

\* Regexes & language of a regex

$r$  a regex, then  $L(r)$  = set of all words that match  $r$ .

\*\* Shorthand notation

In a regex, we will start using the additional symbol  $\Sigma$  to denote any single letter of  $\Sigma$ .

E.g. If  $\Sigma = \{0, 1\}$

$$0 \Sigma 1 = 0(0|1)1$$

i.e.  $\Sigma = (0|1)$

$$0 \Sigma^* 1 = 0(0|1)^* 1$$

E.g. If  $\Sigma = \{a, b, c\}$

then the symbol  $\Sigma$  in a regex is shortcut for  $(a|b|c)$ .

\* Today: The language of a regex

Table of regex constructors vs corresponding language

(1)  $r = \phi$ , then  $L(r) = \phi$

(2)  $r = \epsilon$ , then  $L(r) = \{\epsilon\}$

(3)  $r = a$  for  $a \in \Sigma$ , then  $L(r) = \{a\}$

(4)  $r = r_1 r_2$ , then  $L(r) = L(r_1) \circ L(r_2)$

(5)  $r = r_1 | r_2$  then  $L(r) = L(r_1) \cup L(r_2)$

(6)  $r = (r_1)^*$  then  $L(r) = L(r_1)^*$

(2)

\*\* Rmk: The table tells us how to go from regexes to languages.

But we don't know how to go backwards!

Q1: Given a language  $L \subseteq \Sigma^*$ , is there a regex  $r$  such that  $L(r) = L$ ?

Q2: If yes to Q1, then is there only one such regex?

We haven't yet answered Q1 and Q2.

In fact, the answer to Q2 is NO.

E.g. Let  $L = L(r)$ , then  $L = L(r | \phi)$ .  
 $= L(r) \cup L(\phi)$   
 $= L \cup \phi = L.$

$\Rightarrow$  there usually are multiple regexes that have the same language.

### \*\* Examples

$$(1) \quad \text{E.g. } r = (\Sigma \Sigma 0)^* = ((0|1)(0|1)0)^*$$

$L(r)$  contains; e.g.:  $\epsilon$ ,  $000$ ,  $110$ ,  $010$ ,  $000000$ ,  
 $010110$ ,  $000110$ , ...

= Words in which every third letter is a zero, and the length of the word is a multiple of 3

(~~auto~~ description automatically includes  $\epsilon$ )

(2)  $r = (0 \Sigma^* 0 \mid 1 \Sigma^* 1 \mid 0 \mid 1)$ .

$L(r)$  = Non-empty strings that start & end with the same letter.

(3)  $L = \{w \in \Sigma^* \mid w \text{ contains exactly } 2k \text{ "1"s, for some } k \geq 0\}$ .

Attempt #1 :  $(11)^*$  ,  $(101)^*$  x

Attempt #2 :  $(11)^* \Sigma^* \mid \Sigma^* (11)^*$  x

$(11)^* \mid (11)^* 0^* \mid 0^* (11)^* \mid (10^* 1)^* \mid 0^*$

[Does not cover, e.g. 0101]

Attempt #3 :  $(0^* 1 0^* 1 0^*)^* \mid 0^*$  ✓ [check!]

Attempt #4 :  $0^* (10^* 1)^* 0^*$  x

(000...0) (101 1001 11 100001) 000

(4)  $L = \Sigma^*$

$r = \Sigma^* = (0 \mid 1)^*$

$r_1 = \Sigma^* \mid 0$  ,  $r_2 = \Sigma^* \mid 0 \mid 110$  all

recognise the same language.

\*\* Non-example

~~$L = \{0^n 1^n\}$~~

$L = \{0^n 1^n \mid n \geq 0\}$

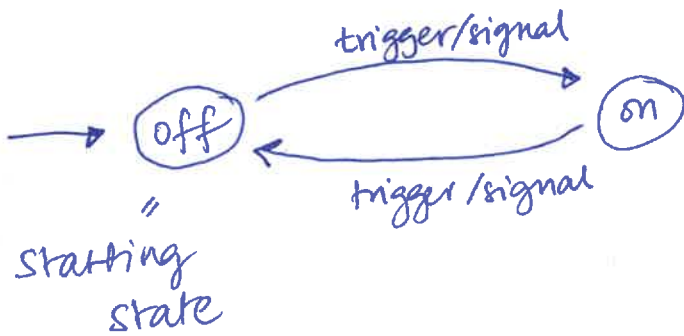
L contains  $\epsilon, 01, 0011, 000111, \dots$

Claim: There is no regular expression  $r$  such that  $L(r) = L$ !

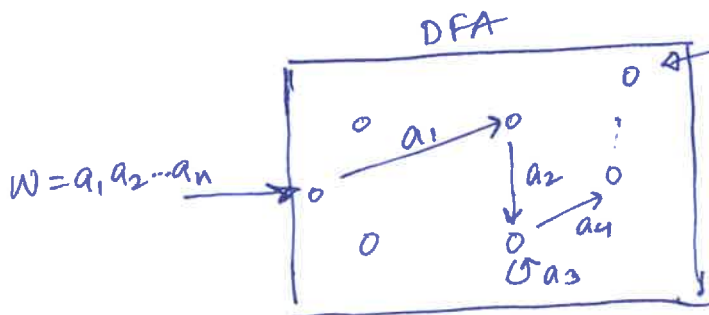
(Prove this later...)

\* Deterministic finite automata (DFAs)

Informal example: sensor tap



DFAs are machines with a finite number of "states", and "transition arrows" between states, labelled by possible inputs.



Final state will either be "accepting" or "rejecting"