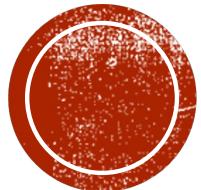


FA TO RE

COMP 4200 – Formal Language





ARDEN'S THEOREM

ARDEN'S THEOREM

- If P and Q are RE over Σ , and if P does not contain \emptyset then the following equation in R is given by $R = Q + RP$ has a unique solution $R = QP^*$.
- Proof:

$$R = Q + RP$$

Replace R with QP^*

$$R = Q + (QP^*)P = Q(\varepsilon + P^*P) = \textcircled{QP^*}$$

Proved



- $$\begin{aligned} R &= Q + RP \\ &= Q + (Q + RP)P \\ &= Q + PQ + (Q+RP)P^2 \\ &= Q + PQ + P^2Q + RP^3 \\ &= Q + PQ + P^2Q + P^3Q + \dots + P^{n+1}R \\ &= Q + PQ + P^2Q + P^3Q + \dots + QP^n + QP^*P^{n+1} \\ &= Q[1 + P + P^2 + P^3 + \dots + P^n + P^*P^{n+1}] \\ &= QP^* \end{aligned}$$

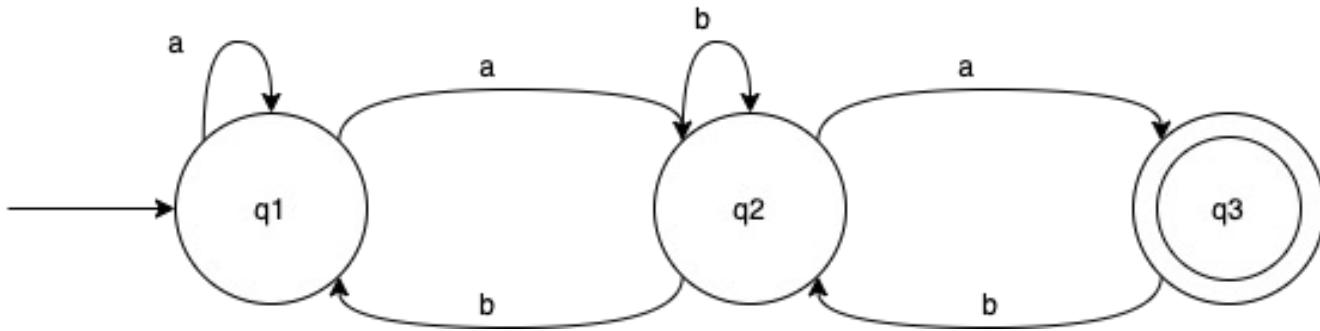


DESIGNING RE FROM LANGUAGE

- Design RE for the following language over {a,b}
 - Language accepting strings of length exactly 2
 - $L = \{aa, ab, ba, bb\}$
 - $R = aa + ab + ba + bb = a(a+b) + b(a+b) = (a+b)(a+b)$
 - Language accepting strings of length atleast 2
 - $L = \{aa, ab, ba, bb, aaa, aab, \dots\}$
 - $R = (a+b)(a+b)(a+b)^*$
 - Language accepting strings of length atmost 2
 - $L = \{\epsilon, a, b, aa, ab, ba, bb\}$
 - $R = \epsilon + a + b + aa + ab + ba + bb = (\epsilon + a + b) (\epsilon + a + b)$



NFA TO RE



Incoming transitions

$$q_3 = q_2a \rightarrow 1$$

$$q_2 = q_2b + q_1a + q_3b \rightarrow 2$$

$$q_1 = \epsilon + q_1a + q_2b \rightarrow 3$$

$$R = Q + RP$$

$$R = QP^*$$



$$\begin{aligned} q_1 &\rightarrow aq_1 + bq_2 + \varepsilon \\ q_2 &\rightarrow aq_1 + bq_2 + bq_3 \\ q_3 &\rightarrow aq_2 \end{aligned}$$

$$q_3 = [q_2 b + q_1 a + q_3 b] a$$

$$q_3 = q_2 ba + q_3 aa + q_3 ba$$

$$\begin{aligned} q_2 &= q_2 b + q_1 a + q_3 b \\ &= q_2 b + q_1 a + q_2 ab \\ q_2 &= q_2 [b + ab] + q_1 a \end{aligned}$$

R R P Q

$$q_2 = q_1 a [b + ab]^*$$

$$q_1 = \varepsilon + q_1 a + [q_1 ab [b + ab]^*]$$

$$q_1 = \varepsilon + q_1 a + q_1 ab (b + ab)^*$$

$$q_1 = q_1 [a + ab (b + ab)^*] + \varepsilon$$

R R P Q

$$q_1 = \varepsilon [a + ab (b + ab)^*]^*$$

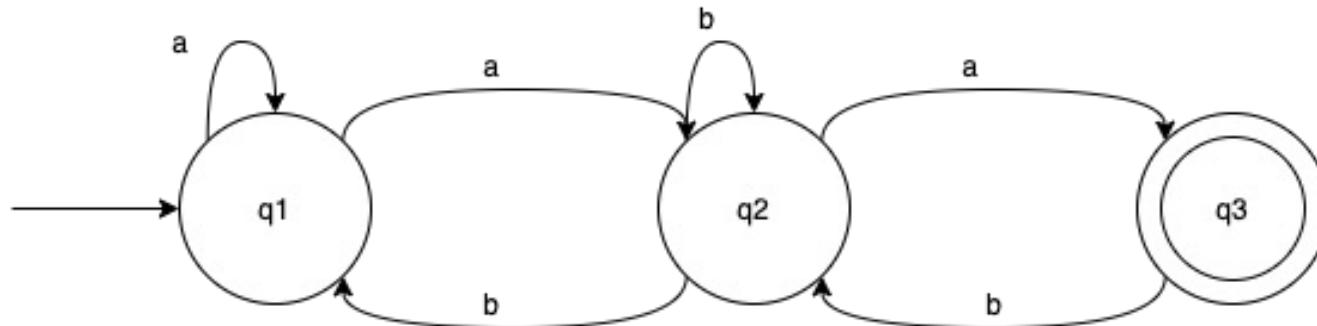
$$q_2 = [a + ab] (b + ab)^* a$$

$$q_3 = [(a + ab) (b + ab)^*]^* a (b + ab)^* a$$

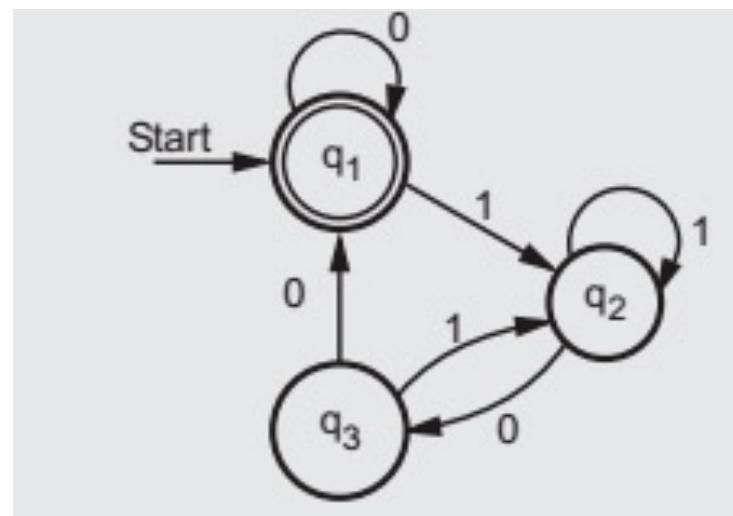
- $q_3 = q_2a$
 - $= [q_1a + q_2b + q_3b]a$
 - $= q_1aa + q_2ba + q_3ba \rightarrow 4$
- $q_2 = q_1a + q_2b + q_3b$
 - putting value of q_3 from eq(1)
 - $= q_1a + q_2b + (q_2a)b$
 - $= q_1a + q_2b + q_2ab$
 - $= q_1a + q_2(b + ab) \rightarrow R = Q + RP \rightarrow R = q_2, Q = q_1a, P = (b + ab) \rightarrow R = QP^*$
 - $= q_1a(b + ab)^* \rightarrow 5$
- $q_1 = \varepsilon + q_1a + q_2b$
 - Putting eq(5) in q_1
 - $= \varepsilon + q_1a + [q_1a(b + ab)^*]b$
 - $= \varepsilon + q_1[a + [a(b + ab)^*]b] \rightarrow R = Q + RP \rightarrow R = q_1, Q = \varepsilon, P = a + [a(b + ab)^*]b \rightarrow R = QP^*$
 - $= \varepsilon[a + [a(b + ab)^*]b]^* \rightarrow \varepsilon.R = R$
 - $= [a + [a(b + ab)^*]b]^* \rightarrow 6$



- Final state, $q_3 = q_2a, q_2$ from eq(5),
 - $= [q_1a(b+ab)^*]a$
 - $= [a + [a(b+ab)^*]b]^* a(b+ab)^*a \rightarrow q_1$ from eq(6)
 - $q_3 = [a + [a(b+ab)^*]b]^* a(b+ab)^*a$
 - Final expression is in terms of input symbols.



EXAMPLE

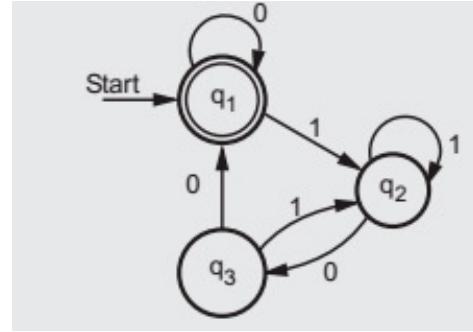


EXAMPLE

- Incoming transitions
 - $q_1 = \epsilon + 0q_1 + 0q_3 \rightarrow 1$
 - $q_2 = 1q_1 + 1q_2 + 1q_3 \rightarrow 2$
 - $q_3 = 0q_2 \rightarrow 3$

- For getting RE we need to solve the final state q_1 .

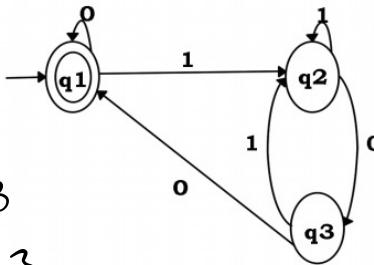
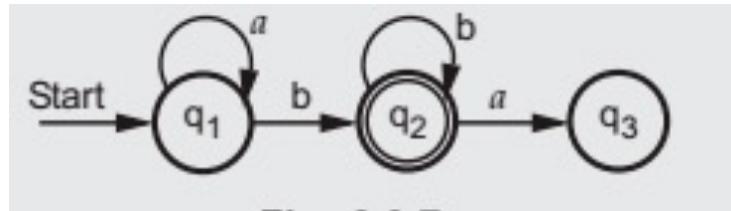
- $q_2 = 1q_1 + 1q_2 + 1q_3$
 - $= 1q_1 + 1q_2 + 1(0q_2) \rightarrow q_3$ from eq(3)
 - $= 1q_1 + q_2(1+10) \rightarrow R = Q + RP \rightarrow R = q_2, Q = 1q_1, P = (1+10) \rightarrow R = QP^*$
 - $= 1q_1(1+10)^* \rightarrow 4$
- $q_1 = \epsilon + 0q_1 + 0q_3$
 - $= \epsilon + 0q_1 + 0(0q_2) \rightarrow q_3$ from eq(3)
 - $= \epsilon + 0q_1 + 00q_2$
 - q_2 from eq(4)
 - $= \epsilon + 0q_1 + 00[1q_1(1+10)^*]$
 - $= \epsilon + 0q_1 + 001q_1(1+10)^*$
 - $= \epsilon + q_1[0 + 001(1+10)^*] \rightarrow R = Q + RP \rightarrow R = q_1, Q = \epsilon, P = [0 + 001(1+10)^*] \rightarrow R = QP^*$
 - $= \epsilon [0 + 001(1+10)^*]^* \rightarrow \epsilon R = R$
 - $[0 + 001(1+10)^*]^*$



$$\begin{aligned}q_1 &= \epsilon + aq_1 \\q_2 &= bq_2 + bq_1 \\q_3 &= aq_2\end{aligned}$$

$$\begin{aligned}P &= RP + Q \\R &= QP^A\end{aligned}$$

TRY YOURSELF!



$$\begin{aligned}q_1 &= \epsilon + 0q_1 + 0q_3 \\q_2 &= 1q_2 + 1q_1 + 1q_3 \\q_3 &= 0q_2\end{aligned}$$

$$\begin{aligned}q_1 &= aq_1 + \epsilon \\&\quad \swarrow R \quad \swarrow P \quad \swarrow Q \\q_2 &= ba^* + b^*q_2 \\&\quad \swarrow Q \quad \swarrow P \quad \swarrow R \\q_2 &= ba^* b^*\end{aligned}$$

$$\begin{aligned}(a) \quad q_1 &= \epsilon a^* \\q_2 &= b a^* + b^* q_2 \\q_2 &= ba^* b^* \\(b) \quad q_1 &= 1a_1 + 10q_2 + 1q_2 \\q_2 &= 1a_1 + a_2 [10+1] \\q_2 &= 1a_1 (10+1)^*\end{aligned}$$



$$q_3 = |q_1 + |q_3 + |q_2$$

$$\therefore q_3 = 0|q_1 + 0|q_3 + 0\left(|q_1|^{(10+1)^*}\right)$$

$$q_3 = q_1 [01 + 011 \left[10 + 1 \right]^{**}] + 0|q_3$$

$$q_3 = q_1 [01 + 011 \left(10 + 1 \right)^*] (01)^*$$

$$q_3 = q_1 [01 + 011 \left(10 + 1 \right)^*] (01)^*$$

$$q_1 = \varepsilon + 0|q_1 + 0\left(q_1 [01 + 011 \left(10 + 1 \right)^*]\right) (01)^*$$

$$q_1 = \varepsilon + \underset{Q}{q_1} \left(0 + 0\left[01 + 011 \left(10 + 1 \right)^* \right] (01)^* \right)$$

$$q_1 = \varepsilon \left(0 + 0\left[01 + 011 \left(10 + 1 \right)^* \right] (01)^* \right)^*$$

EQUIVALENCE OF TWO FA



EQUIVALENCE OF TWO FA

- Two finite automata are equivalent if they accept the same set of strings over L.
- Algorithm:
 - Construct comparison table consisting of $n+1$ columns, where n is number of input symbols .
 - In the **first column**, there will be a pair of vertices (q, q') where $q \in M$ and $q' \in M'$. The **first pair** of vertices will be the **initial states** of the two machines M and M'.
 - The **second column** consists of (q_a, q'_a) , where q_a is **reachable from the initial state** of the machine M for the first input, and q'_a is reachable from the initial state of the machine M' for the first input.
 - The other $n - 2$ columns consist of a pair of vertices from M and M' for $n - 1$ inputs, where $n = 2, 3, \dots, n - 1$.
 - If any new pair of states appear in any of the $n - 1$ next state columns, which were not taken in the first column, take that pair in the present state column and construct subsequent column elements like the first row.



- If a pair of states (q, q') appear in any of the n columns for a pair of states in the present state column, where q is the final state of M and q' is the non-final state of M' or vice versa, terminate the construction and conclude that M and M' are not equivalent.
- If no new pair of states appear, which were not taken in the first column, stop the construction and declare that M and M' are equivalent.
- Simple words:

- For any pair of states (q_i, q_j) , the transition input $a \in \Sigma$ is defined by (q_i, q_j) where,

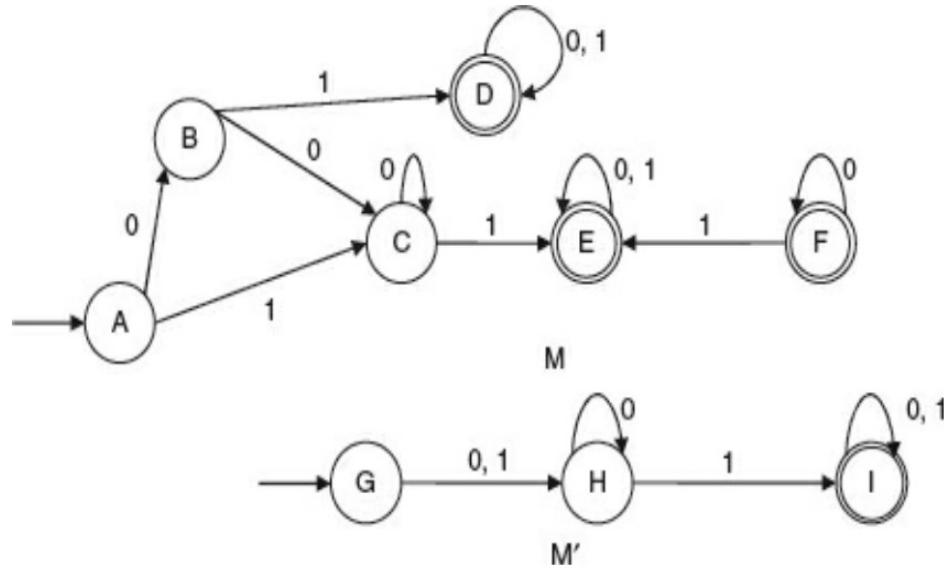
$$\delta(q_i, a) = q_a \text{ and } \delta(q_j, b) = q_b$$

The two automaton are not equivalent if for a pair (q_a, q_b) one is intermediate state and other is final state.

- If the initial state is final state of one automaton, then in second automaton also initial state must be final state for them to be equivalent.



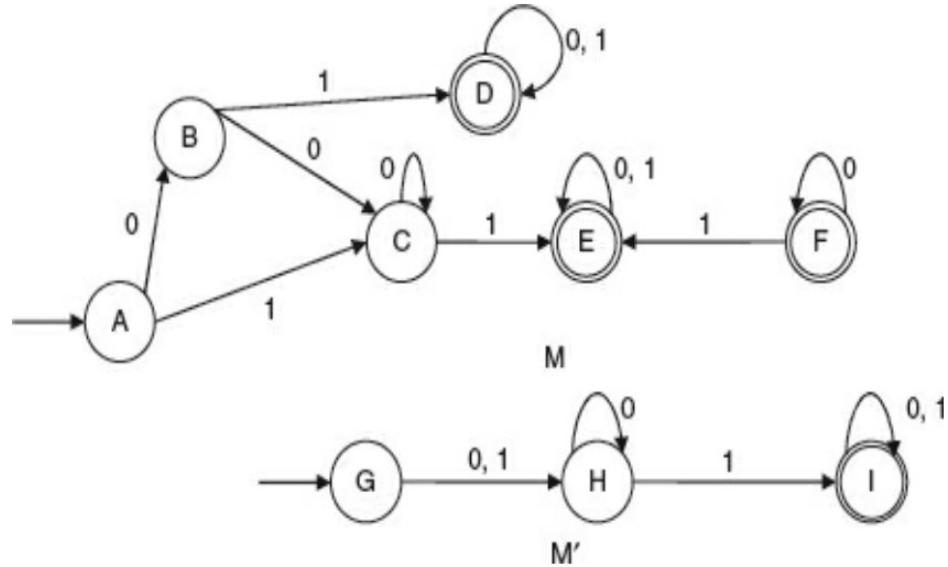
EXAMPLE



States	0	1
(A, G)	(B, H)	(C, I)
(B, H)	(C, H)	(D, I)
(C, H)	(C, H)	(E, I)
(D, I)	(D, I)	(D, I)
(E, I)	(E, I)	(E, I)



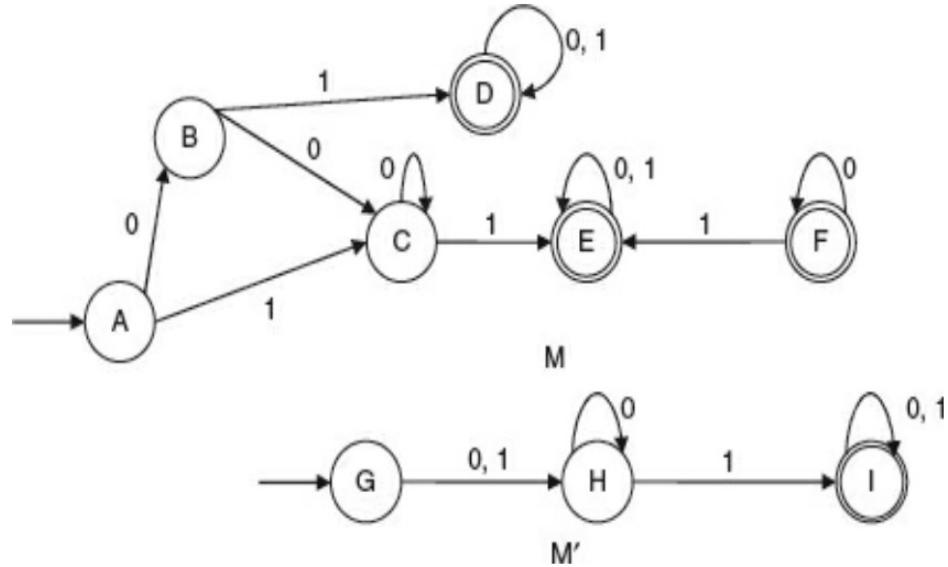
EXAMPLE



States	0	1
(A, G)	(B, H)	(C, H)



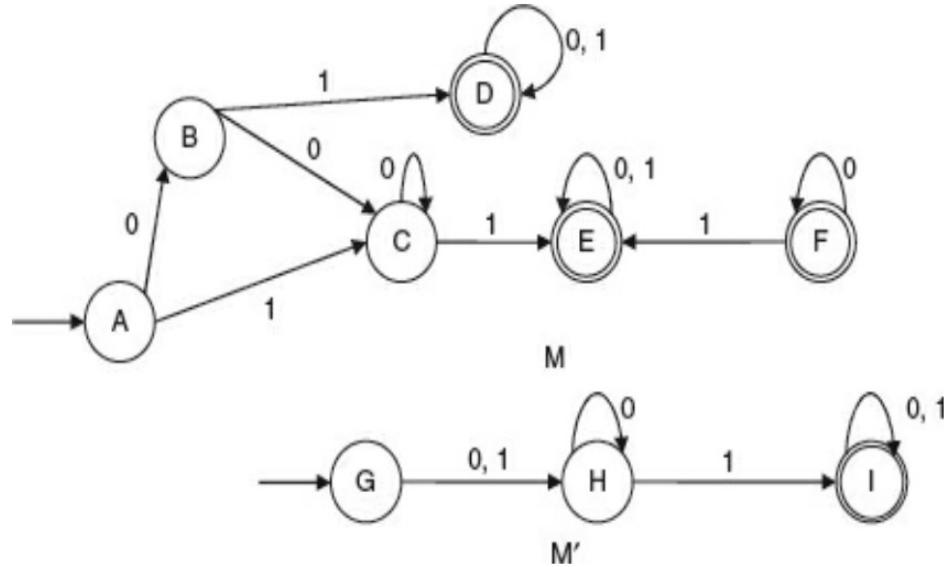
EXAMPLE



States	0	1
(A, G)	(B, H)	(C, H)
(B, H)	(C, H)	(D, I)



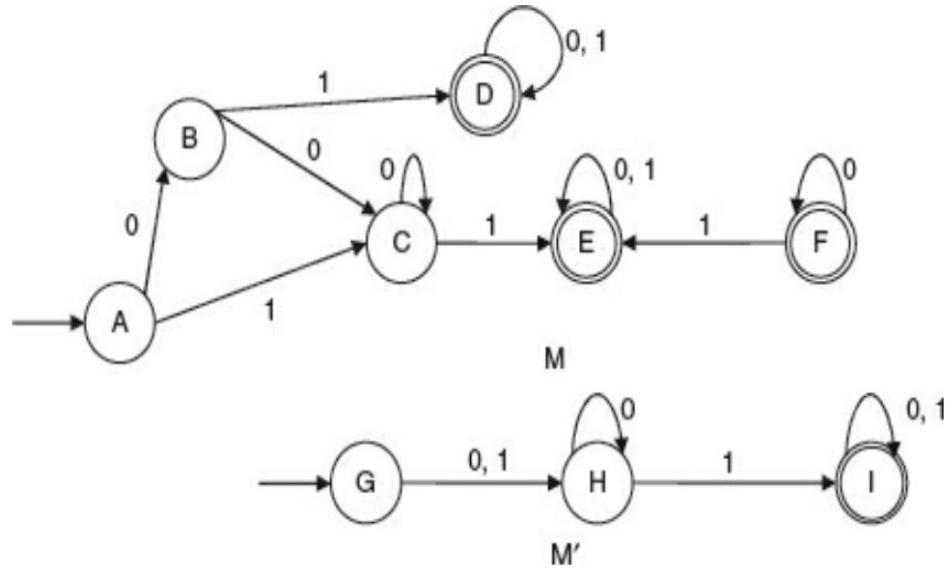
EXAMPLE



States	0	1
(A, G)	(B, H)	(C, H)
(B, H)	(C, H)	(D, I)
(C, H)	(C, H)	(E, I)



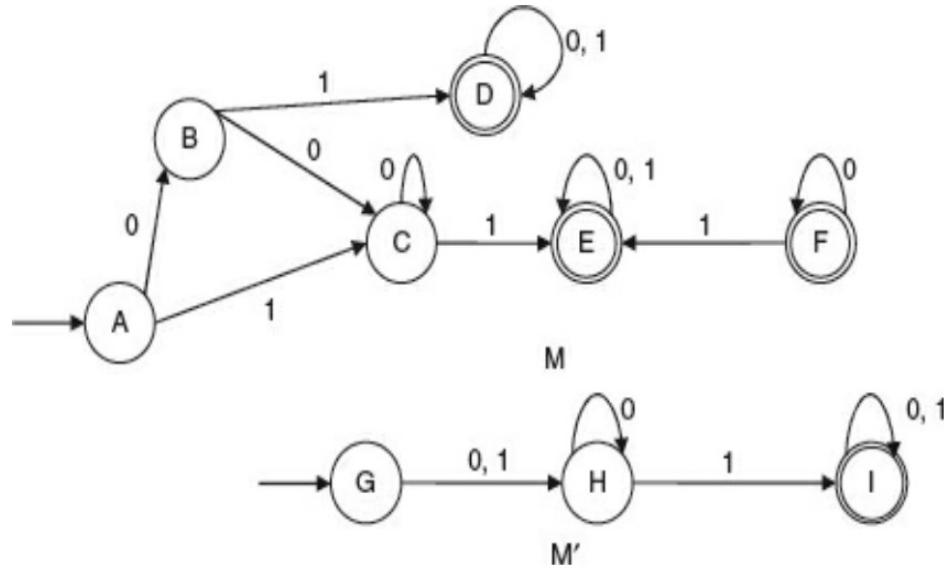
EXAMPLE



States	0	1
(A, G)	(B, H)	(C, H)
(B, H)	(C, H)	(D, I)
(C, H)	(C, H)	(E, I)
(D, I)	(D, I)	(D, I)



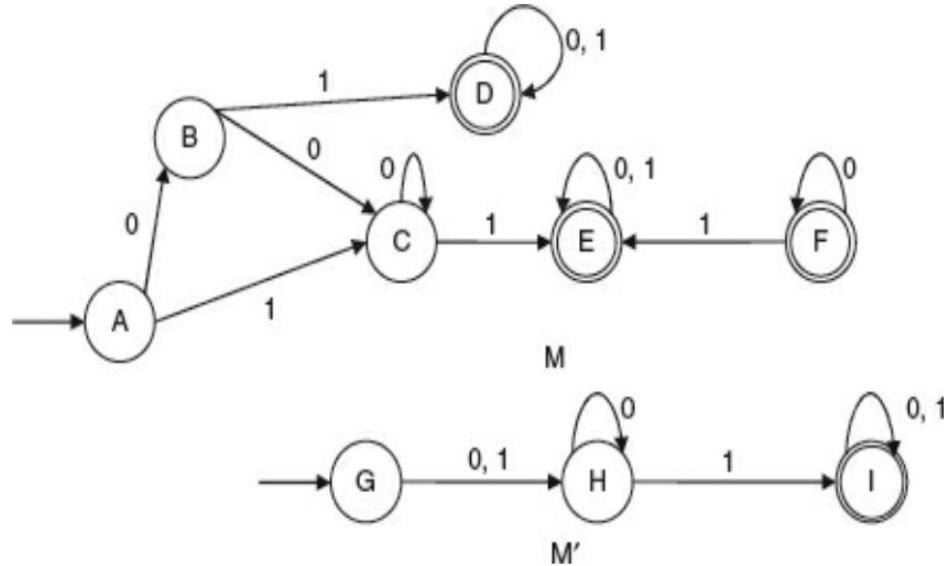
EXAMPLE



States	0	1
(A, G)	(B, H)	(C, H)
(B, H)	(C, H)	(D, I)
(C, H)	(C, H)	(E, I)
(D, I)	(D, I)	(D, I)
(E, I)	(E, I)	(E, I)



EXAMPLE



States	0	1
(A, G)	(B, H)	(C, H)
(B, H)	(C, H)	(D, I)
(C, H)	(C, H)	(E, I)
(D, I)	(D, I)	(D, I)
(E, I)	(E, I)	(E, I)

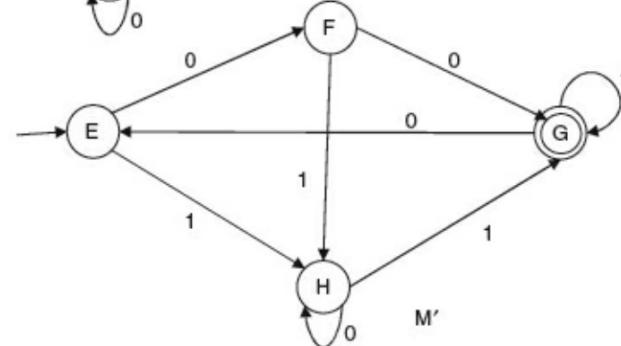
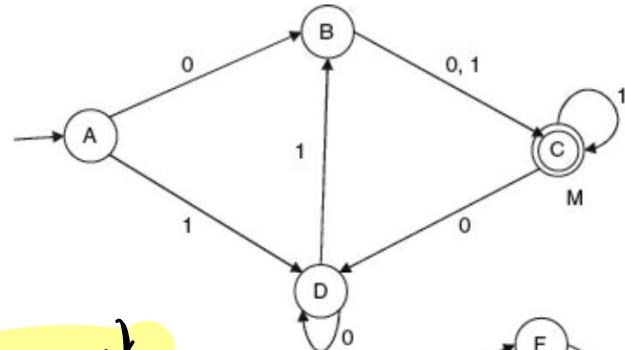
$(B, H) \rightarrow (IS, IS)$
 $(C, H) \rightarrow (IS, IS)$
 $(D, I) \rightarrow (FS, FS)$
 $(E, I) \rightarrow (FS, FS)$

Therefore, the two DFAs are equivalent!!



	0	1
0	(A, E) (B, F)	(D, H)
1	(C, G)	(C, H)

TRY YOURSELF!



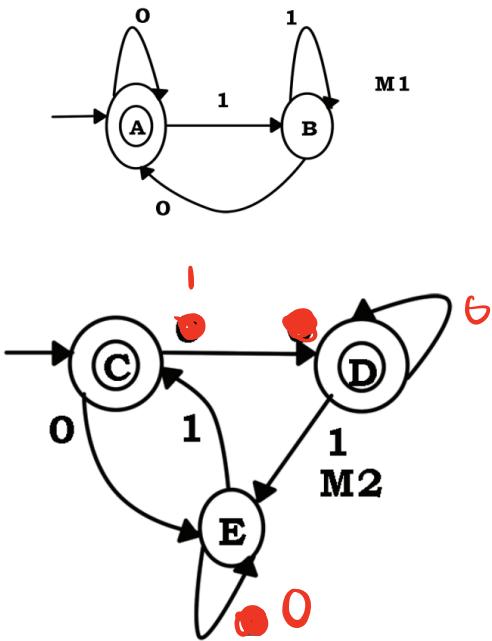
The 2 machines are not equivalent to each other.



	δ	
(A, C)	(A, E)	

TRY YOURSELF!

The 2 machines
are not equivalent to
each other.





PUMPING LEMMA



NON REGULAR LANGUAGE

REGULAR OR NON REGULAR LANGUAGE

$$B = \{0^n 1^n \mid n \geq 0\}.$$

- $C = \{w \mid w \text{ has an equal number of 0s and 1s}\}$, and
- $D = \{w \mid w \text{ has an equal number of occurrences of 01 and 10 as substrings}\}$.



$B = \{0^n1^n \mid n \geq 0\}$.

PUMPING LEMMA

- Technique to prove non regularity → pumping lemma
- The Pumping Lemma states that all regular languages have a special pumping property.
- If a language does not have the pumping property, then it is not regular.
 - So one can use the Pumping Lemma to prove that a given language is not regular.
- The property states that all strings in the language can be “pumped” if they are at least as long as a certain special value, called the *pumping length*.



THEOREM: PUMPING LEMMA

- Let A be a regular language. There is a positive integer p such that any $s \in A$ with $|s| > p$, s can be divided into three pieces, $s = xyz$.
- (p is the pumping length of A .)
- This means that every string $s \in L$ contains a substring that can be repeated any number of times (via a loop).
- The statement "s can be pumped" means that we can write

$s = xyz$, where

1. $xy^i z \in A$ for all $i \geq 0$
2. $|y| > 0$
3. $|xy| \leq p$



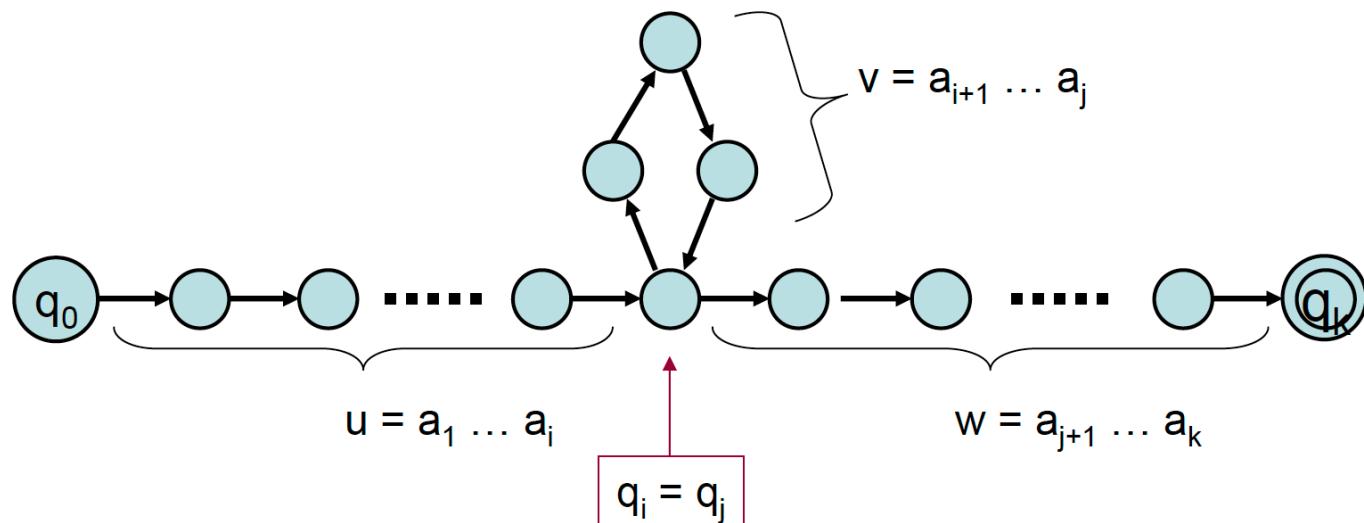
THE PUMPING LEMMA FOR LANGUAGES

- Condition 1: $xy^i z \in A$ for all $i \geq 0$.
- Condition 2: $|y| > 0$.
- Condition 3: $|xy| \leq p$.
- All together, the conditions allow either x or z to be ε , but not both.
- When s is divided into xyz , either x or z may be ε , but condition 2 says that $y \neq \varepsilon$.
- Condition 3 states that the pieces x and y together have length at most p .



THE LOOP

- $u v^m w$ is accepted, since it follows the loop m times (possibly 0 times).



PROOF

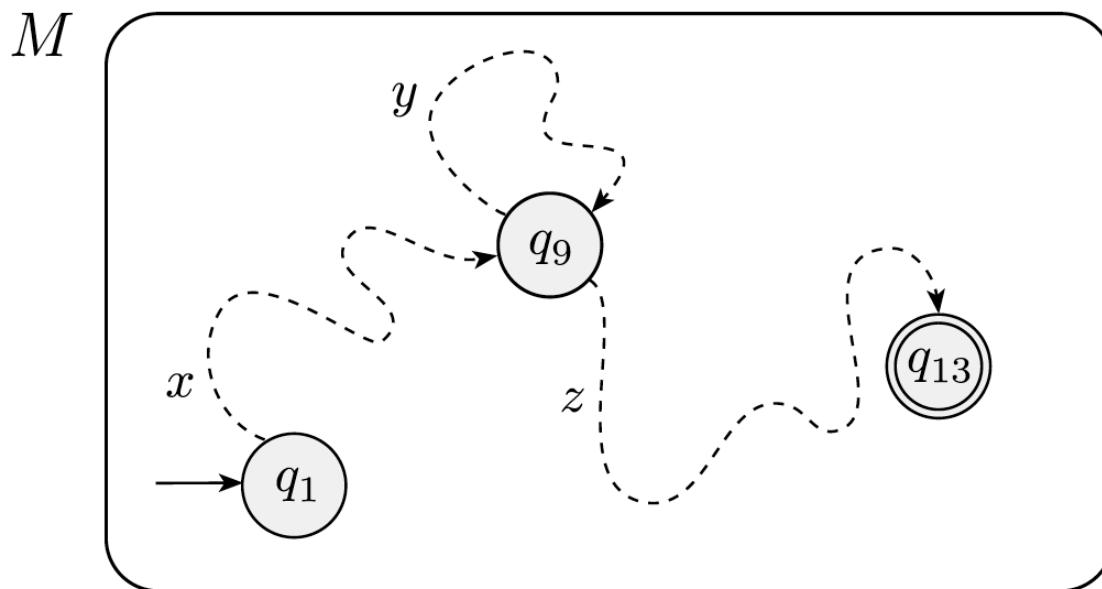
- Let $M = (Q, \Sigma, \delta, q_1, F)$ be a DFA that recognizes A
- Let $p = \text{number of states in the FA}$
- String $s \in A$, length at least p may be broken into the three pieces xyz , i.e., $|s| > p$
- If s in A has length at least p , consider the states that FA goes through for s .
- If we let n be the length of s , then sequence of state is $n+1$.
- Since there are only p states and $|s| > p$, one state must be repeated (via pigeonhole principle)
- So, there is a loop.



$$s = s_1 \ s_2 \ s_3 \ s_4 \ s_5 \ s_6 \ \dots \ s_n$$

$\uparrow \quad \uparrow \quad \uparrow$

$q_1 \quad q_3 \quad q_{20} \quad \textcircled{q}_9 \quad q_{17} \quad \textcircled{q}_9 \quad q_6 \quad \quad \quad q_{35} \quad q_{13}$



STEPS(PROOF BY CONTRADICTION)

- Assume A is regular.
- It should have a pumping length as 'p'.
- All strings longer than 'p' must be pumped $|s| \geq p$.
- Now find a string s in A, such that $|s| \geq p$.
- Divide s into xyz.
- Show that $xy^i z \notin A$, for some i.
- Consider all the ways that s can be divided into xyz.
- Show that none of these can satisfy all the three conditions at same time.
- S cannot be pumped.



B BE THE LANGUAGE $\{0^N 1^N \mid N \geq 0\}$. USE THE PUMPING LEMMA TO PROVE THAT B IS NOT REGULAR.

- Assume to the contrary that B is regular.
- Choose pumping length, p
- Choose string s from the language, $s = 0^p 1^p$, s is a member of B, $|s| \geq p$.
- Split s = xyz, where for any $i \geq 0$ the string $xy^i z$ is in B.
- three cases:
 - string y consists only 0's $\rightarrow xyyz$ has more 0's than 1's. violates the condition 1.
 - String y consists only 1's $\rightarrow xyyz$ has more 1's than 0's. violates the condition 1.
 - String y consist of both 0's and 1's \rightarrow **may have** equal number of 0's and 1's \rightarrow may be out of order hence not member.
- By contradiction, B is not regular.



$C = \{w \mid w \text{ HAS AN EQUAL NUMBER OF } 0\text{'S AND } 1\text{'S}\}.$

$$\begin{array}{c} 0^p 1^p \\ |x|y|z \\ |x|y| \leq p \end{array}$$

- Assume C is regular.
- Choose pumping length, p
- Choose string s from the language, $s = 0^p 1^p$, s is a member of C , $|s| \geq p$.
- Split $s = xyz$, where for any $i \geq 0$ the string $xy^i z$ is in C .
- If we pick $x = z = \varepsilon$ and $y = 0^p 1^p$, can we pump it and have pumped string $xy^i z \in C$?
- Yes! But this choice breaks condition $|xy| \leq p$.
- Suppose we choose x, y, z such that $|xy| \leq p$ and $|y| > 0$. Since $|xy| \leq p$, y consists only of 0's. Hence $xy^i z \notin C$ (too ε many zeros).



$$F = \{ WW \mid W \in \{0, 1\}^*\}.$$

|₀₁

- $F = \{\varepsilon, 00, 11, 0000, 0101, 1010, 1111, \dots\}$
- Assume F is regular.
- Choose pumping length, p
- Choose string s from the language, $s = 0^p 1 0^p 1$, s is a member of F , $|s| \geq p$.
- Split $s = xyz$, where for any $i \geq 0$ the string $xy^i z$ is in F .
- if we let x and z be the empty string, we can pump s . With condition 3 the proof follows because y must consist only of 0s, so $xyyz \notin F$, since 0's separated by 1 must be equal.



$$E = \{0^i 1^j \mid i > j\}$$

- Assume E is regular.
- Let p be the pumping length for E given by the pumping lemma.
- Let s = $0^{p+1}1^p$.
- Then s can be split into xyz, satisfying the conditions of Pumping Lemma.
- By condition 3, y consists only of 0s.
- Let's examine the string xyyz to see whether it can be in E. Adding an extra copy of y increases the number of 0s. But, E contains all strings in 0^*1^* that have more 0s than 1s, so increasing the number of 0s will still give a string in E. No contradiction occurs. We need to try something else.
- Let's consider the string $xy^0z = xz$, Removing string y decreases the number of 0s in s. Recall that s has just one more 0 than 1. Therefore, xz cannot have more 0s than 1s, so it cannot be a member of E. Thus we obtain a contradiction.



$$D = \{ 1^{n^2} \mid n \geq 0 \}$$

- D contains all strings of 1s whose length is a perfect square
- D = $\{\epsilon, 1, 1111, 11111111, \dots\}$ We use the pumping lemma to prove that D is not regular. The proof is by contradiction.
- Let p be the pumping length given by the pumping lemma.
- Let s be the string 1^{p^2}
- Assume we have $xyz \in D$ as per Pumping Lemma.
- By condition 3 of the pumping lemma, $|xy| \leq p$ and thus $|y| \leq p$.
- We have $|xyz| \leq p^2$, then $|xy^2z| \leq p^2 + p$. But $p^2 + p < p^2 + 2p + 1 = (p + 1)^2$. Moreover, condition 2 implies that $|y| > 0$, and so $|xy^2z| > p^2$.
- So the length of xy^2z lies between two consecutive perfect squares p^2 and $(p + 1)^2$, and hence $xy^2z \notin D$



$$A = \{a^n b^n \mid n \geq 0\}$$

- We need to keep the count of number of 'a' to get same number of 'b'
- FA cannot keep count, hence A is not regular.
- Proof:
 - Assume A is regular.
 - Pumping length, p
 - Choose string s, $s = a^p b^p$
 - Divide into 3 parts, xyz
 - Assume $p = 7$, $s = a^7 b^7 = aaaaaaaabbbbbbb$
 - Case 1: The y is in the a part
 - aaaaaaaa bbbbbbb
 - x y z
 - $xy^i z \rightarrow xy^2 z$, here $i = 2$
 - $aaaaaaaaaaaabbbbbbb \neq A$



- Case 2: The y is in the b part

- ~~aaaaaaaa bbbbbbb~~
x y z

- $xy^iz \rightarrow xy^2z$, here $i = 2$

- $aaaaaaaaabbbbbbbbbbb \neq A$

- Case 3: The y is in the a and b part

- ~~aaaaaaaa bbbbbbb~~
x y z

- $xy^iz \rightarrow xy^2z$, here $i = 2$

- $aaaaaaaaabbaabbbbbbb \neq A$

- $|xy| \leq p$, here $p = 7$

- Case 1: $6 < 7$

- Case 2: $13 \nless 7$

- Case 3: $9 \nless 7$



CONTEXT FREE GRAMMAR



OVERVIEW

- **Grammar**
 - Derivation from a grammar
 - Types
- **Context-Free Grammars**
 - Formal Definition
 - Examples of CFG
 - Designing CFG
 - Parse/derivation tree
 - Ambiguity
 - Simplification of CFG
 - Chomsky Normal Form
 - Greibach Normal form
- **Pushdown Automata**
 - Formal Definition
 - Examples of PDA
 - Equivalence with CFGs
- **Non-Context Free Languages**
 - The pumping lemma for context-free languages



GRAMMAR

- Regular grammar → set of rules that we use for proper composition.
- Grammar can be formally described using 4 tuples as $G = \{V, T, S, P\}$
 - $V \rightarrow$ set of variables/non terminal symbols
 - $T \rightarrow$ Terminal symbols
 - $S \rightarrow$ Start symbol
 - $P \rightarrow$ Production rule for terminal and non terminal symbols.
- Production rule, has the form $a \rightarrow \beta$, where a and β are strings on $V \cup T$ and atleast one symbol of ' a ' belongs to V .
- Example: $G = (\{S, A, B\}, \{a, b\}, S, \{S \rightarrow AB, A \rightarrow a, B \rightarrow b\})$



STEPS

- Write down the start variable. It is the variable on the left-hand side of the top rule, unless specified otherwise.
- Find a variable that is written down and a rule that starts with that variable. Replace the written down variable with the right-hand side of that rule.
- Repeat step 2 until no variables remain.



EXAMPLE

$A \rightarrow 0A1$

$A \rightarrow B$

$B \rightarrow \#$

- Grammar G1 generates the string 000#111. The sequence of substitutions to obtain a string is called a derivation.
- A derivation of string 000#111 in grammar G1 is



EXAMPLE

$A \rightarrow 0A1$

$A \rightarrow B$

$B \rightarrow \#$

- Grammar G1 generates the string 000#111. The sequence of substitutions to obtain a string is called a derivation.
- A derivation of string 000#111 in grammar G1 is
- $A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$.



$$G = (\{S, A, B\}, \{A, B\}, S, \{S \rightarrow AB, A \rightarrow a, B \rightarrow b\})$$

- $V = \{S, A, B\}$
- $T = \{a, b\}$
- $S = s$
- $P = S \rightarrow AB, A \rightarrow a, B \rightarrow b$
 - $S \rightarrow AB$, we know that $A \rightarrow a$
 - $S \rightarrow aB$
 - $S \rightarrow ab$



DERIVATION FROM A GRAMMAR

- The set of all the strings that can be derived from a grammar is said to be the Language generated from that grammar.
- Example: $G = (\{S, A\}, \{a, b\}, S, \{S \rightarrow aAb, A \rightarrow aaAbb, A \rightarrow \epsilon\})$
 - $S \rightarrow aAb$, $A \rightarrow aaAbb$
 - Expanding, $S \rightarrow aaaAbb$
 - $S \rightarrow aaaAbb$, $A \rightarrow \epsilon$
 - $S \rightarrow aaabb$, String derived from grammar G
 - $L = a^n b^n$



DERIVATION FROM A GRAMMAR

- The set of all the strings that can be derived from a grammar is said to be the Language generated from that grammar.
- Example: $G = (\{S,A\}, \{a,b\}, S, \{S \rightarrow aAb, A \rightarrow aaAbb, A \rightarrow \epsilon\})$



TYPES OF GRAMMAR

Grammar Type	Grammar Accepted	Language Accepted	Automaton
Type 0	Unrestricted grammar	Recursively enumerable language	Turing Machine
Type 1	Context-sensitive grammar	Context-sensitive language	Linear-bounded automaton
Type 2	Context-free grammar	Context-free language	Pushdown Automaton
Type 3	Regular grammar	Regular language	Finite state automaton



CONTEXT FREE GRAMMAR

- Methods of describing languages:
 - Finite automata
 - Regular expressions
- Some simple languages, such as $\{0^n 1^n \mid n \geq 0\}$, cannot be described by these languages.
- We study Context-Free Grammars, a more powerful method of describing languages. Such grammars can describe certain features that have a recursive structure.
- The collection of languages associated with context-free grammars are called the Context-Free languages
- PDA: a class of machines recognizing the context-free languages A Context Free Grammar is a “machine” that creates a language. A language created by a CFG is called A Context Free Language



CONTEXT FREE LANGUAGE

Consider language $\{ 0^n 1^n \mid n \geq 0 \}$, which is nonregular.

- Start variable S with “substitution rules”: $S \rightarrow 0S1 \quad S \rightarrow \epsilon$
- Rules can yield string $0^k 1^k$ by applying rule “ $S \rightarrow 0S1$ ” k times, followed by rule “ $S \rightarrow \epsilon$ ” once.
- Derivation of string $0^3 1^3$ $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000S111 \Rightarrow 000\epsilon111 = 000111$



EXAMPLE

- Example of CFG Example: Language $\{ 0^n 1^n \mid n \geq 0 \}$ has CFG $G = (V, \Sigma, R, S)$
- Variables $V = \{S\}$
- Terminals $\Sigma = \{0, 1\}$
- Start variable S
- Rules $R: S \rightarrow 0S1 \quad S \rightarrow \epsilon$
- Combine rules with same left-hand side in [Backus-Naur \(or Backus Normal\) Form \(BNF\)](#):

$$S \rightarrow 0S1 \mid \epsilon$$



The following is an example of CFG, which we call G1.

$A \rightarrow 0A1$

$A \rightarrow B$

$B \rightarrow \#$

- A grammar consists of a **collection of rules**. Each rule appears as a line in the grammar, comprising a symbol and a string separated by an arrow.
- The **symbol** is called a **variable**.
- The string consists of variables and other symbols called **terminals**.
- One variable is designated as the **start variable**.
- It usually occurs on the LHS of the topmost rule.
- For example, grammar G1 contains 3 rules. G1's variables are A and B, where A is the start variable. Its terminals are 0, 1, and #.



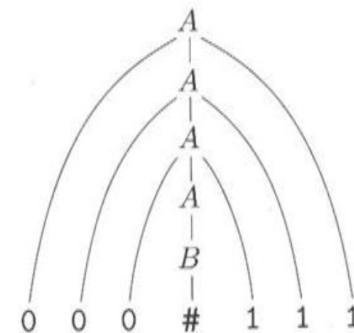
$$\begin{array}{l}
 A \rightarrow 0A1 \\
 A \rightarrow B \\
 B \rightarrow #
 \end{array}
 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \quad A \rightarrow 0A1 \mid B$$

For example, grammar G1 generates the string 000#111. The sequence of substitutions to obtain a string is called a **derivation**.

A derivation of string 000#111 in grammar G1 is

$$A \rightarrow 0 A 1 \rightarrow 0 0 A 1 1 \rightarrow 0 0 0 A 1 1 1 \rightarrow 0 0 0 B 1 1 1 \rightarrow 0 0 0 \# 1 1 1.$$

Language defined by grammar G1 $\{0^n \# 1^n \mid n \geq 0\}$



Parse tree for 000#111 in grammar G1



EXAMPLE

- Consider grammar $G_3 = (\{S\}, \{a, b\}, R, S)$. The set of rules, R , is
 - $S \rightarrow aSb \mid SS \mid \epsilon$
- What are some of the strings generated by this grammar?
- Can you verify if the strings abab, aaabbb, and aababb, belong to the grammar.



EXAMPLE

- CFG $G = (V, \Sigma, R, S)$ with

1. $V = \{S\}$

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

3. Rules $R: S \rightarrow 0S \mid \epsilon$

Then $L(G) = \{ 0^n \mid n \geq 0 \}$.

For example, S derives 0^3

$S \Rightarrow 0S \Rightarrow 00S \Rightarrow 000S \Rightarrow 000\epsilon = 000$

Note that \rightarrow and \Rightarrow are different.

\rightarrow used in defining rules

\Rightarrow used in derivation



EXAMPLE

- CFG $G = (V, \Sigma, R, S)$ with

1. $V = \{S\}$
2. $\Sigma = \{0, 1\}$
3. Rules $R: S \rightarrow 0S \mid 1S \mid \epsilon$

Then $L(G) = \Sigma^*$.

For example, S derives 0100

$$S \Rightarrow 0S \Rightarrow 01S \Rightarrow 010S \Rightarrow 0100S \Rightarrow 0100$$



EXAMPLE

Consider grammar G4 = (V, Σ ,R, <EXPR>)

V is {<EXPR> , <TERM> , <FACTOR>} and Σ is {a, +, x, (,)}.

The rules are

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle$$

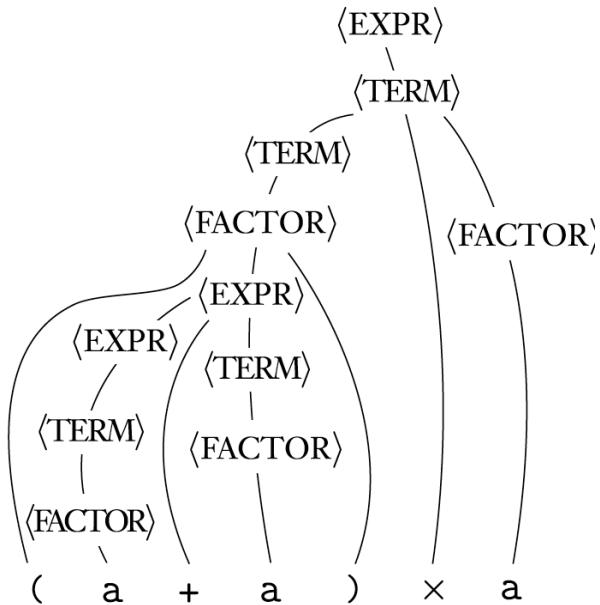
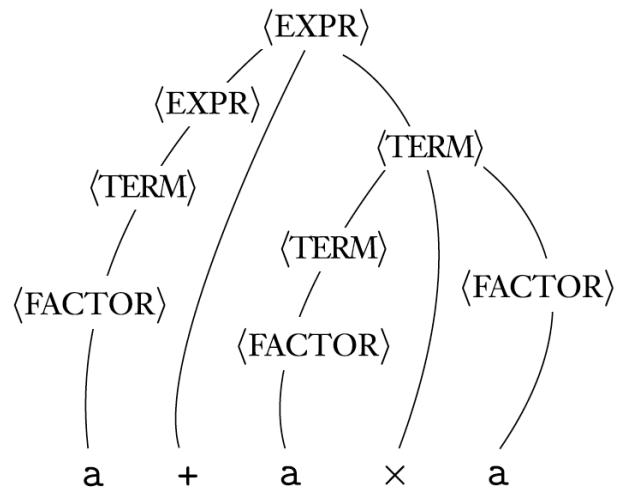
$$\langle \text{TERM} \rangle \rightarrow \langle \text{TERM} \rangle * \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle$$

$$\langle \text{FACTOR} \rangle \rightarrow (\langle \text{EXPR} \rangle) \mid a$$

The two strings a + a * a and (a + a) * a can be generated with grammar G4.



PARSE TREE



Consider grammar G4 = (V, Σ ,R, <EXPR>)

V is {<EXPR> , <TERM> , <FACTOR>} and Σ is {a, +, x, (,)}.

The rules are

<EXPR> → <EXPR > + <TERM> | <TERM>

<TERM> → <TERM> * <FACTOR> | <FACTOR>

<FACTOR> → (<EXPR >) | a

<EXPR> → <EXPR> + <TERM>

→ <TERM> + <TERM>

→ <FACTOR> + <TERM>

→ a + <TERM>

→ a + <FACTOR>

→ a + (<EXPR>)

→ a + (<TERM>)

→ a + (<TERM> * <FACTOR>)

→ a + (<FACTOR> * <FACTOR>)

→ a + (a * <FACTOR>) = a + (a * a)



TRY YOURSELF!

- CFG $G = (V, \Sigma, R, S)$ with $V = \{S\}$, $\Sigma = \{0, 1\}$, Rules $R: S \rightarrow 0S \mid 1S \mid 1$
- CFG $G = (V, \Sigma, R, S)$ with $V = \{S, Z\}$, $\Sigma = \{0, 1\}$, Rules $R: S \rightarrow 0S1 \mid Z Z \rightarrow 0Z \mid \epsilon$
- PALINDROME = { $w \in \Sigma^* \mid w = wR$ }, where $\Sigma = \{a, b\}$. CFG $G = (V, \Sigma, R, S)$ with $V = \{S\}$, $\Sigma = \{a, b\}$, Rules $R: S \rightarrow aSa \mid bSb \mid a \mid b \mid \epsilon$



TRY YOURSELF!

- CFG $G = (V, \Sigma, R, S)$ with $V = \{S\}$, $\Sigma = \{0, 1\}$, Rules $R: S \rightarrow 0S \mid 1S \mid 1$.
Derive 011
- CFG $G = (V, \Sigma, R, S)$ with $V = \{S, Z\}$, $\Sigma = \{0, 1\}$, Rules $R: S \rightarrow 0S1 \mid Z$
 $Z \rightarrow 0Z \mid \epsilon$. $L(G) = \{0^i1^j \mid i \geq j\}$
- PALINDROME = { $w \in \Sigma^*$ | $w = wR$ }, where $\Sigma = \{a, b\}$. CFG $G = (V, \Sigma, R, S)$ with $V = \{S\}$, $\Sigma = \{a, b\}$, Rules $R: S \rightarrow aSa \mid bSb \mid a \mid b \mid \epsilon$
- $A \rightarrow 0A, A \rightarrow 1B \mid 1, B \rightarrow 0A \mid 1B \mid 1$, Derive the string 10011011
- $S \rightarrow 0S1 \mid 1S0 \mid SS \mid \epsilon$, derive the string 00011011.



CFG FOR SIMPLE ARITHMETIC EXPRESSIONS

- CFG $G = (V, \Sigma, R, S)$ with
 1. $V = \{S\}$
 2. $\Sigma = \{ +, -, \times, /, (,), 0, 1, 2, \dots, 9 \}$
 3. Rules $R: S \rightarrow S + S \mid S - S \mid S \times S \mid S/S \mid (S) \mid -S \mid 0 \mid 1 \mid \dots \mid 9$

$L(G)$ is a set of valid arithmetic expressions over single-digit integers.

S derives string $2 \times (3 + 4)$



CFG FOR SIMPLE ARITHMETIC EXPRESSIONS

- CFG $G = (V, \Sigma, R, S)$ with
 1. $V = \{S\}$
 2. $\Sigma = \{ +, -, \times, /, (,), 0, 1, 2, \dots, 9 \}$
 3. Rules $R: S \rightarrow S + S \mid S - S \mid S \times S \mid S/S \mid (S) \mid -S \mid 0 \mid 1 \mid \dots \mid 9$

$L(G)$ is a set of valid arithmetic expressions over single-digit integers.

S derives string $2 \times (3 + 4)$

$S \Rightarrow S \times S \Rightarrow S \times (S) \Rightarrow S \times (S + S) \Rightarrow 2 \times (S + S) \Rightarrow 2 \times (3 + S) \Rightarrow 2 \times (3 + 4)$



DERIVATION TREE

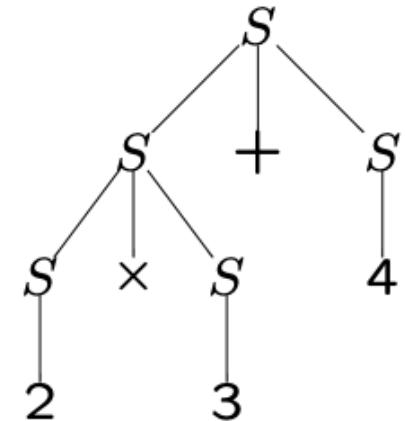
- CFG

$S \rightarrow S + S \mid S - S \mid S \times S \mid S/S \mid (S) \mid -S \mid 0 \mid 1 \mid \dots \mid 9$

- Can generate string $2 \times 3+4$ using derivation

$S \Rightarrow S + S \Rightarrow S \times S + S \Rightarrow 2 \times S + S \Rightarrow 2 \times 3 + S \Rightarrow 2 \times 3+4$

- **Leftmost derivation:** leftmost variable replaced in each step.
- Corresponding **derivation (or parse)** tree



DESIGNING CFG

For example, to get a grammar for the language $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$, first construct the grammar

$S_1 \rightarrow 0 S_1 1 \mid \epsilon$

for the language $\{0^n 1^n \mid n \geq 0\}$ and the grammar

$S_2 \rightarrow 1 S_2 0 \mid \epsilon$

for the language $\{1^n 0^n \mid n \geq 0\}$ and then add the rule

$S \rightarrow S_1 \mid S_2$ to give the grammar

$S \rightarrow S_1 \mid S_2$

$S_1 \rightarrow 0 S_1 1 \mid \epsilon$

$S_2 \rightarrow 1 S_2 0 \mid \epsilon$



EXAMPLE

- Generate CFG for given language $L=\{a^n b^n / n \geq 0\}$
- Let $L=\{a^n b^{2n} / n \geq 1\}$. Find CFG
- $L=\{a^n b a^n / n \geq 1\}$. Find CFG.
- $L=\{wcw^r / w \text{ is element of } (a,b)^*\}$. Find the CFG.



CONTEXT FREE LANGUAGE (CFL)

- The CFG is given as $S \rightarrow aSb/ab$. Find the CFL.



CONTEXT FREE LANGUAGE (CFL)

- The CFG is given as $S \rightarrow aSb/ab$. Find the CFL.
- Let $G = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, \{S\})$
- $S \rightarrow ab$
- $S \rightarrow aSb$
 - $S \rightarrow aabb$
- $S \rightarrow aSb$
 - $S \rightarrow aaSbb$
 - $S \rightarrow aaabbb$
- $L(G) = \{ab, aabb, aaabbb, \dots\}$ Therefore, $L(G) = \{a^n b^n / n \geq 1\}$



EXAMPLE

- Find the CFL for $S \rightarrow aB/bA, A \rightarrow a/aS/bAA, B \rightarrow b/bS/aBB$



TRY YOURSELF!

- Find $L(G)$ $S \rightarrow aSa/bSb/\varepsilon$
- Find $L(G)$ $S \rightarrow aSa/bSa/\varepsilon$
- Find $L(G)$ $S \rightarrow aSa/bSb/a/b$
- Find $L(G)$ $S \rightarrow SS/bS/a$
- Find $L(G)$ for $S \rightarrow aS/bS/a$
- Find $L(G)$ $S \rightarrow aSa/bSb/\varepsilon$



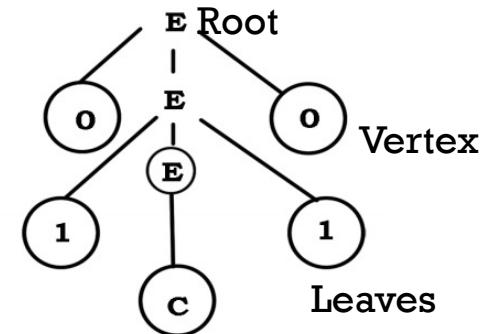
TRY YOURSELF!

- Consider the grammar $G=(\{A,S\}, \{a,b\}, P, S)$ where P consist of $S \rightarrow aAS/a$ $A \rightarrow SbA/SS/ba$ Draw the derivation for string "aabbaa"
- Draw the derivation for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$



PARSE OR DERIVATION TREE

- Parse tree is an ordered rooted tree that graphically presents the semantic information of string derived from CFG.
- Any one production or derivation derived from tree format.
- It is the graph form of representation.
- RULE
 - Start with 'S'. → Root
 - The final answer should be terminal. → Leaves
 - The derivation should be applied from left to right.
 - The intermediate derivation should be terminal or nonterminal. → vertex



- $W = 01C10$, $E \rightarrow 0E0$ $E \rightarrow 1E1$ $E \rightarrow C$

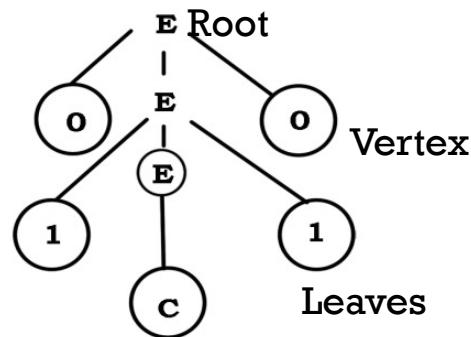
- DERIVATION $E \Rightarrow 0E0$

$\Rightarrow 01E10 :: E \Rightarrow 1E1$

$\Rightarrow 01C10 :: E \Rightarrow C$

$\Rightarrow E \Rightarrow 01C10$

- PARSE TREE (OR) DERIVATION TREE



TYPES OF DERIVATION

- LEFT MOST DERIVATION (LMD): If at each step-in derivation, a production is applied to the left most variable (or) left most non-terminal then the derivation method is called left most derivation.

- Example $w = id + id * id$ $E \rightarrow E+E$, $E \rightarrow E * E$, $E \rightarrow id$

$E \rightarrow E+E$

$E \rightarrow id+E \therefore E=>id$

$E \rightarrow id+E * E \therefore E=>E * E$

$E \rightarrow id+id * E \therefore E=>id$

$E \rightarrow id+id * id$



- **RIGHT MOST DERIVATION (RMD):** A derivation in which the right most variable is replaced at each step then, the derivation method is called right most derivation.

- Example

$E \rightarrow E+E$

$E \rightarrow E+E^*E \therefore E \Rightarrow E^*E$

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

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



EXAMPLE

- Draw derivation for the string abbabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/Baa$ $B \rightarrow b/bS/Abb$
 - $s \rightarrow aB$
 - $s \rightarrow aBB$
 - $s \rightarrow abSB$
 - $s \rightarrow abSbs$
 - $s \rightarrow abbAbs$
 - $s \rightarrow abbAbbA$
 - $s \rightarrow abbabbA$
 - $s \rightarrow abbabba$



EXAMPLE

Draw derivation tree for the string abaaba for the CFG given by G where p is $S \rightarrow aSa/bSb/b/a/\epsilon$

Solution:

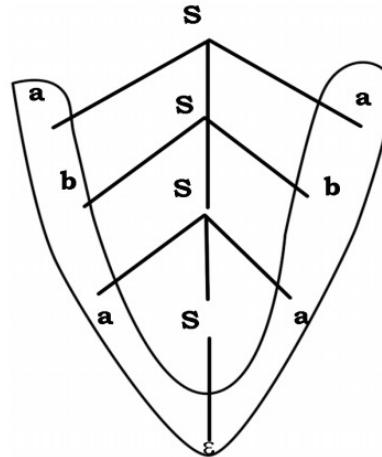
$$S \rightarrow aSa$$

$$S \rightarrow abSba$$

S → abaSaba

$$S \rightarrow aba\varepsilon aba$$

$S \rightarrow abaaba$



TRY YOURSELF!

- Draw derivation tree for the string aabbaabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/bAA$ $B \rightarrow b/bS/aBB$
- Draw the derivation tree for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$.
- Draw the derivation tree for the given graph CFG $G=(v,t,p,s)$ where $p=\{S \rightarrow aSb/\epsilon\}$ and the input string :aabb



EXAMPLE

- Draw derivation for the string abbabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/Baa$ $B \rightarrow b/bS/Abb/BB$
 - $s \rightarrow aB$
 - $s \rightarrow aBB$
 - $s \rightarrow abSB$
 - $s \rightarrow abSbs$
 - $s \rightarrow abbAbS$
 - $s \rightarrow abbAbbA$
 - $s \rightarrow abbabbA$
 - $s \rightarrow abbabba$



TRY YOURSELF!

- Consider the grammar $G=(\{A,S\}, \{a,b\}, P, S)$ where P consist of $S \rightarrow aAS/a$ $A \rightarrow SbA/SS/ba$ Draw the derivation for string "aabbaa"
- Draw the derivation for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$



TYPES OF DERIVATION

- LEFT MOST DERIVATION (LMD): If at each step-in derivation, a production is applied to the left most variable (or) left most non-terminal then the derivation method is called left most derivation.
- Example $w = id + id * id$ $E \rightarrow E+E$, $E \rightarrow E * E$, $E \rightarrow id$

$E \rightarrow E+E$

$E \rightarrow id+E \therefore E=>id$

$E \rightarrow id+E * E \therefore E=>E * E$

$E \rightarrow id+id * E \therefore E=>id$

$E \rightarrow id+id * id$



- **RIGHT MOST DERIVATION (RMD):** A derivation in which the right most variable is replaced at each step then, the derivation method is called right most derivation.

- Example

$E \rightarrow E+E$

$E \rightarrow E+E^*E \therefore E \Rightarrow E^*E$

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

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



TRY YOURSELF!

- Draw derivation tree for the string aabbaabba For thee CFG given G where production is $S \rightarrow bA/aB$ $A \rightarrow a/aS/bAA$ $B \rightarrow b/bS/aBB$
- Draw the derivation tree for the string abb where $S \rightarrow aAB$ $A \rightarrow bBb$ $B \rightarrow A/\epsilon$.
- Draw the derivation tree for the given graph CFG $G=(v,t,p,s)$ where $p=\{S \rightarrow aSb/\epsilon\}$ and the input string :aabb



AMBIGUITY

- CFG G is ambiguous if string $w \in L(G)$ having different parse trees (or equivalently, different leftmost derivations).
- A string is derived ambiguously in a CFG if it has two or more different leftmost derivations
- A grammar that produces **more than one parse tree (or) derivation tree** for some string then the grammar is said to be an ambiguous grammar.
- An ambiguous grammar produces more than one LMD (or) more than 1 RMD then, the given grammar is said to be an ambiguous grammar.
- Leftmost derivation : At every step in the derivation the leftmost variable is replaced
- A grammar is ambiguous if it generates *some* string ambiguously
- Some context free languages are inherently ambiguous, that is, every grammar for the language is ambiguous



- LMD:

$E \Rightarrow id + E :: E \Rightarrow id$

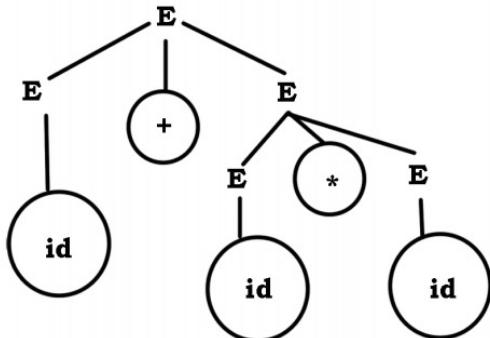
$E \Rightarrow id + E * E$

$E \Rightarrow id + id * E$

$E \Rightarrow id + id * id$

$E \Rightarrow id + id * id$

PARSE TREE



Rules:

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

w = id + id * id

- RMD:

$E \Rightarrow E * E$

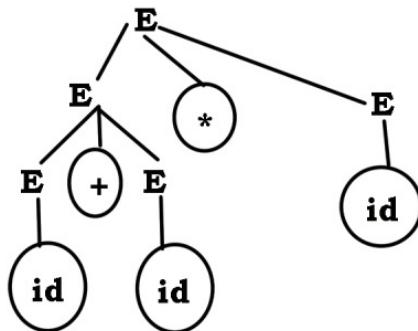
$E \Rightarrow E + E * E$

$E \Rightarrow E + E * id$

$E \Rightarrow E + id * id$

$E \Rightarrow id + id * id$

Therefore, the above grammar is ambiguous.



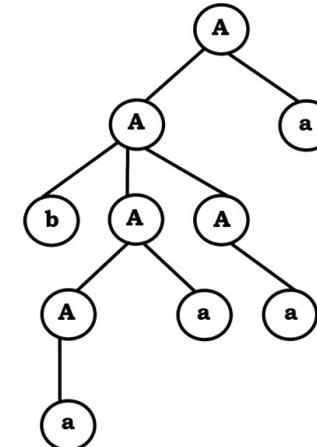
Show that CFG having productions $A \rightarrow a/Aa/bAA/AAb/AbA$ is ambiguous. The string can be baaaa



Show that CFG having productions $A \rightarrow a/Aa/bAA/AAb/AbA$ is ambiguous.

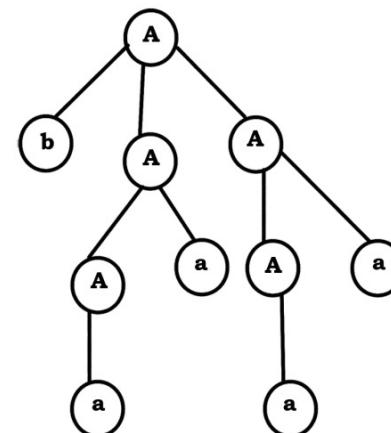
Parse tree 1 $A \rightarrow Aa \Rightarrow bAAa \Rightarrow bAaaa$

$A \rightarrow baaaa$



▪ Parse tree 2 $A \rightarrow bAA \Rightarrow bAaA \Rightarrow bAaAa$

$A \rightarrow baaaa$



▪ Therefore, the above grammar is ambiguous.



TRY YOURSELF!

- Prove that given CFG is ambiguous, $S \rightarrow 0B/1A$ $A \rightarrow 0/0S/1AA$ $B \rightarrow 1/1S/0BB$. The string can be randomly chosen. For reference I have chosen “0011010”
- Show that grammar is $S \Rightarrow aSbS/ bSaS/ \epsilon$ ambiguous.



RESOLVING AMBIGUITY

- Designing unambiguous grammars is tricky and requires planning from the start.
- It's hard to start with an ambiguous grammar and to manually convert it into an unambiguous one.
- Often, must throw the whole thing out and start over.
- We have just seen that this grammar is ambiguous:
$$E \rightarrow E+E|E^*E|id$$
If we take a string $id+id^*id$ or $id+id+id$ we get two parse trees.
- Goals:
 - Eliminate the ambiguity from the grammar.
 - Make the only parse trees for the grammar the ones corresponding to operator precedence.

