

# Compiler Construction 2017

## — Exercise Sheet 2 —

Hand in until May 23rd before the exercise class.

### Exercise 1

(4 Points)

Let  $\alpha_1 = \text{if}$  and  $\alpha_2 = \Sigma(\Sigma|N)^*$ , where  $\Sigma := (a|\dots|z|A|\dots|Z)$ ,  $N := (0|1|\dots|9)$ .

- Construct DFAs  $\mathfrak{A}_i$  for  $\alpha_i$  such that  $\mathcal{L}(\mathfrak{A}_i) = \llbracket \alpha_i \rrbracket$ .
- Construct the product automaton  $\mathfrak{A} = \mathfrak{A}_1 \otimes \mathfrak{A}_2$ .
- Determine the *first match* partitioning of the set of final states in  $\mathfrak{A}$ .  
(The regular expressions are ordered  $(\alpha_1, \alpha_2)$ .)
- Determine the set of reachable and productive states in  $\mathfrak{A}$ .
- Compute the run of the corresponding backtracking DFA for input `ifdef`. Provide the run by giving the corresponding configurations.

### Exercise 2

(6 Points)

On July 20, 2016 the popular website `stackoverflow.com` experienced an unexpected outage due to a malformed post that caused its backtracking Regex engine to consume extraordinary high CPU time.<sup>1</sup> This is an example of a *regular expression denial-of-service* (ReDoS) attack: An attacker tries to paralyze an application by feeding it with input strings that exhibit high computational complexity.

In this exercise, we have a closer look at such "attack strings". To this end, we execute the FLM analysis based on the NFA method from the lecture (lecture 4, slide 10).

Furthermore, we call an NFA  $\mathfrak{A}$  *vulnerable* if and only if the worst-case complexity of the above FLM analysis is exponential in the length of the input string.

**Hint:** In this exercise, we assume that our FLM analysis algorithm is very naive and explores each possible run of an NFA on a word individually.

- Provide an example of a vulnerable NFA  $\mathfrak{A}$  and, for each  $n \geq 1$ , an input string  $w_n$  of length  $n$  such that the runtime of the FLM analysis on  $\mathfrak{A}$  and  $w_n$  is exponential in  $n$ .
- Show that an NFA  $\mathfrak{A} = (Q, \Omega, \delta, q_0, F)$  is vulnerable if there exists a state  $q$  and two distinct paths<sup>2</sup>  $\pi_1, \pi_2$  such that
  - both  $\pi_1$  and  $\pi_2$  start and end at state  $q$ ,
  - $\text{labels}(\pi_1) = \text{labels}(\pi_2)$ ,
  - there is a path  $\pi_p$  from initial state  $q_0$  to  $q$ , and
  - there is a path  $\pi_s$  from  $q$  to a state  $q_r \notin F$ .
- Sketch a procedure that takes an NFA  $\mathfrak{A}$  and one of its states  $q$  as an input and checks whether  $\mathfrak{A}$  is vulnerable for the given state  $q$ .

<sup>1</sup>Details are provided at <http://stackstatus.net/post/147710624694/outage-postmortem-july-20-2016>

<sup>2</sup>A path is a sequence of transitions  $\pi = (q_1, \ell_1, q_2) \dots (q_{m-1}, \ell_{m-1}, q_m)$  such that  $q_i \in Q$ ,  $\ell_i \in \Omega$ , and  $q_{i+1} \in \delta(q_i, \ell_i)$ . Moreover, we set  $\text{labels}(\pi) = \ell_1 \dots \ell_{m-1}$ .