

# Lec2: scanner (lexical analyzer)

- Recognizing words
- Formal grammars
- Regular expressions
- Finite automata

Content is copied from:

- <https://web.stanford.edu/class/cs143/lectures/lecture03.pdf>
- <https://web.stanford.edu/class/cs143/lectures/lecture04.pdf>
- Engineering a Compiler by Cooper and Torczon, 2nd Ed. ch. 1 and sec. 2.1-2.4
- <https://www3.nd.edu/~dthain/compilerbook/chapter3.pdf>

# Review: compiler structure

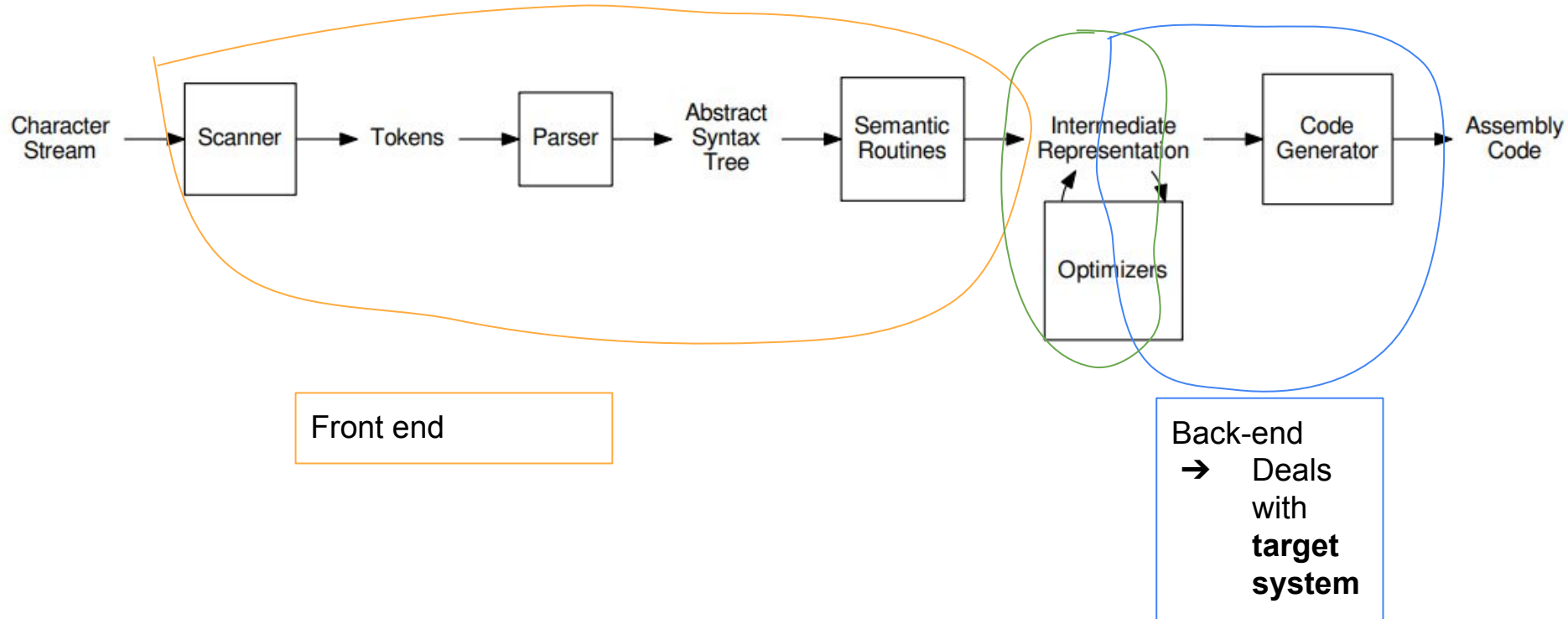


Fig2.2 with some modifications from the book <https://www3.nd.edu/~dthain/compilerbook/chapter2.pdf>

# Scanner (lexical analyzer)

What do we want to do? Example:

```
if (i == j)
    z = 0;
else
    z = 1;
```

→ `\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;`

**Goal:**

Partition this input string into substrings (tokens)

The input is just a string of characters:

→ `\tif (i == j)\n\t\tz = 0;\n\telse\n\t\tz = 1;`

# What is a token?

A syntactic **category**

– In English: noun, verb, adjective, ...

Token categories in any  
Programming Language:

- Keywords
- Identifiers
- Numbers
- Strings
- Comments and whitespace
- ...

Scanner identifies **tokens** from the raw text  
source code of a program.

# What are tokens for?

- Classify program substrings according to role
- Lexical analysis produces a stream of tokens
  - ... which is input to the parser
- Parser relies on token distinctions
  - An identifier is treated differently than a keyword

# Designing a Lexical Analyzer (or designing a new language)

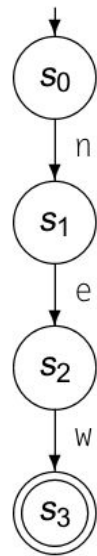
## Step-1: Define a finite set of tokens

- Tokens describe all items of interest
  - Identifiers,
  - integers,
  - keywords(reserved words)
- Choice of tokens depends on
  - language
  - design of parser

## Step-2: Describe which strings belong to each token

# Recognizing words: A hand-made scanner

```
c ← NextChar();  
if (c = 'n')  
  then begin;  
    c ← NextChar();  
    if (c = 'e')  
      then begin;  
        c ← NextChar();  
        if (c = 'w')  
          then report success;  
          else try something else;  
        end;  
      else try something else;  
    end;  
  else try something else;
```

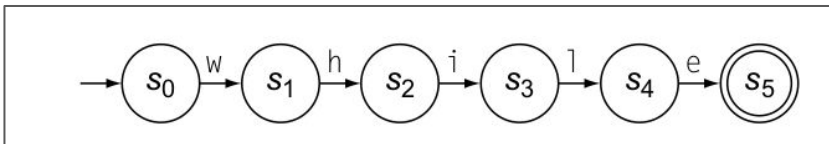


Transition  
to error  
states are  
omitted.

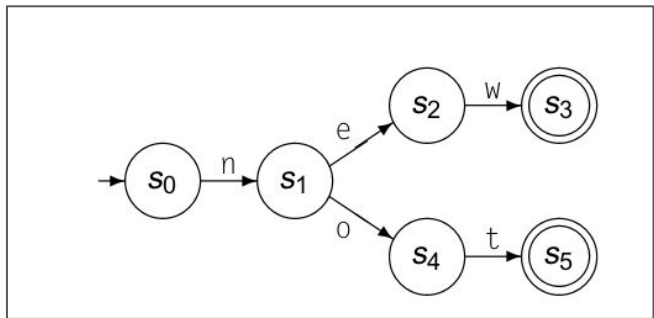
■ **FIGURE 2.1** Code Fragment to Recognize "new".

# Recognizing words: Other examples

while

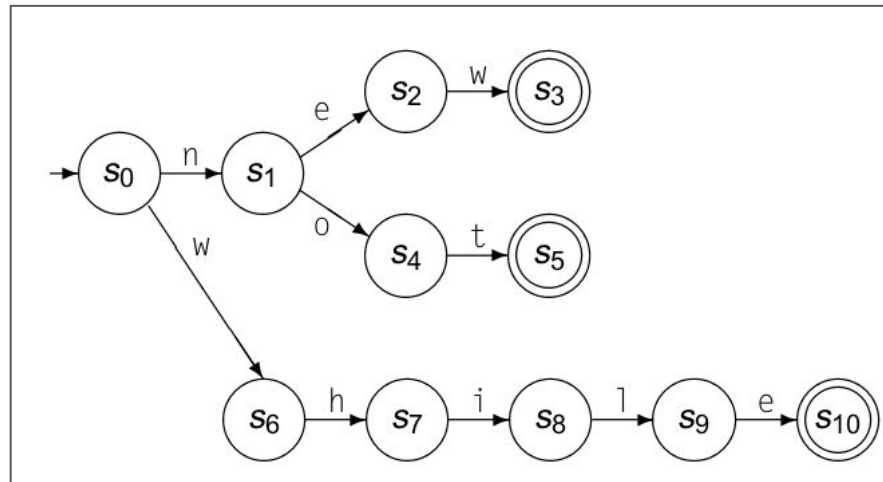


new and not



Combined version: while, new, not

- State s<sub>0</sub> has transitions for n and w





# Handmade scanner in C

```
token_t scan_token( FILE *fp ) {
    int c = fgetc(fp);
    if(c=='*') {
        return TOKEN_MULTIPLY;
    } else if(c=='!') {
        char d = fgetc(fp);
        if(d=='=') {
            return TOKEN_NOT_EQUAL;
        } else {
            ungetc(d,fp);
            return TOKEN_NOT;
        }
    } else if(isalpha(c)) {
        do {
            char d = fgetc(fp);
        } while(isalnum(d));
        ungetc(d,fp);
        return TOKEN_IDENTIFIER;
    } else if ( . . . ) {
        . . .
    }
}
```

**Figure 3.1: A Simple Hand Made Scanner**

# Is it as easy as it sounds?

Sort of... if you do not make it hard!

# Some history

## Lexical Analysis in FORTRAN

FORTRAN rule: Whitespace is insignificant

E.g., **VAR1** is the same as **VA R1**

- ★ A terrible design!
- ★ Historical footnote: FORTRAN Whitespace rule motivated by inaccuracy of punch card operators

Fortran do-loops

DO 5 I = 1,25

DO 5 I = 1.25

From left-to-right, it cannot tell if **DO5I** or **DO stmt .**

“**Lookahead**” may be required to decide where one token ends and the next token begins

- after “,” is reached, it can be determined.

# Lookahead

Even our simple example has lookahead issues

– `i` vs. `if`

– `=` vs. `==`

PL/I keywords are not reserved

**IF ELSE THEN THEN = ELSE; ELSE ELSE = THEN**

PL/I Declarations:

**DECLARE (ARG1,.. . ., ARGN)**

- Cannot tell whether DECLARE is a keyword or array reference until after the ).
  - Requires arbitrary **lookahead!**

# The problems continue today

## Lexical Analysis in C++

C++ template syntax:

```
Foo<Bar>
```

C++ stream syntax:

```
cin >> var;
```

But there is a conflict with nested templates:

```
Foo<Bar<Bazz>>
```

# Goal

The goal of lexical analysis is to

- Partition the input string into lexemes
  - Identify the token of each lexeme
- ★ Left-to-right scan => lookahead sometimes required

We need

- A way to describe the lexemes of each token
- A way to resolve ambiguities
  - Is **if** two variables **i** and **f**?
  - Is **==** two equal signs **= =**?

# Regular Languages

There are several formalisms for specifying tokens

- Regular languages are the most popular
  - Simple and useful theory
  - Easy to understand
  - Efficient implementations



# Languages

## Definition:

**Alphabet( $\Sigma$ )** is a set of characters (symbols).

**String ( $s$ ):** a finite, possibly empty sequence of symbols from an alphabet

**Language ( $L(s)$ )** over  $\Sigma$  is a set of strings of characters drawn from  $\Sigma$ .

Alphabet = English characters

Language = English sentences

★ Not every string of English characters is an English sentence

Alphabet = ASCII

→ Note: ASCII character set is different from English character set

Language = C programs

# Regular expressions

Need some notation for specifying which sets we want

The standard notation for **regular languages** is **regular expressions**.

- **s** is a string
- **L(s)** is the “language of **s**.”

**A regular expression s** is a string which denotes **L(s)**, a set of strings drawn from an alphabet  $\Sigma$ .

# Regular expression base case

- $s$  is a string
- $L(s)$  is the “language of  $s$ .”

**A regular expression  $s$**  is a string which denotes  $L(s)$ , a set of strings drawn from an alphabet  $\Sigma$ .

## Single character

- If  $a \in \Sigma$  then,
  - $a$  is a regular expression
  - and  $L(a) = \{a\}$

## $\epsilon$ is a regular expression

- $L(\epsilon)$  contains only the empty string,  $\{\epsilon\}$

# Regular expression built up rules

## Union (alternation)

$s|t$  is a RE such that

$$L(s|t) = L(s) \cup L(t).$$

Different notation

$$A + B = \{s \mid s \in A \text{ or } s \in B\}$$

## Concatenation

$st$  is a RE such that

- $L(st)$  contains all strings formed by the concatenation of a string in  $L(s)$  followed by a string in  $L(t)$ .

Different notation

$$AB = \{ab \mid a \in A \text{ and } b \in B\}$$

## Closure(Iteration)

The Kleene closure of a set  $A$ , denoted  $A^*$

$$A^* = \bigcup_{i \geq 0} A^i = A^0 \cup A^1 \cup A^2 \cup A^3 \cup A^4 \cup \dots$$

# Examples of closures

## Closure(Iteration)

The Kleene closure of a set A, denoted  $A^*$

$$A^* = \bigcup_{i \geq 0} A^i = A^0 \cup A^1 \cup A^2 \cup A^3 \cup A^4 \cup \dots$$

$$A^0 = \{\epsilon\}$$

$$A^1 = \{A\}$$

$$\{\text{"ab"}, \text{"c"}\}^* =$$

$\{\epsilon, \text{"ab"}, \text{"c"}, \text{"abab"}, \text{"abc"}, \text{"cab"}, \text{"cc"}, \text{"ababab"}, \text{"ababc"}, \text{"abccab"}, \text{"abcc"}, \text{"cabab"}, \text{"cabc"}, \text{"ccab"}, \text{"ccc"}, \dots\}.$

$$\{\text{"a"}, \text{"b"}, \text{"c"}\}^* =$$

$\{\epsilon, \text{"a"}, \text{"b"}, \text{"c"}, \text{"aa"}, \text{"ab"}, \text{"ac"}, \text{"ba"}, \text{"bb"}, \text{"bc"}, \text{"ca"}, \text{"cb"}, \text{"cc"}, \text{"aaa"}, \text{"aab"}, \dots\}.$

$$\emptyset^* = \{\epsilon\}.$$

# examples

| Regular Expression s | Language L(s)  |
|----------------------|--|
| <b>hello</b>         | <b>{hello}</b>   |
| <b>d(o i)g</b>       | <b>{dog, dig}</b>  |
| <b>moo*</b>          | <b>{mo, moo, mooo, ...}</b>                                  |
| <b>(moo)*</b>        | <b>{<math>\epsilon</math>, moo, moomoo, moomoomoo, ... }</b> |
| <b>a(b a)*a</b>      | <b>{aa, aaa, aba, aaaa, aaba, abaa, ... }</b>                |

# Abbreviations and some properties

$$s? = (s|\epsilon)$$

$$s+ = ss^*$$

$$[a-z] = (a|b|\dots|z)$$

$$[\hat{x}] \text{ is } \Sigma - x$$

- Union:  $A + B \equiv A | B$
- Option:  $A + \epsilon \equiv A?$
- Range:  $'a'+ 'b'+ \dots + 'z' \equiv [a-z]$
- Excluded range: complement of  $[a-z] \equiv [\hat{a-z}]$

|                       |                             |
|-----------------------|-----------------------------|
| <b>Associativity:</b> | $a   (b   c) = (a   b)   c$ |
| <b>Commutativity:</b> | $a   b = b   a$             |
| <b>Distribution:</b>  | $a (b   c) = ab   ac$       |
| <b>Idempotency:</b>   | $a ** = a *$                |

# exercises

`[0-9]+(.[0-9]+)?`

Which one matches?

- **123**
  - matches
- **3.14**
  - matches
- **.15**
  - Do not match
- **30.**
  - Do not match

`[A-Z]+([A-Z]|[0-9])*`

Which one matches?

- **PRINT**
  - matches
- **MODE5**
  - matches
- **hello**
  - Do not match
- **4YOU**
  - Do not match

`<[^>]*>`

Which one matches?

- **<tricky part>**
  - matches
- **<<<look left>**
  - matches
- **<this is an <illegal> comment>**
  - Does not match



# examples

| Regular Expression s                    | Language L(s) |
|---|---------------|
| <b>[abc]<sup>+</sup></b>                |               |
| <b>[abc]<sup>*</sup></b>                |               |
| <b>[0-9]<sup>+</sup></b>                |               |
| <b>[1-9][0-9]<sup>*</sup></b>           |               |
| <b>[a-zA-Z][a-zA-Z0-9_]<sup>*</sup></b> |               |

# Regular expression is used to describe grammars

## Keywords

“else” or “if” or “begin” or ...

‘else’ + ‘if’ + ‘begin’ + ...

Abbreviation: ‘else’ = ‘e’ ‘l’ ‘s’ ‘e’

**Integers:** an empty string of digits

```
digit = '0' + '1' + '2' + '3' + '4' + '5' + '6' + '7'  
+ '8' + '9'
```

```
integer = digit digit*
```

Abbreviation:  $A^+ = AA^*$

Abbreviation:  $[0-2] = '0' + '1' + '2'$

# Identifier

Identifier: strings of letters or digits, starting with a letter

letter = 'A' + ... + 'Z' + 'a' + ... + 'z'

identifier = letter (letter + digit)\*

Is **(letter\* + digit\*) = (letter + digit)\*** ?

# Whitespace

Whitespace: a **non-empty sequence of blanks, newlines, and tabs**

`(' ' + '\n' + '\t')+`

## Numeric constants

digit ::= [0-9]

digits ::= digit<sup>+</sup>

number ::= digits ( . digits )?

( [eE] ( + | - )? digits ) ?

## Phone Numbers

Consider (650)-723-3232

$\Sigma$  = digit + { -, (, ), }

exchange = digit<sup>3</sup>

phone = digit<sup>4</sup>

area = digit<sup>3</sup>

phone\_number = '(' area ')' '-' exchange '-' phone

# Recognizing REs

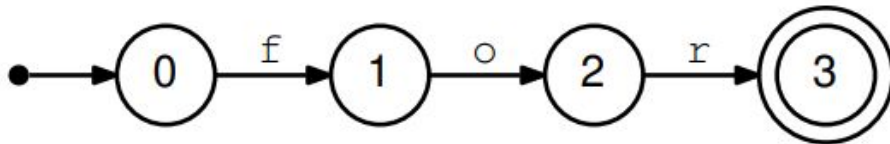
Given string **s** and  
rexp **R**,

is  $s \in L(R)$ ?

**Finite automata** can be used to recognize strings generated by regular expressions

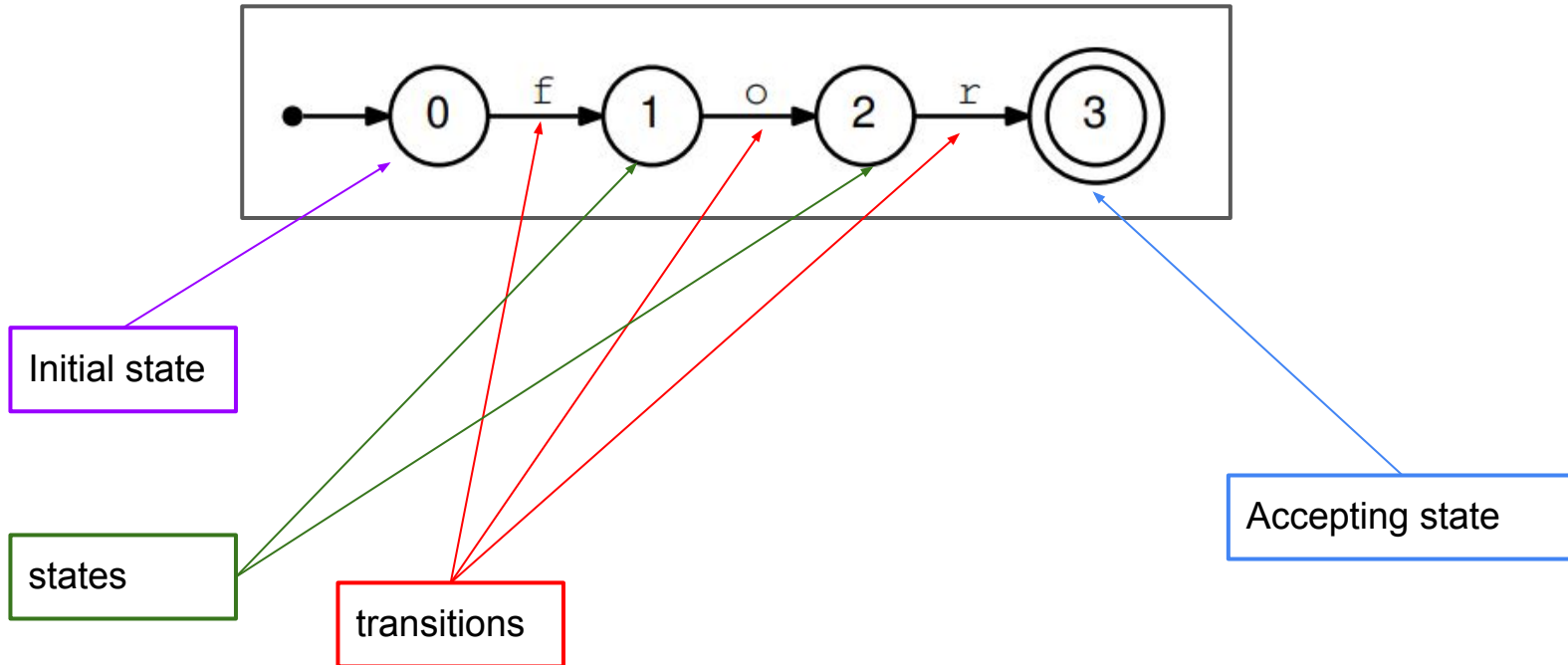
Can build by hand or automatically

- Reasonably straightforward, and can be done systematically
- Tools like **Lex**, **Flex**, **JFlex** etc this automatically, given a set of REs.
- Same techniques used in **grep**, **sed**, and other packages/tools



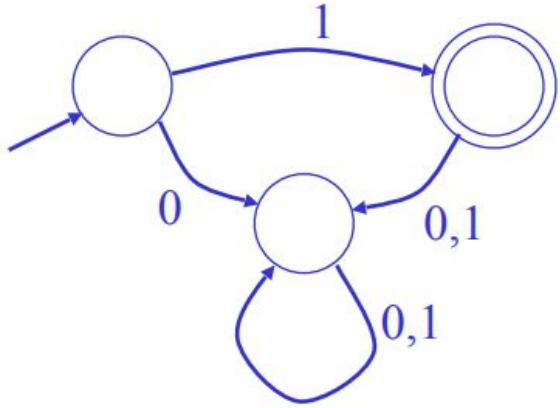
# Finite Automata state graphs

A finite automaton (FA) is an abstract machine that can be used to represent certain forms of computation.



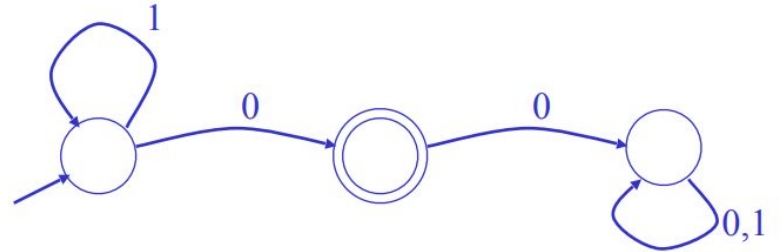
# A simple example

A finite automaton that accepts only “1”



A finite automaton accepting any number of 1's followed by a single 0

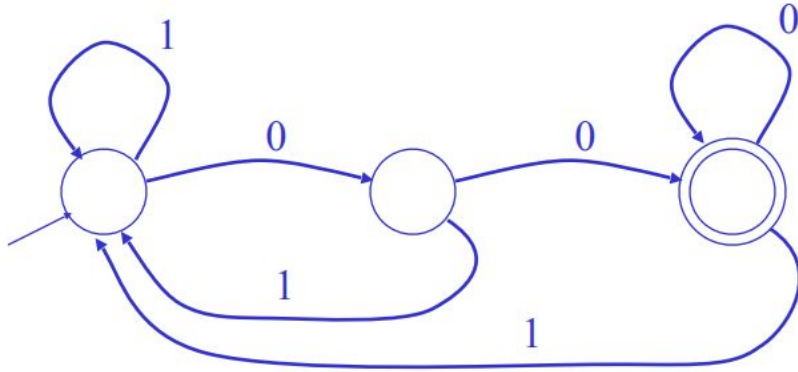
Alphabet:  $\{0,1\}$





Alphabet  $\{0,1\}$

What language does this recognize?



# Deterministic finite automata(DFA)

A DFA is a special case of an FA

- where **every state has no more than one outgoing edge for a given symbol.**

**A DFA has no ambiguity:**

for every combination of state and input symbol, there is exactly one choice of what to do next.

# Deterministic finite automata(DFA)

One integer (c) is needed to keep track of the current state.

The transitions between states are represented by a matrix M

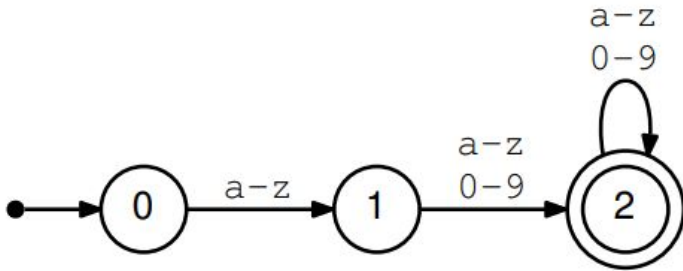
- $M[s, i]$  encodes the next state, given the current state  $s$  and input symbol  $i$ .
- **Error:** (If the transition is not allowed, we mark it with E to indicate an error.)

## Acceptance:

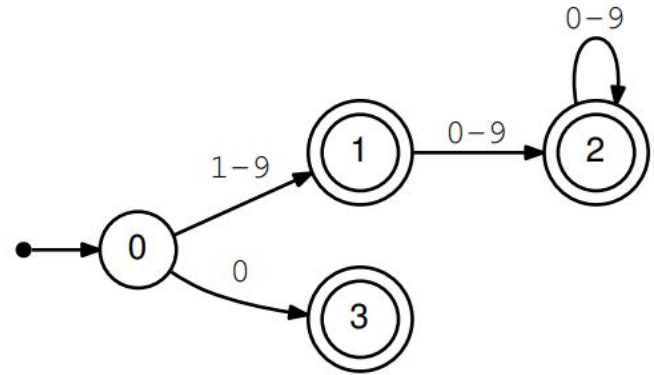
For each symbol,

- We compute  $c = M[s, i]$  until all the input is consumed, or an error state is reached.

[a-z][a-z0-9]+



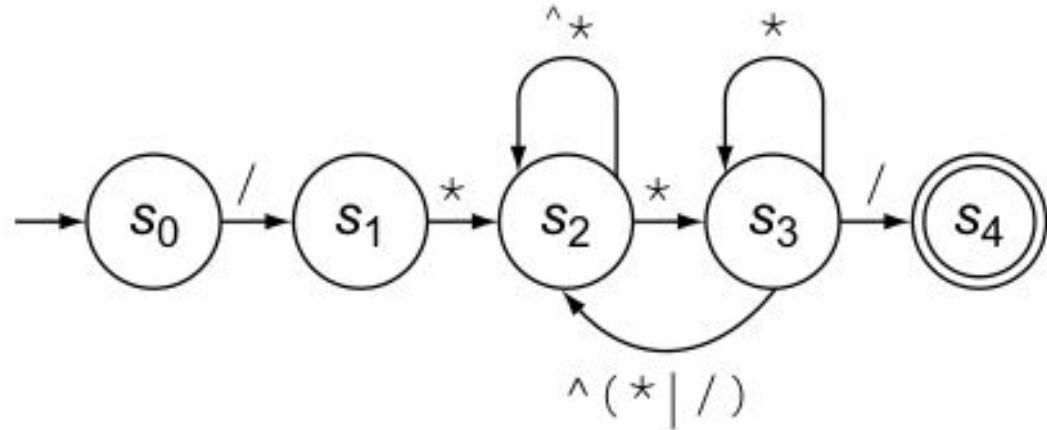
([1-9][0-9]\*)|0



# RE for comments in C++, Java

`/* (^*)* */` without `*`

`/* (^* | * + ^/)* */`

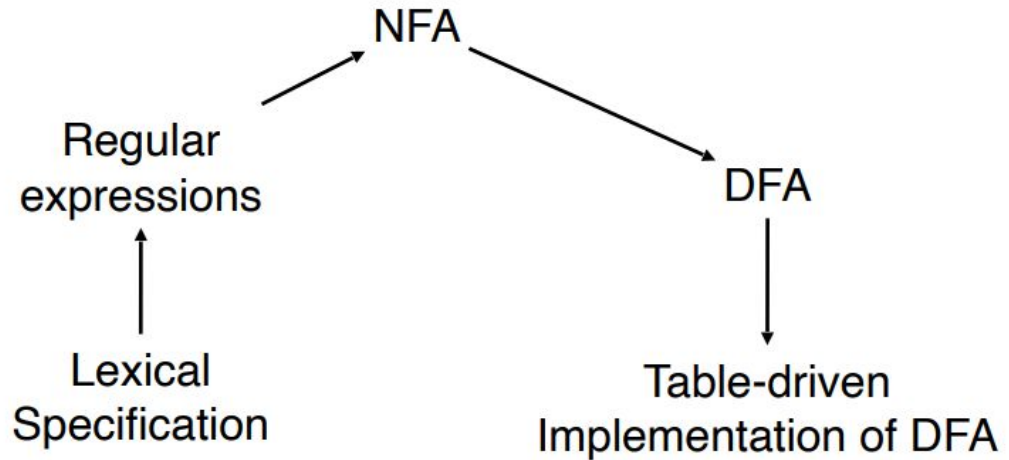


A more correct written form with escape char `\` and without nested comments:

`/\* ([^*]|\\*+[^*/])*\*+\/`

# Convert Regular Expressions to Finite Automata

## High level sketch



# Regular Expressions in Lexical Specification

a specification for the predicate

$$s \in L(R)$$

- **yes/no answer is not enough!**
- Instead: **partition the input into tokens**
- We will adapt regular expressions to this goal

# Lexical Specification → Regex in five steps

## 1. Write a regex for each token

- `Number = digit +`
- `Keyword = 'if' + 'else' + ...`
- `Identifier = letter (letter + digit)*`
- `OpenPar = '('`
- ...

## 2. Construct R, matching all lexemes for all tokens

`R = Keyword + Identifier + Number + ...`  
`= R1 + R2 + ...`

(This step is done automatically by tools like flex)



3. Let input be  $x_1 \dots x_n$

for  $1 \leq i \leq n$  check

$x_1 \dots x_i \in L(R)$

4. If success,

then we know that

$x_1 \dots x_i \in L(R_j)$  for some  $j$

5. Remove  $x_1 \dots x_i$  from input and go to

(3)

# Ambiguity

There are ambiguities in the algorithm

- How much input is used?

What **if**  $x_1 \dots x_i \in L(R)$  **and** also  $x_1 \dots x_k \in L(R)$

- Rule: Pick longest possible string in  $L(R)$ 
  - Pick **k** **if**  $k > i$
  - The “maximal munch”

# Ambiguity

Which token is used?

What **if**  $x_1 \dots x_i \in L(R_j)$  **and**  $x_1 \dots x_i \in L(R_k)$

- Rule: use rule listed first
  - Pick **j** **if**  $j < k$
  - E.g., treat "if" as a keyword, not an identifier

# Error handling

What if **No rule matches a prefix** of input?

- **Problem:** Can't just get stuck ...
- **Solution:**
  - Write a rule matching all "bad" strings
  - Put it **last** (lowest priority)

# Error handling

Regular expressions provide a concise notation for string patterns

- Use in lexical analysis requires small extensions
  - To resolve ambiguities
  - To handle errors
- Good algorithms
  - Require only single pass over the input
  - Few operations per character (table lookup)

# Finite automata again

Regular expressions = specification

Finite automata = implementation

## Deterministic Finite Automata (DFA)

- Exactly one transition per input per state
- No  $\epsilon$ -moves

## Nondeterministic Finite Automata (NFA)

- Can have zero, one, or multiple transitions for one input in a given state
- Can have  $\epsilon$ -moves

# HW1

Posted!