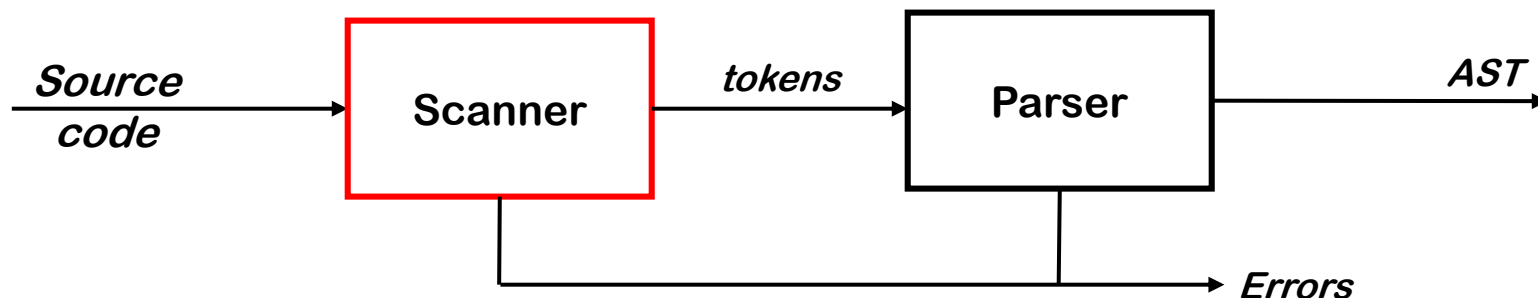# Lexical Analysis

Bernd Burgstaller
Yonsei University

# Outline

- The role of a scanner

- Scanner concepts

  - Tokens, Lexemes, Patterns

- Regular Expressions & Automata

  - Definitions of REs, DFAs and NFAs

  - REs→NFA (Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book)

  - NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon book)

  - DFA→minimal-state DFA (state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)

- Scanner generators

# The Role of the Scanner



Scanner

- Maps stream of characters into words called *tokens*.
  - Basic units of syntax
  - x = y + 21 ; *becomes*

    $\langle \mathbf{id}, x \rangle \; \langle = \rangle \; \langle \mathbf{id}, y \rangle \; \langle + \rangle \; \langle \mathbf{intliteral}, 21 \rangle \; \langle ; \rangle$

- Characters that form a word are its *lexeme.*
- Formally, a token is a tuple <Token_type, value>.
- Informally, we often say "token x" when we mean  $\langle \mathbf{id}, x \rangle$
- Tokens are to programming languages what words are to natural languages.

# Tokens

The tokens of our MiniC language are classified into *token types*:

- identifiers: `i`, `j`, `initial`, `position`, …
- keywords: `if`, `for`, `while`, `int`, `float`, `bool`, …
- operators: `+`  `-`  `*`  `/`  `<=`  `&&` …
- separators: `{`  `}`  `(`  `)`  `[`  `]`  `;`  `,`
- literals
    - integer literals: `0`, `1`, `22`, …
    - float literals: `1.25`  `1.`   `.01`   `1.2e2` …
    - bool literals: `true`, `false`
    - string literals: `"hello\n"`, `"string literal"`, `...`

- The exact token set depends on the given programming language. Pascal and Ada use ":=" for assignment, C uses "=".
- Natural languages also contain different kinds of tokens (words): verbs, nouns, articles, adjectives, … The exact token set depends on the natural language in question.

# Tokens can be described by patterns

- Pattern: a rule describing the set of lexemes of a particular token type.

| Token type | Pattern (informal) | Set of lexemes |
|---|---|---|
| INTLITERAL | a string of decimal digits | {0, 1, 324, …} |
| ID | a string of letters, digits and underscores, beginning with a letter or underscore | {sum, _sum, ptr_1, …} |
| + | the character "+" | {+} |
| if | the characters "i", "f" | {if} |

- The pattern is said to match each string in the set.
- We want a **formal notation** for patterns.
  => allows us to specify the tokens of programming languages.

# Tokens can be described by patterns

- Pattern: a rule describing the set of lexemes (=strings) of a particular token type.

| Pattern name | Pattern | Set of lexemes (strings) |
|---|---|---|
| BinaryDigit | 0 \| 1 | {0, 1} |
| DecimalDigit | 0\|1\|2\|3\|4\|5\|6\|7\|8\|9 | {0,1,2,3,4,5,6,7,8,9} |
| TwoBinaryDigits | (0\|1) (0\|1) | {00, 01, 10, 11} |
| TwoBinaryDigits | BinaryDigit BinaryDigit | {00, 01, 10, 11} |
| Bananana | b a (n a)* | {ba, bana, banana, bananana, …} |
| Register | r (0\|1\|2\|3\|4\|5\|6\|7\|8\|9) | {r0, r1, r2, … r9} |

- The set of lexemes (=strings) described by pattern X is called the language of X.
- We write L(X) to denote the language of pattern X.

# String Concatenation

- If x and y are strings, then xy is the string formed by appending y to x.

- Examples:

| x | y | xy |
|---|---|---|
| key | board | keyboard |
| java | script | javascript |

# Set Operations (Review)

| Operation | Definition |
|---|---|
| Union of $L$ and $M$, written $L \cup M$ | $L \cup M = \{s \mid s \in L \text{ or } s \in M\}$ |
| Concatenation of $L$ and $M$, written $LM$ | $LM = \{st \mid s \in L \text{ and } t \in M\}$ |
| Kleene star (or closure) of $L$, written $L^*$ | $L^* = \bigcup\limits_{0 \leq i \leq \infty} L^i$ <br><br> $L^0 = \{\epsilon\}, \quad L^i = \underbrace{LL \ldots L}_{i \text{ times concatenation}}$ <br><br> $\epsilon$ denotes the empty string |
| Postitive Kleene star (or positive closure) of $L$, written $L^+$ | $L^+ = \bigcup\limits_{1 \leq i \leq \infty} L^i$ |

8

# Examples: Operations on Languages

- $L = \{ a, \ldots z, A, \ldots Z \}$

- $D = \{ 0, \ldots, 9 \}$

| Example | Language |
|---------|----------|
| $L \cup D$ | letters and digits |
| $L^3$ | all 3-letter strings, e.g., abc, xyz, … |
| $LD$ | strings consisting of a letter followed by a digit, e.g., a1, b2, x1 |
| $L^*$ | all strings of letters, including the empty string $\varepsilon$ |
| $L(L \cup D)^*$ | all strings of letters and digits, starting with a letter, e.g., x, y, a1, aa, c999ccc |
| $D^+$ | all strings of one or more digits |

# Regular Expressions (inductive definition)

***Regular expressions*** *(REs)* can be used to describe patterns.

- A regular expression r is a pattern that describes a set of lexemes $L(r)$.

- Lexemes are made from characters of a finite set $\Sigma$ called alphabet. Notation: we <u>underline</u> the characters from $\Sigma$.

> Precedence is *closure*, then *concatenation*, then *alternation*

Inductive definition of regular expressions over alphabet $\Sigma$:

**R1:** $\varepsilon$ is a RE denoting the set $\{\varepsilon\}$

**R2:** If <u>a</u> is in $\Sigma$, then <u>a</u> is a RE denoting $\{\underline{a}\}$

If *x* and *y* are REs denoting *L(x)* and *L(y)* then

      **R3:** *x | y* is a RE denoting $L(x) \cup L(y)$      *called "alternation"*

      **R4:** *xy* is a RE denoting *L(x)L(y)*      *called "concatenation"*

      **R5:** $x^*$ is a RE denoting *L(x)\**      *called "repetition" or "closure"*

# Example REs

- $\sum$ = $\{\underline{0}, \underline{1}\}$

| RE | Language L(RE) |
|----|----------------|
| $\underline{1}$ | $L(\underline{1}) \stackrel{R2}{=} \{\underline{1}\}$ |
| $\underline{0} \mid \underline{1}$ | $L(\underline{0}\mid\underline{1}) \stackrel{R3}{=} L(\underline{0}) \cup L(\underline{1}) \stackrel{R2}{=} \{\underline{0}\} \cup \{\underline{1}\} = \{\underline{0}, \underline{1}\}$ |
| $\underline{1}^*$ | $L(\underline{1}^*) \stackrel{R5}{=} L(\underline{1})^* \stackrel{R2}{=} \{\underline{1}\}^* = \{\varepsilon, \underline{1}, \underline{1}\,\underline{1}, \underline{1}\,\underline{1}\,\underline{1}, \ldots\}$ |
| $\underline{1}^*\underline{1}$ | $L(\underline{1}^*\underline{1}) \stackrel{R4}{=} L(\underline{1}^*)L(\underline{1}) = \ldots = \{\underline{1}, \underline{1}\,\underline{1}, \underline{1}\,\underline{1}\,\underline{1}, \ldots\}$ |
| $\underline{0} \mid \underline{1}^*\underline{0}$ | the set containing $\underline{0}$ and all strings consisting of zero or more $\underline{1}$'s followed by a $\underline{0}$: $\{\underline{0}, \underline{1}\,\underline{0}, \underline{1}\,\underline{1}\,\underline{0}, \ldots\}$ |
| $(\underline{1}\,\underline{1}\,\underline{1})^*$ | the set of strings that contain zero or more sequences of $\underline{111}$: $\{\varepsilon, \underline{1}\,\underline{1}\,\underline{1}, \underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}, \underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\,\underline{1}\ldots\}$ |

11

# Example REs (cont.)

- $\sum$    $= \{\,\underline{a}\,,\underline{b}\,,\underline{c}\,\}$

| RE | Language L(RE) |
|---|---|
| $(\underline{a}\,\underline{b}\,)\,^*$ | zero or more concatenations of the string $\underline{a}\,\underline{b}$ <br> $\{\varepsilon\,,\underline{a}\,\underline{b}\,,\underline{a}\,\underline{b}\,\underline{a}\,\underline{b}\,,\underline{a}\,\underline{b}\,\underline{a}\,\underline{b}\,\underline{a}\,\underline{b}\,,\ldots\}$ |
| $(\underline{a}\mid\underline{b}\,)\,^*$ | zero or more concatenations of $(\underline{a}\mid\underline{b}\,)$ <br> $L((\underline{a}\mid\underline{b}\,)^*) = (L(\underline{a}\mid\underline{b}\,))^* = (L(\underline{a})\cup L(\underline{b}\,))^* = \{\,\underline{a}\,,\underline{b}\,\}\,^* =$ <br> $\{\varepsilon\,,\underline{a}\,,\underline{b}\,,\underline{a}\,\underline{a}\,,\underline{a}\,\underline{b}\,,\underline{b}\,\underline{a}\,,\underline{b}\,\underline{b}\,,\underline{a}\,\underline{a}\,\underline{a}\,,\underline{a}\,\underline{a}\,\underline{b}\,,\underline{a}\,\underline{b}\,\underline{a}\,,\underline{a}\,\underline{b}\,\underline{b}\,,\underline{b}\,\underline{a}\,\underline{a}\,,\underline{b}\,\underline{a}\,\underline{b}\,,\underline{b}\,\underline{b}\,\underline{a}\,,\underline{b}\,\underline{b}\,\underline{b}\,,\underline{a}\,\underline{a}\,\underline{a}\,\underline{a}\,,\ldots\}$ |
| $(\underline{a}\mid\underline{b}\,)\,^*\,\underline{c}$ | strings from above, each string concatenated with character $\underline{c}$ <br> $\{\underline{c}\,,\underline{a}\,\underline{c}\,,\underline{b}\,\underline{c}\,,\underline{a}\,\underline{a}\,\underline{c}\,,\underline{a}\,\underline{b}\,\underline{c}\,,\underline{b}\,\underline{a}\,\underline{c}\,,\underline{b}\,\underline{b}\,\underline{c}\,,\underline{a}\,\underline{a}\,\underline{a}\,\underline{c}\,,\underline{a}\,\underline{a}\,\underline{b}\,\underline{c}\,,\ldots\}$ |

# Precedence Rules

**Precedence rules** are needed to disambiguate regular expressions.

- Like with arithmetic: a + b * c  =  a + (b * c)

- Use parenthesis if you want (a + b) * c

Precedence of regular expression operators is **closure**, then **concatenation**, then **alternation**.

Examples:

ab*        = a(b*)         use parenthesis if you want (ab)*

ab|c       = (ab) | c      use parenthesis if you want a(b|c)

# Example Token Specifications using REs

- **digit**: $\underline{0} \mid \underline{1} \mid \underline{2} \mid \underline{3} \mid \underline{4} \mid \underline{5} \mid \underline{6} \mid \underline{7} \mid \underline{8} \mid \underline{9}$

  Example digits: $\underline{0}, \underline{1}, \underline{2}, ..., \underline{9}$

  > Notation: **digit** and **letter** are defined as symbolic names for REs.

- **letter**: $\underline{a} \mid \underline{b} \mid ... \mid \underline{z} \mid \underline{A} \mid \underline{B} \mid ... \mid \underline{Z}$

  Example letters: $\underline{a}, \underline{b}, \underline{c}, ..., \underline{Z}$

  > Notation: **simple_id** uses the symbolic names for **digit** and **letter** (like a shorthand notation).

- **simple_id**: **letter** (**letter**|**digit**)

  Simple 2-character identifiers, starting with a letter, followed by a letter or digit:

  Examples for simple_id: <u>a1</u>, <u>A1</u>, <u>ab</u>, <u>YD</u>

- **twochar_id** : (**letter** | **digit**) (**letter** | **digit**)

  2-character identifiers, each character can be either a letter or a digit:

  Examples for twochar_id: <u>a1</u>, <u>A1</u>, <u>1A</u>, <u>xY</u>, <u>12</u>, <u>00</u>

  What's the problem with identifiers consisting of digits only?

14

# Example Token Specifications (cont.)

- Identifiers with at least one character, the first character must be a letter, the subsequent characters can be letters or digits:

    $$id : letter\ (letter | digit)*$$

    Examples for id: <u>A</u>, <u>B</u>, <u>a</u>, <u>a1</u>, <u>Bernd</u>, <u>BestVar1y22</u>

- Integer numbers: first character is a digit, followed by zero or more digits:

    $$integer : digit\ digit*$$

    Examples: <u>0</u>, <u>1</u>, <u>2</u>, <u>01</u>, <u>101</u>, <u>999</u>, <u>4711</u>

- Signed integer numbers: an integer that can optionally have "<u>+</u>" or "<u>-</u>" in front:

    $$signed\_integer : (\underline{+} | \underline{-} | \varepsilon)\ integer$$

    Examples: <u>0</u>, <u>1</u>, <u>+1</u>, <u>-1</u>, <u>-4711</u>, <u>+1234</u>

# Writing Regular Expressions

- Write a regular expression for each of the following sets of tokens:

  – Ruby binary literals consisting of "0b" followed by the binary number. Examples: 0b001011, 0b01.

  – Ruby binary literals, with an optional underscore ("_") between a pair of binary digits.
    Examples: 0b0_101 , 0b11_01 , but not 0b_1 and not 0b1_ .

  – Ada identifiers: a letter followed by any number of letters, digits, and underlines. An identifier must not end in an underline or have two underlines in a row.

  – A floating point number: one or more digits (whole-number part) followed by a decimal point (".") and one or more digits (fractional part). Examples: 2.24, 0.1234.

  – Floating point number in scientific notation: same as above, but optionally followed by "e" or "E", and a signed integer exponent. Examples: 1.2e-2, 2.3E+34, 2.3E34.

# Shorthand Notations

- One or more concatenations: $r^+ = rr^*$

    – Denotes the language $(L(r))^+$

    – Same precedence and associativity as *

- Zero or one instance: $r? = \varepsilon | r$

    – Denotes the language $L(r) \cup \{\varepsilon\}$

    – Written as (r)? to indicate grouping, e.g., (12)?

- Character classes:

    $$[\underline{a}-\underline{z}\ \underline{A}-\underline{Z}]$$

# MiniC RE Examples

| Token | RE |
|---|---|
| **letter** | $[\underline{a} - \underline{z}\,\underline{A} - \underline{Z}\,\underline{\_}\,]$ |
| **identifier** | $\textbf{letter}\,(\textbf{letter}\mid\textbf{digit})*$ |
| **integer** | $\textbf{digit}^{+}$ |

- Note that with MiniC, **letter** includes the underscore character "**_**".
- In Java, letters and digits may be drawn from the entire Unicode character set. Examples of Java identifiers are

   abc                              벅스텔라

$\beta\varepsilon\rho\eta\delta$

# The Big Picture

Why are we doing this?

- We want to use regular expressions to specify scanners.

- We want to harness the theory from classes like "Automata".



Goals:

- To simplify specification & implementation of scanners

- To understand the underlying techniques and technologies

# Regular Expressions                    (the whole point)

*We use regular expressions to specify the mapping of lexemes to tokens.*

Using results from automata theory and theory of algorithms, we can automatically build recognizers from regular expressions.

$\Rightarrow$ We study REs and associated theory to automate scanner construction!

$\Rightarrow$ Fortunately the automatic techniques produce fast scanners. Used with text editors, URL filtering software, ...

# Example

Consider the problem of recognizing register names

   **Register**: r (0|1|2| ... | 9) (0|1|2| ... | 9)*

- Allows registers of arbitrary number
- Requires at least one digit
- $\Sigma = \{r,0,1,...,9\}$

RE corresponds to a recognizer (or Deterministic Finite Automaton, DFA)



**Recognizer for Register**

*Transitions on other inputs go to an error state, $s_e$*

# Example (cont.)

DFA operation

- Start in state $S_0$ and take transitions on each input character

- DFA **accepts** input string $\underline{x}$ *iff* $\underline{x}$ leaves it in an accepting state ($S_2$)



**Recognizer for Register**

So,

- $\underline{r17}$ takes it through $s_0$, $s_1$, $s_2$ and accepts

- $\underline{r}$ takes it through $s_0$, $s_1$ and fails

- $\underline{1}$ takes it straight to $s_e$

# Error State

Per convention, the error state, $s_e$, and transitions to it, are not drawn.

Drawing the error state $s_e$, and transitions to it would give us:



Remember: $\Sigma = \{\underline{r}, \underline{0}, \underline{1}, ..., \underline{9}\}$

# Finite Automata (FA)



A FA consists of a 5-tuple

$$\langle \Sigma , S , \delta , F , I \rangle$$

where

- $\Sigma$ is an alphabet
- $S$ is a finite set of states
- $\delta$ is a state transition function $\delta : S \times \Sigma \rightarrow S$
- $F$ is a set of final or accepting states
- $I$ is the start state

Example:

$\Sigma = \{ \underline{r}, \underline{0}, \underline{1}, ..., \underline{9} \}$

$S = \{ s_0 , s_1 , s_2 \}$

$\delta :$ takes a state and a symbol as input and returns the next state (see next slide).

$F = \{ s_2 \}$

$I = \{ s_0 \}$

# State Transition Function and Code

To be useful, recognizer must turn into code

| $\delta$ | $\underline{r}$ | $\underline{0},\underline{1},\underline{2},\underline{3},\underline{4},\underline{5},$ $\underline{6},\underline{7},\underline{8},\underline{9}$ | *All others* |
|---|---|---|---|
| $s_0$ | $s_1$ | $s_e$ | $s_e$ |
| $s_1$ | $s_e$ | $s_2$ | $s_e$ |
| $s_2$ | $s_e$ | $s_2$ | $s_e$ |
| $s_e$ | $s_e$ | $s_e$ | $s_e$ |

*Table encoding RE*

Char $\leftarrow$ *next character*
State $\leftarrow s_0$

while (Char $\neq$ $\underline{EOF}$)
　State $\leftarrow \delta$(State,Char)
　Char $\leftarrow$ *next character*

if (State is a final state )
　then report success
　else  report failure

*Skeleton recognizer*

# Non-deterministic Finite Automata (NFAs)

We can construct an FA for each RE, but...

What about an RE such as $( \underline{a} \mid \underline{b} )^*\, \underline{abb}$ ?



This is a little different than the FA before:

- $S_0$ has a transition on $\varepsilon$

   $\varepsilon$ does not consume any input!

- $S_1$ has two transitions on $\underline{a}$

This is a *non-deterministic finite automaton* (NFA)

Problem: state $s_1$ gives us two choices for character $\underline{a}$: staying in $s_1$ or going to $s_2$. We have to guess the correct transition if we want to reach the accepting state $s_4$. E.g., for string "abb": $s_0$->$s_1$->$s_2$->$s_3$->$s_4$ (accept). Guessing wrongly causes us to reject a valid string. E.g., for string "abb": $s_0$->$s_1$->$s_1$->$s_1$->$s_1$ (fail).

# Non-deterministic Finite Automata

- An NFA accepts a string *x* iff $\exists$ a path though the transition graph from $s_0$ to a final state such that the edge labels spell *x*

- Transitions on $\varepsilon$ consume no input

- To "run" the NFA, start in $s_0$ and *guess* the right transition at each step (if there is more than one transition for a given symbol)
  - Always guess correctly
  - If some sequence of correct guesses accepts x then accept

Why study NFAs?

- They are the key to automating the RE$\rightarrow$DFA construction

- We can glue together NFAs with $\varepsilon$-transitions

$$\boxed{NFA} \xrightarrow{\varepsilon} \boxed{NFA} \quad \textit{becomes an} \quad \boxed{NFA}$$

# Relationship between NFAs and DFAs

DFA is a special case of an NFA

- DFA has no ε transitions

- DFA's transition function is single-valued

  – not possible to have 2 transitions from state s on a symbol <u>a</u>

- Same rules will work

DFA can be simulated with an NFA

  – *Obviously*

NFA can be simulated with a DFA *(less obvious)*

- Simulate sets of possible states

- Possible exponential blowup in the state space

- Still, one state per character in the input stream

# Summary: NFAs and DFAs

A FA is a DFA if

- no state has an $\varepsilon$-transition, i.e., there is no transition on input $\varepsilon$

- for each state s and input symbol a, there is at most one edge labeled a leaving state s.

A FA is an NFA

- if it contains $\varepsilon$-transitions or

- if it has several possible transitions from a state s on a given input symbol a.

# Automating Scanner Construction

To convert a specification to code:

1   Write down the RE for the input language

2   Build a big NFA

3   Build the DFA that simulates the NFA

4   Minimize the number of states in the DFA

5   Generate the scanner code

Scanner generators

•    Lex and Flex work along these lines

•    Algorithms are well-known and well-understood

•    Key issue is interface to parser

•    You could build one in a weekend!

# Outline

- The role of a scanner ✔

- Scanner concepts

  - Tokens, Lexemes, Patterns ✔

- Regular Expressions & Automata

  - Definitions of REs, DFAs and NFAs ✔

  - REs→NFA (Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book)

  - NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon)

  - DFA→minimal-state DFA (state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)

- Scanner generators

# Automating Scanner Construction

RE$\rightarrow$ NFA  *(Thompson's construction)*

- Build an NFA for $\varepsilon$ and each symbol $\underline{a} \in \Sigma$ occurring in the RE

- Combine them with $\varepsilon$-moves

NFA $\rightarrow$ DFA *(subset construction)*

- Build the simulation of the NFA

*The Cycle of Constructions*

RE $\longrightarrow$ NFA $\longrightarrow$ DFA $\longrightarrow$ *minimal* **DFA**

DFA $\rightarrow$ Minimal DFA

- Hopcroft's algorithm

DFA $\rightarrow$RE *(Not part of the scanner construction)*

- All pairs, all paths problem
- Take the union of all paths from $s_0$ to an accepting state

# RE → NFA using Thompson's Construction

Key idea:

• construct NFA for each symbol and each operator of the RE

• Join them with ε moves in precedence order of RE operators

**NFA for a**

**NFA for b**

**NFA for ab**

**NFA for a | b**

**NFA for a\***

Note: when joining two NFAs, we renumber states so that $s_0$ is again the start state, and all state names are unique.

# Example of Thompson's Construction

Let's try **a ( b | c )***

1.  a, b, and c:

$s_0$ —**a**→ $s_1$    $s_0$ —**b**→ $s_1$    $s_0$ —**c**→ $s_1$

2.  b | c

$s_0$ —ε→ $s_1$ —**b**→ $s_2$ —ε→ $s_5$
$s_0$ —ε→ $s_3$ —**c**→ $s_4$ —ε→ $s_5$

3.  ( b | c )*

$s_0$ —ε→ $s_1$ —ε→ $s_2$ —**b**→ $s_3$ —ε→ $s_6$ —ε→ $s_7$
$s_1$ —ε→ $s_4$ —**c**→ $s_5$ —ε→ $s_6$
(ε from $s_6$ back to $s_1$)
(ε from $s_0$ to $s_7$)

34

# Example of Thompson's Construction (*cont.*)

4. <u>a</u> ( <u>b</u> | <u>c</u> )*



Of course, a human would design something simpler ...



> **But, we can automate the production of the more complex automaton ...**

# Thompson's Construction

- Syntax-driven
    - follows the structure of REs

- <span style="color:red">Inductive:</span>
    - the cases in the construction of the NFA follow the cases in the definition of REs.
    - larger NFAs are built from smaller constituent NFAs

- <span style="color:red">Important:</span> if a symbol <u>a</u> occurs several times in a RE r, then a separate NFA is constructed for each occurrence of <u>a</u>.

# Thompson's Construction

Inductive Base:

- For $\varepsilon$, construct the NFA

$$\boxed{s_0} \xrightarrow{\varepsilon} \boxed{\boxed{s_1}}$$

- For $\underline{a} \in \Sigma$ , construct the NFA

$$\boxed{s_0} \xrightarrow{\underline{a}} \boxed{\boxed{s_1}}$$

Thompson's Construction for automata follows the inductive definition of regular expressions from page 10.

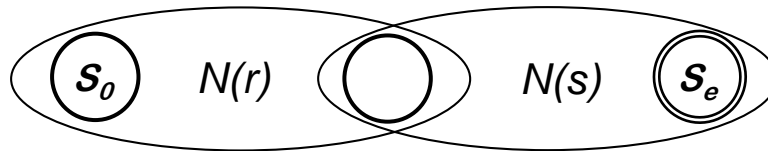Inductive step: suppose *N(r)* and *N(s)* are NFAs for REs *r* and *s*.
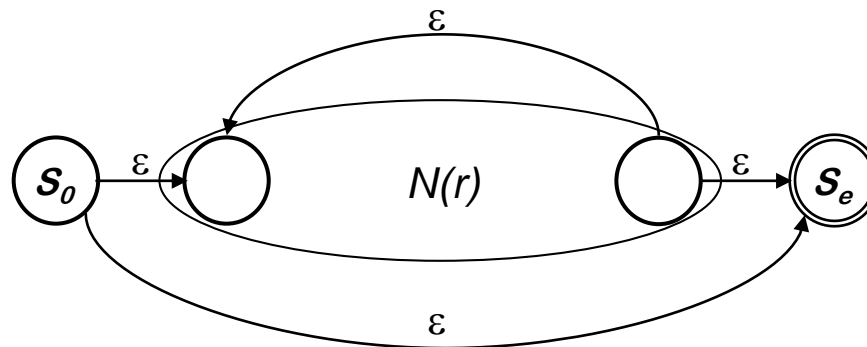
Then...  (continued on next slide).

# Thompson's Construction
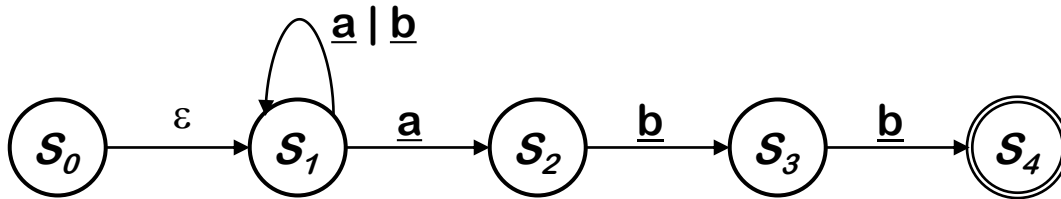
RE *r | s*:



RE *r s*:



RE *r\**:



Thompson's Construction for automata follows the inductive definition of regular expressions from page 10.

# Outline

- The role of a scanner ✔

- Scanner concepts

    - Tokens, Lexemes, Patterns ✔

- Regular Expressions & Automata

    - Definitions of REs, DFAs and NFAs ✔

    - REs→NFA (Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book) ✔

    - NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon)

    - DFA→minimal-state DFA (state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)

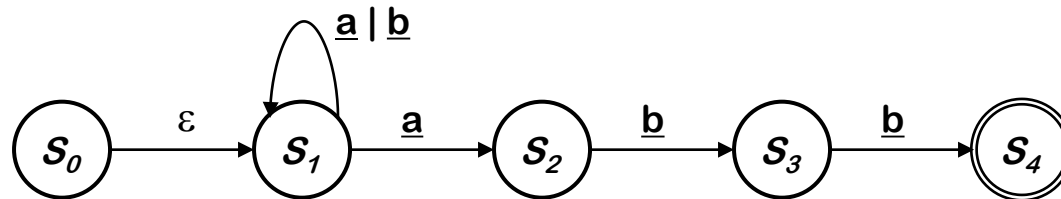- Scanner generators

# Subset Construction (NFA→DFA)



## RECALL:

Instead of guessing in state $s_1$ on symbol <u>a</u>, we can follow both transitions ($s_1 \to s_1$ and $s_1 \to s_2$) in parallel.

Problem: state $s_1$ gives us two choices for character <u>a</u>: staying in $s_1$ or going to $s_2$. We have to <u>guess</u> the correct transition if we want to reach the accepting state $s_4$. E.g., for string "abb": $s_0 \to s_1 \to s_2 \to s_3 \to s_4$ (accept). <u>Guessing wrongly</u> causes us to reject a valid string. E.g., for string "abb": $s_0 \to s_1 \to s_1 \to s_1 \to s_1$ (fail).

- We introduce a "virtual state" that combines the states that we reach from $s_1$ on symbol <u>a</u>.

- The new virtual state contains states $s_1$ and $s_2$, we write $\{s_1, s_2\}$ for the virtual state.

- **A question:** if we are in the virtual state $\{s_1, s_2\}$, where can we go when we read symbol <u>b</u>?

- **Answer:** the virtual state takes us anywhere that one of its member states takes us on symbol <u>b</u> (either $s_1$ or $s_3$ in the above NFA). So we introduce a second virtual state $\{s_1, s_3\}$.                (continued on next slide.)

# Subset Construction (NFA→DFA)



- A "virtual state" is a subset of the set of states S={$s_0,s_1,s_2,s_3,s_4$} of the NFA.

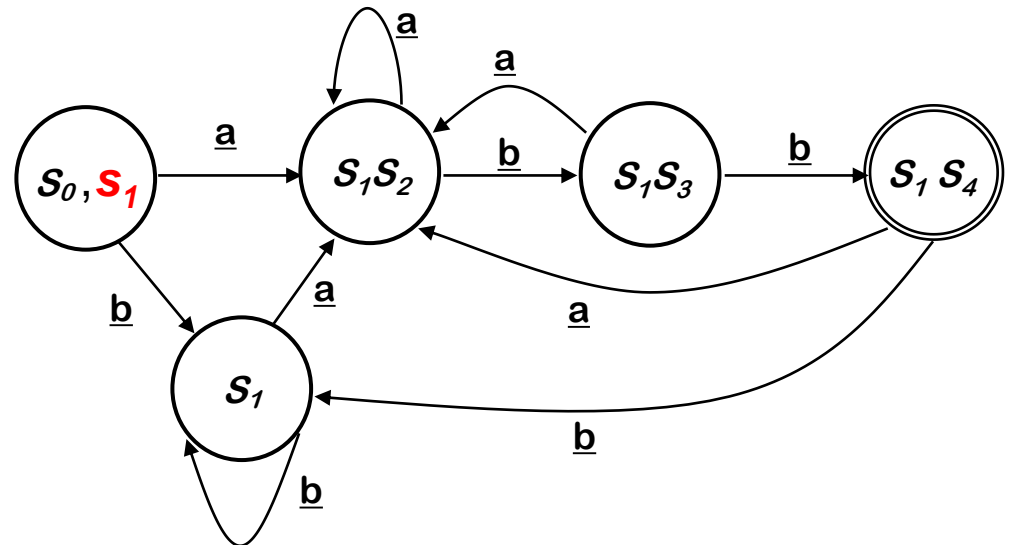| $\delta$ | $\underline{a}$ | $\underline{b}$ |
|---|---|---|
| {$s_0$, **$s_1$**} | {$s_1$, $s_2$} | {$s_1$} |
| {$s_1$, $s_2$} | {$s_1$, $s_2$} | {$s_1$, $s_3$} |
| {$s_1$} | {$s_1$, $s_2$} | {$s_1$} |
| {$s_1$, $s_3$} | {$s_1$, $s_2$} | {$s_1$, $s_4$} |
| {$s_1$, $s_4$} | {$s_1$, $s_2$} | {$s_1$} |

*Table encoding DFA*

41

# Algorithm: NFA $\rightarrow$ DFA with Subset Construction

We need to build a simulation of the NFA

Two key functions

- *Move($s_i$, a)* : gives set of states reachable from set $s_i$ by a

- $\varepsilon$-*closure($s_i$)* : gives set of states reachable from set $s_i$ by $\varepsilon$

The algorithm:

- Start state derived from $s_0$ of the NFA:

  – Take its $\varepsilon$-closure $S_0 = \varepsilon$-closure($\{s_0\}$)

- Take the image of $S_0$, Move($S_0$, $\alpha$), for each $\alpha \in \Sigma$, and take its $\varepsilon$-closure

- Iterate until no more states are added

Sounds more complex than it is...

...see also the algorithm animation in YSCEC...

# Algorithm: NFA $\rightarrow$ DFA with Subset Construction

**The algorithm:**

*Dstates $\leftarrow$ {};*

**add** $\varepsilon$-*closure($s_0$ ) as an*

  *unmarked state to Dstates;*

**while** *( there is an unmarked*
  *state T in Dstates ) {*

 **mark** *T;*
 **for each** $\alpha \in \Sigma$ {
    *U$\leftarrow$ $\varepsilon$-closure(Move(T,$\alpha$))*
    **if** *( U $\notin$ Dstates )* **then**
       **add** *U as an unmarked*
            *state to Dstates;*
       $\delta$ *[T,$\alpha$] $\leftarrow$ U*
 *}*
*}*

*Let's think about why this works*

**The algorithm terminates:**

1. ***Dstates*** **contains no duplicates (test before adding)**
2. $2^{|S|}$ **is finite**
3. **while loop adds to *Dstates*, but does not remove from *Dstates (monotone)***

$\Rightarrow$ **the loop halts**

**Dstates contains all the reachable NFA states:**

*It tries each character $\alpha \in \Sigma$ in each state T.*

*It builds every possible* **NFA** *configuration.*

$\Rightarrow$ ***Dstates and*** $\delta$ ***form the*** **DFA**

Any DFA state containing a final state from the NFA becomes a final state in the DFA.

43

# NFA →DFA with Subset Construction

Example of a *fixed-point* computation

• Monotone construction of some finite set

• Terminates when it stops adding to the set

• Proofs of termination & correctness are similar
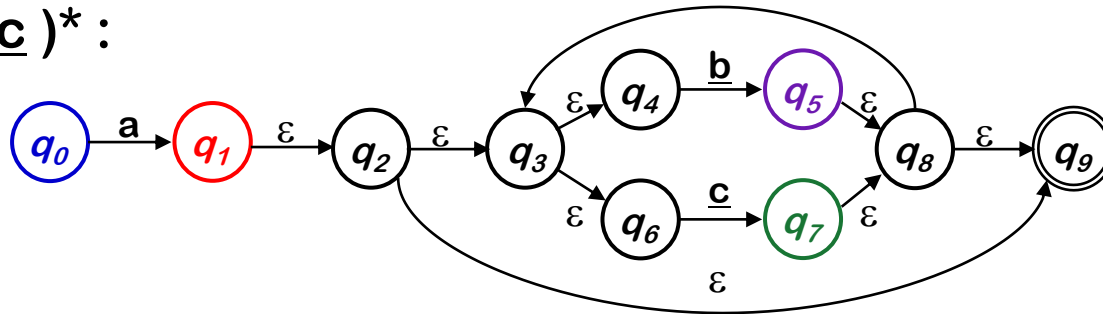
• These computations arise in many contexts

Other fixed-point computations

• Classic data-flow analysis (& Gaussian Elimination)

– Solving sets of simultaneous equations

*We will see many more fixed-point computations*

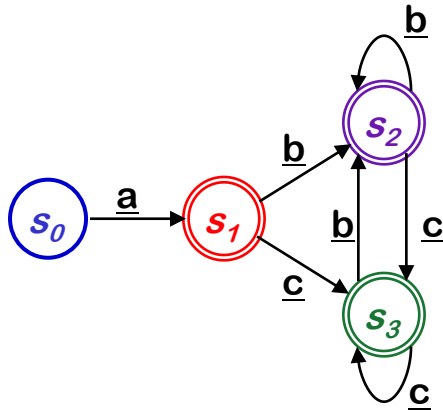# Example: NFA $\rightarrow$ DFA with Subset Construction

$\underline{a} ( \underline{b} | \underline{c} )^* :$



Applying the subset construction:

|  | NFA states | $\varepsilon$-closure (Move (s,*)) | | |
|---|---|---|---|---|
|  |  | **a** | **b** | **c** |
| $s_0$ | $q_0$ | $q_1, q_2, q_3,$ $q_4, q_6, q_9$ | none | none |
| $s_1$ | $q_1, q_2, q_3,$ $q_4, q_6, q_9$ | none | $q_5, q_8, q_9,$ $q_3, q_4, q_6$ | $q_7, q_8, q_9,$ $q_3, q_4, q_6$ |
| $s_2$ | $q_5, q_8, q_9,$ $q_3, q_4, q_6$ | none | $s_2$ | $s_3$ |
| $s_3$ | $q_7, q_8, q_9,$ $q_3, q_4, q_6$ | none | $s_2$ | $s_3$ |

**Final states**

45

# Example: NFA →DFA with Subset Construction

The DFA for $\underline{a}$ ( $\underline{b}$ | $\underline{c}$ )$^*$



| $\delta$ | $\underline{a}$ | $\underline{b}$ | $\underline{c}$ |
|---|---|---|---|
| $s_0$ | $s_1$ | - | - |
| $s_1$ | - | $s_2$ | $s_3$ |
| $s_2$ | - | $s_2$ | $s_3$ |
| $s_3$ | - | $s_2$ | $s_3$ |

- Ends up smaller than the NFA

- All transitions are deterministic

- Use same code skeleton as before

But, remember our goal:



46

# Outline

- The role of a scanner ✓

- Scanner concepts

  - Tokens, Lexemes, Patterns ✓

- Regular Expressions & Automata

  - Definitions of REs, DFAs and NFAs ✓

  - REs→NFA (Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book) ✓

  - NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon) ✓

  - DFA→minimal-state DFA (state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book)
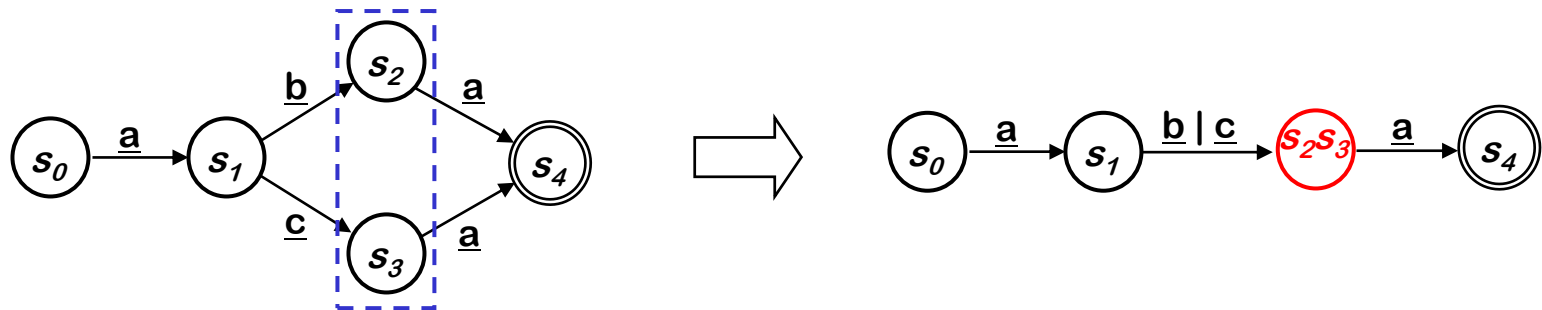
- Scanner generators

# DFA Minimization Overview

The Big Picture:

- Discover distinguishable states
- States that cannot be distinguished can be represented by a single state.

Two states p and q are distinguishable by input string **w** iff the DFA starting in p accepts w but starting in q does not accept w.

Two states p and q are distinguishable if they are distinguishable by some input string w.



In this example, $s_2$ and $s_3$ cannot be distinguished, therefore they can be represented by a single state $s_2 s_3$.

48

# DFA Minimization Overview

The set of states is divided into subsets of states that cannot be distinguished.

We say that the set of states $S$ is partitioned into partition $P$:

- Each state $s \in S$ is in exactly one set $p_i \in P$
- States in the same set have <u>not been distinguished yet</u>.
- States from different sets <u>are known to be distinguishable</u>.

**Step 1:**
Initially the partition $P$ consists of 2 sets:

- the accepting states $F$ and the non-accepting states $S - F$.
- States from F and S-F are distinguishable by $\varepsilon$.
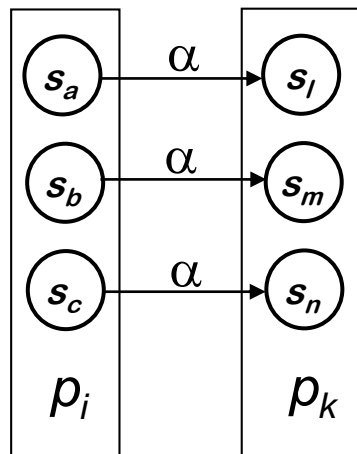
**Step 2:**
The minimization algorithm repeatedly picks a set $p_i \in P$ and tries to distinguish between its states by some symbol $\alpha \in \Sigma$.

$\rightarrow$ split $p_i$ along $\alpha$ (see next slide).

# Splitting of a State Set along Symbol $\alpha$

The minimization algorithm repeatedly picks a set $p_i \in P$ and tries
to distinguish between its states by some symbol $\alpha \in \Sigma$:



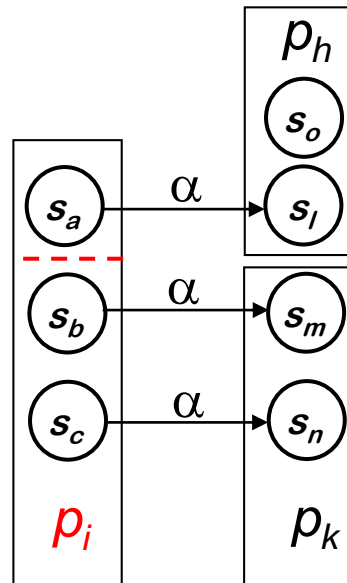$\alpha$ does not split $p_i$

$\alpha$ splits $p_i$ into $\{s_a\}$ and $\{s_{b,} s_c\}$

$\{s_a\}$ and $\{s_{b,} s_c\}$ are distinguishable by $\alpha$

Eventually no more set can be split and the algorithm terminates.

# DFA Minimization

The algorithm:

$P \leftarrow \{ F, S - F \}$

**while** *( P is still changing )*

   $T \leftarrow \{ \}$

   **for** each set $p \in P$

      $T \leftarrow T \cup$ Split *(p)*

   $P \leftarrow T$


Split *(p):*

   **for** each $\alpha \in \Sigma$

   {

      **if** $\alpha$ splits *p* into *p1* and *p2*

      **then return**  *{p1, p2}*

   }

   **return** *p;*

*This is a fixed-point algorithm!*

Why does this **terminate**?

- $p \in 2^S$ *(powerset)*
- Start off with 2 subsets of S:
  *F and S-F*
- **While** **loop** takes $P_i \rightarrow P_{i+1}$ by splitting 1 or more sets
- $P_{i+1}$ is at least one step closer to the partition with |S| sets
- Maximum of |S| splits

Note that

- Partitions are <u>never</u> combined, only split.
- → algorithm eventually terminates

51

# DFA Minimization

The algorithm:

$P \leftarrow \{ F, S\text{-}F \}$

**while** *( P is still changing )*
　$T \leftarrow \{ \}$
　**for** each set $p \in P$
　　$T \leftarrow T \cup \text{Split } (p)$
　$P \leftarrow T$

Split *(p)*
　**for** each $\alpha \in \Sigma$
　{
　　**if** $\alpha$ splits *p* into *p1* and *p2*
　　**then return**　*{p1, p2}*
　}
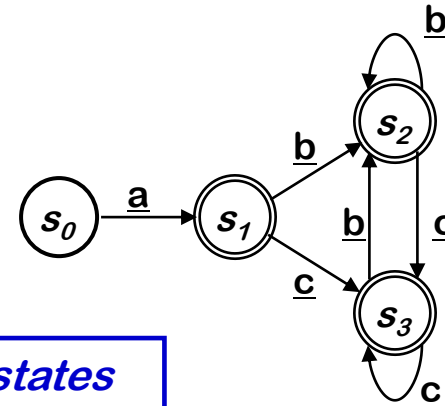　**return** *p;*

*This is a fixed-point algorithm!*

Why does this **work**?

- *The algorithm maintains*
  *2 invariants:*

  1) States remaining in the same set
     have not been distinguished yet
     by any string.

  2) States winding up in different sets
     are distinguishable by some
     string.

- The sets in the final partition
  contain the sets of distinguishable
  states of the DFA.

- Proof sketch: Dragon book, 2nd
  ed., Section 3.9.6.

# DFA Minimization Example
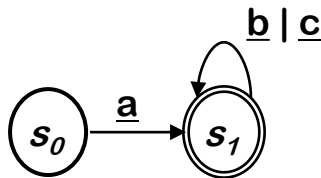
Applying the minimization algorithm...

| | | Split on | | |
|---|---|---|---|---|
| | Current Partition | a | b | c |
| $P_0$ | {$s_1$, $s_2$, $s_3$} {$s_0$} | none | none | none |



accepting states

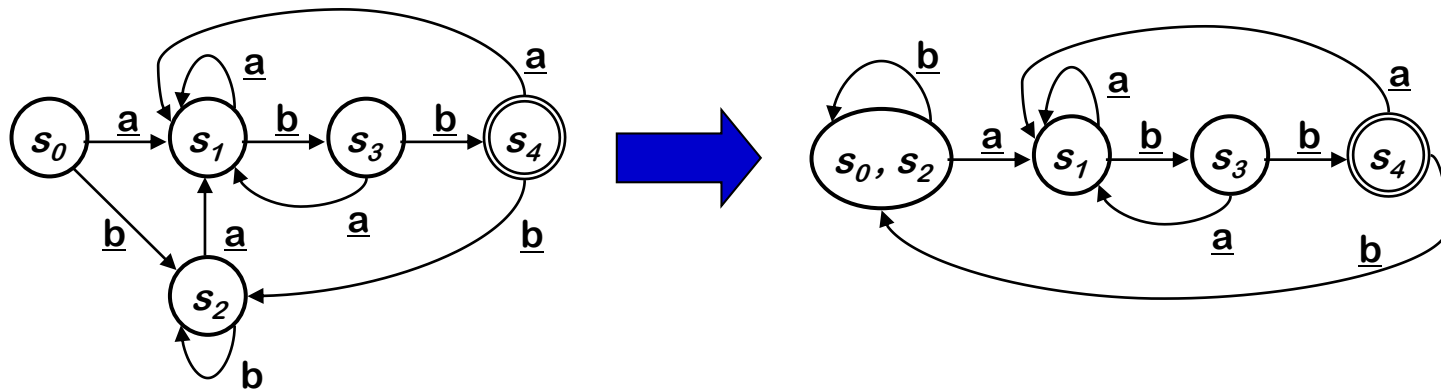...to produce the minimal DFA



In lecture 5, we observed that a human would design a simpler automaton than Thompson's construction & the subset construction did.

Minimizing that DFA produces the one that a human would design.

Theoretical result: every RE language can be recognized by a minimal-state DFA that is unique up to state names.

A detailed example: minimizing the DFA for **(a|b)*abb**

| Partition | Set | Split on <u>a</u> | Split on <u>b</u> |
|---|---|---|---|
| **P$_0$** | {s$_0$,s$_1$,s$_2$,s$_3$} | none | {s$_0$,s$_1$,s$_2$},{s$_3$} |
| | {s$_4$} | --- | --- |
| **P$_1$** | {s$_0$,s$_1$,s$_2$} | none | {s$_0$,s$_2$}, {s$_1$} |
| | {s$_3$} | --- | --- |
| | {s$_4$} | --- | --- |
| **P$_2$** | {s$_0$,s$_2$}, | none | none |
| | {s$_1$} | --- | --- |
| | {s$_3$} | --- | --- |
| | {s$_4$} | --- | --- |



54

# Outline

- The role of a scanner ✔

- Scanner concepts

    – Tokens, Lexemes, Patterns ✔

- Regular Expressions & Automata

    – Definitions of REs, DFAs and NFAs ✔

    – REs→NFA (Thompson's construction, Algorithm 3.3, Red Dragon book, Algorithm 3.23, Purple Dragon book) ✔

    – NFA→DFA (subset construction, Algorithm 3.2, Red Dragon book, Algorithm 3.20, Purple Dragon) ✔

    – DFA→minimal-state DFA (state minimization, Algorithm 3.6, Red Dragon book, Algorithm 3.39, Purple Dragon book) ✔

- Scanner generators

# Scanner Generators

Scanners generated in C:

•  lex (UNIX)

•  flex (GNU's f̲ast lex, UNIX)

•  mks lex (MS-DOS, Windows, OS/2)

Scanners generated in Java:

•  JLex (Princeton University)

•  JavaCC (Oracle)

# The Scanner Spec in JLex
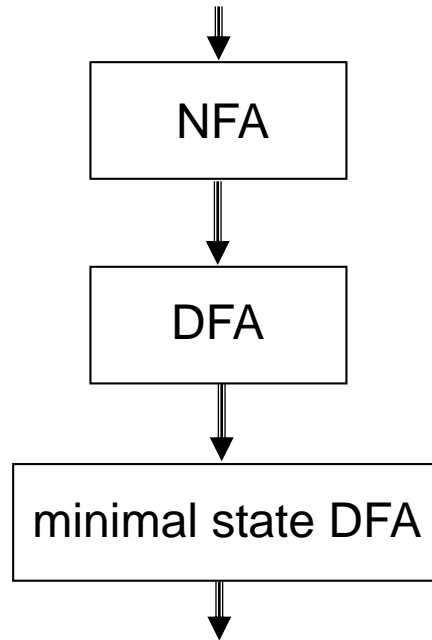
user code `// copied verbatim to the scanner file`

`%%`

`JLEX directives (declarations)`

`%%`

`regular expression rules`

# How a Scanner Generator Works

Specification of tokens using REs

```
┌─────────────────────┐
│        NFA          │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│        DFA          │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│  minimal state DFA  │
└─────────────────────┘
          │
          ▼
```

One of two possible DFA representations:

- Table-driven code (JLex), simulates a DFA on an input

- Hard-wired code (like our hand-crafted scanner for Assignment 1)

See Slide Add-on: Table-driven vs. handcrafted scanners.

# JLex Example Spec

```
%%
LETTER=[a-zA-Z_]
DIGIT=[0-9]
%%
"if" { return new Token(Token.IF, "if", src_pos); }
"<"  { return new Token(Token.LESS, "<", src_pos); }
"<="
  { return new Token(Token.LESS_EQ, "<=", src_pos); }
{LETTER}({LETTER}|{DIGIT})*
  { return new Token.ID, "spelling", src_pos); }
```
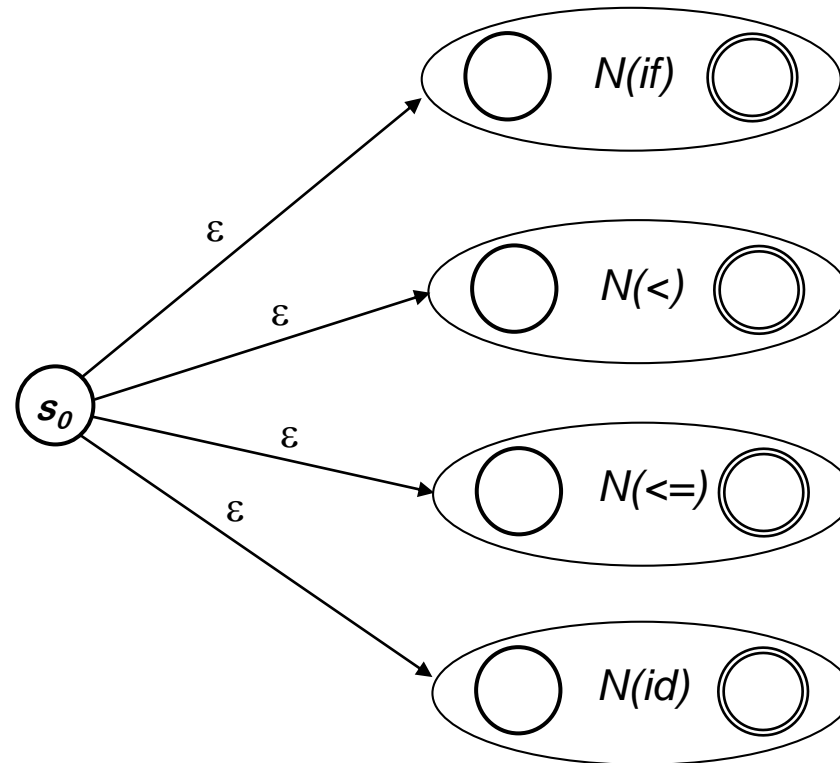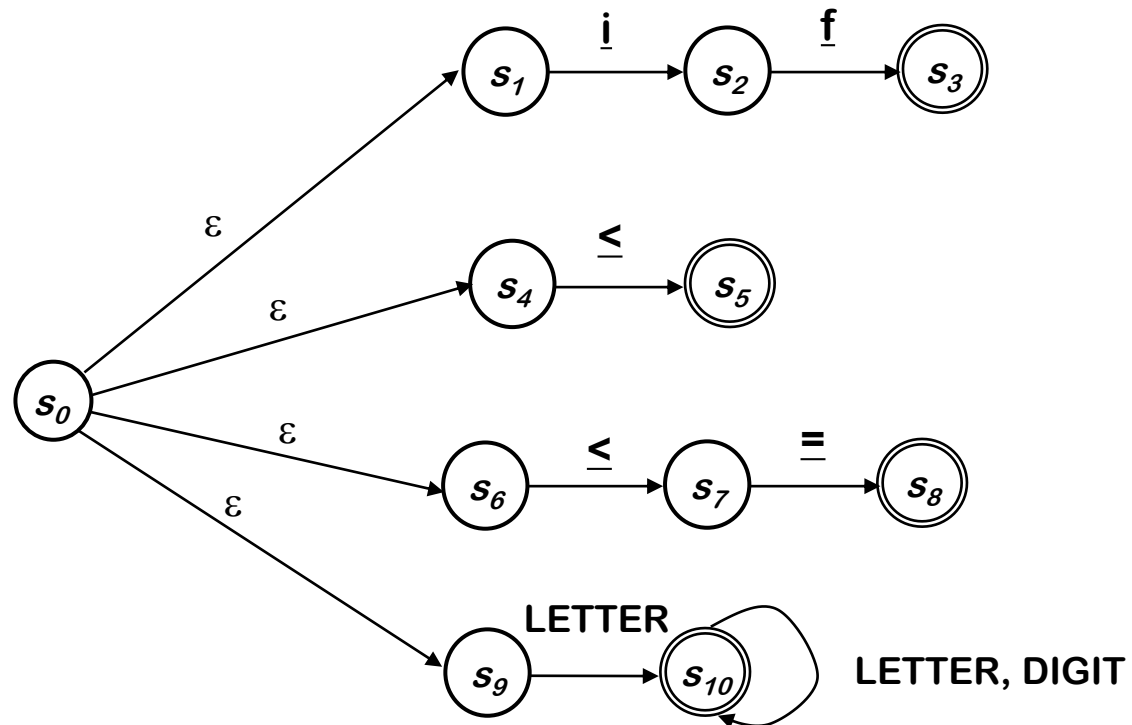
Two rules:

- The first pattern is used if more than one pattern matches.
  Example: `if` as a keyword and not as an ID.

- The longest possible match is taken.
  Example: "<=" as one token.

# NFA for JLex Example Spec



- The NFAs for the different REs are combined as above.
- Instead of an NFA, a DFA can also be used for each pattern.

# NFA for JLex Example Spec

# Running JLex on a Sample Scanner Spec

Scanner.l: the spec passed to the scanner generator JLex

`java JLex.Main Scanner.l`

```
Processing first section -- user code.
Processing second section -- JLex declarations.
Processing third section -- lexical rules.
Creating NFA machine representation.
NFA comprised of 315 states.
Working on character classes.::.::::.:..::..::.:...:.:...
NFA has 46 distinct character classes.
Creating DFA transition table.
Working on DFA states...............................
Minimizing DFA transition table.
109 states after removal of redundant states.
Outputting lexical analyzer code.
```

Scanner.l.java: the generated scanner

`javac Scanner.l.java`

# Limitations of Regular Languages

Advantages of Regular Expressions

- Simple & powerful notation for specifying patterns

- Automatic construction of fast recognizers (scanners)

- Many patterns can be specified with REs

Example ─ an expression grammar

> **Term:**    [a-zA-Z] ([a-zA-z] | [0-9])$^*$
>
> **Op:**      + | - | * | /
>
> **Expr :**   ( Term Op )$^*$ Term

Of course, this would generate a DFA ...

If REs are so useful ...

> *Why not use them for everything?*

# Limitations of Regular Languages (cont.)

Regular expressions are limited in what can be expressed.

– Example: it is not possible to specify a regular expression for $\underline{0^n1^n}$ (the set of strings of N zeroes followed by N ones).

For the same reason, we cannot specify the set of arithmetic expressions for which left and right parentheses match.
e.g., (a), ((a)), (((a))), ((((a))))

Regular expressions only usable to specify keywords, identifiers, literals, operators and punctuation characters of a programming language. For everything else (expressions, statements, nested statements, …) we need a more powerful mechanism.

# Regular Expressions and Context-Free Grammars

- The "languages" that can be defined by REs and Context-free grammars (CFGs) have been extensively studied by theoretical computer scientists. These are some important conclusions / terminologies:

  - RE are a "weaker" formalism than CFGs: Any language expressible by a RE can be expressed by a CFG **but not the other way around!**

  - The languages expressible as RE are called regular languages

  - Generally: a language that exhibits "self embedding" cannot be expressed by a RE.

  - Programming languages exhibit self embedding.
    Example: an expression can contain another expression

    ```
    expr ::=  id | integer | - expr | ( expr ) | expr op expr
    op   ::= + | - | * | /
    ```

How many expressions can be derived? Infinitely many.