

University of Camerino
School of Science and Technologies



Formal Languages and Compilers
Project Paper

Prof. Andrea Polini

The project to be developed asks students to derive a compiler for a regular expression grammar that generate a program able to distinguish words belonging to the language generated by the defined Regular Expression (RExp). In order to derive the program students can refer to the regular automaton defined by Thompson as corresponding acceptance automaton for a RExp. Such rules need to be codified in a programming language (e.g. Java). The result will be a program in the selected language that then will be compiled to provide the needed functionality.

For instance the program generated by the compiler, given the following generic RE on the alphabet $\Sigma = \{a, b, c, \dots, z, A, B, C, \dots, Z, 0, 1, 2, \dots, 9, ., -, _\}$ should reply “OK” or “KO” as specified in the following:

RE	word	Result
a+b	a	OK
	b	OK
	ac	KO
a*b	aab	OK
	b	OK
	aaaab	OK
	abb	KO

In order to derive the compiler the following grammar for RE should be considered:

<i>RE</i>	::=	<i>union</i> <i>simpleRE</i>
<i>union</i>	::=	<i>simpleRE</i> + <i>RE</i>
<i>simpleRE</i>	::=	<i>concatenation</i> <i>basicRE</i>
<i>concatenation</i>	::=	<i>basicRE</i> <i>simpleRE</i>
<i>basicRE</i>	::=	<i>group</i> <i>any</i> <i>char</i>
<i>group</i>	::=	(<i>RE</i>) (<i>RE</i>)* (<i>RE</i>)+
<i>any</i>	::=	?
<i>char</i>	::=	<i>a</i> <i>b</i> <i>c</i> ... <i>z</i> <i>A</i> <i>B</i> <i>C</i> ... <i>Z</i> <i>0</i> <i>1</i> <i>2</i> ... <i>9</i> <i>.</i> <i>-</i> <i>_</i>

Suggestion: before starting the ANTLR project can be useful to write a simple program able to behave according to a statically defined regular automaton with ϵ -transitions. For such an automaton actions are driven by a transition table that need to be suitably represented. In particular it will be necessary to define a data structure that can be used to represent the transition table of a NDFSA.