

# 15-453

## FORMAL LANGUAGES, AUTOMATA AND COMPUTABILITY

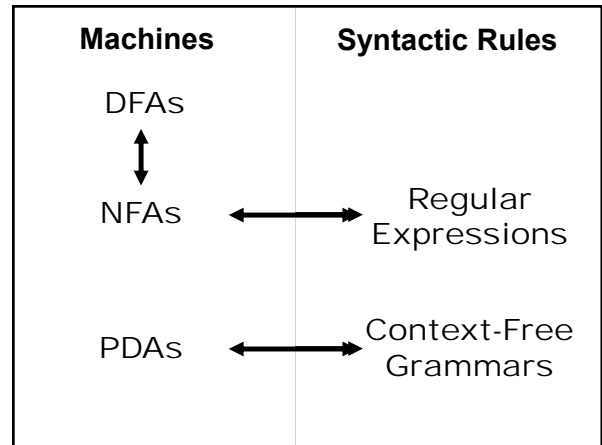
REVIEW for MIDTERM 1  
THURSDAY Feb 6

**Midterm 1 will cover everything we  
have seen so far**

The PROBLEMS will be from Sipser,  
Chapters 1, 2, 3

**It will be Closed-Book,  
Closed-Everything**

- 1. Deterministic Finite Automata and Regular Languages
- 2. Non-Deterministic Finite Automata
- 3. Pumping Lemma for Regular Languages; Regular Expressions
- 4. Minimizing DFAs
- 5. PDAs, CFGs; Pumping Lemma for CFLs
- 6. Equivalence of PDAs and CFGs
- 7. Chomsky Normal Form
- 8. Turing Machines



### THE REGULAR OPERATIONS

Union:  $A \cup B = \{ w \mid w \in A \text{ or } w \in B \}$

Intersection:  $A \cap B = \{ w \mid w \in A \text{ and } w \in B \}$

Negation:  $\neg A = \{ w \in \Sigma^* \mid w \notin A \}$

Reverse:  $A^R = \{ w_1 \dots w_k \mid w_k \dots w_1 \in A \}$

Concatenation:  $A \cdot B = \{ vw \mid v \in A \text{ and } w \in B \}$

Star:  $A^* = \{ s_1 \dots s_k \mid k \geq 0 \text{ and each } s_i \in A \}$

### REGULAR EXPRESSIONS

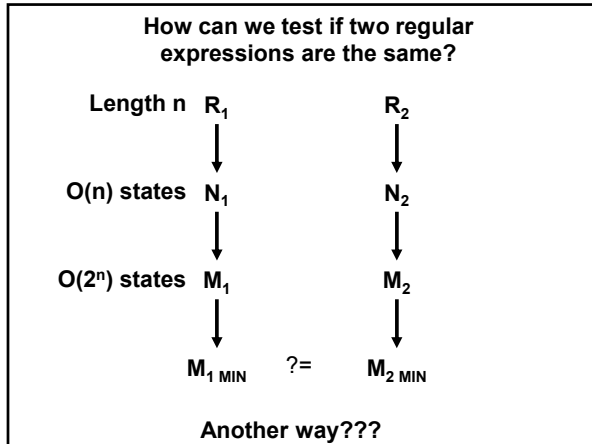
$\sigma$  is a regexp representing  $\{\sigma\}$

$\epsilon$  is a regexp representing  $\{\epsilon\}$

$\emptyset$  is a regexp representing  $\emptyset$

If  $R_1$  and  $R_2$  are regular expressions representing  $L_1$  and  $L_2$  then:

- $(R_1 R_2)$  represents  $L_1 \cdot L_2$
- $(R_1 \cup R_2)$  represents  $L_1 \cup L_2$
- $(R_1)^*$  represents  $L_1^*$



THEOREMS  
and  
CONSTRUCTIONS

THE PUMPING LEMMA  
(for Regular Languages)

Let  $L$  be a regular language with  $|L| = \infty$

Then there is an integer  $P$  such that  
if  $w \in L$  and  $|w| \geq P$   
then can write  $w = xyz$ , where:

- $|y| > 0$
- $|xy| \leq P$
- $xy^iz \in L$  for any  $i \geq 0$

THE PUMPING LEMMA  
(for Context Free Grammars)

Let  $L$  be a context-free language with  $|L| = \infty$

Then there is an integer  $P$  such that  
if  $w \in L$  and  $|w| \geq P$   
then can write  $w = uvxyz$ , where:

- $|vy| > 0$
- $|vxy| \leq P$
- $uv^ixy^iz \in L$ ,  
for any  $i \geq 0$

CONVERTING NFAs TO DFAs

Input: NFA  $N = (Q, \Sigma, \delta, Q_0, F)$

Output: DFA  $M = (Q', \Sigma, \delta', q_0', F')$

$Q' = 2^Q$

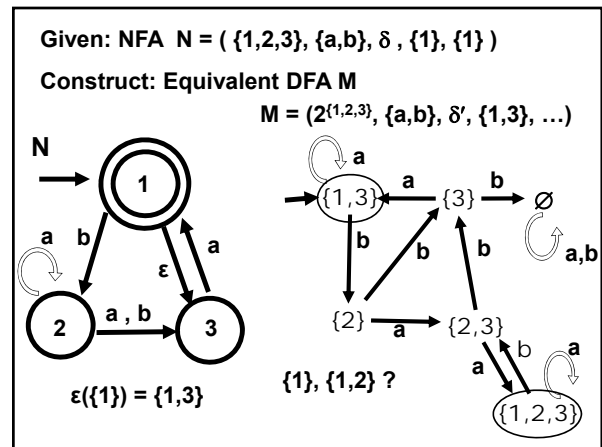
$\delta' : Q' \times \Sigma \rightarrow Q'$

$\delta'(R, \sigma) = \bigcup_{r \in R} \epsilon(\delta(r, \sigma))$

$q_0' = \epsilon(Q_0)$

$F' = \{ R \in Q' \mid f \in R \text{ for some } f \in F \}$

\* For  $R \subseteq Q$ , the  $\epsilon$ -closure of  $R$ ,  $\epsilon(R) = \{q \text{ that can be reached from some } r \in R \text{ by traveling along zero or more } \epsilon \text{ arrows}\}$



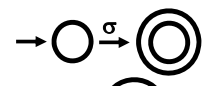
EQUIVALENCE

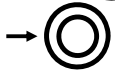
L can be represented by a regexp  
 $\Leftrightarrow$   
 L is a regular language

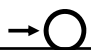
✓ L can be represented by a regexp  
 $\Rightarrow$   
 L is a regular language

Induction on the length of R:

Base Cases (R has length 1):

$R = \sigma \rightarrow$  

$R = \epsilon \rightarrow$  

$R = \emptyset \rightarrow$  

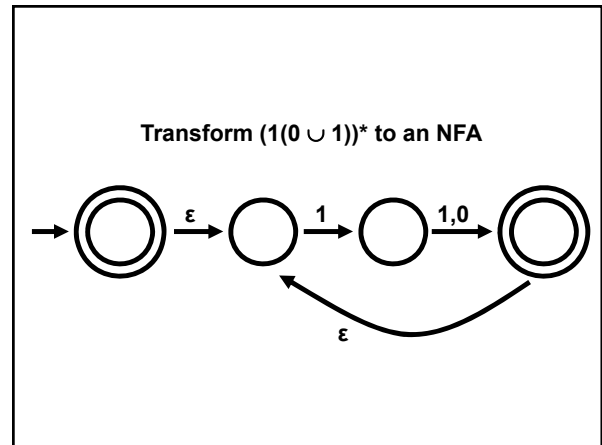
Inductive Step:

Assume R has length  $k > 1$ ,  
 and that every regexp of length  $< k$   
 represents a regular language

Three possibilities for what R can be:

$R = R_1 \cup R_2$  (Closure under Union)  
 $R = R_1 R_2$  (Closure under Concat.)  
 $R = (R_1)^*$  (Closure under Star)


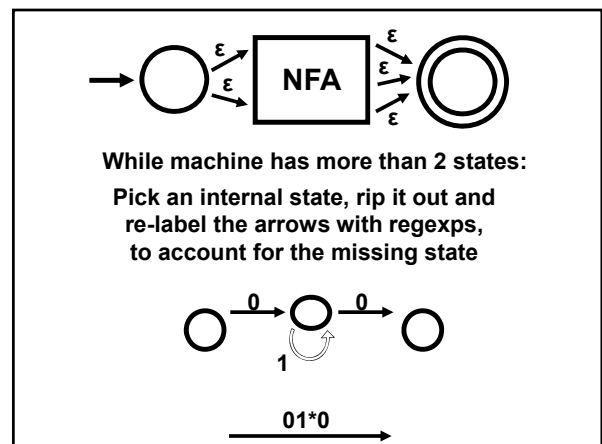
Therefore: L can be represented by a regexp  
 $\Rightarrow$  L is regular

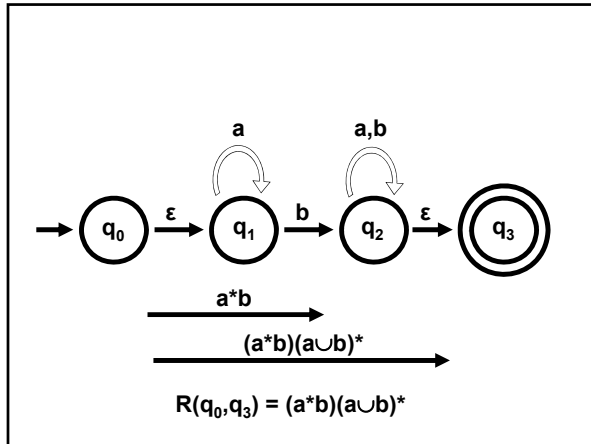


✓ L is a regular language  $\Rightarrow$   
 L can be represented by a regexp

Proof idea: Transform an NFA for L into a  
 regular expression by removing states and re-  
 labeling the arrows with regular expressions

Add unique and distinct start and accept states



**THEOREM**

For every regular language  $L$ , there exists a **UNIQUE** (up to re-labeling of the states) minimal DFA  $M$  such that  $L = L(M)$

**EXTENDING  $\delta$**

Given DFA  $M = (Q, \Sigma, \delta, q_0, F)$ , extend  $\delta$  to  $\hat{\delta} : Q \times \Sigma^* \rightarrow Q$  as follows:

- $\hat{\delta}(q, \epsilon) = q$
- $\hat{\delta}(q, \sigma) = \delta(q, \sigma)$
- $\hat{\delta}(q, w_1 \dots w_{k+1}) = \delta(\hat{\delta}(q, w_1 \dots w_k), w_{k+1})$

Note:  $\delta(q_0, w) \in F \Leftrightarrow M$  accepts  $w$

String  $w \in \Sigma^*$  distinguishes states  $q_1$  and  $q_2$  iff exactly ONE of  $\hat{\delta}(q_1, w), \hat{\delta}(q_2, w)$  is a final state

Fix  $M = (Q, \Sigma, \delta, q_0, F)$  and let  $p, q, r \in Q$

**Definition:**

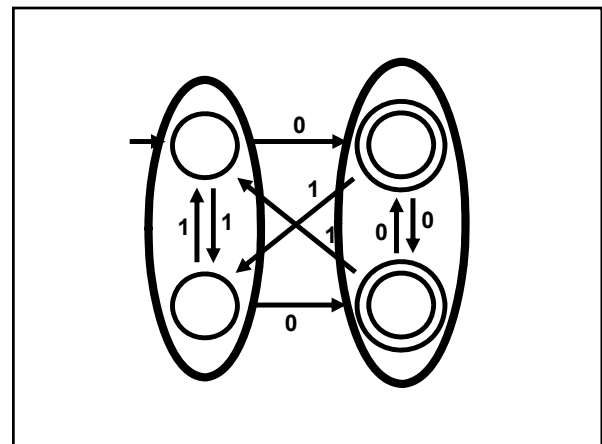
- $p \sim q$  iff  $p$  is indistinguishable from  $q$
- $p \not\sim q$  iff  $p$  is distinguishable from  $q$

**Proposition:**  $\sim$  is an equivalence relation

- $p \sim p$  (reflexive)
- $p \sim q \Rightarrow q \sim p$  (symmetric)
- $p \sim q$  and  $q \sim r \Rightarrow p \sim r$  (transitive)

**Proposition:**  $\sim$  is an equivalence relation  
so  $\sim$  partitions the set of states of  $M$  into disjoint equivalence classes

$[q] = \{ p \mid p \sim q \}$





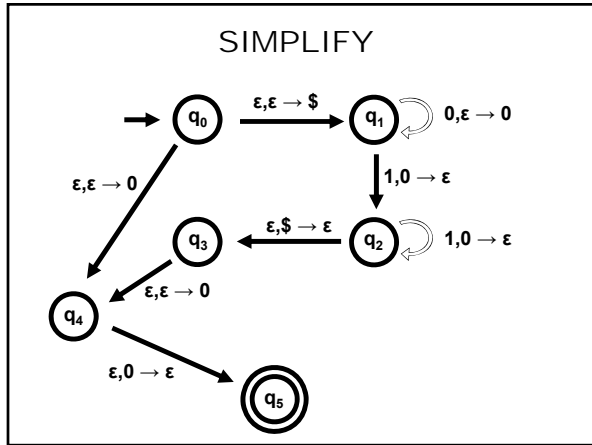
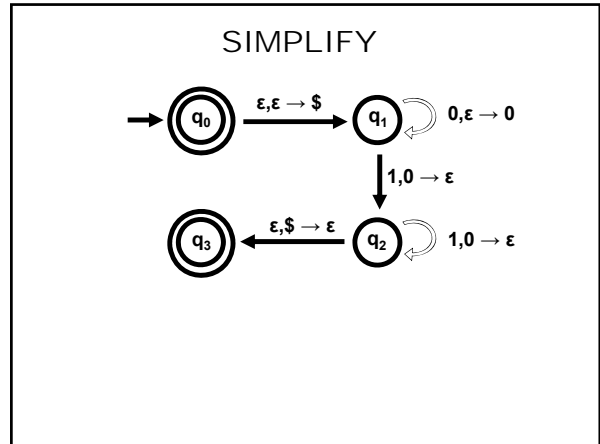
A Language  $L$  is generated by a CFG  
 $\Leftarrow$   
 $L$  is recognized by a PDA

Given PDA  $P = (Q, \Sigma, \Gamma, \delta, q, F)$

Construct a CFG  $G = (V, \Sigma, R, S)$  such that  $L(G) = L(P)$

First, simplify  $P$  to have the following form:

- (1) It has a single accept state,  $q_{\text{accept}}$
- (2) It empties the stack before accepting
- (3) Each transition either pushes a symbol or pops a symbol, but not both at the same time

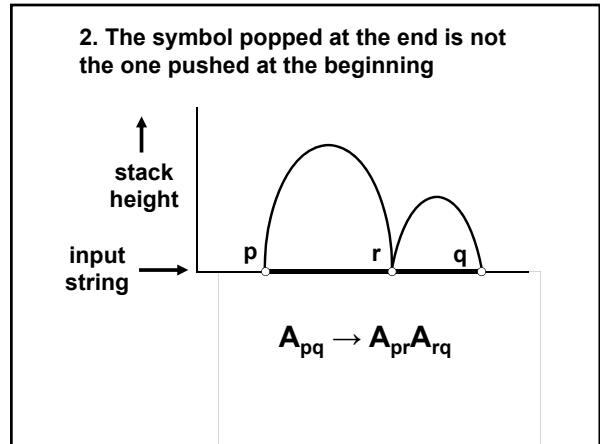
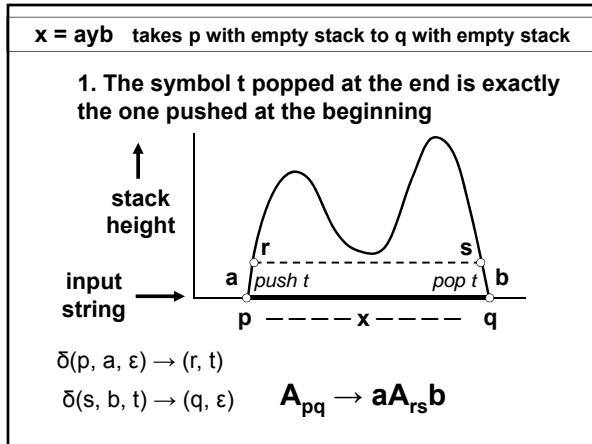


Idea For Our Grammar  $G$ :  
 For every pair of states  $p$  and  $q$  in PDA  $P$ ,

$G$  will have a variable  $A_{pq}$  which generates all strings  $x$  that can take:

$P$  from  $p$  with an empty stack to  $q$  with an empty stack

$$V = \{A_{pq} \mid p, q \in Q\}$$

$$S = A_{q_0 q_{\text{acc}}}$$


**Formally:**  
 $V = \{A_{pq} \mid p, q \in Q\}$   
 $S = A_{q_0 q_{acc}}$

For every  $p, q, r, s \in Q, t \in \Gamma$  and  $a, b \in \Sigma_\epsilon$   
 If  $(r, t) \in \delta(p, a, \epsilon)$  and  $(q, \epsilon) \in \delta(s, b, t)$   
 Then add the rule  $A_{pq} \rightarrow aA_{rs}b$


For every  $p, q, r \in Q$ ,  
 add the rule  $A_{pq} \rightarrow A_{pr}A_{rq}$

For every  $p \in Q$ ,  
 add the rule  $A_{pp} \rightarrow \epsilon$


**THE CHOMSKY NORMAL FORM**

A context-free grammar is in Chomsky normal form if every rule is of the form:

$A \rightarrow BC$  B, C are variables (not the start var)  
 $A \rightarrow a$  a is a terminal  
 $S \rightarrow \epsilon$  S is the start variable

  $S \rightarrow 0S1$   
 $S \rightarrow TT$   
 $T \rightarrow \epsilon$

$S_0 \rightarrow TU \mid \epsilon$   
 $T \rightarrow 0$   
 $U \rightarrow SV \mid 1$   
 $S \rightarrow TU$   
 $V \rightarrow 1$



**Theorem:** If G is in CNF,  $w \in L(G)$  and  $|w| > 0$ , then any derivation of w in G has length  $2|w| - 1$

**Theorem:** Any context-free language can be generated by a context-free grammar in Chomsky normal form

**“Can transform any CFG into Chomsky normal form”**

**Theorem:** Any CFL can be generated by a CFG in Chomsky normal form

**Algorithm:**

1. Add a new start variable ( $S_0 \rightarrow S$ )
2. Eliminate all  $A \rightarrow \epsilon$  rules:  
 For each occurrence of A on the RHS of a rule, add a new rule that removes that occurrence (unless this new rule was previously removed)
3. Eliminate all  $A \rightarrow B$  rules:  
 For each rule with B on LHS of a rule, add a new rule that puts A on the LHS instead (unless this new rule was previously removed)
4. Convert  $A \rightarrow u_1 u_2 \dots u_k$  to  $A \rightarrow u_1 A_1, A_1 \rightarrow u_2 A_2, \dots$   
 If  $u_i$  is a terminal, replace  $u_i$  with  $U_i$  and add  $U_i \rightarrow u_i$

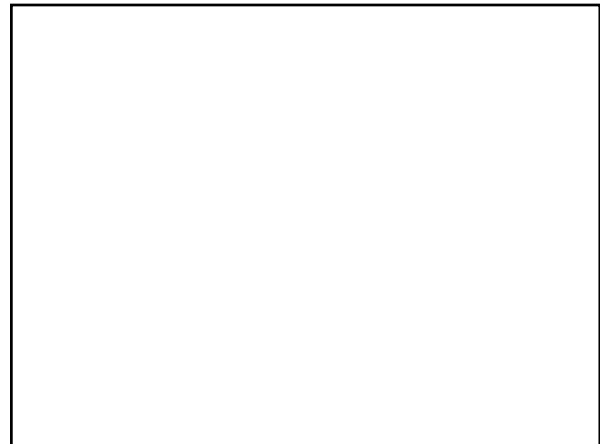
Convert the following into Chomsky normal form:

$A \rightarrow BAB \mid B \mid \epsilon$   
 $B \rightarrow 00 \mid \epsilon$

$S_0 \rightarrow A$        $S_0 \rightarrow A \mid \epsilon$   
 $A \rightarrow BAB \mid B \mid \epsilon$        $A \rightarrow BAB \mid B \mid BB \mid AB \mid BA$   
 $B \rightarrow 00 \mid \epsilon$        $B \rightarrow 00$

$S_0 \rightarrow BAB \mid 00 \mid BB \mid AB \mid BA \mid \epsilon$   
 $A \rightarrow BAB \mid 00 \mid BB \mid AB \mid BA$   
 $B \rightarrow 00$

$S_0 \rightarrow BC \mid DD \mid BB \mid AB \mid BA \mid \epsilon, C \rightarrow AB,$   
 $A \rightarrow BC \mid DD \mid BB \mid AB \mid BA, B \rightarrow DD, D \rightarrow 0$



FORMAL DEFINITIONS

deterministic DFA  
**A finite automaton is a 5-tuple  $M = (Q, \Sigma, \delta, q_0, F)$**

**Q is the set of states (finite)**

**$\Sigma$  is the alphabet (finite)**

**$\delta : Q \times \Sigma \rightarrow Q$  is the transition function**

**$q_0 \in Q$  is the start state**

**$F \subseteq Q$  is the set of accept states**

Let  $w_1, \dots, w_n \in \Sigma$  and  $w = w_1 \dots w_n \in \Sigma^*$   
 Then **M accepts w** if there are  $r_0, r_1, \dots, r_n \in Q$ , s.t.

1.  $r_0 = q_0$
2.  $\delta(r_i, w_{i+1}) = r_{i+1}$ , for  $i = 0, \dots, n-1$ , and
3.  $r_n \in F$

**A non-deterministic finite automaton (NFA)** is a 5-tuple  $N = (Q, \Sigma, \delta, Q_0, F)$

**Q is the set of states**

**$\Sigma$  is the alphabet**

**$\delta : Q \times \Sigma_\epsilon \rightarrow 2^Q$  is the transition function**

**$Q_0 \subseteq Q$  is the set of start states**

**$F \subseteq Q$  is the set of accept states**

$2^Q$  is the set of all possible subsets of Q  
 $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$

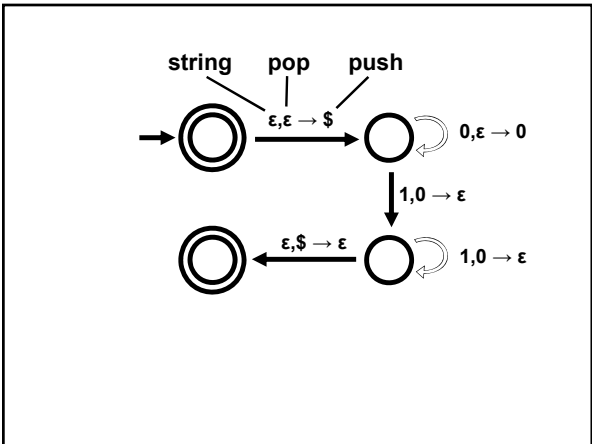
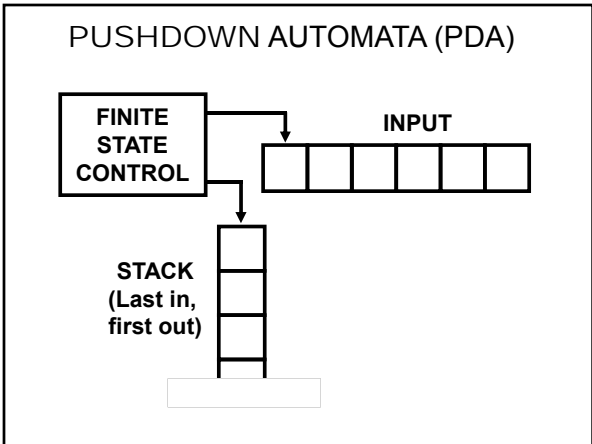
Let  $w \in \Sigma^*$  and suppose w can be written as  $w_1 \dots w_n$  where  $w_i \in \Sigma_\epsilon$  ( $\epsilon$  = empty string)

Then **N accepts w** if there are  $r_0, r_1, \dots, r_n \in Q$  such that

1.  $r_0 \in Q_0$
2.  $r_{i+1} \in \delta(r_i, w_{i+1})$  for  $i = 0, \dots, n-1$ , and
3.  $r_n \in F$

**$L(N)$  = the language recognized by N**  
 = set of all strings machine N accepts

**A language L is recognized by an NFA N**  
 if  $L = L(N)$ .





**Definition:** A (non-deterministic) PDA is a tuple  $P = (Q, \Sigma, \Gamma, \delta, q_0, F)$ , where:

- $Q$  is a finite set of states
- $\Sigma$  is the input alphabet
- $\Gamma$  is the stack alphabet
- $\delta : Q \times \Sigma_\epsilon \times \Gamma_\epsilon \rightarrow 2^{Q \times \Gamma_\epsilon}$
- $q_0 \in Q$  is the start state
- $F \subseteq Q$  is the set of accept states

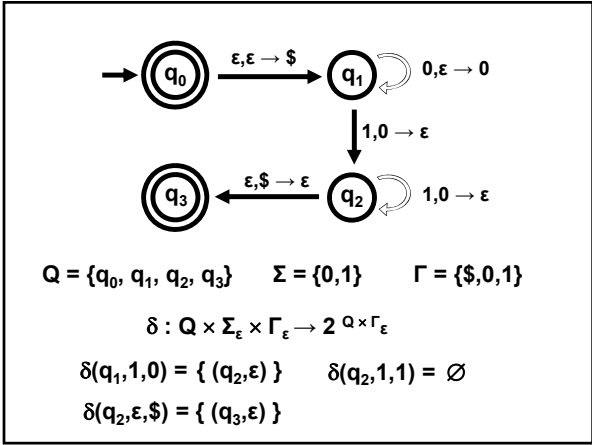
$2^Q$  is the set of subsets of  $Q$  and  $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$

Let  $w \in \Sigma^*$  and suppose  $w$  can be written as  $w_1 \dots w_n$  where  $w_i \in \Sigma_\epsilon$  (recall  $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$ )

Then  $P$  accepts  $w$  if there are

- $r_0, r_1, \dots, r_n \in Q$  and
- $s_0, s_1, \dots, s_n \in \Gamma^*$  (sequence of stacks) such that

1.  $r_0 = q_0$  and  $s_0 = \epsilon$  ( $P$  starts in  $q_0$  with empty stack)
2. For  $i = 0, \dots, n-1$ :  
 $(r_{i+1}, b) \in \delta(r_i, w_{i+1}, a)$ , where  $s_i = at$  and  $s_{i+1} = bt$  for some  $a, b \in \Gamma_\epsilon$  and  $t \in \Gamma^*$   
 ( $P$  moves correctly according to state, stack and symbol read)
3.  $r_n \in F$  ( $P$  is in an accept state at the end of its input)



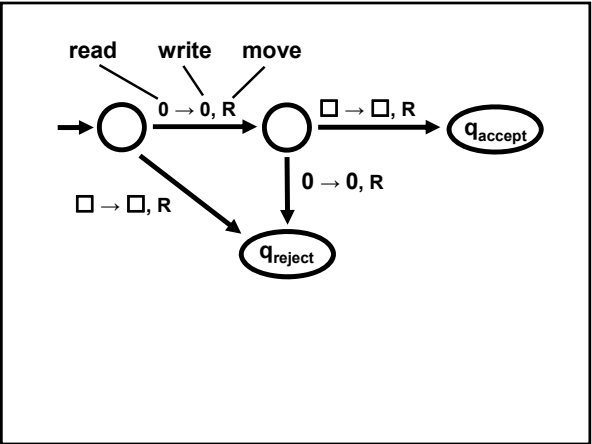
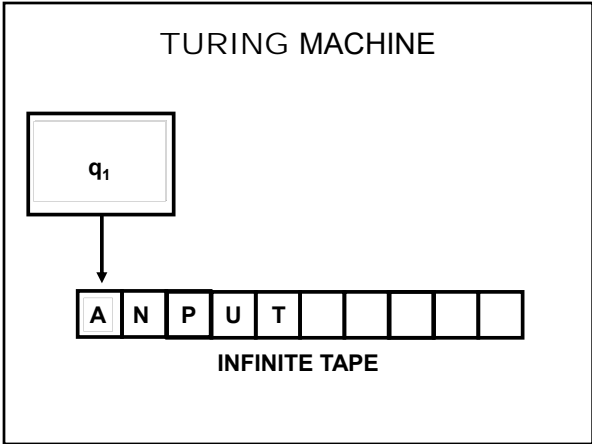
**CONTEXT-FREE GRAMMARS**

A context-free grammar (CFG) is a tuple  $G = (V, \Sigma, R, S)$ , where:

- $V$  is a finite set of variables
- $\Sigma$  is a finite set of terminals (disjoint from  $V$ )
- $R$  is set of production rules of the form  $A \rightarrow W$ , where  $A \in V$  and  $W \in (V \cup \Sigma)^*$
- $S \in V$  is the start variable
- $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$  Strings Generated by  $G$

$G = \{ \{S\}, \{0, 1\}, R, S \}$      $R = \{ S \rightarrow 0S1, S \rightarrow \epsilon \}$

$L(G) = \{ 0^n 1^n \mid n \geq 0 \}$  Strings Generated by  $G$



**Definition: A Turing Machine is a 7-tuple**  
 $T = (Q, \Sigma, \Gamma, \delta, q_0, q_{\text{accept}}, q_{\text{reject}})$ , where:

**Q is a finite set of states**

**$\Sigma$  is the input alphabet, where  $\square \notin \Sigma$**

**$\Gamma$  is the tape alphabet, where  $\square \in \Gamma$  and  $\Sigma \subseteq \Gamma$**

**$\delta : Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$**

**$q_0 \in Q$  is the start state**

**$q_{\text{accept}} \in Q$  is the accept state**

**$q_{\text{reject}} \in Q$  is the reject state, and  $q_{\text{reject}} \neq q_{\text{accept}}$**

A Turing Machine M accepts input w if there is a sequence of configurations  $C_1, \dots, C_k$  such that

1.  $C_1$  is a *start* configuration of M on input w, ie  $C_1$  is  $q_0w$
2. each  $C_i$  *yields*  $C_{i+1}$ , ie M can legally go from  $C_i$  to  $C_{i+1}$  in a single step

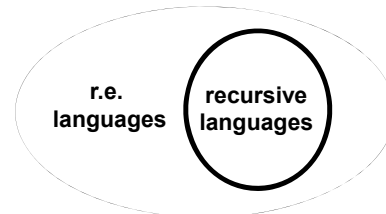
$u a q_i b v \text{ yields } u q_j a c v \text{ if } \delta(q_i, b) = (q_j, c, L)$   
 $u a q_i b v \text{ yields } u a c q_j v \text{ if } \delta(q_i, b) = (q_j, c, R)$

A Turing Machine M accepts input w if there is a sequence of configurations  $C_1, \dots, C_k$  such that

1.  $C_1$  is a *start* configuration of M on input w, ie  $C_1$  is  $q_0w$
2. each  $C_i$  *yields*  $C_{i+1}$ , ie M can legally go from  $C_i$  to  $C_{i+1}$  in a single step
3.  $C_k$  is an *accepting* configuration, ie the state of the configuration is  $q_{\text{accept}}$

**A language is called Turing-recognizable or recursively enumerable (r.e.) or semi-decidable if some TM recognizes it**

**A language is called decidable or recursive if some TM decides it**



**Theorem: If A and  $\neg A$  are r.e. then A is recursive**

$A_{\text{DFA}} = \{ (B, w) \mid B \text{ is a DFA that accepts string } w \}$

**Theorem:  $A_{\text{DFA}}$  is decidable**

**Proof Idea: Simulate B on w**

$A_{\text{NFA}} = \{ (B, w) \mid B \text{ is an NFA that accepts string } w \}$

**Theorem:  $A_{\text{NFA}}$  is decidable**

$A_{\text{CFG}} = \{ (G, w) \mid G \text{ is a CFG that generates string } w \}$

**Theorem:  $A_{\text{CFG}}$  is decidable**

**Proof Idea: Transform G into Chomsky Normal Form. Try all derivations of length up to  $2|w|-1$**

WWW.FLAC.WS

Happy studying!