

# Formal Language Theory

There is nondeterminism in how the string  $w$  might be broken.

$$A \cdot B = \{w : w \text{ can be broken as } w =: x \cdot y \text{ such that } x \in A \text{ \& } y \in B\}$$

$$\alpha \cdot \beta = \{w : w \text{ can be broken as } w =: x \cdot y \text{ such that } x \in A \text{ \& } y \in B\}$$

**!**  $A \cdot B \neq B \cdot A$  in general (as with matrices)

$$A^2 =_{\text{def}} A \cdot A = \{x \cdot y : x \in A \text{ \& } y \in A\} \text{ NOT } \{x \cdot x : x \in A\}$$

$$A \cdot A = A \cdot A = \{w : w \text{ can be broken as } w =: x \cdot y \text{ st } x \in A \text{ \& } y \in A\}$$

$$A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \cup \dots = \bigcup_{i=0}^{\infty} A^i$$

Rule: For any language  $A$ , even  $A = \emptyset$ ,  $A^0 = \{\epsilon\}$ . Like  $a^0 = 1$ . in math.

Rule: For all languages  $A$ ,  $A \cdot \emptyset = \emptyset$  and  $A \cdot \{\epsilon\} = A$  (char)  $\emptyset \cdot A = \emptyset$  and  $\{\epsilon\} \cdot A = A$

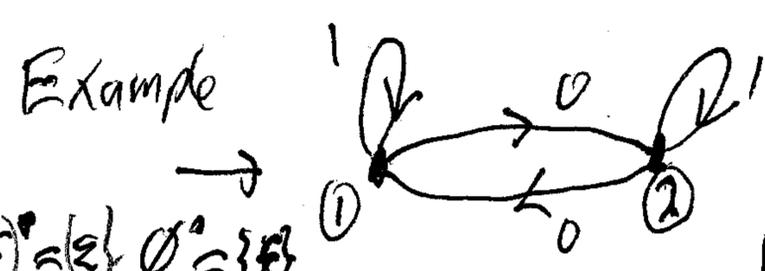
For all languages  $A, B$ , and  $C$ :

Distributive Law:  $A \cdot (B \cup C) = A \cdot B \cup A \cdot C$  and  $(A \cup B) \cdot C = A \cdot C \cup B \cdot C$

With regexps  $d, \beta, \gamma$   $\alpha \cdot (\beta + \gamma) = \alpha \cdot \beta + \alpha \cdot \gamma$   $(\alpha + \beta) \cdot \gamma = \alpha \cdot \gamma + \beta \cdot \gamma$

With regexps  $r, t, u$   $r(t+u) = rt + ru$  Main Difference:  $(r+r) = r$   $A \cup A = A$

An example of a totally new kind of rule:  $(\epsilon + r)^* = r^*$   $(\{\epsilon\} \cup A)^* = A^*$

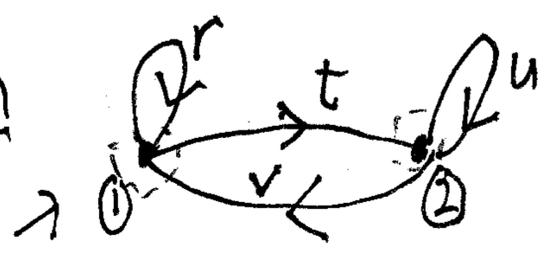


$L_{11} = \{x \in \{0,1\}^* : \#0(x) \text{ is even}\}$   $L_{12} = \{x \in \{0,1\}^* : \#0(x) \text{ is odd}\}$   
 $L_{11} = \left( \begin{matrix} \text{"once"} \\ L_{11} \end{matrix} \right)^*$

Rule: For any state  $p$ ,  $L_{pp} = r^*$  for some regexp  $r$ .

$L_{11}^{\text{once}} = (1 + 01^*0)$   $\therefore L_{11} = (1 + 01^*0)^*$  is equiv to  $1^*(01^*01^*)^*$   
 $L_{12} = L_{11}^{\text{once}}$   $01^* = (1 + 01^*0)^* 01^* = L_{12} \cdot L_{22} = 1^*0 \cdot \left( \begin{matrix} \text{"once"} \\ L_{22} \end{matrix} \right)^*$

The Abstract  
2-State GNFA

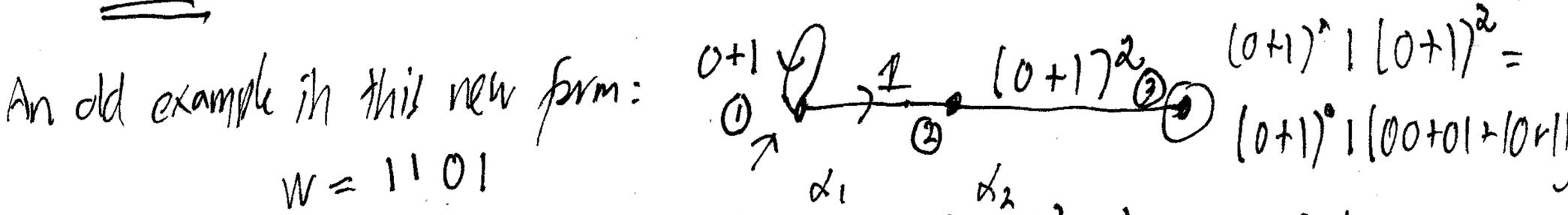


$L_{11} = (r + tu^*v)^*$   $L_{12} = L_{11}tu^* = (r + tu^*v)^*tu^*$   
 $L_{22} = (u + vr^*t)^*$   $L_{21} = L_{22}vr^* = (u + vr^*t)^*vr^*$

Def<sup>n</sup>: A Generalized NFA (GNFA) is a 5-tuple  $N = (Q, \Sigma, \delta, s, F)$  which is like an NFA except  $\delta \subseteq Q \times \text{Regexps}(\Sigma) \times Q$ .  $(p, \alpha, q)$  where  $\alpha$  is a regexp.

In fact, an NFA is just the case  $\delta \subseteq Q \times \text{BasicRegexps}(\Sigma) \times Q$  where  $\alpha$  is a regexp.

Def<sup>n</sup>: A computation path from a state  $p$  to a state  $q$  is a sequence  $(p, \alpha_1, q_1, \alpha_2, q_2, \alpha_3, q_3, \dots, q_{m-1}, \alpha_m, q_m = q)$  such that for all  $i, 1 \leq i \leq m, (q_{i-1}, \alpha_i, q_i)$  is an instruction in  $\delta$ .  
 (This computation path of  $N$  can process a string  $w \in \Sigma^*$  from  $p$  to  $q$  if  $w$  can be broken into  $m$  substrings  $w =: u_1 \cdot u_2 \cdot u_3 \dots u_m$  such that for each  $i, 1 \leq i \leq m$   $u_i \in L(\alpha_i)$  ( $u_i$  matches  $\alpha_i$ )  
 Finally,  $L_{pq} = \{w : N \text{ can process } w \text{ from } p \text{ to } q\}$ .



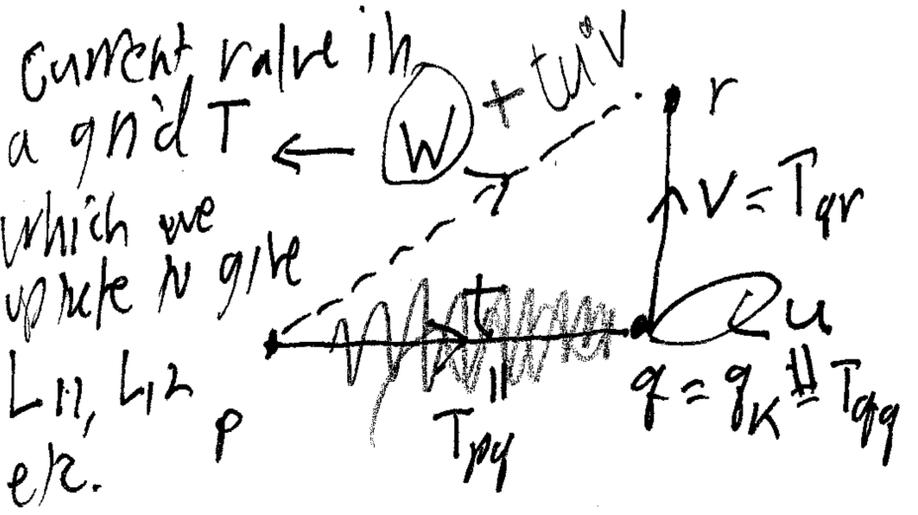
To witness this defn, if I try path =  $(1, 1, 2, (0+1)^2, 3)$  then I fail.  
 $w = u_1 \cdot u_2$  with  $u_1 = 1$  and  $u_2 = 101 \notin L((0+1)^2)$ .  
 Good path:  $(1, 0+1, 1, 1, 2, (0+1)^2, 3)$   
 Matched with  $w = \underline{1} \cdot \underline{1} \cdot \underline{01}$   
 $u_1 = 1, u_2 = 1, u_3 = 01$

$$L(N) = \bigcup_{q \in F} L_{sq} = \text{just } L_{13}$$

Theorem: Given any DFA or NFA or GNFA  $N$ , we can compute a regexp  $\alpha$  such that  $L(N) = L(\alpha)$ .  
 Text target  $s \rightarrow q$

Proof: We have already proved this for  $K \leq 2$ , where  $N = (Q, \Sigma, \delta, s, F)$  and  $K = |Q|$ .  
 Given  $N$  with  $K \geq 2$ , if  $N$  has more than one accepting state  $q \neq s$ , then add a new accepting state  $f$  and  $\epsilon$ -arcs from all  $q \in F$  to  $f$ .  
 Well grounded?

Hence we may suppose  $N$  is in well-graded form. <sup>(3) (maybe old)</sup>  
 Number  $s=1$ , and the final states  $2, \{1, 2\}$ , or just  $1$  if  $F=\{s\}$ .  
 Thus states  $3, \dots, K$  are non-accepting. We will eliminate them one-by-one.

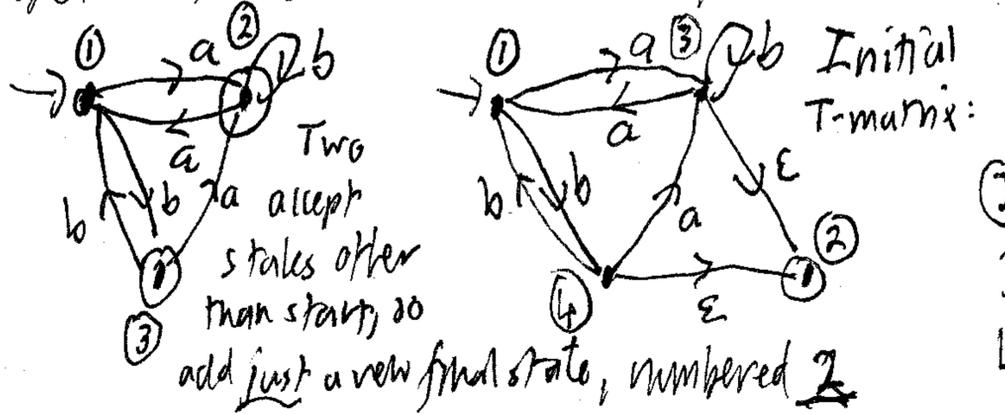


To eliminate the highest-numbered state  $q$ , we bypass all arcs  $(p, t, q)$  into  $q$ .  
 To bypass  $(p, t, q)$ , for all outgoing arcs  $(q, v, r)$ , we update the direct path  $p \rightarrow r$

by  $T_{pr} += t u^r v$  Then we can delete  $(p, t, q)$ .  
 New  $T_{pr} = W + t u^r v$ .  
 When all incoming arcs to  $q$  have been bypassed, then we delete  $q$ .

for  $(int\ j = K; j \geq 3; j--)$  { // eliminate state  $j$  ie state  $q_j$   
 for  $(int\ i = 1; i \leq j-1; i++)$  {  
 if (there is an instruction  $(i, t, j)$  in  $\delta$ ) {  
 for  $(int\ h = 1; h \leq j-1; h++)$  {  
 if (there is an instruction  $(j, v, h)$ ) {  
 $T_{i,h} += T_{i,t} T_{t,j}^* T_{j,v} T_{j,h}$ .  
 }  
 }  
 }  
 }  
 (that is, new  $T_{i,h} = old\ T_{i,h} \cup T_{i,t} T_{t,j}^* T_{j,v} T_{j,h}$ )  
 Read off final answer from  $K=2$  case after outer loop exits.

**Extra:** You can regard  $T$  as a "matrix of regular expressions" if you think better "in code" than "in pictures." E.g. for my version of the text, example on p76:



Initial T-matrix:

	1	2	3	4
1	$\emptyset$	$\emptyset$	$a$	$b$
2	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
3	$a$	$\epsilon$	$b$	$\emptyset$
4	$b$	$\epsilon$	$a$	$\emptyset$

$T_{11}^* T_{12} = \text{Final answer.}$

Example (text w/o adding new start state, still remember by new dest state <sup>(4)</sup> as 2)

Elim state 4: Incoming: Just (1, b, 4)

Outgoing: (4, b, 1) ∴ Update (1, 1) (1, —, 1)

$\epsilon^r = \epsilon$  and  $\emptyset^o = \epsilon$

(4, a, 3) Update (1, —, 3)  
 (4, ε, 2)  $T(4,4) = \emptyset$  but  $T(4,4)^o = \emptyset^o = \epsilon$ .

New  $T(1,1) = \text{old } T(1,1) \cup T(1,4)T(4,4)^*T(4,1)$

You can also initialize  $T(1,1) = \epsilon$ , likewise  $T(i,i)$  for all  $i$ .  
 $\emptyset \cup b \cdot \epsilon \cdot b = bb$   
 $\epsilon \cup b \cdot \epsilon \cdot b = \epsilon + bb$

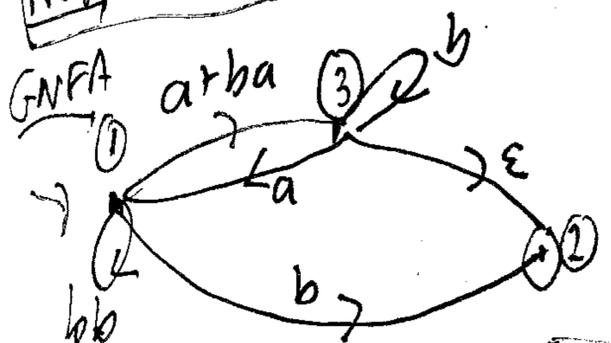
Bypass (1, b, 4) to 3 expression will have  $T(1,1)^*$  and  $(bb)^* = (\epsilon + bb)^*$  as a component term.

New  $T(1,3) = \text{old } T(1,3) \cup T(1,4)T(4,4)^*T(4,3)$

Bypass (1, b, 4) to 2.  $a \cup b \cdot \epsilon \cdot a = a + ba$

New  $T(1,2) = \text{old } T(1,2) \cup T(1,4)T(4,4)^*T(4,2) = b$

Now elim (4)



We cannot replace this  $\emptyset$  by  $\epsilon$  because dest state 2  $\neq$  origin 1.

Elim 3: Incoming: only (1, —, 3)

Outgoing: (3, —, 1) ∴ Update  $T(1,1)$  again  
 (3, ε, 2) Update  $T(1,2)$

Bypass (1, a+ba, 3) back to 1

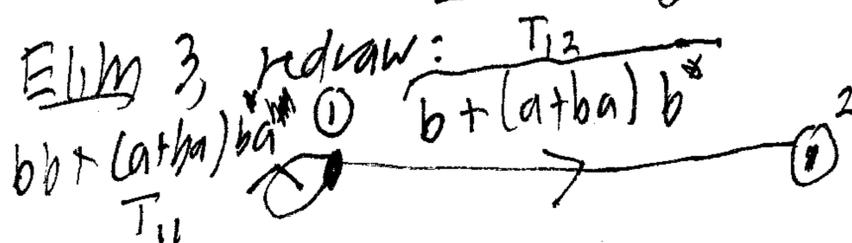
New  $T(1,1) = \text{old } T(1,1) \cup T(1,3)T(3,3)^*T(3,1)$

$bb + (a+ba)b^a$  or  $\epsilon + bb + (a+ba)b^a$

Bypass (1, a+ba, 3) to 2

New  $T(1,2) = \text{old } T(1,2) \cup T(1,3)T(3,3)^*T(3,2)$

$b + (a+ba)b^a$  putting  $\epsilon$  here makes no difference



$L_{11} = T_{11}^* = (r^{\text{as}}) = (bb + (a+ba)b^a)^*$   
 $L(M) = L_{12} = L_{11} \cdot T_{12} = (bb + (a+ba)b^a)^* \cdot (b + (a+ba)b^a)$

The last 3 1/2 lectures have proved a theorem first obtained by Stephen Kleene in the mid-1950s:

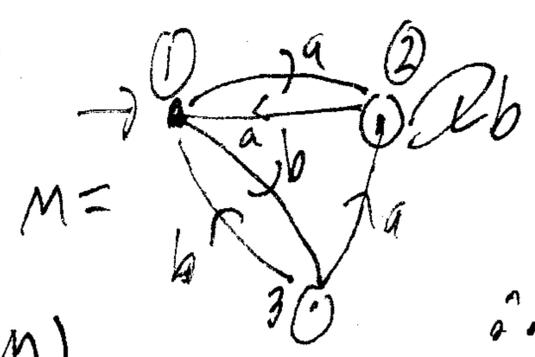
Theorem: For any language  $L \subseteq \Sigma^*$  (any  $\Sigma$ ):

- (a) there is a DFA  $M$  st.  $L = L(M)$
  - (b) there is an NFA  $N$  st.  $L = L(N)$  (GNFA ditto)
  - (c) there is a regular expression  $r$  st.  $L = L(r)$
- }  $L$  is a regular language if any holds.

When is a language  $L$  (non-) regular?

Related: Why is a regular  $L$  regular?

Let's revisit the original DFA



Note that a and ba both get processed from 1 to 2.

Let  $L = L(M)$

★ If  $x$  and  $y$  are any two strings such that for some  $z \in \Sigma^*$

$xz \in L \wedge yz \notin L$   
 or  $xz \notin L \wedge yz \in L$

∴ For any suffix string  $z$

$az$  and  $baz$  get processed to the same state.

∴  $(\forall z \in \Sigma^*) az \in L(M) \Leftrightarrow baz \in L(M)$

If  $z = \epsilon$ , or if  $z = ab$ , they're both in  $L(M)$

If  $z = a$  they're both out: aa, ba  $\notin L(M)$ .

$xz \in L \not\Leftrightarrow yz \in L$

ie.  $xz \in L$  XOR  $yz \in L$ , then  $x$  and  $y$  CANNOT be processed to the same state

ie.  $L(xz) \neq L(yz)$  I.e. they must be processed to different states

∴ Hence,  $M$  must have at least 2 states.

Suppose  $S$  is a set of strings such that for any distinct  $x, y \in S$  there exists  $z$  st.  $L(xz) \neq L(yz)$ .

Then for any  $x, y \in S$  ( $x \neq y$ ), any DFA  $M_n$  must process  $x, y$  to different states (st.  $L(M) = L$ ).

- ∴ If  $S$  has  $K$  strings, then any DFA  $M_n$  must have at least  $K$  states.
- ∴ If  $S$  has  $\infty$ -many strings, then any DFA  $M$  st.  $L = L(M)$  must have at least  $\infty$ -many states. Hence  $L$  has no finite DFA, so  $L$  is not regular.

∴ Theorem : Suppose  $L$  is a language and  $S$  is a set of strings st.  $S$  is infinite, and  
 (Half of the Myhill-Nerode Thm) • for all  $x, y \in S$  ( $x \neq y$ ) ( $\exists z \in \Sigma^*$ )  $L(xz) \neq L(yz)$ .

Then  $L$  is not regular.

Example =  $L = \{a^n b^n : n \geq 0\}$ . To <sup>prove</sup> show  $L$  is not regular:

Take  $S = a^* = \{a^n : n \geq 0\}$ . Clearly  $S$  is infinite

Let any  $x, y \in S$ ,  $x \neq y$ , be given. Then there are numbers

$m, n \geq 0$ ,  $m \neq n$ , s.t.  $x = a^m$  and  $y = a^n$ .

Take  $z = b^m$ .

Then  $xz = a^m b^m \in L$ , but  $yz = a^n b^m \notin L$  since  $m \neq n$ .

∴  $L(xz) \neq L(yz)$ , and since  $x, y \in S$  are arbitrary,

$S$  is an infinite PD set for  $L$ , ∴  $L$  is not regular.  $\square$

Extra: The other half of the theorem: If  $L$  is not regular, then there is an infinite set  $S$  such that  $(\forall x, y \in S, x \neq y) (\exists z \in \Sigma^*) L(xz) \neq L(yz)$  i.e. such that  $S$  is PD (for  $L$ .)  
 This means that in principle every nonregular language has some kind of Myhill-Nerode proof.

The contrapositive of this is: If every PD set  $S$  for  $L$  is finite, then  $L$  is regular.  
 We can prove this - which completes both halves of the theorem - without going into quite as much detail as the text in problem 1-52. Consider the following unbounded "allocation" process:

For string  $x = \epsilon, 0, 1, 00, 01, 10, 11, 000, \dots$  in order, (having initialized  $S = \emptyset$ )  $\left\{ \begin{array}{l} \text{If } (\exists y < x) (\forall z \in \Sigma^*) L(yz) = L(xz) \text{ then define } \text{State}(x) = \text{State}(y). \\ \text{Else } \text{allocate } \text{State}(x) \text{ as a new state and put } S := S \cup \{x\}; // \text{Invariant: } S \text{ is PD for } L. \end{array} \right.$  Note:  $S$  always gets  $\epsilon$  right away.  
 If  $(x = 1^n)$  then for all  $w$  of length  $n-1$  and  $c \in \Sigma$ , define  $\text{delta}(w, c) = \text{State}(wc)$ ;  
 } By hypothesis,  $S$  stays finite, and so  $(S, \Sigma, \text{delta}, \text{State}(\epsilon), F)$  is a DFA  $M$  st.  $L(M) = L$ , where  $F = \{\text{State}(x) : x \in L\}$ !