



中山大學
SUN YAT-SEN UNIVERSITY

计算机学院（软件学院）

SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

Compilation Principle

编译原理

第2讲：词法分析(2)

张献伟

xianweiz.github.io

DCS290, 3/5/2024



中山大學
SUN YAT-SEN UNIVERSITY



Review Questions

- Q1: input and output of lexical analysis?
Input: source code/char stream, output: tokens
- Q2: lexical analysis of “while i>=1)”?
(keyword, ‘while’), (id, ‘i’), (sym, ‘>=’), (num, ‘1’), (sym, ‘)’)
- Q3: $\Sigma = \{a, b\}$, $L_1 = \{aa\}$, $L_2 = \{bbb\}$. What are $L_1 \mid L_2$ and L_1L_2 ?
 $L_3 = L_1 \mid L_2 = \{aa\} \mid \{bbb\} = \{aa, bbb\}$, $L_4 = L_1L_2 = \{aabbbb\}$
- Q4: L_3^2 ?
 $L_3^2 = L_3L_3 = \{aa, bbb\}\{aa, bbb\} = \{aaaa, aabbbb, bbbbaa, bbbbbbb\}$
- Q5: describe the meaning of $L_1^* \mid L_2^*$?
A language composed of ‘a’s and ‘b’s of length $2N$ and $3N$, respectively, including ϵ
- Q6: RE of identifiers in C language?
 $(_letter)(_letter \mid digit)^*$

Compound REs[组合表达式]

- **Compound**

- Large REs built from smaller ones

- Suppose r and s are REs denoting languages $L(r)$ and $L(s)$

- (r) is a RE denoting the language $L(r)$

- We can add additional $()$ around expressions without changing the language they denote

- $(r)|(s)$ is a RE denoting the language $L(r) \cup L(s)$

- $(r)(s)$ is a RE denoting the language $L(r)L(s)$

- $(r)^*$ is a RE denoting the language $(L(r))^*$

- REs often contain unnecessary $()$, which could be dropped

- $(\mathbf{A}) \equiv \mathbf{A}$: A is a RE

- $(a)|((b)^*(c)) \equiv a|b^*c$

Operator Precedence[优先级]

- RE operator precedence

- (A)

- A^*

- AB

- $A|B$

- Example: $ab^*c|d$

- $a(b^*)c|d$

- $(a(b^*))c|d$

- $((a(b^*))c)|d$

Common REs[常用表达]

- **At least one:** $A^+ \equiv AA^*$
- **Option:** $A? \equiv A \mid \varepsilon$
- **Characters:** $[a_1a_2\dots a_n] \equiv a_1 \mid a_2 \mid \dots \mid a_n$
- **Range:** 'a' + 'b' + ... + 'z' $\equiv [a-z]$
- **Excluded range:** complement of $[a-z] \equiv [^a-z]$

RE Examples

Regular Expression	Explanation
a^*	0 or more a's (ϵ , a, aa, aaa, aaaa, ...)
a^+	1 or more a's (a, aa, aaa, aaaa, ...)
$(a b)(a b)$	(aa, ab, ba, bb)
$(a b)^*$	all strings of a's and b's (including ϵ)
$(aa ab ba bb)^*$	all strings of a's and b's of even length
$[a-zA-Z]$	shorthand for "a b ...z A B ... Z"
$[0-9]$	shorthand for "0 1 2 ... 9"
$0([0-9])^*0$	numbers that start and end with 0
$1^*(0 \epsilon)1^*$	binary strings that contain at most one zero
$(0 1)^*00(0 1)^*$	all binary strings that contain '00' as substring

- Q: are $(a|b)^*$ and $(a^*b^*)^*$ equivalent?

Different REs of the Same Language

- $(a|b)^* = ?$

- $(L(a|b))^* = (L(a) \cup L(b))^* = (\{a\} \cup \{b\})^* = \{a, b\}^*$

- $= \{a, b\}^0 + \{a, b\}^1 + \{a, b\}^2 + \dots$

- $= \{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$

- $(a^*b^*)^* = ?$

- $(L(a^*b^*))^* = (L(a^*)L(b^*))^*$

- $= L(\{\epsilon, a, aa, \dots\}\{\epsilon, b, bb, \dots\})^*$

- $= L(\{\epsilon, a, b, aa, ab, bb, \dots\})^*$

- $= \epsilon + \{\epsilon, a, b, aa, ab, bb, \dots\} + \{\epsilon, a, b, aa, ab, bb, \dots\}^2 + \{\epsilon, a, b, aa, ab, bb, \dots\}^3 + \dots$

More Examples

- Keywords: 'if' or 'else' or 'then' or 'for' ...
 - RE = 'i''f' + 'e''l''s''e' + ... = 'if' + 'else' + 'then' + ...
- Numbers: a non-empty string of digits
 - digit = '0' + '1' + '2' + '3' + '4' + '5' + '6' + '7' + '8' + '9'
 - integer = digit digit*
 - Q: is '000' an integer?
- Identifier: strings of letters or digits, starting with a letter
 - letter = 'a' + 'b' + ... 'z' + 'A' + 'B' + ... + 'Z' = [a-zA-Z]
 - RE = letter(letter + digit)*
 - Q: is the RE valid for identifiers in C?
- Whitespace: a non-empty sequence of blanks, newline, tabs
 - ('' + '\n' + '\t')+

'+' == '|'

REs in Programming Language

Symbol	Meaning		
<code>\d</code>	Any decimal digit, i.e. [0-9]		
<code>\D</code>	Any non-digit char, i.e., [^0-9]		
<code>\s</code>	Any whitespace char, i.e., [\t\n\r\f\v]		
<code>\S</code>	Any non-whitespace char, i.e., [^ \t\n\r\f\v]		
<code>\w</code>	Any alphanumeric char, i.e., [a-zA-Z0-9_]		
<code>\W</code>	Any non-alphanumeric char, i.e., [^a-zA-Z0-9_]		
<code>.</code>	Any char	<code>\.</code>	Matching “.”
<code>[a-f]</code>	Char range	<code>[^a-f]</code>	Exclude range
<code>^</code>	Matching string start	<code>\$</code>	Matching string end
<code>(...)</code>	Capture matches		

<https://docs.python.org/3/howto/regex.html>

Lexical Specification of a Language

- **S0**: write a regex for the lexemes of each token class
 - Numbers = digit^+
 - Keywords = 'if' + 'else' + ...
 - Identifiers = $\text{letter}(\text{letter} + \text{digit})^*$
- **S1**: construct R , matching all lexemes for all tokens
 - $R = \text{numbers} + \text{keywords} + \text{identifiers} + \dots = R_1 + R_2 + R_3 + \dots$
- **S2**: let input be $x_1 \dots x_n$, for $1 \leq i \leq n$, check $x_1 \dots x_i \in L(R)$
- **S3**: if successful, then we know $x_1 \dots x_i \in L(R_j)$ for some j
 - E.g., an identifier or a number ...
- **S4**: remove $x_1 \dots x_i$ from input and go to step S2

Lexical Spec. of a Language(cont.)

- How much input is used?
 - $x_1 \dots x_i \in L(R), x_1 \dots x_j \in L(R), i \neq j$
 - Which one do we want? (e.g., '==' or '=')
 - Maximal match: always choose the longer one[最长匹配]
- Which token is used if more than one matches?
 - $x_1 \dots x_i \in L(R)$ where $R = R_1 + R_2 + \dots + R_n$
 - $x_1 \dots x_i \in L(R_m), x_1 \dots x_i \in L(R_n), m \neq n$
 - E.g., keywords = 'if', identifier = letter(letter+digit)*
 - Keyword has higher priority
 - Rule of thumb: choose the one listed first[次序]
- What if no rule matches?
 - $x_1 \dots x_i \notin L(R) \rightarrow$ Error

Summary: RE

- We have learnt how to specify tokens for lexical analysis[定义]
 - Regular expressions
 - Concise notations for the string patterns
- Used in lexical analysis with some extensions[适度扩展]
 - To resolve ambiguities
 - To handle errors
- RE is only a language specification[只是定义了语言]
 - An implementation is still needed
 - Next: to construct a token recognizer for languages given by regular expressions – by using **finite automata**[有穷自动机]

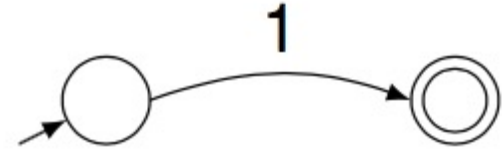
Impl. of Lexical Analyzer[实现]

- How do we go from specification to implementation?
 - RE \rightarrow finite automata (FA)
- **Solution 1:** to implement using a tool — Lex (for C), Flex (for C++), Jlex (for java)
 - Programmer specifies tokens using REs
 - The tool generates the source code from the given REs
 - The Lex tool essentially does the following translation: REs (specification) \Rightarrow FAs (implementation)
- **Solution 2:** to write the code yourself
 - More freedom; even tokens not expressible through REs
 - But difficult to verify; not self-documenting; not portable; usually not efficient

Transition Diagram[转换图]

- REs \rightarrow transition diagrams

- By hand
- Automatic



- Node[节点]: state

- Each state represents a condition that may occur in the process
- Initial state (Start): only one, circle marked with ‘start \rightarrow ’
- Final state (Accepting): may have multiple, double circle

- Edge[边]: directed, labeled with symbol(s)

- From one state to another on the input

Finite Automata[有穷自动机]

- **Regular Expression** = **specification**[正则表达是定义]
- **Finite Automata** = **implementation**[自动机是实现]
- Automaton (pl. automata): a machine or program
- **Finite automaton (FA)**: a program with a finite number of states
- Finite Automata are similar to transition diagrams
 - They have states and labelled edges
 - There are one unique start state and one or more than one final states

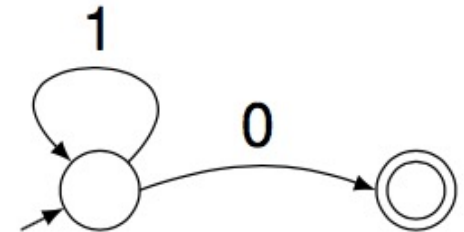
FA: Language

- An FA is a program for classifying strings (accept, reject)
 - In other words, a program for recognizing a language
 - The Lex tool essentially does the following translation: REs (specification) \Rightarrow FAs (implementation)
 - For a given string 'x', if there is transition sequence for 'x' to move from start state to certain accepting state, then we say 'x' is accepted by the FA
 - Otherwise, rejected
- Language of FA = set of strings accepted by that FA
 - $L(\text{FA}) \equiv L(\text{RE})$

Example

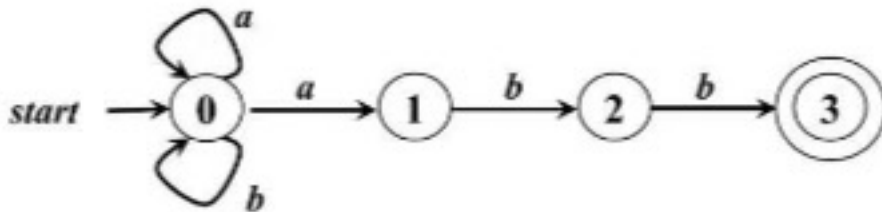
- Are the following strings acceptable?

- 0 ✓
- 1 ✗
- 11110 ✓
- 11101 ✗
- 11100 ✗
- 1111110 ✓



- What language does the state graph recognize? $\Sigma = \{0, 1\}$

Any number of '1's followed by a single 0



L(FA): all strings of $\Sigma\{a, b\}$, ending with 'abb'

L(RE) = $(a|b)^*abb$

DFA and NFA

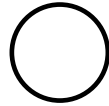
- Deterministic Finite Automata (**DFA**): the machine can exist in only one state at any given time[确定]
 - One transition per input per state
 - No ϵ -moves
 - Takes only one path through the state graph
- Nondeterministic Finite Automata (**NFA**): the machine can exist in multiple states at the same time[非确定]
 - Can have multiple transitions for one input in a given state
 - Can have ϵ -moves
 - Can choose which path to take
 - An NFA accepts if some of these paths lead to accepting state at the end of input

State Graph

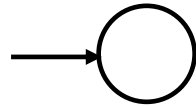
- 5 components $(\Sigma, S, n, F, \delta)$

- An input alphabet Σ

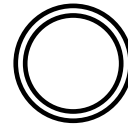
- A set of states S



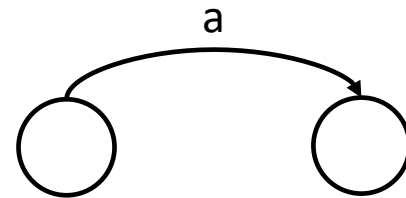
- A start state $n \in S$



- A set of accepting states $F \subseteq S$



- A set of transitions $\delta: S_a \xrightarrow{\text{input}} S_b$

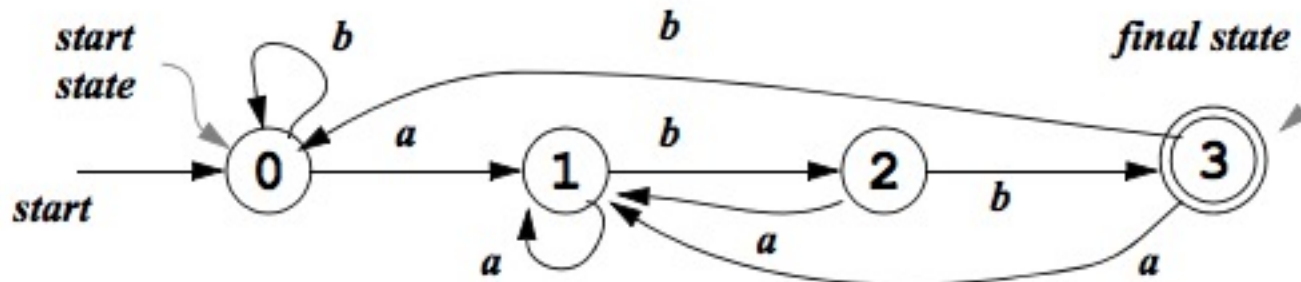


Example: DFA

- There is **only one** possible sequence of moves --- either lead to a final state and accept or the input string is rejected

– Input string: **aabb**

– Successful sequence:



A DFA accepts $(a|b)^*abb$

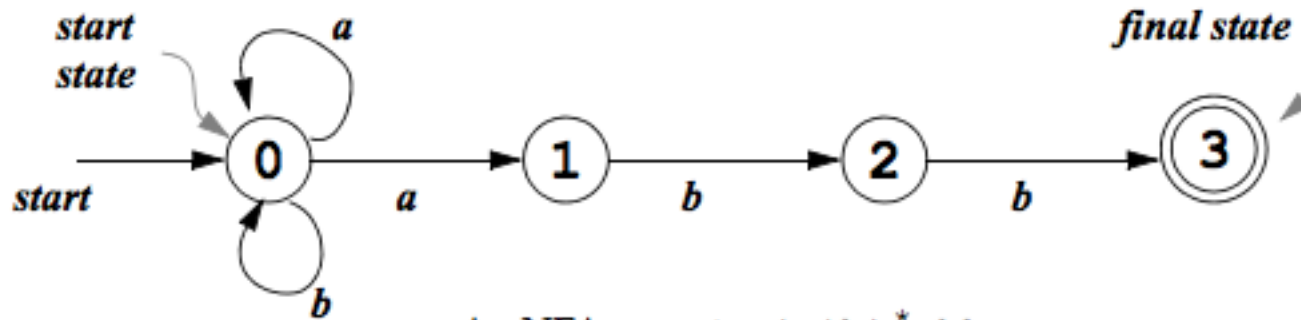
Example: NFA

- There are **many possible** moves: to accept a string, we only need one sequence of moves that lead to a final state

– Input string: **aabb**

– Successful sequence: $0 \xrightarrow{a} 0 \xrightarrow{a} 1 \xrightarrow{b} 2 \xrightarrow{b} 3$

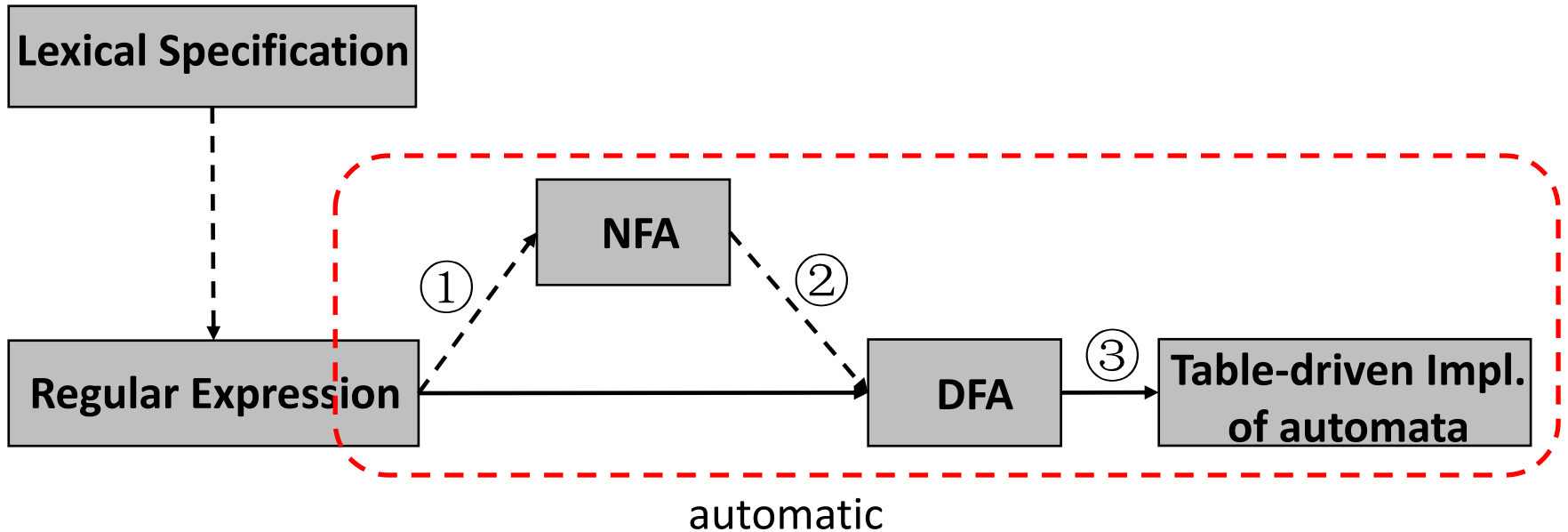
– Unsuccessful sequence: $0 \xrightarrow{a} 0 \xrightarrow{a} 0 \xrightarrow{b} 0 \xrightarrow{b} 0$



An NFA accepts $(a|b)^*abb$

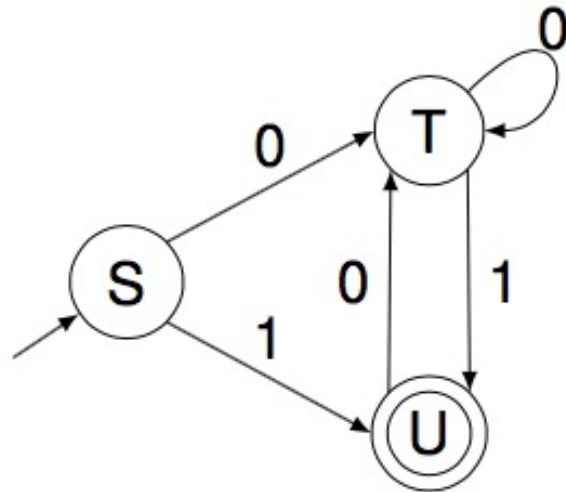
Conversion Flow[转换流程]

- Outline: RE \rightarrow NFA \rightarrow DFA \rightarrow Table-drive Implementation
 - ③ Converting DFAs to table-driven implementations
 - ① Converting REs to NFAs
 - ② Converting NFAs to DFAs



DFA \rightarrow Table

- FA can also be represented using transition table



alphabet \rightarrow

state \downarrow

	0	1
S	T	U
T	T	U
U	T	x

Table-driven Code:

```
DFA() {
    state = "S";
    while (!done) {
        ch = fetch_input();
        state = Table[state][ch];
        if (state == "x")
            print("reject");
    }
    if (state  $\in$  F)
        printf("accept");
    else
        printf("reject");
}
```

Q: which is/are accepted?

111

000

001

More on Table

- Implementation is efficient[表格是一种高效实现]
 - Table can be automatically generated
 - Need finite memory $O(S \times \Sigma)$
 - Size of transition table
 - Need finite time $O(\text{input length})$
 - Number of state transitions
- Pros and cons of table[表格实现的优劣]
 - **Pro**: can easily find the transitions on a given state and input
 - **Con**: takes a lot of space, when the input alphabet is large, yet most states do not have any moves on most of the input symbols