

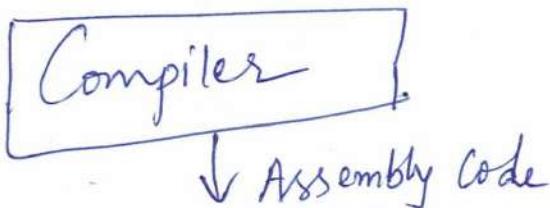
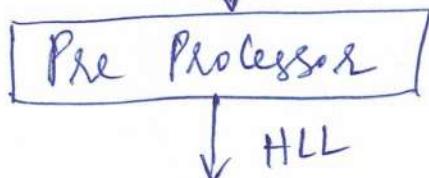
①

# Language Processing System

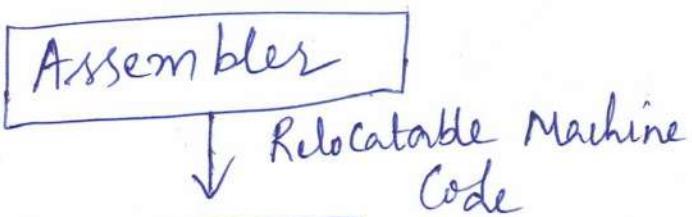


High Level Language(HLL)

① → Linker:-



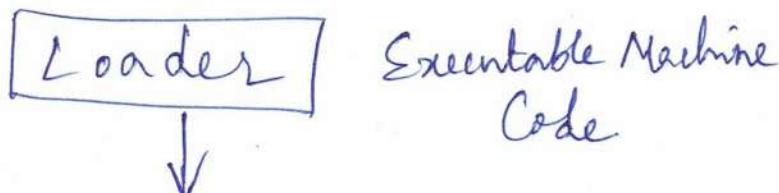
② → Loader :-



③ → Cross Compiler :-



④ → Console to Console Compiler :-



P

# Language Processing System

①-②



High Level Language(HLL)

↓  
Pre Processor

↓ HLL

↓ Compiler

↓ Assembly Code

↓ Assembler

↓ Relocatable Machine  
Code

↓ Linker

↓ Executable Machine  
Code

↓ Loader

① → Linker:- It links and merges various object files together in order to make an executable file.

② → Loader:- It is responsible for loading executable files into memory & execute them.

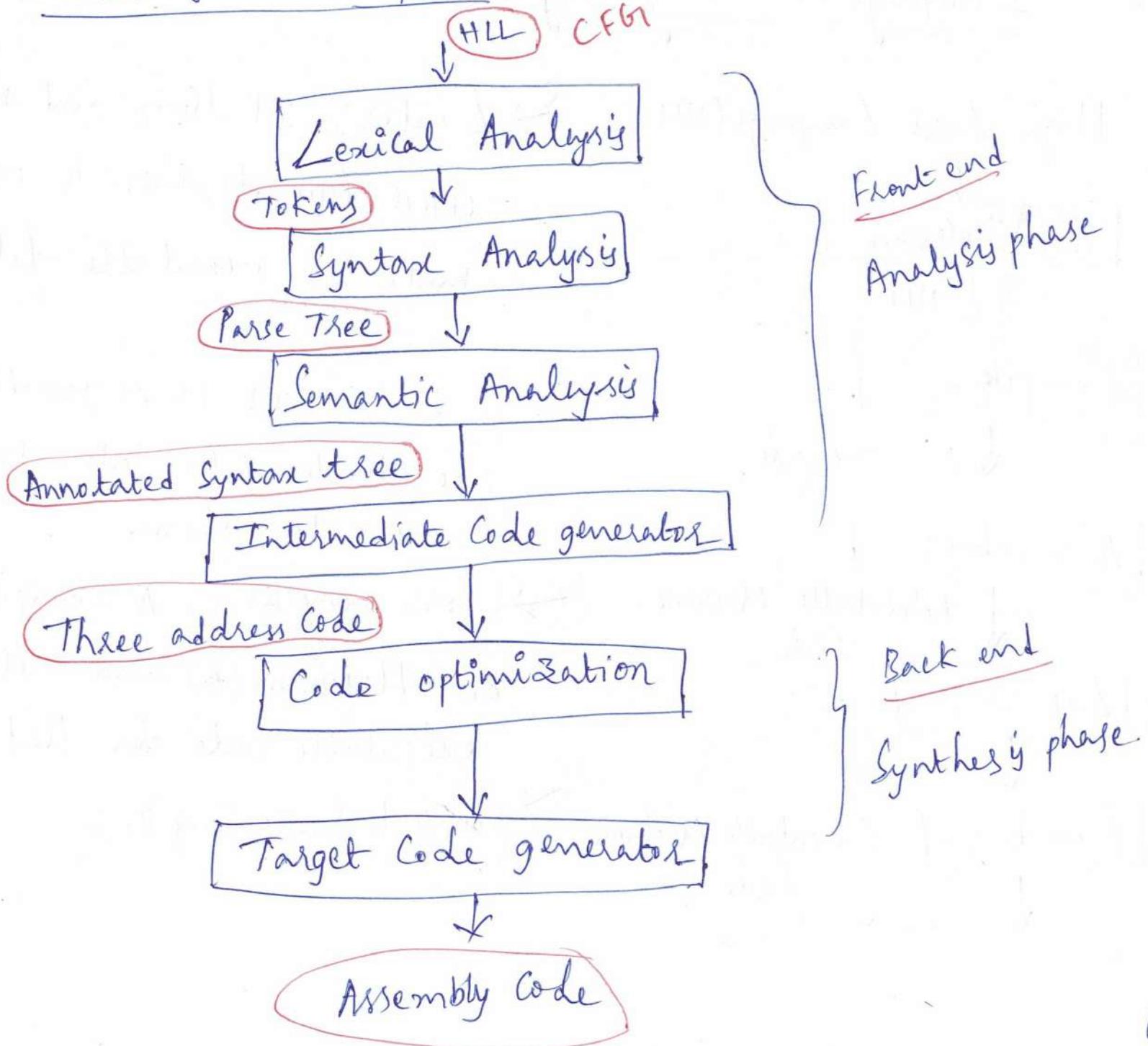
③ → Cross Compiler:- A Compiler that runs on Platform (A) and able to generate executable code for Platform B.

④ → Source to Source Compiler:-

P

# Phases of a Compiler

Symbol Table



Front-end  
Analysis phase

Back-end  
Synthesis phase

Q

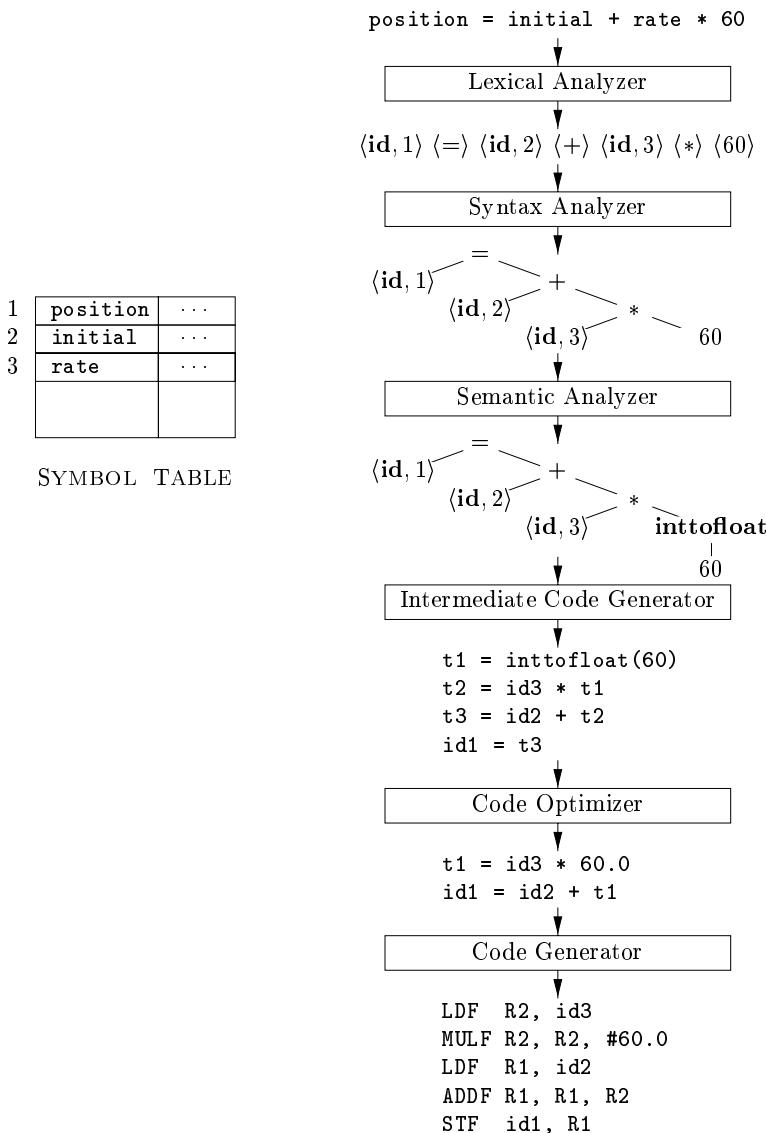


Figure 1.7: Translation of an assignment statement

# Deterministic Finite State Machine (DFA)

④

→ 5-tuple  $(Q, \Sigma, \delta, q_0, F)$

→ ① A finite set of states -  $Q$

→ ② A " " of input symbols -  $\Sigma$

→ ③ A transition function -  $\delta: Q \times \Sigma \rightarrow Q$

Important

→ ④ An initial or start state -  $q_0 \in Q$

→ ⑤ A set of accept states -  $F \subseteq Q$

---

$\Sigma = \{0, 1\}$  - Binary

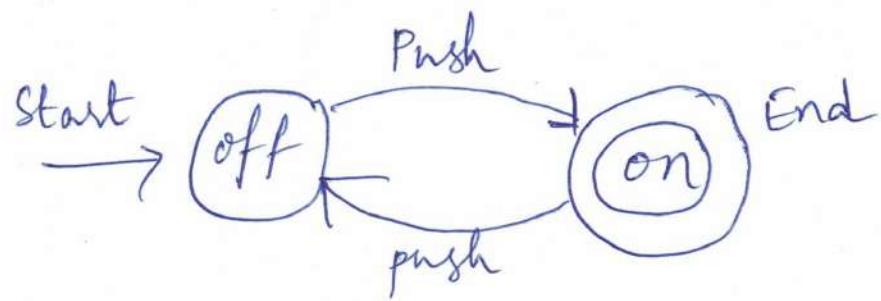
$\Sigma$  = Set of all ASCII values

$\Sigma = \{a, b, \dots, z\}$  - Alphabet

P

# Example of A finite Automate modeling on/off switch.

⑤



$$Q = \{ \text{off}, \text{on} \}$$

$$\Sigma = \{ \text{push} \}$$

$$q_0 = \{ \text{off} \}$$

$$F = \{ \text{on} \}$$

$$\delta(\text{on}, \text{push}) = \text{off}$$

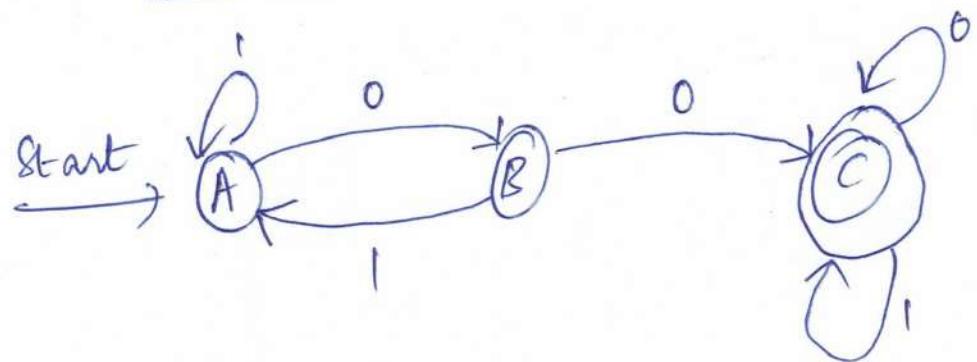
$$\delta(\text{off}, \text{push}) = \text{on}$$

$\delta$	push
on	<del>off</del> off
off	on

Q

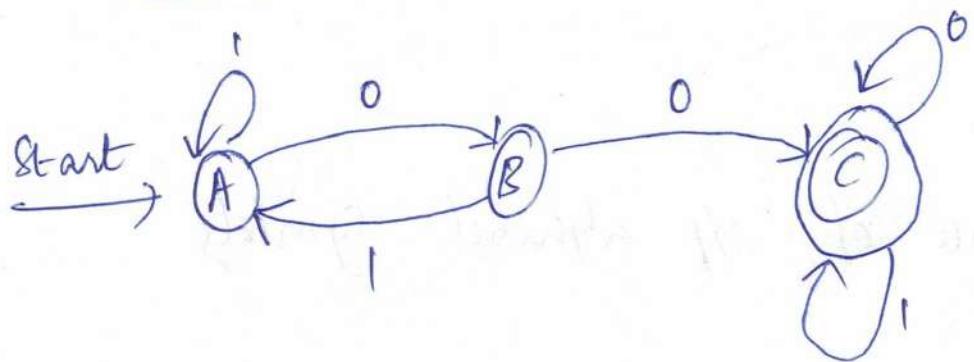
(6)

Another example of DFA



P

## Another example of DFA



$$Q = \{A, B, C\}$$

$$\Sigma = \{0, 1\}$$

$$q_0 = \{A\}$$

$$F = \{C\}$$

$\Sigma$	0	1
A	B	A
B	C	A
C	C	C

State transition table

P

Alphabet

$$\Sigma = \{0, 1\}$$

Strings/Word:- A finite sequence of i/p alphabet symbols

Empty string:-  $\epsilon \rightarrow \text{Epsilon}$   
zero occurrence

$$w = w\epsilon = \epsilon w$$

$$001 = 001\epsilon = \epsilon 001 = 0\epsilon 01$$

$$x = a_1, a_2, \dots, a_n, a_i \in \Sigma$$

$$y = b_1, b_2, \dots, b_m, b_i \in \Sigma$$

$$|xy| = |a_1, a_2, \dots, a_n, b_1, b_2, \dots, b_m| = m+n$$

} Concatenation.

Q

⑨

## Power of Alphabet

$$\Sigma = \text{Input alphabet} = \{a, b\}$$

$$\Sigma^k = \{(x_1, x_2, \dots, x_k) \mid x_i \in \Sigma\}$$

$$\Sigma^0 = \epsilon$$

$$\Sigma^1 = \Sigma = \{a, b\}$$

$$\Sigma^2 = \{aa, ab, ba, bb\}$$

$$\Sigma^3 = \{aaa, aab, aba, abb, \dots\}$$

$$\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots$$

$$\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots$$

$$\Sigma^* = \Sigma^+ \cup \{\epsilon\}$$

⑩

(10)

Language  $L \subseteq \Sigma^*$ 

→ Any subset of  $\Sigma^*$  is called a language over  $\Sigma$ .

Example 1

$L =$  language of strings consisting  $n$  0's followed by  $n$  1's.  
 $n \geq 0$

$$= \{ \epsilon, 01, 0011, 000111, \dots \}$$

$$= \{ 0^n 1^n \mid n \geq 0 \} \subseteq \Sigma^*$$

Example 2:-  $L =$  set of all strings having equal 0's and 1's

$$= \{ \epsilon, 01, 10, 0011, 0101, 0110, 1100, \dots \}$$

(Q)

Language  $L \subseteq \Sigma^*$

May contain an infinite number of strings  
drawn from a finite i/p alphabet set  $\Sigma$ .

Problem:- Decision problem

$\Sigma \rightarrow$  An alphabet

$L \subseteq \Sigma^*$

Given a string  $w \in \underline{\Sigma^*}$

decide whether  $\underline{w \in L}$  or not ??

## How a DFA processing Works

(12)

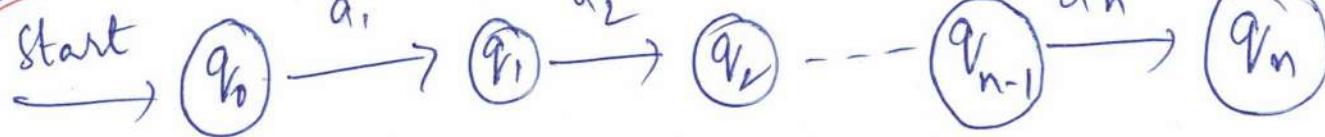
~~Given~~

$$A = \{Q, \Sigma, \delta, q_0, F\}$$

$$\delta: Q \times \Sigma \rightarrow Q$$

$$w = a_1 a_2 \dots a_n, a_i \in \Sigma$$

~~Processing~~



If  $q_n \in F \Leftrightarrow w$  is accepted by DFA.

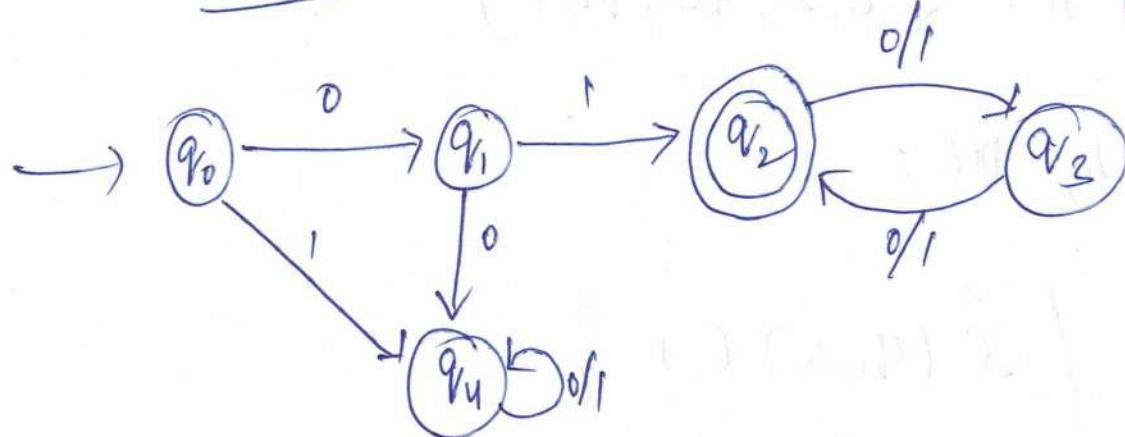
$q_n \notin F \Leftrightarrow w$  is rejected.

P.

## Example

(13)

$$Q = \{q_0, q_1, q_2, q_3, q_4\} \quad F = \{q_2\}$$



$$w = 011101$$

$s$	0	1
$q_0$	$q_1$	$q_4$
$q_1$	$q_4$	$q_2$
$q_2$	$q_3$	$q_3$
$q_3$	$q_2$	$q_2$
$q_4$	$q_4$	$q_4$

$\stackrel{s}{\Delta}(q_0, w) \in F ??$

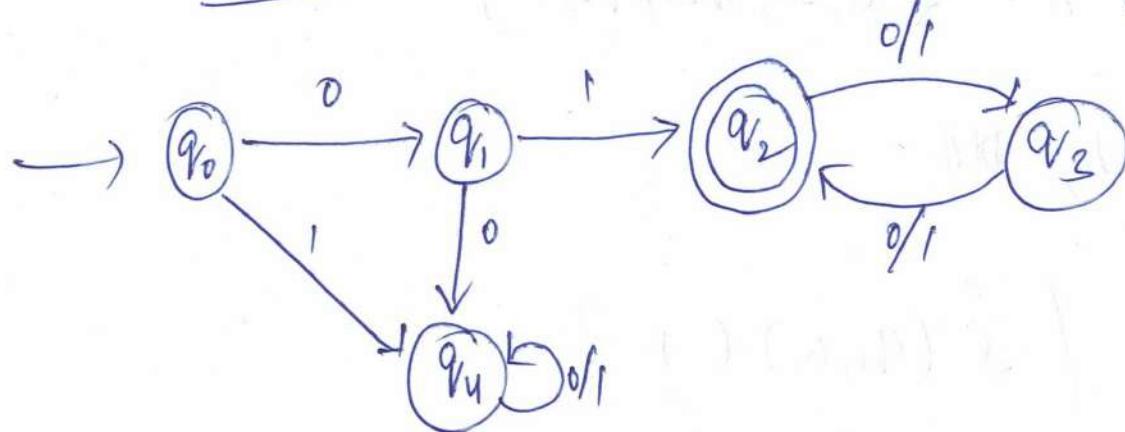
w accepted or not ??

(P)

## Example

$$Q = \{q_0, q_1, q_2, q_3, q_4\} \quad F = \{q_2\}$$

(13) (14)

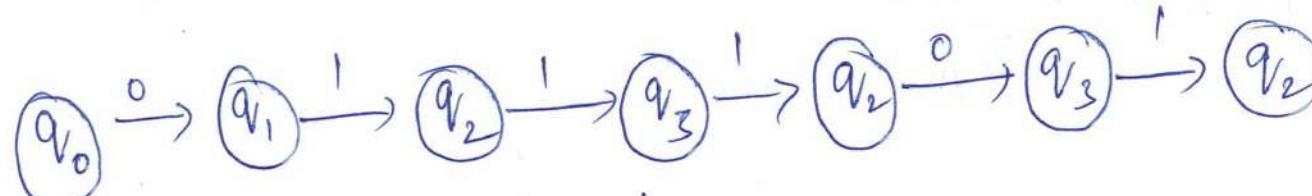


$$w = 011101$$

$S$	0	1
$q_0$	$q_1$	$q_4$
$q_1$	$q_4$	$q_2$
$q_2$	$q_3$	$q_3$
$q_3$	$q_2$	$q_2$
$q_4$	$q_4$	$q_4$

$\stackrel{?}{S}(q_0, w) \in F$  ??

w accepted or not ??



Processing.

This  $q_2 \in F$   
 $\therefore w$  is accepted by DFA.

P

(15) (Q)

$$\delta^1(q_0, w) \in F \quad A = \{ Q, \Sigma, \delta, q_0, F \}$$

→  $w$  is accepted by DFA

$$L(A) = \{ w \in \Sigma^* \mid \delta^1(q_0, w) \in F \}$$

⇒ Language of the DFA

---

$\Sigma \rightarrow$  I/P alphabet

$$L \subseteq \Sigma^*$$

→ Regular Language

$$\text{If } L = L(A)$$

DFA.

P

## Building DFA

$$\Sigma = \{0, 1\}$$

$L = \{w \in \Sigma^* \mid w \text{ is binary string with odd number of } 1's\}$

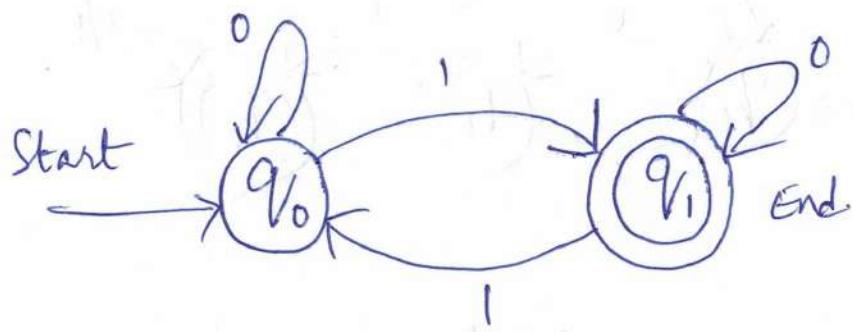
$1 \in L, 01 \in L, 001 \in L$   
 $111 \in L, 010101 \in L - \text{etc.}$

Building DFA

$$\Sigma = \{0, 1\}$$

$L = \{w \in \Sigma^* \mid n \text{ is binary string with odd number of } 1's\}$

$1 \in L, 01 \in L, 001 \in L,$   
 $111 \in L, 010101 \in L - \text{etc.}$



Trial and error Method.

$L = L(A)$ , A is a DFA

$\therefore L$  is a Regular Language.

$$\Sigma = \{0, 1\}$$

$L = \{w \mid w \text{ is a binary string containing '00' as substring}\}$

$$= \{x00y \mid x, y \in \Sigma^*\}$$

$00 \in L, 0010 \in L, 100 \in L, 001 \in L \dots \text{etc}$

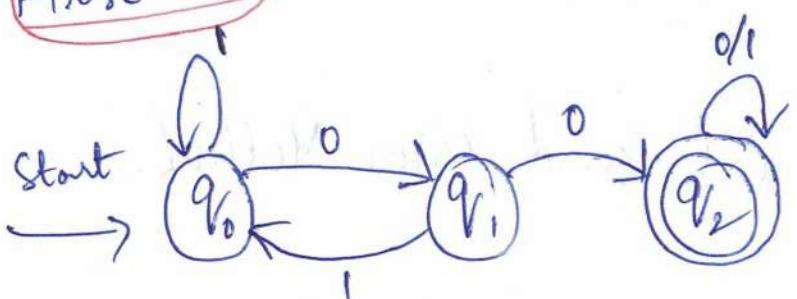
$$\Sigma = \{0, 1\}$$

$L = \{w \mid w \text{ is a binary string containing '00' as substring}\}$

$$= \{x00y \mid x, y \in \Sigma^*\}$$

$00 \in L, 0010 \in L, 100 \in L, 001 \in L \dots \text{etc}$

First DFA



DFA

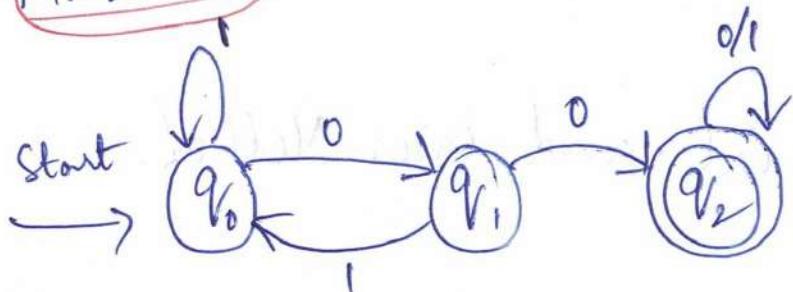
$$\Sigma = \{0, 1\}$$

$L = \{w \mid w \text{ is a binary string containing '00' as substring}\}$

$$= \{x00y \mid x, y \in \Sigma^*\}$$

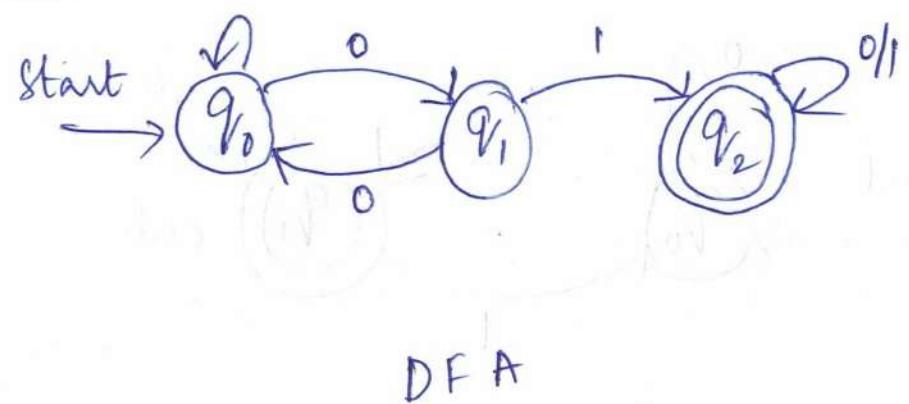
$00 \in L$ ,  $0010 \in L$ ,  $100 \in L$ ,  $001 \in L$  etc

First DFA



DFA

Second DFA



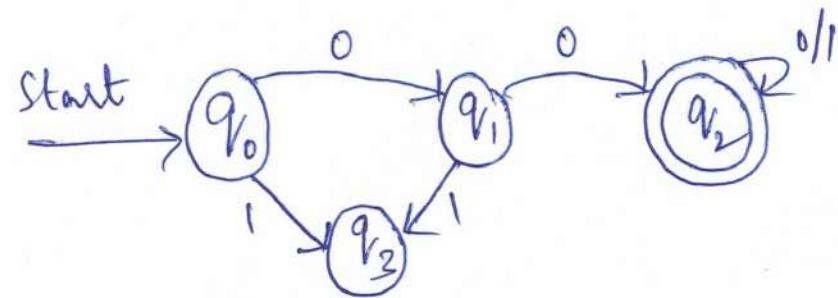
DFA

For both DFA's

$$L = L(A) \quad \therefore L \text{ is Regular Language.}$$

Example ①  $\Sigma = \{0, 1\}$  (21)

DFA that accepts all binary strings with starting '00'!

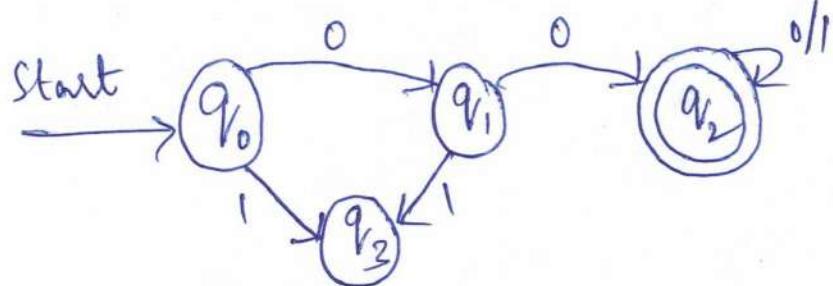


Example ② Example ① + even length  
Accept all binary strings starting '00' and length is even.

(P)

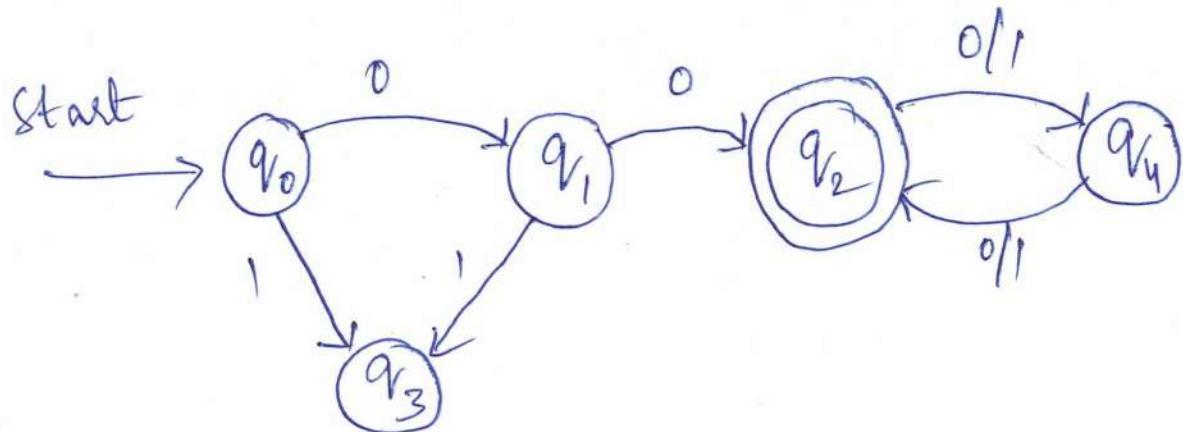
Example ①  $\Sigma = \{0, 1\}$  (21) (22)

DFA that accepts all binary strings with starting '00'.



$00 \in L$ ,  $001 \in L$ ,  $0000 \in L$  ... etc.

Example ② Example ① + even length  
Accept all binary strings starting '00' and length is even.



$0000 \in L$ ,  $0001 \notin L$   
 $0010 \in L$  ... etc.

(P)

DFA that accept all binary strings where 0's and 1's  
are alternate.

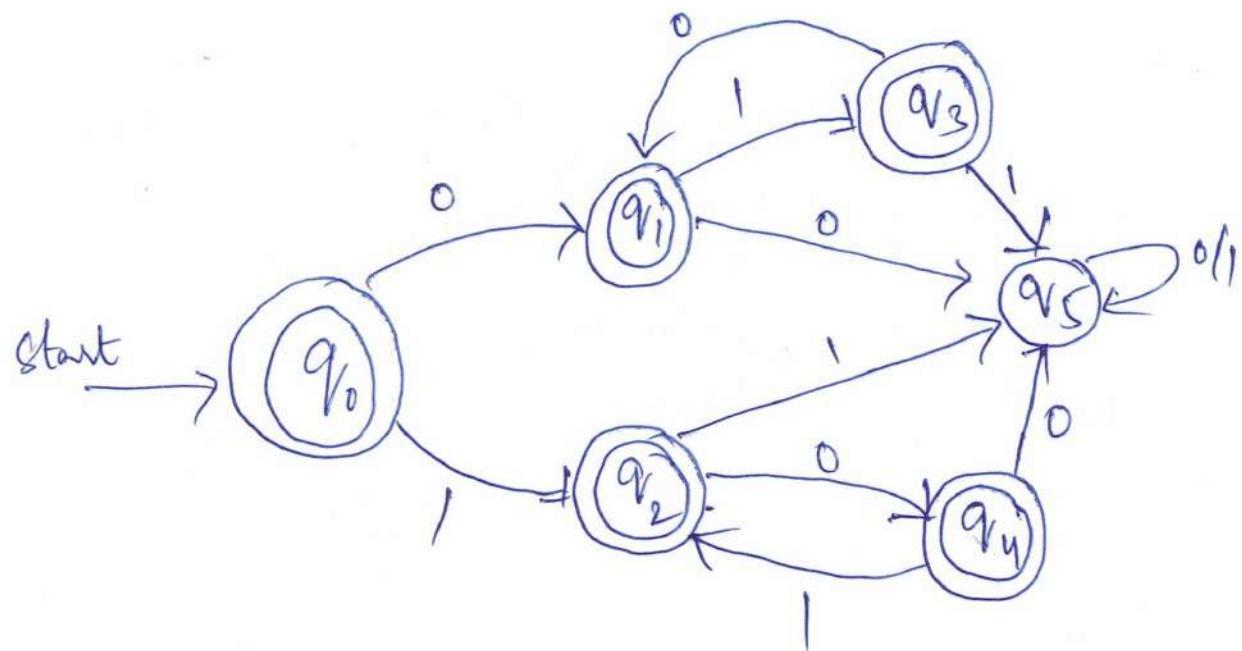
$\epsilon \in L$ ,  $0 \in L$ ,  $1 \in L$ ,  $01 \in L$ ,  $10 \in L$ ,  $010 \in L$ ,  $101 \in L$  --- etc.

$L = \{ \epsilon, 0, 1, 01, 10, 010, 101, 0101, 1010, \dots \}$ .

DFA that accept all binary strings where 0's and 1's are alternate.

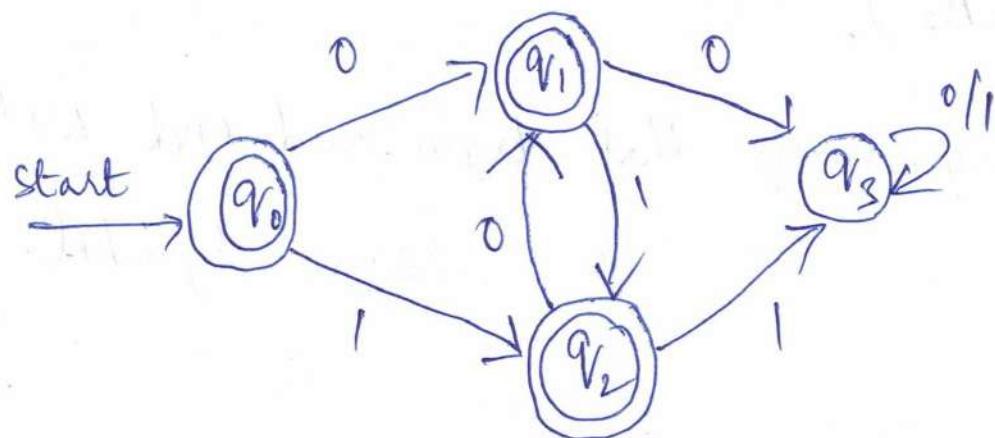
$\epsilon \in L$ ,  $0 \in L$ ,  $1 \in L$ ,  $01 \in L$ ,  $10 \in L$ ,  $010 \in L$ ,  $101 \in L$  --- etc.

$$L = \{ \epsilon, 0, 1, 01, 10, 010, 101, 0101, 1010, \dots \}.$$



$$L = \{ \epsilon, 0, 1, 01, 10, 010, \dots \}$$

DFA that accept all binary strings where 0's and 1's are alternate.



$$\delta: Q \times \Sigma \rightarrow Q \quad q \xrightarrow{a} p$$

$$\forall q \in Q$$

$$\forall a \in \Sigma$$

$$A = (Q, \Sigma, \delta, q_0, F)$$

$$Q = \{q_0, q_1, q_2, q_3\}$$

$$F = \{q_0, q_1, q_2\}$$

$q$	0	1
$q_0$	$q_1$	$q_2$
$q_1$	$q_3$	$q_2$
$q_2$	$q_1$	$q_3$
$q_3$	$q_3$	$q_3$

p

Note:-

- 1) Given a DFA  $\Rightarrow$  Unique  $L(A)$
- 2) Given a regular language  $L \Rightarrow A_1, A_2, \dots, A_n$   
 $L(A_1) = L = L(A_2)$

Example: DFA that accept all binary strings that begin and end with same symbol.

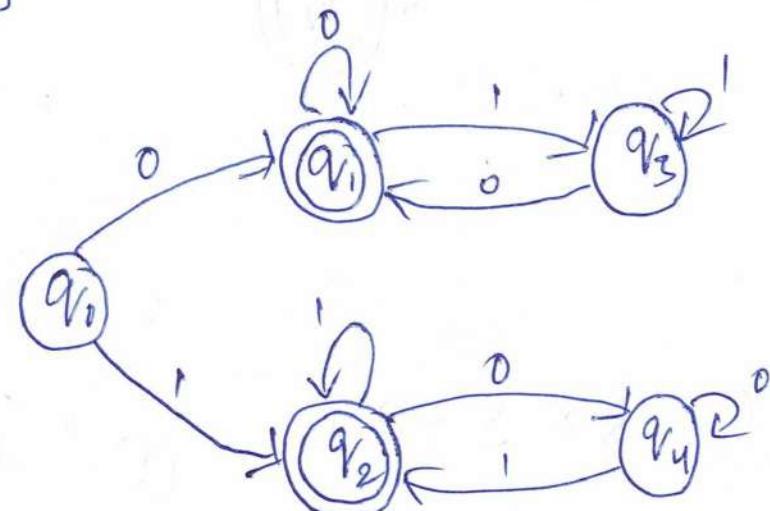
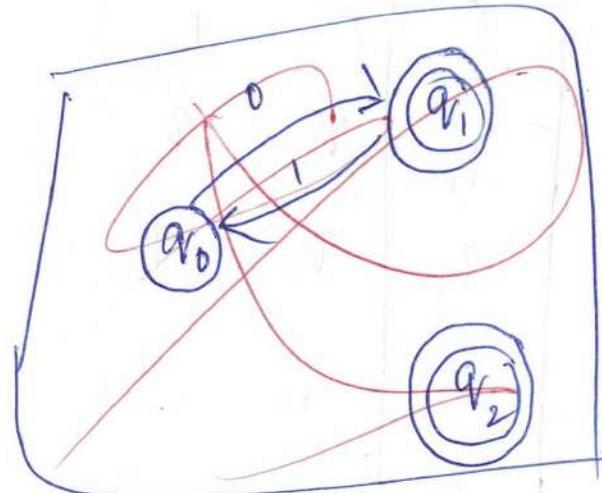
Note:-

- 1) Given a DFA  $\Rightarrow$  Unique  $L(A)$
- 2) Given a regular language  $L \Rightarrow A_1, A_2, \dots, A_n$

$$L(A_1) = L = L(A_2)$$

Example: DFA that accept all binary strings that begin and end with same symbol.

$$L = \{axa / a \in \Sigma, x \in \Sigma^*\}$$



(28) (29)

$$L = \{0^n, 1^n \mid n \geq 0\}$$

Regular ?  $\Rightarrow$  DFA ?

Ans: No

Q

$$L = \{0^n, 1^n \mid n \geq 0\}$$

(28) (29)

Regular ?  $\Rightarrow$  DFA ?

Ans: No

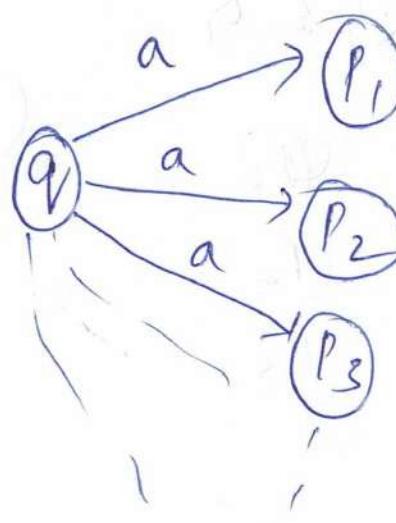
NFA  $\rightarrow$  Non deterministic Finite Automata

$$(Q, \Sigma, \delta, q_0, F)$$



$$\delta(q, a) = P$$

$$\text{DFA: } \delta: Q \times \Sigma \rightarrow Q$$



$$\delta(q, a) = \{p_1, p_2, p_3, \dots, p_k\} \subseteq Q$$

$$\text{NFA} = \delta: Q \times \Sigma \rightarrow 2^Q = \text{pow}(Q)$$

NFA

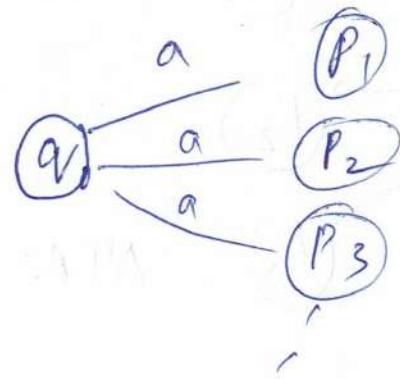
$$N = \{Q, \Sigma, \delta, q_0, F\}$$

$Q \rightarrow$  Set of possible states

Finite

$\Sigma \rightarrow$  Input alphabet

$$\delta(q, a) = \{p_1, p_2, \dots, p_k\} \subseteq Q$$



$q_0 \rightarrow$  starting state

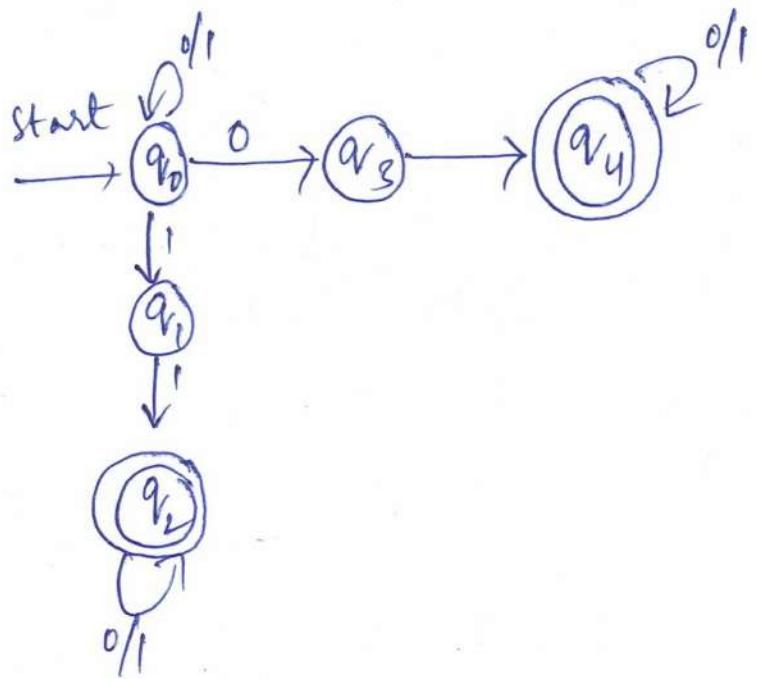
$F \rightarrow$  Final state  $F \subseteq Q$ .

$$\delta: Q \times \Sigma \rightarrow \mathcal{P}(Q)$$

= set of all possible  
subset of  $Q$ .

$\Sigma = \{0, 1\}$  Example of NFA

(3)



$$N = (Q, \Sigma, \delta, q_0, F)$$

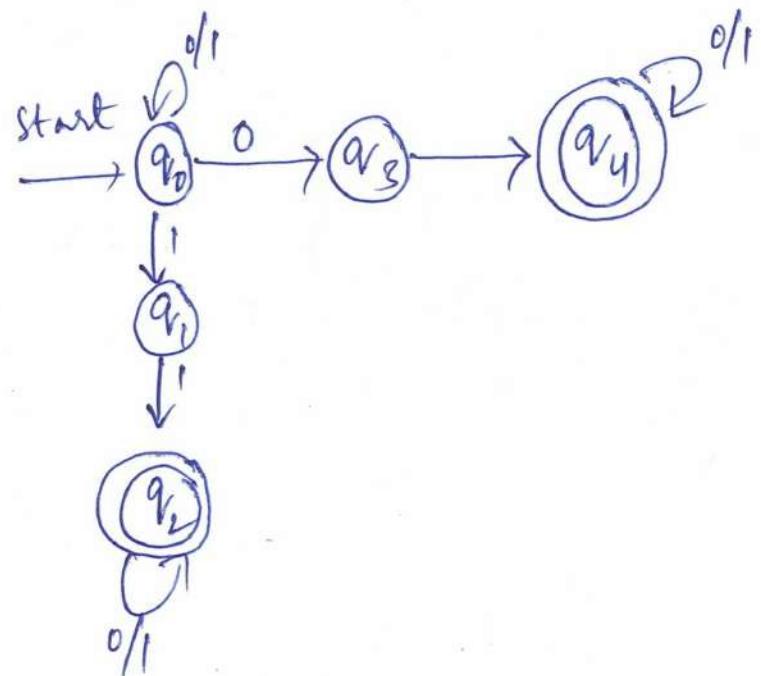
$$Q = \{q_0, q_1, q_2, q_3, q_4\}$$

$$F = \{q_2, q_4\}$$

P

$\Sigma = \{0, 1\}$  Example of NFA

(31)-32



$$N = (Q, \Sigma, \delta, q_0, F)$$

$$Q = \{q_0, q_1, q_2, q_3, q_4\}$$

$$F = \{q_2, q_4\}$$

$\Sigma$	0	1
$q_0$	$\{q_0, q_3\}^*$	$\{q_0, q_1\}^*$
$q_1$	$\emptyset^2$	$\{q_2\}$
$q_2$	$\{q_2\}$	$\{q_2\}$
$q_3$	$\{q_4\}$	$\emptyset^2$
$q_4$	$\{q_4\}$	$\{q_4\}$

- ① Multiple  
② Empty

R

## Building NFA

$$\Sigma = \{0, 1\}$$

$$L \subseteq \Sigma^*$$



$$N \rightarrow \text{NFA} \quad L(N) = L$$

$$L(N) = L$$

→ \* Every DFA is NFA.

Q

## Building NFA

Example:-

$$\Sigma = \{0, 1\}$$

{ NFA accepting all binary strings that end with  
Pattern 101 }

$$L = \{ x101 / x \in \Sigma^* \}$$

$$= \{ 101, 0101, 1101, 00101, \dots \}$$

P

## Building NFA

(34) (35)

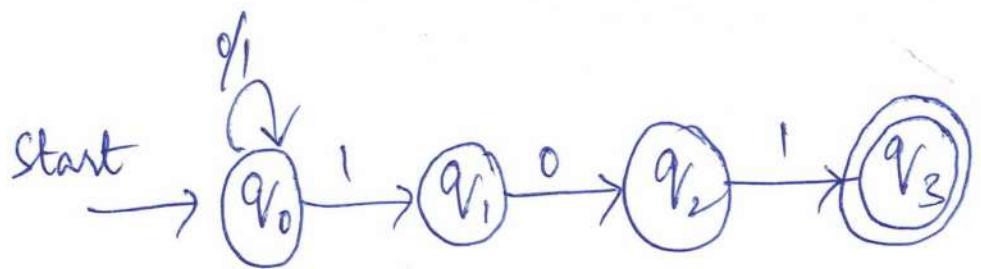
Example:-

$$\Sigma = \{0, 1\}$$

{ NFA accepting all binary strings that end with  
Pattern 101 }

$$L = \{ x101 / x \in \Sigma^* \}$$

$$= \{ 101, 0101, 1101, 00101, \dots \}$$



R

## Building NFA

Example:-  $\Sigma = \{0, 1\}$

Two patterns simultaneously

$$L = 0^* + (01)^*$$

$$L = \left\{ 0, 00, 000, \dots \right. \\ \left. 01, 0101, 010101, \dots \right\}$$

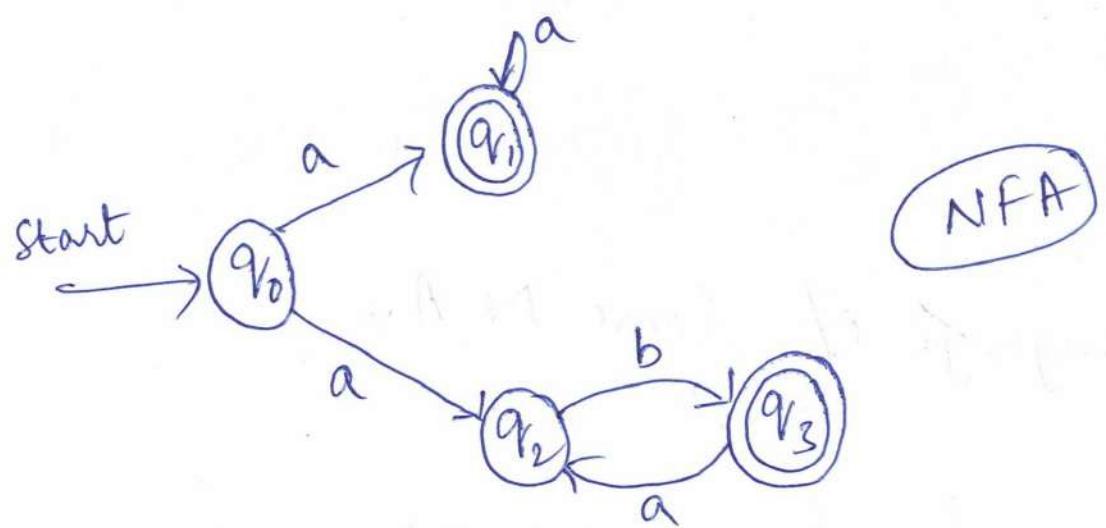
## Building NFA

Example:-  $\Sigma = \{0, 1\}$

Two patterns simultaneously

$$L = 0^* + (01)^*$$

$$L = \left\{ \begin{array}{l} 0, 00, 000, \dots \\ 01, 0101, 010101, \dots \end{array} \right\}$$



## The equivalence of DFA and NFA

- DFA can be treated as NFA



$$(\mathcal{Q}, \Sigma, \delta, q_0, F)$$



$$(\mathcal{Q}, \Sigma, \delta, q_0, F)$$

$$\delta: \mathcal{Q} \times \Sigma \rightarrow \mathcal{P}$$

$$\delta: \mathcal{Q} \times \Sigma \rightarrow 2^{\mathcal{Q}}$$

$$\delta(q, a) = P = \{p\} \in 2^{\mathcal{Q}}$$

→ Language of a NFA

$L(N)$  is also a language of some DFA.



NFA accepts only regular languages.

Regular Language

## NFA to DFA

→ Given NFA  $N = (Q_N, \Sigma, \delta_N, q_0, F_N)$

Design:-  $D = (Q_D, \Sigma, \delta_D, \{q_0\}, F_D)$  such that  $L(N) = L(D)$ .

$$Q_N = \{q_0, q_1, \dots, q_n\}$$

$$Q_D = 2^{Q_N}$$

$$|Q_D| = 2^n$$

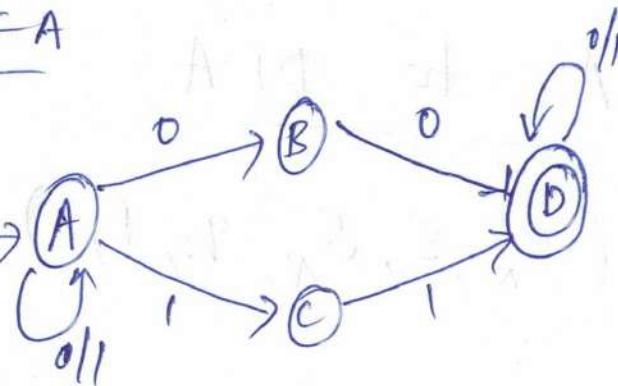
$$\left\{ \begin{array}{l} \{q_0\} \{q_1\} \dots \{q_n\} \\ \{q_0, q_1\} \{q_0, q_2\} \dots \\ \vdots \end{array} \right\}$$

$$Q_D \subseteq 2^{Q_N}$$

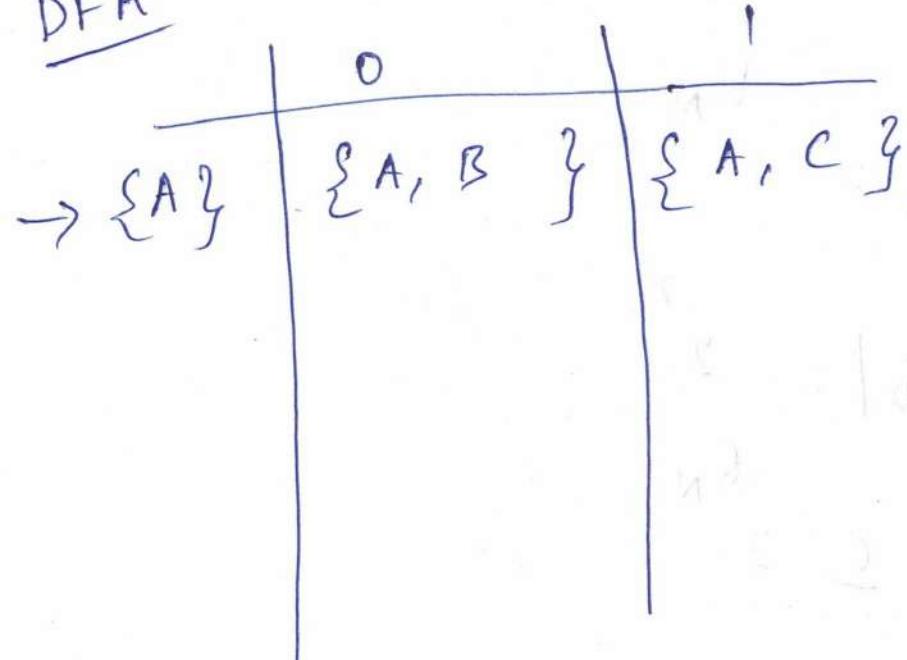
(40)

NFA to DFAExample:-  $Q_N = \{A, B, C, D\}$ 

Start



NFA

DFA

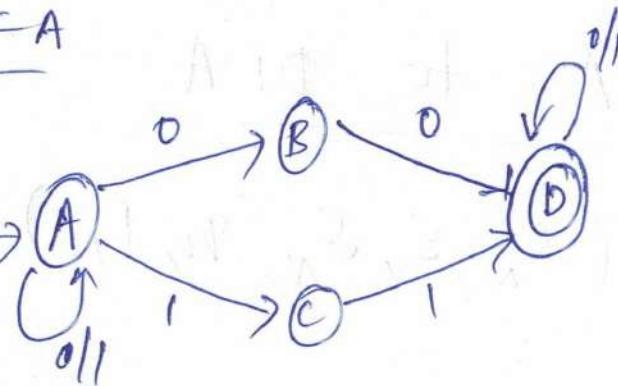
P.

(40)  
41

NFA to DFA

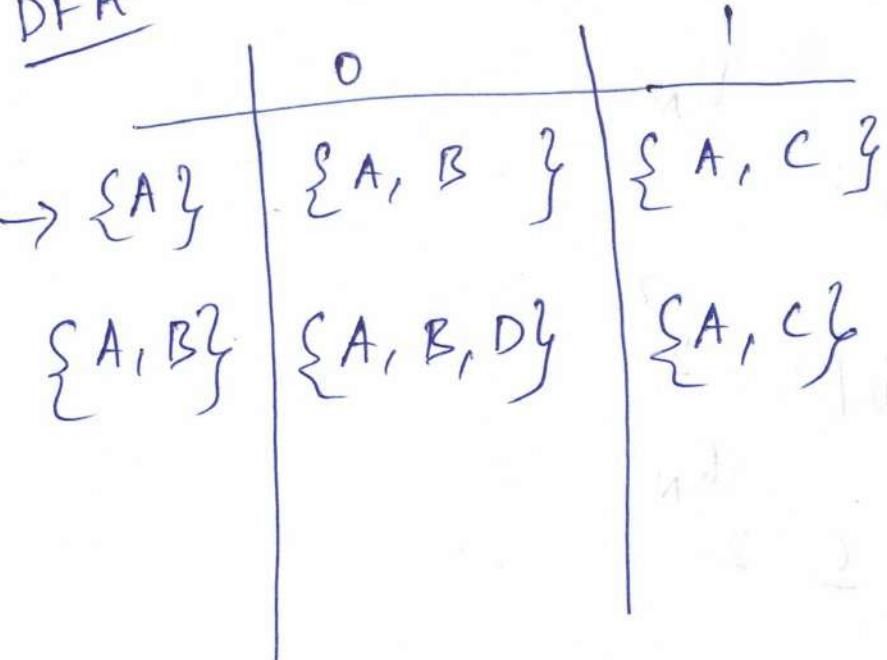
Example:-  $Q_N = \{A, B, C, D\}$

Start



NFA

DFA



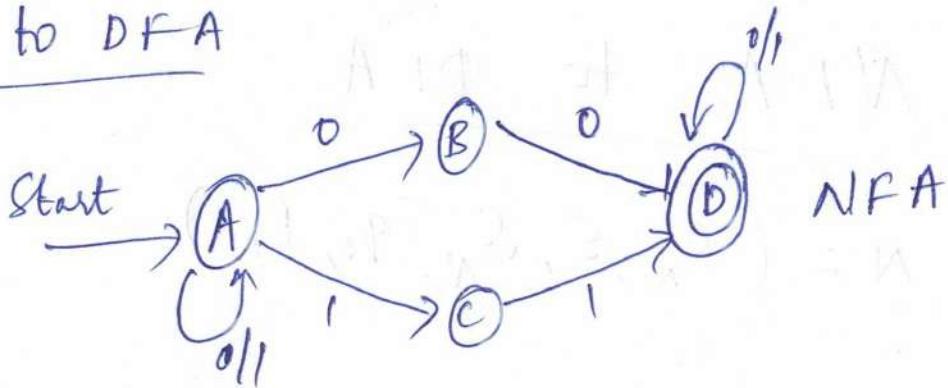
P

(40)  
(41)  
(42)

NFA to DFA

Example:-  $Q_N = \{A, B, C, D\}$

Start



DFA

$\rightarrow \{A\}$	$\{A, B\}$	$\{A, C\}$
$\{A, B\}$	$\{A, B, D\}$	$\{A, C\}$
$\{A, C\}$	$\{A, B\}$	$\{A, C, D\}$

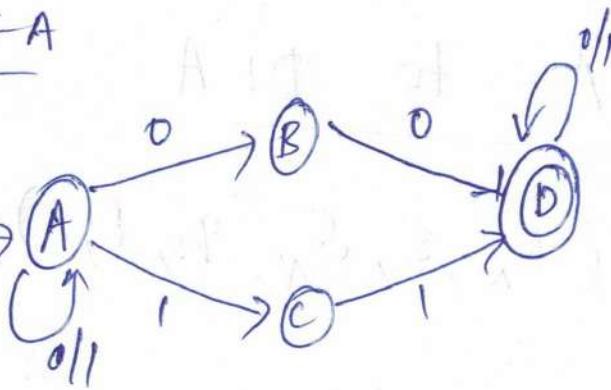
R

40  
41  
M2  
43

NFA to DFA

Example:-  $Q_N = \{A, B, C, D\}$

Start



NFA

DFA

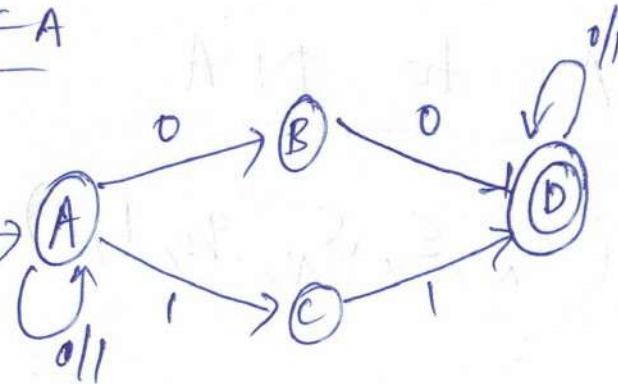
	0	1
$\rightarrow \{A\}$	$\{A, B\}$	$\{A, C\}$
$\{A, B\}$	$\{A, B, D\}$	$\{A, C\}$
$\{A, C\}$	$\{A, B\}$	$\{A, C, D\}$
* $\{A, B, D\}$	$\{A, B, D\}$	$\{A, C, D\}$
$\{A, C, D\}$	$\{A, B, D\}$	$\{A, C, D\}$

P.

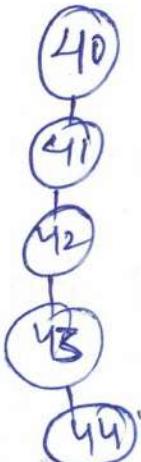
NFA to DFA

Example:-  $Q_N = \{A, B, C, D\}$

Start



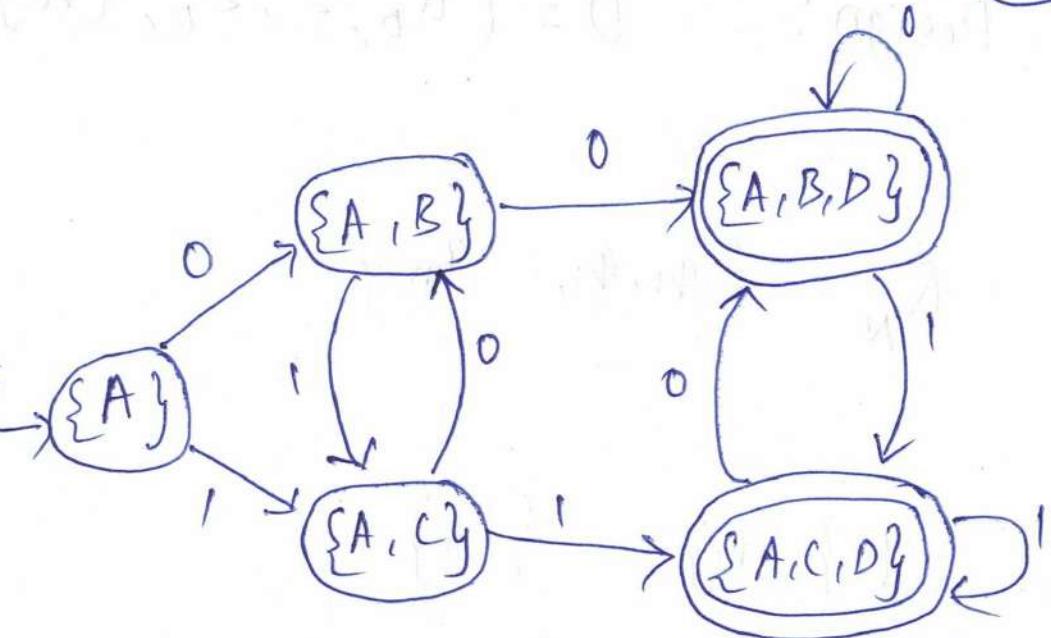
NFA



DFA

	0	1
$\rightarrow \{A\}$	$\{A, B\}$	$\{A, C\}$
$\{A, B\}$	$\{A, B, D\}$	$\{A, C\}$
$\{A, C\}$	$\{A, B\}$	$\{A, C, D\}$
* $\{A, B, D\}$	$\{A, B, D\}$	$\{A, C, D\}$
* $\{A, C, D\}$	$\{A, B, D\}$	$\{A, C, D\}$

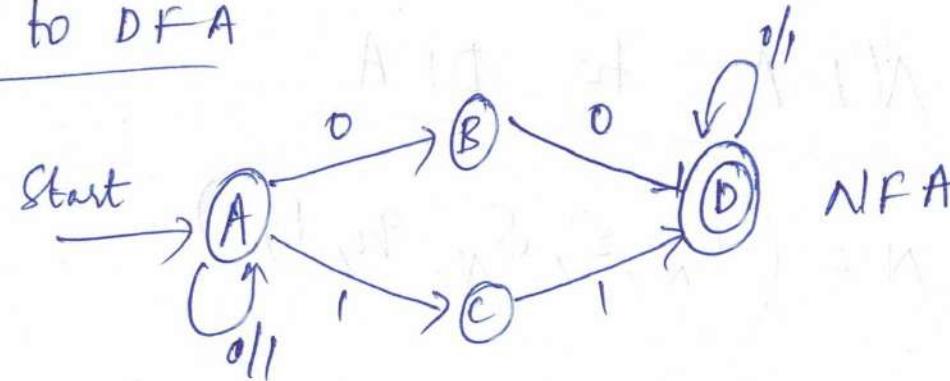
start



P

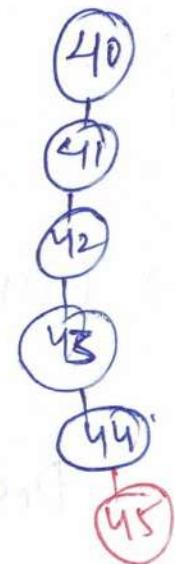
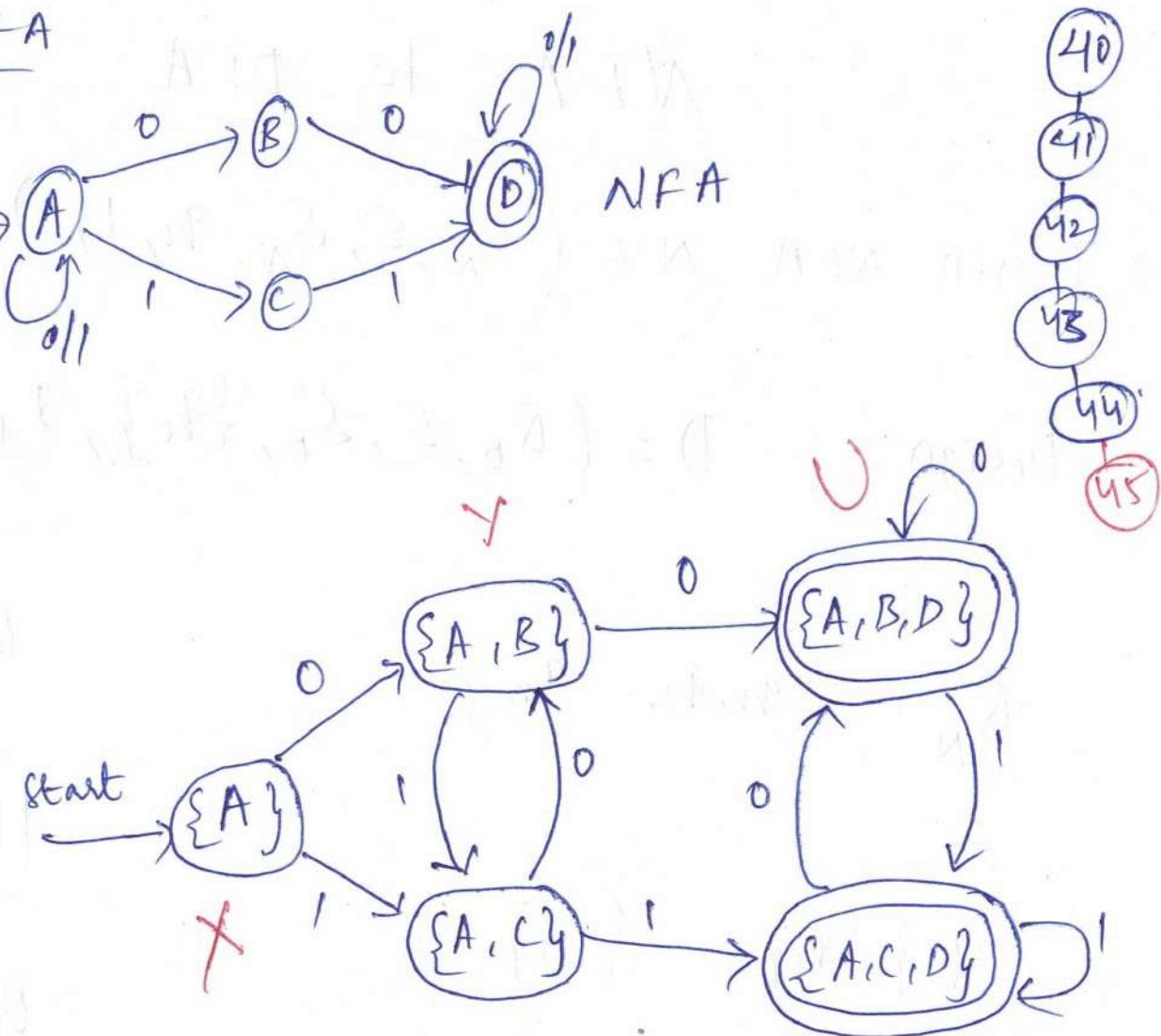
NFA to DFA

Example:-  $Q_N = \{A, B, C, D\}$



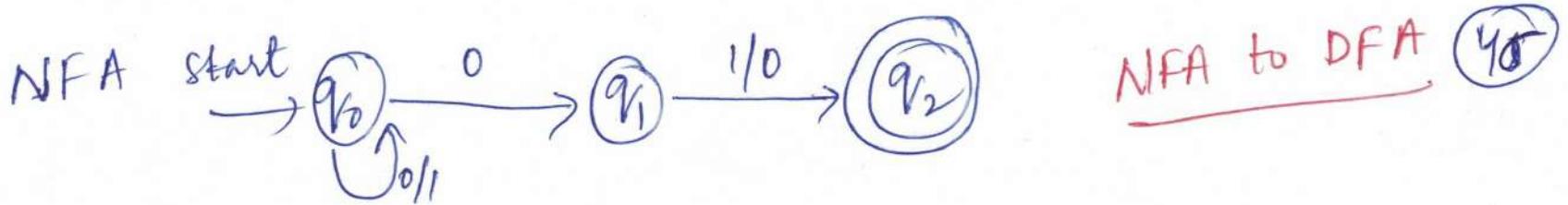
DFA

	0	1
$\rightarrow \{A\}$	$\{A, B\}$	$\{A, C\}$
$\{A, B\}$	$\{A, B, D\}$	$\{A, C\}$
$\{A, C\}$	$\{A, B\}$	$\{A, C, D\}$
* $\{A, B, D\}$	$\{A, B, D\}$	$\{A, C, D\}$
** $\{A, C, D\}$	$\{A, B, D\}$	$\{A, C, D\}$



P

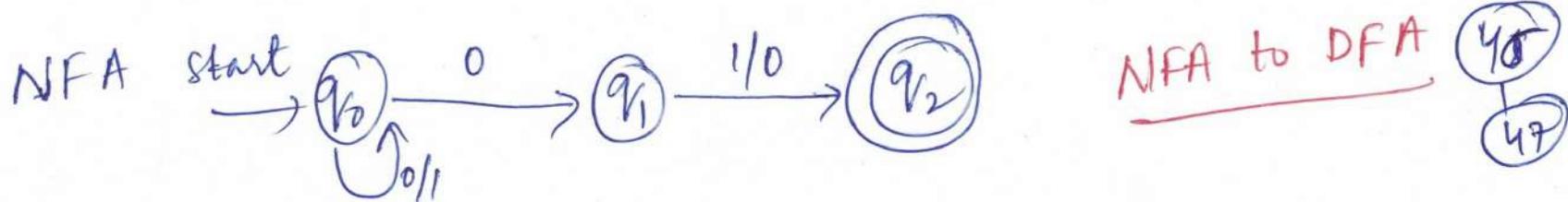
Example:-



$$S_D(s, a) = \bigcup_{p \in s} S_N(p, a)$$

$S_D$	0	1
$\{q_0\}$	$\{q_0\}$	$\{q_0, q_1\}$

Example:



NFA to DFA

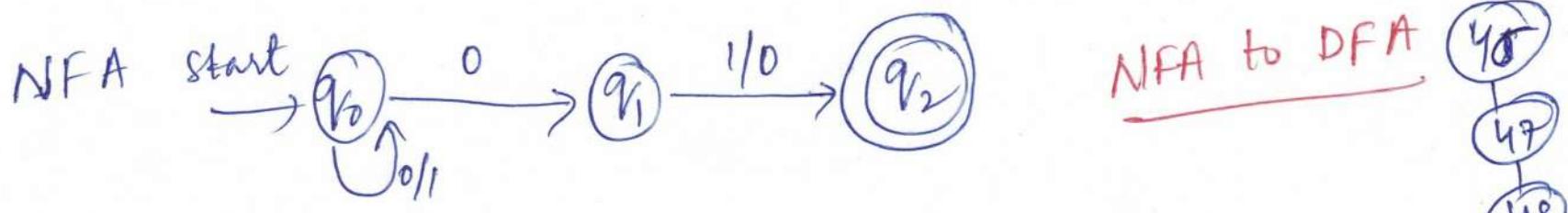


$$S_D(\xi, a) = \bigcup_{p \in \xi} S_N(p, a)$$

$S_D$	0	1
$\{q_0\}$	$\{q_0\}$	$\{q_0, q_1\}$
$\{q_0, q_1\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$

P

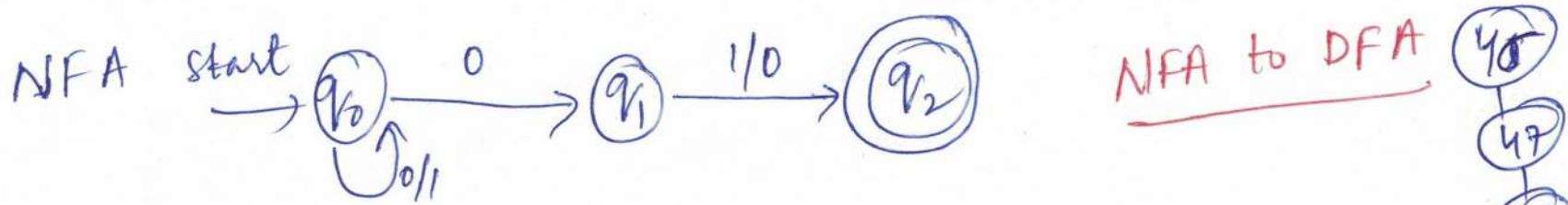
Example:-



$$S_D(\xi, a) = \bigcup_{p \in \xi} S_N(p, a)$$

$S_D$	0	1
$\{q_0\}$	$\{q_0\}$	$\{q_0, q_1\}$
$\{q_0, q_1\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$
$\{q_0, q_2\}$	$\{q_0\}$	$\{q_0, q_1\}$
$\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$

Example:

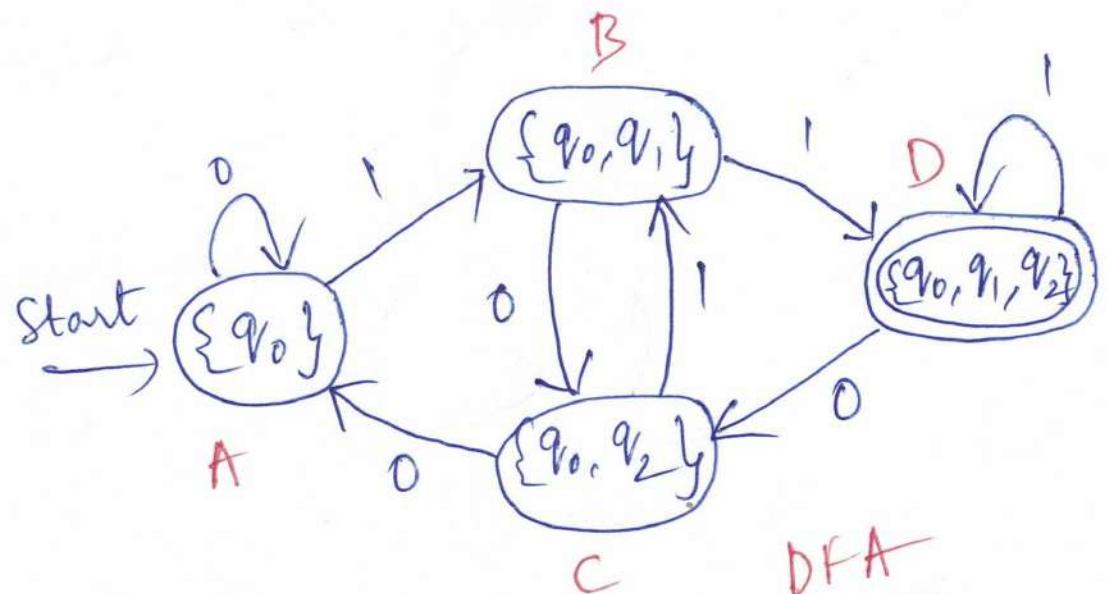


NFA to DFA



$$S_D(S, a) = \bigcup_{p \in S} S_N(p, a)$$

$S_D$	0	1
$\{q_0\}$	$\{q_0\}$	$\{q_0, q_1\}$
$\{q_0, q_1\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$
$\{q_0, q_2\}$	$\{q_0\}$	$\{q_0, q_1\}$
$\{q_0, q_1, q_2\}$	$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$



NFA      DFA  
 $L(N) = L(D)$

DFA ✓

P.

# Finite Automata with $\epsilon$ -moves

( $\epsilon$ -NFA)

$(Q, \Sigma, \delta, q_0, F)$

Transition

$$\boxed{\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow \text{Pow}(Q)}$$

$Q \rightarrow$  Set of all states

$\Sigma \rightarrow$  alphabet

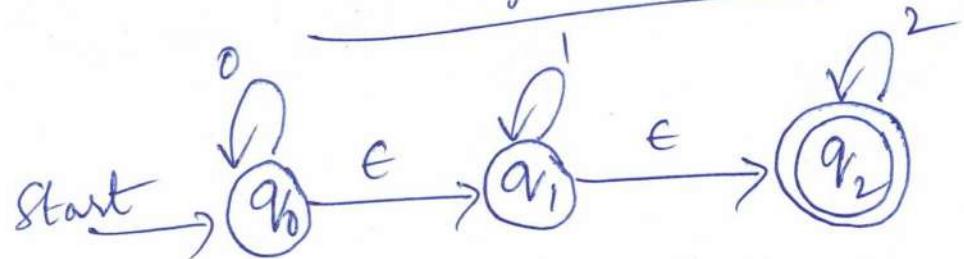
$q_0 \rightarrow$  Initial

$F \rightarrow$  Final / acceptable

P.

(51)

### Example of $\epsilon$ -NFA



Transition function ( $\delta$ )

State	0	1	2	$\epsilon$
$q_0$	$\{q_0\}$	$\emptyset$	$\emptyset$	$\{q_1\}$
$q_1$	$\emptyset$	$\{q_1\}$	$\emptyset$	$\{q_2\}$
$q_2$	$\emptyset$	$\emptyset$	$\{q_2\}$	$\emptyset$

$$Q = \{q_0, q_1, q_2\}$$

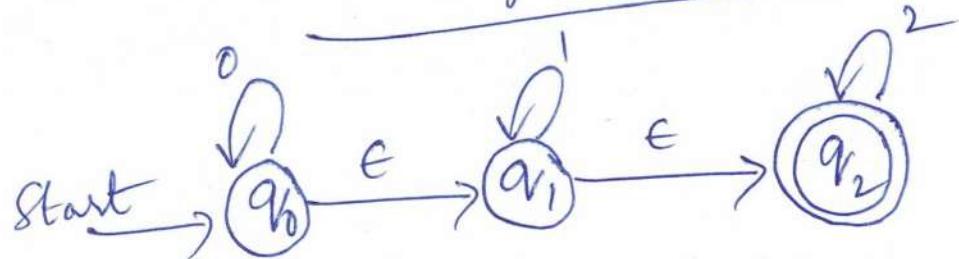
$$\Sigma = \{0, 1, 2\}$$

$$F = \{q_2\}$$

Q

(51)  
(52)

### Example of $\epsilon$ -NFA



Transition function ( $\Sigma$ )

State	0	1	2	$\epsilon$
$q_0$	$\{q_0\}$	$\emptyset$	$\emptyset$	$\{q_1\}$
$q_1$	$\emptyset$	$\{q_1\}$	$\emptyset$	$\{q_2\}$
$q_2$	$\emptyset$	$\emptyset$	$\{q_2\}$	$\emptyset$

$$Q = \{q_0, q_1, q_2\}$$

$$\Sigma = \{0, 1, 2\}$$

$$F = \{q_2\}$$

$$L = \begin{matrix} * & * & * \\ 0 & 1 & 2 \end{matrix}$$

$$0^* = \{\epsilon, 0, 00, 000, \dots\}$$

$$1^* = \{\epsilon, 1, 11, 111, \dots\}$$

$$2^* = \{\epsilon, 2, 22, 222, \dots\}$$

$$L = \{\epsilon, 0, 1, 2, \dots\}$$

Q

53

## $\epsilon$ -NFA to NFA

$$E = (Q, \Sigma, \delta, q_0, F)$$

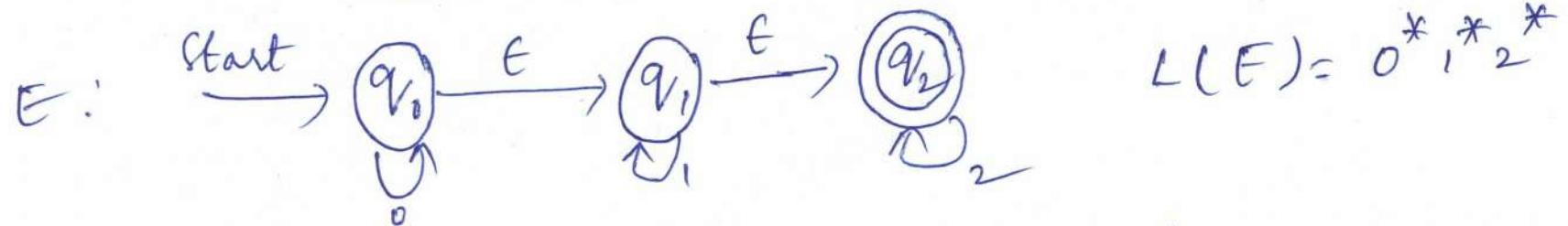
$$N = (Q, \Sigma, \delta', q_0, F')$$

$$\delta'(q, a) = \overline{\delta}(q, a)$$

$$F' = \begin{cases} F \cup \{q_0\}, & \text{if } \epsilon\text{-closure}(q_0) \cap F \neq \emptyset \\ F & \text{otherwise} \end{cases}$$

P

54

 $\epsilon$ -NFA to NFA

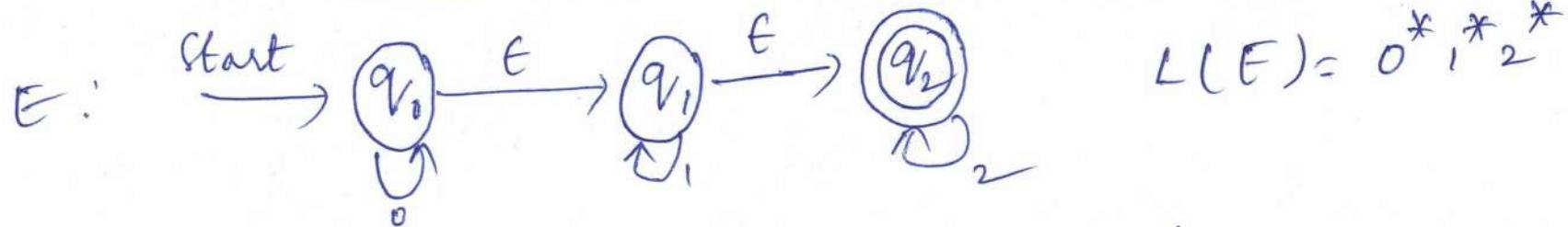
$$N = (Q, \Sigma, \delta', q_0, F') \quad \delta'(q_0, 0) = \hat{\delta}(q_0, \epsilon 0 \epsilon) = ?$$

$\delta'$	0	1	2
$\rightarrow q_0$			
$q_1$			
$* q_2$			

P

54  
CP

### $\epsilon$ -NFA to NFA



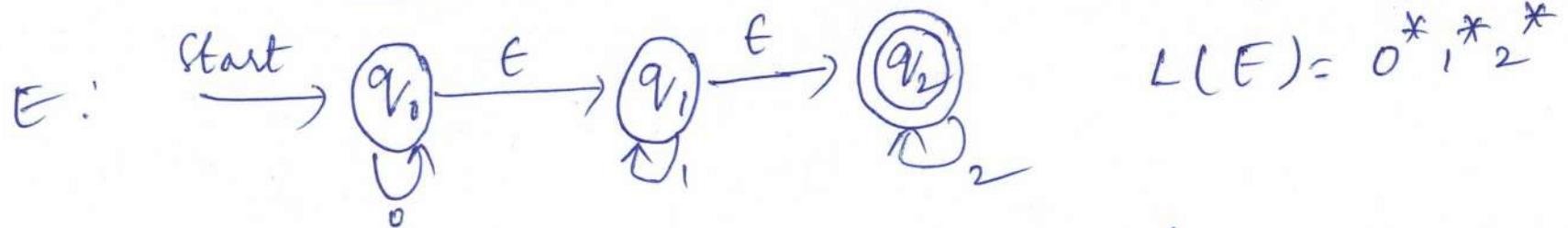
$$N = (Q, \Sigma, \delta', q_0, F')$$

$$\delta'(q_0, 0) = \hat{\delta}(q_0, \epsilon \circ \epsilon) = ? (q_0, q_1, q_2)$$

$\delta'$	0	1	2
$\rightarrow q_0$	$\{q_0, q_1, q_2\}$	$\{q_1, q_2\}$	$\{q_2\}$
$q_1$			
$* q_2$			

P

$\epsilon$ -NFA to NFA



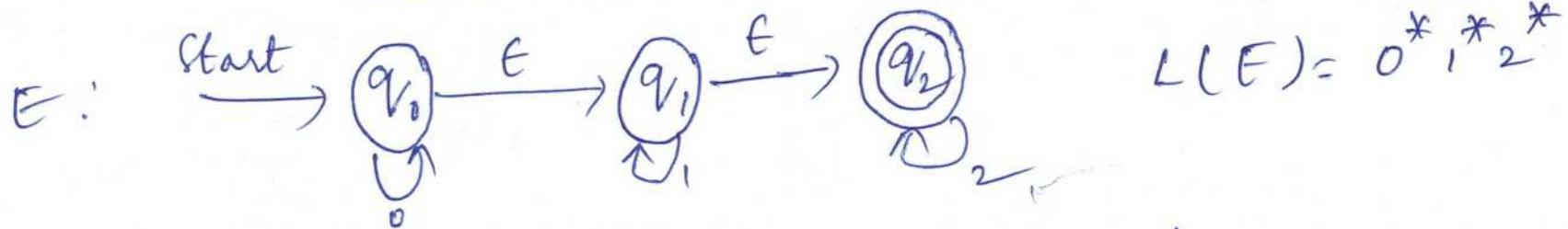
$$N = (Q, \Sigma, \delta', q_0, F') \quad \delta'(q_0, 0) = \hat{\delta}(q_0, \epsilon \circ \epsilon) = ? (q_0, q_1, q_2)$$

$\delta'$	0	1	2
$\rightarrow q_0$	$\{q_0, q_1, q_2\}$	$\{q_1, q_2\}$	$\{q_2\}$
$q_1$	$\emptyset$	$\{q_1, q_2\}$	$\{q_2\}$
$* q_2$	$\emptyset$	$\emptyset$	$\{q_2\}$

54  
55  
56

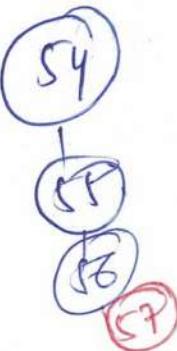
P

# E-NFA to NFA

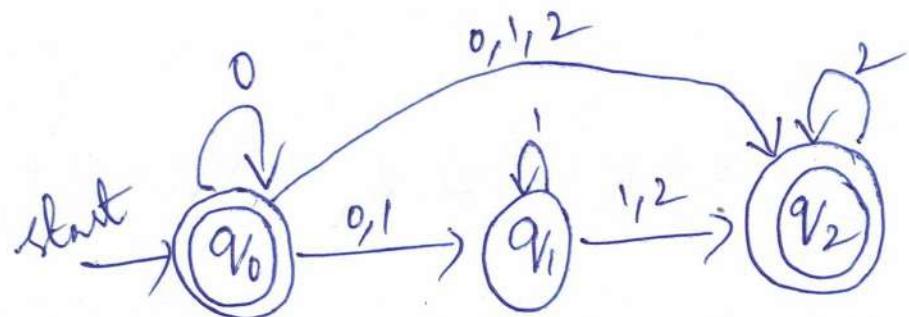


$$N = (Q, \Sigma, \delta', q_0, F')$$

$$\delta'(q_0, 0) = \hat{\delta}(q_0, \epsilon \circ \epsilon) = ? (q_0, q_1, q_2)$$



$\delta'$	0	1	2
$\rightarrow q_0$	$\{q_0, q_1, q_2\}$	$\{q_1, q_2\}$	$\{q_2\}$
$q_1$	$\emptyset$	$\{q_1, q_2\}$	$\{q_2\}$
$* q_2$	$\emptyset$	$\emptyset$	$\{q_2\}$



NFA

P

## Regular Expression

$\Sigma$  = Set of alphabet =  $\{0, 1\}$  or  $\{a, b, c, \dots, z\}$

$\Sigma^*$  = Set of all possible strings over  $\Sigma$

3 operations -  $(\cdot, +, *)$

$L_1, L_2 \rightarrow$  2 languages

$L = L_1 \cdot L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}$

For eg:-  $L_1 = \{10, 1\}, L_2 = \{011, 11\}$

Regular Expression

$\Sigma$  = Set of alphabet =  $\{0, 1\}$  or  $\{a, b, c, \dots, z\}$

$\Sigma^*$  = Set of all possible strings over  $\Sigma$

3 operations -  $(\cdot, +, *)$

$L_1, L_2 \rightarrow$  2 languages

$L = L_1 \cdot L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}$

Exg:-  $L_1 = \{10, 1\}, L_2 = \{011, 11\}$

$L = L_1 \cdot L_2 = \{10011, 1011, 1011, 111\}$ .

P

(60)

$$L^* = \bigcup_{i=0}^{\infty} L^i$$

↓

Kleene closure

$$L^+ = L^*/\{\epsilon\}$$

$$= \bigcup_{i=1}^{\infty} L^i$$

$$= L \cup L^r \cup L^3$$

positive closure

Eg:-  $L = \{10, 11\}$

$$L^* = \bigcup_{i=0}^{\infty} L^i$$

$$L^0 = \{\epsilon\}, \quad L^1 = L = \{10, 11\}, \quad L^r = L \cdot L = \{1010, 1011, 1110, 1111\}$$

$$L^3 = L \cdot L^r =$$

etc.

$L^+$  will not have  $L^0$

Q

## Regular Expression

$r, s$ , concatenation ( $\cdot$ ).

$r, s \rightarrow$  regular expressions

$r+s, r.s, r^*$  will also be regular expressions.

Denotations related to regular expressions (RE)

(1)  $\phi$  is a RE and its language  $L(\phi)$  is empty set.

(2)  $\epsilon$  is a RE  $L(\epsilon) = \{\epsilon\}$ .

(3)  $a \in \Sigma$  is a RE  $L(a) = \{a\}$

(4) If 'i' and 'S' are RE, then  $i+S, iS, i^*$  are RE

$$S = \{0, 1\} \quad 0+1, 01, 0^*$$

(62)

$$\lambda \rightarrow R = L(\lambda)$$

$$\varsigma \rightarrow S = L(\varsigma)$$

$$(i) \quad \lambda + \varsigma = R \cup S = L(\lambda) \cup L(\varsigma) = L(\lambda + \varsigma)$$

$$(ii) \quad L(\lambda\varsigma) = L(\lambda)L(\varsigma) = RS$$

$$(iii) \quad L(\lambda^*) = (L(\lambda))^* = R^* = \bigcup_{i=0}^{\infty} R^i$$

Eg:-  $\Sigma = \{0, 1\}$      $\lambda = \emptyset$  and  $\varsigma = 1$

$$L(\emptyset + 1) = L(\emptyset) \cup L(1) = \{0, 1\} = \Sigma \checkmark$$

$$L(01) = L(0) \cdot L(1) = \{01\} \checkmark$$

$$0^* = (L(0))^* = \{\epsilon, 0, 00, 000, \dots\}$$

P

(63)

$$(0+1)^* = \Sigma^*$$
$$= \{ \epsilon, 0, 1, 00, 01, 10, 11, \dots \}$$

$$\therefore L((0+1)^*) = \overline{\Sigma^*}$$

---

(63) (64)

$$(0+1)^* = \Sigma^* \\ = \{ \epsilon, 0, 1, 00, 01, 10, 11, \dots \}.$$

$$\therefore L((0+1)^*) = \Sigma^*$$

$$\Sigma = \{a, b\}$$

$$L = (a+b)^* aab (a+b)^*$$

$L(L)$  = All binary strings containing  $aab$  as substring

$$= \{aab, aaab, aaab, \dots\}.$$

Q.

(65)

$$\Sigma = \{0, 1\}$$

$$h = (1+01)$$

$$h \xrightarrow{*} ?$$

Q

(65) (66)

$$\Sigma = \{0, 1\}$$

$$L = (1+01)$$

$$L^* \rightarrow ?$$

$$L(x) = x_1 + x_2 = 1+01 = \{1, 01\}$$

$$L^* = (L(x))^* = \{1, 01\}^*$$

$$= \{e, 1, 01, 101, 011, \dots\}$$

Q

## Precedence of the operators

(67)

+, Concatenation ( $\cdot$ ),  $*$

$*$  > Concatenation > +

Eg!  $((0(1^*)) + 0) = 01^* + 1$

$$= \{0, 1, 01, 011, 0111, \dots\}$$

P

## Precedence of the operators

+, Concatenation( $\cdot$ ), \*

\* > Concatenation > +

$$\text{Eg} \quad ((0(1^*)) + 0) = 01^* + 1$$

$$= \{0, 1, 01, 011, 0111, \dots\}$$

Eg<sub>2</sub>: -  $0^* 1^* 2^*$   $\rightarrow$  Set of all strings with any 0's followed by any 1's followed by any 2's.

## Algebraic laws of the RE operators

(1) Commutativity of '+':  $r+s = s+r$

$$\text{Eg } 0+1 = 1+0$$

(2) Associativity of '+':  $r, s, t$

$$(r+s)+t = r+(s+t)$$

(3) Associativity of concatenation (''):  $(rs)t = r(st)$

$$\{xyz / x \in \Sigma, y \in \Sigma, z \in \Sigma\}$$

(70)

(4) Is Concatenation Commutative?

Eg.  $\lambda = 0$  and  $S = 1$ 

$$L(\lambda S) = 01 = \{xy \mid x \in \lambda, y \in S\}.$$

$$L(S\lambda) = 10 = \{xy \mid x \in S, y \in \lambda\}.$$

$L(\lambda S) \neq L(S\lambda)$  *∴ Not Commutative.*

$$(5) 0 + \lambda = \lambda + \emptyset = \lambda$$

$$L(\emptyset) \cup L(\lambda) = \emptyset \cup L(\lambda) = L(\lambda)$$

$$(6) \epsilon \lambda = \lambda \epsilon = \lambda$$

$$L(\epsilon \lambda) = L(\lambda \epsilon) = L(\lambda).$$

(P)

## Distributive laws

(71)

$$\lambda(s+t) = \lambda s + \lambda t$$

$$= \{ xy \mid x \in L(s), y \in L(s+t) \} \\ \quad \quad \quad y \in L(s) \text{ or } y \in L(t) \}$$

$$= \{ xt_1 \mid x \in L(s), t_1 \in L(s) \}$$

$$\cup \{ xt_2 \mid x \in L(s), t_2 \in L(t) \}$$

$$= L(\lambda s) \cup L(\lambda t)$$

Q

## distributive laws (Contd...)

$$\underline{(r+s)t} = rt + st$$

$$L(r) \cup L(s)$$

$$= \{x \in \{x \in L(r) \text{ or } x \in L(s) \text{ and } z \in L(t)\}\}$$

$$= \{x \in \{x \in L(r) \text{ and } z \in L(t)\}\}$$

$$\cup \{x \in \{x \in L(s) \text{ and } z \in L(t)\}\}$$

$$= rt + st$$

Q

Idempotent law :-

$$\lambda + \lambda = \lambda$$

$$= L(\lambda) \cup L(\lambda) = L(\lambda)$$

Laws on \* :-

$$(\lambda^*)^* = \lambda^*$$

$$\text{Eg:- } (0^*)^* = 0^*$$

$$\text{Eg2:- } ((01)^*)^* = (01)^*$$

$$\epsilon^* = \epsilon$$

$$L(\epsilon^*) = \bigcup_{i=0}^{\infty} (L(\epsilon))^i \\ = \epsilon$$

$$\phi^* = \epsilon$$

$$L(\phi) = \phi$$

$$L(\phi^*) = \bigcup_{i=0}^{\infty} (L(\phi))^i$$

$$= \epsilon \cup \phi \cup \phi \dots = \epsilon$$

Q.

$$(1) (\lambda + \Sigma)^* = (\lambda^* \Sigma^*)^*$$

$$(2) \lambda^* = \lambda^* \lambda^*$$

$$(3) \lambda + \Sigma \lambda \neq (\lambda + \Sigma) \cdot \lambda$$

Eg1:-  $(0+1)^* = (0^*, 1^*)^* = \Sigma^*$        $\Sigma = \{0, 1\}$

All Binary strings.

Eg2:-  $0^* = 0^* \cdot 0^*$

Eg3:-  $0+10 \neq (0+1) \cdot 0$

$$\textcircled{1} \rightarrow 0 \cdot 1^* + 2 = \lambda$$

$$L(\lambda) =$$

$$\textcircled{2} \rightarrow (0^*, 1^*)^* = \lambda$$

$$L(\lambda) =$$

$$\textcircled{3} \rightarrow (0^*, 1^* + 0^*, 1^*)^* = \lambda$$

$$L(\lambda) =$$

$$\textcircled{4} \rightarrow (0+1)^* \cdot (0+1)^* = \lambda$$

$$L(\lambda) =$$

$$\textcircled{5} \rightarrow (0+1) \cdot 001 = \lambda$$

$$L(\lambda) =$$

$$\textcircled{1} \rightarrow 0 \cdot 1^* + 2 = \lambda$$

(75) (76)

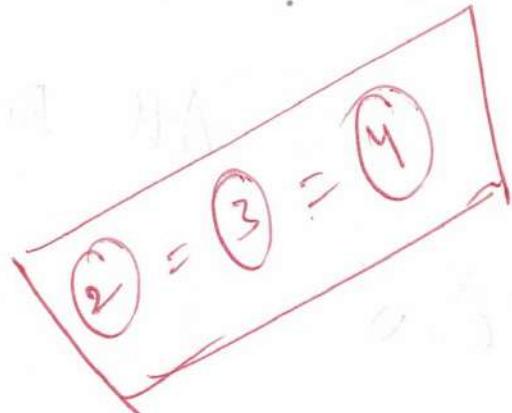
$$L(\lambda) = \{0, 2, 01, 011, 0111, \dots\}$$

$$\textcircled{2} \rightarrow (0^*, 1^*)^* = \lambda$$

$$L(\lambda) = \text{Set of all binary strings} = \{\epsilon, 0, 1, 00, 01, 10, 11, \dots\}$$

$$\textcircled{3} \rightarrow (0^*, 1^* + 0^*, 1^*)^* = \lambda$$

$$L(\lambda) = \text{Same as } \textcircled{2}$$



$$\textcircled{4} \rightarrow (0+1)^* \cdot (0+1)^* = \lambda$$

$$L(\lambda) = \text{Same as } \textcircled{2}$$

$$\textcircled{5} \rightarrow (0+1) \cdot 001 = \lambda$$

$$L(\lambda) = \{0001, 1001\}$$

Q.

## Equivalent between Regular Expression and $\epsilon$ -NFA

(77)

$$r \rightarrow L(r)$$

$$\epsilon\text{-NFA} \rightarrow E$$

$$L(E) = L(r)$$

Theorem:- Let  $r$  be a RE, then  $\exists \epsilon\text{-NFA}, E$  that accept  $L(r)$

$$L(E) = L(r)$$

Proof :-

Q

## Equivalent between Regular Expression and $\epsilon$ -NFA

77

78

$$r \rightarrow L(r)$$

$$\epsilon\text{-NFA} \rightarrow E$$

$$L(E) = L(r)$$

Theorem:- Let  $r$  be a RE, then  $\exists \epsilon\text{-NFA}, E$  that accept  $L(r)$

$$L(E) = L(r)$$

Proof:- By Induction on all operators in  $r$ , we can show that there is a  $\epsilon\text{-NFA}$   $E$ , having one final state and no transition out from this final state s.t.  $L(r) = L(E)$

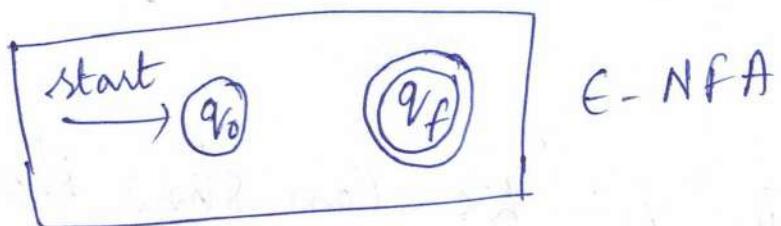
Q

## Base (Zero operators)

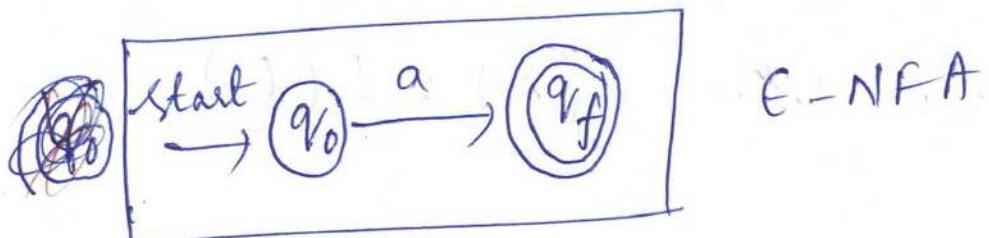
1)  $\lambda = \epsilon \Rightarrow L(\lambda) = \{\epsilon\}$



2)  $\lambda = \emptyset \Rightarrow L(\lambda) = \emptyset$



3)  $\lambda = a, a \in \Sigma \Rightarrow L(\lambda) = \{a\}$



P.

## Induction (one or more operators)

(80)

I.H

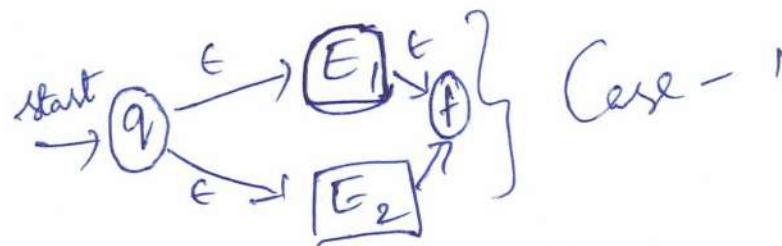
Assume that the theorem is true for RE with fewer than  
i operators,  $i \geq 1$

$\lambda = i$  operator

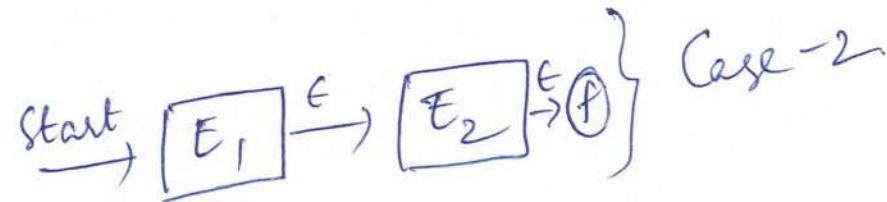
$$L(\lambda_1) = L(E_1)$$

$$L(\lambda_2) = L(E_2)$$

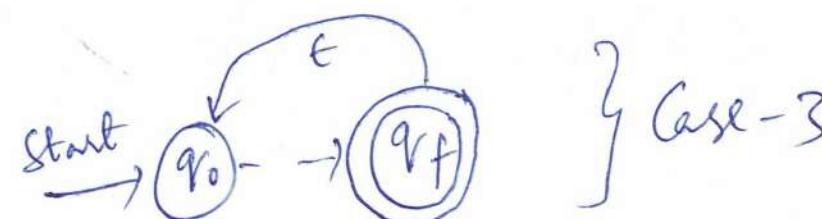
Case-1  $\lambda = \lambda_1 + \lambda_2$



Case-2  $\lambda = \lambda_1 \cdot \lambda_2$



Case-3  $\lambda = \lambda_1^*$



$$E_1 = \text{start} \rightarrow q_0 \dots \rightarrow q_f$$

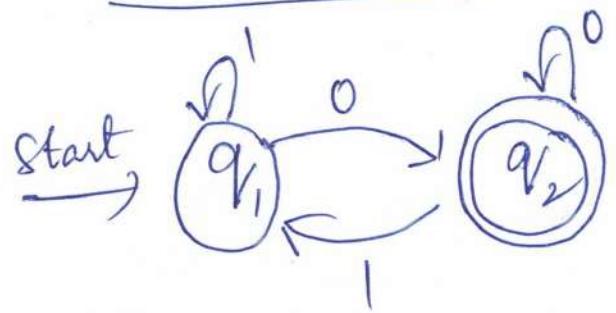
∴ Every RE will have  $\epsilon$ -NFA

Q

(81)

## DFA to Regular Expression

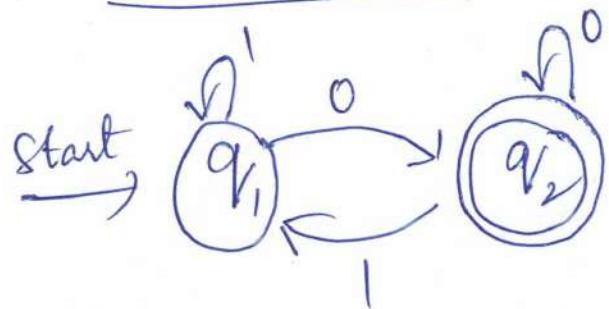
Eg:-



P

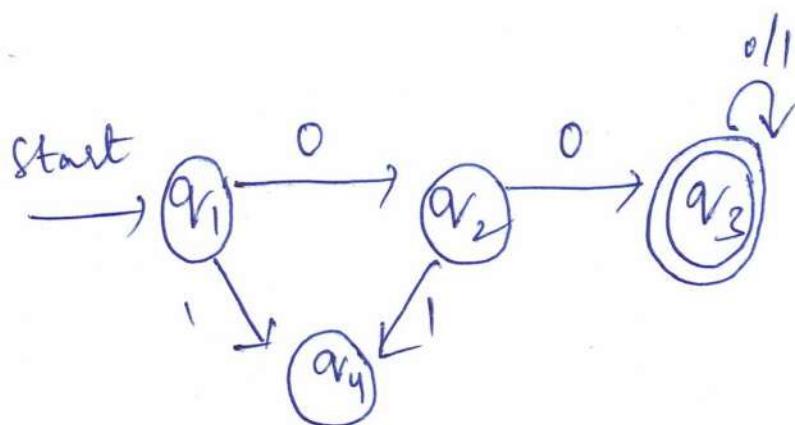
(81) (82)

## DFA to Regular Expression

Eg 1:-

$$L(\lambda) = \{x_0 \mid x \in \Sigma^*\}$$

$$\lambda = (0+1)^* 0$$

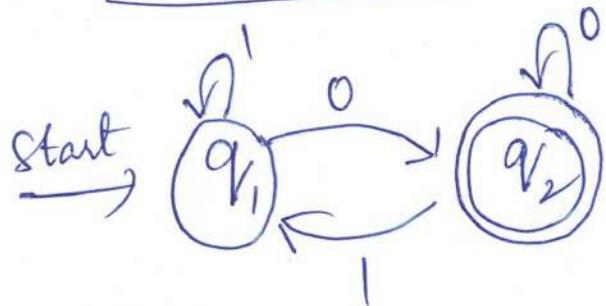
Eg 2:-

P

## DFA to Regular Expression

(81) (82)  
(83)

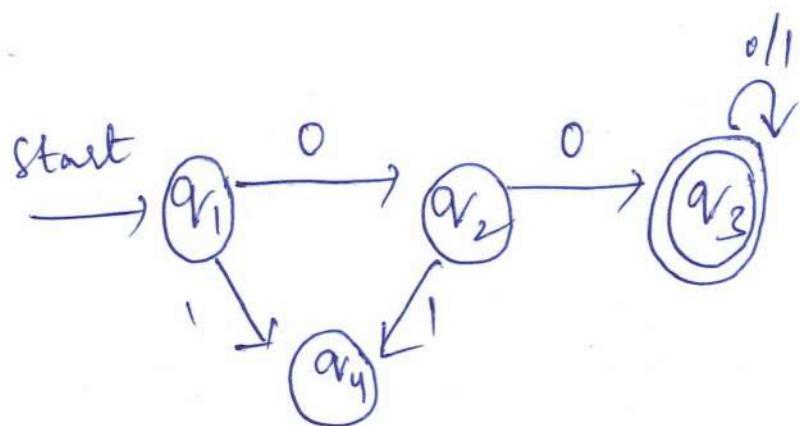
Eg 1:-



$$L(\lambda) = \{x_0 \mid x \in \Sigma^*\}$$

$$\lambda = (0+1)^* 0$$

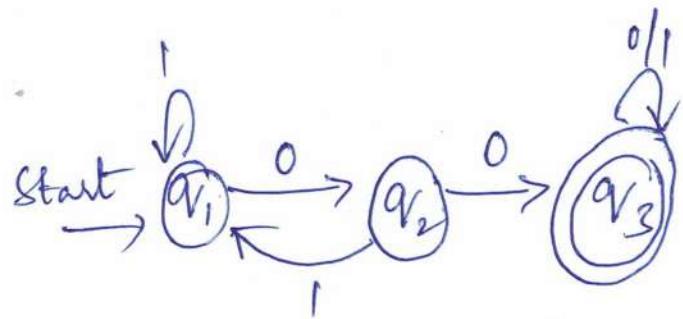
Eg 2:-



$$L^{(x)} = \{00x \mid x \in \Sigma^*\}$$

$$x = 00 (0+1)^*$$

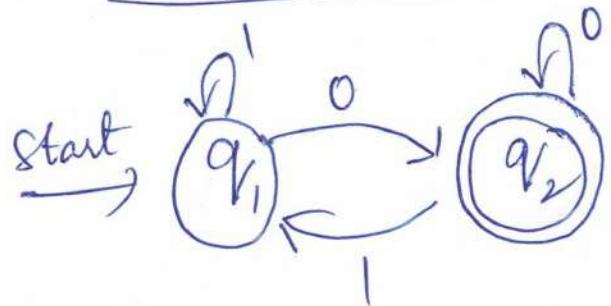
Eg 3:-



P

## DFA to Regular Expression

Eg 1:-

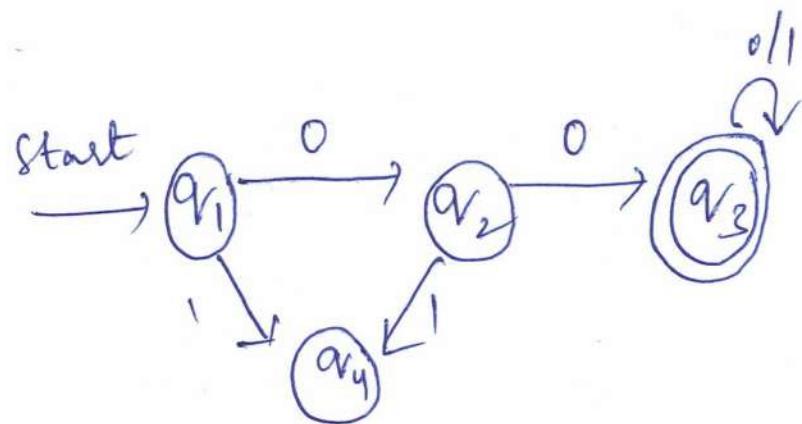


$$L(\lambda) = \{x_0 \mid x \in \Sigma^*\}$$

$$\lambda = (0+1)^* 0$$

(81)-(82)  
(83)-(84)

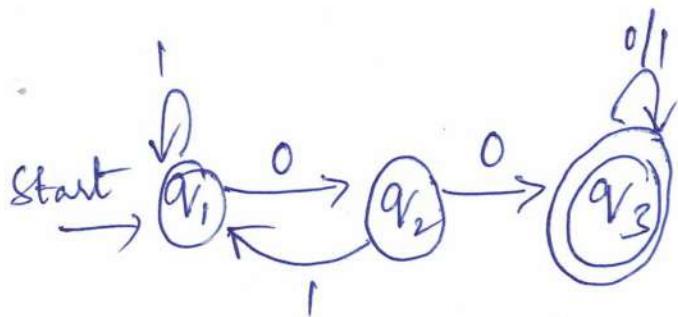
Eg 2:-



$$L^{(1)} = \{00x \mid x \in \Sigma^*\}$$

$$\lambda = 00 (0+1)^*$$

Eg 3:-



$$L(\lambda) = \{x00x \mid x \in \Sigma^*\}$$

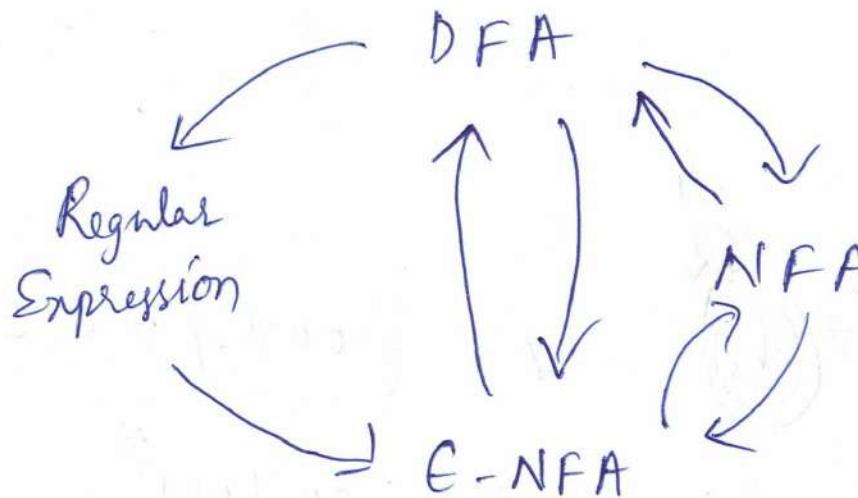
$$\lambda = (0+1)^* 00 (0+1)^*$$

P

# Summary      Diagram

(85)

Regular Language



R

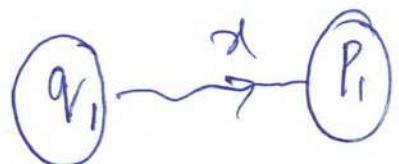
## Minimization of FA

$(Q, \Sigma, \delta, q_0, F)$  DFA

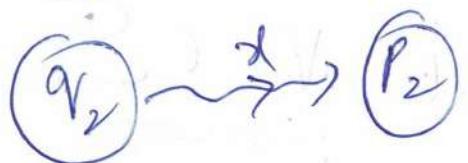
$q_1, q_2$

$q_1 \equiv q_2$

if  $\overset{\uparrow}{\delta}(q_1, a)$  and  $\overset{\uparrow}{\delta}(q_2, a)$  are both  $\in F$   
or both  $\notin F$



either  $p_1, p_2 \in F$  (final)



or  
 $p_1, p_2 \notin F$  (final).

Q.

(1)  $P \equiv P$

Equivivalence (reflexive, symmetric, transitive) (87)

(2)  $P \equiv q \Leftrightarrow q \equiv P$

(3)  $P \equiv q$  and  $q \equiv r \Rightarrow P \equiv r$

Recursive Construction of  $(k+1)^{th}$  equivalence.

(i)  $q_1, q_2$   $(k+1)^{th}$  equivalent  $\Rightarrow$   $q_1, q_2$  must be  $k$ -equivalent

(ii)  $q_1, q_2$   $(k+1)^{th}$  equivalent if

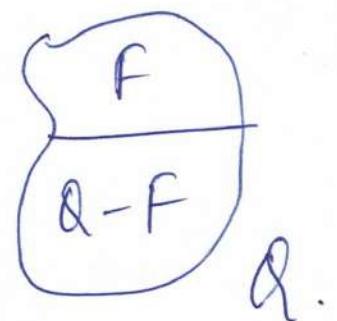
$\delta(q_1, a), \delta(q_2, a)$  are  $k$ -equivalent  $\forall a \in \Sigma$

## Construction of Minimum DFA

(88)

1.  $\Pi_0 \rightarrow 0^{\text{th}}$  equivalent class

$$\Pi_0 = (\mathbb{Q}^0, \mathcal{Q}_2) \text{ i.e. } \mathbb{Q}^0 = F, \mathbb{Q}_2^0 = \mathbb{Q} - F$$



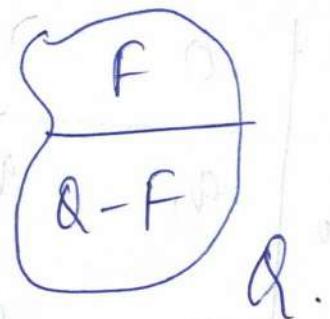
P

## Construction of Minimum DFA

(88-89)

1.  $\Pi_0 \rightarrow 0^{\text{th}}$  equivalent class

$$\Pi_0 = (\overset{\circ}{Q_1}, Q_2) \text{ i.e. } \overset{\circ}{Q_1} = F, \overset{\circ}{Q_2} = Q - F$$



2. Construction of  $\Pi_{k+1}$  from  $\Pi_k$

$$q_1, q_2 \in \overset{k}{Q_i} \text{ and}$$

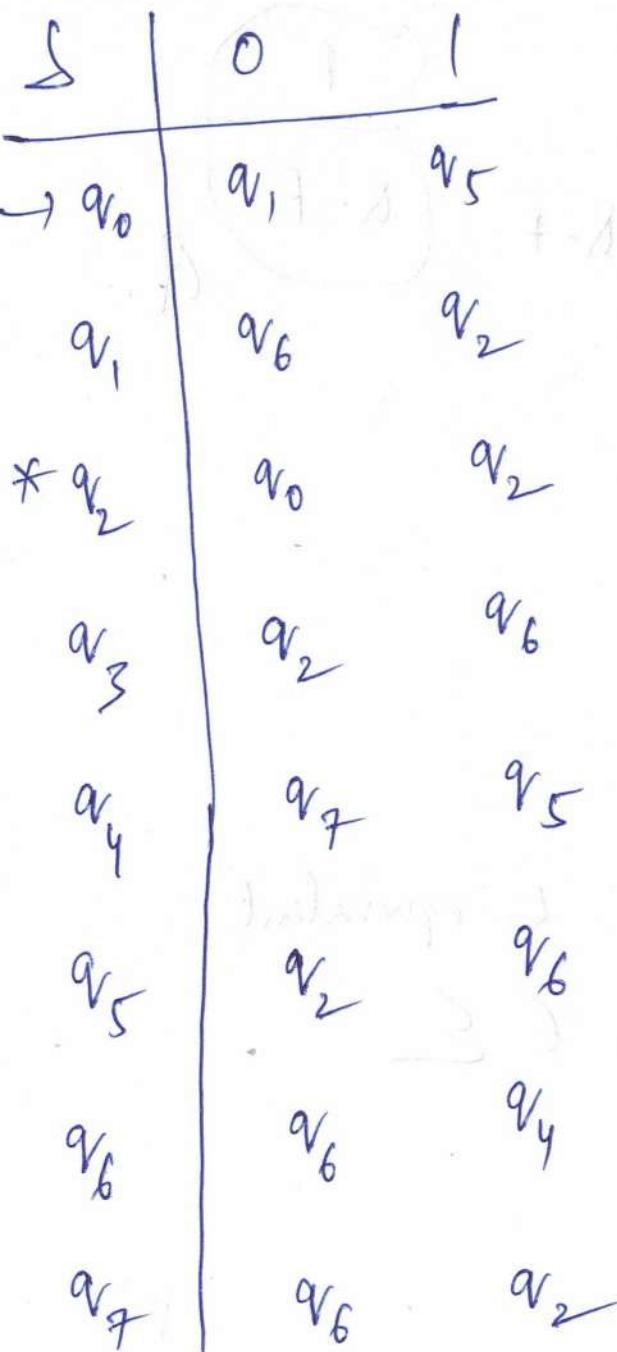
$q_1, q_2 \in \overset{k+1}{Q_j}$  if  $l(q_1, a), s(q_2, a)$  are  $k$ -equivalent  
+  $a \in \Sigma$

Repeat until

$$\Pi_n = \Pi_{n+1} \checkmark$$

P

Example:-



$$\Pi_0 = \{ \{q_2\}, \{q_0, q_1, q_3, q_4, q_5, q_6, q_7\} \} \quad (90)$$

$$\Pi_1 = ??$$

$$q_0 \equiv q_1 ??$$

$$\delta(q_{0,0}) = q_1 \quad \delta(q_{0,1}) = q_1 \notin F \quad \therefore q_0 \equiv q_1$$

$$\delta(q_{1,0}) = q_6 \quad \delta(q_{1,1}) = \underline{q_2} \in F$$

$$q_0 \equiv q_4 ???$$

$$\delta(q_{0,0}) = q_1 \quad \delta(q_{0,1}) = q_7$$

$$\delta(q_{1,0}) = q_5 \quad \delta(q_{1,1}) = q_5$$

Example:-

$\delta$	0	1
$\rightarrow v_0$	$v_1$	$v_5$
$v_1$	$v_6$	$v_2$
$*v_2$	$v_0$	$v_2$
$v_3$	$v_2$	$v_6$
$v_4$	$v_7$	$v_5$
$v_5$	$v_2$	$v_6$
$v_6$	$v_6$	$v_4$
$v_7$	$v_6$	$v_2$

$$\Pi_0 = \left\{ \{v_2\}, \{v_0, v_1, v_3, v_4, v_5, v_6, v_7\} \right\}$$

(90)-(91)

$$\Pi_1 = ??$$

$$v_0 = v_1 ??$$

$$\delta(v_0, 0) = v_1 \quad \delta(v_0, 1) = v_1 \notin F \quad \therefore v_0 = v_1$$

$$\delta(v_1, 0) = v_6 \quad \delta(v_1, 1) = \underline{v_2} \in F$$

$$v_0 = v_4 ???$$

$$\delta(v_0, 0) = v_1 \quad \delta(v_0, 1) = v_7$$

$$\delta(v_1, 0) = v_5 \quad \delta(v_1, 1) = v_5$$

$$\Pi_1 = \left\{ \{v_2\}, \{v_0, v_4, v_6\}, \{v_1, v_7\}, \{v_3, v_5\} \right\}$$

$$\Pi_2 = ??$$

$$\delta(v_4, 0) = v_7 \} \text{ belongs to diff. classes}$$

$$\delta(v_6, 0) = v_6 \} \text{ in } \Pi_1$$

Example:-

(90)-(91)-(92)

$$\Pi_0 = \left\{ \{q_2\}, \{q_0, q_1, q_3, q_4, q_5, q_6, q_7\} \right\}$$

$\delta$	0	1
$\rightarrow q_0$	$q_1$	$q_5$
$q_1$	$q_6$	$q_2$
$* q_2$	$q_0$	$q_2$
$q_3$	$q_2$	$q_6$
$q_4$	$q_7$	$q_5$
$q_5$	$q_2$	$q_6$
$q_6$	$q_6$	$q_4$
$q_7$	$q_6$	$q_2$

$$\Pi_1 = ??$$

$$q_0 \equiv q_1 ??$$

$$\delta(q_{0,0}) = q_1 \quad \delta(q_{0,1}) = q_1 \notin F \quad \therefore q_0 \equiv q_1$$

$$\delta(q_{1,0}) = q_6 \quad \delta(q_{1,1}) = q_2 \in F$$

$$q_0 \equiv q_4 ???$$

$$\delta(q_{0,0}) = q_1 \quad \delta(q_{0,1}) = q_7$$

$$\delta(q_{0,1}) = q_5 \quad \delta(q_{0,1}) = q_5$$

$$\Pi_1 = \left\{ \{q_2\}, \{q_0, q_4, q_6\}, \{q_1, q_7\}, \{q_3, q_5\} \right\}$$

$$\Pi_2 = ??$$

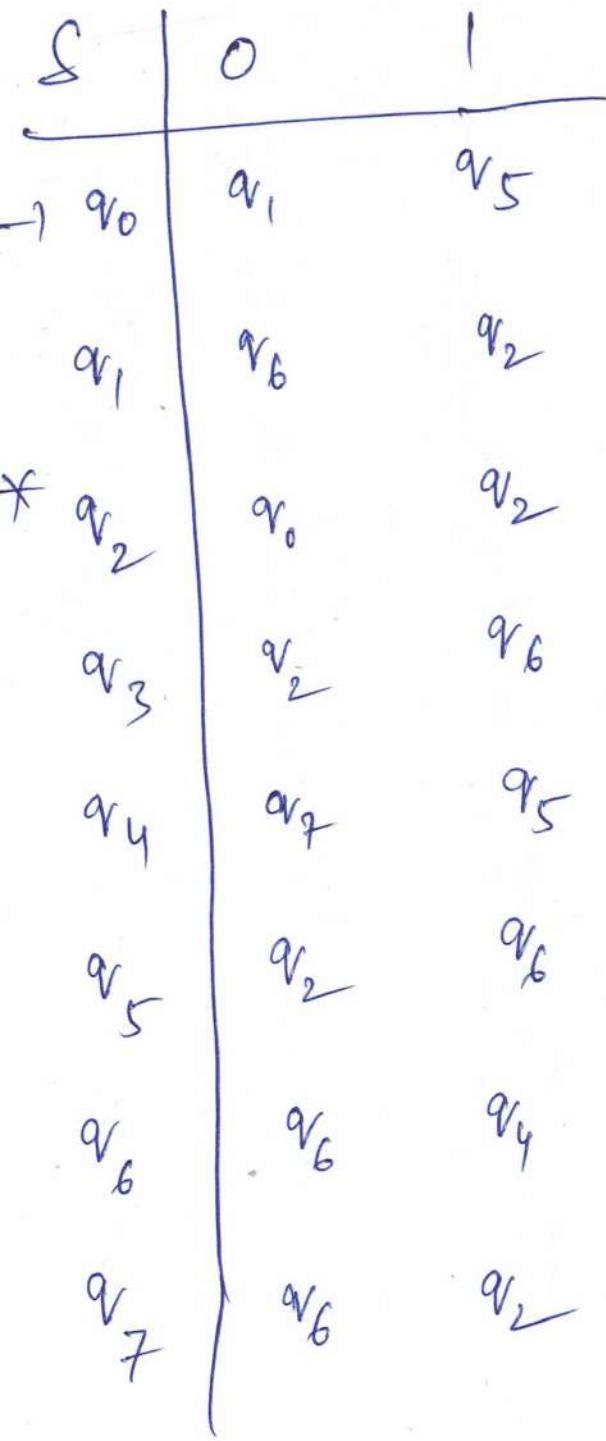
$$\delta(q_{0,0}) = q_7 \quad \text{belongs to diff. classes}$$

$$\delta(q_{0,0}) = q_6 \quad \text{in } \Pi_1$$

$$\Pi_2 = \left\{ \{q_2\}, \{q_6\}, \{q_0, q_4\}, \{q_1, q_7\}, \{q_3, q_5\} \right\}$$

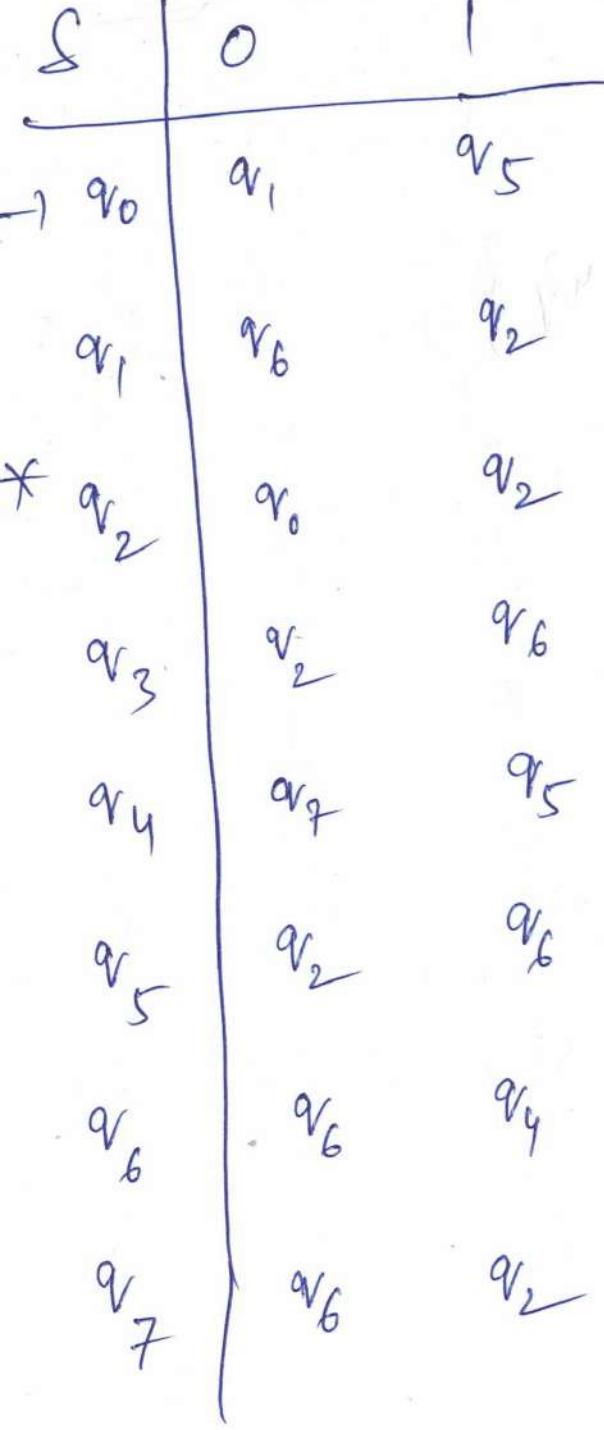
$$\Pi_3 = \underline{\Pi_2}$$

Q



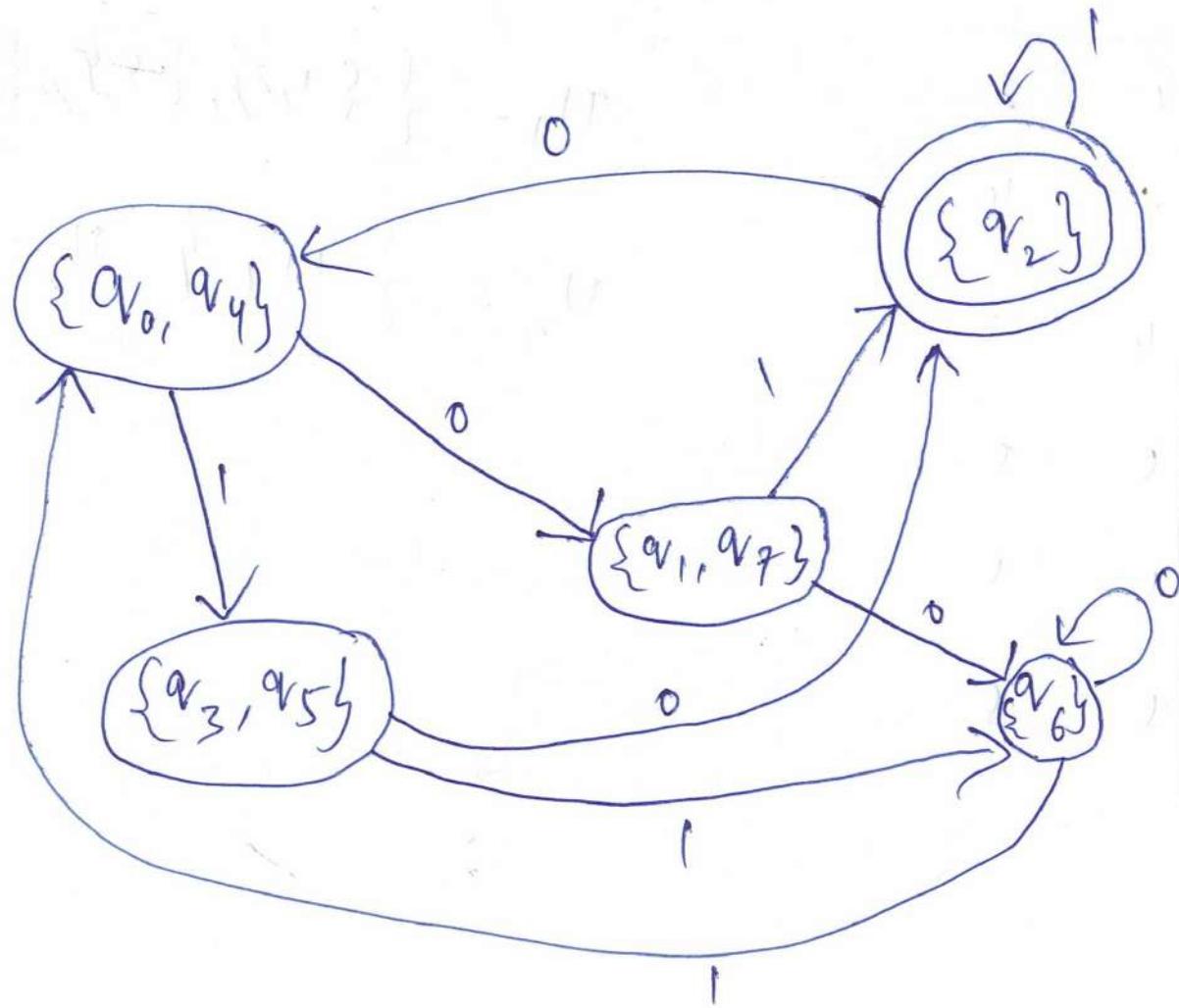
$$\Pi_3 = \Pi_2 = \{\{v_2\}, \{v_6\}, \{v_0, v_4\}, \{v_1, v_7\}, \{v_3, v_5\}\} \quad (93)$$

(P)



$$\Pi_3 = \Pi_2 = \{\{v_2\}, \{v_6\}, \{v_0, v_4\}, \{v_1, v_7\}, \{v_3, v_5\}\}$$

(93) (94)



Q

Example  $\Sigma = \{a, b\}$

Q5

S	a	b
*	0	1 2
*	1	3 4
*	2	4 3
3	5 5	
4	5 5	
*	5	5 5

$$\Pi_0 = \{\{1, 2, 5\}, \{0, 3, 4\}\}$$

$$\Pi_1 = \{\{1, 2\}, \{5\}, \{0\}, \{3, 4\}\}$$

$\Pi_2 = \cancel{\{\Pi_1\}}$   $\Pi_2 = \Pi_1$

Q

Example  $\Sigma = \{a, b\}$

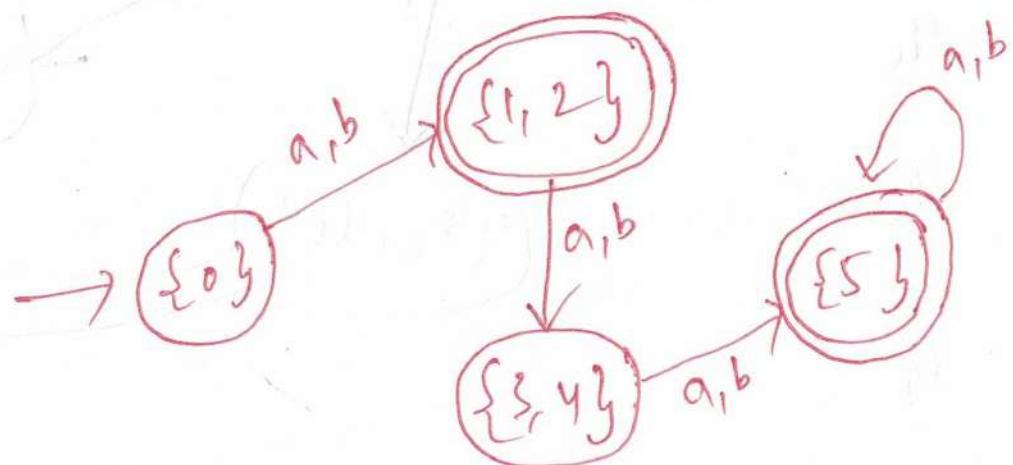
95 96

$S$	a	b
0	1	2
1	3	4
2	4	3
3	5	5
4	5	5
5	5	5

$$\Pi_0 = \{\{1, 2, 5\}, \{0, 3, 4\}\}$$

$$\Pi_1 = \{\{1, 2\}, \{5\}, \{0\}, \{3, 4\}\}$$

$\Pi_2 = \{\{1, 2\}, \{5\}\}$   $\Pi_2 = \Pi_1$



$$1 \approx 2, 3 \approx 4$$

Q

## Content Free Grammar (CFG)

$$G_1 = (V, T, P, S)$$

③  $P$  = finite set of productions  
rules.

①  $V \rightarrow$  Set of Variables    ②  $T =$  Set of terminals

$$V \neq \emptyset$$

$$T \neq \emptyset$$

$$|V| = \text{finite}$$

$$|T| = \text{finite}$$

$$V \cap T = \emptyset$$

④  $S =$  Starting Variable

$$S \in V$$

P

① Example:-  $G_1 = (\{E\}, \{+, *\}, (,), \text{id}\}, P, E)$  ② (CFL)

P:-

- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$
- $E \rightarrow \text{id}$

P

Example:-  $G_1 = (\{E\}, \{+, *\}, \{(), (), id\}, P, E)$  (CPU) ② ③

P:-  $E \rightarrow E + E$   
 $E \rightarrow E * E$   
 $E \rightarrow (E)$   
 $E \rightarrow id$

$$\begin{aligned} E &\Rightarrow E * E \\ &\Rightarrow (E) * E \\ &\Rightarrow (E) * id \\ &\Rightarrow (E + E) * id \\ &\Rightarrow \cancel{(E + id)} * id \\ &\Rightarrow \cancel{(id + id)} * id \end{aligned}$$

P

Example:-

Palindromes

(4)

$\epsilon, 0, 1$

$\epsilon, 0, 1, 00, 11, 010, 0110, 101, 1001 \dots$  etc.

$$G_{\text{pal}} = (\{\epsilon\}, \{0, 1\}, P, \epsilon)$$

P :- ? ?

P

Example:-

## Palindromes

(4) (5)

$\epsilon, 0, 1$

$\epsilon, 0, 1, 00, 11, 010, 0110, 101, 1001 \dots$  etc.

$$G_{\text{pal}} = (\{\epsilon\}, \{0, 1\}, P, \epsilon)$$

P:- ??

$$\underline{P}:- E \rightarrow \epsilon / 0 / 1 / 0 E 0 / 1 E 1 \quad \checkmark$$

$$E \Rightarrow \epsilon$$

$$E \rightarrow 0$$

$$E \rightarrow 1$$

$$E \rightarrow 0 E 0$$

$$E \rightarrow 1 E 1$$

P

3.11  
G<sub>1</sub> = (V, T, P, S)

(CFG)

A → R A α / ε

A ⇒ R A α ⇒ R R A α α ⇒ R R α α

A  $\xrightarrow[G_1]{*}$  R R α α

6  
G = (V, T, P, S)

L(G) = {w ∈ T<sup>\*</sup> | s  $\xrightarrow[G]{*}$  w}

↳ Language of the CFG G.

Q

Example:-  $G = (V, T, P, S)$

(P)

$$V = \{S\} \quad T = \{a, b\}$$

$$P = \{ \hookrightarrow asb \mid ab \}$$

$$L(G) = ? ?$$

(P)

Example:-  $G = (V, T, P, S)$

(P) 8

$$V = \{S\} \quad T = \{a, b\}$$

$$P = \{S \rightarrow aSb \mid ab\}$$

$$\begin{aligned} S &\Rightarrow ab \\ S &\xrightarrow{*} a^nb^n \\ S &\xrightarrow{*} a^3b^3 \\ &\quad \vdots \end{aligned}$$

$$L(G) = ? ?$$

$$L(G) = \{ab, a^nb^n, a^3b^3, \dots a^n b^n\}_{n \geq 1}$$

$$\therefore L(G) = \{a^n b^n \mid n \geq 1\}$$

P

Example find a CFG for all binary strings with ⑨

$L \{ \text{even number of } 0's \}$        $\Sigma = \{0, 1\}$

$L = \{ \epsilon, 00, 010, 001, \overset{100}{\dots} - \text{etc} \}$

Solut:-

Case 1) 1st symbol start with 1, followed by even number of 0's

2) 1st symbol start with 0, then followed by even 0's.

P

(9)  
(10)

Example find a CFG for all binary strings with  
 $L \{ \text{even number of } 0's \}$        $\Sigma = \{0, 1\}$

$$L = \{ \epsilon, \overset{\circ}{0}0, \overset{\circ}{0}10, \overset{\circ}{0}01, \overset{\circ}{0}00, \dots \text{etc} \}$$

Soln:-

Case 1) 1st symbol start with 0 followed by even number of 0's

2) 1st symbol start with 1 and having then followed by even 0's.

$$G_1 = (V, T, P, S)$$

$$S \rightarrow 1S / 0A0S / \epsilon$$

$$A \rightarrow 1A / \epsilon$$

P

$L = \{w \in \{0,1\}^* / w \text{ has even number of } 0's\}$

⑪

$G_1 = (\{\}, \{S, A\}, \{0, 1\}, P_1, S)$

$P_1: S \rightarrow 1S / 0A0S / \epsilon$

$A \rightarrow 1A / \epsilon$

⑫

$L = \{ \omega \in (0,1)^* / \omega \text{ has even number of } 0's \}$

(11)-12

$$G_{11} = (\{S, A\}, \{0, 1\}, P_1, S)$$

$$P_1: S \rightarrow 1S / 0A0S / \epsilon$$

$$A \rightarrow 1A / \epsilon$$

$$G_{12} = (\{S, T\}, \{0, 1\}, P_2, S)$$

$$P_2: \begin{array}{l} S \rightarrow 1S / 0T / \epsilon \\ T \rightarrow 1T / 0S \end{array}$$

$$L(G_{11}) = L = L(G_{12})$$

Two Equivalent CFGs.

Q

(13)

## Content Free Grammar

$$G = (V, T, P, S)$$

$$P: \{A \rightarrow \alpha\}$$

$$A \in V, \quad \alpha \in (V \cup T)^*$$

$$A \in (V \cup T)^*$$

Q

Example:  $L(G) = \{0^n 1^n 2^n \mid n \geq 1\}$  (14)

$$G_1 = (\{S, A_1\}, \{0, 1, 2\}, P, S)$$

$$P: S \rightarrow \text{OSA}^2$$

$$S \rightarrow 012$$

$$2A \rightarrow A2$$

$$1A \rightarrow 11$$

Example:  $L(G) = \{0^n 1^n 2^n \mid n \geq 1\}$

(14)  
(15)

$$G_1 = (\{S, A_0\}, \{0, 1, 2\}, P, S)$$

$$P: S \rightarrow \text{OSA}^2$$

$$S \rightarrow 012$$

$$2A \rightarrow A2$$

$$1A \rightarrow 11$$

$$S \rightarrow 012 \in L(G)$$

$$S \rightarrow 0\cancel{S} A2$$

$$S \rightarrow 001\cancel{2} A2$$

$$S \rightarrow 001\cancel{A}22$$

$$S \rightarrow 001122$$

P

## Derivation tree / Parse tree

$G_1 = (V, T, P, S)$ , a Parse tree for  $G_1$  ~~for~~ if :-

(i) Every vertex has label which is symbol of  $V \cup T \cup \{\epsilon\}$

(ii) Root is  $S$ .

(iii) If a vertex is interior and has label by  $A \in V$ ,

(iv)  $A \rightarrow x_1 x_2 \dots x_n$

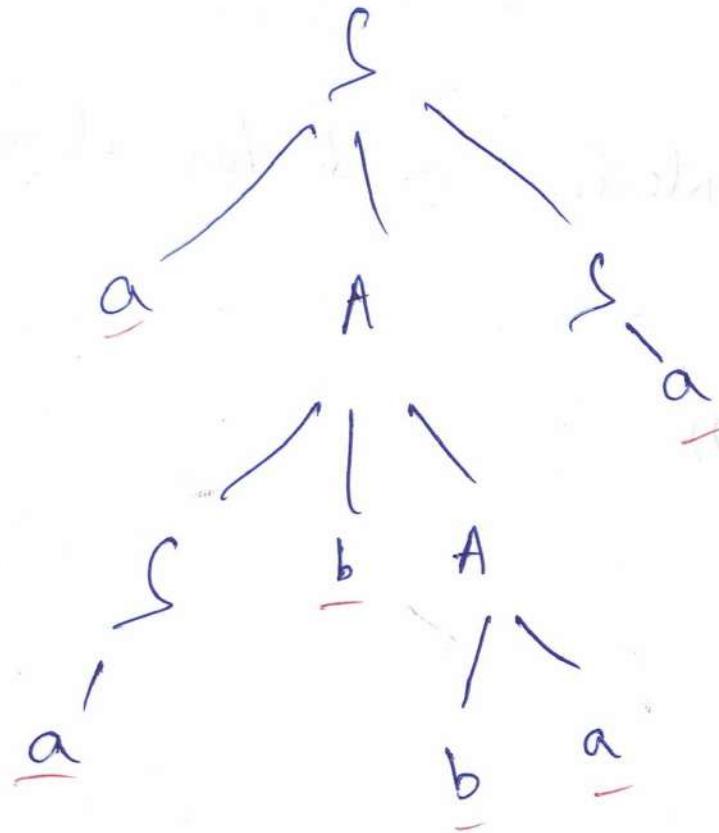
Example  $G_1 = (\{S, A\}, \{a, b\}, P, S)$

(17)

$P: S \rightarrow aAS/a$

$A \rightarrow SbA/ss/ba$

Parse tree



yields

aabbaa

~~parse tree of this string~~

~~aabbaa~~  
 $(V \cup T)^*$

yield string ∈

Q

(18)

$$\alpha \in (V \cup T)^*$$

$\frac{S}{G} \xrightarrow{*} \alpha$  iff we have a derivation tree  
that yield  $\alpha$ .

Example: P:  $S \rightarrow aAS/a$

$$A \rightarrow SbA/SS/ba$$

$$S \Rightarrow aAS \Rightarrow aSSS \Rightarrow aSbas \Rightarrow a\{bAbas$$

$$\alpha = a\underbrace{\{bAbas}_{\in (V \cup T)^*}$$

P

Example

$$G_1 = (\{S, A\}, \{a, b\}, P, S)$$

(19)

$$P: S \rightarrow aA\$ / a$$

$$A \rightarrow \$bA / \$\$ / ba$$

$$\alpha = \underline{\underline{aabbaa}}$$

Leftmost derivation ??

PR

Example

$$G_1 = (\{S, A\}, \{a, b\}, P, S)$$

(19)-20

P:  $S \rightarrow aAS/a$   
 $A \rightarrow SbA/S/ba$

$$\alpha = \underline{\underline{aabbaa}}$$

Leftmost derivation ?

$$S \Rightarrow a \underset{=} {AS} \Rightarrow a \underset{=} {SbAS} \Rightarrow aab \underset{=} {AS} \Rightarrow aabba \underset{=} {}$$

aabbba



$\alpha$ .

$$S \xrightarrow{*} aabbba$$

G

Q

Example:-  $G_1 = (\{S, A\}, \{a, b\}, P, S)$

(21)

$$P: S \rightarrow aAS/a$$
$$A \rightarrow SbA/Ss/ba$$

$$\alpha = aabbbaa$$

Right most derivation ??

P

Example:-  $G_1 = (\{S, A\}, \{a, b\}, P, S)$

(21) (22)

P:  $S \rightarrow aAS/a$        $\alpha = aabbbaa$

A  $\rightarrow SbA/Ss/ba$

Right most derivation ??

$$S \Rightarrow aA\underline{S} \Rightarrow a\underline{A}\underline{a} \Rightarrow a\underline{SbAa} \Rightarrow a\underline{S}\underline{bbaa}$$

aabbbaa  
 $\alpha.$

$\begin{matrix} * \\ S \Rightarrow \\ G_1 \end{matrix} aabbbaa$

P

~~Sample~~  $S \rightarrow 0B/1A, A \rightarrow 0/0S/1AA, B \rightarrow 1/1S/0BB$  (23)

$\alpha = 00110101$

$G = (\{S, A, B\}, \{0, 1, P, S\})$

Derivation tree for  $\alpha = 00110101$

Left most derivation. ??

P.

~~Sample~~  $S \rightarrow 0B/1A, A \rightarrow 0/0s/1AA, B \rightarrow 1/1s/0BB$  (23-24)

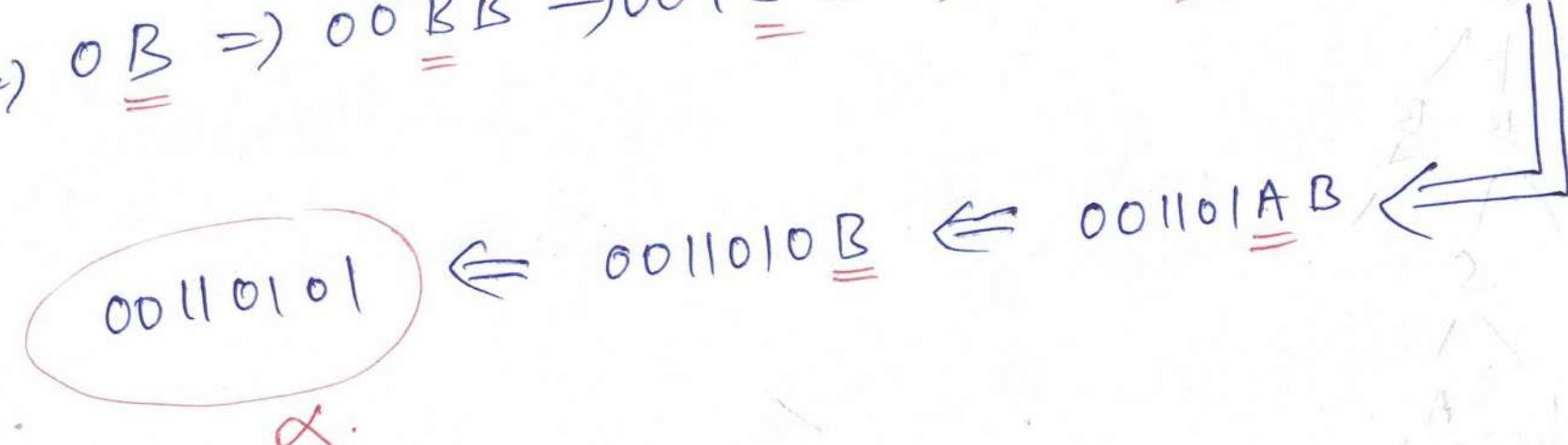
$$\alpha = 00110101$$

$$G = (\{S, A, B\}, \{0, 1\}, P, S)$$

Derivation tree for  $\alpha = 00110101$

Left most derivation. ??

$$S \Rightarrow 0\underline{B} \Rightarrow 00\underline{B}B \Rightarrow 001\underline{S}B \Rightarrow 0011\underline{A}B \Rightarrow 00110\underline{S}B$$



P.

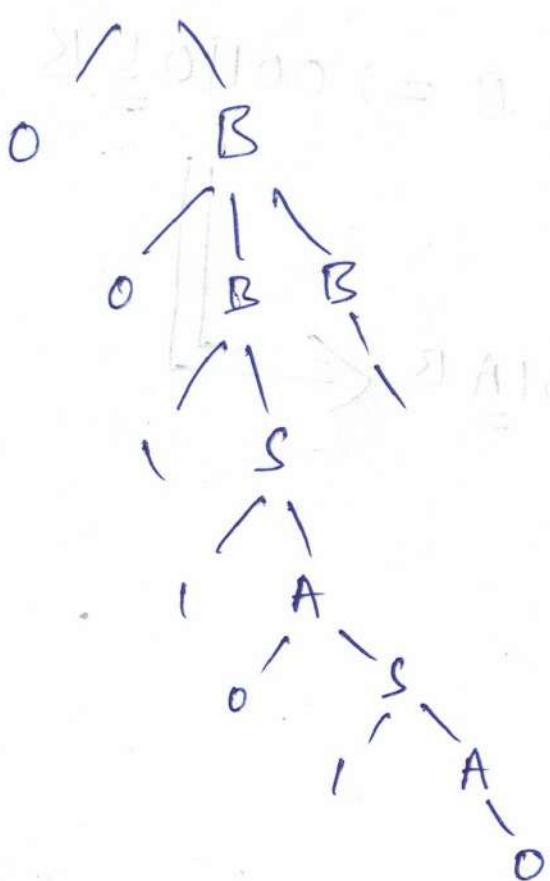
$S \rightarrow 0B/1A, A \rightarrow 0/0S/1AA, B \rightarrow 1/1S/0BB$

(25)

$$\alpha = 00110101$$

derivation tree for  $\alpha = 00110101$

Left most



26

$$S \rightarrow 0B / IA, \quad A \rightarrow 0/0S / CA, \quad B \rightarrow 1 / 1S / 0BB$$

$$\alpha = 00110101$$

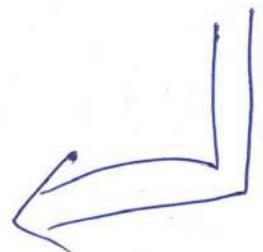
Right most derivation ??

$$S \Rightarrow 0\underline{B} \Rightarrow 00\underline{B}\underline{B} \Rightarrow 00\underline{B}\underline{1} \Rightarrow 001\underline{S} \Rightarrow 0011\underline{A}\underline{1}$$

00110101  
 $\cancel{\alpha}$

$\leftarrow 001101\underline{A}\underline{1}$

$\leftarrow 00110\underline{S}\underline{1}$



Q

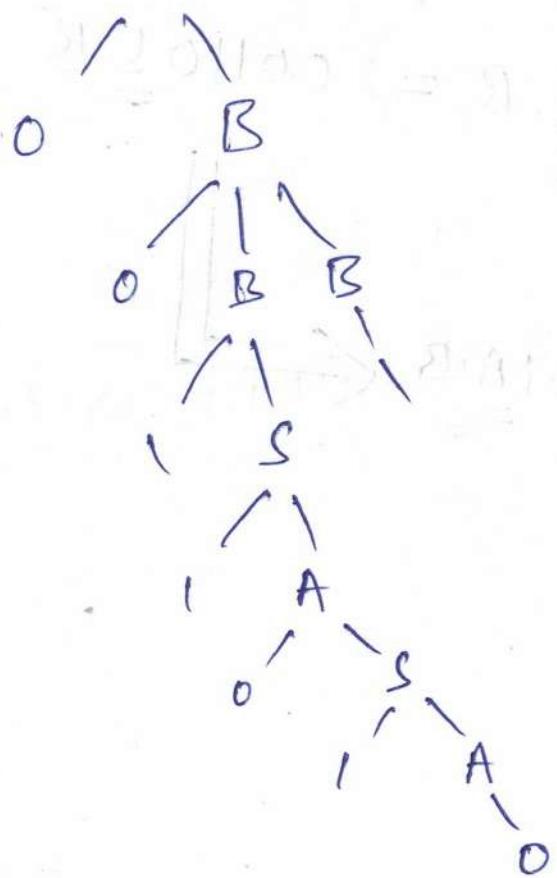
$S \rightarrow 0B/1A, A \rightarrow 0/0S/1AA, B \rightarrow 1/1S/0BB$

$$\alpha = 00110101$$

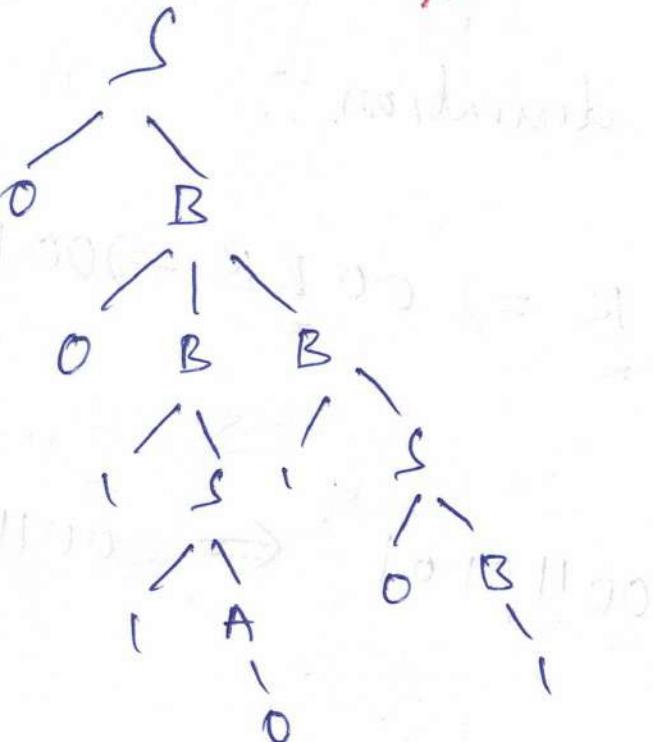
(25)  
(27)

derivation tree for  $\alpha = 00110101$

Left most



Right most



Another Left most derivation tree.

Q

$S \rightarrow OB / IA, A \rightarrow O/OS/IAA, B \rightarrow I/IS/OBB$

(28)

$\alpha = 00110101$

Another left most derivation for  $\alpha = 00110101$  ??

$S \Rightarrow OB \Rightarrow OO \underline{BB} \Rightarrow OO \underline{I} \underline{B} \Rightarrow OO \underline{II} \underline{AB}$

$\alpha = 00110101$

Same way Another Rightmost derivation also exist for  
Same derivation tree.

## Ambiguity in CFG.

Example:- "In books selected information is given."

Selected  $\vdash$  books  
Information.

Definition! - A terminal string  $w \in L(G)$  is ambiguous if

$\exists$  two or more derivation tree for  $w$ .

( $\exists$  two or more left most derivation of  $w$ )

- A CFG  $G$  is ambiguous if  $\exists$  some  $w \in L(G)$  which is ambiguous.

- A CFL for which every CFG is ambiguous then CFL is called inherently ambiguous.

Example :-  $G_1 = (\{S\}, \{a, b, +, *\}, P, S)$  (30)

P:  $S \rightarrow S+S \mid S*S \mid a \mid b$   $\omega = a+a*b.$

one derivation tree

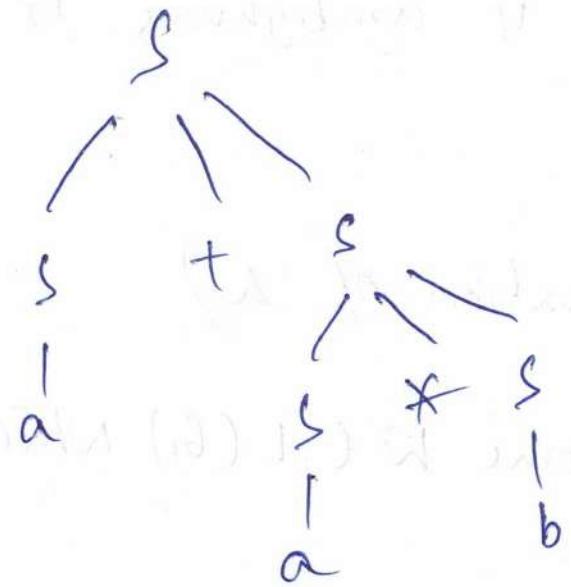
P

Example :-  $G_1 = (\{S\}, \{a, b, +, *\}, P, S)$

(30) (31)

P:  $S \rightarrow S + S \mid S * S \mid a \mid b$        $w = a + a * b.$

one derivation tree.



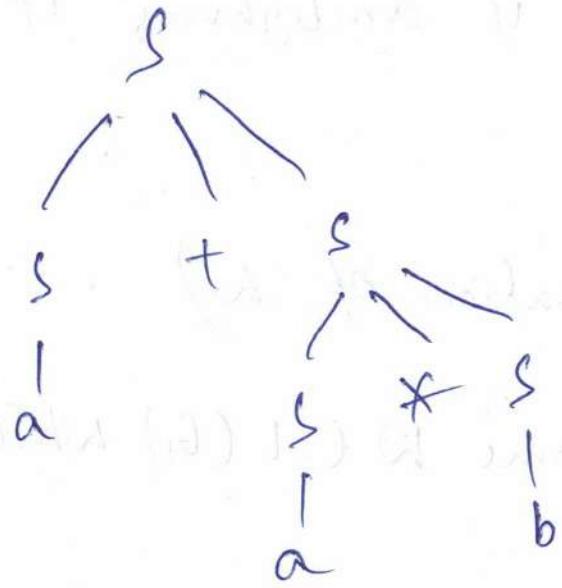
P

Example :-  $G_1 = (\{S\}, \{a, b, +, *\}, P, S)$

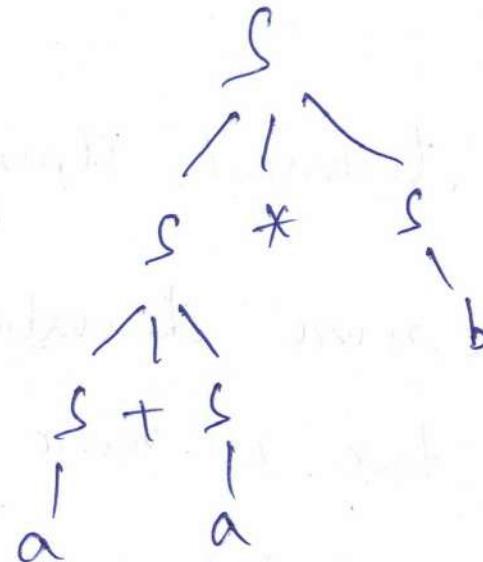
P:  $S \rightarrow S+S \mid S*S \mid a \mid b$        $w = a+a*b.$

(30) (31)  
(32)

one derivation tree:



Another derivation tree



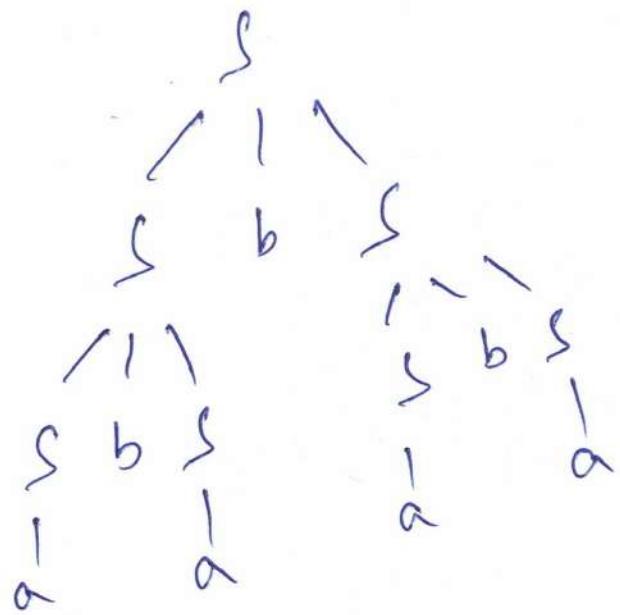
Two derivation trees exist so  $w = a+a*b$  is ambiguous.  
and the  $G_1$  is also ambiguous.

P

Another Example:  $G_1 = (\{s\}, \{a, b\}, P, S)$

$P: S \rightarrow S^* / a \quad w = ababab a$

One derivation tree

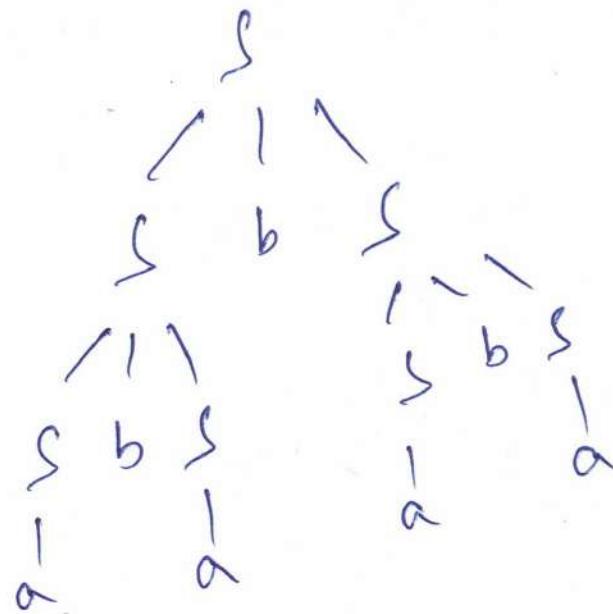


P

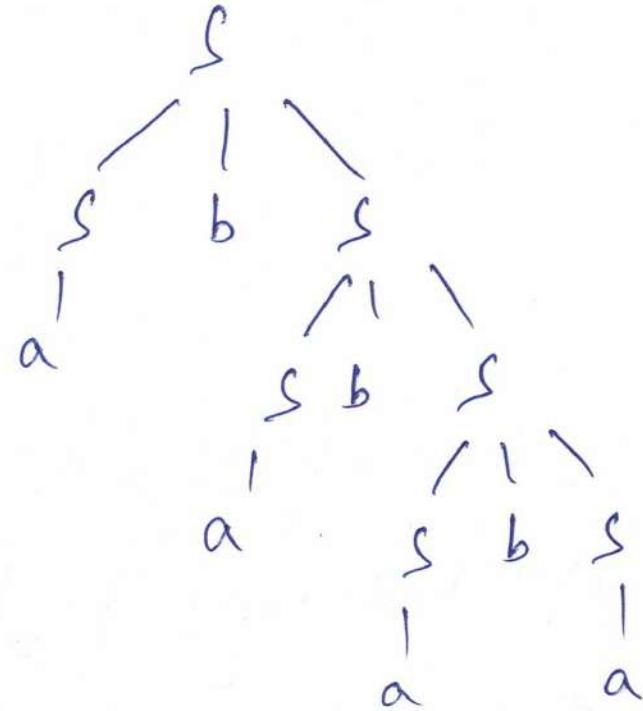
Another Example:-  $G = (\{S\}, \{a, b\}, P, S)$

$P: S \rightarrow SbS/a \quad w = ababab a$

One derivation tree



Another derivation tree



$w = ababab a$  is Ambiguous  $\Rightarrow G$  is Ambiguous

(P)

## Simplification of CFG

$G_1 = (V, T, P, S)$ , starting variable  $S \in V$

↓      ↓  
Variable    terminal

Productions

$$L(G_1) = \{ w \in T^* \mid S \xrightarrow[G_1]{*} w \}$$

↳ May be not all symbols vut are used

There is a possibility for simplification.

Example:  $G_1 = (\{S, A, B, D, E\}, \{a, b, c\}, P, S)$

$$P = \{S \rightarrow AB, A \rightarrow a, B \rightarrow b, B \rightarrow D, E \rightarrow c/e\}$$

$$L(G_1) = \{w \in T^* \mid \xrightarrow[G_1]{S} w\} = \{ab\}$$

P

Example:  $G_1 = (\{S, A, B, D, E\}, \{a, b, c\}, P, S)$

$P = \{S \rightarrow AB, A \rightarrow a, B \rightarrow b, \boxed{B \rightarrow D}, \boxed{E \rightarrow c/e}\}$

$$L(G_1) = \{w \in T^* \mid \xrightarrow[G_1]{S} w\} = \{ab\}$$

derivation

$$\xrightarrow[S]{} AB \Rightarrow A B \Rightarrow ab$$

~~Remove Revamp~~  $V, T, P$  which are not useful.

$$G'_1 = (V', T', P', S)$$

$$V' = \{S, A, B\}, T' = \{a, b\}$$

$P' = \text{remove red colored } \emptyset \text{ from } P$

$$= \{S \rightarrow AB, A \rightarrow a, B \rightarrow b\}.$$

Q

# Syntax Analysis:

Context-free Grammars, Pushdown Automata and Parsing  
Part - 1

Y.N. Srikant

Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis?
- Specification of programming languages: context-free grammars
- Parsing context-free languages: push-down automata
- Top-down parsing: LL(1) and recursive-descent parsing
- Bottom-up parsing: LR-parsing

- Every programming language has precise grammar rules that describe the syntactic structure of well-formed programs
  - In C, the rules state how functions are made out of parameter lists, declarations, and statements; how statements are made of expressions, etc.
- Grammars are easy to understand, and parsers for programming languages can be constructed automatically from certain classes of grammars
- Parsers or syntax analyzers are generated *for* a particular grammar
- Context-free grammars are usually used for syntax specification of programming languages

# What is Parsing or Syntax Analysis?

- A parser for a grammar of a programming language
  - verifies that the string of tokens for a program in that language can indeed be generated from that grammar
  - reports any syntax errors in the program
  - constructs a parse tree representation of the program (not necessarily explicit)
  - usually calls the lexical analyzer to supply a token to it when necessary
  - could be hand-written or automatically generated
  - is based on *context-free* grammars
- Grammars are generative mechanisms like regular expressions
- Pushdown automata are machines recognizing context-free languages (like FSA for RL)

# Context-free Grammars

- A CFG is denoted as  $G = (N, T, P, S)$ 
  - $N$ : Finite set of non-terminals
  - $T$ : Finite set of terminals
  - $S \in N$ : The start symbol
  - $P$ : Finite set of productions, each of the form  $A \rightarrow \alpha$ , where  $A \in N$  and  $\alpha \in (N \cup T)^*$
- Usually, only  $P$  is specified and the first production corresponds to that of the start symbol
- Examples

(1)	(2)	(3)	(4)
$E \rightarrow E + E$	$S \rightarrow 0S0$	$S \rightarrow aSb$	$S \rightarrow aB \mid bA$
$E \rightarrow E * E$	$S \rightarrow 1S1$	$S \rightarrow \epsilon$	$A \rightarrow a \mid aS \mid bAA$
$E \rightarrow (E)$	$S \rightarrow 0$		$B \rightarrow b \mid bS \mid aBB$
$E \rightarrow id$	$S \rightarrow 1$		
	$S \rightarrow \epsilon$		

# Derivations

- $E \xrightarrow{E \rightarrow E+E} E + E \xrightarrow{E \rightarrow id} id + E \xrightarrow{E \rightarrow id} id + id$   
is a derivation of the terminal string  $id + id$  from  $E$
- In a derivation, a production is applied at each step, to replace a nonterminal by the right-hand side of the corresponding production
- In the above example, the productions  $E \rightarrow E + E$ ,  $E \rightarrow id$ , and  $E \rightarrow id$ , are applied at steps 1,2, and, 3 respectively
- The above derivation is represented in short as,  
 $E \Rightarrow^* id + id$ , and is read as **S derives**  $id + id$

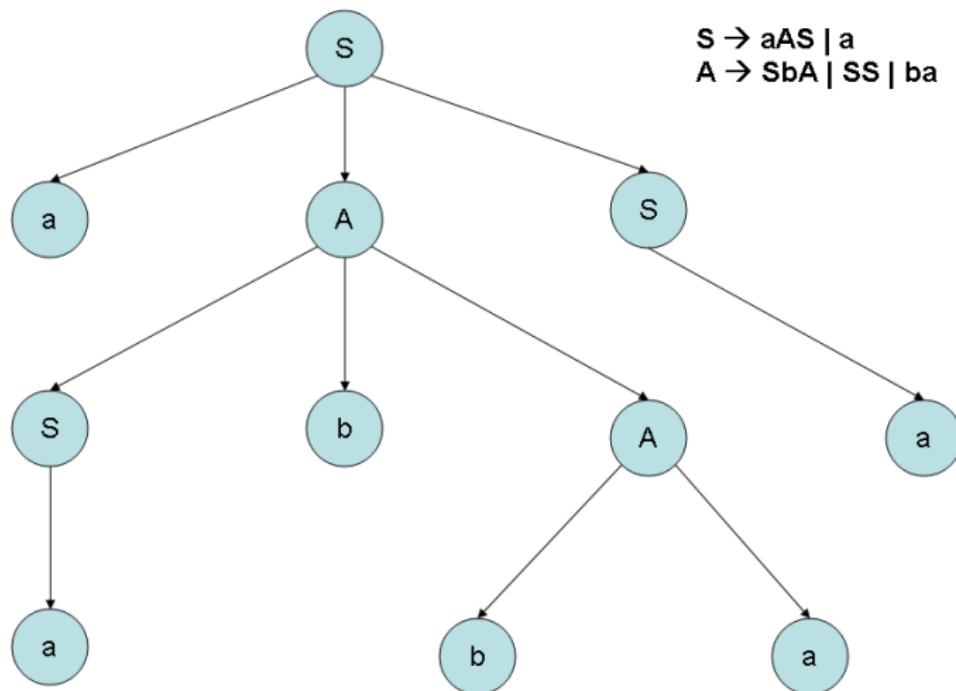
# Context-free Languages

- Context-free grammars generate context-free languages (grammar and language resp.)
- The *language generated by G*, denoted  $L(G)$ , is  
$$L(G) = \{w \mid w \in T^*, \text{ and } S \Rightarrow^* w\}$$
i.e., a string is in  $L(G)$ , if
  - ① the string consists solely of terminals
  - ② the string can be derived from  $S$
- Examples
  - ①  $L(G_1)$  = Set of all expressions with  $+$ ,  $*$ , names, and balanced ' $($ ' and ' $)$ '
  - ②  $L(G_2)$  = Set of palindromes over 0 and 1
  - ③  $L(G_3) = \{a^n b^n \mid n \geq 0\}$
  - ④  $L(G_4) = \{x \mid x \text{ has equal no. of } a's \text{ and } b's\}$
- A string  $\alpha \in (N \cup T)^*$  is a **sentential form** if  $S \Rightarrow^* \alpha$
- Two grammars  $G_1$  and  $G_2$  are equivalent, if  $L(G_1) = L(G_2)$

# Derivation Trees

- Derivations can be displayed as trees
- The internal nodes of the tree are all nonterminals and the leaves are all terminals
- Corresponding to each internal node A, there exists a production  $\in P$ , with the RHS of the production being the list of children of A, read from left to right
- The **yield** of a derivation tree is the list of the labels of all the leaves read from left to right
- If  $\alpha$  is the yield of some derivation tree for a grammar  $G$ , then  $S \Rightarrow^* \alpha$  and conversely

# Derivation Tree Example

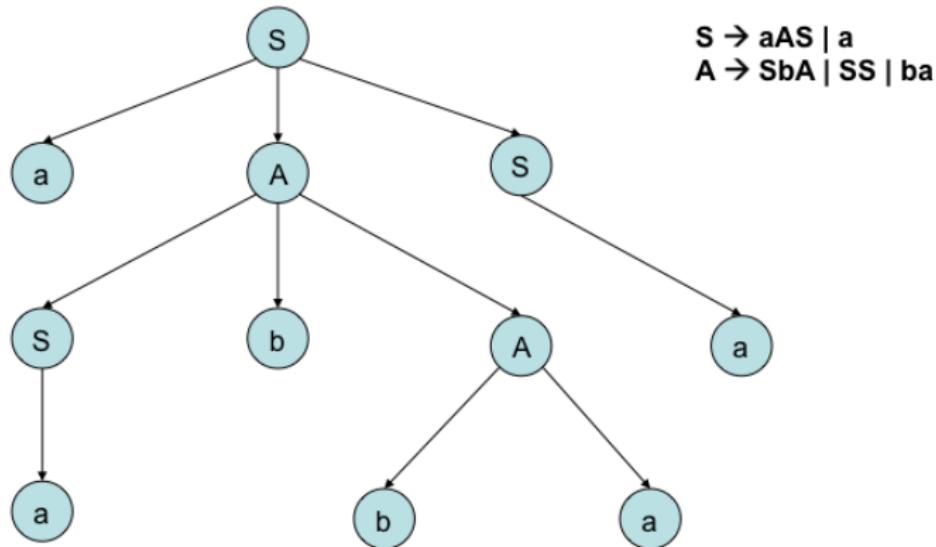


$S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbaS \Rightarrow aabbaaa$

# Leftmost and Rightmost Derivations

- If at each step in a derivation, a production is applied to the leftmost nonterminal, then the derivation is said to be **leftmost**. Similarly **rightmost derivation**.
- If  $w \in L(G)$  for some  $G$ , then  $w$  has at least one parse tree and corresponding to a parse tree,  $w$  has unique leftmost and rightmost derivations
- If some word  $w$  in  $L(G)$  has two or more parse trees, then  $G$  is said to be **ambiguous**
- A CFL for which every  $G$  is ambiguous, is said to be an **inherently ambiguous** CFL

# Leftmost and Rightmost Derivations: An Example



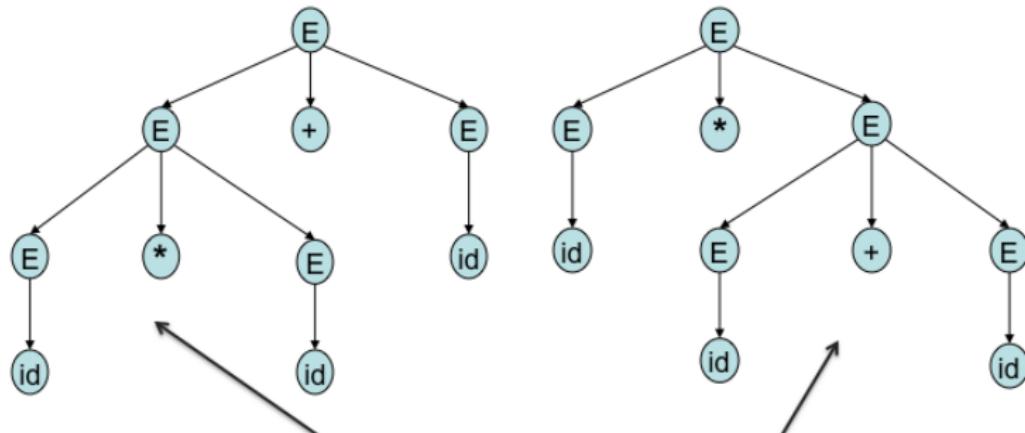
Leftmost derivation:  $S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbaS \Rightarrow aabbaa$

Rightmost derivation:  $S \Rightarrow aAS \Rightarrow aAa \Rightarrow aSbAa \Rightarrow aSbbAa \Rightarrow aSbbbaa \Rightarrow aabbaa$

# Ambiguous Grammar Examples

- The grammar,  $E \rightarrow E + E | E * E | (E) | id$  is ambiguous, but the following grammar for the same language is unambiguous  
 $E \rightarrow E + T | T, T \rightarrow T * F | F, F \rightarrow (E) | id$
- The grammar,  
 $stmt \rightarrow IF\ expr\ stmt | IF\ expr\ stmt\ ELSE\ stmt | other\_stmt$  is ambiguous, but the following equivalent grammar is not  
 $stmt \rightarrow IF\ expr\ stmt | IF\ expr\ matched\_stmt\ ELSE\ stmt$   
 $matched\_stmt \rightarrow IF\ expr\ matched\_stmt\ ELSE\ matched\_stmt | other\_stmt$
- The language,  
 $L = \{a^n b^n c^m d^m \mid n, m \geq 1\} \cup \{a^n b^m c^m d^n \mid n, m \geq 1\},$  is inherently ambiguous

# Ambiguity Example 1

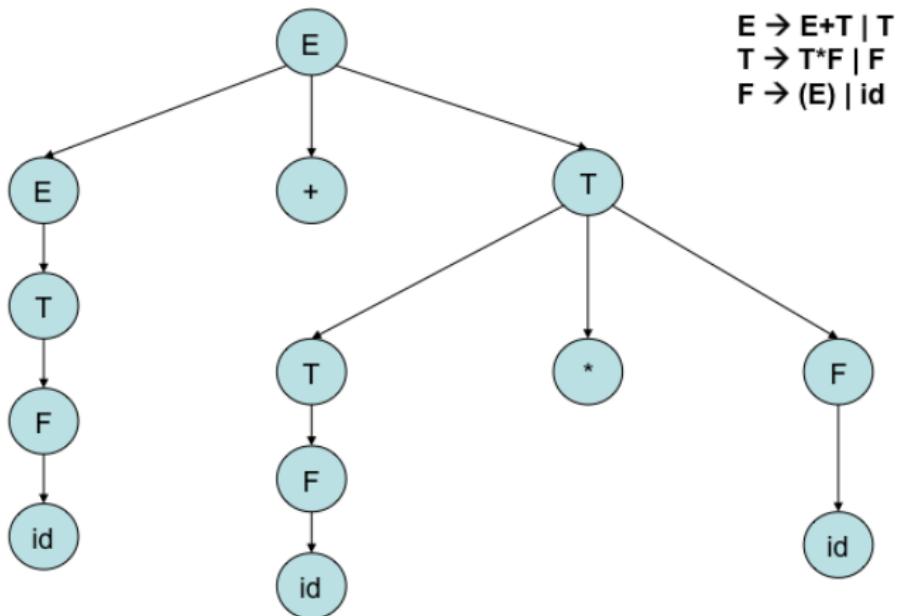


$E \Rightarrow E+E \Rightarrow E^*E+E \Rightarrow id^*E+E \Rightarrow id^*id+E \Rightarrow id^*id+id$

$E \Rightarrow E^*E \Rightarrow id^*E \Rightarrow id^*E+E \Rightarrow id^*id+E \Rightarrow id^*id+id$

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

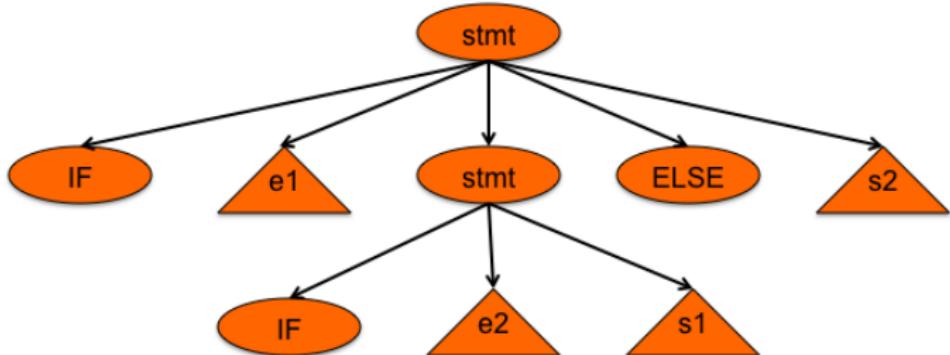
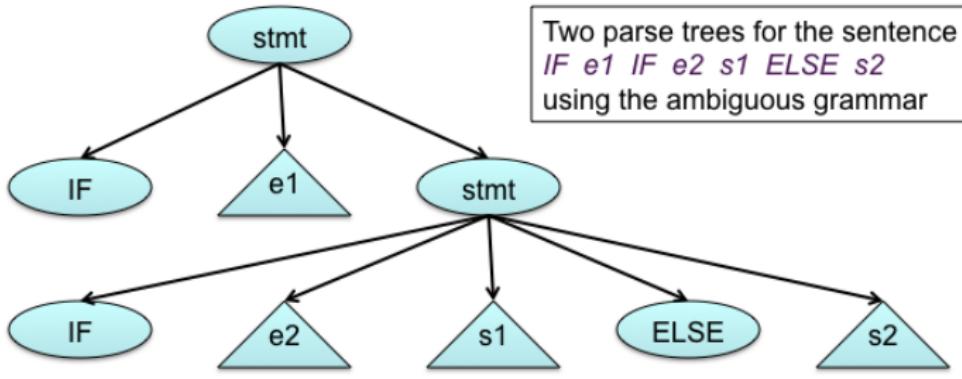
# Equivalent Unambiguous Grammar



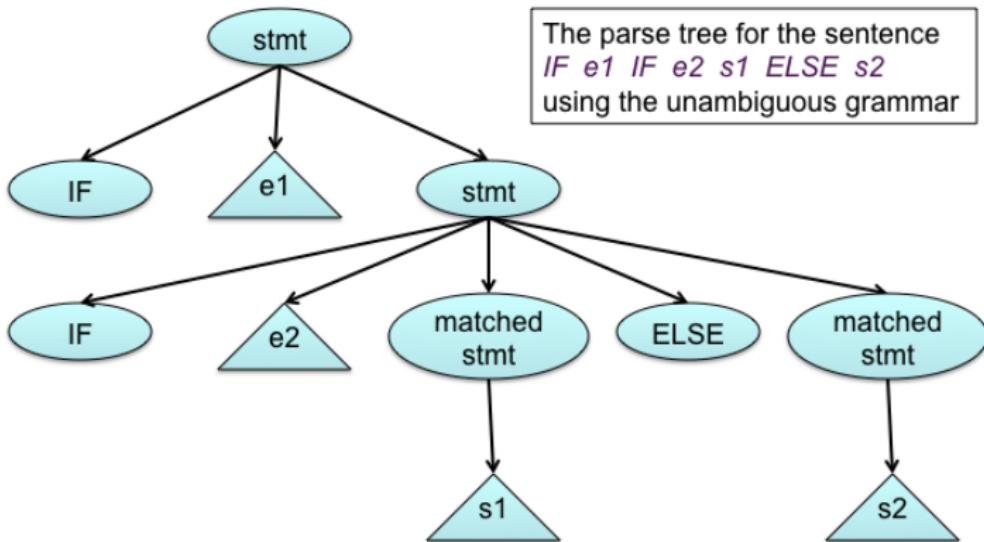
$E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow id+T \Rightarrow id+T^*F \Rightarrow id+F^*F \Rightarrow id+id^*F \Rightarrow id+id^*id$

$E \Rightarrow T^*F \Rightarrow F^*F \Rightarrow (E)^*F \Rightarrow (E+T)^*F \Rightarrow (T+T)^*F \Rightarrow (F+T)^*F \Rightarrow (id+T)^*F$   
 $\Rightarrow (id+F)^*id \Rightarrow (id+id)^*F \Rightarrow (id+id)^*id$

# Ambiguity Example 2



## Ambiguity Example 2 (contd.)



```
s → IF e s | IF e ms ELSE s  
ms → IF e ms ELSE ms | other_s
```

# Fragment of C-Grammar (Statements)

```
program --> VOID MAIN '(' ')' compound_stmt
compound_stmt --> '{}' | '{' stmt_list '}'
                  | '{' declaration_list stmt_list '}'
stmt_list --> stmt | stmt_list stmt
stmt --> compound_stmt | expression_stmt
                  | if_stmt | while_stmt
expression_stmt --> ';' | expression ';'
if_stmt --> IF '(' expression ')' stmt
          | IF '(' expression ')' stmt ELSE stmt
while_stmt --> WHILE '(' expression ')' stmt
expression --> assignment_expr
                  | expression ',' assignment_expr
```

# Fragment of C-Grammar (Expressions)

```
assignment_expr --> logical_or_expr
    | unary_expr assign_op assignment_expr
assign_op --> '=' | MUL_ASSIGN | DIV_ASSIGN
    | ADD_ASSIGN | SUB_ASSIGN
    | AND_ASSIGN | OR_ASSIGN
unary_expr --> primary_expr
    | unary_operator unary_expr
unary_operator --> '+' | '-' | '!'
primary_expr --> ID | NUM | '(' expression ')'
logical_or_expr --> logical_and_expr
    | logical_or_expr OR_OP logical_and_expr
logical_and_expr --> equality_expr
    | logical_and_expr AND_OP equality_expr
equality_expr --> relational_expr
    | equality_expr EQ_OP relational_expr
    | equality_expr NE_OP relational_expr
```

# Fragment of C-Grammar (Expressions and Declarations)

```
relational_expr --> add_expr
                    | relational_expr '<' add_expr
                    | relational_expr '>' add_expr
                    | relational_expr LE_OP add_expr
                    | relational_expr GE_OP add_expr
add_expr --> mult_expr | add_expr '+' mult_expr
                    | add_expr '-' mult_expr
mult_expr --> unary_expr | mult_expr '*' unary_expr
                    | mult_expr '/' unary_expr
declarationlist --> declaration
                    | declarationlist declaration
declaration --> type idlist ';'
idlist --> idlist ',' ID | ID
type --> INT_TYPE | FLOAT_TYPE | CHAR_TYPE
```

# Pushdown Automata

A PDA  $M$  is a system  $(Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$ , where

- $Q$  is a finite set of states
- $\Sigma$  is the input alphabet
- $\Gamma$  is the stack alphabet
- $q_0 \in Q$  is the start state
- $z_0 \in \Gamma$  is the start symbol on stack (initialization)
- $F \subseteq Q$  is the set of final states
- $\delta$  is the transition function,  $Q \times \Sigma \cup \{\epsilon\} \times \Gamma$  to finite subsets of  $Q \times \Gamma^*$

A typical entry of  $\delta$  is given by

$$\delta(q, a, z) = \{(p_1, \gamma_1), ((p_2, \gamma_2), \dots, (p_m, \gamma_m))\}$$

The PDA in state  $q$ , with input symbol  $a$  and top-of-stack symbol  $z$ , can enter any of the states  $p_i$ , replace the symbol  $z$  by the string  $\gamma_i$ , and advance the input head by one symbol.

## Pushdown Automata (contd.)

- The leftmost symbol of  $\gamma_i$  will be the new top of stack
- $a$  in the above function  $\delta$  could be  $\epsilon$ , in which case, the input symbol is not used and the input head is not advanced
- For a PDA  $M$ , we define  $L(M)$ , the language accepted by  **$M$  by final state**, to be

$$L(M) = \{w \mid (q_0, w, Z_0) \vdash^* (p, \epsilon, \gamma), \text{ for some } p \in F \text{ and } \gamma \in \Gamma^*\}$$

- We define  $N(M)$ , the language accepted by  **$M$  by empty stack**, to be

$$N(M) = \{w \mid (q_0, w, Z_0) \vdash^* (p, \epsilon, \epsilon), \text{ for some } p \in Q\}$$

- When acceptance is by empty stack, the set of final states is irrelevant, and usually, we set  $F = \emptyset$

# PDA - Examples

- $L = \{0^n 1^n \mid n \geq 0\}$   
 $M = (\{q_0, q_1, q_2, q_3\}, \{0, 1\}, \{Z, 0\}, \delta, q_0, Z, \{q_0\})$ , where  $\delta$  is defined as follows  
 $\delta(q_0, 0, Z) = \{(q_1, 0Z)\}, \delta(q_1, 0, 0) = \{(q_1, 00)\},$   
 $\delta(q_1, 1, 0) = \{(q_2, \epsilon)\}, \delta(q_2, 1, 0) = \{(q_2, \epsilon)\},$   
 $\delta(q_2, \epsilon, Z) = \{(q_0, \epsilon)\}$
- $(q_0, 0011, Z) \vdash (q_1, 011, 0Z) \vdash (q_1, 11, 00Z) \vdash (q_2, 1, 0Z) \vdash (q_2, \epsilon, Z) \vdash (q_0, \epsilon, \epsilon)$
- $(q_0, 001, Z) \vdash (q_1, 01, 0Z) \vdash (q_1, 1, 00Z) \vdash (q_2, \epsilon, 0Z) \vdash error$
- $(q_0, 010, Z) \vdash (q_1, 10, 0Z) \vdash (q_2, 0, Z) \vdash error$

# Syntax Analysis:

Context-free Grammars, Pushdown Automata and Parsing  
Part - 2

Y.N. Srikant

Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis? (covered in lecture 1)
- Specification of programming languages: context-free grammars (covered in lecture 1)
- Parsing context-free languages: push-down automata
- Top-down parsing: LL(1) and recursive-descent parsing
- Bottom-up parsing: LR-parsing

# Pushdown Automata

A PDA  $M$  is a system  $(Q, \Sigma, \Gamma, \delta, q_0, z_0, F)$ , where

- $Q$  is a finite set of states
- $\Sigma$  is the input alphabet
- $\Gamma$  is the stack alphabet
- $q_0 \in Q$  is the start state
- $z_0 \in \Gamma$  is the start symbol on stack (initialization)
- $F \subseteq Q$  is the set of final states
- $\delta$  is the transition function,  $Q \times \Sigma \cup \{\epsilon\} \times \Gamma$  to finite subsets of  $Q \times \Gamma^*$

A typical entry of  $\delta$  is given by

$$\delta(q, a, z) = \{(p_1, \gamma_1), ((p_2, \gamma_2), \dots, (p_m, \gamma_m))\}$$

The PDA in state  $q$ , with input symbol  $a$  and top-of-stack symbol  $z$ , can enter any of the states  $p_i$ , replace the symbol  $z$  by the string  $\gamma_i$ , and advance the input head by one symbol.

## Pushdown Automata (contd.)

- The leftmost symbol of  $\gamma_i$  will be the new top of stack
- $a$  in the above function  $\delta$  could be  $\epsilon$ , in which case, the input symbol is not used and the input head is not advanced
- For a PDA  $M$ , we define  $L(M)$ , the language accepted by  **$M$  by final state**, to be

$$L(M) = \{w \mid (q_0, w, Z_0) \vdash^* (p, \epsilon, \gamma), \text{ for some } p \in F \text{ and } \gamma \in \Gamma^*\}$$

- We define  $N(M)$ , the language accepted by  **$M$  by empty stack**, to be

$$N(M) = \{w \mid (q_0, w, Z_0) \vdash^* (p, \epsilon, \epsilon), \text{ for some } p \in Q\}$$

- When acceptance is by empty stack, the set of final states is irrelevant, and usually, we set  $F = \emptyset$

# PDA - Examples

- $L = \{0^n 1^n \mid n \geq 0\}$   
 $M = (\{q_0, q_1, q_2, q_3\}, \{0, 1\}, \{Z, 0\}, \delta, q_0, Z, \{q_0\})$ , where  $\delta$  is defined as follows  
 $\delta(q_0, 0, Z) = \{(q_1, 0Z)\}, \delta(q_1, 0, 0) = \{(q_1, 00)\},$   
 $\delta(q_1, 1, 0) = \{(q_2, \epsilon)\}, \delta(q_2, 1, 0) = \{(q_2, \epsilon)\},$   
 $\delta(q_2, \epsilon, Z) = \{(q_0, \epsilon)\}$
- $(q_0, 0011, Z) \vdash (q_1, 011, 0Z) \vdash (q_1, 11, 00Z) \vdash (q_2, 1, 0Z) \vdash (q_2, \epsilon, Z) \vdash (q_0, \epsilon, \epsilon)$
- $(q_0, 001, Z) \vdash (q_1, 01, 0Z) \vdash (q_1, 1, 00Z) \vdash (q_2, \epsilon, 0Z) \vdash error$
- $(q_0, 010, Z) \vdash (q_1, 10, 0Z) \vdash (q_2, 0, Z) \vdash error$

## PDA - Examples (contd.)

- $L = \{ww^R \mid w \in \{a, b\}^+\}$

$M = (\{q_0, q_1, q_2\}, \{a, b\}, \{Z, a, b\}, \delta, q_0, Z, \{q_2\})$ , where  $\delta$  is defined as follows

$$\delta(q_0, a, Z) = \{(q_0, aZ)\}, \quad \delta(q_0, b, Z) = \{(q_0, bZ)\},$$

$$\delta(q_0, a, a) = \{(q_0, aa), (q_1, \epsilon)\}, \quad \delta(q_0, a, b) = \{(q_0, ab)\},$$

$$\delta(q_0, b, a) = \{(q_0, ba)\}, \quad \delta(q_0, b, b) = \{(q_0, bb), (q_1, \epsilon)\},$$

$$\delta(q_1, a, a) = \{(q_1, \epsilon)\}, \quad \delta(q_1, b, b) = \{(q_1, \epsilon)\},$$

$$\delta(q_1, \epsilon, Z) = \{(q_2, \epsilon)\}$$

- $(q_0, abba, Z) \vdash (q_0, bba, aZ) \vdash (q_0, ba, baZ) \vdash (q_1, a, aZ) \vdash (q_1, \epsilon, Z) \vdash (q_2, \epsilon, \epsilon)$
- $(q_0, aaa, Z) \vdash (q_0, aa, aZ) \vdash (q_0, a, aaZ) \vdash (q_1, \epsilon, aZ) \vdash \text{error}$
- $(q_0, aaa, Z) \vdash (q_0, aa, aZ) \vdash (q_1, a, Z) \vdash \text{error}$

# Nondeterministic and Deterministic PDA

- Just as in the case of NFA and DFA, PDA also have two versions: NPDA and DPDA
- However, NPDA are strictly more powerful than the DPDA
- For example, the language,  $L = \{ww^R \mid w \in \{a, b\}^+\}$  can be recognized only by an NPDA and not by any DPDA
- In the same breath, the language,  $L = \{wcw^R \mid w \in \{a, b\}^+\}$ , can be recognized by a DPDA
- In practice we need DPDA, since they have exactly one possible move at any instant
- Our parsers are all DPDA

- Parsing is the process of constructing a parse tree for a sentence generated by a given grammar
- If there are no restrictions on the language and the form of grammar used, parsers for context-free languages require  $O(n^3)$  time ( $n$  being the length of the string parsed)
  - Cocke-Younger-Kasami's algorithm
  - Earley's algorithm
- Subsets of context-free languages typically require  $O(n)$  time
  - Predictive parsing using  $LL(1)$  grammars (top-down parsing method)
  - Shift-Reduce parsing using  $LR(1)$  grammars (bottom-up parsing method)

# Top-Down Parsing using LL Grammars

- Top-down parsing using predictive parsing, traces the left-most derivation of the string while constructing the parse tree
- Starts from the start symbol of the grammar, and “predicts” the next production used in the derivation
- Such “prediction” is aided by parsing tables (constructed off-line)
- The next production to be used in the derivation is determined using the next input symbol to lookup the parsing table (look-ahead symbol)
- Placing restrictions on the grammar ensures that no slot in the parsing table contains more than one production
- At the time of parsing table construction, if two productions become eligible to be placed in the same slot of the parsing table, the grammar is declared unfit for predictive parsing

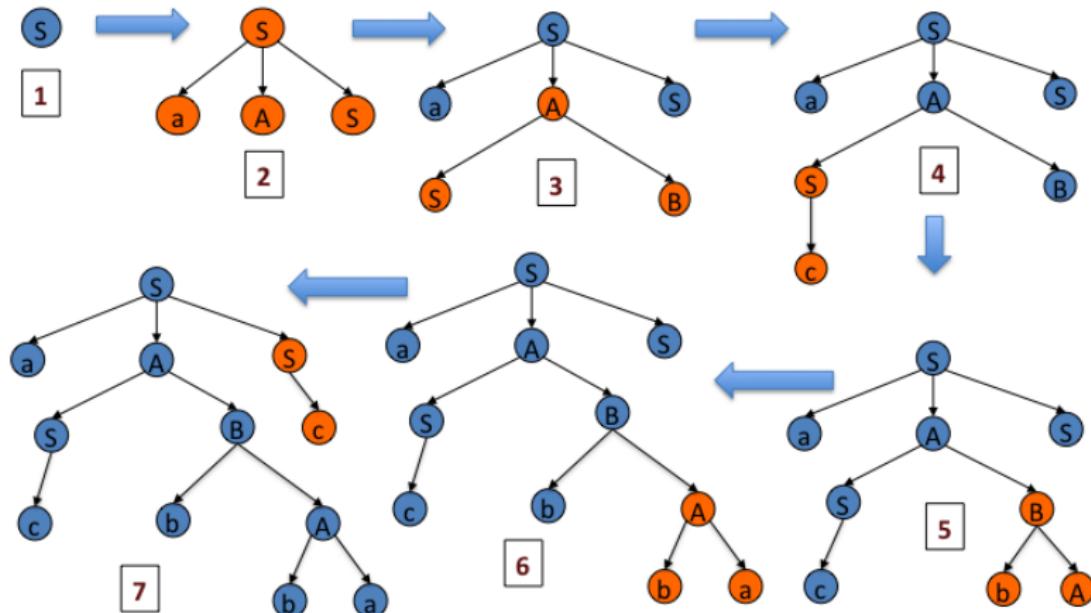
# Top-Down LL-Parsing Example

$S \rightarrow aAS \mid c$   
 $A \rightarrow ba \mid SB$   
 $B \rightarrow bA \mid S$

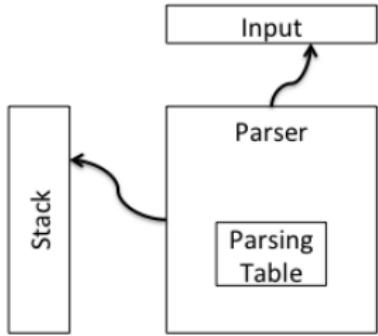
Leftmost derivation of the string *acbbac*

$S \Rightarrow aAS \Rightarrow aSBS \Rightarrow acBS \Rightarrow acbAS \Rightarrow acbbaS \Rightarrow acbbac$

1      2      3      4      5      6      7



# LL(1) Parsing Algorithm



Initial configuration: Stack =  $S$ , Input =  $w\$$ ,  
where,  $S$  = start symbol,  $\$$  = end of file marker  
repeat {  
    let  $X$  be the top stack symbol;  
    let  $a$  be the next input symbol /\*may be  $\$$ \*;/  
    if  $X$  is a terminal symbol or  $\$$  then  
        if  $X == a$  then {  
            pop  $X$  from Stack;  
            remove  $a$  from input;  
        } else ERROR();  
    else /\*  $X$  is a non-terminal symbol \*/  
        if  $M[X,a] == X \rightarrow Y_1Y_2\dots Y_k$  then {  
            pop  $X$  from Stack;  
            push  $Y_k, Y_{k-1}, \dots, Y_1$  onto Stack;  
            ( $Y_1$  on top)  
        }  
    }  
} until Stack has emptied;

# LL(1) Parsing Algorithm Example

Grammar

$$S' \rightarrow S\$$$

$$S \rightarrow aAS \mid c$$

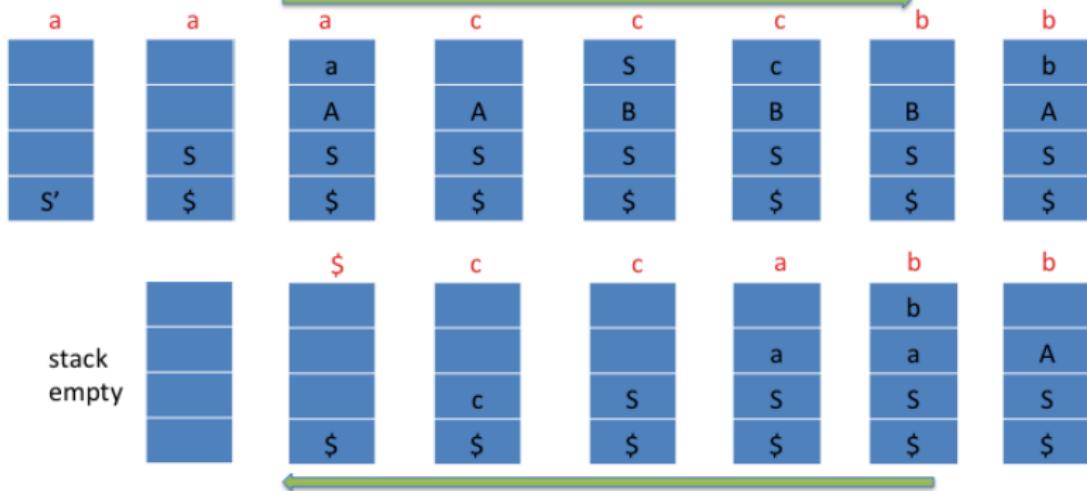
$$A \rightarrow ba \mid SB$$

$$B \rightarrow bA \mid S$$

string: acbbac

LL(1) Parsing Table

	a	b	c	\$
S'	$S' \rightarrow S\$$		$S' \rightarrow S\$$	
S	$S \rightarrow aAS$		$S \rightarrow c$	
A	$A \rightarrow SB$	$A \rightarrow ba$	$A \rightarrow SB$	
B	$B \rightarrow S$	$B \rightarrow bA$	$B \rightarrow S$	



# Strong LL(k) Grammars

Let the given grammar be G

- Input is extended with  $k$  symbols,  $\$^k$ ,  $k$  is the lookahead of the grammar
- Introduce a new nonterminal  $S'$ , and a production,  $S' \rightarrow S\$^k$ , where  $S$  is the start symbol of the given grammar
- Consider leftmost derivations only and assume that the grammar has no *useless symbols*
- A production  $A \rightarrow \alpha$  in  $G$  is called a *strong LL(k)* production, if in G
$$S' \Rightarrow^* wA\gamma \Rightarrow w\alpha\gamma \Rightarrow^* wz\gamma \\ S' \Rightarrow^* w'A\delta \Rightarrow w'\beta\delta \Rightarrow^* w'zx \\ |z| = k, z \in \Sigma^*, w \text{ and } w' \in \Sigma^*, \text{ then } \alpha = \beta$$
- A grammar (nonterminal) is strong LL( $k$ ) if all its productions are strong LL( $k$ )

## Strong LL(k) Grammars (contd.)

- Strong LL(k) grammars do not allow different productions of the same nonterminal to be used even in two different derivations, if the first  $k$  symbols of the strings produced by  $\alpha\gamma$  and  $\beta\delta$  are the same
- Example:  $S \rightarrow Abc|aAcb, A \rightarrow \epsilon|b|c$   
 $S$  is a strong LL(1) nonterminal
  - $S' \Rightarrow S\$ \Rightarrow Abc\$ \Rightarrow bc\$, bbc\$,$  and  $cbc\$,$  on application of the productions,  $A \rightarrow \epsilon,$   $A \rightarrow b,$  and,  $A \rightarrow c,$  respectively.  
 $z = b,$   $b,$  or  $c,$  respectively
  - $S' \Rightarrow S\$ \Rightarrow aAcb\$ \Rightarrow acb\$, abcb\$,$  and  $accb\$,$  on application of the productions,  $A \rightarrow \epsilon,$   $A \rightarrow b,$  and,  $A \rightarrow c,$  respectively.  $z = a,$  in all three cases
  - In this case,  $w = w' = \epsilon,$   $\alpha = Abc,$   $\beta = aAcb,$  but  $z$  is different in the two derivations, in all the derived strings
  - Hence the nonterminal  $S$  is strong LL(1)

# Strong LL(k) Grammars (contd.)

$A$  is not strong LL(1)

- $S' \Rightarrow^* Abc\$ \Rightarrow \underline{bc}\$, w = \epsilon, z = b, \alpha = \epsilon (A \rightarrow \epsilon)$   
 $S' \Rightarrow^* Abc\$ \Rightarrow \underline{bbc}\$, w' = \epsilon, z = b, \beta = b (A \rightarrow b)$
- Even though the lookaheads are the same ( $z = b$ ),  $\alpha \neq \beta$ , and therefore, the grammar is not strong LL(1)

$A$  is not strong LL(2)

- $S' \Rightarrow^* Abc\$ \Rightarrow \underline{bc}\$, w = \epsilon, z = bc, \alpha = \epsilon (A \rightarrow \epsilon)$   
 $S' \Rightarrow^* aAcb\$ \Rightarrow \underline{abcb}\$, w' = a, z = bc, \beta = b (A \rightarrow b)$
- Even though the lookaheads are the same ( $z = bc$ ),  $\alpha \neq \beta$ , and therefore, the grammar is not strong LL(2)

$A$  is strong LL(3) because all the six strings ( $bc\$, bbc, cbc, cb\$, bcb, ccb$ ) can be distinguished using 3-symbol lookahead  
(details are for home work)

# Testable Conditions for LL(1)

- We call strong LL(1) as LL(1) from now on and we will not consider lookaheads longer than 1
- The classical condition for LL(1) property uses  $FIRST$  and  $FOLLOW$  sets
- If  $\alpha$  is any string of grammar symbols ( $\alpha \in (N \cup T)^*$ ), then  
 $FIRST(\alpha) = \{a \mid a \in T, \text{ and } \alpha \Rightarrow^* ax, x \in T^*\}$   
 $FIRST(\epsilon) = \{\epsilon\}$
- If  $A$  is any nonterminal, then  
 $FOLLOW(A) = \{a \mid S \Rightarrow^* \alpha A a \beta, \alpha, \beta \in (N \cup T)^*, a \in T \cup \{\$\}$
- $FIRST(\alpha)$  is determined by  $\alpha$  alone, but  $FOLLOW(A)$  is determined by the “context” of  $A$ , i.e., the derivations in which  $A$  occurs

# *FIRST* and *FOLLOW* Computation Example

- Consider the following grammar

$$S' \rightarrow S\$, \quad S \rightarrow aAS \mid c, \quad A \rightarrow ba \mid SB, \quad B \rightarrow bA \mid S$$

- $\text{FIRST}(S') = \text{FIRST}(S) = \{a, c\}$  because

$$S' \Rightarrow S\$ \Rightarrow \underline{c\$}, \text{ and } S' \Rightarrow S\$ \Rightarrow \underline{aAS\$} \Rightarrow \underline{abaS\$} \Rightarrow \underline{abac\$}$$

- $\text{FIRST}(A) = \{a, b, c\}$  because

$A \Rightarrow \underline{ba}$ , and  $A \Rightarrow SB$ , and therefore all symbols in  $\text{FIRST}(S)$  are in  $\text{FIRST}(A)$

- $\text{FOLLOW}(S) = \{a, b, c, \$\}$  because

$$S' \Rightarrow \underline{S\$},$$
$$S' \Rightarrow^* aAS\$ \Rightarrow a\underline{SBS\$} \Rightarrow aS\underline{bAS\$},$$
$$S' \Rightarrow^* a\underline{SBS\$} \Rightarrow a\underline{SSS\$} \Rightarrow aS\underline{aASS\$},$$
$$S' \Rightarrow^* a\underline{SSS\$} \Rightarrow aS\underline{cS\$}$$

- $\text{FOLLOW}(A) = \{a, c\}$  because

$$S' \Rightarrow^* a\underline{AS\$} \Rightarrow aA\underline{aAS\$},$$
$$S' \Rightarrow^* a\underline{AS\$} \Rightarrow aA\underline{c}$$

# Computation of FIRST: Terminals and Nonterminals

```
{  
    for each ( $a \in T$ )  $\text{FIRST}(a) = \{a\}$ ;  $\text{FIRST}(\epsilon) = \{\epsilon\}$ ;  
    for each ( $A \in N$ )  $\text{FIRST}(A) = \emptyset$ ;  
    while ( $\text{FIRST}$  sets are still changing) {  
        for each production  $p$  {  
            Let  $p$  be the production  $A \rightarrow X_1 X_2 \dots X_n$ ;  
             $\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(X_1) - \{\epsilon\})$ ;  
             $i = 1$ ;  
            while ( $\epsilon \in \text{FIRST}(X_i)$   $\&&$   $i \leq n - 1$ ) {  
                 $\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(X_{i+1}) - \{\epsilon\})$ ;  $i++$ ;  
            }  
            if ( $i == n$ )  $\&&$  ( $\epsilon \in \text{FIRST}(X_n)$ )  
                 $\text{FIRST}(A) = \text{FIRST}(A) \cup \{\epsilon\}$   
        }  
    }  
}
```

# Syntax Analysis:

Context-free Grammars, Pushdown Automata and Parsing  
Part - 3

Y.N. Srikant

Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis? (covered in lecture 1)
- Specification of programming languages: context-free grammars (covered in lecture 1)
- Parsing context-free languages: push-down automata (covered in lectures 1 and 2)
- Top-down parsing: LL(1) and recursive-descent parsing
- Bottom-up parsing: LR-parsing

# Testable Conditions for LL(1)

- We call strong LL(1) as LL(1) from now on and we will not consider lookaheads longer than 1
- The classical condition for LL(1) property uses  $FIRST$  and  $FOLLOW$  sets
- If  $\alpha$  is any string of grammar symbols ( $\alpha \in (N \cup T)^*$ ), then  
 $FIRST(\alpha) = \{a \mid a \in T, \text{ and } \alpha \Rightarrow^* ax, x \in T^*\}$   
 $FIRST(\epsilon) = \{\epsilon\}$
- If  $A$  is any nonterminal, then  
 $FOLLOW(A) = \{a \mid S \Rightarrow^* \alpha A a \beta, \alpha, \beta \in (N \cup T)^*, a \in T \cup \{\$\}$
- $FIRST(\alpha)$  is determined by  $\alpha$  alone, but  $FOLLOW(A)$  is determined by the “context” of  $A$ , i.e., the derivations in which  $A$  occurs

# *FIRST* and *FOLLOW* Computation Example

- Consider the following grammar

$$S' \rightarrow S\$, \quad S \rightarrow aAS \mid c, \quad A \rightarrow ba \mid SB, \quad B \rightarrow bA \mid S$$

- $\text{FIRST}(S') = \text{FIRST}(S) = \{a, c\}$  because

$$S' \Rightarrow S\$ \Rightarrow \underline{c\$}, \text{ and } S' \Rightarrow S\$ \Rightarrow \underline{aAS\$} \Rightarrow \underline{abaS\$} \Rightarrow \underline{abac\$}$$

- $\text{FIRST}(A) = \{a, b, c\}$  because

$A \Rightarrow \underline{ba}$ , and  $A \Rightarrow SB$ , and therefore all symbols in  $\text{FIRST}(S)$  are in  $\text{FIRST}(A)$

- $\text{FOLLOW}(S) = \{a, b, c, \$\}$  because

$$S' \Rightarrow \underline{S\$},$$
$$S' \Rightarrow^* aAS\$ \Rightarrow a\underline{SBS\$} \Rightarrow aS\underline{bAS\$},$$
$$S' \Rightarrow^* a\underline{SBS\$} \Rightarrow a\underline{SSS\$} \Rightarrow aS\underline{aASS\$},$$
$$S' \Rightarrow^* a\underline{SSS\$} \Rightarrow aS\underline{cS\$}$$

- $\text{FOLLOW}(A) = \{a, c\}$  because

$$S' \Rightarrow^* a\underline{AS\$} \Rightarrow aA\underline{aAS\$},$$
$$S' \Rightarrow^* a\underline{AS\$} \Rightarrow aA\underline{c}$$

# Computation of FIRST: Terminals and Nonterminals

```
{  
    for each ( $a \in T$ )  $\text{FIRST}(a) = \{a\}$ ;  $\text{FIRST}(\epsilon) = \{\epsilon\}$ ;  
    for each ( $A \in N$ )  $\text{FIRST}(A) = \emptyset$ ;  
    while ( $\text{FIRST}$  sets are still changing) {  
        for each production  $p$  {  
            Let  $p$  be the production  $A \rightarrow X_1 X_2 \dots X_n$ ;  
             $\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(X_1) - \{\epsilon\})$ ;  
             $i = 1$ ;  
            while ( $\epsilon \in \text{FIRST}(X_i)$   $\&&$   $i \leq n - 1$ ) {  
                 $\text{FIRST}(A) = \text{FIRST}(A) \cup (\text{FIRST}(X_{i+1}) - \{\epsilon\})$ ;  $i++$ ;  
            }  
            if ( $i == n$ )  $\&&$  ( $\epsilon \in \text{FIRST}(X_n)$ )  
                 $\text{FIRST}(A) = \text{FIRST}(A) \cup \{\epsilon\}$   
        }  
    }  
}
```

# Computation of $FIRST(\beta)$ : $\beta$ , a string of Grammar Symbols

```
{ /* It is assumed that FIRST sets of terminals and nonterminals  
   are already available */  
FIRST( $\beta$ ) =  $\emptyset$ ;  
while (FIRST sets are still changing) {  
    Let  $\beta$  be the string  $X_1 X_2 \dots X_n$ ;  
    FIRST( $\beta$ ) = FIRST( $\beta$ )  $\cup$  (FIRST( $X_1$ ) -  $\{\epsilon\}$ );  
     $i = 1$ ;  
    while ( $\epsilon \in FIRST(X_i)$  &&  $i \leq n - 1$ ) {  
        FIRST( $\beta$ ) = FIRST( $\beta$ )  $\cup$  (FIRST( $X_{i+1}$  -  $\{\epsilon\}$ ));  $i++$ ;  
    }  
    if ( $i == n$ ) && ( $\epsilon \in FIRST(X_n)$ )  
        FIRST( $\beta$ ) = FIRST( $\beta$ )  $\cup \{\epsilon\}$   
    }  
}
```

# FIRST Computation: Algorithm Trace - 1

- Consider the following grammar

$$S' \rightarrow S\$, \quad S \rightarrow aAS \mid \epsilon, \quad A \rightarrow ba \mid SB, \quad B \rightarrow cA \mid S$$

- Initially,  $\text{FIRST}(S) = \text{FIRST}(A) = \text{FIRST}(B) = \emptyset$

- Iteration 1

- $\text{FIRST}(S) = \{a, \epsilon\}$  from the productions  $S \rightarrow aAS \mid \epsilon$
- $\text{FIRST}(A) = \{b\} \cup \text{FIRST}(S) - \{\epsilon\} \cup \text{FIRST}(B) - \{\epsilon\} = \{b, a\}$   
from the productions  $A \rightarrow ba \mid SB$   
(since  $\epsilon \in \text{FIRST}(S)$ ,  $\text{FIRST}(B)$  is also included;  
since  $\text{FIRST}(B)=\phi$ ,  $\epsilon$  is not included)
- $\text{FIRST}(B) = \{c\} \cup \text{FIRST}(S) - \{\epsilon\} \cup \{\epsilon\} = \{c, a, \epsilon\}$   
from the productions  $B \rightarrow cA \mid S$   
( $\epsilon$  is included because  $\epsilon \in \text{FIRST}(S)$ )

# FIRST Computation: Algorithm Trace - 2

- The grammar is  
 $S' \rightarrow S\$, S \rightarrow aAS \mid \epsilon, A \rightarrow ba \mid SB, B \rightarrow cA \mid S$
- From the first iteration,  
 $\text{FIRST}(S) = \{a, \epsilon\}, \text{FIRST}(A) = \{b, a\}, \text{FIRST}(B) = \{c, a, \epsilon\}$
- Iteration 2  
(values stabilize and do not change in iteration 3)
  - $\text{FIRST}(S) = \{a, \epsilon\}$  (no change from iteration 1)
  - $\text{FIRST}(A) = \{b\} \cup \text{FIRST}(S) - \{\epsilon\} \cup \text{FIRST}(B) - \{\epsilon\} \cup \{\epsilon\}$   
 $= \{b, a, c, \epsilon\}$  (changed!)
  - $\text{FIRST}(B) = \{c, a, \epsilon\}$  (no change from iteration 1)

# Computation of FOLLOW

```
{ for each ( $X \in N \cup T$ )  $\text{FOLLOW}(X) = \emptyset$ ;  
     $\text{FOLLOW}(S) = \{\$\}$ ; /*  $S$  is the start symbol of the grammar */  
repeat {  
    for each production  $A \rightarrow X_1 X_2 \dots X_n$  /*  $X_i \neq \epsilon$  */  
         $\text{FOLLOW}(X_n) = \text{FOLLOW}(X_n) \cup \text{FOLLOW}(A)$ ;  
        REST =  $\text{FOLLOW}(A)$ ;  
        for  $i = n$  downto 2 {  
            if ( $\epsilon \in \text{FIRST}(X_i)$ ) {  $\text{FOLLOW}(X_{i-1}) =$   
                 $\text{FOLLOW}(X_{i-1}) \cup (\text{FIRST}(X_i) - \{\epsilon\}) \cup \text{REST}$ ;  
                REST =  $\text{FOLLOW}(X_{i-1})$ ;  
            } else {  $\text{FOLLOW}(X_{i-1}) = \text{FOLLOW}(X_{i-1}) \cup \text{FIRST}(X_i)$  ;  
                REST =  $\text{FOLLOW}(X_{i-1})$ ;  
            }  
        }  
    }  
} until no FOLLOW set has changed  
}
```

# FOLLOW Computation: Algorithm Trace

- Consider the following grammar

$$S' \rightarrow S\$, \quad S \rightarrow aAS \mid \epsilon, \quad A \rightarrow ba \mid SB, \quad B \rightarrow cA \mid S$$

- Initially,  $follow(S) = \{\$\}$ ;  $follow(A) = follow(B) = \emptyset$   
 $first(S) = \{a, \epsilon\}$ ;  $first(A) = \{a, b, c, \epsilon\}$ ;  $first(B) = \{a, c, \epsilon\}$ ;
- Iteration 1 /\* In the following,  $x \cup = y$  means  $x = x \cup y$  \*/

- $S \rightarrow aAS$ :  $follow(S) \cup = \{\$\}$ ;  $rest = follow(S) = \{\$\}$   
 $follow(A) \cup = (first(S) - \{\epsilon\}) \cup rest = \{a, \$\}$
- $A \rightarrow SB$ :  $follow(B) \cup = follow(A) = \{a, \$\}$   
 $rest = follow(A) = \{a, \$\}$   
 $follow(S) \cup = (first(B) - \{\epsilon\}) \cup rest = \{a, c, \$\}$
- $B \rightarrow cA$ :  $follow(A) \cup = follow(B) = \{a, \$\}$
- $B \rightarrow S$ :  $follow(S) \cup = follow(B) = \{a, c, \$\}$
- At the end of iteration 1  
 $follow(S) = \{a, c, \$\}$ ;  $follow(A) = follow(B) = \{a, \$\}$

## FOLLOW Computation: Algorithm Trace (contd.)

- $first(S) = \{a, \epsilon\}$ ;  $first(A) = \{a, b, c, \epsilon\}$ ;  $first(B) = \{a, c, \epsilon\}$ ;
- At the end of iteration 1  
 $follow(S) = \{a, c, \$\}$ ;  $follow(A) = follow(B) = \{a, \$\}$
- Iteration 2
  - $S \rightarrow aAS$ :  $follow(S) \cup = \{a, c, \$\}$ ;  
 $rest = follow(S) = \{a, c, \$\}$   
 $follow(A) \cup = (first(S) - \{\epsilon\}) \cup rest = \{a, c, \$\}$  (changed!)
  - $A \rightarrow SB$ :  $follow(B) \cup = follow(A) = \{a, c, \$\}$  (changed!)  
 $rest = follow(A) = \{a, c, \$\}$   
 $follow(S) \cup = (first(B) - \{\epsilon\}) \cup rest = \{a, c, \$\}$  (no change)
- At the end of iteration 2  
 $follow(S) = follow(A) = follow(B) = \{a, c, \$\}$ ;
- The  $follow$  sets do not change any further

# LL(1) Conditions

- Let  $G$  be a context-free grammar
- $G$  is LL(1) iff for every pair of productions  $A \rightarrow \alpha$  and  $A \rightarrow \beta$ , the following condition holds
  - $\text{dirsymbol}(\alpha) \cap \text{dirsymbol}(\beta) = \emptyset$ , where  
 $\text{dirsymbol}(\gamma) = \text{if } (\epsilon \in \text{first}(\gamma)) \text{ then } ((\text{first}(\gamma) - \{\epsilon\}) \cup \text{follow}(A)) \text{ else } \text{first}(\gamma)$   
( $\gamma$  stands for  $\alpha$  or  $\beta$ )
  - $\text{dirsymbol}$  stands for “direction symbol set”
- An equivalent formulation (as in ALSU’s book) is as below
  - $\text{first}(\alpha.\text{follow}(A)) \cap \text{first}(\beta.\text{follow}(A)) = \emptyset$
- Construction of the LL(1) parsing table

for each production  $A \rightarrow \alpha$

for each symbol  $s \in \text{dirsymbol}(\alpha)$

/\*  $s$  may be either a terminal symbol or  $\$$  \*/

add  $A \rightarrow \alpha$  to  $LLPT[A, s]$

Make each undefined entry of  $LLPT$  as *error*

# LL(1) Table Construction using *FIRST* and *FOLLOW*

for each production  $A \rightarrow \alpha$

    for each terminal symbol  $a \in \text{first}(\alpha)$

        add  $A \rightarrow \alpha$  to  $\text{LLPT}[A, a]$

    if  $\epsilon \in \text{first}(\alpha)$  {

        for each terminal symbol  $b \in \text{follow}(A)$

            add  $A \rightarrow \alpha$  to  $\text{LLPT}[A, b]$

        if  $\$ \in \text{follow}(A)$

            add  $A \rightarrow \alpha$  to  $\text{LLPT}[A, \$]$

}

Make each undefined entry of  $\text{LLPT}$  as *error*

- After the construction of the LL(1) table is complete (following any of the two methods), if any slot in the LL(1) table has two or more productions, then the grammar is NOT LL(1)

# Simple Example of LL(1) Grammar

- P1:  $S \rightarrow \text{if } (a) S \text{ else } S \mid \text{while } (a) S \mid \text{begin } SL \text{ end}$   
P2:  $SL \rightarrow S \ S'$   
P3:  $S' \rightarrow ; \ SL \mid \epsilon$
- {if, while, begin, end, a, (, ), ;} are all terminal symbols
- Clearly, all alternatives of P1 start with distinct symbols and hence create no problem
- P2 has no choices
- Regarding P3,  $\text{dirsymbol}(;SL) = \{;\}$ , and  $\text{dirsymbol}(\epsilon) = \{\text{end}\}$ , and the two have no common symbols
- Hence the grammar is LL(1)

# LL(1) Table Construction Example 1

LL(1) Parsing Table for the original grammar

	if	id	else	a	\$
S'	$S' \rightarrow S\$$			$S' \rightarrow S\$$	
S	$S \rightarrow \text{if id } S$ $S \rightarrow \text{if id } S \text{ else } S$			$S \rightarrow a$	

Original Grammar

Grammar is not LL(1)

$S' \rightarrow S\$$   
 $S \rightarrow \text{if id } S \mid$   
 $\quad \text{if id } S \text{ else } S \mid$   
 $\quad a$

tokens: if, id, else, a

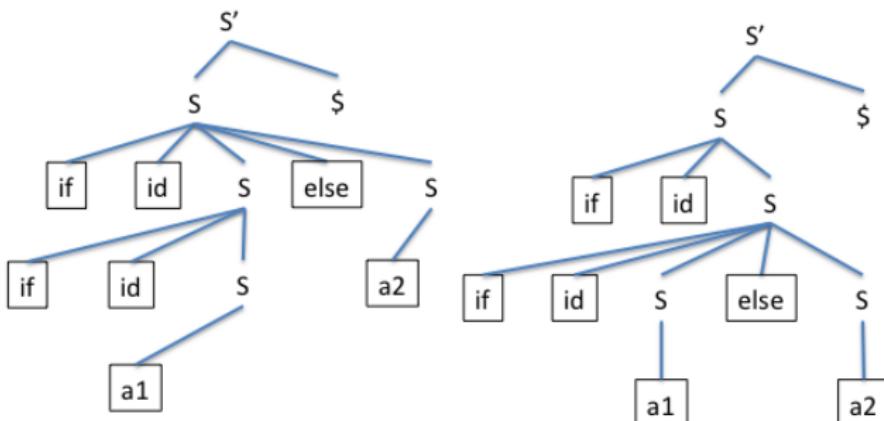
$\text{dirsymbol}(S\$) = \{\text{if}, \text{id}\}; \text{dirsymbol}(a) = \{a\}$   
 $\text{dirsymbol}(\text{if id } S) = \{\text{if}\}$   
 $\text{dirsymbol}(\text{if id } S \text{ else } S) = \{\text{if}\}$

$$\text{dirsymbol}(\text{if id } S) \cap \text{dirsymbol}(a) = \emptyset$$

$$\text{dirsymbol}(\text{if id } S \text{ else } S) \cap \text{dirsymbol}(a) = \emptyset$$

$$\text{dirsymbol}(\text{if id } S) \cap \text{dirsymbol}(\text{if id } S \text{ else } S) \neq \emptyset$$

# LL(1) Table Problem Example 1



string: if id ( if id a1 ) else a2

parentheses are not part of the string

string: if id (if id a1 else a2)

parentheses are not part of the string

# LL(1) Table Construction Example 2

Original Grammar		LL(1) Parsing Table for modified grammar			
		if	else	a	\$
$S' \rightarrow S\$$					
$S \rightarrow \text{if id } S \mid$	$S' \rightarrow S\$$			$S' \rightarrow S\$$	
$\text{if id } S \text{ else } S \mid$	$S \rightarrow \text{if id } S \ S1$			$S \rightarrow a$	
$a$			$S1 \rightarrow \epsilon$ $S1 \rightarrow \text{else } S$		$S1 \rightarrow \epsilon$

$\text{dirsymbol}(S\$) = \{\text{if}, \text{a}\}; \text{dirsymbol}(a) = \{a\}$   
 $\text{dirsymbol}(\text{if id } S \ S1) = \{\text{if}\}$   
 $\text{dirsymbol}(\text{else } S) = \{\text{else}\}$   
 $\text{dirsymbol}(\epsilon) = \{\text{else, \$}\}$

Grammar is not LL(1)

Left-Factored Grammar

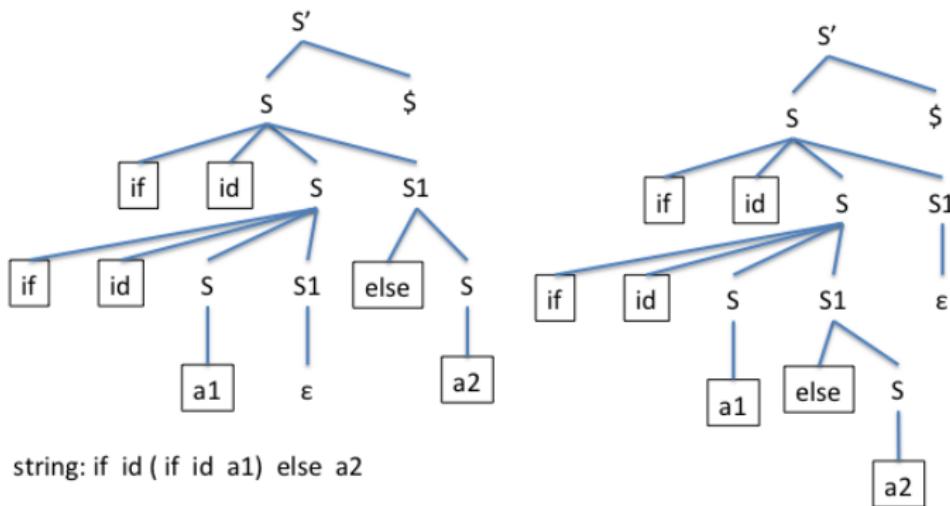
$S' \rightarrow S\$$   
 $S \rightarrow \text{if id } S \ S1 \mid a$   
 $S1 \rightarrow \epsilon \mid \text{else } S$

tokens: if, id, else, a

$$\text{dirsymbol}(\text{if id } S \ S1) \cap \text{dirsymbol}(a) = \emptyset$$

$$\text{dirsymbol}(\epsilon) \cap \text{dirsymbol}(\text{else } S) \neq \emptyset$$

# LL(1) Table Problem Example 2



string: if id ( if id a1) else a2

parentheses are not part of the string

string: if id (if id a1 else a2)

parentheses are not part of the string

# LL(1) Table Construction Example 3

$S' \rightarrow S\$$

$S \rightarrow aAS \mid c$

$A \rightarrow ba \mid SB$

$B \rightarrow bA \mid S$

Grammar is LL(1)

	a	b	c	\$
S'	$S' \rightarrow S\$$		$S' \rightarrow S\$$	
S	$S \rightarrow aAS$		$S \rightarrow c$	
A	$A \rightarrow SB$	$A \rightarrow ba$	$A \rightarrow SB$	
B	$B \rightarrow S$	$B \rightarrow bA$	$B \rightarrow S$	

$\text{first}(S) = \{a, c\}$

$\text{first}(A) = \{a, b, c\}$

$\text{first}(B) = \{a, b, c\}$

$\text{dirsymbol}(aAS) \cap \text{dirsymbol}(c) = \emptyset$

$\text{dirsymbol}(ba) \cap \text{dirsymbol}(SB) = \emptyset$

$\text{dirsymbol}(bA) \cap \text{dirsymbol}(S) = \emptyset$

$\text{follow}(S) = \{a, b, c, \$\}$

$\text{follow}(A) = \{a, c\}$

$\text{follow}(B) = \{a, c\}$

$\text{dirsymbol}(S\$) = \{a, c\}$

$\text{dirsymbol}(aAS) = \{a\}$

$\text{dirsymbol}(c) = \{c\}$

$\text{dirsymbol}(ba) = \{b\}$

$\text{dirsymbol}(SB) = \{a, c\}$

$\text{dirsymbol}(bA) = \{b\}$

$\text{dirsymbol}(S) = \{a, c\}$

# LL(1) Table Construction Example 4

Left-Recursive Grammar  
for Statement List

$$\begin{aligned}S' &\rightarrow SL \$ \\SL &\rightarrow SL S \mid S \\S &\rightarrow a\end{aligned}$$

$$\begin{aligned}\text{dirsymbol}(SL \$) &= \{a\} \\ \text{dirsymbol } (a) &= \{a\} \\ \text{dirsymbol}(SL S) &= \{a\} \\ \text{dirsymbol}(S) &= \{a\}\end{aligned}$$

$$\text{dirsymbol}(SL S) \cap \text{dirsymbol}(S) \neq \emptyset$$

$$\begin{aligned}\text{dirsymbol}(SL \$) &= \{a\} \\ \text{dirsymbol } (a) &= \{a\} \\ \text{dirsymbol}(S A) &= \{a\} \\ \text{dirsymbol}(\epsilon) &= \{\$\}\end{aligned}$$

LL(1) Parsing Table for  
Left-Recursive Grammar

	a
S'	$S' \rightarrow SL \$$
SL	$SL \rightarrow SL S$ $SL \rightarrow S$
S	$S \rightarrow a$

Grammar is not LL(1)

Right-Recursive Grammar  
for Statement List

$$\begin{aligned}S' &\rightarrow SL \$ \\SL &\rightarrow S A \\A &\rightarrow S A \mid \epsilon \\A &\rightarrow a\end{aligned}$$

LL(1) Parsing Table for  
Right-Recursive Grammar

	a	\$
S'	$S' \rightarrow SL \$$	
SL	$SL \rightarrow S A$	
A	$A \rightarrow S A$	$A \rightarrow \epsilon$
S	$S \rightarrow a$	

$$\text{dirsymbol}(S A) \cap \text{dirsymbol}(\epsilon) = \emptyset$$

# Elimination of Useless Symbols

Now we study the *grammar transformations*, elimination of useless symbols, elimination of left recursion and left factoring

- Given a grammar  $G = (N, T, P, S)$ , a non-terminal  $X$  is *useful* if  $S \Rightarrow^* \alpha X \beta \Rightarrow^* w$ , where,  $w \in T^*$   
Otherwise,  $X$  is useless
- Two conditions have to be met to ensure that  $X$  is useful
  - $X \Rightarrow^* w$ ,  $w \in T^*$  ( $X$  derives some terminal string)
  - $S \Rightarrow^* \alpha X \beta$  ( $X$  occurs in some string derivable from  $S$ )
- Example:  $S \rightarrow AB \mid CA$ ,  $B \rightarrow BC \mid AB$ ,  $A \rightarrow a$ ,  
 $C \rightarrow aB \mid b$ ,  $D \rightarrow d$ 
  - $A \rightarrow a$ ,  $C \rightarrow b$ ,  $D \rightarrow d$ ,  $S \rightarrow CA$
  - $S \rightarrow CA$ ,  $A \rightarrow a$ ,  $C \rightarrow b$

# Testing for $X \Rightarrow^* w$

$G' = (N', T', P', S')$  is the new grammar

$N_{OLD} = \phi;$

$N_{NEW} = \{X \mid X \rightarrow w, w \in T^*\}$

while  $N_{OLD} \neq N_{NEW}$  do {

$N_{OLD} = N_{NEW};$

$N_{NEW} = N_{OLD} \cup \{X \mid X \rightarrow \alpha, \alpha \in (T \cup N_{OLD})^*\}$

}

$N' = N_{NEW}; T' = T; S' = S;$

$P' = \{p \mid \text{all symbols of } p \text{ are in } N' \cup T'\}$

# Testing for $S \Rightarrow^* \alpha X \beta$

$G' = (N', T', P', S')$  is the new grammar

$N' = \{S\}$ ;

Repeat {

for each production  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$  with  $A \in N'$  do

    add all nonterminals of  $\alpha_1, \alpha_2, \dots, \alpha_n$  to  $N'$  and

    all terminals of  $\alpha_1, \alpha_2, \dots, \alpha_n$  to  $T'$

} until there is no change in  $N'$  and  $T'$

$P' = \{p \mid \text{all symbols of } p \text{ are in } N' \cup T'\}; S' = S$

# Syntax Analysis:

Context-free Grammars, Pushdown Automata and Parsing  
Part - 4

Y.N. Srikant

Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis? (covered in lecture 1)
- Specification of programming languages: context-free grammars (covered in lecture 1)
- Parsing context-free languages: push-down automata (covered in lectures 1 and 2)
- Top-down parsing: LL(1) parsing (covered in lectures 2 and 3)
- Recursive-descent parsing
- Bottom-up parsing: LR-parsing

# Elimination of Left Recursion

- A *left-recursive* grammar has a non-terminal  $A$  such that  $A \Rightarrow^+ A\alpha$
- Top-down parsing methods (LL(1) and RD) cannot handle left-recursive grammars
- Left-recursion in grammars can be eliminated by transformations
- A simpler case is that of grammars with *immediate left recursion*, where there is a production of the form  $A \rightarrow A\alpha$ 
  - Two productions  $A \rightarrow A\alpha \mid \beta$  can be transformed to  $A \rightarrow \beta A', A' \rightarrow \alpha A' \mid \epsilon$
  - In general, a group of productions:  
$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$
can be transformed to  
$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A', A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

# Left Recursion Elimination - An Example

$$A \rightarrow A\alpha \mid \beta \Rightarrow A \rightarrow \beta A', A' \rightarrow \alpha A' \mid \epsilon$$

- The following grammar for regular expressions is ambiguous:

$$E \rightarrow E + E \mid E E \mid E^* \mid (E) \mid a \mid b$$

- Equivalent left-recursive but unambiguous grammar is:

$$E \rightarrow E + T \mid T, T \rightarrow T F \mid F, F \rightarrow F^* \mid P, P \rightarrow (E) \mid a \mid b$$

- Equivalent non-left-recursive grammar is:

$$\begin{aligned} E &\rightarrow T E', E' \rightarrow +T E' \mid \epsilon, T \rightarrow F T', T' \rightarrow F T' \mid \epsilon, \\ F &\rightarrow P F', F' \rightarrow *F' \mid \epsilon, P \rightarrow (E) \mid a \mid b \end{aligned}$$

# Left Factoring

- If two alternatives of a production begin with the same string, then the grammar is not LL(1)
- Example:  $S \rightarrow 0S1 \mid 01$  is not LL(1)
  - After left factoring:  $S \rightarrow 0S'$ ,  $S' \rightarrow S1 \mid 1$  is LL(1)
- General method:  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \Rightarrow A \rightarrow \alpha A'$ ,  $A' \rightarrow \beta_1 \mid \beta_2$
- Another example: a grammar for logical expressions is given below

$E \rightarrow T \text{ or } E \mid T$ ,  $T \rightarrow F \text{ and } T \mid F$ ,  
 $F \rightarrow \text{not } F \mid (E) \mid \text{true} \mid \text{false}$

- This grammar is not LL(1) but becomes LL(1) after left factoring
- $E \rightarrow TE'$ ,  $E' \rightarrow \text{or } E \mid \epsilon$ ,  $T \rightarrow FT'$ ,  $T' \rightarrow \text{and } T \mid \epsilon$ ,  
 $F \rightarrow \text{not } F \mid (E) \mid \text{true} \mid \text{false}$

# Grammar Transformations may not help!

Original Grammar

$$\begin{aligned}S' &\rightarrow S\$ \\S &\rightarrow \text{if id } S \mid \\&\quad \text{if id } S \text{ else } S \mid \\&\quad a\end{aligned}$$

LL(1) Parsing Table for modified grammar

	if	else	a	\$
S'	$S' \rightarrow S\$$		$S' \rightarrow S\$$	
S	$S \rightarrow \text{if id } S S1$		$S \rightarrow a$	
S1		$S1 \rightarrow \epsilon$ $S1 \rightarrow \text{else } S$		$S1 \rightarrow \epsilon$

$\text{dirsymbol}(S\$) = \{\text{if, a}\}; \text{dirsymbol}(a) = \{a\}$   
 $\text{dirsymbol}(\text{if id } S S1) = \{\text{if}\}$   
 $\text{dirsymbol}(\text{else } S) = \{\text{else}\}$   
 $\text{dirsymbol}(\epsilon) = \{\text{else, \$}\}$

Grammar is not LL(1)

Left-Factored Grammar

$$\begin{aligned}S' &\rightarrow S\$ \\S &\rightarrow \text{if id } S S1 \mid a \\S1 &\rightarrow \epsilon \mid \text{else } S\end{aligned}$$

tokens: if, id, else, a

$\text{dirsymbol}(\text{if id } S S1) \cap \text{dirsymbol}(a) = \emptyset$

$\text{dirsymbol}(\epsilon) \cap \text{dirsymbol}(\text{else } S) \neq \emptyset$

Choose  $S1 \rightarrow \text{else } S$  instead of  $S1 \rightarrow \epsilon$  on lookahead *else*.

This resolves the conflict. Associates *else* with the innermost *if*

# Recursive-Descent Parsing

- Top-down parsing strategy
- One function/procedure for each nonterminal
- Functions call each other recursively, based on the grammar
- Recursion stack handles the tasks of LL(1) parser stack
- LL(1) conditions to be satisfied for the grammar
- Can be automatically generated from the grammar
- Hand-coding is also easy
- Error recovery is superior

# An Example

Grammar:  $S' \rightarrow S\$, S \rightarrow aAS \mid c, A \rightarrow ba \mid SB, B \rightarrow bA \mid S$

```
/* function for nonterminal S' */
void main() /* S' --> \$ */
    fS(); if (token == eof) accept();
        else error();
}
/* function for nonterminal S */
void fS() /* S --> aAS | c */
switch token {
    case a : get_token(); fA(); fS();
                break;
    case c : get_token(); break;
    others : error();
}
}
```

## An Example (contd.)

```
void fA() /* A --> ba | SB */
    switch token {
        case b : get_token();
                    if (token == a) get_token();
                    else error(); break;
        case a,c : fS(); fB(); break;
        others : error();
    }
}

void fB() /* B --> bA | S */
    switch token {
        case b : get_token(); fA(); break;
        case a,c : fS(); break;
        others : error();
    }
}
```

- Scheme is based on structure of productions
- Grammar must satisfy LL(1) conditions
- function `get_token()` obtains the next token from the lexical analyzer and places it in the global variable `token`
- function `error()` prints out a suitable error message
- In the next slide, for each grammar component, the code that must be generated is shown

# Automatic Generation of RD Parsers (contd.)

- ①  $\epsilon : ;$
- ②  $a \in T : \text{if } (\text{token} == a) \text{ get\_token(); else error();}$
- ③  $A \in N : fA(); /* \text{function call for nonterminal } A */$
- ④  $\alpha_1 | \alpha_2 | \dots | \alpha_n :$ 

```
switch token {
    case dirsym( $\alpha_1$ ): program_segment( $\alpha_1$ ); break;
    case dirsym( $\alpha_2$ ): program_segment( $\alpha_2$ ); break;
    ...
    others: error();
}
```
- ⑤  $\alpha_1 \alpha_2 \dots \alpha_n :$ 

```
program_segment( $\alpha_1$ ); program_segment( $\alpha_2$ ); ... ;
program_segment( $\alpha_n$ );
```
- ⑥  $A \rightarrow \alpha : \text{void } fA() \{ \text{program\_segment}(\alpha); \}$

# Bottom-Up Parsing

- Begin at the leaves, build the parse tree in small segments, combine the small trees to make bigger trees, until the root is reached
- This process is called *reduction* of the sentence to the start symbol of the grammar
- One of the ways of “reducing” a sentence is to follow the rightmost derivation of the sentence in *reverse*
  - *Shift-Reduce* parsing implements such a strategy
  - It uses the concept of a *handle* to detect when to perform reductions

- **Handle:** A *handle* of a right sentential form  $\gamma$ , is a production  $A \rightarrow \beta$  and a position in  $\gamma$ , where the string  $\beta$  may be found and replaced by  $A$ , to produce the previous right sentential form in a rightmost derivation of  $\gamma$ .  
That is, if  $S \Rightarrow_{rm}^* \alpha Aw \Rightarrow_{rm} \alpha \beta w$ , then  $A \rightarrow \beta$  in the position following  $\alpha$  is a handle of  $\alpha \beta w$
- A handle will always eventually appear on the top of the stack, never submerged inside the stack
- In S-R parsing, we locate the handle and reduce it by the LHS of the production repeatedly, to reach the start symbol
- These reductions, in fact, trace out a rightmost derivation of the sentence in reverse. This process is called handle pruning
- *LR-Parsing* is a method of shift-reduce parsing

# Examples

- 1  $S \rightarrow aAcBe, A \rightarrow Ab \mid b, B \rightarrow d$

For the string =  $abbcde$ , the rightmost derivation marked with handles is shown below

$S \Rightarrow \underline{aAcBe}$  ( $aAcBe, S \rightarrow aAcBe$ )  
 $\Rightarrow \underline{aAcde}$  ( $d, B \rightarrow d$ )  
 $\Rightarrow \underline{aAbcde}$  ( $Ab, A \rightarrow Ab$ )  
 $\Rightarrow \underline{abbcde}$  ( $b, A \rightarrow b$ )

The handle is unique if the grammar is unambiguous!

## Examples (contd.)

②  $S \rightarrow aAS \mid c, A \rightarrow ba \mid SB, B \rightarrow bA \mid S$

For the string =  $acbbac$ , the rightmost derivation marked with handles is shown below

$$\begin{aligned} S &\Rightarrow \underline{aAS} \text{ (aAS, } S \rightarrow aAS) \\ &\Rightarrow \underline{aAc} \text{ (c, } S \rightarrow c) \\ &\Rightarrow \underline{aSBc} \text{ (SB, } A \rightarrow SB) \\ &\Rightarrow \underline{aSbAc} \text{ (bA, } B \rightarrow bA) \\ &\Rightarrow \underline{aSbbac} \text{ (ba, } A \rightarrow ba) \\ &\Rightarrow \underline{acbbac} \text{ (c, } S \rightarrow c) \end{aligned}$$

## Examples (contd.)

③  $E \rightarrow E + E, E \rightarrow E * E, E \rightarrow (E), E \rightarrow id$

For the string =  $id + id * id$ , two rightmost derivation marked with handles are shown below

$$E \Rightarrow \underline{E + E} \quad (E + E, \quad E \rightarrow E + E)$$

$$\Rightarrow E + \underline{E * E} \quad (E * E, \quad E \rightarrow E * E)$$

$$\Rightarrow E + E * \underline{id} \quad (id, \quad E \rightarrow id)$$

$$\Rightarrow E + \underline{id * id} \quad (id, \quad E \rightarrow id)$$

$$\Rightarrow \underline{id + id * id} \quad (id, \quad E \rightarrow id)$$

$$E \Rightarrow \underline{E * E} \quad (E * E, \quad E \rightarrow E * E)$$

$$\Rightarrow E * \underline{id} \quad (id, \quad E \rightarrow id)$$

$$\Rightarrow \underline{E + E * id} \quad (E + E, \quad E \rightarrow E + E)$$

$$\Rightarrow E + \underline{id * id} \quad (id, \quad E \rightarrow id)$$

$$\Rightarrow \underline{id + id * id} \quad (id, \quad E \rightarrow id)$$

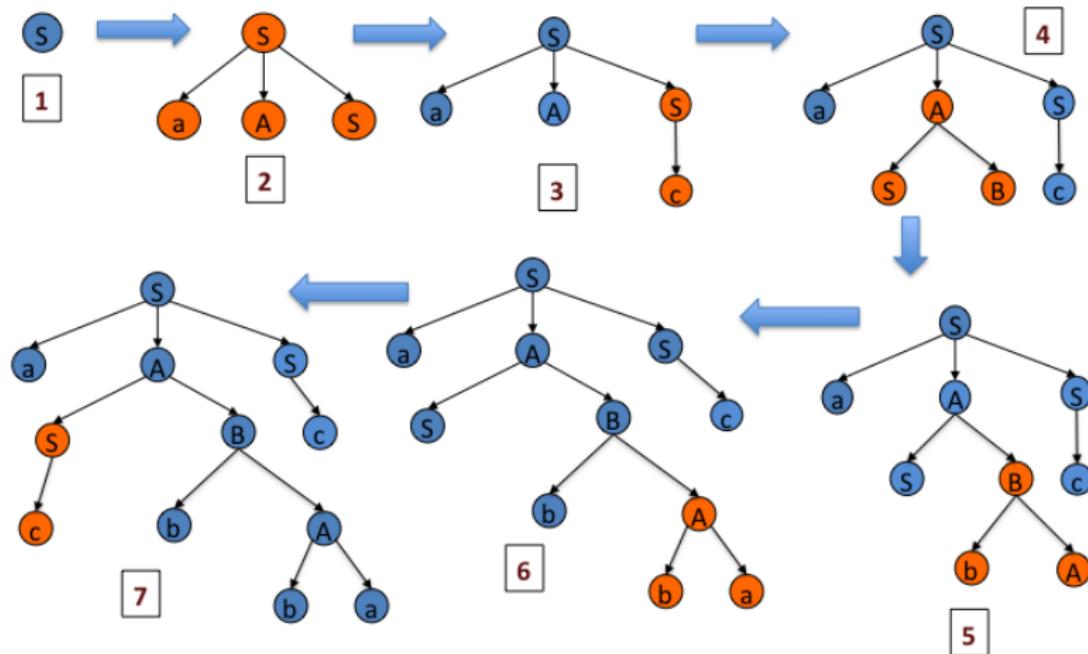
# Rightmost Derivation and Bottom-UP Parsing

$$\begin{array}{l} S \rightarrow aAS \mid c \\ A \rightarrow ba \mid SB \\ B \rightarrow bA \mid S \end{array}$$

Rightmost derivation of the string *acbbac*

$$S \Rightarrow \underline{aAS} \Rightarrow a\underline{Ac} \Rightarrow a\underline{S}bc \Rightarrow aS\underline{bA}c \Rightarrow aS\underline{bb}ac \Rightarrow a\underline{c}bbac$$

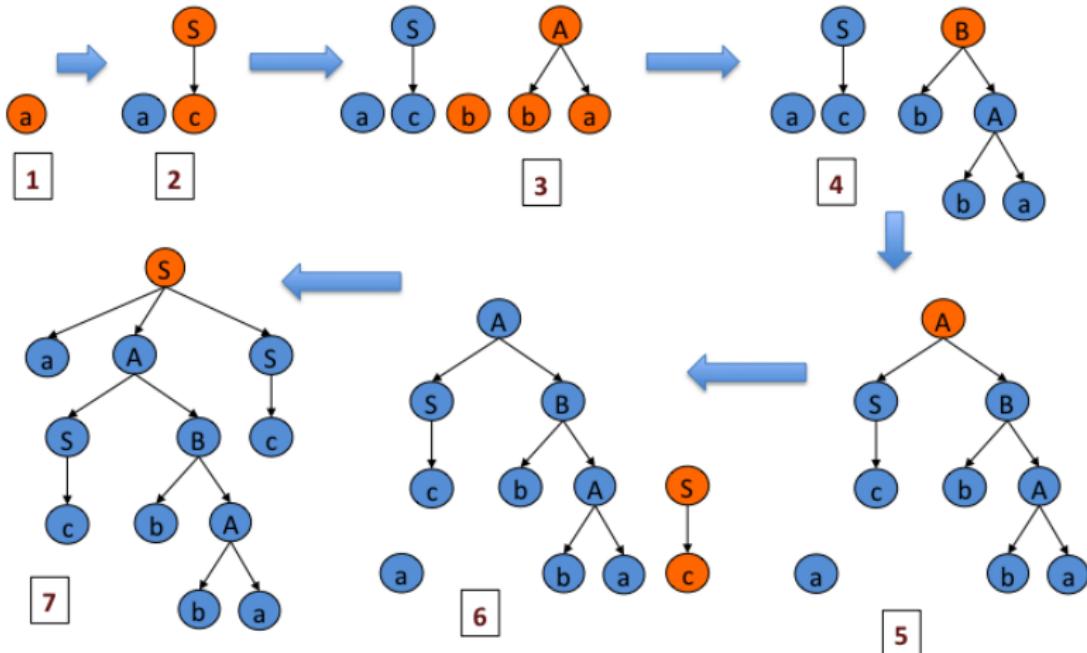
1      2      3      4      5      6      7



# Rightmost Derivation and Bottom-UP Parsing (contd.)

$$\begin{array}{l} S \rightarrow aAS \mid c \\ A \rightarrow ba \mid SB \\ B \rightarrow bA \mid S \end{array}$$

Rightmost derivation of the string *acbbac* in reverse  
 $S \leq_a \underline{aAS} \leq_a \underline{aAc} \leq_a \underline{aSBC} \leq_a \underline{aSbAC} \leq_a \underline{aSbbAC} \leq_a \underline{aCbbac}$   
7      6      5      4      3      2      1



# Shift-Reduce Parsing Algorithm

- How do we locate a handle in a right sentential form?
  - An LR parser uses a DFA to detect the condition that a handle is now on the stack
- Which production to use, in case there is more than one with the same RHS?
  - An LR parser uses a parsing table similar to an LL parsing table, to choose the production
- A stack is used to implement an S-R parser, The parser has four actions
  - ① **shift**: the next input symbol is shifted to the top of stack
  - ② **reduce**: the right end of the handle is the top of stack; locates the left end of the handle inside the stack and replaces the handle by the LHS of an appropriate production
  - ③ **accept**: announces successful completion of parsing
  - ④ **error**: syntax error, error recovery routine is called

# S-R Parsing Example 1

$\$$  marks the bottom of stack and the right end of the input

Stack	Input	Action
$\$$	$acbbac\$$	shift
$\$ a$	$cbbac\$$	shift
$\$ ac$	$bbac\$$	reduce by $S \rightarrow c$
$\$ aS$	$bbac\$$	shift
$\$ aSb$	$bac\$$	shift
$\$ aSbb$	$ac\$$	shift
$\$ aSba$	$c\$$	reduce by $A \rightarrow ba$
$\$ aSbA$	$c\$$	reduce by $B \rightarrow bA$
$\$ aSB$	$c\$$	reduce by $A \rightarrow SB$
$\$ aA$	$c\$$	shift
$\$ aAc$	$\$$	reduce by $S \rightarrow c$
$\$ aAS$	$\$$	reduce by $S \rightarrow aAS$
$\$ S$	$\$$	accept

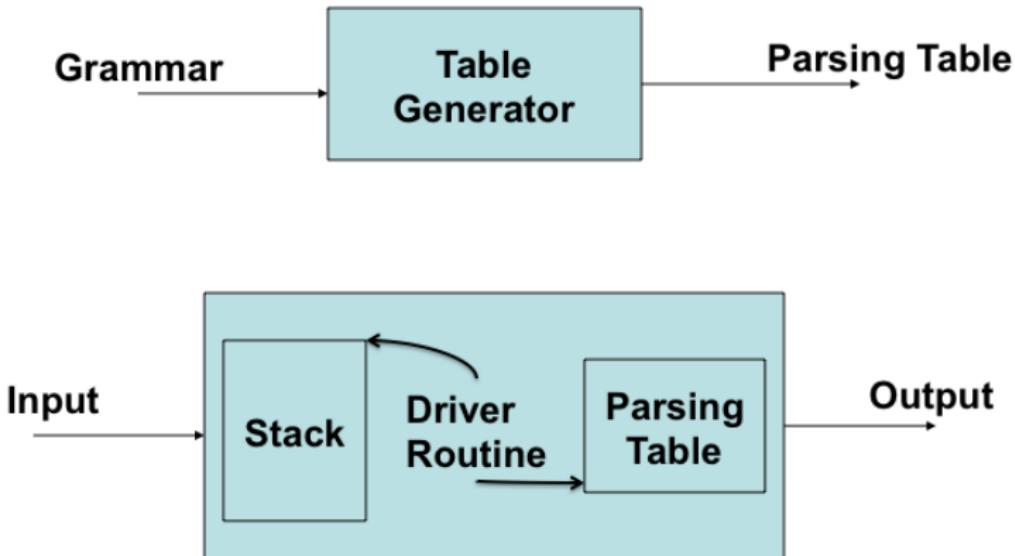
## S-R Parsing Example 2

$\$$  marks the bottom of stack and the right end of the input

Stack	Input	Action
$\$$	$id_1 + id_2 * id_3 \$$	shift
$\$ id_1$	$+ id_2 * id_3 \$$	reduce by $E \rightarrow id$
$\$ E$	$+ id_2 * id_3 \$$	shift
$\$ E +$	$id_2 * id_3 \$$	shift
$\$ E + id_2$	$* id_3 \$$	reduce by $E \rightarrow id$
$\$ E + E$	$* id_3 \$$	shift
$\$ E + E *$	$id_3 \$$	shift
$\$ E + E * id_3$	$\$$	reduce by $E \rightarrow id$
$\$ E + E * E$	$\$$	reduce by $E \rightarrow E * E$
$\$ E + E$	$\$$	reduce by $E \rightarrow E + E$
$\$ E$	$\$$	accept

- LR( $k$ ) - Left to right scanning with Rightmost derivation in reverse,  $k$  being the number of lookahead tokens
  - $k = 0, 1$  are of practical interest
- LR parsers are also automatically generated using parser generators
- LR grammars are a subset of CFGs for which LR parsers can be constructed
- LR(1) grammars can be written quite easily for practically all programming language constructs for which CFGs can be written
- LR parsing is the most general non-backtracking shift-reduce parsing method (known today)
- LL grammars are a strict subset of LR grammars - an LL( $k$ ) grammar is also LR( $k$ ), but not vice-versa

# LR Parser Generation

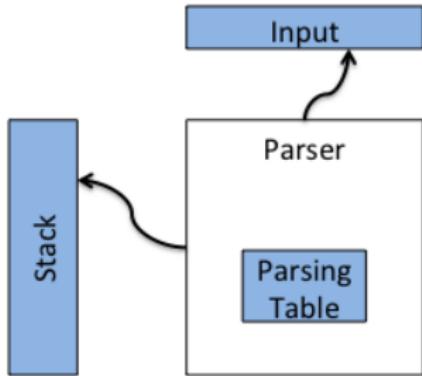


**LR Parser Generator**

# LR Parser Configuration

- A configuration of an LR parser is:  
 $(s_0 X_1 s_2 X_2 \dots X_m s_m, \quad a_i a_{i+1} \dots a_n \$)$ , where,  
**stack            unexpended input**  
 $s_0, s_1, \dots, s_m$ , are the states of the parser, and  $X_1, X_2, \dots, X_m$ , are grammar symbols (terminals or nonterminals)
- Starting configuration of the parser:  $(s_0, a_1 a_2 \dots a_n \$)$ , where,  $s_0$  is the initial state of the parser, and  $a_1 a_2 \dots a_n$  is the string to be parsed
- Two parts in the parsing table: *ACTION* and *GOTO*
  - The *ACTION* table can have four types of entries: **shift**, **reduce**, **accept**, or **error**
  - The *GOTO* table provides the next state information to be used after a *reduce* move

# LR Parsing Algorithm



```
Initial configuration: Stack = state 0, Input = w$,  
a = first input symbol;  
repeat {  
    let s be the top stack state;  
    let a be the next input symbol;  
    if ( ACTION[s, a] == shift p ) {  
        push a and p onto the stack (in that order);  
        advance input pointer;  
    } else if ( ACTION[s, a] == reduce A → α ) then {  
        pop 2*/α/* symbols off the stack;  
        let s' be the top of stack state now;  
        push A and GOTO[s', A] onto the stack  
        (in that order);  
    } else if ( ACTION[s, a] == accept ) break;  
    /* parsng is over */  
    else error();  
} until true; /* for ever */
```

# LR Parsing Example 1 - Parsing Table

STATE	ACTION				GOTO		
	a	b	c	\$	S	A	B
0	S2		S3		1		
1				R1 acc			
2	S2	S6	S3		8	4	
3	R3	R3	R3	R3			
4	S2		S3		5		
5	R2	R2	R2	R2			
6	S7						
7	R4	R4	R4	R4			
8	S2	S10	S3		12		9
9	R5	R5	R5	R5			
10	S2	S6	S3		8	11	
11	R6	R6	R6	R6			
12	R7	R7	R7	R7			

1.  $S' \rightarrow S$
2.  $S \rightarrow aAS$
3.  $S \rightarrow c$
4.  $A \rightarrow ba$
5.  $A \rightarrow SB$
6.  $B \rightarrow bA$
7.  $B \rightarrow S$

# Syntax Analysis:

Context-free Grammars, Pushdown Automata and Parsing  
Part - 5

Y.N. Srikant

Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

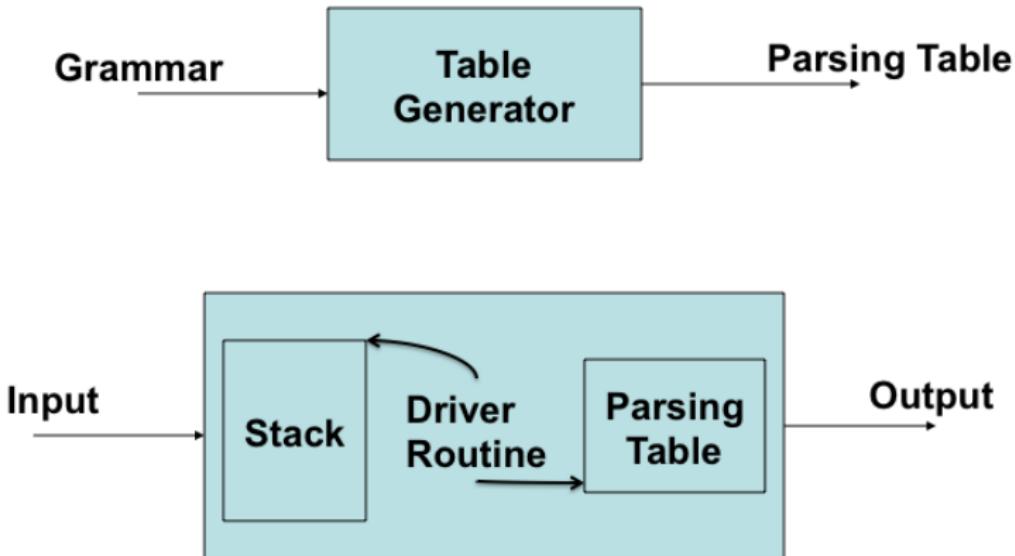
NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis? (covered in lecture 1)
- Specification of programming languages: context-free grammars (covered in lecture 1)
- Parsing context-free languages: push-down automata (covered in lectures 1 and 2)
- Top-down parsing: LL(1) parsing  
(covered in lectures 2 and 3)
- Recursive-descent parsing (covered in lecture 4)
- Bottom-up parsing: LR-parsing

- LR( $k$ ) - Left to right scanning with Rightmost derivation in reverse,  $k$  being the number of lookahead tokens
  - $k = 0, 1$  are of practical interest
- LR parsers are also automatically generated using parser generators
- LR grammars are a subset of CFGs for which LR parsers can be constructed
- LR(1) grammars can be written quite easily for practically all programming language constructs for which CFGs can be written
- LR parsing is the most general non-backtracking shift-reduce parsing method (known today)
- LL grammars are a strict subset of LR grammars - an LL( $k$ ) grammar is also LR( $k$ ), but not vice-versa

# LR Parser Generation

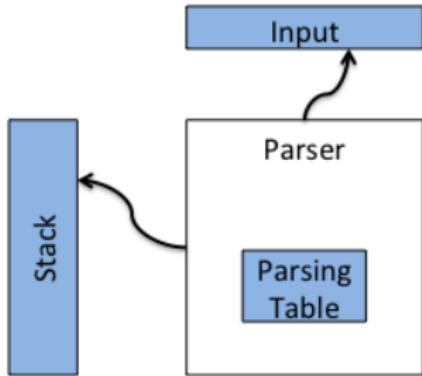


**LR Parser Generator**

# LR Parser Configuration

- A configuration of an LR parser is:  
 $(s_0 X_1 s_2 X_2 \dots X_m s_m, \quad a_i a_{i+1} \dots a_n \$)$ , where,  
**stack            unexpended input**  
 $s_0, s_1, \dots, s_m$ , are the states of the parser, and  $X_1, X_2, \dots, X_m$ , are grammar symbols (terminals or nonterminals)
- Starting configuration of the parser:  $(s_0, a_1 a_2 \dots a_n \$)$ , where,  $s_0$  is the initial state of the parser, and  $a_1 a_2 \dots a_n$  is the string to be parsed
- Two parts in the parsing table: *ACTION* and *GOTO*
  - The *ACTION* table can have four types of entries: **shift**, **reduce**, **accept**, or **error**
  - The *GOTO* table provides the next state information to be used after a *reduce* move

# LR Parsing Algorithm



```
Initial configuration: Stack = state 0, Input = w$,  
a = first input symbol;  
repeat {  
    let s be the top stack state;  
    let a be the next input symbol;  
    if ( ACTION[s, a] == shift p) {  
        push a and p onto the stack (in that order);  
        advance input pointer;  
    } else if (ACTION[s, a] == reduce A → α) then {  
        pop 2*/α/ symbols off the stack;  
        let s' be the top of stack state now;  
        push A and GOTO[s', A] onto the stack  
        (in that order);  
    } else if (ACTION[s, a] == accept) break;  
    /* parsng is over */  
    else error();  
} until true; /* for ever */
```

# LR Parsing Example 1 - Parsing Table

STATE	ACTION				GOTO		
	a	b	c	\$	S	A	B
0	S2		S3		1		
1				R1 acc			
2	S2	S6	S3		8	4	
3	R3	R3	R3	R3			
4	S2		S3		5		
5	R2	R2	R2	R2			
6	S7						
7	R4	R4	R4	R4			
8	S2	S10	S3		12		9
9	R5	R5	R5	R5			
10	S2	S6	S3		8	11	
11	R6	R6	R6	R6			
12	R7	R7	R7	R7			

1.  $S' \rightarrow S$
2.  $S \rightarrow aAS$
3.  $S \rightarrow c$
4.  $A \rightarrow ba$
5.  $A \rightarrow SB$
6.  $B \rightarrow bA$
7.  $B \rightarrow S$

## LR Parsing Example 1 (contd.)

Stack	Input	Action
0	acbbac\$	S2
0a2	cbbac\$	S3
0a2c3	bbac\$	R3 ( $S \rightarrow c$ , goto(2,S) = 8)
0a2S8	bbac\$	S10
0a2S8b10	bac\$	S6
0a2S8b10b6	ac\$	S7
0a2S8b10b6a7	c\$	R4 ( $A \rightarrow ba$ , goto(10,A) = 11)
0a2S8b10A11	c\$	R6 ( $B \rightarrow bA$ , goto(8,B) = 9)
0a2S8B9	c\$	R5 ( $A \rightarrow SB$ , goto(2,A) = 4)
0a2A4	c\$	S3
0a2A4c3	\$	R3 ( $S \rightarrow c$ , goto(4,S) = 5)
0a2A4S5	\$	R2 ( $S \rightarrow aAS$ , goto(0,S) = 1)
0S1	\$	R1 ( $S' \rightarrow S$ ), and accept

# LR Parsing Example 2 - Parsing Table

STATE	ACTION						GOTO		
	id	+	*	(	)	\$	E	T	F
0	S5			S4			1	2	3
1		S6				R7 acc			
2		R2	S7		R2	R2			
3		R4	R4		R4	R4			
4	S5			S4			8	2	3
5		R6	R6		R6	R6			
6	S5			S4			9	3	
7	S5			S4					10
8		S6			S11				
9		R1	S7		R1	R1			
10		R3	R3		R3	R3			
11		R5	R5		R5	R5			

1.  $E \rightarrow E + T$
2.  $E \rightarrow T$
3.  $T \rightarrow T^* F$
4.  $T \rightarrow F$
5.  $F \rightarrow (E)$
6.  $F \rightarrow id$
7.  $S \rightarrow E$

## LR Parsing Example 2(contd.)

Stack	Input	Action
0	$id + id * id \$$	S5
0 $id$ 5	$+ id * id \$$	R6 ( $F \rightarrow id$ , $G(0,F) = 3$ )
0 $F$ 3	$+ id * id \$$	R4 ( $T \rightarrow F$ , $G(0,T) = 2$ )
0 $T$ 2	$+ id * id \$$	R2 ( $E \rightarrow T$ , $G(0,E) = 1$ )
0 $E$ 1	$+ id * id \$$	S6
0 $E$ 1 $+$ 6	$id * id \$$	S5
0 $E$ 1 $+$ 6 $id$ 5	$* id \$$	R6 ( $F \rightarrow id$ , $G(6,F) = 3$ )
0 $E$ 1 $+$ 6 $F$ 3	$* id \$$	R4 ( $T \rightarrow F$ , $G(6,T) = 9$ )
0 $E$ 1 $+$ 6 $T$ 9	$* id \$$	S7
0 $E$ 1 $+$ 6 $T$ 9 $*$ 7	$id \$$	S5
0 $E$ 1 $+$ 6 $T$ 9 $*$ 7 $id$ 5	$\$$	R6 ( $F \rightarrow id$ , $G(7,F) = 10$ )
0 $E$ 1 $+$ 6 $T$ 9 $*$ 7 $F$ 10	$\$$	R3 ( $T \rightarrow T * F$ , $G(6,T) = 9$ )
0 $E$ 1 $+$ 6 $T$ 9	$\$$	R1 ( $E \rightarrow E + T$ , $G(0,E) = 1$ )
0 $E$ 1	$\$$	R7 ( $S \rightarrow E$ ) and accept

- Consider a rightmost derivation:  
 $S \Rightarrow_{rm}^* \phi Bt \Rightarrow_{rm} \phi \beta t$ ,  
where the production  $B \rightarrow \beta$  has been applied
- A grammar is said to be **LR(k)**, if for any given input string, at each step of any rightmost derivation, the handle  $\beta$  can be detected by examining the string  $\phi\beta$  and scanning *at most*, first  $k$  symbols of the unused input string  $t$

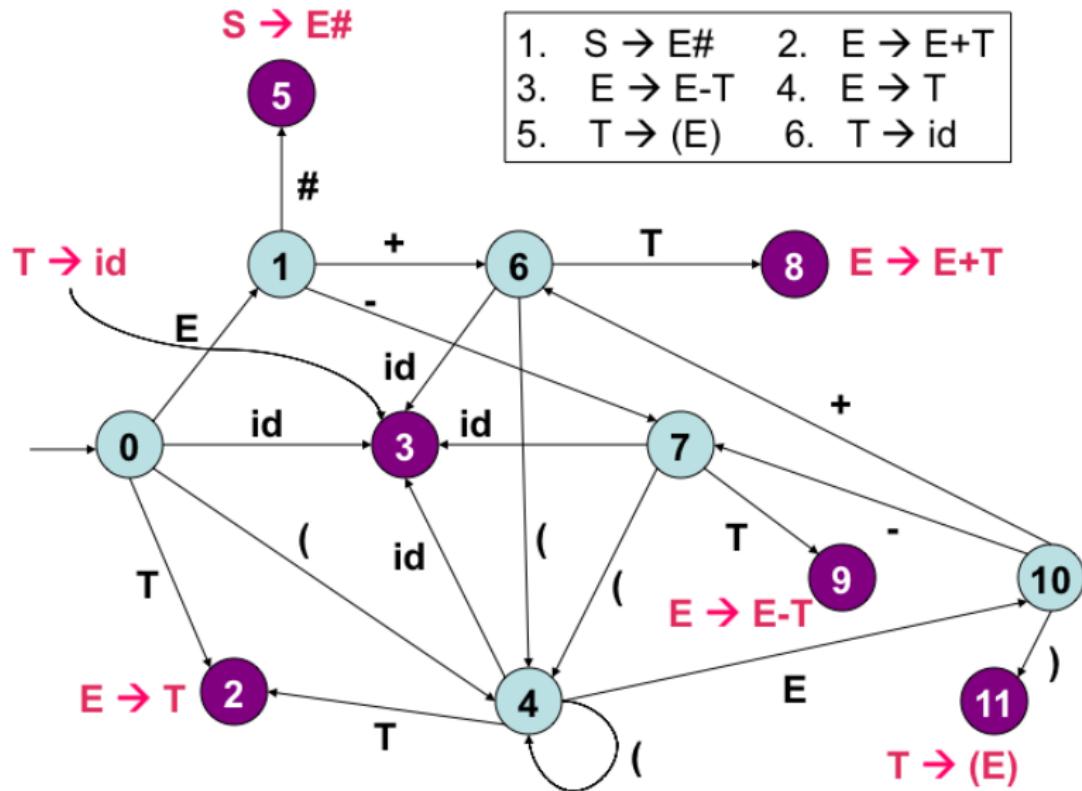
# LR Grammars (contd.)

- Example: The grammar,  
 $\{S \rightarrow E, E \rightarrow E + E \mid E * E \mid id\}$ , is not LR(2)
  - $S \Rightarrow^1 \underline{E} \Rightarrow^2 \underline{E + E} \Rightarrow^3 E + \underline{E * E} \Rightarrow^4 E + E * \underline{id} \Rightarrow^5 E + \underline{id * id} \Rightarrow^6 \underline{id} + id * id$
  - $S \Rightarrow^{1'} \underline{E} \Rightarrow^{2'} \underline{E * E} \Rightarrow^{3'} E * \underline{id} \Rightarrow^{4'} \underline{E + E * id} \Rightarrow^{5'} E + \underline{id * id} \Rightarrow^{6'} \underline{id} + id * id$
- In the above two derivations, the handle at steps 6 & 6' and at steps 5 & 5', is  $E \rightarrow id$ , and the position is underlined (with the same lookahead of two symbols,  $id+$  and  $+id$ )
- However, the handles at step 4 and at step 4' are different ( $E \rightarrow id$  and  $E \rightarrow E + E$ ), even though the lookahead of 2 symbols is the same ( $*id$ ), and the stack is also the same ( $\phi = E + E$ )
- That means that the handle cannot be determined using the lookahead

- A **viable prefix** of a sentential form  $\phi\beta t$ , where  $\beta$  denotes the handle, is any prefix of  $\phi\beta$ . A viable prefix cannot contain symbols to the right of the handle
- Example:  $S \rightarrow E\#$ ,  $E \rightarrow E + T \mid E - T \mid T$ ,  $T \rightarrow id \mid (E)$   
 $S \Rightarrow E\# \Rightarrow E + T\# \Rightarrow E + (E)\# \Rightarrow E + (T)\#$   
 $\Rightarrow E + (id)\#$   
 $E$ ,  $E+$ ,  $E + ($ , and  $E + (id)$ , are all viable prefixes of the right sentential form  $E + (id)\#$
- It is always possible to add appropriate terminal symbols to the end of a viable prefix to get a right-sentential form
- Viable prefixes characterize the prefixes of sentential forms that can occur on the stack of an LR parser

- **Theorem:** The set of all viable prefixes of all the right sentential forms of a grammar is a regular language
- The DFA of this regular language can detect handles during LR parsing
- When this DFA reaches a “reduction state”, the corresponding viable prefix cannot grow further and thus signals a reduction
- This DFA can be constructed by the compiler using the grammar
- All LR parsers have such a DFA incorporated in them
- We construct an augmented grammar for which we construct the DFA
  - If  $S$  is the start symbol of  $G$ , then  $G'$  contains all productions of  $G$  and also a new production  $S' \rightarrow S$
  - This enables the parser to halt as soon as  $S'$  appears on the stack

# DFA for Viable Prefixes - LR(0) Automaton



# Items and *Valid* Items

- A finite set of *items* is associated with each state of DFA
  - An *item* is a marked production of the form  $[A \rightarrow \alpha_1.\alpha_2]$ , where  $A \rightarrow \alpha_1\alpha_2$  is a production and '.' denotes the mark
  - Many items may be associated with a production e.g., the items  $[E \rightarrow .E + T]$ ,  $[E \rightarrow E.+T]$ ,  $[E \rightarrow E+ .T]$ , and  $[E \rightarrow E + T.]$  are associated with the production  $E \rightarrow E + T$
- An item  $[A \rightarrow \alpha_1.\alpha_2]$  is *valid* for some viable prefix  $\phi\alpha_1$ , iff, there exists some rightmost derivation  $S \Rightarrow^* \phi At \Rightarrow \phi\alpha_1\alpha_2 t$ , where  $t \in \Sigma^*$
- There may be several items valid for a viable prefix
  - The items  $[E \rightarrow E - .T]$ ,  $[T \rightarrow .id]$ , and  $[T \rightarrow .(E)]$  are all valid for the viable prefix “ $E -$ ” as shown below  
 $S \Rightarrow E\# \Rightarrow E - T\#$ ,  $S \Rightarrow E\# \Rightarrow E - T\# \Rightarrow E - id\#$ ,  
 $S \Rightarrow E\# \Rightarrow E - T\# \Rightarrow E - (E)\#$

# Valid Items and States of LR(0) DFA

- An item indicates how much of a production has already been seen and how much remains to be seen
  - $[E \rightarrow E - . T]$  indicates that we have already seen a string derivable from " $E -$ " and that we hope to see next, a string derivable from  $T$
- Each state of an LR(0) DFA contains only those items that are valid for the *same set of viable prefixes*
  - All items in state 7 are valid for the viable prefixes " $E -$ " and " $(E -$ " (and many more)
  - All items in state 4 are valid for the viable prefix "(" (and many more)
  - In fact, the set of all viable prefixes for which the items in a state  $s$  are valid is the set of strings that can take us from state 0 (initial) to state  $s$
- Constructing the LR(0) DFA using sets of items is very simple

# Closure of a Set of Items

```
Itemset closure(I){ /* I is a set of items */
    while (more items can be added to I) {
        for each item  $[A \rightarrow \alpha.B\beta] \in I$  {
            /* note that B is a nonterminal and is right after the “.” */
            for each production  $B \rightarrow \gamma \in G$ 
                if (item  $[B \rightarrow .\gamma] \notin I$ ) add item  $[B \rightarrow .\gamma]$  to I
        }
        return I
    }
}
```

<u>State 0</u>	<u>State 1</u>	<u>State 7</u>	<u>State 2</u>
$S \rightarrow .E\#$	$S \rightarrow E.\#$	$E \rightarrow E-.T$	$E \rightarrow T.$
$E \rightarrow .E+T$	$E \rightarrow E.+T$	$T \rightarrow .(E)$	
$E \rightarrow .E-T$	$E \rightarrow E.-T$	$T \rightarrow .id$	
$E \rightarrow .T$			
$T \rightarrow .(E)$			
$T \rightarrow .id$			



indicates closure items

# GOTO set computation

Itemset  $GOTO(I, X)\{ /* I is a set of items$   
 $X$  is a grammar symbol, a terminal or a nonterminal \*/  
Let  $I' = \{[A \rightarrow \alpha X . \beta] \mid [A \rightarrow \alpha . X \beta] \in I\};$   
return ( $closure(I')$ )  
}

<u>State 0</u>	<u>State 1</u>	<u>State 7</u>
$S \rightarrow .E\#$	$S \rightarrow E.\#$	$E \rightarrow E-.T$
$E \rightarrow .E+T$	$E \rightarrow E.+T$	$T \rightarrow .(E)$
$E \rightarrow .E-T$	$E \rightarrow E.-T$	$T \rightarrow .id$
$E \rightarrow .T$		
$T \rightarrow .(E)$	● indicates closure items	
$T \rightarrow .id$	$GOTO(0, E) = 1$ $GOTO(1, -) = 7$	

# Intuition behind *closure* and *GOTO*

- If an item  $[A \rightarrow \alpha.B\delta]$  is in a state (i.e., item set  $I$ ), then, some time in the future, we expect to see in the input, a string derivable from  $B\delta$ 
  - This implies a string derivable from  $B$  as well
  - Therefore, we add an item  $[B \rightarrow .\beta]$  corresponding to each production  $B \rightarrow \beta$  of  $B$ , to the state (i.e., item set  $I$ )
- If  $I$  is the set of items valid for a viable prefix  $\gamma$ 
  - All the items in  $\text{closure}(I)$  are also valid for  $\gamma$
  - $\text{GOTO}(I, X)$  is the set items valid for the viable prefix  $\gamma X$ 
    - If  $[A \rightarrow \alpha.B\delta]$  (in item set  $I$ ) is valid for the viable prefix  $\phi\alpha$ , and  $B \rightarrow \beta$  is a production, we have
$$S \Rightarrow^* \phi At \Rightarrow \phi\alpha B\delta t \Rightarrow^* \phi\alpha Bxt \Rightarrow \phi\alpha\beta xt$$
demonstrating that the item  $[B \rightarrow .\beta]$  (in the closure of  $I$ ) is valid for  $\phi\alpha$
    - The above derivation also shows that the item  $[A \rightarrow \alpha B.\delta]$  (in  $\text{GOTO}(I, B)$ ) is valid for the viable prefix  $\phi\alpha B$

# Construction of Sets of Canonical LR(0) Items

```
void Set_of_item_sets(G'){ /* G' is the augmented grammar */  
    C = {closure({S' → .S})};/* C is a set of item sets */  
    while (more item sets can be added to C) {  
        for each item set I ∈ C and each grammar symbol X  
            /* X is a grammar symbol, a terminal or a nonterminal */  
            if ((GOTO(I, X) ≠ ∅) && (GOTO(I, X) ∉ C))  
                C = C ∪ GOTO(I, X)  
    }  
}
```

- Each set in  $C$  (above) corresponds to a state of a DFA (LR(0) DFA)
- This is the DFA that recognizes viable prefixes

# Construction of an LR(0) Automaton - Example 1

## State 0

$S \rightarrow .E\#$

$E \rightarrow .E+T$

$E \rightarrow .E-T$

$E \rightarrow .T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 1

$S \rightarrow E.\#$

$E \rightarrow E.+T$

$E \rightarrow E.-T$

## State 2

$E \rightarrow T.$

## State 3

$T \rightarrow id.$

## State 6

$E \rightarrow E+.T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 4

$T \rightarrow (.E)$

$E \rightarrow .E+T$

$E \rightarrow .E-T$

$E \rightarrow .T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 7

$E \rightarrow E.-T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 9

$E \rightarrow E-T.$

## State 10

$T \rightarrow (E.)$

$E \rightarrow E.+T$

$E \rightarrow E.-T$

## State 8

$E \rightarrow E+T.$

## State 11

$T \rightarrow (E).$

## State 5

$S \rightarrow E\#.$



indicates closure items



indicates kernel items

# Shift and Reduce Actions

- If a state contains an item of the form  $[A \rightarrow \alpha.]$  ("reduce item"), then a reduction by the production  $A \rightarrow \alpha$  is the action in that state
- If there are no "reduce items" in a state, then shift is the appropriate action
- There could be shift-reduce conflicts or reduce-reduce conflicts in a state
  - Both shift and reduce items are present in the same state (S-R conflict), or
  - More than one reduce item is present in a state (R-R conflict)
  - It is normal to have more than one shift item in a state (no shift-shift conflicts are possible)
- If there are no S-R or R-R conflicts in any state of an LR(0) DFA, then the grammar is LR(0), otherwise, it is not LR(0)

# Syntax Analysis:

Context-free Grammars, Pushdown Automata and Parsing  
Part - 6

Y.N. Srikant

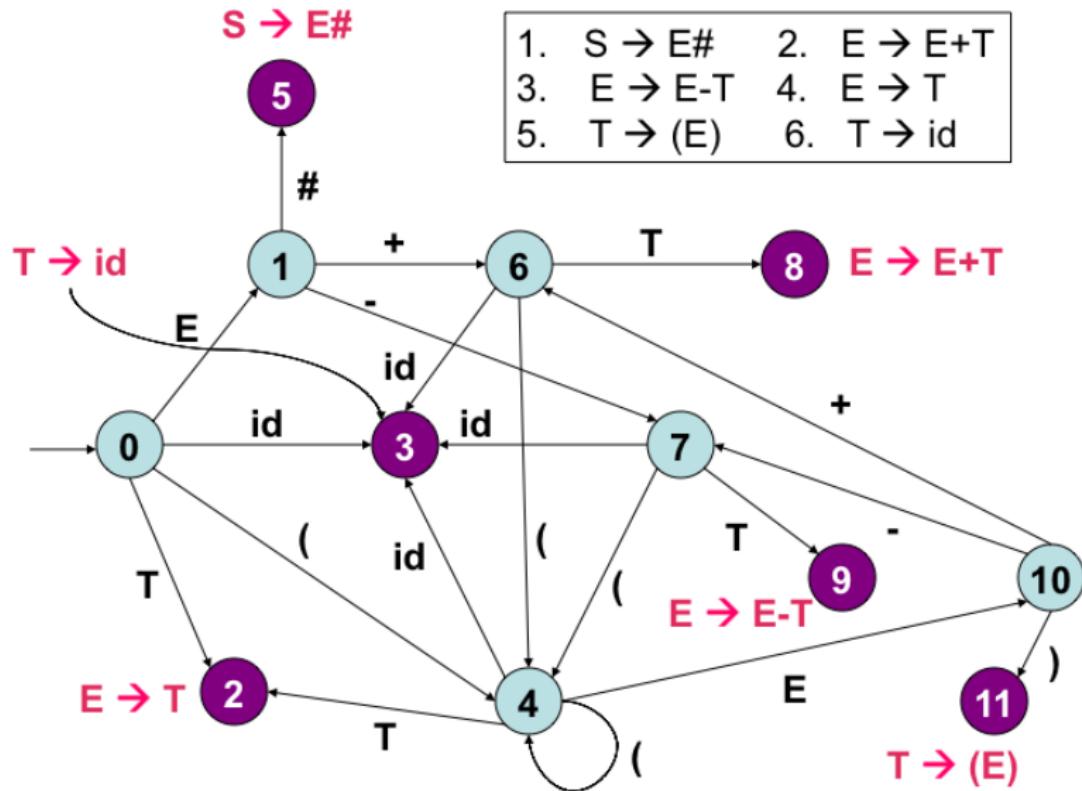
Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- What is syntax analysis? (covered in lecture 1)
- Specification of programming languages: context-free grammars (covered in lecture 1)
- Parsing context-free languages: push-down automata (covered in lectures 1 and 2)
- Top-down parsing: LL(1) parsing (covered in lectures 2 and 3)
- Recursive-descent parsing (covered in lecture 4)
- Bottom-up parsing: LR-parsing (continued)

# DFA for Viable Prefixes - LR(0) Automaton



# Construction of Sets of Canonical LR(0) Items

```
void Set_of_item_sets(G'){ /* G' is the augmented grammar */  
    C = {closure({S' → .S})};/* C is a set of item sets */  
    while (more item sets can be added to C) {  
        for each item set I ∈ C and each grammar symbol X  
            /* X is a grammar symbol, a terminal or a nonterminal */  
            if ((GOTO(I, X) ≠ ∅) && (GOTO(I, X) ∉ C))  
                C = C ∪ GOTO(I, X)  
    }  
}
```

- Each set in  $C$  (above) corresponds to a state of a DFA (LR(0) DFA)
- This is the DFA that recognizes viable prefixes

# Construction of an LR(0) Automaton - Example 1

## State 0

$S \rightarrow .E\#$

$E \rightarrow .E+T$

$E \rightarrow .E-T$

$E \rightarrow .T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 1

$S \rightarrow E.\#$

$E \rightarrow E.+T$

$E \rightarrow E.-T$

## State 2

$E \rightarrow T.$

## State 3

$T \rightarrow id.$

## State 4

$T \rightarrow (.E)$

$E \rightarrow .E+T$

$E \rightarrow .E-T$

$E \rightarrow .T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 5

$S \rightarrow E\#.$

## State 6

$E \rightarrow E+.T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 7

$E \rightarrow E-.T$

$T \rightarrow .(E)$

$T \rightarrow .id$

## State 8

$E \rightarrow E+T.$

## State 9

$E \rightarrow E-T.$

## State 10

$T \rightarrow (E.)$

$E \rightarrow E.+T$

$E \rightarrow E.-T$

## State 11

$T \rightarrow (E).$

● indicates closure items

● indicates kernel items

# Shift and Reduce Actions

- If a state contains an item of the form  $[A \rightarrow \alpha.]$  ("reduce item"), then a reduction by the production  $A \rightarrow \alpha$  is the action in that state
- If there are no "reduce items" in a state, then shift is the appropriate action
- There could be shift-reduce conflicts or reduce-reduce conflicts in a state
  - Both shift and reduce items are present in the same state (S-R conflict), or
  - More than one reduce item is present in a state (R-R conflict)
  - It is normal to have more than one shift item in a state (no shift-shift conflicts are possible)
- If there are no S-R or R-R conflicts in any state of an LR(0) DFA, then the grammar is LR(0), otherwise, it is not LR(0)

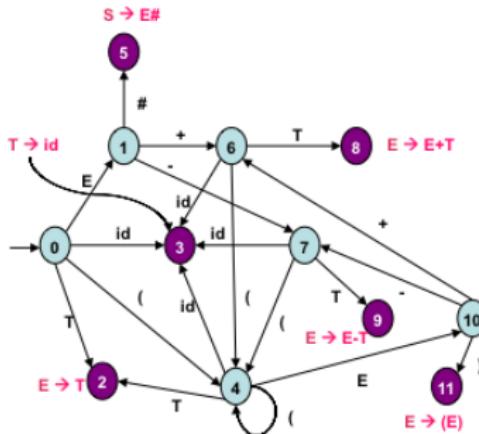
# LR(0) Parser Table - Example 1

STATE	ACTION						GOTO		
	+	-	(	)	id	#	S	E	T
0			S4		S3			1	2
1	S6	S7				S5			
2	R4	R4	R4	R4	R4	R4			
3	R6	R6	R6	R6	R6	R6			
4			S4		S3			10	2
5	R1 acc	R1 acc	R1 acc	R1 acc	R1 acc	R1 acc			
6			S4		S3				8
7			S4		S3				9
8	R2	R2	R2	R2	R2	R2			
9	R3	R3	R3	R3	R3	R3			
10	S6	S7		S11					
11	R5	R5	R5	R5	R5	R5			

1.  $S \rightarrow E\#$
2.  $E \rightarrow E+T$
3.  $E \rightarrow E-T$
4.  $E \rightarrow T$
5.  $T \rightarrow (E)$
6.  $T \rightarrow id$

# Construction of an LR(0) Parser Table - Example 1

STATE	ACTION					GOTO			
	+	-	(	)	id	#	S	E	T
0			S4		S3			1	2
1	S6	S7				S5			
2	R4	R4	R4	R4	R4	R4			
3	R6	R6	R6	R6	R6	R6			
4			S4		S3			10	2
5	R1 acc	R1 acc	R1 acc	R1 acc	R1 acc	R1 acc			
6			S4		S3				8
7			S4		S3				9
8	R2	R2	R2	R2	R2	R2			
9	R3	R3	R3	R3	R3	R3			
10	S6	S7		S1 1					
11	R5	R5	R5	R5	R5	R5			



1.  $S \rightarrow E\#$
2.  $E \rightarrow E+T$
3.  $E \rightarrow E-T$
4.  $E \rightarrow T$
5.  $T \rightarrow (E)$
6.  $T \rightarrow id$

State 0	State 2	State 4	State 7	State 10	State 11
$S \rightarrow .E\#$	$E \rightarrow T.$	$T \rightarrow (.E)$	$E \rightarrow E-.T$	$T \rightarrow (E.)$	$T \rightarrow (E).$
$E \rightarrow .E+T$		$E \rightarrow .E+T$		$E \rightarrow E.+T$	
$E \rightarrow .E-T$	$E \rightarrow .E-T$	$E \rightarrow .E-T$	$T \rightarrow .id$		$E \rightarrow E.-T$
$E \rightarrow .T$	$E \rightarrow T.$	$E \rightarrow T.$			
$T \rightarrow .(E)$		$T \rightarrow .(E)$			
$T \rightarrow .id$		$T \rightarrow .id$			

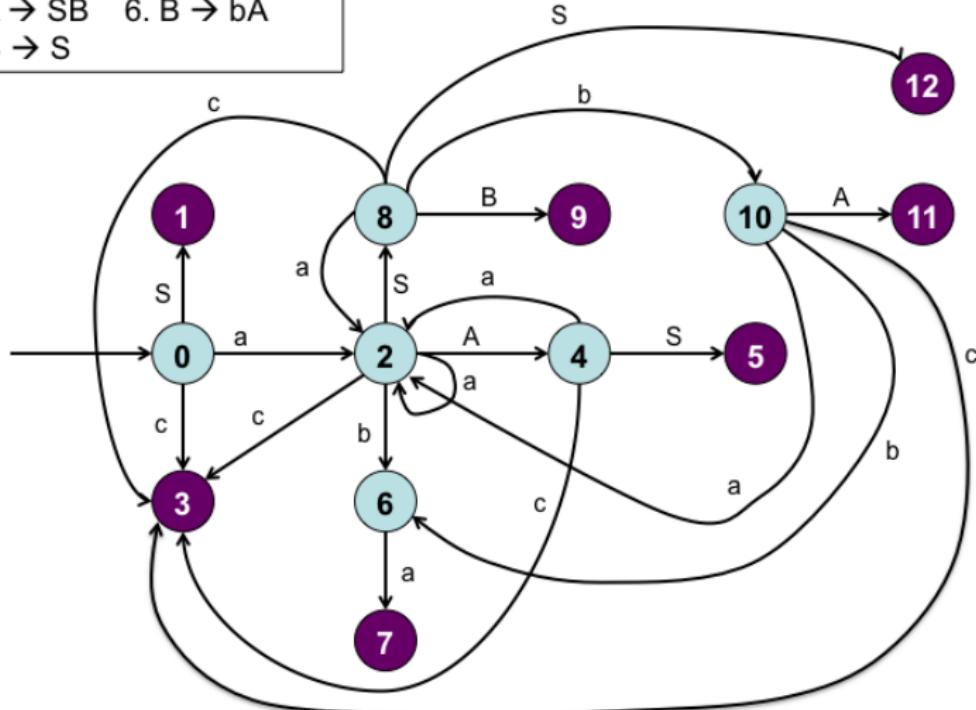
  

State 1	State 6	State 5	State 9
$S \rightarrow E. \#$	$E \rightarrow E+.T$	$S \rightarrow E. \#$	$E \rightarrow E-T.$
$E \rightarrow E.+T$	$T \rightarrow .(E)$		
$E \rightarrow E.-T$	$T \rightarrow .id$		

● indicates closure items  
 ● indicates kernel items

# LR(0) Automaton - Example 2

- 1.  $S' \rightarrow S$
- 2.  $S \rightarrow aAS$
- 3.  $S \rightarrow c$
- 4.  $A \rightarrow ba$
- 5.  $A \rightarrow SB$
- 6.  $B \rightarrow bA$
- 7.  $B \rightarrow S$



# Construction of an LR(0) Automaton - Example 2

<u>State 0</u> $S' \rightarrow .S$ $S \rightarrow .aAS$ $S \rightarrow .c$	<u>State 3</u> $S \rightarrow c.$	<u>State 7</u> $A \rightarrow ba.$	<u>State 10</u> $B \rightarrow b.A$ $A \rightarrow .ba$ $A \rightarrow .SB$ $S \rightarrow .aAS$ $S \rightarrow .c$
<u>State 1</u> $S' \rightarrow S.$	<u>State 4</u> $S \rightarrow aA.S$ $S \rightarrow .aAS$ $S \rightarrow .c$	<u>State 8</u> $A \rightarrow S.B$ $B \rightarrow .bA$ $B \rightarrow .S$	
<u>State 2</u> $S \rightarrow a.AS$ $A \rightarrow .ba$ $A \rightarrow .SB$ $S \rightarrow .aAS$ $S \rightarrow .c$	<u>State 5</u> $S \rightarrow aAS.$ $S \rightarrow .aAS$	$S \rightarrow .aAS$ $S \rightarrow .c$	<u>State 11</u> $B \rightarrow bA.$
	<u>State 6</u> $A \rightarrow b.a$	<u>State 9</u> $A \rightarrow SB.$	<u>State 12</u> $B \rightarrow S.$



indicates closure items



indicates kernel items

1.  $S' \rightarrow S$
2.  $S \rightarrow aAS$
3.  $S \rightarrow c$
4.  $A \rightarrow ba$
5.  $A \rightarrow SB$
6.  $B \rightarrow bA$
7.  $B \rightarrow S$

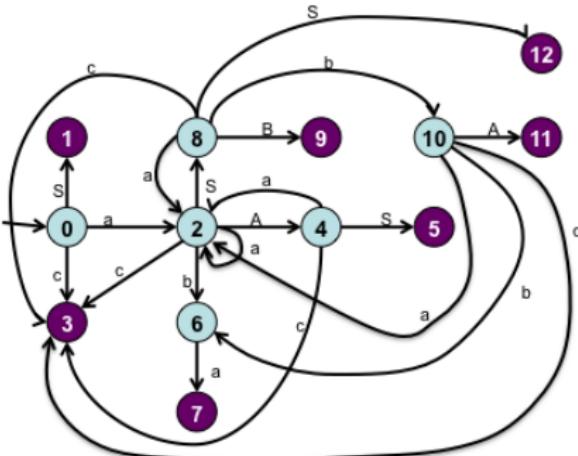
# LR(0) Parser Table - Example 2

STATE	ACTION				GOTO		
	a	b	c	\$	S	A	B
0	S2		S3		1		
1				R1 acc			
2	S2	S6	S3		8	4	
3	R3	R3	R3	R3			
4	S2		S3		5		
5	R2	R2	R2	R2			
6	S7						
7	R4	R4	R4	R4			
8	S2	S10	S3		12		9
9	R5	R5	R5	R5			
10	S2	S6	S3		8	11	
11	R6	R6	R6	R6			
12	R7	R7	R7	R7			

1.  $S' \rightarrow S$
2.  $S \rightarrow aAS$
3.  $S \rightarrow c$
4.  $A \rightarrow ba$
5.  $A \rightarrow SB$
6.  $B \rightarrow bA$
7.  $B \rightarrow S$

# Construction of an LR(0) Parser Table - Example 2

STATE	ACTION				GOTO		
	a	b	c	\$	S	A	B
0	S2		S3		1		
1				R1 acc			
2	S2	S6	S3		8	4	
3	R3	R3	R3	R3			
4	S2		S3		5		
5	R2	R2	R2	R2			
6	S7						
7	R4	R4	R4	R4			
8	S2	S1 0	S3		12		9
9	R5	R5	R5	R5			
10	S2	S6	S3		8	11	
11	R6	R6	R6	R6			
12	R7	R7	R7	R7			



1.  $S' \rightarrow S$   
 2.  $S \rightarrow aAS$   
 3.  $S \rightarrow c$   
 4.  $A \rightarrow ba$   
 5.  $A \rightarrow SB$   
 6.  $B \rightarrow bA$   
 7.  $B \rightarrow S$

<u>State 0</u>	<u>State 2</u>	<u>State 4</u>	<u>State 6</u>	<u>State 8</u>	<u>State 10</u>
$S' \rightarrow .S$	$S \rightarrow .aAS$	$S \rightarrow a.AS$	$A \rightarrow b.a$	$A \rightarrow S.B$	$B \rightarrow b.A$
$S \rightarrow .c$	$A \rightarrow .ba$	$S \rightarrow .aAS$	$B \rightarrow .bA$	$B \rightarrow .S$	$A \rightarrow .ba$
	$A \rightarrow .SB$	$S \rightarrow .c$	$S \rightarrow .aAS$	$A \rightarrow .SB$	$S \rightarrow .aAS$
	$S \rightarrow .aAS$	<u>State 5</u>	<u>State 7</u>	$B \rightarrow .c$	$S \rightarrow .c$
	$S \rightarrow .c$	$S \rightarrow aAS.$	$A \rightarrow ba.$	$S \rightarrow .c$	$S \rightarrow .c$
		<u>State 3</u>	<u>State 9</u>	<u>State 11</u>	<u>State 12</u>
	$S \rightarrow c.$	$S \rightarrow .c$	$A \rightarrow SB.$	$B \rightarrow bA.$	$B \rightarrow S.$

● indicates closure items

● indicates kernel items

# A Grammar that is not LR(0) - Example 1

State 0  
 $S \rightarrow .E$   
 $E \rightarrow .E+T$   
 $E \rightarrow .E-T$   
 $E \rightarrow .T$   
 $T \rightarrow .(E)$   
 $T \rightarrow .id$

State 2  
 $E \rightarrow T.$

State 5  
 $E \rightarrow E+.T$   
 $T \rightarrow .(E)$   
 $T \rightarrow .id$

State 8  
 $E \rightarrow E-T.$

State 3  
 $T \rightarrow id.$

State 6  
 $E \rightarrow E-.T$   
 $T \rightarrow .(E)$   
 $T \rightarrow .id$

State 9  
 $T \rightarrow (E.)$   
 $E \rightarrow E.+T$   
 $E \rightarrow E.-T$

State 1  
 $S \rightarrow E.$   
 $E \rightarrow E.+T$   
 $E \rightarrow E.-T$

State 4  
 $T \rightarrow (.E)$   
 $E \rightarrow .E+T$   
 $E \rightarrow .E-T$   
 $E \rightarrow .T$   
 $T \rightarrow .(E)$   
 $T \rightarrow .id$

State 7  
 $E \rightarrow E+T.$

State 10  
 $T \rightarrow (E).$

shift-reduce  
conflicts in  
state 1

● indicates closure items  
● indicates kernel items

$\text{follow}(S) = \{\$\}$ , where \$ is EOF  
Reduction on \$, and shifts on + and -, will resolve the conflicts  
This is similar to having an end marker such as #

Grammar is  
not LR(0), but  
is SLR(1)

# SLR(1) Parsers

- If the grammar is not LR(0), we try to resolve conflicts in the states using one look-ahead symbol
- Example: The expression grammar that is not LR(0)  
The state containing the items  $[T \rightarrow F.]$  and  $[T \rightarrow F.* T]$  has S-R conflicts
  - Consider the reduce item  $[T \rightarrow F.]$  and the symbols in  $FOLLOW(T)$
  - $FOLLOW(T) = \{+, ), \$\}$ , and reduction by  $T \rightarrow F$  can be performed on seeing one of these symbols in the input (look-ahead), since shift requires seeing  $*$  in the input
  - Recall from the definition of  $FOLLOW(T)$  that symbols in  $FOLLOW(T)$  are the only symbols that can legally follow  $T$  in any sentential form, and hence reduction by  $T \rightarrow F$  when one of these symbols is seen, is correct
  - If the S-R conflicts can be resolved using the FOLLOW set, the grammar is said to be SLR(1)

# A Grammar that is not LR(0) - Example 2

## State 0

$S \rightarrow .E$   
 $E \rightarrow .E+T$   
 $E \rightarrow .T$   
 $T \rightarrow .F^*T$   
 $T \rightarrow .F$   
 $F \rightarrow .(E)$   
 $F \rightarrow .id$

## State 2

$E \rightarrow T.$

## State 5

$F \rightarrow id.$   
 $E \rightarrow E+.T$   
 $T \rightarrow F$   
Shift-reduce  
conflict

## State 8

$F \rightarrow (E.)$   
 $E \rightarrow E.+T$   
 $E \rightarrow E+T.$

## State 1

$S \rightarrow E.$   
 $E \rightarrow E.+T$   
Shift-reduce  
conflict

## State 4

$F \rightarrow (.E)$   
 $E \rightarrow .E+T$   
 $E \rightarrow .T$   
 $T \rightarrow .F^*T$   
 $T \rightarrow .F$   
 $F \rightarrow .(E)$   
 $F \rightarrow .id$

## State 7

$T \rightarrow F^*.T$   
 $T \rightarrow .F^*T$   
 $T \rightarrow .F$   
 $F \rightarrow .(E)$   
 $F \rightarrow .id$

## State 10

$E \rightarrow F^*T.$

## State 11

$F \rightarrow (E).$

$\text{follow}(S) = \{\$\}$ , Reduction on \$ and shift on +, eliminates conflicts  
 $\text{follow}(T) = \{\$, , +\}$ , where \$ is EOF

Reduction on \$, , and +, and shift on \*, eliminates conflicts

Grammar is  
not LR(0), but  
is SLR(1)

# Construction of an SLR(1) Parsing Table

Let  $C = \{I_0, I_1, \dots, I_i, \dots, I_n\}$  be the canonical LR(0) collection of items, with the corresponding states of the parser being  $0, 1, \dots, i, \dots, n$ . Without loss of generality, let  $0$  be the initial state of the parser (containing the item  $[S' \rightarrow .S]$ )

Parsing actions for state  $i$  are determined as follows

1. If  $([A \rightarrow \alpha.a\beta] \in I_i) \text{ && } ([A \rightarrow \alpha a.\beta] \in I_j)$   
set ACTION[i, a] = *shift j* /\*  $a$  is a terminal symbol \*/
2. If  $([A \rightarrow \alpha.] \in I_i)$   
set ACTION[i, a] = *reduce A  $\rightarrow \alpha$ , for all  $a \in follow(A)$*
3. If  $([S' \rightarrow S.] \in I_i)$  set ACTION[i, \$] = *accept*

S-R or R-R conflicts in the table imply grammar is not SLR(1)

4. If  $([A \rightarrow \alpha.A\beta] \in I_i) \text{ && } ([A \rightarrow \alpha A.\beta] \in I_j)$   
set GOTO[i, A] =  $j$  /\*  $A$  is a nonterminal symbol \*/

All other entries not defined by the rules above are made *error*

# A Grammar that is not LR(0) - Example 3

Grammar

$$S' \rightarrow S, S \rightarrow aSb, S \rightarrow \epsilon$$

State 0

$$S' \rightarrow .S$$

$$S \rightarrow .aSb$$

$$S \rightarrow .$$

State 3

$$S \rightarrow aS.b$$

State 1

$$S' \rightarrow S.$$

State 4

$$S \rightarrow aSb.$$

State 2

$$S \rightarrow a.Sb$$

$$S \rightarrow .aSb$$

$$S \rightarrow .$$

shift-reduce  
conflict in  
states 0, 2



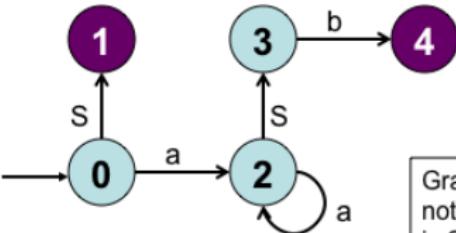
indicates closure items



indicates kernel items

$$\text{follow}(S) = \{\$, b\}$$

	a	b	\$	S
0	S2	reduce $S \rightarrow \epsilon$	reduce $S \rightarrow \epsilon$	1
1			accept	
2	S2	reduce $S \rightarrow \epsilon$	reduce $S \rightarrow \epsilon$	3
3		S4		
4		reduce $S \rightarrow aSb$	reduce $S \rightarrow aSb$	



Grammar is  
not LR(0), but  
is SLR(1)

# A Grammar that is not SLR(1) - Example 1

Grammar:  $S' \rightarrow S$ ,  
 $S \rightarrow aSb$ ,  $S \rightarrow ab$ ,  $S \rightarrow \epsilon$

$\text{follow}(S) = \{\$, b\}$   
**State 0:** Reduction on  $\$$  and  $b$ , by  $S \rightarrow \epsilon$ , and shift on  $a$  resolves conflicts  
**State 2:** S-R conflict on  $b$  still remains

**State 0**  
 $S' \rightarrow .S$   
 $S \rightarrow .aSb$   
 $S \rightarrow .ab$   
 $S \rightarrow .$

**State 3**  
 $S \rightarrow aS.b$

**State 4**  
 $S \rightarrow aS.b.$

**State 1**  
 $S' \rightarrow S.$

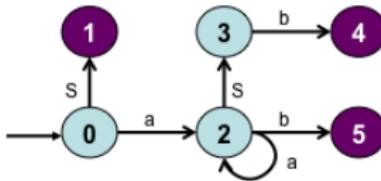
**State 5**  
 $S \rightarrow ab.$

**State 2**  
 $S \rightarrow a.Sb$   
 $S \rightarrow a.b$   
 $S \rightarrow .aSb$   
 $S \rightarrow .ab$   
 $S \rightarrow .$

shift-reduce  
conflict in  
states 0, 2

Grammar is  
neither LR(0)  
nor SLR(1)

	a	b	\$	S
0	S2	R: $S \rightarrow \epsilon$	R: $S \rightarrow \epsilon$	1
1			accept	
2	S2	<b>S5, R: <math>S \rightarrow \epsilon</math></b>	R: $S \rightarrow \epsilon$	3
3		S4		
4		R: $S \rightarrow aSb$	R: $S \rightarrow aSb$	
5		R: $S \rightarrow ab$	R: $S \rightarrow ab$	



# A Grammar that is not SLR(1) - Example 2

<u>Grammar</u>	<u>State 0</u>	<u>State 2</u>	<u>State 6</u>
$S' \rightarrow S$	$S' \rightarrow .S$	$S \rightarrow L = R$	$S \rightarrow L = .R$
$S \rightarrow L=R$	$S \rightarrow .L=R$	$R \rightarrow L.$	$R \rightarrow .L$
$S \rightarrow R$	$S \rightarrow .R$	shift-reduce	$L \rightarrow .*R$
$L \rightarrow *R$	$L \rightarrow .*R$	conflict	$L \rightarrow .id$
$L \rightarrow id$	$L \rightarrow .id$		
$R \rightarrow L$	$R \rightarrow .L$		
		<u>State 4</u>	<u>State 7</u>
		$L \rightarrow *.R$	$L \rightarrow *R.$
	<u>State 1</u>	$R \rightarrow .L$	
	$S' \rightarrow S.$	$L \rightarrow .*R$	
		$L \rightarrow .id$	
			<u>State 8</u>
			$R \rightarrow L.$
	<u>State 3</u>		
	$S \rightarrow R.$		
		<u>State 5</u>	
		$L \rightarrow id.$	
			<u>State 9</u>
			$S \rightarrow L=R.$

Grammar is  
neither LR(0)  
nor SLR(1)

Follow(R) = {\$,=} does not resolve S-R conflict

# The Problem with SLR(1) Parsers

- SLR(1) parser construction process does not remember enough left context to resolve conflicts
  - In the " $L = R$ " grammar (previous slide), the symbol '=' got into follow( $R$ ) because of the following derivation:  
 $S' \Rightarrow S \Rightarrow L = R \Rightarrow L = L \Rightarrow L = id \Rightarrow *R = id \Rightarrow \dots$
  - The production used is  $L \rightarrow *R$
  - The following rightmost derivation in *reverse* does not exist (and hence reduction by  $R \rightarrow L$  on '=' in state 2 is illegal)  
 $id = id \Leftarrow L = id \Leftarrow R = id \dots$
- Generalization of the above example
  - In some situations, when a state  $i$  appears on top of the stack, a viable prefix  $\beta\alpha$  may be on the stack such that  $\beta A$  cannot be followed by 'a' in any right sentential form
  - Thus, the reduction by  $A \rightarrow \alpha$  would be invalid on 'a'
  - In the above example,  $\beta = \epsilon$ ,  $\alpha = L$ , and  $A = R$ ;  $L$  cannot be reduced to  $R$  on '=', since it would lead to the above illegal derivation sequence

# LR(1) Parsers

- LR(1) items are of the form  $[A \rightarrow \alpha.\beta, a]$ ,  $a$  being the “lookahead” symbol
- Lookahead symbols have no part to play in shift items, but in reduce items of the form  $[A \rightarrow \alpha., a]$ , reduction by  $A \rightarrow \alpha$  is valid only if the next input symbol is ‘ $a$ ’
- An LR(1) item  $[A \rightarrow \alpha.\beta, a]$  is *valid* for a viable prefix  $\gamma$ , if there is a derivation  $S \Rightarrow_{rm}^* \delta A w \Rightarrow_{rm} \delta \alpha \beta w$ , where,  $\gamma = \delta \alpha$ ,  $a = \text{first}(w)$  or  $w = \epsilon$  and  $a = \$$
- Consider the grammar:  $S' \rightarrow S, S \rightarrow aSb \mid \epsilon$ 
  - $[S \rightarrow a.Sb, \$]$  is valid for the VP  $a$ ,  $S' \Rightarrow S \Rightarrow aSb$
  - $[S \rightarrow a.Sb, b]$  is valid for the VP  $aa$ ,  
 $S' \Rightarrow S \Rightarrow aSb \Rightarrow aaSbb$
  - $[S \rightarrow .., \$]$  is valid for the VP  $\epsilon$ ,  $S' \Rightarrow S \Rightarrow \epsilon$
  - $[S \rightarrow aSb., b]$  is valid for the VP  $aaSb$ ,  
 $S' \Rightarrow S \Rightarrow aSb \Rightarrow aaSbb$

# LR(1) Grammar - Example 1

Grammar

$$S' \rightarrow S, S \rightarrow aSb, S \rightarrow \epsilon$$

State 0

$$S' \rightarrow .S, \$$$

$$S \rightarrow .aSb, \$$$

$$S \rightarrow ., \$$$

State 4

$$S \rightarrow a.Sb, b$$

$$S \rightarrow .aSb, b$$

$$S \rightarrow ., b$$

State 1

$$S' \rightarrow S., \$$$

State 5

$$S \rightarrow aSb., \$$$

State 2

$$S \rightarrow a.Sb, \$$$

$$S \rightarrow .aSb, b$$

$$S \rightarrow ., b$$

State 6

$$S \rightarrow aS.b, b$$

State 7

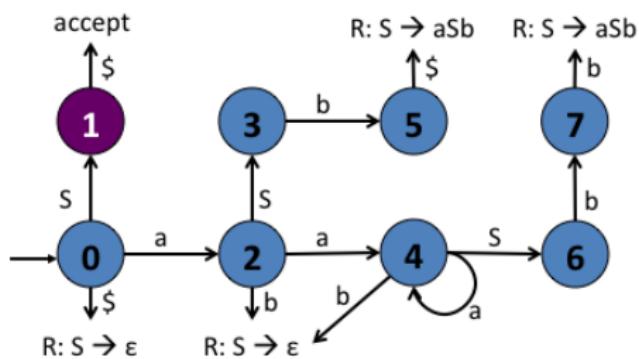
$$S \rightarrow aSb., b$$

State 3

$$S \rightarrow aS.b, \$$$

Grammar is  
LR(1)

	a	b	\$	S
0	S2		R: $S \rightarrow \epsilon$	1
1			accept	
2	S4	R: $S \rightarrow \epsilon$		3
3		S5		
4	S4	R: $S \rightarrow \epsilon$		6
5			R: $S \rightarrow aSb$	
6		S7		
7		R: $S \rightarrow aSb$		



# Closure of a Set of LR(1) Items

```
Itemset closure(I){ /* I is a set of LR(1) items */
    while (more items can be added to I) {
        for each item [A → α.Bβ, a] ∈ I {
            for each production B → γ ∈ G
                for each symbol b ∈ first(βa)
                    if (item [B → .γ, b] ∉ I) add item [B → .γ, b] to I
        }
    }
    return I
}
```

Grammar $S' \rightarrow S$ $S \rightarrow aSb \mid \epsilon$	State 0 $S' \rightarrow .S, \$$ $S \rightarrow .aSb, \$$ $S \rightarrow ., \$$	State 3 $S \rightarrow aS.b, \$$	State 4 $S \rightarrow a.Sb, b$ $S \rightarrow .aSb, b$ $S \rightarrow ., b$	State 7 $S \rightarrow aSb., b$
--	---	-------------------------------------	---	------------------------------------

# GOTO set computation

```
Itemset GOTO(I, X){ /* I is a set of LR(1) items
X is a grammar symbol, a terminal or a nonterminal */
Let I' = {[A → αX.β, a] | [A → α.Xβ, a] ∈ I};
return (closure(I'))
}
```

Grammar $S' \rightarrow S$ $S \rightarrow aSb \mid \epsilon$	<u>State 0</u> $S' \rightarrow .S, \$$ $S \rightarrow .aSb, \$$ $S \rightarrow ., \$$	<u>State 1</u> $S' \rightarrow S., \$$	<u>State 2</u> $S \rightarrow a.Sb, \$$ $S \rightarrow .aSb, b$ $S \rightarrow ., b$	<u>State 4</u> $S \rightarrow a.Sb, b$ $S \rightarrow .aSb, b$ $S \rightarrow ., b$
--	--	---	---	--

$$\text{GOTO}(0, S) = 1, \quad \text{GOTO}(0, a) = 2, \quad \text{GOTO}(2, a) = 4$$

# Construction of Sets of Canonical of LR(1) Items

```
void Set_of_item_sets(G){ /* G' is the augmented grammar */  
    C = {closure({S' → .S, $})};/* C is a set of LR(1) item sets */  
    while (more item sets can be added to C) {  
        for each item set I ∈ C and each grammar symbol X  
            /* X is a grammar symbol, a terminal or a nonterminal */  
            if ((GOTO(I, X) ≠ ∅) && (GOTO(I, X) ∉ C))  
                C = C ∪ GOTO(I, X)  
    }  
}
```

- Each set in  $C$  (above) corresponds to a state of a DFA (LR(1) DFA)
- This is the DFA that recognizes viable prefixes

# LR(1) DFA Construction - Example 1

Grammar

$S' \rightarrow S, S \rightarrow aSb, S \rightarrow \epsilon$

State 0

$S' \rightarrow .S, \$$

$S \rightarrow .aSb, \$$

$S \rightarrow ., \$$

State 4

$S \rightarrow a.Sb, b$

$S \rightarrow .aSb, b$

$S \rightarrow ., b$

State 1

$S' \rightarrow S., \$$

State 5

$S \rightarrow aSb., \$$

State 2

$S \rightarrow a.Sb, \$$

$S \rightarrow .aSb, b$

$S \rightarrow ., b$

State 6

$S \rightarrow aS.b, b$

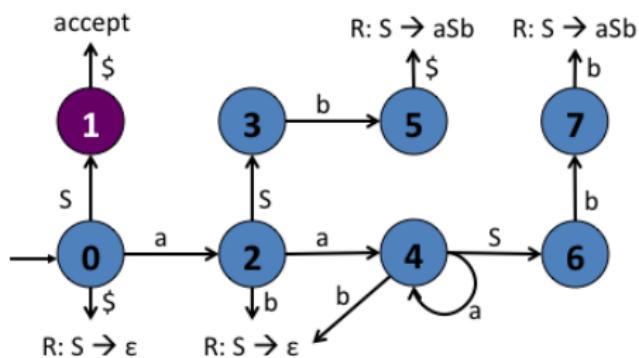
State 7  
 $S \rightarrow aSb., b$

State 3

$S \rightarrow aS.b, \$$

Grammar is  
LR(1)

	a	b	\$	S
0	S2		R: $S \rightarrow \epsilon$	1
1			accept	
2	S4	R: $S \rightarrow \epsilon$		3
3		S5		
4	S4	R: $S \rightarrow \epsilon$		6
5			R: $S \rightarrow aSb$	
6		S7		
7		R: $S \rightarrow aSb$		



# Construction of an LR(1) Parsing Table

Let  $C = \{I_0, I_1, \dots, I_i, \dots, I_n\}$  be the canonical LR(1) collection of items, with the corresponding states of the parser being  $0, 1, \dots, i, \dots, n$ . Without loss of generality, let  $0$  be the initial state of the parser (containing the item  $[S' \rightarrow .S, \$]$ )

Parsing actions for state  $i$  are determined as follows

1. If  $([A \rightarrow \alpha.a\beta, b] \in I_i) \& ([A \rightarrow \alpha a.\beta, b] \in I_j)$   
set ACTION[i, a] = *shift j* /\*  $a$  is a terminal symbol \*/

2. If  $([A \rightarrow \alpha., a] \in I_i)$   
set ACTION[i, a] = *reduce A  $\rightarrow \alpha$*

3. If  $([S' \rightarrow S., \$] \in I_i)$  set ACTION[i, \$] = *accept*

S-R or R-R conflicts in the table imply grammar is not LR(1)

4. If  $([A \rightarrow \alpha.A\beta, a] \in I_i) \& ([A \rightarrow \alpha A.\beta, a] \in I_j)$   
set GOTO[i, A] =  $j$  /\*  $A$  is a nonterminal symbol \*/

All other entries not defined by the rules above are made *error*

# Semantic Analysis with Attribute Grammars

## Part 1

Y.N. Srikant

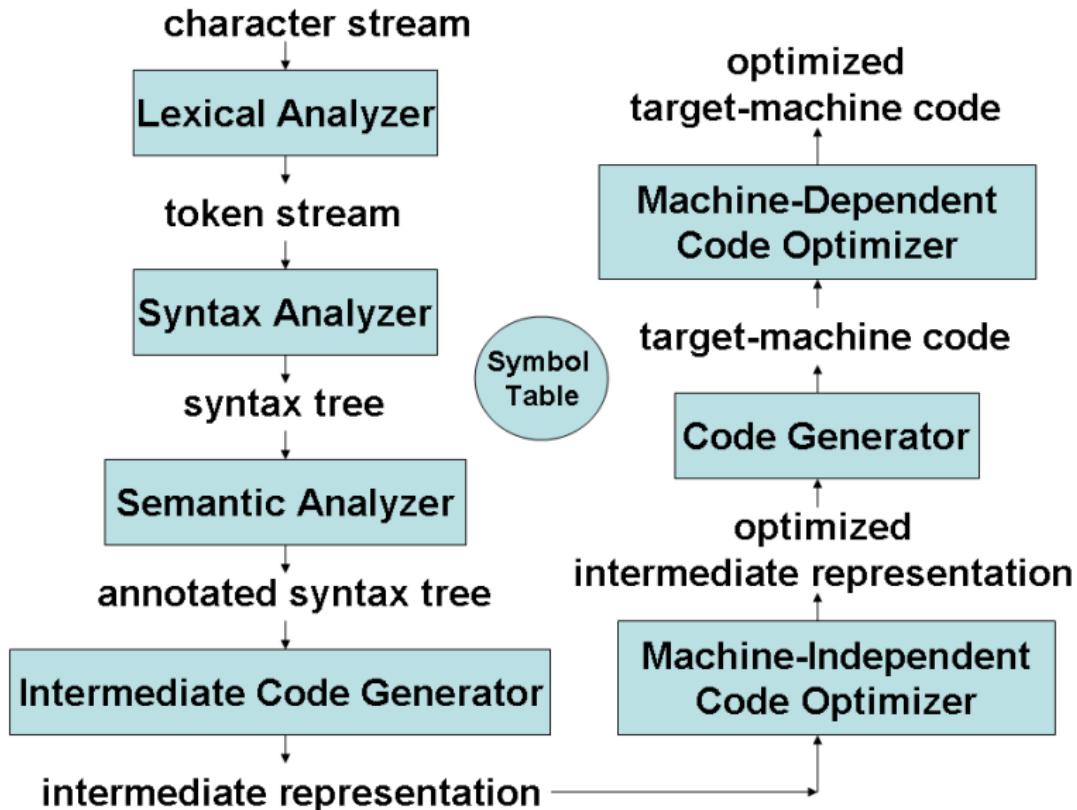
Department of Computer Science and Automation  
Indian Institute of Science  
Bangalore 560 012

NPTEL Course on Principles of Compiler Design

# Outline of the Lecture

- Introduction
- Attribute grammars
- Attributed translation grammars
- Semantic analysis with attributed translation grammars

# Compiler Overview



# Semantic Analysis

- Semantic consistency that cannot be handled at the parsing stage is handled here
- Parsers cannot handle context-sensitive features of programming languages
- These are *static semantics* of programming languages and can be checked by the semantic analyzer
  - Variables are declared before use
  - Types match on both sides of assignments
  - Parameter types and number match in declaration and use
- Compilers can only generate code to check *dynamic semantics* of programming languages at runtime
  - whether an overflow will occur during an arithmetic operation
  - whether array limits will be crossed during execution
  - whether recursion will cross stack limits
  - whether heap memory will be insufficient

# Static Semantics

```
int dot_prod(int x[], int y[]){
    int d, i; d = 0;
    for (i=0; i<10; i++) d += x[i]*y[i];
    return d;
}
main() {
    int p; int a[10], b[10];
    p = dot_prod(a,b);
}
```

Samples of static semantic checks in *main*

- Types of *p* and return type of *dot\_prod* match
- Number and type of the parameters of *dot\_prod* are the same in both its declaration and use
- *p* is declared before use, same for *a* and *b*

# Static Semantics: Errors given by gcc Compiler

```
int dot_product(int a[], int b[]) { ... }

1 main(){int a[10]={1,2,3,4,5,6,7,8,9,10};
2 int b[10]={1,2,3,4,5,6,7,8,9,10};
3 printf("%d", dot_product(b));
4 printf("%d", dot_product(a,b,a));
5 int p[10]; p=dotproduct(a,b); printf("%d",p); }
```

In function 'main':

error in 3: too few arguments to fn 'dot\_product'  
error in 4: too many arguments to fn 'dot\_product'  
error in 5: incompatible types in assignment  
warning in 5: format '%d' expects type 'int', but  
argument 2 has type 'int \*'

# Static Semantics

```
int dot_prod(int x[], int y[]){
    int d, i; d = 0;
    for (i=0; i<10; i++) d += x[i]*y[i];
    return d;
}
main() {
    int p; int a[10], b[10];
    p = dot_prod(a,b);
}
```

Samples of static semantic checks in *dot\_prod*

- *d* and *i* are declared before use
- Type of *d* matches the return type of *dot\_prod*
- Type of *d* matches the result type of “\*”
- Elements of arrays *x* and *y* are compatible with “\*”

# Dynamic Semantics

```
int dot_prod(int x[], int y[]){
    int d, i; d = 0;
    for (i=0; i<10; i++) d += x[i]*y[i];
    return d;
}
main() {
    int p; int a[10], b[10];
    p = dot_prod(a,b);
}
```

Samples of dynamic semantic checks in *dot\_prod*

- Value of *i* does not exceed the declared range of arrays *x* and *y* (both lower and upper)
- There are no overflows during the operations of “\*” and “+” in *d* += *x*[*i*]\**y*[*i*]

# Dynamic Semantics

```
int fact(int n) {  
    if (n==0) return 1;  
    else return (n*fact(n-1));  
}  
main(){int p; p = fact(10); }
```

Samples of dynamic semantic checks in *fact*

- Program stack does not overflow due to recursion
- There is no overflow due to “\*” in  $n * \text{fact}(n-1)$

# Semantic Analysis

- Type information is stored in the symbol table or the syntax tree
  - Types of variables, function parameters, array dimensions, etc.
  - Used not only for semantic validation but also for subsequent phases of compilation
- If declarations need not appear before use (as in C++), semantic analysis needs more than one pass
- Static semantics of PL can be specified using attribute grammars
- Semantic analyzers can be generated semi-automatically from attribute grammars
- Attribute grammars are extensions of context-free grammars

# Attribute Grammars

- Let  $G = (N, T, P, S)$  be a CFG and let  $V = N \cup T$ .
- Every symbol  $X$  of  $V$  has associated with it a set of *attributes* (denoted by  $X.a$ ,  $X.b$ , etc.)
- Two types of attributes: *inherited* (denoted by  $AI(X)$ ) and *synthesized* (denoted by  $AS(X)$ )
- Each attribute takes values from a specified domain (finite or infinite), which is its *type*
  - Typical domains of attributes are, integers, reals, characters, strings, booleans, structures, etc.
  - New domains can be constructed from given domains by mathematical operations such as *cross product*, *map*, etc.
  - array*: a map,  $\mathcal{N} \rightarrow \mathcal{D}$ , where,  $\mathcal{N}$  and  $\mathcal{D}$  are domains of natural numbers and the given objects, respectively
  - structure*: a cross-product,  $A_1 \times A_2 \times \dots \times A_n$ , where  $n$  is the number of fields in the structure, and  $A_i$  is the domain of the  $i^{th}$  field

# Attribute Computation Rules

- A production  $p \in P$  has a set of attribute computation rules (functions)
- Rules are provided for the computation of
  - Synthesized attributes of the LHS non-terminal of  $p$
  - Inherited attributes of the RHS non-terminals of  $p$
- These rules can use attributes of symbols from the production  $p$  only
  - Rules are strictly local to the production  $p$  (no side effects)
- Restrictions on the rules define different types of attribute grammars
  - L-attribute grammars, S-attribute grammars, ordered attribute grammars, absolutely non-circular attribute grammars, circular attribute grammars, etc.

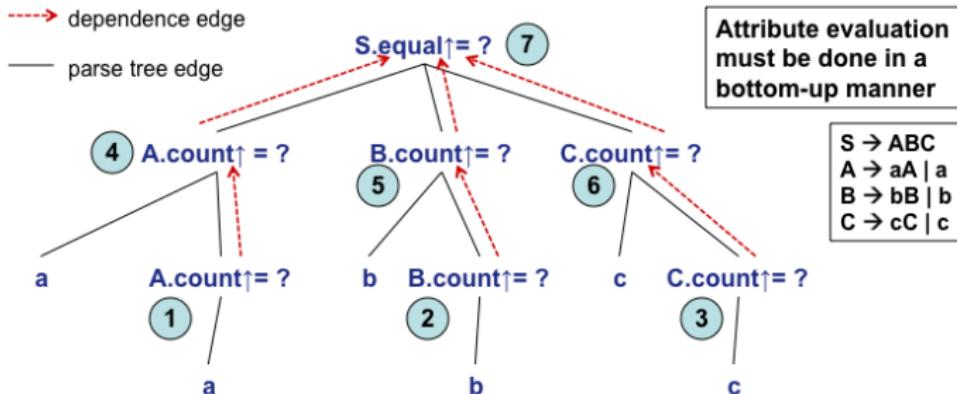
# Synthesized and Inherited Attributes

- An attribute cannot be both synthesized and inherited, but a symbol can have both types of attributes
- Attributes of symbols are evaluated over a parse tree by making passes over the parse tree
- Synthesized attributes are computed in a bottom-up fashion from the leaves upwards
  - Always synthesized from the attribute values of the children of the node
  - Leaf nodes (terminals) have synthesized attributes initialized by the lexical analyzer and cannot be modified
  - An AG with only synthesized attributes is an *S-attributed grammar (SAG)*
  - YACC permits only SAGs
- Inherited attributes flow down from the parent or siblings to the node in question

# Attribute Grammar - Example 1

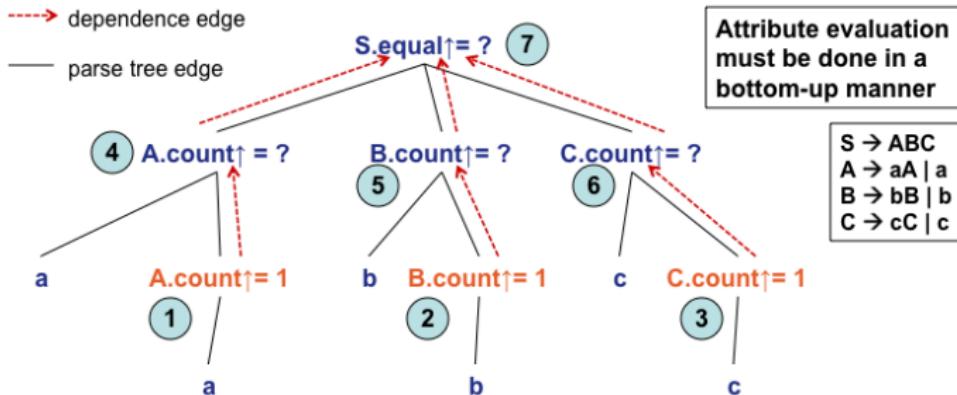
- The following CFG  
 $S \rightarrow A\ B\ C, A \rightarrow aA \mid a, B \rightarrow bB \mid b, C \rightarrow cC \mid c$   
generates:  $L(G) = \{a^m b^n c^p \mid m, n, p \geq 1\}$
- We define an AG (attribute grammar) based on this CFG to generate  $L = \{a^n b^n c^n \mid n \geq 1\}$
- All the non-terminals will have only synthesized attributes
  - $AS(S) = \{equal \uparrow: \{T, F\}\}$
  - $AS(A) = AS(B) = AS(C) = \{count \uparrow: integer\}$

# Attribute Grammar - Example 1 (contd.)



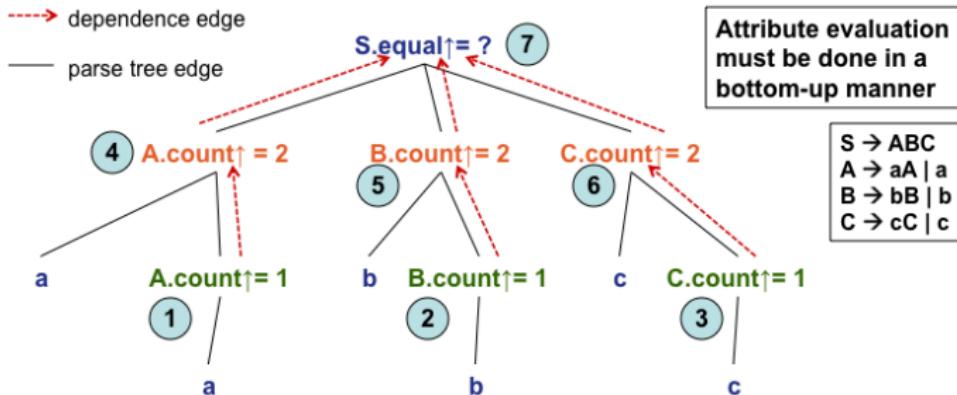
- 1  $S \rightarrow ABC \{S.\text{equal} \uparrow := \text{if } A.\text{count} \uparrow = B.\text{count} \uparrow \& B.\text{count} \uparrow = C.\text{count} \uparrow \text{ then } T \text{ else } F\}$
- 2  $A_1 \rightarrow aA_2 \{A_1.\text{count} \uparrow := A_2.\text{count} \uparrow + 1\}$
- 3  $A \rightarrow a \{A.\text{count} \uparrow := 1\}$
- 4  $B_1 \rightarrow bB_2 \{B_1.\text{count} \uparrow := B_2.\text{count} \uparrow + 1\}$
- 5  $B \rightarrow b \{B.\text{count} \uparrow := 1\}$
- 6  $C_1 \rightarrow cC_2 \{C_1.\text{count} \uparrow := C_2.\text{count} \uparrow + 1\}$
- 7  $C \rightarrow c \{C.\text{count} \uparrow := 1\}$

# Attribute Grammar - Example 1 (contd.)



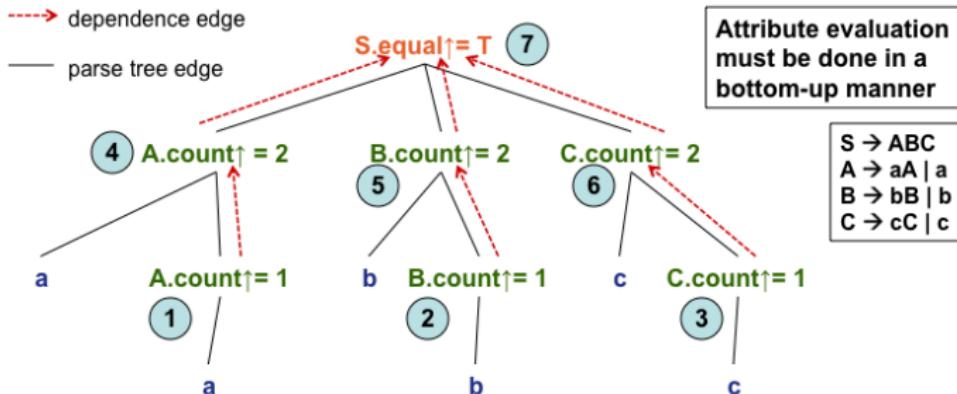
- 1  $S \rightarrow ABC \{S.\text{equal} \uparrow := \text{if } A.\text{count} \uparrow = B.\text{count} \uparrow \& B.\text{count} \uparrow = C.\text{count} \uparrow \text{ then } T \text{ else } F\}$
- 2  $A_1 \rightarrow aA_2 \{A_1.\text{count} \uparrow := A_2.\text{count} \uparrow + 1\}$
- 3  $A \rightarrow a \{A.\text{count} \uparrow := 1\}$
- 4  $B_1 \rightarrow bB_2 \{B_1.\text{count} \uparrow := B_2.\text{count} \uparrow + 1\}$
- 5  $B \rightarrow b \{B.\text{count} \uparrow := 1\}$
- 6  $C_1 \rightarrow cC_2 \{C_1.\text{count} \uparrow := C_2.\text{count} \uparrow + 1\}$
- 7  $C \rightarrow c \{C.\text{count} \uparrow := 1\}$

# Attribute Grammar - Example 1 (contd.)



- 1  $S \rightarrow ABC \{S.\text{equal} \uparrow := \text{if } A.\text{count} \uparrow = B.\text{count} \uparrow \& B.\text{count} \uparrow = C.\text{count} \uparrow \text{ then } T \text{ else } F\}$
- 2  $A_1 \rightarrow aA_2 \{A_1.\text{count} \uparrow := A_2.\text{count} \uparrow + 1\}$
- 3  $A \rightarrow a \{A.\text{count} \uparrow := 1\}$
- 4  $B_1 \rightarrow bB_2 \{B_1.\text{count} \uparrow := B_2.\text{count} \uparrow + 1\}$
- 5  $B \rightarrow b \{B.\text{count} \uparrow := 1\}$
- 6  $C_1 \rightarrow cC_2 \{C_1.\text{count} \uparrow := C_2.\text{count} \uparrow + 1\}$
- 7  $C \rightarrow c \{C.\text{count} \uparrow := 1\}$

# Attribute Grammar - Example 1 (contd.)



- 1  $S \rightarrow ABC \{S.\text{equal} \uparrow := \text{if } A.\text{count} \uparrow = B.\text{count} \uparrow \& B.\text{count} \uparrow = C.\text{count} \uparrow \text{ then } T \text{ else } F\}$
- 2  $A_1 \rightarrow aA_2 \{A_1.\text{count} \uparrow := A_2.\text{count} \uparrow + 1\}$
- 3  $A \rightarrow a \{A.\text{count} \uparrow := 1\}$
- 4  $B_1 \rightarrow bB_2 \{B_1.\text{count} \uparrow := B_2.\text{count} \uparrow + 1\}$
- 5  $B \rightarrow b \{B.\text{count} \uparrow := 1\}$
- 6  $C_1 \rightarrow cC_2 \{C_1.\text{count} \uparrow := C_2.\text{count} \uparrow + 1\}$
- 7  $C \rightarrow c \{C.\text{count} \uparrow := 1\}$

# Attribute Dependence Graph

- Let  $T$  be a parse tree generated by the CFG of an AG,  $G$ .
- The *attribute dependence graph* (dependence graph for short) for  $T$  is the directed graph,  $DG(T) = (V, E)$ , where

$V = \{b | b \text{ is an attribute instance of some tree node}\}$ , and

$E = \{(b, c) | b, c \in V, b \text{ and } c \text{ are attributes of grammar symbols in the same production } p \text{ of } B, \text{ and the value of } b \text{ is used for computing the value of } c \text{ in an attribute computation rule associated with production } p\}$

# Attribute Dependence Graph

- An AG  $G$  is *non-circular*, iff for all trees  $T$  derived from  $G$ ,  $DG(T)$  is acyclic
  - Non-circularity is very expensive to determine (exponential in the size of the grammar)
  - Therefore, our interest will be in subclasses of AGs whose non-circularity can be determined efficiently
- Assigning consistent values to the attribute instances in  $DG(T)$  is *attribute evaluation*

# Attribute Evaluation Strategy

- Construct the parse tree
- Construct the dependence graph
- Perform topological sort on the dependence graph and obtain an evaluation order
- Evaluate attributes according to this order using the corresponding attribute evaluation rules attached to the respective productions
- Multiple attributes at a node in the *parse tree* may result in that node to be visited multiple number of times
  - Each visit resulting in the evaluation of at least one attribute

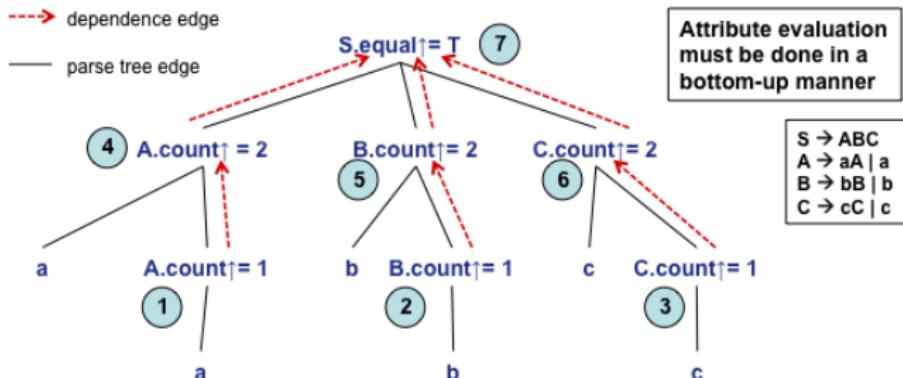
# Attribute Evaluation Algorithm

**Input:** A parse tree  $T$  with unevaluated attribute instances

**Output:**  $T$  with consistent attribute values

```
{ Let  $(V, E) = DG(T)$ ;  
Let  $W = \{b \mid b \in V \text{ & } \text{indegree}(b) = 0\}$ ;  
while  $W \neq \emptyset$  do  
{ remove some  $b$  from  $W$ ;  
   $\text{value}(b) :=$  value defined by appropriate attribute  
    computation rule;  
  for all  $(b, c) \in E$  do  
  {  $\text{indegree}(c) := \text{indegree}(c) - 1$ ;  
    if  $\text{indegree}(c) = 0$  then  $W := W \cup \{c\}$ ;  
  }  
}  
}
```

# Dependence Graph for Example 1



1,2,3,4,5,6,7 and 2,3,6,5,1,4,7 are two possible evaluation orders. 1,4,2,5,3,6,7 can be used with LR-parsing. The right-most derivation is below (its reverse is LR-parsing order)

$S \Rightarrow ABC \Rightarrow ABcC \Rightarrow ABcc \Rightarrow AbBcc \Rightarrow Abbcc \Rightarrow aAbbcc \Rightarrow aabbcc$

1. A.count = 1 { $A \rightarrow a$ , {A.count := 1}}
4. A.count = 2 { $A_1 \rightarrow aA_2$ , {A<sub>1</sub>.count := A<sub>2</sub>.count + 1}}
2. B.count = 1 { $B \rightarrow b$ , {B.count := 1}}
5. B.count = 2 { $B_1 \rightarrow bB_2$ , {B<sub>1</sub>.count := B<sub>2</sub>.count + 1}}
3. C.count = 1 { $C \rightarrow c$ , {C.count := 1}}
6. C.count = 2 { $C_1 \rightarrow cC_2$ , {C<sub>1</sub>.count := C<sub>2</sub>.count + 1}}
7. S.equal = 1 { $S \rightarrow ABC$ , {S.equal := if A.count = B.count & B.count = C.count then T else F}}

# Attribute Grammar - Example 2

- AG for the evaluation of a real number from its bit-string representation

Example: 110.101 = 6.625

- $N \rightarrow L.R, L \rightarrow BL \mid B, R \rightarrow BR \mid B, B \rightarrow 0 \mid 1$

- $AS(N) = AS(R) = AS(B) = \{value \uparrow: real\},$   
 $AS(L) = \{length \uparrow: integer, value \uparrow: real\}$

- ①  $N \rightarrow L.R \{N.value \uparrow := L.value \uparrow + R.value \uparrow\}$
- ②  $L \rightarrow B \{L.value \uparrow := B.value \uparrow; L.length \uparrow := 1\}$
- ③  $L_1 \rightarrow BL_2 \{L_1.length \uparrow := L_2.length \uparrow + 1;$   
 $L_1.value \uparrow := B.value \uparrow * 2^{L_2.length \uparrow} + L_2.value \uparrow\}$
- ④  $R \rightarrow B \{R.value \uparrow := B.value \uparrow / 2\}$
- ⑤  $R_1 \rightarrow BR_2 \{R_1.value \uparrow := (B.value \uparrow + R_2.value \uparrow) / 2\}$
- ⑥  $B \rightarrow 0 \{B.value \uparrow := 0\}$
- ⑦  $B \rightarrow 1 \{B.value \uparrow := 1\}$