

**FORMAL LANGUAGES AND AUTOMATA, 2024 FALL SEMESTER**

# Lec 04. Regular expression

**Eunjung Kim**

# REGULAR EXPRESSION

We want to compactly describe the 'pattern' of the following languages using union, concatenation and kleene star operations.

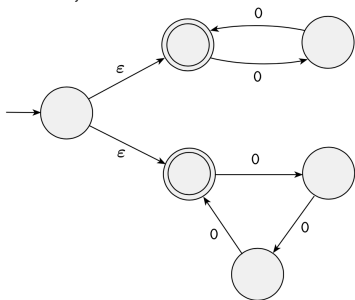


Figure 1.33, Sipser 2012.

- 1 The set of all 0, 1 strings with exactly one 1's.
- 2 The set of all 0, 1 strings with at least one 1's.
- 3 The set of strings over  $\Sigma$  of even length.

# FORMAL DEFINITION OF REGULAR EXPRESSION

## REGULAR EXPRESSION OVER A FINITE ALPHABET $\Sigma$

Regular expression over  $\Sigma$  is a string consisting of symbols of  $\Sigma$ , parenthesis ( ), and the operators  $\cup, \circ, *$  that can be generated as follows.

- Each symbol  $x \in \Sigma \cup \{\epsilon\}$  is a regular expression.
- $\emptyset$  is a regular expression.
- $(R_1 \cup R_2)$  is a regular expression if  $R_1$  and  $R_2$  are regular expressions.
- $(R_1 \circ R_2)$  is a regular expression if  $R_1$  and  $R_2$  are regular expressions.
- $R^*$  is a regular expression if  $R$  is a regular expression

# EXAMPLES OF REGULAR EXPRESSION

Assume  $\Sigma = \{0, 1\}$ . Which language does the regular expression describe?

1  $0^*10^*$ .

2  $\Sigma^*1\Sigma^*$ .

3  $1^*(01^+)^*$ .

4  $(\Sigma\Sigma)^*$ .

5  $0\Sigma^*0 \cup 1\Sigma^*1 \cup 0 \cup 1$ .

6  $1^*\emptyset = \emptyset$ .

7  $\emptyset^* = \{\epsilon\}$ .

# REGULAR LANGUAGE

## VALUE OF A REGULAR EXPRESSION

For a regular expression  $R$ , the set of all strings which can be generated following the expression is denoted by  $\mathcal{L}(R)$ , said to be the language of  $R$ .  $\mathcal{L}(R)$  is also called the value of  $R$ , or the language described by  $R$ .

For two regular expressions  $R_1$  and  $R_2$ , the value of union / concatenation / kleene star is...

- $L(R_1 \cup R_2) = L(R_1) \cup L(R_2)$ .
- $L(R_1 \circ R_2) = L(R_1) \circ L(R_2)$ .
- $L(R_1^*) = L(R_1)^*$ .

# EQUIVALENCE OF REGULAR EXPRESSION AND FINITE AUTOMATA

## REGULAR EXPRESSION=NFA

A language  $A \subseteq \Sigma^*$  is described by a regular expression if and only if it is a regular language, i.e. recognized by some NFA.

One direction can be proved easily using what we learnt in the previous lecture. Which direction is it?

# EQUIVALENCE PROOF: EASY DIRECTION

## EQUIVALENCE THEOREM, PART I

If a language  $A \subseteq \Sigma^*$  is described by a regular expression, then there is a (nondeterminist) finite automata  $M$  such that  $L(M) = A$ .

Proof idea: inductively build an NFA accepting each symbol and each regular operations (union, concatenation, Kleene star) applied to smaller regular expressions.

# EQUIVALENCE PROOF: EASY DIRECTION

Constructing NFA recognizing the language  $(\{\epsilon\} \cup \{a\} \cup \{ab\})^*$ .



# EQUIV PROOF: THE OTHER DIRECTION

## EQUIVALENCE THEOREM, PART II

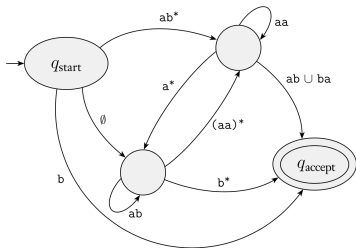
If a language  $A \subseteq \Sigma^*$  is recognized by a finite automata  $A$ , then there is a regular expression  $R$  such that  $\mathcal{L}(R) = A$ .

# EQUIV PROOF: THE OTHER DIRECTION

## GENERALIZED NONDETERMINISTIC FINITE AUTOMATA

- Generalized NFA is a NFA in which each arc carries a regular expression as a label.
- There are a unique source  $q_{start}$  as an initial state and a unique sink  $q_{accept}$  as an accept state.
- Between any other states there are arcs in both ways including a loop.

Figure 1.61, Sipser 2012.



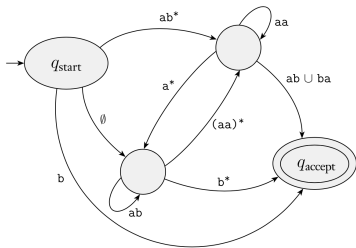
# EQUIV PROOF: THE OTHER DIRECTION

## A FORMAL DESCRIPTION OF GNFA

- Generalized NFA is a 5-tuple  $(Q, \Sigma, \delta, q_{start}, q_{accept})$ , where
- the transition function maps  $(Q \setminus q_{accept}) \times (Q \setminus q_{start})$  to the set of all regular expressions over  $\Sigma$ .

In the transition diagram, we often omit an arc which carries  $\emptyset$ .

Figure 1.61, Sipser 2012.



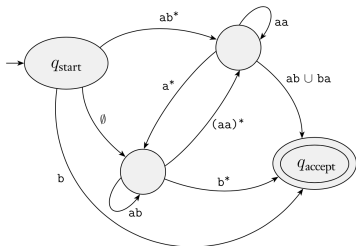
# EQUIV PROOF: THE OTHER DIRECTION

## A GNFA ACCEPTS A STRING $w$

If  $w$  can be written as  $w_1 w_2 \cdots w_\ell$  and there is a sequence of states  $r_0, \dots, r_\ell$  such that

- $r_0$  is the initial state of GNFA,
- $w_i$  is in the language described by  $\delta(r_{i-1}, r_i)$ , i.e.  $w_i \in L(\delta(r_{i-1}, r_i))$ , and
- $r_\ell$  is the accept state of GNFA.

Figure 1.61, Sipser 2012.



# EQUIV PROOF: THE OTHER DIRECTION

Proof idea.

- Initial NFA can be seen as GNFA. Reduce the number of states of GNFA inductively by eliminating a state of  $Q - \{q_{start}, q_{accept}\}$  one by one, each elimination yielding an equivalent GNFA.
- The final GNFA with two states  $q_{start}$  and  $q_{accept}$  carries a single regular expression, a desired end product.
- How to eliminate a state  $q_k$  and update the label on  $(q_i, q_j)$ :

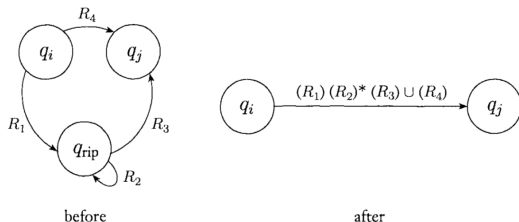


Figure 1.63, Sipser 2012.

# EXAMPLE: FROM NFA TO REGULAR EXPRESSION

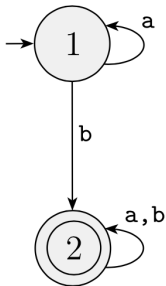


Figure 1.67 (a), Sipser 2012.

## EQUIV PROOF: THE OTHER DIRECTION

Proof by induction. Let  $G'$  be a GNFA obtained by eliminating state  $q_k$  from GNFA  $G$ . It suffices to prove that  $G$  accepts a string  $w$  if and only if  $G'$  does.

- Let  $r_0, \dots, r_\ell$  be a sequence of states appearing in the accepting computation history for  $w = w_1 \cdots w_\ell$  in  $G$ .
- If  $q_k$  does not appear in this sequence, done.

## EQUIV PROOF: THE OTHER DIRECTION

Proof by induction. Let  $G'$  be a GNFA obtained by eliminating state  $q_k$  from GNFA  $G$ . It suffices to prove that  $G$  accepts a string  $w$  if and only if  $G'$  does.

- Let  $r_0, \dots, r_\ell$  be a sequence of states appearing in the accepting computation history for  $w = w_1 \cdots w_\ell$  in  $G$ .
- If  $q_k$  does not appear in this sequence, done.
- If not, consider a maximal subsequence  $r_a, \dots, r_b$  of contiguous occurrences of  $q_k$ .



## EQUIV PROOF: THE OTHER DIRECTION

Proof by induction. Let  $G'$  be a GNFA obtained by eliminating state  $q_k$  from GNFA  $G$ . It suffices to prove that  $G$  accepts a string  $w$  if and only if  $G'$  does.

- Let  $r_0, \dots, r_\ell$  be a sequence of states appearing in the accepting computation history for  $w = w_1 \cdots w_\ell$  in  $G$ .
- If  $q_k$  does not appear in this sequence, done.
- If not, consider a maximal subsequence  $r_a, \dots, r_b$  of contiguous occurrences of  $q_k$ .
- Observe:  $w_a \in L(\delta(r_{a-1}, r_a))$ ,  $w_{i+1} \in L(\delta(r_i, r_{i+1})) = L(\delta(q_k, q_k))$  for every  $i \in [a, b-1]$ , and  $w_{b+1} \in L(\delta(r_b, r_{b+1}))$ .
- That is:  $w_a \cdots w_{b+1}$  is described by

$$\delta(r_{a-1}, r_a) \cdot \delta(q_k, q_k)^* \cdot \delta(q_k, r_{b+1})$$

which is in the union expression  $\delta'(r_{a-1}, r_{b+1})$  in the  $G'$ .

# EQUIV PROOF: THE OTHER DIRECTION

Conversely,

- Let  $r_0, \dots, r_\ell$  be a sequence of states appearing in the accepting computation history for  $w = w_1 \cdots w_\ell$  in  $G'$ .
  - Each  $w_i$  in is the language of  $\delta'(r_{i-1}, r_i)$ .
  - $\delta'(r_{i-1}, r_i) = \delta(r_{i-1}, r_i) \cup \delta(r_{i-1}, q_k) \cdot \delta(q_k, q_k)^* \cdot \delta(q_k, r_i)$ .
  - Therefore,  $w_i$  can be written as  $x_1 \cdots x_m$  such that
    - 1  $x_1 \in \delta(r_{i-1}, q_k)$
    - 2  $x_i \in \delta(q_k, q_k)$  for each  $i < m$
    - 3  $x_m \in \delta(q_k, r_i)$
- or
- 1  $x_1 \in \delta(r_{i-1}, q_k)$
  - 2  $x_2 \in \delta(q_k, r_i)$
- or  $w_i \in \delta(r_{i-1}, r_i)$ .
- Now replace the computation history in  $G'$  by..... to obtain a computation history in  $G$ .

# EXAMPLE: FROM NFA TO REGULAR EXPRESSION

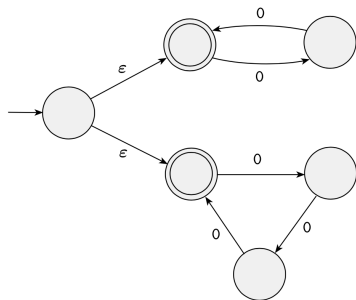


Figure 1.33, Sipser 2012.