



中山大學
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

Compilation Principle 编译原理

第4讲：词法分析(4)

张献伟

xianweiz.github.io

DCS290, 3/3/2022



中山大學
SUN YAT-SEN UNIVERSITY



Review Questions

- Q1: how does RE, NFA, DFA relate to each other?

$$L(RE) \equiv L(NFA) \equiv L(DFA)$$

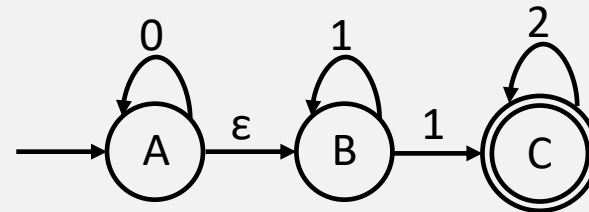
- Q2: NFA \rightarrow DFA, time and space complexity

$$\text{Time: NFA-}O(|X| * N^2) \text{ DFA-}O(|X|)$$

$$\text{Space: NFA-}O(N), \text{ DFA-}O(2^N)$$

- Q3: RE of the NFA?

$$0^*1+2^*$$

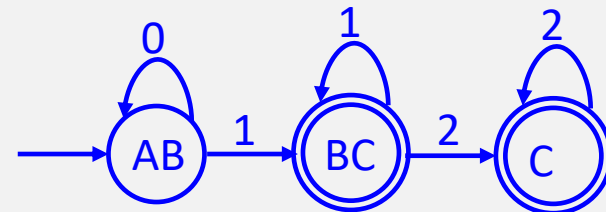


- Q4: start state of the equivalent DFA? $\epsilon\text{-closure}(A) = \{A, B\}$

	0	1	2
AB	AB	BC	
BC		BC	C
C			C

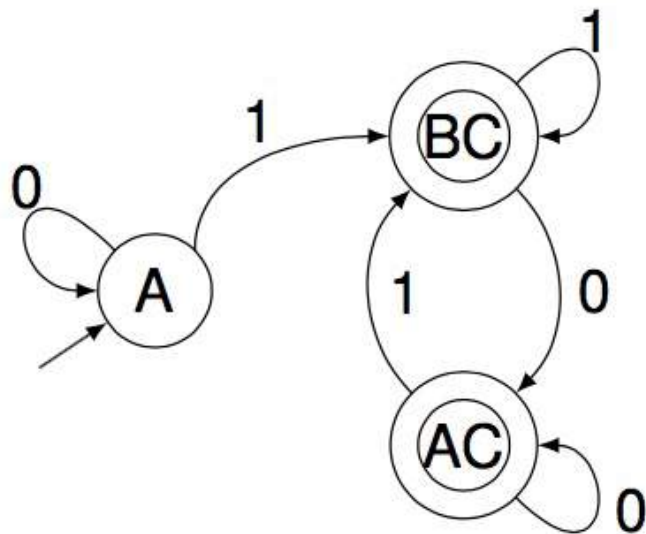
$$\epsilon\text{-closure}(\text{move}(\{AB\}, 0)) = \epsilon\text{-closure}(\{A\}) \Rightarrow \{A, B\}$$

$$\epsilon\text{-closure}(\text{move}(\{AB\}, 1)) = \epsilon\text{-closure}(\{B, C\}) \Rightarrow \{B, C\}$$



NFA \rightarrow DFA: Minimization[最小化]

- Any DFA can be converted to its minimum-state equivalent DFA
 - Discover sets of equivalent states
 - Represent each such set with just one state
- Two states are equivalent if and only if[等价]:
 - $\forall \alpha \in \Sigma$, transitions on α lead to equivalent states
 - α -transitions to distinct sets \Rightarrow states must be in distinct sets



Initial: {A}, {BC, AC}

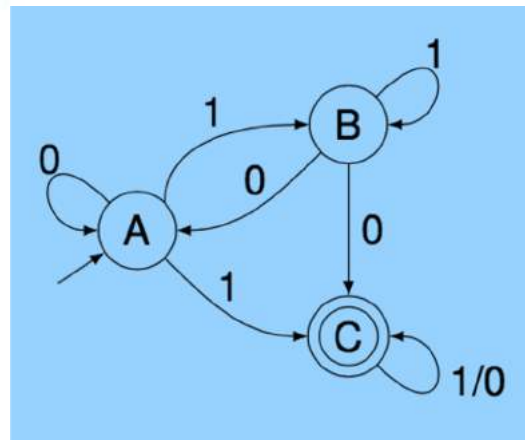
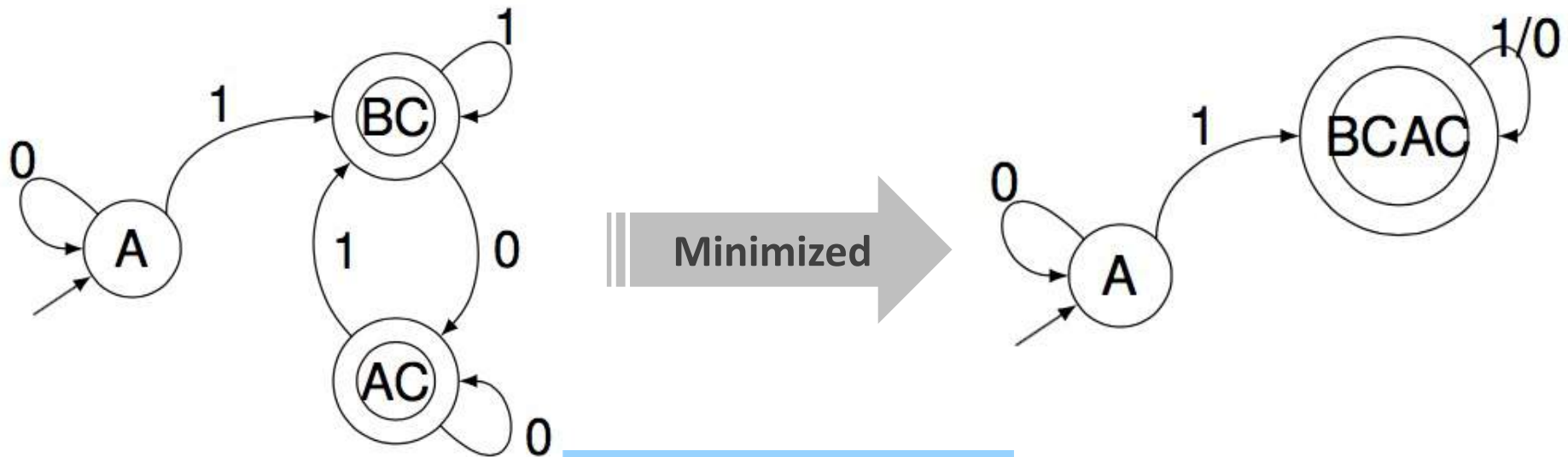
For {BC, AC} **Initial sets:** {non-accepting states}, {accepting states}

- BC on '0' \rightarrow AC, AC on '0' \rightarrow AC
- BC on '1' \rightarrow BC, AC on '1' \rightarrow BC
- No way to distinguish BC from AC on any string starting with '0' or '1'

Final: {A}, {BCAC}

NFA \rightarrow DFA: Minimization (cont.)

- States *BC* and *AC* do not need differentiation
 - Should be merged into one



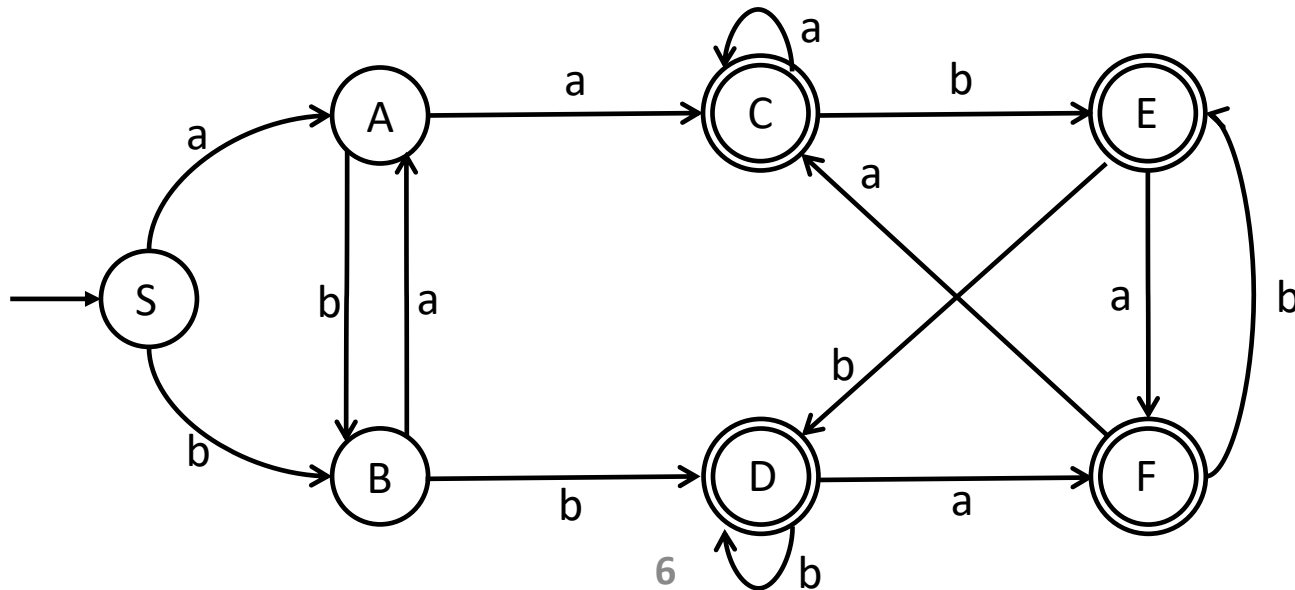
Minimization Algorithm

- The algorithm
 - Partitioning the states of a DFA into groups of states that **cannot be distinguished (i.e., equivalent)** [组内状态不可区分]
 - Each groups of states is then merged into a single state of the min-state DFA
- For a DFA $(\Sigma, S, n, F, \delta)$
 - The initial partition P_0 , has two sets $\{F\}$ and $\{S-F\}$
 - Splitting a set (i.e., partitioning a set s by input symbol α)
 - Assume q_a and $q_b \in s$, and $\delta(q_a, \alpha) = q_x$ and $\delta(q_b, \alpha) = q_y$
 - If q_x and q_y are not in the same set, then s must be split (i.e., α splits s)
 - One state in the final DFA cannot have two transitions on α

```
P ← {F}, {S-F}
while (P is still changing)
  T ← {}
  for each state s ∈ P
    for each α ∈ Σ
      partition s by α into s1 and s2
      T ← T ∪ s1 ∪ s2
  if T ≠ P then
    P ← T
```

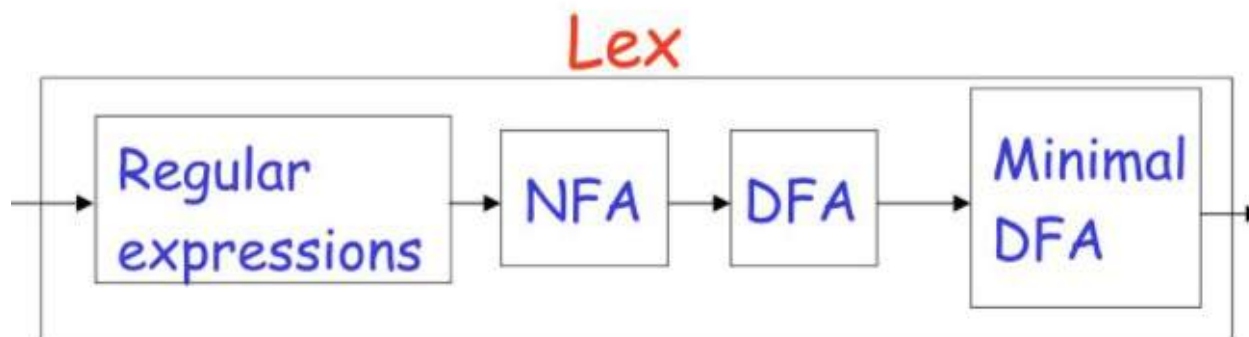
Example

- P0: $s_1 = \{S, A, B\}$, $s_2 = \{C, D, E, F\}$
- For s_1 , further splits into $\{S\}$, $\{A\}$, $\{B\}$
 - a: $S \rightarrow A \in s_1$, $A \rightarrow C \in s_2$, $B \rightarrow A \in s_1 \Rightarrow$ a distincts s_1
 - b: $S \rightarrow B \in s_1$, $A \rightarrow B \in s_1$, $B \rightarrow D \in s_2 \Rightarrow$ b distincts s_1
- For s_2 , all states are equivalent
 - a: $C \rightarrow C \in s_2$, $D \rightarrow F \in s_2$, $E \rightarrow F \in s_2$, $F \rightarrow C \in s_2 \Rightarrow$ a doesn't
 - b: $C \rightarrow E \in s_2$, $D \rightarrow D \in s_2$, $E \rightarrow D \in s_2$, $F \rightarrow E \in s_2 \Rightarrow$ b doesn't



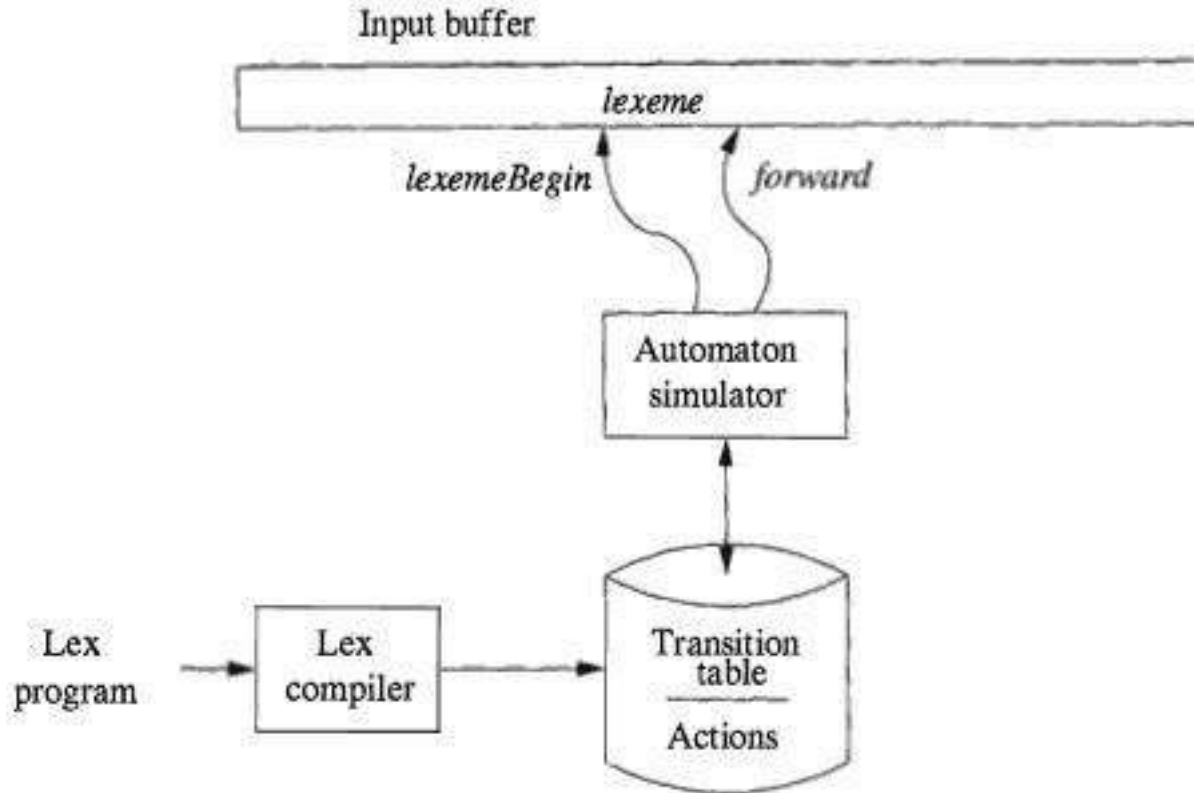
Implementation in Practice[实际实现]

- Lex: RE \rightarrow NFA \rightarrow DFA \rightarrow Table
 - Converts regular expressions to NFA
 - Converts NFA to DFA
 - Performs DFA state minimization to reduce space
 - Generate the transition table from DFA
 - Performs table compression to further reduce space
- Most other automated lexers also choose DFA over NFA
 - Trade off space for speed



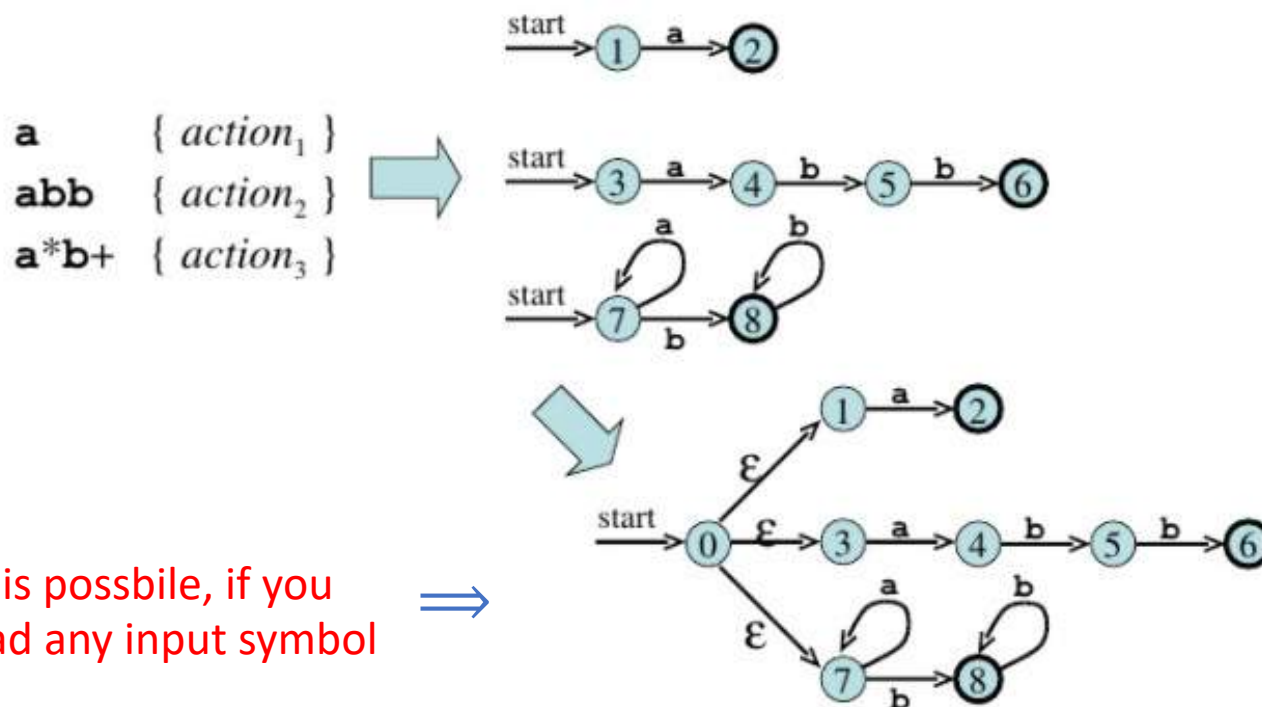
Lexical Analyzer Generated by Lex

- A Lex program is turned into a transition table and actions, which are used by a FA simulator
- Automaton recognizes matching any of the patterns



Lex: Example

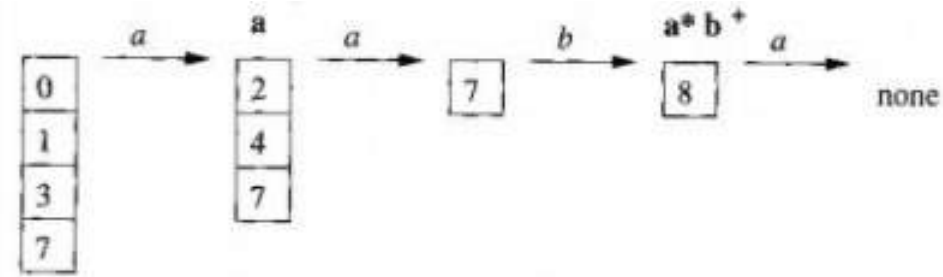
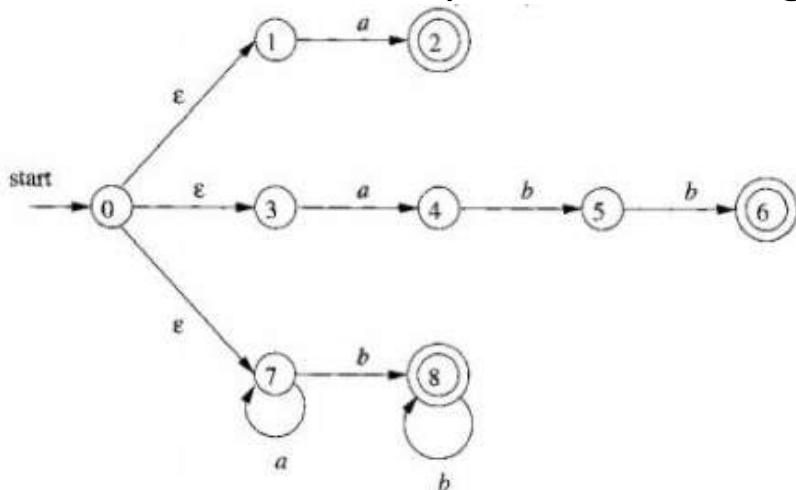
- Three patterns, three NFAs
- Combine three NFAs into a single NFA
 - Add start state 0 and ϵ -transitions



Any one is possible, if you haven't read any input symbol

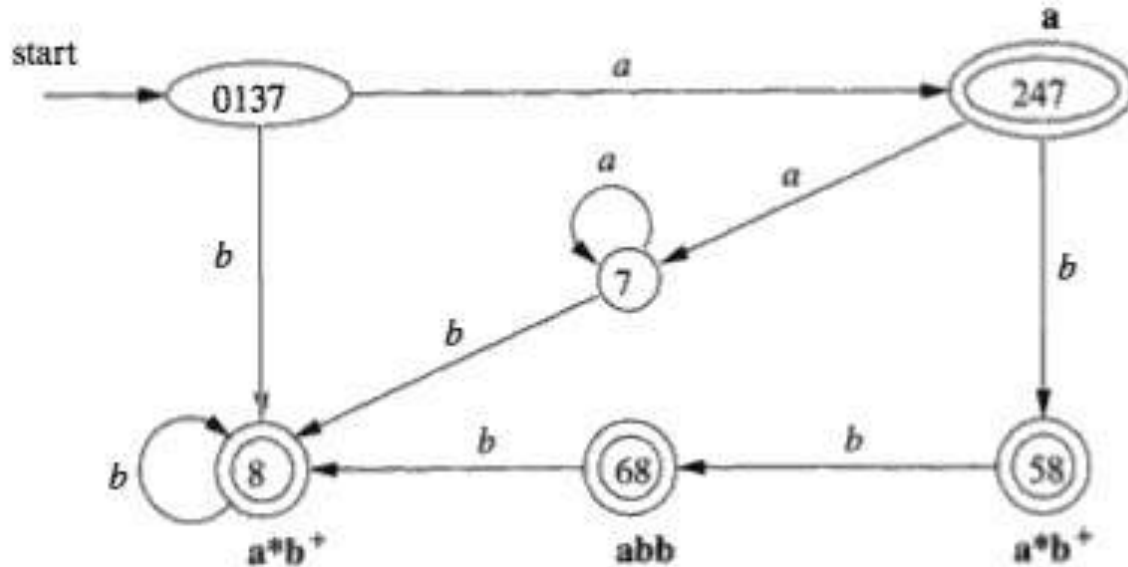
Lex: Example (cont.)

- NFA's for lexical analyzer
- Input: **aaba**
 - ϵ -closure(0) = {0, 1, 3, 7}
 - Empty states after reading the fourth input symbol
 - There are no transitions out of state 8
 - Back up, looking for a set of states that include an accepting state
 - State 8: a^*b^+ has been matched
 - Select **aab** as the lexeme, execute action₃
 - Return to parser indicating that token w/ pattern $p_3=a^*b^+$ has been found



Lex: Example (cont.)

- DFA's for lexical analyzer
- Input: **abba**
 - Sequence of states entered: $0137 \rightarrow 247 \rightarrow 58 \rightarrow 68$
 - At the final a , there is no transition out of state 68
 - 68 itself is an accepting state that reports pattern $p_2 = \mathbf{abb}$



How Much Should We Match? [匹配多少]

- In general, find the **longest match** possible
 - We have seen examples
 - One more example: input string **aabbb ...**
 - Have many prefixes that match the third pattern
 - Continue reading *b*'s until another *a* is met
 - Report the lexeme to be the initial *a*'s followed by as many *b*'s as there are
- If same length, rule appearing first takes precedence
 - String **abb** matches both the second and third
 - We consider it as a lexeme for p_2 , since that pattern listed first

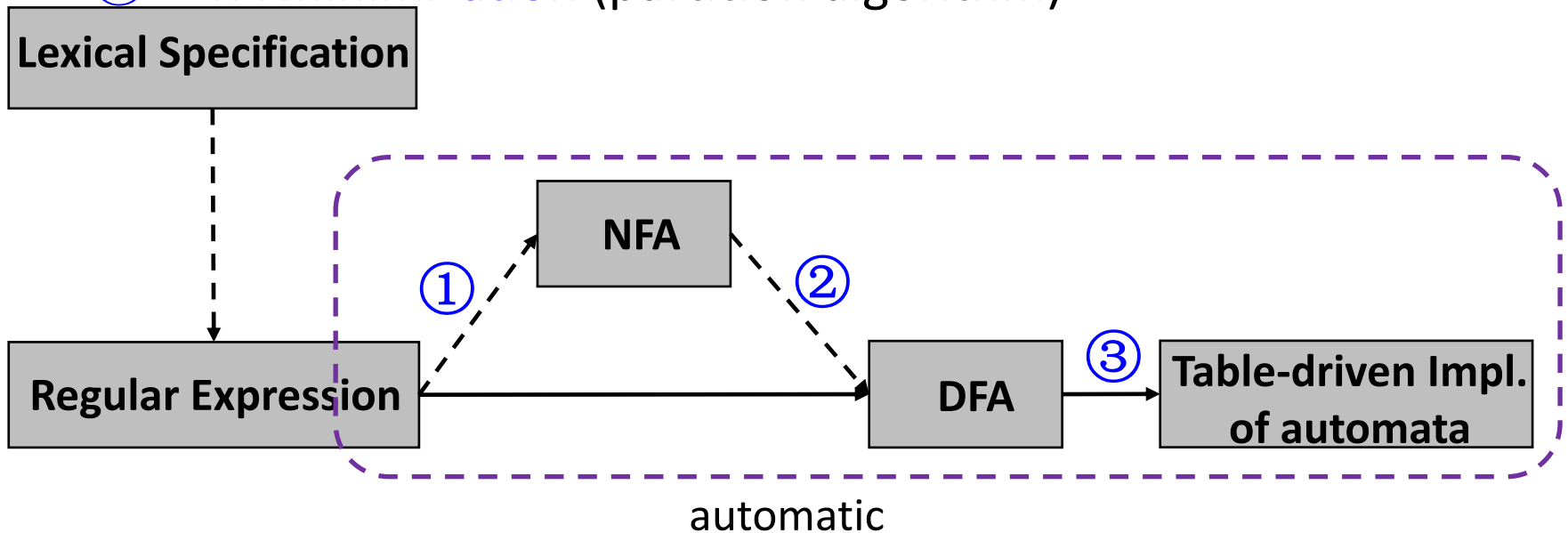
a	{ <i>action</i> ₁ }
abb	{ <i>action</i> ₂ }
a*b+	{ <i>action</i> ₃ }

How to Match Keywords? [匹配关键字]

- Example: to recognize the following tokens
 - Identifiers: letter(letter|digit)*
 - Keywords: if, then, else
- **Approach 1:** make REs for keywords and place them before REs for identifiers so that they will take precedence
 - Will result in more bloated finite state machine
- **Approach 2:** recognize keywords and identifiers using same RE but differentiate using special keyword table
 - Will result in more streamlined finite state machine
 - But extra table lookup is required
- Usually approach 2 is more efficient than 1, but you can implement approach 1 in your projects for simplicity

The Conversion Flow[转换流程]

- Outline: RE \rightarrow NFA \rightarrow DFA \rightarrow Table-driven Implementation
 - ③ Converting DFAs to table-driven implementations
 - ① Converting REs to NFAs (M-Y-T algorithm)
 - ② Converting NFAs to DFAs (subset construction)
 - ③' DFA minimization (partition algorithm)

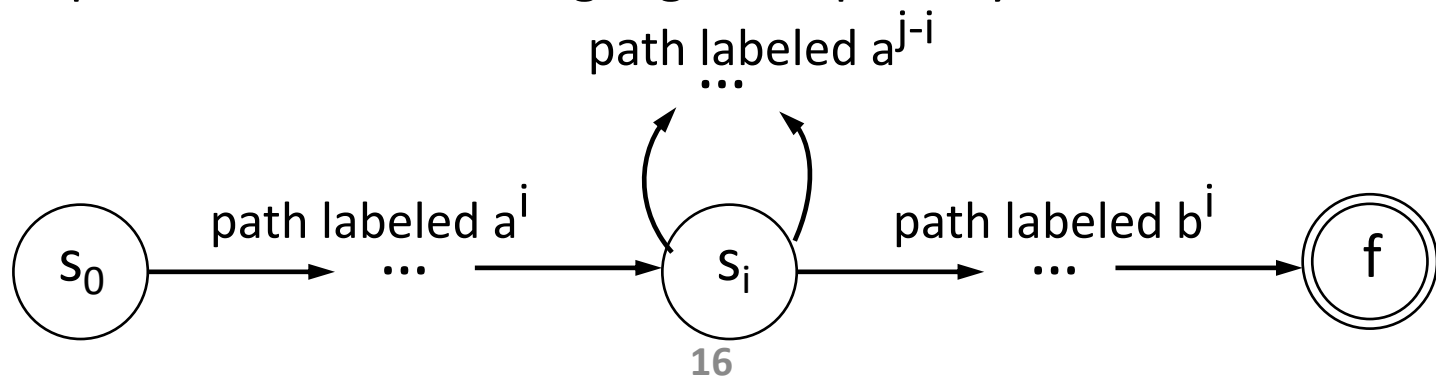


Beyond Regular Languages

- Regular languages are **expressive enough for tokens**
 - Can express identifiers, strings, comments, etc.
- However, it is the **weakest** (least expressive) language
 - Many languages are not regular
 - C programming language is not
 - The language matching braces “{{{...}}}” is also not
 - FA cannot count # of times char encountered
 - $L = \{a^n b^n \mid n \geq 1\}$
 - Crucial for analyzing languages with nested structures (e.g. nested for loop in C language)
- We need a more powerful language for parsing
 - Later, we will discuss context-free languages (CFGs)

RE/FA is NOT Powerful Enough

- $L = \{a^n b^n \mid n \geq 1\}$ is **NOT** a Regular Language
 - Suppose L were the language defined by regular expression
 - Then we could construct a DFA D with k states to accept L
 - Since D has only k states, for an input beginning with more than k a 's, D must enter some state twice, say s_j
 - Suppose that the path from s_j back to itself is labeled with a^{j-i}
 - Since $a^i b^i$ is in L , there must be a path labeled b^i from s_j to an accepting state f
 - But, there is also a path from s_0 through s_j to f labelled $a^i b^i$
 - Thus, D also accepts $a^i b^i$, which is not in L , contradicting the assumption that L is the language accepted by D



RE/FA is NOT Powerful Enough(cont.)

- $L = \{a^n b^n \mid n \geq 1\}$ is not a Regular Language
 - Proof \rightarrow Pumping Lemma (泵引理)
 - FA does not have any memory (FA cannot count)
 - The above L requires to keep count of a's before seeing b's
- Matching parenthesis is not a RL
- Any language with nested structure is not a RL
 - if ... if ... else ... else
- Regular Languages
 - Weakest formal languages that are widely used



中山大學
SUN YAT-SEN UNIVERSITY



国家超级计算广州中心
NATIONAL SUPERCOMPUTER CENTER IN GUANGZHOU

Compilation Principle 编译原理

第4讲：语法分析(1)

张献伟

xianweiz.github.io

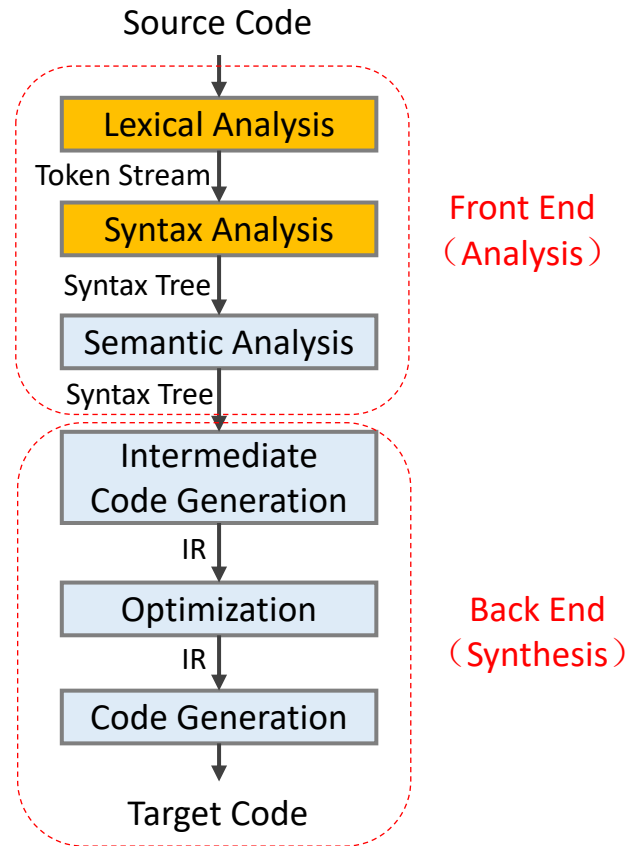
DCS290, 3/3/2022



中山大學
SUN YAT-SEN UNIVERSITY



Compilation Phases[编译阶段]

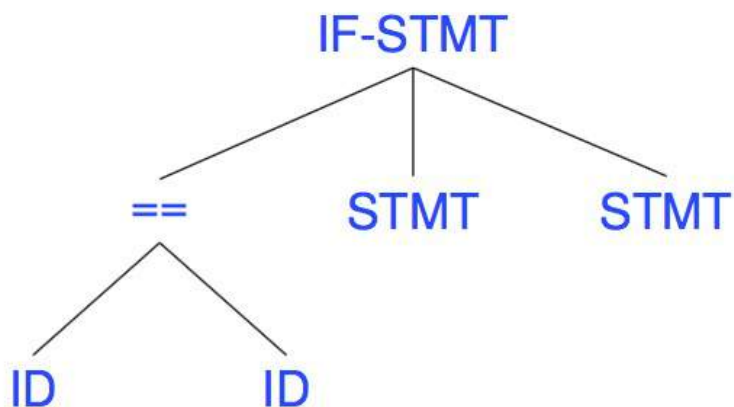


Syntax Analysis[语法分析]

- Second phase of compilation[第二阶段]
 - Also called as **parser**
- Parser obtains a string of tokens from the lexical analyzer[以token作为输入]
 - **Lexical analyzer** reads the chars of the source program, groups them into lexically meaningful units called **lexemes**
 - and produces as output **tokens** representing these lexemes
 - Token: <token name, attribute value>
 - Token names are used by parser for syntax analysis
 - tokens → parse tree/AST
- Parse tree[分析树]
 - Graphically represent the syntax structure of the token stream

Parsing Example

- Input: `if(x==y) ... else ...` [源程序输入]
- Parser input (Lexical output) [语法分析输入]
KEY(IF) '(' ID(x) OP('==') ID(y) ')' ... KEY(ELSE) ...
- Parser output [语法分析输出]



Parsing Example (cont.)

- Example: `<id, x> <op, *> <op, %>`
 - Is it a valid token stream in C language? **YES**
 - Is it a valid statement in C language (`x *%`)? **NO**
- Not every string of tokens are valid
 - Parser must distinguish between valid and invalid token strings
- We need a method to describe what is valid string?
 - To specify the syntax of a programming language

How to Specify Syntax?

- How can we specify a syntax with nested structures?
 - Is it possible to use RE/FA?
 - $L(\text{Regular Expression}) \equiv L(\text{Finite Automata})$
- RE/FA is **not powerful enough**
 - $L = \{a^n b^n \mid n \geq 1\}$ is not a Regular Language
- Example: matching parenthesis: # of '(' == # of ')'

– $(x+y)^*z$	✓
– $((x+y)+y)^*z$	✓
– $(\dots(((x+y)+y)+y)\dots)$	✓
– $((x+y)+y)+y)^*z$	✗

What Language Do We Need?

- C-language syntax: **Context Free Language (CFL)**[上下文无关语言]
e.g., 'else' is always 'else', wherever you place it
 - A broader category of languages that includes languages with nested structures
- Before discussing CFL, we need to learn a more general way of specifying languages than RE, called **Grammars**[文法]
 - Can specify both RL and CFL
 - and more ...
- Everything that can be described by a regular expression can also be described by a grammar
 - Grammars are most useful for describing nested structures

Concepts

- **Language**[语言]
 - Set of strings over alphabet
 - *String*: finite sequence of symbols
 - *Alphabet*: finite set of symbols
- **Grammar**[文法]
 - To systematically describe the syntax of programming language constructs like expressions and statements
- **Syntax**[语法]
 - Describes the proper form of the programs
 - Specified by grammar

Grammar[文法]

- Formal definition[形式化定义]: 4 components $\{T, N, s, \delta\}$
- **T**: set of terminal symbols[终结符]
 - Basic symbols from which strings are formed
 - Essentially tokens - leaves in the parse tree
- **N**: set of non-terminal symbols[非终结符]
 - Each represents a set of strings of terminals – internal nodes
 - E.g.: declaration, statement, loop, ...
- **s**: start symbol[开始符号]
 - One of the non-terminals
- σ : set of productions[产生式]
 - Specify the manner in which the terminals and non-terminals can be combined to form strings
 - “LHS \rightarrow RHS”: left-hand-side produces right-hand-side

Grammar (cont.)

- Usually, we can only write the σ [简写]
- Merge rules sharing the same LHS[规则合并]
 - $\alpha \rightarrow \beta_1, \alpha \rightarrow \beta_2, \dots, \alpha \rightarrow \beta_n$
 - $\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

$G = (\{id, +, *, (,)\}, \{E\}, E, P)$
 $P = \{ E \rightarrow E + E,$
 $E \rightarrow E * E,$
 $E \rightarrow (E),$
 $E \rightarrow id \}$

$G: E \rightarrow E + E,$
 $E \rightarrow E * E,$
 $E \rightarrow (E),$
 $E \rightarrow id \}$

$E \rightarrow E + E \mid E * E \mid (E) \mid id$