

# THEORY OF COMPUTATION (II)

Yijia Chen  
Fudan University

Review

A **language**  $L$  is a subset of strings over an **alphabet**  $\Sigma$ .

Our goal is to identify those languages that can be recognized by one of the simplest computation models – **finite automata**.

We will also show that those languages can be formed by some simple syntactic manipulations – **regular expressions**.

# Finite automata

## Definition

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

1.  $Q$  is a finite set called the states,
2.  $\Sigma$  is a finite set called the alphabet,
3.  $\delta : Q \times \Sigma \rightarrow Q$  is the transition function,
4.  $q_0 \in Q$  is the start state, and
5.  $F \subseteq Q$  is the set of accept states.

## Formal definition of computation

Let  $M = (Q, \Sigma, \delta, q_0, F)$  be a finite automaton and let  $w = w_1 w_2 \cdots w_n$  be a string with  $w_i \in \Sigma$  for all  $i \in [n]$ . Then  $M$  accepts  $w$  if a sequence of states  $r_0, r_1, \dots, r_n$  in  $Q$  exists with:

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

We say that  $M$  recognizes  $A$  if

$$A = \{w \mid M \text{ accepts } w\}.$$

# Regular languages

## Definition

A language is called regular if some finite automaton recognizes it.

## The regular operators

### Definition

Let  $A$  and  $B$  be languages. We define the regular operations union, concatenation, and star as follows:

- ▶ Union:  $A \cup B = \{x \mid x \in A \text{ or } x \in B\}$ .
- ▶ Concatenation:  $A \circ B = \{xy \mid x \in A \text{ and } y \in B\}$ .
- ▶ Star:  $A^* = \{x_1x_2 \dots x_k \mid k \geq 0 \text{ and each } x_i \in A\}$ .

## Closure under union

### Theorem

*The class of regular languages is closed under the union operation.*

*In other words, if  $A_1$  and  $A_2$  are regular languages, so is  $A_1 \cup A_2$ .*



## Proof (1)

For  $i \in [2]$  let  $M_i = (Q_i, \Sigma_i, \delta_i, q_i, F_i)$  recognize  $A_i$ . We can assume without loss of generality  $\Sigma_1 = \Sigma_2$ :

- ▶ Let  $a \in \Sigma_2 - \Sigma_1$ .
- ▶ We add  $\delta_1(r, a) = r_{\text{trap}}$ , where  $r_{\text{trap}}$  is a new state with

$$\delta_1(r_{\text{trap}}, w) = r_{\text{trap}}$$

for every  $w$ .

## Proof (2)

We construct  $M = (Q, \Sigma, \delta, q_0, F)$  to recognize  $A_1 \cup A_2$ :

1.  $Q = Q_1 \times Q_2 = \{(r_1, r_2) \mid r_1 \in Q_1 \text{ and } r_2 \in Q_2\}$ .
2.  $\Sigma = \Sigma_1 = \Sigma_2$ .
3. For each  $(r_1, r_2) \in Q$  and  $a \in \Sigma$  we let

$$\delta((r_1, r_2), a) = (\delta_1(r_1, a), \delta_2(r_2, a)).$$

4.  $q_0 = (q_1, q_2)$ .
5.  $F = (F_1 \times Q_2) \cup (Q_1 \times F_2) = \{(r_1, r_2) \mid r_1 \in F_1 \text{ or } r_2 \in F_2\}$ .

## Closure under concatenation

### Theorem

*The class of regular languages is closed under the concatenation operation.*

*In other words, if  $A_1$  and  $A_2$  are regular languages, so is  $A_1 \circ A_2$ .*

We prove the above theorem by **nondeterministic finite automata** which

- ▶ have the same power as deterministic finite automata,
- ▶ are much easier to manipulate.

# Nondeterminism

## Definition

A nondeterministic finite automaton (NFA) is a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$ , where

1.  $Q$  is a finite set of states,
2.  $\Sigma$  is a finite alphabet,
3.  $\delta : Q \times \Sigma_\epsilon \rightarrow \mathcal{P}(Q)$  is the transition function, where  $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$ ,
4.  $q_0 \in Q$  is the start state, and
5.  $F \subseteq Q$  is the set of accept states.

## Formal definition of computation

Let  $N = (Q, \Sigma, \delta, q_0, F)$  be an NFA and let  $w = y_1 y_2 \cdots y_m$  be a string with  $y_i \in \Sigma_\varepsilon$  for all  $i \in [m]$ . Then  $N$  accepts  $w$  if a sequence of states  $r_0, r_1, \dots, r_m$  in  $Q$  exists with:

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

## Equivalence of NFAs and DFAs

### Theorem

*Every NFA has an equivalent DFA, i.e., they recognize the same language.*

## Corollary

*A language is regular if and only if some nondeterministic finite automaton recognizes it.*

## Second proof of the closure under union

For  $i \in [2]$  let  $N_i = (Q_i, \Sigma_i, \delta_i, q_i, F_i)$  recognize  $A_i$ . We construct an  $N = (Q, \Sigma, \delta, q_0, F)$  to recognize  $A_1 \cup A_2$ :

1.  $Q = \{q_0\} \cup Q_1 \cup Q_2$ .
2.  $q_0$  is the start state.
3.  $F = F_1 \cup F_2$ .
4. For any  $q \in Q$  and any  $a \in \Sigma_\varepsilon$

$$\delta(q, a) = \begin{cases} \delta_1(q, a) & q \in Q_1 \\ \delta_2(q, a) & q \in Q_2 \\ \{q_1, q_2\} & q = q_0 \text{ and } a = \varepsilon \\ \emptyset & q = q_0 \text{ and } a \neq \varepsilon. \end{cases}$$



## Closure under concatenation

### Theorem

*The class of regular languages is closed under the concatenation operation.*

## Proof

For  $i \in [2]$  let  $N_i = (Q_i, \Sigma_i, \delta_i, q_i, F_i)$  recognize  $A_i$ . We construct an  $N = (Q, \Sigma, \delta, q_1, F_2)$  to recognize  $A_1 \circ A_2$ :

1.  $Q = Q_1 \cup Q_2$ .
2. The start state  $q_1$  is the same as the start state of  $N_1$ .
3. The accept states  $F_2$  are the same as the accept states of  $N_2$ .
4. For any  $q \in Q$  and any  $a \in \Sigma_\varepsilon$

$$\delta(q, a) = \begin{cases} \delta_1(q, a) & q \in Q_1 - F_1 \\ \delta_1(q, a) & q \in F_1 \text{ and } a \neq \varepsilon \\ \delta_1(q, a) \cup \{q_2\} & q \in F_1 \text{ and } a = \varepsilon \\ \delta_2(q, a) & q \in Q_2. \end{cases}$$

## Closure under star

### Theorem

*The class of regular languages is closed under the star operation.*

## Proof

Let  $N_1 = (Q_1, \Sigma, \delta_1, q_1, F_1)$  recognize  $A_i$ . We construct an  $N = (Q, \Sigma, \delta, q_0, F)$  to recognize  $A_1^*$ :

1.  $Q = \{q_0\} \cup Q_1$ .
2. The start state  $q_0$  is the new start state.
3.  $F = \{q_0\} \cup F_1$ .
4. For any  $q \in Q$  and any  $a \in \Sigma_\epsilon$

$$\delta(q, a) = \begin{cases} \delta_1(q, a) & q \in Q_1 - F_1 \\ \delta_1(q, a) & q \in F_1 \text{ and } a \neq \epsilon \\ \delta_1(q, a) \cup \{q_1\} & q \in F_1 \text{ and } a = \epsilon \\ \{q_1\} & q = q_0 \text{ and } a = \epsilon \\ \emptyset & q = q_0 \text{ and } a \neq \epsilon. \end{cases}$$

# Regular expression

## Definition

We say that  $R$  is a regular expression if  $R$  is

1.  $a$  for some  $a \in \Sigma$ ,
2.  $\varepsilon$ ,
3.  $\emptyset$ ,
4.  $(R_1 \cup R_2)$ , where  $R_1$  and  $R_2$  are regular expressions,
5.  $(R_1 \circ R_2)$ , where  $R_1$  and  $R_2$  are regular expressions,
6.  $(R_1^*)$ , where  $R_1$  is a regular expressions.

We often write  $R_1R_2$  instead of  $(R_1 \circ R_2)$  if no confusion arises.

## Language defined by regular expressions

regular expression $R$	language $L(R)$
$a$	$\{a\}$
$\varepsilon$	$\{\varepsilon\}$
$\emptyset$	$\emptyset$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 \circ R_2)$	$L(R_1) \circ L(R_2)$
$(R_1^*)$	$L(R_1)^*$

## Equivalence with finite automata

### Theorem

*A language is regular if and only if some regular expression describes it.*

## The languages defined by regular expressions are regular

1.  $R = a$ : Let  $N = (\{q_1, q_2\}, \Sigma, \delta, q_1, \{q_2\})$ , where  $\delta(q_1, a) = \{q_2\}$  and  $\delta(r, b) = \emptyset$  for all  $r \neq q_1$  or  $b \neq a$ .
2.  $R = \varepsilon$ : Let  $N = (\{q_1\}, \Sigma, \delta, q_1, \{q_1\})$ , where  $\delta(r, b) = \emptyset$  for all  $r$  and  $b$ .
3.  $R = \emptyset$ : Let  $N = (\{q_1\}, \Sigma, \delta, q_1, \emptyset)$ , where  $\delta(r, b) = \emptyset$  for all  $r$  and  $b$ .
4.  $R = R_1 \cup R_2$ :  $L(R) = L(R_1) \cup L(R_2)$ .
5.  $R = R_1 \circ R_2$ :  $L(R) = L(R_1) \circ L(R_2)$ .
6.  $R = R_1^*$ :  $L(R) = L(R_1)^*$ .



To prove the other direction we need to further generalize NFA to **generalized nondeterministic finite automata**, which might be viewed as a mixture of automata and regular expressions.

# Generalized nondeterministic finite automata

## Definition

A GNFA is a 5-tuple  $(Q, \Sigma, \delta, q_{\text{start}}, q_{\text{accept}})$ , where

1.  $Q$  is a finite set of states,
2.  $\Sigma$  is a finite alphabet,
3.  $\delta : (Q - \{q_{\text{accept}}\}) \times (Q - \{q_{\text{start}}\}) \rightarrow \mathcal{R}$  is the transition function, where  $\mathcal{R}$  is the set of regular expressions,
4.  $q_{\text{start}}$  is the start state, and
5.  $q_{\text{accept}}$  is the accept state.

## Formal definition of computation

A GNFA accepts a string  $w \in \Sigma^*$  if  $w = w_1 w_2 \dots w_k$ , where each  $w_i \in \Sigma^*$  and a sequence of states  $q_0, q_1, \dots, q_k$  exists such that

1.  $q_0 = q_{\text{start}}$  is the start state,
2.  $q_k = q_{\text{accept}}$  is the accept state, and
3. for each  $i \in [k]$ , we have  $w_i \in L(R_i)$ , where  $R_i = \delta(q_{i-1}, q_i)$ .

## Regular languages can be defined by regular expressions

Let  $M$  be the DFA for language  $A$ .

- ▶ We convert  $M$  to a GNFA  $G$  by adding a new start state and a new accept state and additional transition arrows as necessary.
  1. The start state has transition arrows going to every other state but no arrows coming in from any other state.
  2. There is only a single accept state, and it has arrows coming in from every other state but no arrows going to any other state. Furthermore, the accept state is not the same as the start state.
  3. Except for the start and accept states, one arrow goes from every state to every other state and also from each state to itself.
- ▶ Then we use a procedure **CONVERT** on  $G$  to return an equivalent regular expression.

## CONVERT( $G$ ):

1. Let  $k$  be the number of states of  $G$ .
2. If  $k = 2$ , then return the regular expression  $R$  labelling the arrow from  $q_{\text{start}}$  to  $q_{\text{accept}}$ .
3. If  $k > 2$ , we select any state  $q_{\text{rip}} \in Q - \{q_{\text{start}}, q_{\text{accept}}\}$  and let  $G' = (Q', \Sigma, \delta', q_{\text{start}}, q_{\text{accept}})$  be the GNFA, where

$$Q' = Q - \{q_{\text{rip}}\},$$

and for any  $q_i \in Q' - \{q_{\text{accept}}\}$  and  $q_j \in Q' - \{q_{\text{start}}\}$ , let

$$\delta'(q_i, q_j) = (R_1)(R_2)^*(R_3) \cup (R_4),$$

for  $R_1 = \delta(q_i, q_{\text{rip}})$ ,  $R_2 = \delta(q_{\text{rip}}, q_{\text{rip}})$ ,  $R_3 = \delta(q_{\text{rip}}, q_j)$ , and  $R_4 = \delta(q_i, q_j)$ .

4. Compute  $\text{CONVERT}(G')$  and return this value.

## Languages need counting

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

### Lemma

*D is regular.*

## The pumping lemma for regular languages

### Lemma

*If  $A$  is a regular language, then there is a number  $p$  (i.e., the pumping length) where if  $s$  is any string in  $A$  of length at least  $p$ , then  $s$  may be divided into three pieces,  $s = xyz$ , satisfying the following conditions:*

1. *for each  $i \geq 0$ , we have  $xy^iz \in A$ ,*
2.  *$|y| > 0$ , and*
3.  *$|xy| \leq p$ .*

Any string  $xyz$  in  $A$  can be pumped along  $y$ .

## Proof

Let  $M = (Q, \Sigma, \delta, q_1, F)$  be a DFA recognizing  $A$  and  $p := |Q|$ .

Let  $s = s_1s_2 \cdots s_n$  be a string in  $A$  with  $n \geq p$ . Let  $r_1, \dots, r_{n+1}$  be the sequence of states that  $M$  enters while processing  $s$ , i.e.,

$$r_{i+1} = \delta(r_i, s_i)$$

for  $i \in [n]$ .

Among the first  $p + 1$  states in the sequence, two must be the same, say  $r_j$  and  $r_\ell$  with  $j < \ell \leq p + 1$ . We define

$$x = s_1 \cdots s_{j-1}, y = s_j \cdots s_{\ell-1}, \text{ and } z = s_\ell \cdots s_n.$$



## Example (I)

### Example

The language  $\{0^n 1^n \mid n \geq 0\}$  is not regular.

### Proof.

Choose  $p$  be the pumping length and consider  $s = 0^p 1^p$ . By the Pumping Lemma,  $s = xyz$  with  $xy^i z \in \{0^n 1^n \mid n \geq 0\}$  for all  $i \geq 0$ .

1.  $y \in 00^*$ , then  $xyyz$  has more 0s than 1s, a contradiction.
2.  $y \in 11^*$ , then  $xyyz$  has more 1s than 0s, again a contradiction.
3.  $y$  consists of both 0s and 1s, then  $xyyz$  have 0 and 1 interleaved.



## Example (II)

### Example

The language  $\{w \mid w \text{ has an equal number of 0s and 1s}\}$  is not regular.

### Proof.

Choose  $p$  be the pumping length and consider  $s = 0^p 1^p$ . By the Pumping Lemma,  $s = xyz$  with  $|xy| \leq p$  and

$$xy^i z \in \{w \mid w \text{ has an equal number of 0s and 1s}\}$$

for all  $i \geq 0$ . Thus  $xy \in 00^*$  and the contradiction follows easily.



# Context-Free Languages

## An example

The grammar:

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

A derivation:

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000\#111.$$

## A second example (1)

⟨SENTENCE⟩ → ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩

⟨NOUN-PHRASE⟩ → ⟨CMPLX-NOUN⟩|⟨CMPLX-NOUN⟩⟨PREP-PHRASE⟩

⟨VERB-PHRASE⟩ → ⟨CMPLX-VERB⟩|⟨CMPLX-VERB⟩⟨PREP-PHRASE⟩

⟨PREP-PHRASE⟩ → ⟨PREP⟩⟨CMPLX-NOUN⟩

⟨CMPLX-NOUN⟩ → ⟨ARTICLE⟩⟨NOUN⟩

⟨CMPLX-VERB⟩ → ⟨VERB⟩|⟨VERB⟩⟨NOUN-PHRASE⟩

⟨ARTICLE⟩ → **a|the**

⟨NOUN⟩ → **boy|girl|flower**

⟨VERB⟩ → **touches|likes|sees**

⟨PREP⟩ → **with**

## A second example (2)

$\langle \text{SENTENCE} \rangle \Rightarrow \langle \text{NOUN-PHRASE} \rangle \langle \text{VERB-PHRASE} \rangle$   
 $\Rightarrow \langle \text{CMPLX-NOUN} \rangle \langle \text{VERB-PHRASE} \rangle$   
 $\Rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle \langle \text{VERB-PHRASE} \rangle$   
 $\Rightarrow \mathbf{a} \langle \text{NOUN} \rangle \langle \text{VERB-PHRASE} \rangle$   
 $\Rightarrow \mathbf{a} \text{ boy} \langle \text{VERB-PHRASE} \rangle$   
 $\Rightarrow \mathbf{a} \text{ boy} \langle \text{CMPLX-VERB} \rangle$   
 $\Rightarrow \mathbf{a} \text{ boy} \langle \text{VERB} \rangle$   
 $\Rightarrow \mathbf{a} \text{ boy sees.}$

## Context-free grammars

### Definition

A context-free grammar (CFL) is a 4-tuple  $(V, \Sigma, R, S)$ , where

1.  $V$  is a finite set called the variables,
2.  $\Sigma$  is a finite set, disjoint from  $V$ , called the terminals.
3.  $R$  is a finite set of rules, with each rule being a variable and a string of variables and terminals, and
4.  $S \in V$  is the start variable.

## Derivations

Let  $u, v, w$  be strings of variables and terminals, and

$$A \rightarrow w$$

is a rule of the grammar. We say that  $uAv$  yields  $uwv$ , written  $uAv \Rightarrow uwv$ .

Say that  $u$  derives  $v$ , written  $u \xRightarrow{*} v$ , if  $u = v$  or if a sequence  $u_1, u_2, \dots, u_k$  exists for  $k \geq 0$  and

$$u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_k \Rightarrow v.$$

The language of the grammar is

$$\{w \in \Sigma^* \mid S \xRightarrow{*} w\},$$

which is a context-free language (CFL).



## Examples

1. The language  $\{0^n 1^n \mid n \geq 0\}$  has a grammar:

$$S_1 \rightarrow 0S_11 \mid \varepsilon.$$

2. The language  $\{1^n 0^n \mid n \geq 0\}$  has a grammar:

$$S_2 \rightarrow 1S_20 \mid \varepsilon.$$

3. The language

$$\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$$

has a grammar:

$$S \rightarrow S_1 \mid S_2$$

$$S_1 \rightarrow 0S_11 \mid \varepsilon$$

$$S_2 \rightarrow 1S_20 \mid \varepsilon.$$