

# Introduction to Language Theory and Compilation: Exercises

## Session 2: Regular expressions



# Regular expressions (RE)

- Finite automata are an equivalent formalism to regular languages (for each regular language, there exists at least one FA that recognizes it, and each FA recognizes a regular language).
- Regular expressions are another formalism defined inductively just as regular languages.
- It can be proven that regular expressions and regular languages are equivalent (see lecture notes).

# Regular expressions (RE) (ctd.)

Base cases:

RE	language
$\phi$	$\emptyset$
$\varepsilon$	$\{\varepsilon\}$
$a \quad (\forall a \in \Sigma)$	$\{a\}$

If  $p$  and  $q$  are regular expressions representing the languages  $P$  and  $Q$  respectively, then:

RE	language
$p + q$	$P \cup Q$
$pq$ (or $p \cdot q$ )	$P \cdot Q$
$p^*$	$P^*$

Extended regular expression example:  $p^+ \equiv pp^*$

- $0 + 1$  denotes the language  $\{0, 1\}$
- $a(b + c)$  denotes the language  $\{a\} \cdot \{b, c\} = \{ab, ac\}$ 
  - ... which could also be denoted by  $ab + ac$
- $x^*$  denotes  $\{x\}^*$ 
  - ... which could also be denoted by  $\varepsilon + x + xxx^*$
- A regular expression is equivalent to *one and one only* regular language, but a regular language can have more than one corresponding regular expression.

# Exercise 1

For each of the following languages (defined on the alphabet  $\Sigma = \{0, 1\}$ ), design a RE that recognizes it:

- 1 The set of strings ending with 00.
- 2 The set of strings whose 10<sup>th</sup> symbol, counted from the end of the string, is a 1.
- 3 The set of strings where each pair of zeroes is followed by a pair of ones.
- 4 The set of strings not containing 101.
- 5 The set of binary numbers divisible by 4.

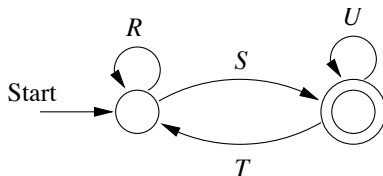
# State elimination method

Given a DFA  $M$ , we can craft a corresponding regular expression using the *state elimination* method. The general idea is to label transitions in the automaton using regular expressions, pick a final state, then remove all other states step by step to finally reach a simple automaton which can then be used to easily determine a regular expression. There are two possible cases:

- 1 The start state of  $M$  is not final ( $q_0 \notin F$ )
- 2 The start state of  $M$  is final ( $q_0 \in F$ )

# State elimination method (ctd.)

In the case where  $q_0 \notin F$ , we reach a two state automaton:

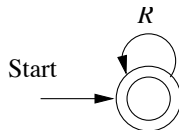


The corresponding regular expression is:

$$(R + SU^*T)^*SU^*$$

# State elimination method (ctd.)

In the case where  $q_0 \in F$ , we reach a single state automaton:



The corresponding regular expression is:

$$R^*$$



# State elimination method (ctd.)

For each final state  $q^F \in F$ , one has to build such a simple automaton to derive a regular expression  $RE(q^F)$  that expresses all possible inputs that are accepted when  $M$  stops in  $q^F$ . The actual regular expression that describes the language  $L(M)$  of the automaton  $M$  then simply becomes:

$$RE(q_1^F) + RE(q_2^F) + \dots + RE(q_k^F) \quad \text{where } \{q_1^F, \dots, q_k^F\} = F$$

# Algorithm

First, preprocess by labeling all transitions by a RE.

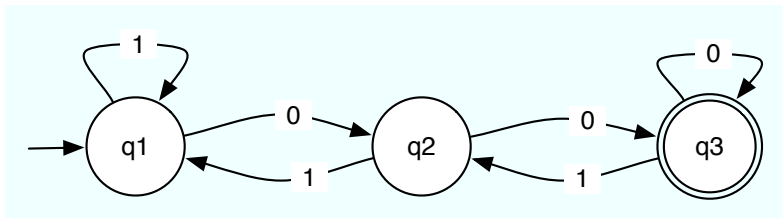
Then, for each state  $S_x$  to be eliminated, consider each transition  $(S_a, S_x)$ ,  $(S_x, S_b)$  or  $(S_x, S_x)$  with respective labels  $A$ ,  $B$  and  $X$ .

The transition  $(S_a, S_b)$  labeled by  $E$  becomes the absorbing transition  $E + (AX * B)$  and remove  $A$ ,  $B$ ,  $X$  and  $E$ .

**Note:** some transitions can not exist. In that case, does not consider the transition. For instance, if  $E = (S_a, S_b)$  cannot be generated by  $\delta$  (the transition function, see definition), then the absorption transition will be  $AX * B$ .

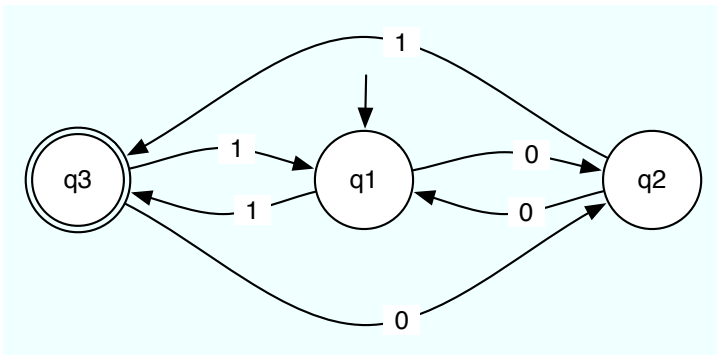
## Exercise 2.1

Design a RE accepting the same language as:



## Exercise 2.2

Design a RE accepting the same language as:



# Exercise 3

Convert the following REs into  $\varepsilon$ -NFAs:

- ①  $01^*$
- ②  $(0 + 1)01$
- ③  $00(0 + 1)^*$

# Extended regular expressions (ERE)

- Very popular on UNIX-like tools (`grep`, `find`, *etc.*)
- Grant more flexibility than traditional regular expressions
- Typically used by scanner generators such as `lex`

# ERE syntax

Expression	Accepted language
$r^*$	0 or more $rs$
$r^+$	1 or more $rs$
$r?$	0 or 1 $r$
$[abc]$	$a$ or $b$ or $c$
$[a-z]$	Any character in the interval $a \dots z$
$.$	Any character except $\backslash n$
$[\^s]$	Any character but those in $s$
$r\{m,n\}$	Between $m$ and $n$ occurrences of $r$
$r_1 r_2$	The concatenation of $r_1$ and $r_2$

# ERE syntax (ctd.)

Expression	Accepted language
$r1 \mid r2$	$r1$ or $r2$
$(r)$	$r$
$\wedge r$	$r$ if it starts a line
$r\$$	$r$ if it ends a line
" $s$ "	The string $s$
$\backslash c$	The character $c$
$r1(?:r2)$	$r1$ when it's followed by $r2$



# Examples

Expression	Accepted language
<code>[a-zA-Z]</code>	Any letter (upper or lower case)
<code>[0-9]</code>	Any digit
<code>a[^A-Za-z]b</code>	An a followed by a non-alphabetical character and a b
<code>^Silly</code>	Silly if it starts a line
<code>[a-zA-Z]([a-zA-Z]   [0-9])*</code>	An identifier in the Pascal language

# Exercise 4

- ① Give an extended regular expression (ERE) that targets any sequence of 5 characters, including the newline character `\n`
- ② Give an ERE that targets any string starting with an arbitrary number of `\` followed by any number of `*`
- ③ UNIX-like shells (such as `bash`) allow the user to write *batch* files in which comments can be added. A line is defined to be a comment if it starts with a `#` sign. What ERE accepts such comments?
- ④ Design an ERE that accepts numbers in scientific notation. Such a number must contain at least one digit and has two optional parts:
  - A "decimal" part : a dot followed by a sequence of digits
  - An "exponential" part: an `E` followed by an integer that may be prefixed by `+` or `-`
  - Examples : 42, 66.4E-5, 8E17, ...

# Exercise 4

- 5 Design an ERE that accepts "correct" phrases that fulfill the following criteria:
- The first word must start with a capital letter
  - The phrase must end with a full stop .
  - The phrase must be made of one or more words (made of the characters a . . . z and A . . . Z) separated by a single space
  - There cannot be two phrases on the same line.

Punctuation signs other than a full stop are not allowed.

- 6 Craft an ERE that accepts old school DOS-style filenames (8 characters in a . . . z, A . . . Z and \_) whose extension is .ext and that begin with the string abcde. We ask that the ERE only accept the filename *without the extension!*
- Example: on abcdeLOL.ext, the ERE must accept abcdeLOL