a \\ y \rightarrow b \end{array} \quad \left\{ \begin{array}{l} \text{new variable introduce} \end{array} \right.$$

$$S \rightarrow AB | XAA | XA | a | b | YB | YBB | \in$$

$$S \rightarrow AB \mid AA \mid XA \mid a \mid b \mid YB \mid bBB \mid c$$

$$X \rightarrow a$$

$$Y \rightarrow b$$

$$A \rightarrow XAA \mid XA \mid a$$

$$B \rightarrow b \mid YBB \mid YB$$

right hand side \in either
 terminal \in more than 2 or
 variable \in

⑥

$$S \rightarrow XAA$$

$$\left\{ \begin{array}{l} P \rightarrow XA \\ S \rightarrow PA \end{array} \right.$$

$$S \rightarrow ABCD$$

$$\hookrightarrow S \rightarrow PQ$$

$$P \rightarrow AB$$

$$Q \rightarrow CD$$

$$S \rightarrow ABCD$$

$$\hookrightarrow S \rightarrow AP$$

$$P \rightarrow BQ$$

$$Q \rightarrow CD$$

$$S \rightarrow AB \mid AA \mid XA \mid a \mid b \mid YB \mid bBB \quad | \in$$
$$X \rightarrow a$$
$$Y \rightarrow b$$
$$A \rightarrow P A \mid XA \mid a$$
$$B \rightarrow b \mid Q B \mid YB$$
$$P \rightarrow XA$$
$$Q \rightarrow YB$$

$S \rightarrow S S \mid (S) \mid \in$ → valid parenthesis string recognize

CNF:

① done

② $S' \rightarrow S$

$S \rightarrow S S \mid (S) \mid \in$

③ $S' \rightarrow S \mid \in$

$S \rightarrow S S \mid (S) \mid S \mid ()$



$S \rightarrow S S \mid (S) \mid ()$

$S \rightarrow S$ কোনো rule নাই,
বাদ দিবো

④

$S' \rightarrow S S \mid (S) \mid () \mid \in$

$S \rightarrow S S \mid (S) \mid ()$

⑤

$L \rightarrow ($

$R \rightarrow)$

$$S' \rightarrow SS | LS R | LR | \epsilon$$
$$S \rightarrow SS | LS R | LR$$

⑥

$$S' \rightarrow SS | XR | LR | \epsilon$$
$$S \rightarrow SS | XR | LR$$
$$X \rightarrow LS$$
$$L \rightarrow ($$
$$R \rightarrow)$$

Eliminating Useless Symbols

- Two things a symbol has to be able to do to be useful
 - We say X is generating if $X \xrightarrow{*} w$ for some terminal string w . Note that every terminal is generating since w can be that terminal itself
 - We say X is reachable if there is a derivation $S \xrightarrow{*} \alpha X \beta$ for some α and β
- If a symbol is not useful, it is *useless*

Example

- Consider the grammar

$$\begin{aligned}S &\rightarrow AB \mid a \\A &\rightarrow b\end{aligned}$$

- Find generating and reachable symbols using induction
- B is not generating

$$\begin{aligned}S &\rightarrow a \\A &\rightarrow b\end{aligned}$$

- A is not reachable

$$S \rightarrow a.$$

Eliminating ϵ -productions

- Discover variables that are *nullable*
 - A variable A is nullable if $A \xrightarrow{*} \epsilon$
- If A is nullable, then whenever A appears in a production body, say $B \rightarrow CAD$, A might or might not derive ϵ . We make two versions of the production
 - one without A in the body $B \rightarrow CD$ which corresponds to the case where A would have been used to derive ϵ
 - and the other with A still present $B \rightarrow CAD$
- If language contains ϵ , add $S \rightarrow \epsilon$ where S is the start symbol

Example

- Consider the grammar

$$\begin{aligned}S &\rightarrow AB \\A &\rightarrow aAA \mid \epsilon \\B &\rightarrow bBB \mid \epsilon\end{aligned}$$

- A, B and S are nullable
- Production 1 becomes

$$S \rightarrow AB \mid A \mid B$$

- Production 2 becomes

$$A \rightarrow aAA \mid aA \mid aA \mid a$$

Example

- Similarly

$$B \rightarrow bBB \mid bB \mid b$$

- So, the grammar after eliminating ϵ -productions is

$$\begin{aligned}S &\rightarrow AB \mid A \mid B \\A &\rightarrow aAA \mid aA \mid a \\B &\rightarrow bBB \mid bB \mid b\end{aligned}$$

- Since S is nullable add $S \rightarrow \epsilon$

Eliminating unit productions

- A unit production is a production of the form $A \rightarrow B$ where both A and B are variables
- Identify *unit pairs*
 - A pair (A, B) is called unit pair if $A \xrightarrow{*} B$ using only unit productions
- For each unit pair (A, B) , add all the productions $A \rightarrow \alpha$, where $B \rightarrow \alpha$ is a nonunit production. Note that $A = B$ is possible in that way. Only the non-unit productions remain

Example

- Consider the grammar

$$\begin{array}{lcl}
 I & \rightarrow & a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\
 F & \rightarrow & I \mid (E) \\
 T & \rightarrow & F \mid T * F \\
 E & \rightarrow & T \mid E + T
 \end{array}$$

- Find the unit pairs. $(E, E), (T, T), (F, F), (I, I)$ are unit pairs by zero steps

- (E, E) and the production $E \rightarrow T$ gives us unit pair (E, T) .
- (E, T) and the production $T \rightarrow F$ gives us unit pair (E, F) .
- (E, F) and the production $F \rightarrow I$ gives us unit pair (E, I) .
- (T, T) and the production $T \rightarrow F$ gives us unit pair (T, F) .
- (T, F) and the production $F \rightarrow I$ gives us unit pair (T, I) .
- (F, F) and the production $F \rightarrow I$ gives us unit pair (F, I) .

Example

- The productions to be added/kept

Pair	Productions
(E, E)	$E \rightarrow E + T$
(E, T)	$E \rightarrow T * F$
(E, F)	$E \rightarrow (E)$
(E, I)	$E \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$
(T, T)	$T \rightarrow T * F$
(T, F)	$T \rightarrow (E)$
(T, I)	$T \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$
(F, F)	$F \rightarrow (E)$
(F, I)	$F \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$
(I, I)	$I \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$

Example

- The resulting grammar

$$\begin{aligned} E &\rightarrow E + T \mid T * F \mid (E) \mid a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\ T &\rightarrow T * F \mid (E) \mid a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\ F &\rightarrow (E) \mid a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\ I &\rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \end{aligned}$$

Converting to CNF

- The grammar has had its ϵ -productions, unit productions and useless symbols removed
- Our tasks are to
 - Arrange that all bodies of length 2 or more consist only of variables
 - Break bodies of length 3 or more into a cascade of productions, each with a body consisting of two variables

Example

- Consider the grammar

$$\begin{aligned} E &\rightarrow E + T \mid T * F \mid (E) \mid a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\ T &\rightarrow T * F \mid (E) \mid a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\ F &\rightarrow (E) \mid a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \\ I &\rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1 \end{aligned}$$

Example

- Eight terminals $a, b, 0, 1, +, *, (,$ and $),$ appears in a body that is not a single terminal
- We must introduce eight new variables, corresponding to these terminals, and eight productions in which the new variable is replaced by its terminal

$$\begin{array}{llll} A \rightarrow a & B \rightarrow b & Z \rightarrow 0 & O \rightarrow 1 \\ P \rightarrow + & M \rightarrow * & L \rightarrow (& R \rightarrow) \end{array}$$

Example

- We introduce these productions, and replace every terminal in a body that is other than a single terminal by the corresponding variable

$E \rightarrow EPT \mid TMF \mid LER \mid a \mid b \mid IA \mid IB \mid IZ \mid IO$
 $T \rightarrow TMF \mid LER \mid a \mid b \mid IA \mid IB \mid IZ \mid IO$
 $F \rightarrow LER \mid a \mid b \mid IA \mid IB \mid IZ \mid IO$
 $I \rightarrow a \mid b \mid IA \mid IB \mid IZ \mid IO$
 $A \rightarrow a$
 $B \rightarrow b$
 $Z \rightarrow 0$
 $O \rightarrow 1$
 $P \rightarrow +$
 $M \rightarrow *$
 $L \rightarrow ($
 $R \rightarrow)$

Example

- Introduce variables to break bodies of length 3 or more

E	\rightarrow	$EC_1 \mid TC_2 \mid LC_3 \mid a \mid b \mid IA \mid IB \mid IZ \mid IO$
T	\rightarrow	$TC_2 \mid LC_3 \mid a \mid b \mid IA \mid IB \mid IZ \mid IO$
F	\rightarrow	$LC_3 \mid a \mid b \mid IA \mid IB \mid IZ \mid IO$
I	\rightarrow	$a \mid b \mid IA \mid IB \mid IZ \mid IO$
A	\rightarrow	a
B	\rightarrow	b
Z	\rightarrow	0
O	\rightarrow	1
P	\rightarrow	$+$
M	\rightarrow	$*$
L	\rightarrow	$($
R	\rightarrow	$)$
C_1	\rightarrow	PT
C_2	\rightarrow	MF
C_3	\rightarrow	ER