

CS 131 Compilers: Discussion 2: (Non)Deterministic Finite Automata

杨易为 吴凌云 樊雨鑫

{yangyw,wuly2,fanyx}@shanghaitech.edu.cn

2022 年 2 月 23 日

1 DFA and NFA

1.1 Introduction

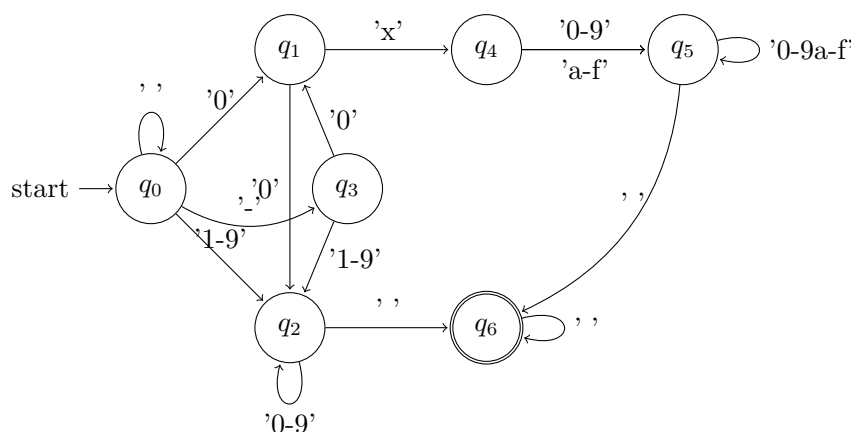
Finite state automata (FSA) are abstract machines that feature states and guarded transitions from one state to another. An FSA can only be in one state at a time, and its total number of states is finite. The machine takes transitions in response to inputs it receives sequentially; if an input matches the guard of a transition that departs from the current state that transition is said to be enabled; only enabled transitions can be taken. FSA that have an accepting state provide the machinery to determine whether an input string is in a regular language. If no transitions are enabled by a given input then the entire input string gets rejected. If all inputs are processed and the FSA is in an accepting state then the entire input string is accepted. Deterministic finite state automata (DFA) can only have one enabled transition at a time while a non-deterministic finite state automata (NFA) can have multiple.

Before NLP, we utilize state machine to accept languages. Two different perspectives can be considered, namely symbolic and random. Chomsky's formal language theory embodies a symbolic approach. Based on this view, a language contains a sequence of symbols that must follow the syntactic rules of its generative grammar. This view reduces the structure of a language to a set of well-defined rules that allow for the decomposition of each sentence and word into structural components.

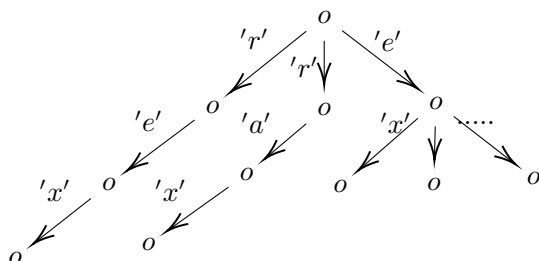
1.2 Examples of accepting x86 Assembly

https://github.com/yangminz/bcst_csapp

There are three main types of terminators for assembly instructions: Immediately and Scale for numbers, Reg and Op for characters, and other symbols, including spaces (' '), commas (','), and parentheses ('(',')'). For numeric terminators, we can use Deterministic Finite Automaton (DFA) to determine and parse them according to Regular Grammar, see below Figure. It should be noted that the numbers in assembly instructions include positive and negative integers in hexadecimal and decimal, and We need to check whether the value overflows at each state transfer step.



For character type terminators, i.e. operands and registers, we can take a simple string matching (Matching). Here, we have a variety of data structures to employ. The most common approach is to use an array or list that stores all possible string values, e.g., {"rax", "eax"} and {"mov", "ld", "sd"}, and iterate through this array and match each element. With such a data structure, the time complexity of finding and matching is $O(LE[strlen])$, where $E[strlen]$ is the average length of the string. Another strategy is to build a dictionary tree (Trie) or a prefix tree for the above array of strings.



By using prefix trees, searching and matching can be done without traversing the entire array, and the time complexity is reduced to $O(E[\textit{strlen}])$. In the implementation, a prefix tree can be created for all instruction operator names and register names before executing the CPU instruction cycle, so that the matching time can be reduced during the execution of the instruction cycle. Each node of the prefix tree usually includes a `char→(void *)` mapped data structure, where `char` is the character and `(void *)` is the address of the subtree. As you can see, the prefix tree is very similar to the DFA, except that the prefix tree automatically builds the state transfer (so it is not streamlined enough), while the DFA requires us to calculate the state transfer in advance.

1.3 Examples of accepting yaml

We cannot do it, because of indent can only be accepted with stack machine.

```

1 // Identity definitions
2 <whitespace> ::= " " | "\t"
3 <in-line-whitespace> ::= <whitespace> | <whitespace> <in-line-whitespace>
4 <digit> ::= "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9"
5 <integer-number> ::= <digit> | <digit> <integer-number>
6 <end-of-line> ::= "\r\n" | "\r" | "\n"
7 <non-whitespace> ::=
8 <non-whitespace-word> ::= <non-whitespace> | <non-whitespace> <non-whitespace-word>
9
10 // Basic document structure
11 <yaml-stream> ::= <yaml-document-list>
12 <yaml-document-list> ::= <yaml-document> | <yaml-document> <yaml-document-list>
13 <yaml-document> ::= <yaml-directive-list> <yaml-directives-end> <yaml-content-list>

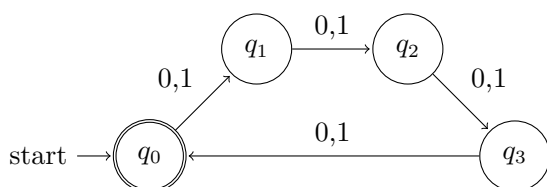
```

```

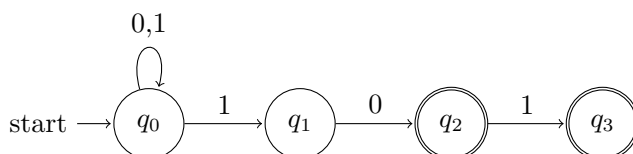
14      <yaml-document-end>
15 <yaml-directive-list> ::= <yaml-directive> <yaml-directive-list> | ""
16 <yaml-directive> ::= <yaml-directive-yaml> | <yaml-directive-tag>
17 <yaml-directive-yaml> ::= "%YAML" <in-line-whitespace> <integer-number> "." <integer-number>
18      <end-of-line> | ""
19 <yaml-directive-tag> ::= "%TAG" <in-line-whitespace> <yaml-tag-handle> <in-line-whitespace>
20      <yaml-tag-prefix> <end-of-line> <yaml-directive-tag> | ""
21 <yaml-tag-handle> ::= "!" | "!!" | "!" <non-whitespace-word> "!"
22 <yaml-tag-prefix> ::= "!" <non-whitespace-word> | <non-whitespace-word>
23 <yaml-directives-end> ::= "---" | ""
24 <yaml-content-list> ::= <yaml-content> <yaml-content-list> | <yaml-content> <end-of-line>
25 <yaml-document-end> ::= "... " | ""
26
27 // Internals of an individual document
28 <yaml-content> ::=

```

1.4 What language is accepted by the following DFA?



1.5 What language is accepted by the following NFA?



1.6 Construct a DFA that accepts binary strings of any length, except 3.

1.7 Construct the NFA that accepts

1.7.1 $x(zy?(yz)^*)$

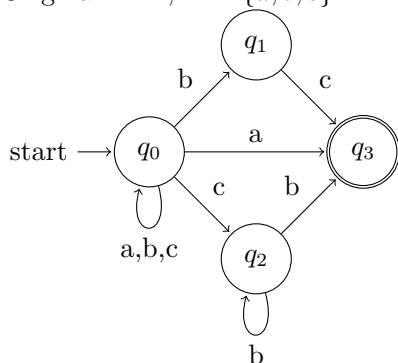
Answer:

1.7.2 $(10)^* 1^* 01^* 0$

Answer:

1.8 Minimized the NFA

1. Original NFA, $\Sigma = \{a, b, c\}$:



Answer:

2 Pumping Lemma

Pumping Lemma: For any DFA (or NFA or regular expression) that accepts an infinite number of strings, there is some minimum length, M , such that any string with length greater than or equal to M that the machine accepts must have the form uxv , where u , x , and v are strings, x is not empty, the length of ux is $\leq M$, and the machine accepts all strings of the form $ux^n v$, for all $n \geq 0$.

2.1 Proof of not regular

Let $A = \{1^j z \mid z \in \{0, 1\}^* \text{ and } z \text{ contains at most } j \text{ 1's, for any } j \geq 1\}$. Prove, by the pumping lemma, that A is not regular.

Answer: