# Principles of Compiler Construction

DFA         iff.
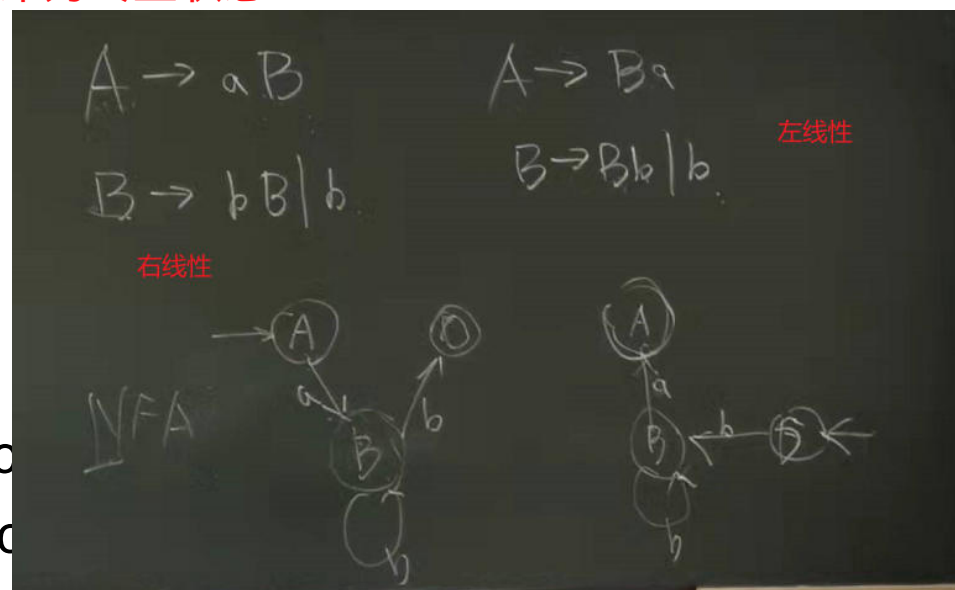
A->Ba    L(B)      a        A

**Prof. Wen-jun LI**

School of Computer Sc

lnslwj@mail.sysu.edu.c



右线性

左线性

eg. reduced DFA

lexical
    [ˈleksɪk(ə)l]        [ˈleksɪkl]
adj.

# Lecture 3. Lexical Analysis

scanner

1. Introduction

2. Scanner Construction

3. Lexical Specification

4. Finite Automata

5. Transformation and Equivalence

6. Limits of Regular Languages
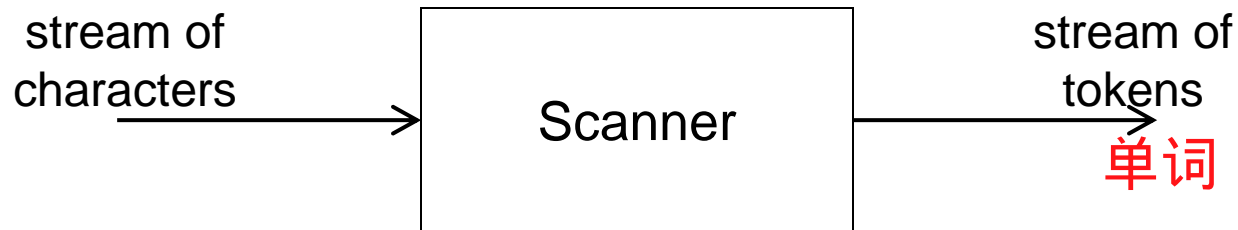
7. Lexical Analysis in Practice

# 1. Introduction

○ Software Architecture: Pipes and Filters
n.

stream of
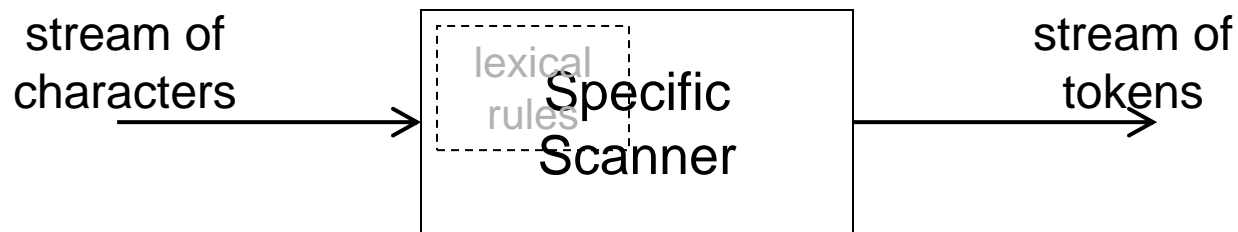characters → | Scanner | → stream of
tokens

# Structure of a Scanner

○ Implicit lexical rules

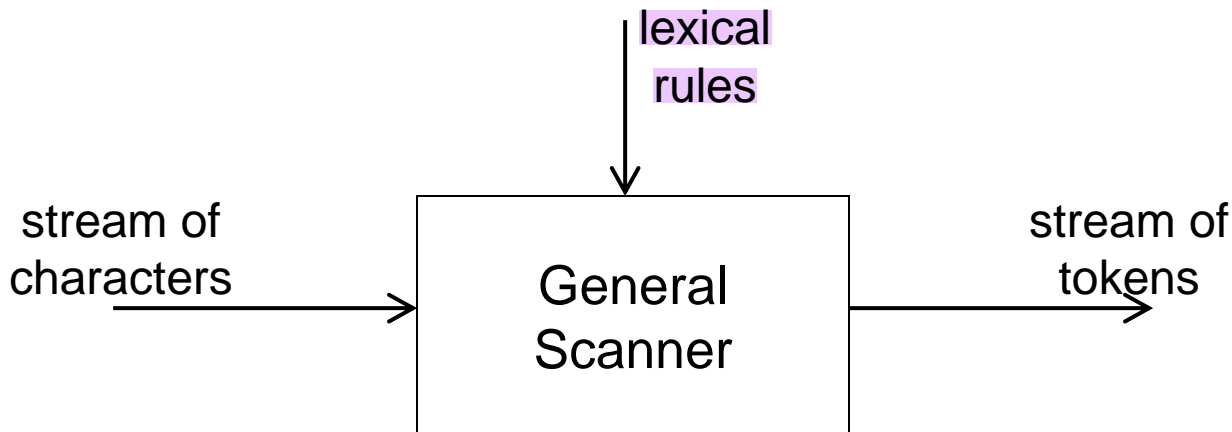| stream of characters | → | lexical rules | Specific Scanner | → | stream of tokens |

Specific to some predefined language.

interpretation
[ɪnˌtɜːprəˈteɪʃn]     [ɪnˌtɜːrprəˈteɪʃn]
n.

# Structure of a Scanner (cont')

○ Explicit lexical rules (interpretation model)

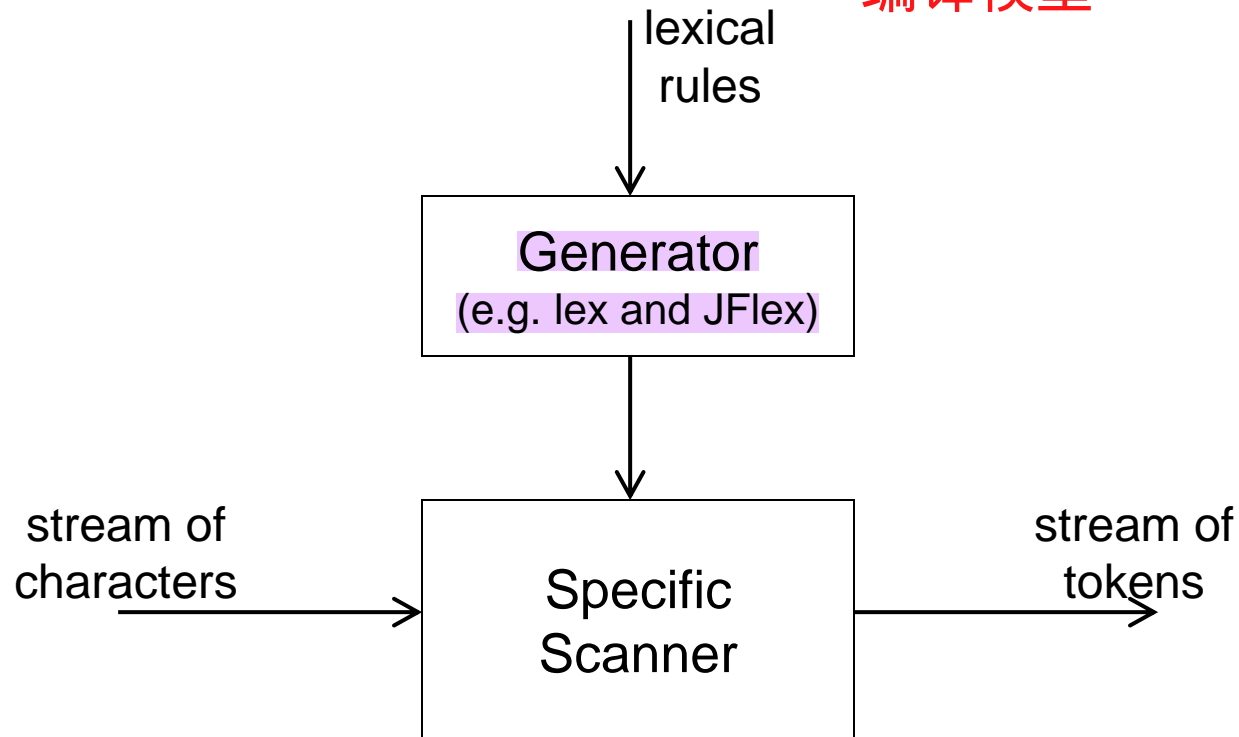lexical
rules

stream of
characters → 

General
Scanner

→ stream of
tokens

rules

No hard-coding of language-specific code.

# Structure of a Scanner (cont')

○ Explicit lexical rules (compilation model)



lexical
rules

↓

**Generator**
(e.g. lex and JFlex)

↓

stream of
characters →

**Specific
Scanner**

→ stream of
tokens

# Input Buffering

- Most runtime of the front end of a compiler is spent on scanning.
  - But the most time consuming stage is optimization.
- Buffer Pairs
  - 2-buffer scheme: each buffer is of size N.
  - Avoid overwriting: $|lexeme| + |lookaheads| \leq N$

### 3.2.1 缓冲区对

由于在编译一个大型源程序时需要处理大量的字符，处理这些字符需要很多的时间，因此开发了一些特殊的缓冲技术来减少用于处理单个输入字符的时间开销。一种重要的机制就是利用两个交替读入的缓冲区，如图 3-3 所示。
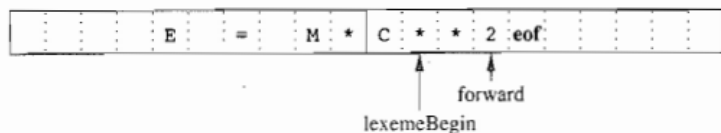


图 3-3 使用一对输入缓冲区

每个缓冲区的容量都是 $N$ 个字符，通常 $N$ 是一个磁盘块的大小，如 4096 字节。我们可以使用系统读取命令一次将 $N$ 个字符读入到缓冲区中，而不是每读入一个字符调用一次系统读取命令。如果输入文件中的剩余字符不足 $N$ 个，那么就会有一个特殊字符（用 eof 表示）来标记源文件的结束。这个特殊字符不同于任何可能出现在源程序中的字符。

程序为输入维护了两个指针：

1) lexemeBegin 指针：该指针指向当前词素的开始处。当前我们正试图确定这个词素的结尾。

2) forward 指针：它一直向前扫描，直到发现某个模式被匹配为止。做出这个决定所依据的策略将在本章的其余部分中讨论。

一旦确定了下一个词素，forward 指针将指向该词素结尾的字符。词法分析器将这个词素作为某个返回给语法分析器的词法单元的属性值记录下来。然后使 lexemeBegin 指针指向刚刚找到的词素之后的第一个字符。在图 3-3 中，我们看到，forward 指针已经越过下一个词素 ** （Fortran 的指数运算符）。在处理完这个词素后，它将会被左移一个位置。

将 forward 指针前移要求我们首先检查是否已经到达某个缓冲区的末尾。如果是，我们必须将 $N$ 个新字符读到另一个缓冲区中，且将 forward 指针指向这个新载入字符的缓冲区的头部。只要我们从不需要越过实际的词素向前看很远，以至于这个词素的长度加上我们向前看的距离大于 $N$，我们就决不会在识别这个词素之前覆盖掉这个尚在缓冲区中的词素。

# Output Tokens

- 5 kinds for most modern programming languages
  - Identifiers: getBalance, weight, …
  - Reserved words: IF, ELSE, WHILE, …
  - Constants: 10, 3.14, -1.26E-5, 'a', "abc", …
  - Operators: +, -, *, /, <<, …
  - Punctuation: (, ), ;, :, …
- Output
  - A Pair of <kind, associatedAttributeValues>
    - Some values are pointers to the symbol table.
  - Symbol (string) table
    - Array, LinkedList, HashSet, TreeSet, …

key words

reserved words

# Interaction with Parser

- ○ **Aim at reducing passes**

- ○ 2 levels of abstraction

  - ● **Logically**
    - ○ **A pipe** through which tokens are transferred.

  - ● **Physically**
    - ○ A disk file of the token sequence
    - ○ Concurrent threads or co-routines
    - ○ **Method invocation**                    caller&callee

invocation
    [ˌɪnvəˈkeɪʃn]          [ˌɪnvəˈkeɪʃn]

n.

# 2. Scanner Construction

○ Definition of lexical rules: 3 equivalent notations
$L(RE)=L(RG)$

- Expressions *describe*
  - ○ Regular expression
  - ○ Regular definition
- Grammars *generate*
  - ○ Regular grammar
  - ○ Left/right linear grammar
- Finite automata *recognize*
  - ○ Deterministic Finite Automata (DFA)
  - ○ Nondeterministic Finite Automata (NFA)
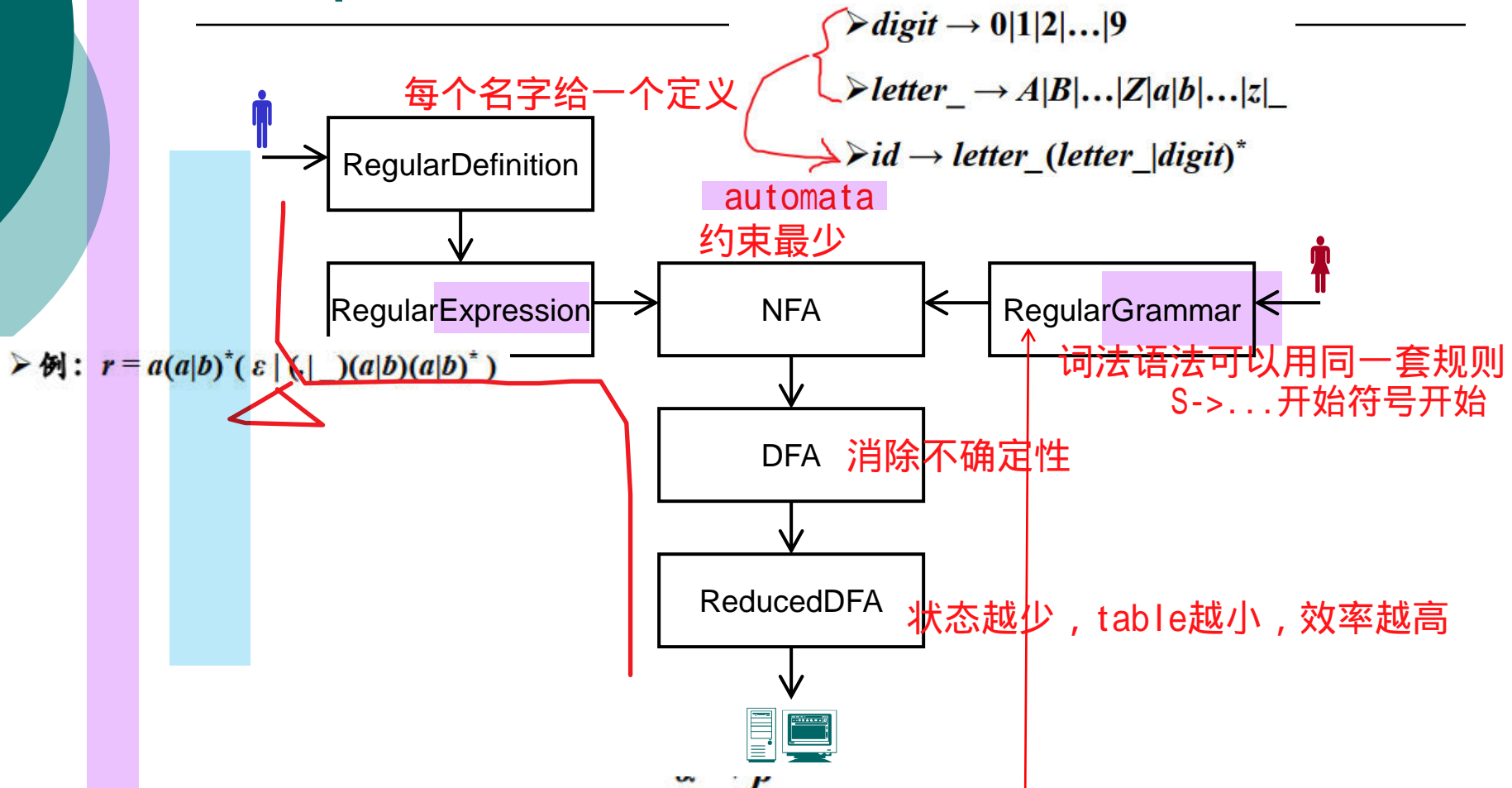
# Begin with Formal Specification

○ Specification of tokens
- Regular definition
  - ○ Hierarchical and brief, suitable for **humans**.
  - ○ Then transferred to regular expressions.
  - ○ Then transferred to finite automata.
- Regular grammar
  - ○ Compliant with the spec. of syntax rules.
  - ○ Then transferred to finite automata.

○ Transferring finite automata
- NFA → DFA → Reduced DFA
- Suitable for implementation on a **machine**.

compliant
  [kəmˈplaɪənt]        [kəmˈplaɪənt]
adj.

# Steps of Construction

> $digit \rightarrow 0|1|2|\ldots|9$

> $letter\_ \rightarrow A|B|\ldots|Z|a|b|\ldots|z|\_$

> $id \rightarrow letter\_(letter\_|digit)^*$

automata

RegularDefinition

> 例: $r = a(a|b)^*(\ \varepsilon\ |\ (\ |\ \ )(a|b)(a|b)^*\ )$

RegularExpression → NFA ← RegularGrammar

S->...

DFA

table

ReducedDFA

> 正则文法 (Regular Grammar, RG)

种分类 > 右线性(Right Linear)文法: $A \rightarrow wB$ 或 $A \rightarrow w$ 右侧最多1个非终结符，在同一侧

> 左线性(Left Linear) 文法: $A \rightarrow Bw$ 或 $A \rightarrow w$

# Programming a Scanner

- Implicit transition diagram
  - A manual approach
  - State transitions are hard-coded in the program.
- Explicit transition diagram
  - A table-driven approach
  - No hard-coding of specific lexical rules
  - Lead to automatic scanner generation.

# 3. Lexical Specification

- Regular Expression
- Regular Definition
- Regular Grammar
- Transformation and Equivalence

# Regular Expression

- Regular expression: constructively defined
  - Basis:                          L:
    - $\varepsilon$ is a regular expr; $L(\varepsilon) = \{\varepsilon\}$.
    - **a** is a regular expr if a $\in \Sigma$; $L(\textbf{a}) = \{a\}$.
  - Induction: if **r** and **s** are regular expressions,
    - **r | s** is a regular expr; $L(\textbf{r | s}) = L(\textbf{r}) \cup L(\textbf{s})$.
    - **r s** is a regular expr; $L(\textbf{r s}) = L(\textbf{r})\,L(\textbf{s})$.
    - **r\*** is a regular expr; $L(\textbf{r}^*) = (L(\textbf{r}))^*$.
    - **(r)** is a regular expr; $L(\textbf{(r)}) = L(\textbf{r})$.
  - Extensions:
    - $\textbf{r}^+ = \textbf{r}\,\textbf{r}^* = \textbf{r}^*\,\textbf{r}$          $\textbf{r}^* = \textbf{r}^+ \,|\, \varepsilon$  optional
    - $\textbf{r}^? = \textbf{r} \,|\, \varepsilon$
    - **[ace] = a | c | e**      **[A–Z] = A | B | … | Z**

                                            RE

# Regular Expression Language

- Regular expression: a language
  - Syntax
    - **a (a | b) b**
    - **a (a | b)***
  - Semantics
    - {aab, abb}
    - {a, aa, ab, aaa, aab, aba, abb, …}
- Language is an alternative approach to problem solving.

# Discussions

- ○ Signature of an operator
  - • true: $\rightarrow$ bool
  - • +: int $\times$ int $\rightarrow$ int
  - • +: real $\times$ real $\rightarrow$ real  (overloading)
- ○ What is the signature of the semantic function L in the previous slides ?
  - • Syntactic category
  - • Semantic category
  - • Mapping

# Alphabet: Lexical vs. Syntax

○ Both scanner and parser are based over an alphabet $\Sigma$.

○ But the meaning (elements) of $\Sigma$ is quite different.

- Scanner: elements in $\Sigma$ are characters in source programs.
  - ASCII, EBCDIC, Unicode, …
- Parser: elements in $\Sigma$ are tokens generated and passed by the scanner.

# Algebraic Laws for Regular Expressions

- r | s = s | r
- r | (s | t) = (r | s) | t
- r | r = r
- r(st) = (rs)t
- r(s | t) = rs | rt
- (s | t)r = sr | tr
- rε = εr = r
- r*r* = r*

- r* = ε | r | rr | …
- (r*)* = r*
- rr* = r*r
- (r* | s*)* = (r*s*)*
- (r*s*)* = (r | s)*
- (rs)*r = r(sr)*
- (r | s)* = (r*s)*r*

**Discussion: what does it mean that two regular expressions are equivalent ?**

# Binding to Practice

○ Identifiers in Pascal

$(a \mid b \mid ... \mid z \mid A \mid B \mid ... \mid Z) ( (a \mid b \mid ... \mid z \mid A \mid B \mid ... \mid Z) \mid (0 \mid 1 \mid ... \mid 9) )^*$

○ Hard to read and write

# Regular Definition

○ In the form of

- $d_1 \rightarrow r_1$
- $d_2 \rightarrow r_2$
- …
- $d_n \rightarrow r_n$

  ○ where $d_i \notin \Sigma \wedge (d_i = d_j \Rightarrow i = j)$, and
  ○ $r_i$ is a regular expr over $\Sigma \cup \{d_1, d_2, …, d_{i-1}\}$

# Binding to Practice

- ○ Identifiers in Pascal (extended)
  - • letter_ → A | B | … | Z | a | b | … | z | _
  - • digit → 0 | 1 | … | 9
  - • id → letter_ (letter_ | digit)*

- ○ Real numbers in Pascal
  - • digit → 0 | 1 | … | 9
  - • digits → digit digit*
  - • optionalFraction → . digits | ε
  - • optionalExponent → (E (+ | - | ε) digits) | ε
  - • real → digits optionalFraction optionalExponent

- ○ Hierarchical representation of regular expression

# Binding to Practice (cont')

○ Writing a regular definition <mark>with extended operators</mark>

- digit $\rightarrow$ [0–9]
- digits $\rightarrow$ digit$^+$
- real $\rightarrow$ digits (. digits)$^?$ (E [+-]$^?$ digits)$^?$

**Note the difference between '–' and '-'**

# Transforming between Regular Definitions and Regular Expressions

- Regular definition → regular expression
  - Top-down, stepwise substitution
  - **Discussion**: does the substitution process always terminate ?
- Regular expression → regular definition
  - Simply add a new symbol and "→".
- Equivalence
  - Defined by the language they denote
  - Bi-direction

# Transforming Regular Definitions to Regular Grammars

- ○ Difficulties arise from the closure operator
  - Transform to a left/right recursive production.
- ○ In regular definition
  - id $\rightarrow$ letter (letter | digit)$^*$
- ○ Introduce a right part
  - id $\rightarrow$ letter rid
  - rid $\rightarrow$ rid (letter | digit) | $\varepsilon$ (not permitted in Grammar)
    that is, rid $\rightarrow$ rid letter | rid digit | $\varepsilon$
  - or right-linear, rid $\rightarrow$ letter rid | digit rid | $\varepsilon$

# Class 3 Grammars

○ Right Linear (Regular) Grammar

- A $\rightarrow$ aB or A $\rightarrow$ a, where A, B $\in$ N $\wedge$ a $\in$ $\Sigma$ $\bigcup$ {$\varepsilon$}.

○ Left Linear Grammar

- A $\rightarrow$ Ba or A $\rightarrow$ a, where A, B $\in$ N $\wedge$ a $\in$ $\Sigma$ $\bigcup$ {$\varepsilon$}.

○ Extended Forms

- Right Linear (Regular) Grammar
  A $\rightarrow$ $\alpha$B or A $\rightarrow$ $\alpha$, where A, B $\in$ N $\wedge$ $\alpha$ $\in$ $\Sigma^*$.

- Left Linear Grammar
  A $\rightarrow$ B$\alpha$ or A $\rightarrow$ $\alpha$, where A, B $\in$ N $\wedge$ $\alpha$ $\in$ $\Sigma^*$.

# 4. Finite Automata *recognize*

- 3 types of finite automata
  - DFA (Deterministic Finite Automata)
  - NFA (Nondeterministic Finite Automata)
  - ε-NFA (Nondeterministic Finite Automata with empty transitions)
- Languages defined by a finite automaton
- Representation of finite automata
  - For humans: Transition Diagrams
  - For computers/machines: Transition Tables

# Formal Definition

- A finite automaton is a 5-tuple: $M = (\Sigma, S, \Delta, s_0, F)$, where
  - $\Sigma$ is the input alphabet.
  - $S \cap \Sigma = \varnothing$ is a **finite** set of states.
  - $\Delta$ is a transition function.
  - $s_0 \in S$ is the start (initial) state.
  - $F \subseteq S$ is a set of final (accepting) states.
- DFA, NFA, and $\varepsilon$-NFA only differ in
  - For a DFA, $\qquad\qquad \delta : S \times \Sigma \to S$
  - For an NFA, $\qquad\qquad \delta : S \times \Sigma \to 2^S$
  - For an $\varepsilon$-NFA, $\qquad\quad \delta : S \times (\Sigma \cup \{\varepsilon\}) \to 2^S$

# Language Defined by an FA

- **Overloading** the transition function
  - For any $a \in \Sigma$, $s \in S$ and $\omega \in \Sigma^*$,
    - $\delta(s, \varepsilon) = s$
    - $\delta(s, \omega a) = \delta(\delta(s, \omega), a)$

- Given a finite automaton M,
  - $L(M) = \{ \omega \mid \omega \in \Sigma^* \wedge \exists s \in F. \ \delta(s_0, \omega) = s \}$

# Transition Diagrams

○ Graphical representation of FA

# Commonly Used in Practice

○ Description of the state transition of a stateful object



○ UML (Unified Modeling Language)
- De facto standard for object-oriented modeling and design.
- Statechart of a class.

# Internal Representation

- Transition table
  - Rows
    - DFA, NFA, and $\varepsilon$-NFA:  State S
  - Columns
    - DFA and NFA: Alphabet $\Sigma$
    - $\varepsilon$-NFA:  $\Sigma \cup \{\varepsilon\}$
  - Cells
    - DFA: an element of S
    - NFA and $\varepsilon$-NFA:  a subset of S
- That is why NFA & $\varepsilon$-NFA are not suitable for implementation.

# 5. Transformation and Equivalence

- Regular Grammar to $\varepsilon$-NFA
- Regular Expression to $\varepsilon$-NFA
- Determination (of $\varepsilon$-NFA)
- Reduction (of DFA)

# Regular Grammar to ε-NFA

○ Algorithm

  • Add a new final state $f$

  • For any a $\in \Sigma \cup \{\varepsilon\}$ and A, B $\in$ N,

    ○ If A $\rightarrow$ a where a $\in \Sigma \cup \{\varepsilon\}$, let $\delta$ (A, a) = $f$

    ○ If A $\rightarrow$ aB, let $\delta$ (A, a) = B

# An Example

- Given a regular grammar

$$A \rightarrow 0 \mid 0B \mid 1D$$

$$B \rightarrow 0D \mid 1C$$

$$C \rightarrow 0 \mid 0B \mid 1D$$

$$D \rightarrow 0D \mid 1D$$

- We have an equivalent $\varepsilon$-NFA

# Regular Expression to ε-NFA

○ Thompson's construction (McNaughton-Yamada-Thompson algorithm)

- A constructive approach corresponding with the constructive definition of regular expressions
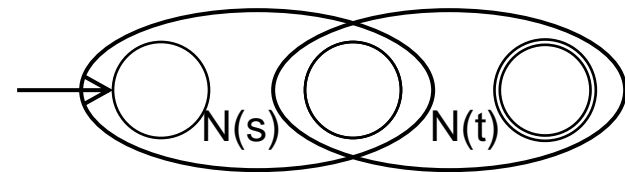
- Basis:

  ○ ε

  ○ a ∈ Σ

- Induction:
  - r = s | t
    2 new states
    4 new ε transitions

  - r = s t
    No new states
    or transitions

  - r = s*
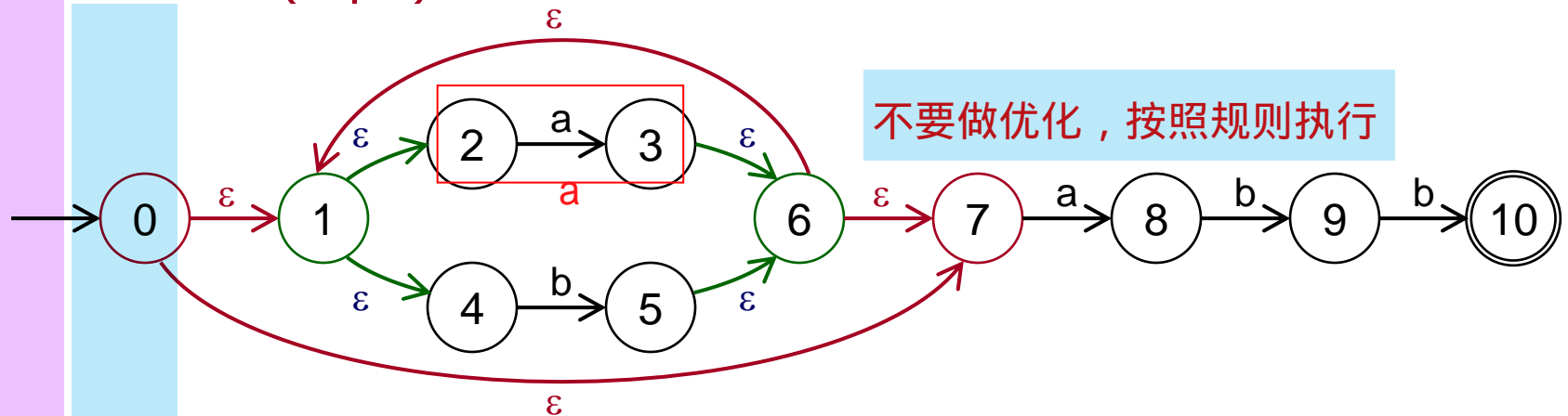    2 new states
    4 new ε transitions

# An Example

- Given regular expression
  - (a | b)*abb
- We have
  - a | b
  - (a | b)*

- NFA

# ε-NFA to DFA

- Both 2 approaches are applied
  - Subset construction
  - ε-closure construction
- Subset
  - A state in the new DFA corresponds to a subset of states in the original NFA.
- ε-closure
  - The empty transition does not consume any input.

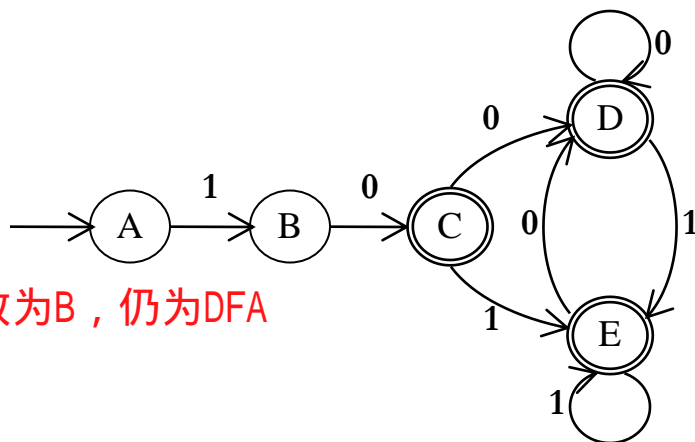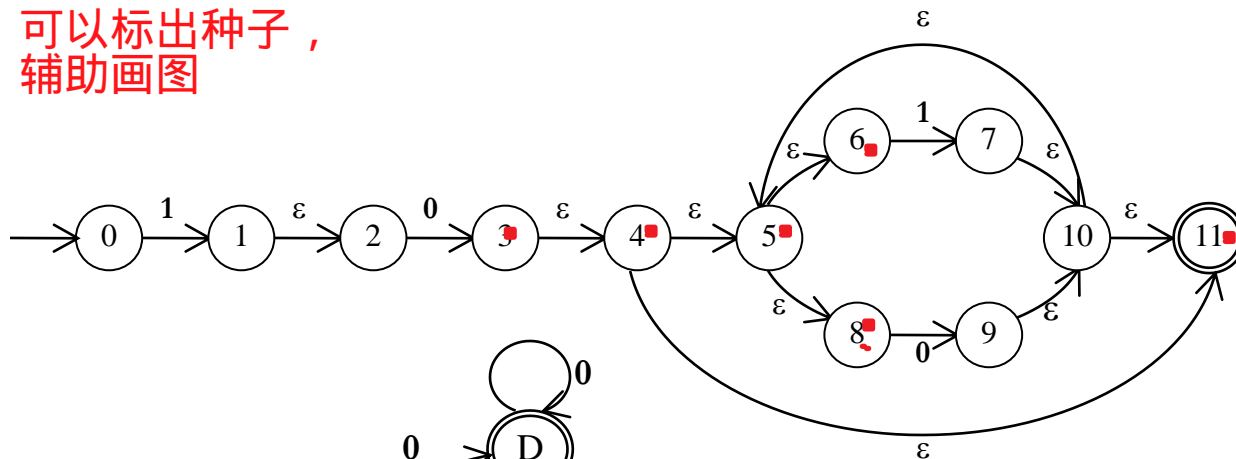# An Example

A = {0}, B = {1, 2}, C = {3, 4, 5, 6, 8, **11**}
D = {9, 10, **11**, 5, 6, 8}, E = {7, 10, **11**, 5, 6, 8}

# Calculate the ε-closure

- Let ε-closure({0}) = {0} = A, then
  - $\delta$ (A, 1) = ε-closure({1}) = {1, 2} = B
  - $\delta$ (B, 0) = ε-closure({3}) = {3, 4, 5, 6, 8, 11} = C
  - $\delta$ (C, 0) = ε-closure({9}) = {5, 6, 8, 9, 10, 11} = D
  - $\delta$ (C, 1) = ε-closure({7}) = {5, 6, 7, 8, 10, 11} = E
  - $\delta$ (D, 0) = ε-closure({9}) = D
  - $\delta$ (D, 1) = ε-closure({7}) = E
  - $\delta$ (E, 0) = ε-closure({9}) = D
  - $\delta$ (E, 1) = ε-closure({7}) = E

# Reduction of DFA

○ Find out equivalent classes of states
  - Approach 1:
    ○ Use a table to represent the equivalent relationship
    ○ Only a triangle is needed.
  - Approach 2:
    ○ Use partitions of the set of states
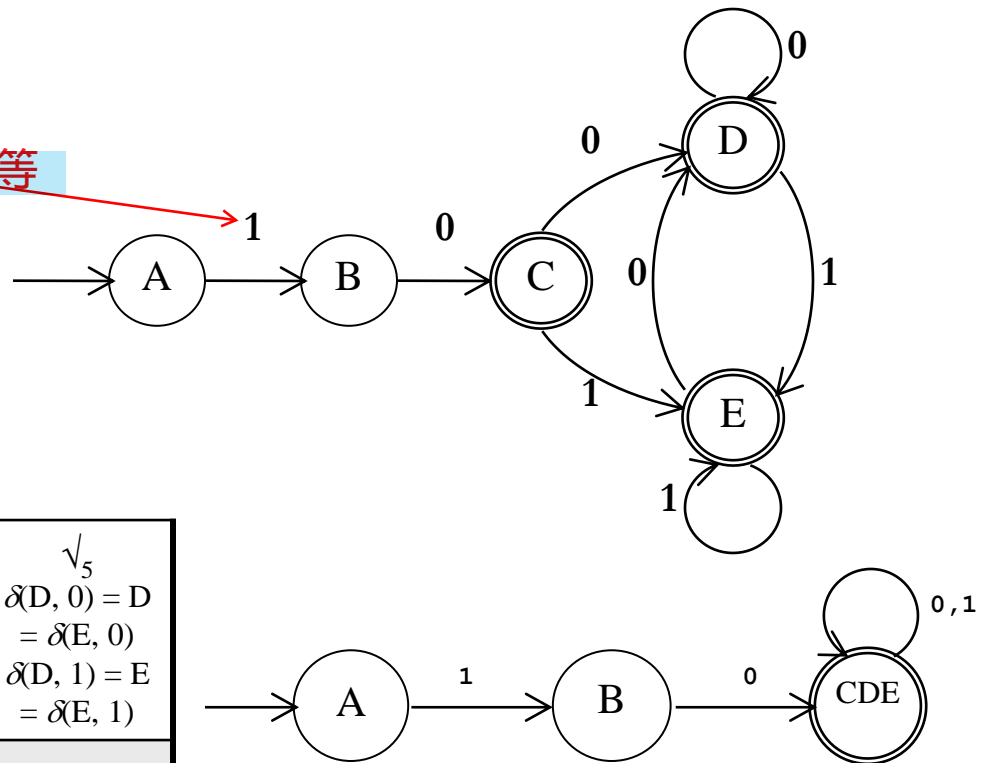
# An Example

| | | | | |
|---|---|---|---|---|
| B | $\times_2$ $\delta(A, 1) = B$ $\delta(B, 1) = \perp$ | | | |
| C | $\times_1$ Final | $\times_1$ Final | | |
| D | $\times_1$ Final | $\times_1$ Final | $\sqrt{}_3$ $\delta(C, 0) = D$ $= \delta(D, 0)$ $\delta(C, 1) = E$ $= \delta(D, 1)$ | |
| E | $\times_1$ Final | $\times_1$ Final | $\sqrt{}_4$ $\delta(C, 0) = D$ $= \delta(E, 0)$ $\delta(C, 1) = E$ $= \delta(E, 1)$ | $\sqrt{}_5$ $\delta(D, 0) = D$ $= \delta(E, 0)$ $\delta(D, 1) = E$ $= \delta(E, 1)$ |
| | A | B | C | D |

CDE
i fuxi lou
AB



yi fuxi lou

# More Example

- {1, 2, 4, 5, 6, 7, 8} {3}  (final and non-final, that is from $\varepsilon$)
- {1, 4, 5, 6, 7} {2, 8} {3}  (from 1 to final and non-final)
- {1, 5, 7} {4, 6} {2, 8} {3}  (from 0 to final and non-final)
- {1, 5} {7} {4, 6} {2, 8} {3}  (from 0 to state {2, 8} and {7})
  {1, 5} {7} {4, 6} {2, 8} {3}  (or from 1 to state {6} and {5})
- No difference found, terminate.

# 6. Limits of Regular Languages

scanni ng

- Regular languages are not able to handle infinite counting and matching
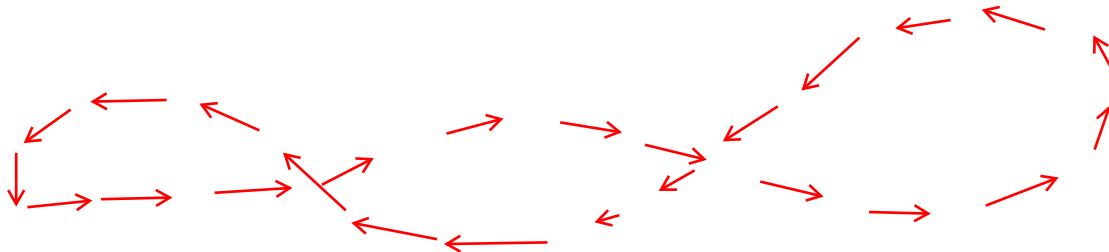  - e.g. $\{a^n b c^n \mid n \geq 1\}$ is NOT a regular set.

    begi n    end    a    c

- The matching of parentheses in most programming languages
  - Can not be described by regular expressions
  - Can not be generated by regular grammars
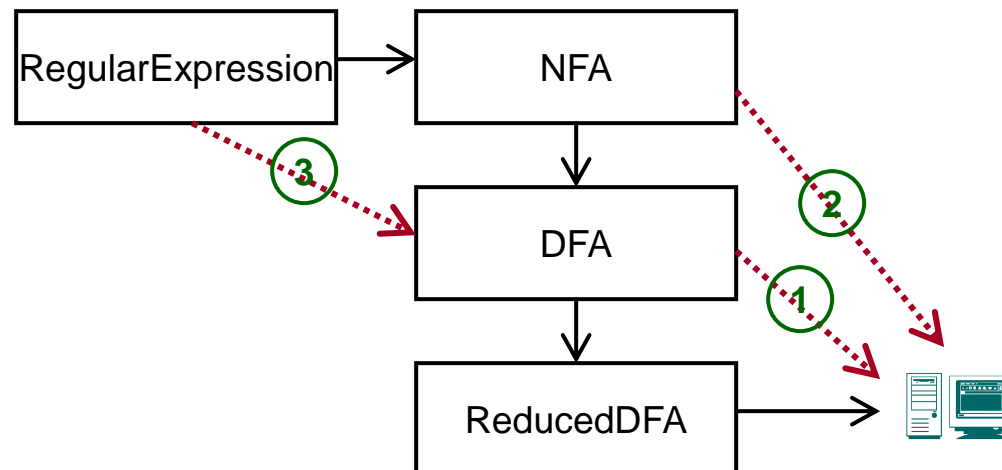  - Can not be recognized by finite automata

# 7. Lexical Analysis in Practice

- Scanner in practice: trained in our lab-time
  - Construct a scanner manually
    - Define the lexical rules
    - Construct a finite automaton
    - Write a scanner based on the blueprint of the finite automaton
  - Construct a scanner with a scanner generator
    - lex, Flex, and JFlex
      C                Java

# String Processing

○ Regular expression is a common way to describe a pattern in a string.
- Trade-off and consequence in the implementation of regular expressions
  - ○ Recognizing with an NFA
  - ○ Regular expression to DFA directly

# Exercise 3.1

- Give the recognized tokens of the following program in Pascal.

```
function max(i, j: integer): integer;
{return the maximum of integers i and j}
begin
  if i > j then max := i else max := j
end;
```

# Exercise 3.2

- (DBv2, Ch.3, pp.125, ex.3.3.2) Describe the languages denoted by the following regular expressions:
  - a (a | b)* a
  - a* b a* b a* b a*

# **Exercise 3.3**

- (DBv2, Ch.3, pp.125, ex.3.3.4) Most Languages are case sensitive, so keywords can be written only one way, and the regular expressions describing their lexemes are very simple.

- However, some languages, like Pascal and SQL, are case insensitive. For example, the SQL keyword `SELECT` can also be written `select`, `Select`, or `sELEcT`.

- Show how to write a regular expression for a keyword in a case insensitive language. Illustrate your idea by writing the expression for `SELECT` in SQL.

# Exercise 3.4

- Given the following regular expression

$$1^*(0 \mid 01)^*$$

(1) Transform it to an equivalent finite automaton.

(2) Construct an equivalent DFA for the result of exercise (1).

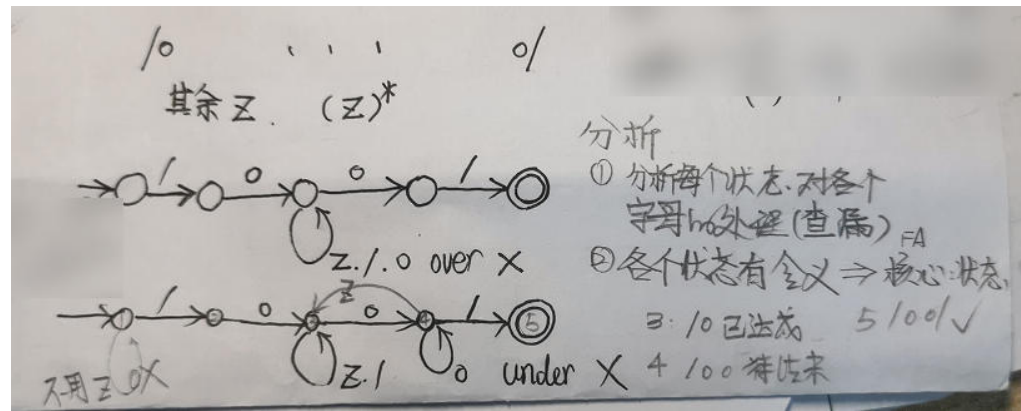(3) Reduce the result of (2) and get a reduced DFA.

# Exercise 3.5**

c

content          *

- Given the alphabet $\Sigma$ = { z, o, / }, a comment in a program over $\Sigma$ begins with "/o" and ends with "o/". Embedded comments are not permitted.

    (1) Draw a DFA that recognizes nothing but all the comments in the source programs.

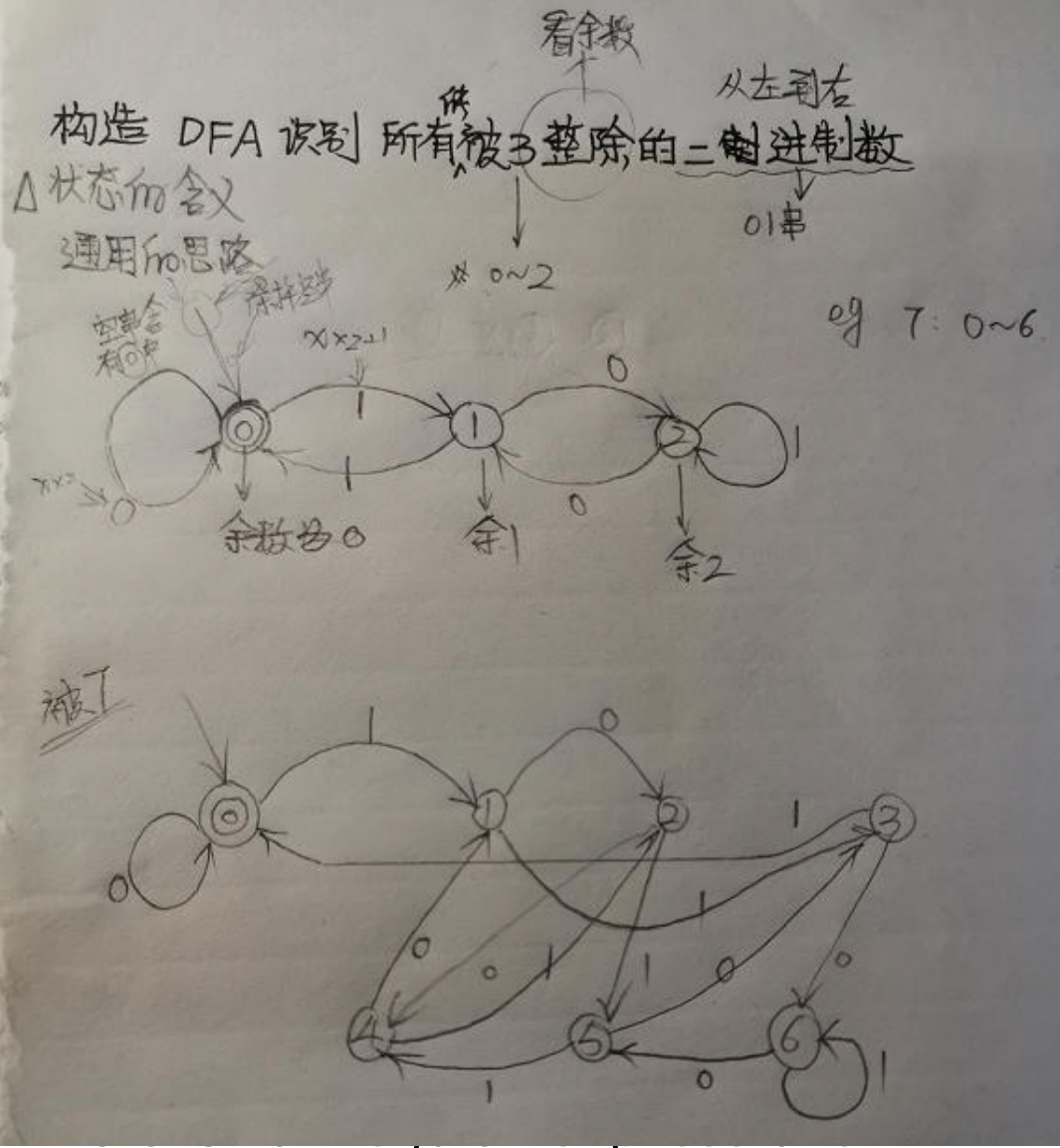    (2) Write a single regular expression that exactly describes all the comments in the source programs.

over specification
under ...

# **Further Rea**

- Dragon Book,
  - Comprehensiv
    - Section 2.6
      scanner.
    - Section 3.3
      definitions.
    - Section 3.6
      related tran
  - Skip reading:
    - Section 3.4
    - Section 3.9
      directly to DFAs.

# Enjoy the Course!