LEX5: Regexps to NFA

Lexical Analysis

CMPT 379: Compilers Instructor: Anoop Sarkar anoopsarkar.github.io/compilers-class

Building a Lexical Analyzer

- Token \Rightarrow Pattern
- Pattern \Rightarrow Regular Expression
- Regular Expression \Rightarrow NFA
- NFA \Rightarrow DFA
- DFA \Rightarrow Table-driven implementation of DFA

- Converts regexps to equivalent NFA
- Six simple rules
 - Empty language
 - Symbols (Σ)
 - Empty String (ε)
 - Alternation $(r_1 \text{ or } r_2)$
 - Concatenation (*r*₁ followed by *r*₂)
 - Repetition (r_1^*)

Used by Ken Thompson for pattern-based search in text editor QED (1968)

- For the empty language φ
- (optionally include a *sinkhole* state)



• For each symbol x of the alphabet, there is a NFA that accepts it



• There is an NFA that accepts only ε



• Given 2 NFAs r_1 , r_2 , there is a NFA that accepts $r_1 | r_2$





• Given 2 NFAs r_1 , r_2 , there is a NFA that accepts $r_1 | r_2$



• Given 2 NFAs r_1 , r_2 , there is a NFA that accepts r_1r_2





• Given 2 NFAs r_1 , r_2 , there is a NFA that accepts r_1r_2



• Given 2 NFAs r_1 , r_2 , there is a NFA that accepts r_1r_2



• Given an NFA for r, there is an NFA that accepts r*



Given an NFA for r, there is an NFA that accepts r*



Example

- Set of all binary strings that are divisible by four (include 0 in this set)
- Defined by the regexp: ((0|1)*00) | 0
- Apply Thompson's Rules to create an NFA





0|1 ((0|1)*00) | 0

ε ε (0|1)* ((0|1)*00) | 0



(0|1)*00 ((0|1)*00) | 0





(0|1)*00 ((0|1)*00) | 0























































n7= nfa(n5, n6, .)

40

Q: Use Thompson's construction to build an NFA for (0|1)(0|1)*

Thompson's construction -

(a(a|b))c

