

# REGULAR EXPRESSIONS

Recall: NFAs and DFAs both recognize the class of regular languages

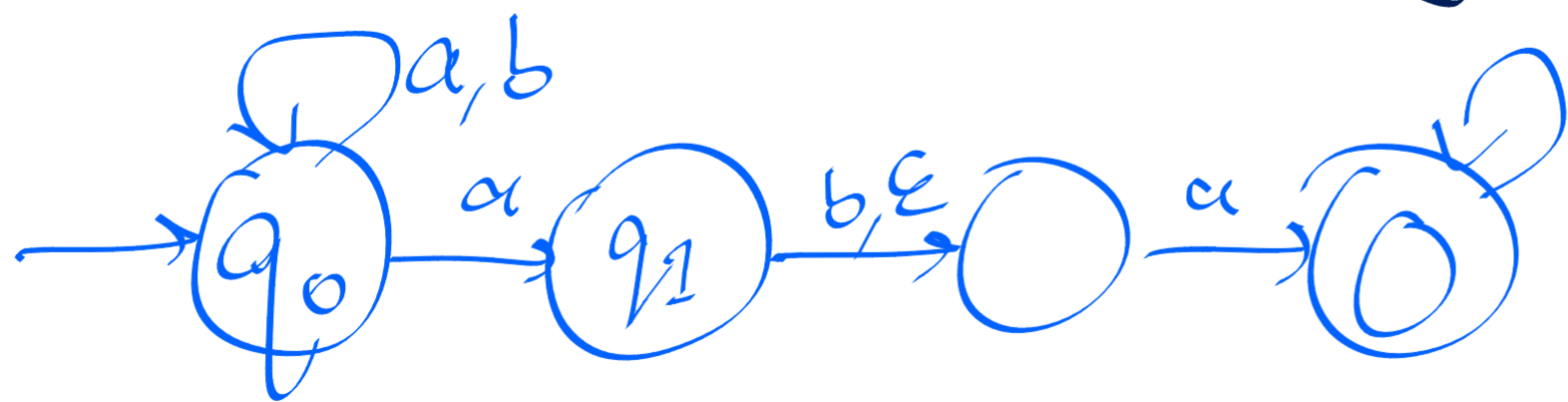
Today: Another way to specify regular languages

When we move from an NFA to a DFA, there is potentially an exponential blow-up in the number of states.

Consider, over  $\Sigma = \{a, b\}$ , the following language

$L = \{s \mid s \text{ contains } aba \text{ or } aa \text{ as a substring}\}$

What is a finite automaton recognizing  $L$ ?



But what if I have to ask my text editor to find all such words?

Is there a less informal way to do it? Regular expressions

What are the basic regular expressions?

- $a$ , for every  $a \in \Sigma$
- $\epsilon$
- $\phi$

Suppose  $\varphi$  and  $\psi$  are regular expressions. Then, so are:

- $\varphi \cdot \psi$
  - $\varphi + \psi$
  - $\varphi^*$
  - $\sim \varphi$
- every word either matches  $\varphi$  or matches  $\psi$
- every word matches some iterations of  $\varphi$
- every word avoids  $\varphi$

every word is of the form  $\omega_1 \omega_2$ , where  $\omega_1$  matches  $\varphi$ ,  $\omega_2$  matches  $\psi$

$$\mathcal{L} = \{s \mid s \text{ contains } aba \text{ or } aa \text{ as a substring}\}$$

What is a regular expression that represents  $\mathcal{L}$ ?

$$\Sigma = \{a, b\}$$

$$(a+b)^* aba (a+b)^* + (a+b)^* aa (a+b)^*$$

$$\Sigma^* aba \Sigma^* + \Sigma^* aa \Sigma^*$$

Is there a connection between this and the earlier automaton?

$\mathcal{L}$ : strings over  $\Sigma = \{a, b\}$  with an odd number of  $a$ 's

$\mathcal{L}$  = strings over  $\Sigma = \{a, b\}$  ending in  $b$  and  
not containing  $aa$

$$\sim \left( \sim (\Sigma^* b) + \Sigma^* aa \Sigma^* \right)$$

**Reg** is the class of languages expressible as regexes

We know **Reg** is the class recognized by finite-state automata

So we show that regexes are "equivalent" to FSAs.

For each regex expressing  $L$ , there is an NFA  $M$  st.  $L(M) = L$ .

We have automata accepting each of the regex patterns

Single letter, empty string, empty language

Union (+), concatenation ( $\cdot$ ), star (\*), complement ( $\sim$ )

Proof by induction On what?

For each NFA  $M$ , there exists a regex representing  $L(M)$ .

Basic idea: Keep track of the patterns tracked by each accepting path

Enumerate the states of  $M$ .  $\{1, \dots, n\}$  initial state

Maintain a set  $L_{ij}^k$  of all strings that take  $M$  from  $i$  to  $j$ , s.t. any state (other than  $i$  &  $j$ ) encountered on this path is  $\leq k$ .

$$L_{ij}^0 = \begin{cases} \{a \mid \delta(i, a) = j\} & i \neq j \\ \{a \mid \delta(i, a) = j\} \cup \{\epsilon\} & i = j \end{cases}$$

$$L_{ij}^k = L_{ij}^{k-1} \cup L_{ik}^{k-1} \cdot (L_{kk}^{k-1})^* \cdot L_{kj}^{k-1}$$

$$L(M) = \bigcup_{s \in F} L_{1s}^n$$