

# Compilation

0368-3133 (Semester A, 2013/14)

Lecture 2: Lexical Analysis

Modern Compiler Design: Chapter 2.1

Noam Rinetzky

1

# Admin

- Slides: [http://www.cs.tau.ac.il/~maon/...](http://www.cs.tau.ac.il/~maon/)
  - All info: [Moodle](#)
- Class room: Trubowicz 101 (Law school)
  - Except 5 Nov 2013
- Mobiles ...

2

# What is a Compiler?

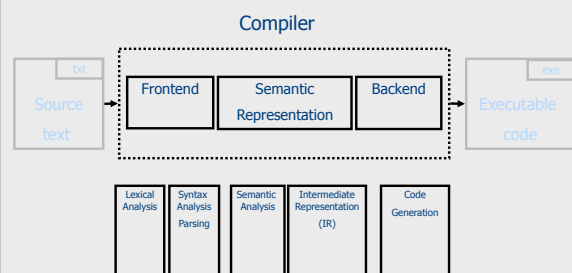
“A compiler is a **computer program** that **transforms** source code written in a programming language (**source language**) into another language (**target language**).

The most common reason for wanting to transform source code is to create an **executable program**.”

--Wikipedia

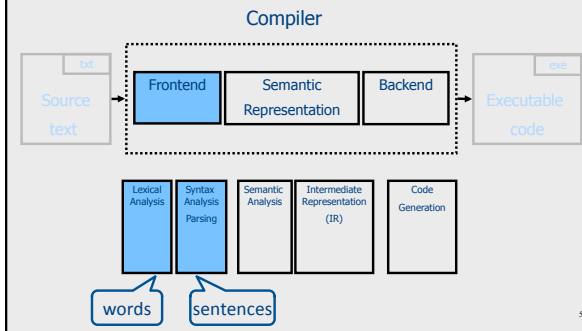
3

# Conceptual Structure of a Compiler



4

## Conceptual Structure of a Compiler



## What does Lexical Analysis do?

- Language: fully parenthesized expressions

$\text{Expr} \rightarrow \text{Num} \mid \text{LP Expr Op Expr RP}$

$\text{Num} \rightarrow \text{Dig} \mid \text{Dig Num}$

$\text{Dig} \rightarrow '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

$\text{LP} \rightarrow '('$

$\text{RP} \rightarrow ')'$

$\text{Op} \rightarrow '+' \mid '*'$

`( ( 23 + 7 ) * 19 )`

## What does Lexical Analysis do?

- Language: fully parenthesized expressions

Context free language  $\text{Expr} \rightarrow \text{Num} \mid \text{LP Expr Op Expr RP}$

Regular languages  $\text{Num} \rightarrow \text{Dig} \mid \text{Dig Num}$

Regular languages  $\text{Dig} \rightarrow '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

Regular languages  $\text{LP} \rightarrow '('$

Regular languages  $\text{RP} \rightarrow ')'$

Regular languages  $\text{Op} \rightarrow '+' \mid '*'$

`( ( 23 + 7 ) * 19 )`

## What does Lexical Analysis do?

- Language: fully parenthesized expressions

Context free language  $\text{Expr} \rightarrow \text{Num} \mid \text{LP Expr Op Expr RP}$

$\text{Num} \rightarrow \text{Dig} \mid \text{Dig Num}$

$\text{Dig} \rightarrow '0' \mid '1' \mid '2' \mid '3' \mid '4' \mid '5' \mid '6' \mid '7' \mid '8' \mid '9'$

$\text{LP} \rightarrow '('$

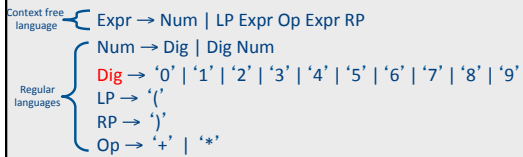
$\text{RP} \rightarrow ')'$

$\text{Op} \rightarrow '+' \mid '*'$

`( ( 23 + 7 ) * 19 )`

## What does Lexical Analysis do?

- Language: fully parenthesized expressions

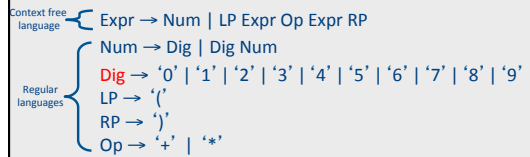


(	(	23	+	7	)	*	19	)
---	---	----	---	---	---	---	----	---

9

## What does Lexical Analysis do?

- Language: fully parenthesized expressions

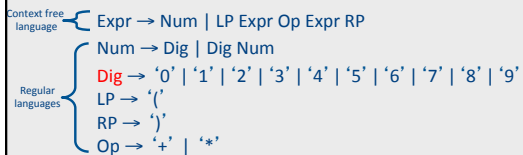


(	(	23	+	7	)	*	19	)
LP	LP	Num	Op	Num	RP	Op	Num	RP

10

## What does Lexical Analysis do?

- Language: fully parenthesized expressions

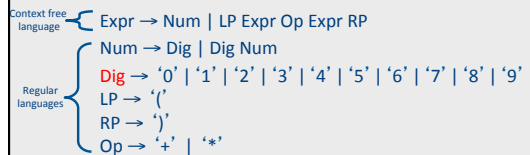


Value	(	(	23	+	7	)	*	19	)
Kind	LP	LP	Num	Op	Num	RP	Op	Num	RP

11

## What does Lexical Analysis do?

- Language: fully parenthesized expressions



Value	(	(	23	+	7	)	*	19	)
Kind	LP	LP	Num	Op	Num	RP	Op	Num	RP

12

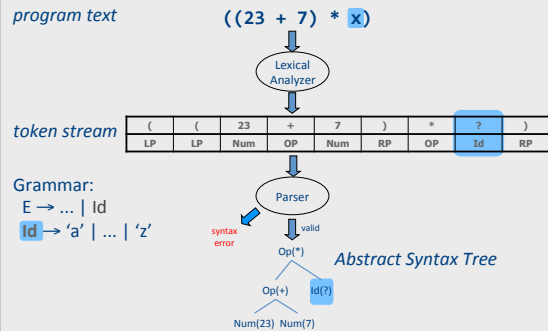
## What does Lexical Analysis do?

- Partitions the input into stream of **tokens**
  - Numbers
  - Identifiers
  - Keywords
  - Punctuation
- Usually represented as (kind, value) pairs
  - (Num, 23)
  - (Op, '\*')

• "word" in the source language  
• "meaningful" to the syntactical analysis

13

## From scanning to parsing



## Why Lexical Analysis?

- Simplifies the syntax analysis (parsing)
  - And language definition
- Modularity
- Reusability
- Efficiency

15

## Lecture goals

- Understand role & place of lexical analysis
- Lexical analysis theory
- Using program generating tools

16

## Lecture Outline

- ✓ Role & place of lexical analysis
- What is a token?
- Regular languages
- Lexical analysis
- Error handling
- Automatic creation of lexical analyzers

17

## What is a token? (Intuitively)

- A “word” in the source language
  - Anything that should appear in the input to syntax analysis
    - Identifiers
    - Values
    - Language keywords
- Usually, represented as a pair of (kind, value)

18

## Example Tokens

Type	Examples
ID	foo, n_14, last
NUM	73, 00, 517, 082
REAL	66.1, .5, 5.5e-10
IF	if
COMMA	,
NOTEQ	!=
LPAREN	(
RPAREN	)

19

## Example Non Tokens

Type	Examples
comment	/* ignored */
preprocessor directive	#include <foo.h>
	#define NUMS 5.6
macro	NUMS
whitespace	\t, \n, \b, ‘ ‘

20

## Some basic terminology

- **Lexeme** (aka symbol) - a series of letters separated from the rest of the program according to a convention (space, semi-column, comma, etc.)
- **Pattern** - a rule specifying a set of strings.  
Example: "an identifier is a string that starts with a letter and continues with letters and digits"  
– (Usually) a regular expression
- **Token** - a pair of (pattern, attributes)

21

## Example

```
void match0(char *s) /* find a zero */
{
    if (!strcmp(s, "0.0", 3))
        return 0. ;
}
```

```
VOID ID(match0) LPAREN CHAR Deref ID(s) RPAREN
LBRACE
IF LPAREN NOT ID(strcmp) LPAREN ID(s) COMMA STRING(0.0)
COMMA NUM(3) RPAREN RPAREN
RETURN REAL(0.0) SEMI
RBRACE
EOF
```

22

## Example Non Tokens

Type	Examples
comment	/* ignored */
preprocessor directive	#include <foo.h> #define NUMS 5.6
macro	NUMS
whitespace	\t, \n, \b, ' '

- Lexemes that are recognized but get consumed rather than transmitted to parser
  - if
  - i/\*comment\*/f

23

## How can we define tokens?

- Keywords – easy!
  - if, then, else, for, while, ...
- Identifiers?
- Numerical Values?
- Strings?
- Characterize **unbounded sets of values** using a **bounded description**?

24

## Lecture Outline

- ✓ Role & place of lexical analysis
- ✓ What is a token?
- Regular languages
- Lexical analysis
- Error handling
- Automatic creation of lexical analyzers

25

## Regular languages

- Formal languages
  - $\Sigma$  = finite set of letters
  - Word = sequence of letter
  - Language = set of words
- Regular languages defined equivalently by
  - Regular expressions
  - Finite-state automata

26

## Common format for reg-exps

Basic Patterns	Matching
x	The character x
.	Any character, usually except a new line
[xyz]	Any of the characters x,y,z
^x	Any character except x
<b>Repetition Operators</b>	
R?	An R or nothing (=optionally an R)
R*	Zero or more occurrences of R
R+	One or more occurrences of R
<b>Composition Operators</b>	
R1R2	An R1 followed by R2
R1 R2	Either an R1 or R2
<b>Grouping</b>	
(R)	R itself

27

## Examples

- $ab^*|cd?$  =
- $(a|b)^*$  =
- $(0|1|2|3|4|5|6|7|8|9)^*$  =

28

## Escape characters

- What is the expression for one or more + symbols?
  - `(+)+` won't work
  - `(\+)+` will
- backslash `\` before an operator turns it to standard character
  - `\*`, `\?`, `\+`, `a\b\+\*`, `(a\b\+\*)+`, ...
- backslash double quotes surrounds text
  - `"a(b+*)"`, `"a(b+*)"`

29

## Shorthands

- Use names for expressions
  - `letter = a | b | ... | z | A | B | ... | Z`
  - `letter_ = letter | _`
  - `digit = 0 | 1 | 2 | ... | 9`
  - `id = letter_ (letter_ | digit)*`
- Use hyphen to denote a range
  - `letter = a-z | A-Z`
  - `digit = 0-9`

30

## Examples

- `if = if`
- `then = then`
- `relop = < | > | <= | >= | = | <>`
  
- `digit = 0-9`
- `digits = digit+`

31

## Example

- A number is  

```
number = ( 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 )+
          ( ε | . ( 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 )+
            ( ε | E ( 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 )+
              )
          )
```
- Using shorthands it can be written as  

```
number = digits ( ε | .digits ( ε | E (ε|+|-) digits ) )
```

32



## Exercise 1 - Question

- Language of Java identifiers
  - Identifiers start with either an underscore '\_' or a letter
  - Continue with either underscore, letter, or digit

33

## Exercise 1 - Answer

- Language of Java identifiers
  - Identifiers start with either an underscore '\_' or a letter
  - Continue with either underscore, letter, or digit
- $(\_|a|b|\dots|z|A|\dots|Z)(\_|a|b|\dots|z|A|\dots|Z|0|\dots|9)^*$
- Using shorthand macros
  - First =  $\_|a|b|\dots|z|A|\dots|Z$
  - Next =  $\text{First}|0|\dots|9$
  - R =  $\text{First Next}^*$

34

## Exercise 1 - Question

- Language of rational numbers in decimal representation (no leading, ending zeros)
  - 0
  - 123.757
  - .933333
  - Not 007
  - Not 0.30

35

## Exercise 1 - Answer

- Language of rational numbers in decimal representation (no leading, ending zeros)
  - Digit =  $1|2|\dots|9$
  - Digit0 =  $0|Digit$
  - Num =  $Digit Digit0^*$
  - Frac =  $Digit0^* Digit$
  - Pos =  $Num | .Frac | 0.Frac | Num.Frac$
  - PosOrNeg =  $(\epsilon|-)Pos$
  - R =  $0 | PosOrNeg$

36

## Exercise 2 - Question

- Equal number of opening and closing parenthesis:  $[^n]^n = [], [[]], [[[]]], \dots$

37

## Exercise 2 - Answer

- Equal number of opening and closing parenthesis:  $[^n]^n = [], [[]], [[[]]], \dots$
- Not regular
- Context-free
- Grammar:  $S ::= [] \mid [S]$

38

## Challenge: Ambiguity

- `if = if`
- `id = letter_ (letter_ | digit)*`
- “if” is a valid word in the language of identifiers... so what should it be?
- How about the identifier “iffy”?
- Solution
  - Always find longest matching token
  - Break ties using order of definitions... first definition wins ( $\Rightarrow$  list rules for keywords before identifiers)

39

## Creating a lexical analyzer

- Given a list of token definitions (pattern name, regex), write a program such that
  - Input: String to be analyzed
  - Output: List of tokens
- How do we build an analyzer?

40

## Lecture Outline

- ✓ Role & place of lexical analysis
- ✓ What is a token?
- ✓ Regular languages
  - Lexical analysis
  - Error handling
  - Automatic creation of lexical analyzers

41

## A Simplified Scanner for C

```
Token nextToken()
{
    char c ;
    loop: c = getchar();
    switch (c){
        case '\n': goto loop ;
        case ';': return SemiColumn;
        case '+':
            c = getchar() ;
            switch (c) {
                case '+': return PlusPlus ;
                case '=': return PlusEqual;
                default: ungetc(c); return Plus;
            };
        case '<': ...
        case 'w': ...
    }
}
```

42

## But we have a much better way!

- Generate a lexical analyzer **automatically** from token definitions
- **Main idea:** Use finite-state automata to match regular expressions
  - Regular languages defined equivalently by
    - Regular expressions
    - Finite-state automata

43

## Reg-exp vs. automata

- **Regular expressions** are declarative
  - Offer compact way to define a regular language *by humans*
  - Don't offer direct way to check whether a given word is in the language
- **Automata** are operative
  - Define an *algorithm* for deciding whether a given word is in a regular language
  - Not a natural notation for humans

44

## Overview

- Define tokens using regular expressions
- Construct a nondeterministic finite-state automaton (NFA) from regular expression
- Determinize the NFA into a deterministic finite-state automaton (DFA)
- DFA can be directly used to identify tokens

45

## Finite-State Automata

- **Deterministic** automaton
- $M = (\Sigma, Q, \delta, q_0, F)$ 
  - $\Sigma$  - alphabet
  - $Q$  – finite set of state
  - $q_0 \in Q$  – initial state
  - $F \subseteq Q$  – final states
  - $\delta : Q \times \Sigma \rightarrow Q$  - transition function

46

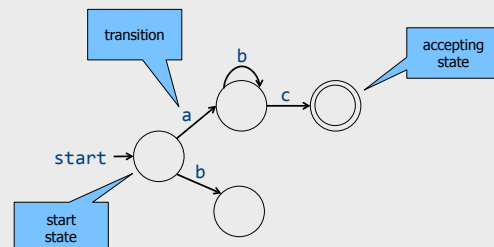
## Finite-State Automata

- **Non-Deterministic** automaton
- $M = (\Sigma, Q, \delta, q_0, F)$ 
  - $\Sigma$  - alphabet
  - $Q$  – finite set of state
  - $q_0 \in Q$  – initial state
  - $F \subseteq Q$  – final states
  - $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$  - transition function
- Possible  $\epsilon$ -transitions
- For a word  $w$ ,  $M$  can reach a number of states or get stuck. If some state reached is final,  $M$  accepts  $w$ .

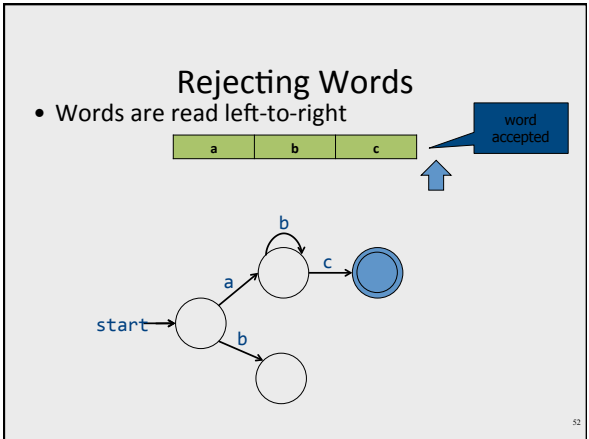
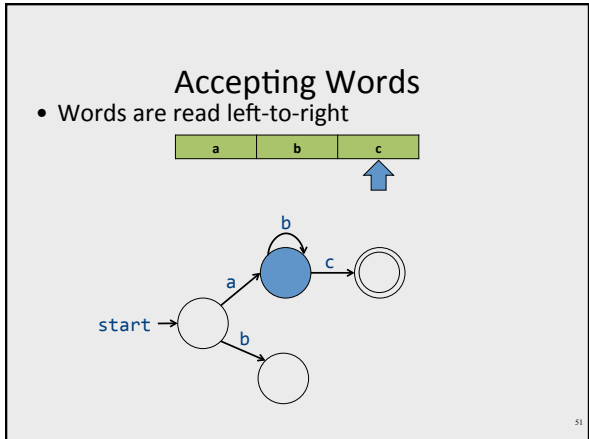
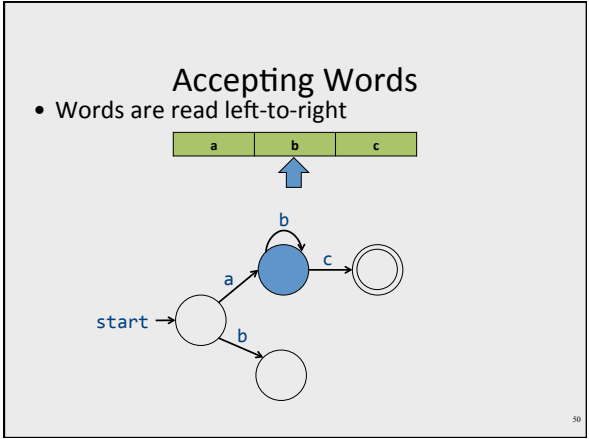
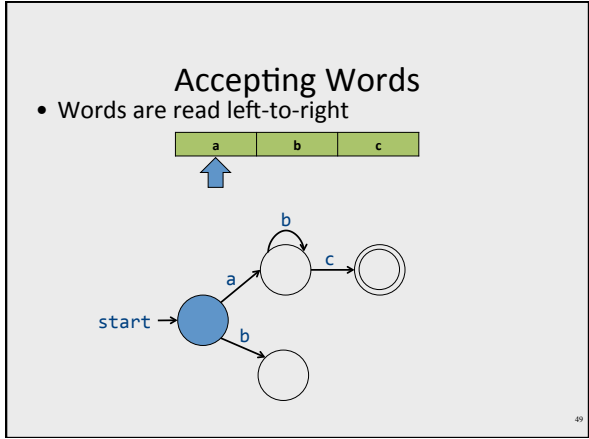
47

## Deterministic Finite automata

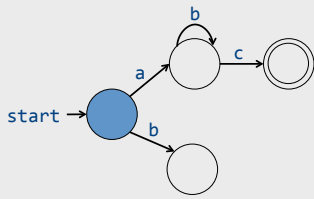
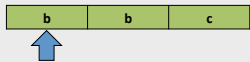
- An automaton is defined by states and transitions



48



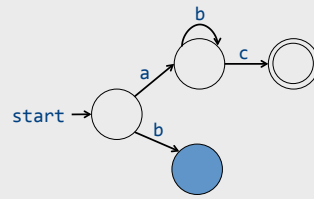
### Rejecting Words



53

### Rejecting Words

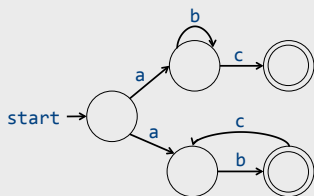
- Missing transition means non-acceptance



54

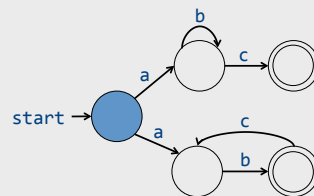
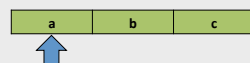
### Non-deterministic automata

- Allow multiple transitions from given state labeled by same letter

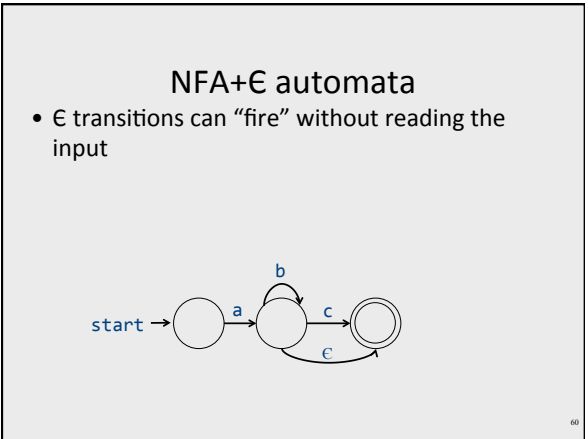
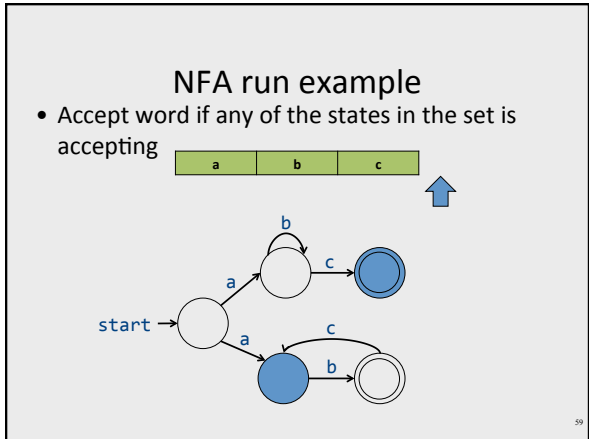
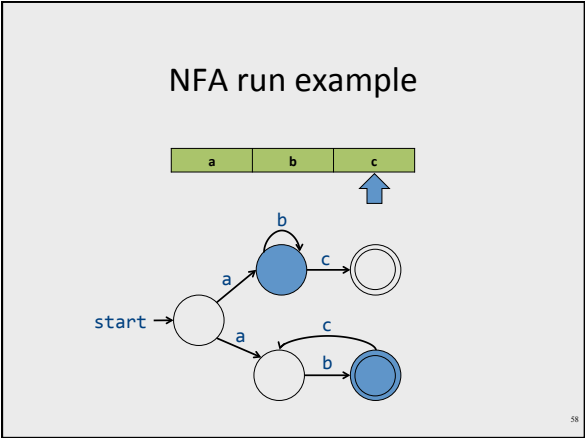
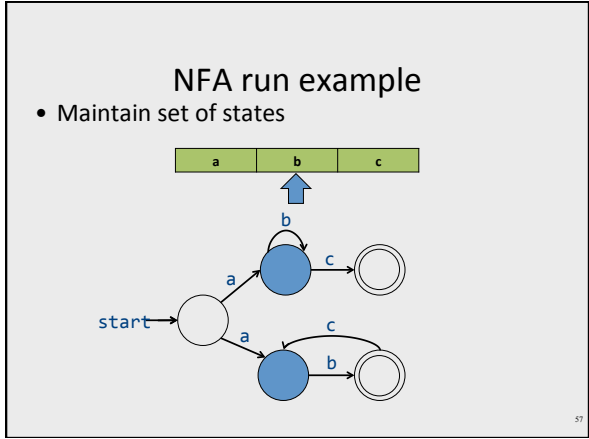


55

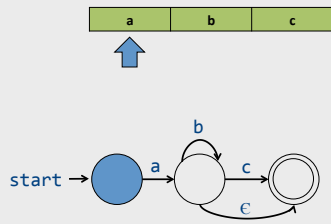
### NFA run example



56



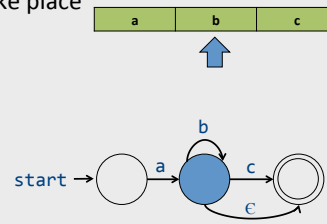
### NFA+ $\epsilon$ run example



61

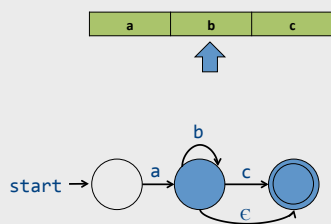
### NFA+ $\epsilon$ run example

- Now  $\epsilon$  transition can non-deterministically take place



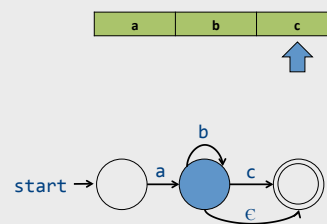
62

### NFA+ $\epsilon$ run example



63

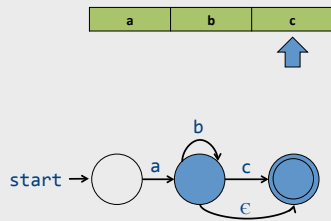
### NFA+ $\epsilon$ run example



64



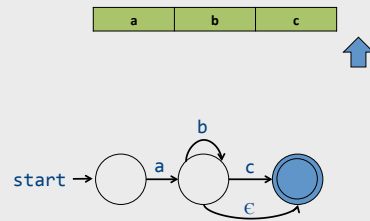
### NFA+ $\epsilon$ run example



65

### NFA+ $\epsilon$ run example

- Word accepted



66

### From regular expressions to NFA

- Step 1: assign expression names and obtain pure regular expressions  $R_1 \dots R_m$
- Step 2: construct an NFA  $M_i$  for each regular expression  $R_i$
- Step 3: combine all  $M_i$  into a single NFA
- Ambiguity resolution: prefer longest accepting word

67

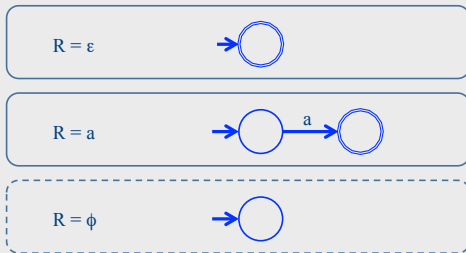
### From reg. exp. to automata

- Theorem: *there is an algorithm to build an NFA + $\epsilon$  automaton for any regular expression*
- Proof: *by induction on the structure of the regular expression*



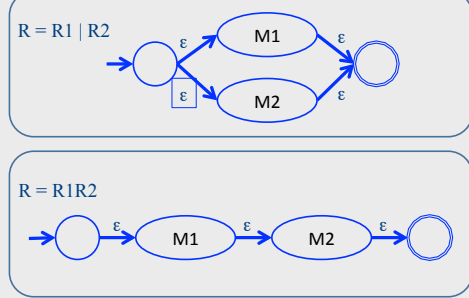
68

### Basic constructs



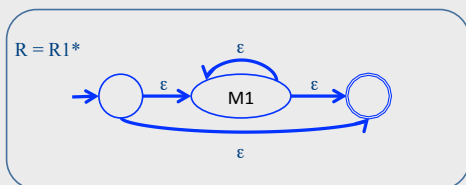
69

### Composition



70

### Repetition

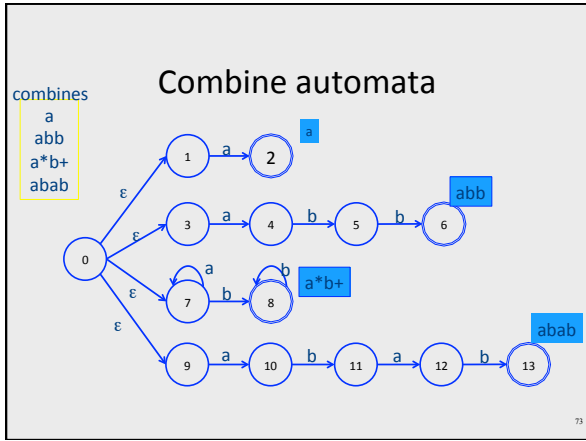


71

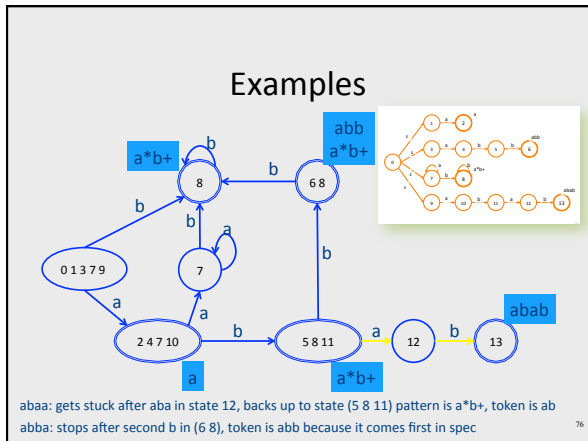
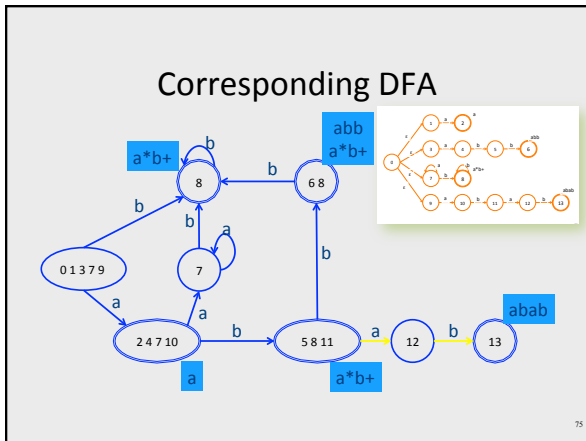
### What now?

- Naïve approach: try each automaton separately
- Given a word  $w$ :
  - Try  $M_1(w)$
  - Try  $M_2(w)$
  - ...
  - Try  $M_n(w)$
- Requires resetting after every attempt

72



- ### Ambiguity resolution
- Recall...
  - Longest word
  - Tie-breaker based on **order of rules** when words have same length
  - Recipe
    - Turn NFA to DFA
    - Run until stuck, remember last accepting state, this is the token to be returned
- 74



## Summary of Construction

- Describes tokens as regular expressions
  - and decides which attributes (values) are saved for each token
- Regular expressions turned into a DFA
  - describes expressions and specifies which attributes (values) to keep
- Lexical analyzer simulates the run of an automata with the given transition table on any input string

77

## A Few Remarks

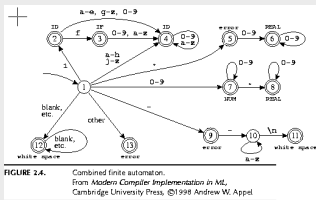
- Turning an NFA to a DFA is expensive, but
  - Exponential in the worst case
  - In practice, works fine
- The construction is done once per-language
  - At Compiler construction time
  - **Not** at compilation time

78

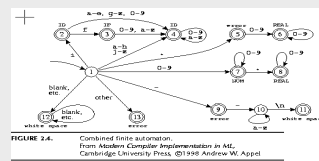
## Implementation by Example

```

if [a-z][a-z0-9]* { return IF; }
[a-z][a-z0-9]* { return ID; }
[0-9]+ { return NUM; }
[0-9]*"."[0-9]*|[0-9]*"."[0-9]+ { return REAL; }
(\\-\\-[a-z]*\\n)|(“ ”|\\n|\\t) { error(); }
    
```



79



```

int edges[][256]= {
/* state 0 */ {0, ..., 0, 0, 0, 0, 0, ..., 0, 0, 0, 0, 0, 0},
/* state 1 */ {13, ..., 7, 7, 7, 7, ..., 9, 4, 4, 4, 4, 2, 4, ..., 13, 13},
/* state 2 */ {0, ..., 4, 4, 4, 4, ..., 0, 4, 3, 4, 4, 4, ..., 0, 0},
/* state 3 */ {0, ..., 4, 4, 4, 4, ..., 0, 4, 4, 4, 4, 4, ..., 0, 0},
/* state 4 */ {0, ..., 4, 4, 4, 4, ..., 0, 4, 4, 4, 4, 4, ..., 0, 0},
/* state 5 */ {0, ..., 6, 6, 6, 6, ..., 0, 0, 0, 0, 0, 0, ..., 0, 0},
/* state 6 */ {0, ..., 6, 6, 6, 6, ..., 0, 0, 0, 0, 0, 0, ..., 0, 0},
/* state 7 */ {0, ..., 6, 6, 6, 6, ..., 0, 0, 0, 0, 0, 0, ..., 0, 0},
/* state ... */ {0, ..., 0, 0, 0, 0, 0, ..., 0, 0, 0, 0, 0, 0, ..., 0, 0},
/* state 13 */ {0, ..., 0, 0, 0, 0, 0, ..., 0, 0, 0, 0, 0, 0, ..., 0, 0}
};
    
```

80

## Pseudo Code for Scanner

```

Token nextToken()
{
    lastFinal = 0;
    currentState = 1;
    inputPositionAtLastFinal = input;
    currentPosition = input;
    while (not(isDead(currentState))) {
        nextState = edges[currentState][currentPosition];
        if (isFinal(nextState)) {
            lastFinal = nextState;
            inputPositionAtLastFinal = currentPosition;
        }
        currentState = nextState;
        advance currentPosition;
    }
    input = inputPositionAtLastFinal;
    return action[lastFinal];
}
    
```

81

## Example

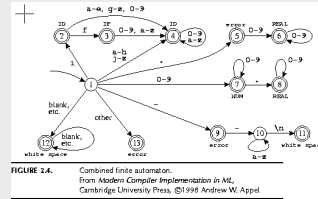


FIGURE 2.4. Combined finite automaton. From Modern Compiler Implementation in ML, Cambridge University Press, ©1998 Andrew W. Appel

Input: "if --not-a-com"

82

final	state	input
0	1	if --not-a-com
2	2	if --not-a-com
3	3	if --not-a-com
3	0	if --not-a-com

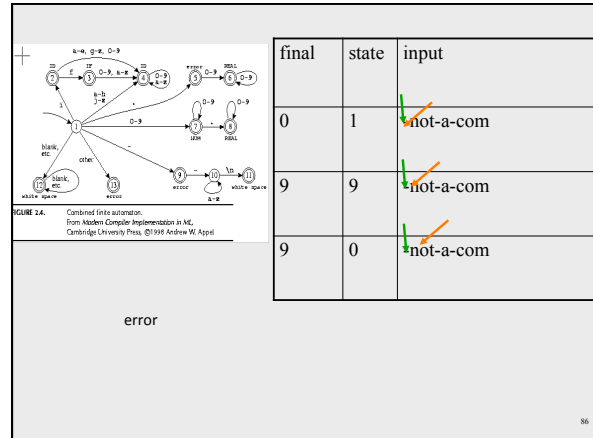
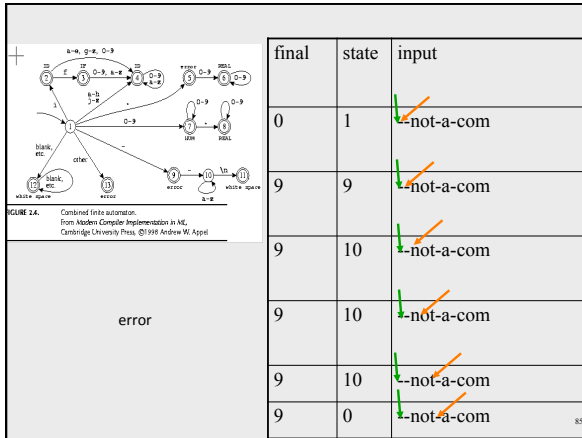
return IF

83

final	state	input
0	1	--not-a-com
12	12	--not-a-com
12	0	--not-a-com

found whitespace

84



### Efficient Scanners

- Efficient state representation
- Input buffering
- Using switch and gotos instead of tables

87

### DFSM from Specification

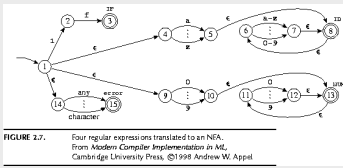
- Create a non-deterministic automaton (NFA) from every regular expression
- Merge all the automata using epsilon moves (like the | construction)
- Construct a deterministic finite automaton (DFA)
  - State priority
- Minimize the automaton starting
  - separate accepting states by token kinds

88

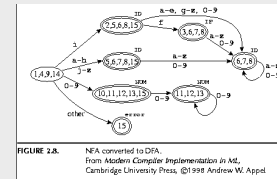
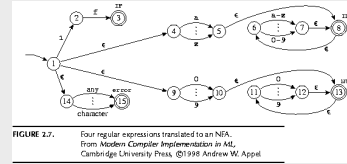
## NFA Construction

```

if { return IF; }
[a-z][a-z0-9]* { return ID; }
[0-9]+ { return NUM; }
    
```

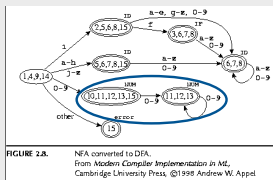


89



90

## Minimization



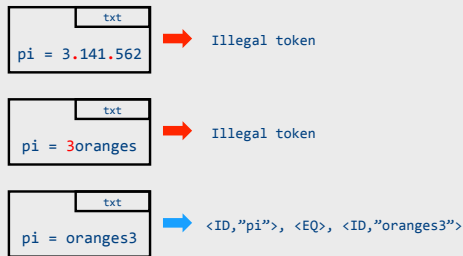
91

## Lecture Outline

- ✓ Role & place of lexical analysis
- ✓ What is a token?
- ✓ Regular languages
- ✓ Lexical analysis
  - Error handling
  - Automatic creation of lexical analyzers

92

## Errors in lexical analysis



93

## Error Handling

- Many errors cannot be identified at this stage
- Example: "fi (a==f(x))". Should "fi" be "if"? Or is it a routine name?
  - We will discover this later in the analysis
  - At this point, we just create an identifier token
- Sometimes the lexeme does not match any pattern
  - Easiest: eliminate letters until the beginning of a legitimate lexeme
  - Alternatives: eliminate/add/replace one letter, replace order of two adjacent letters, etc.
- Goal: allow the compilation to continue
- Problem: errors that spread all over

94

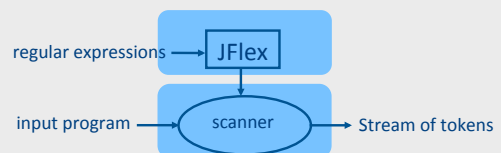
## Lecture Outline

- ✓ Role & place of lexical analysis
- ✓ What is a token?
- ✓ Regular languages
- ✓ Lexical analysis
- ✓ Error handling
- Automatic creation of lexical analyzers

95

## Use of Program-Generating Tools

- Automatically generated lexical analyzer
  - Specification → Part of compiler
  - Compiler-Compiler

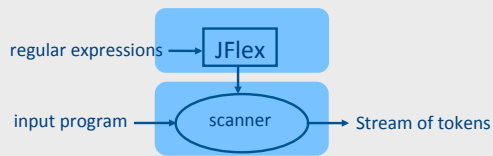


96



## Use of Program-Generating Tools

- Input: regular expressions and actions
  - Action = Java code
- Output: a scanner program that
  - Produces a stream of tokens
  - Invoke actions when pattern is matched



97

## Line Counting Example

- Create a program that counts the number of lines in a given input text file

98

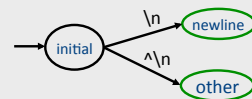
## Creating a Scanner using Flex

```
int num_lines = 0;
%%
\n      ++num_lines;
.       ;
%%
main() {
  yylex();
  printf( "# of lines = %d\n", num_lines);
}
```

99

## Creating a Scanner using Flex

```
int num_lines = 0;
%%
\n      ++num_lines;
.       ;
%%
main() {
  yylex();
  printf( "# of lines = %d\n", num_lines);
}
```



100

## JFlex Spec File

User code: Copied directly to Java file

%%

JFlex directives: macros, state names

%%

Lexical analysis rules:

- Optional state, regular expression, action
- How to break input to tokens
- Action when token matched

Possible source of javac errors down the road

```
DIGIT= [0-9]
LETTER= [a-zA-Z]
```

```
YYINITIAL
```

```
{LETTER}
({LETTER}|{DIGIT})*
```

101

## Creating a Scanner using JFlex

```
import java_cup.runtime.*;
%%
%cup
%{
    private int lineNumber = 0;
}%

%eofval{
    System.out.println("line number=" + lineNumber);
    return new Symbol(sym.EOF);
%eofval}

NEWLINE=\n
%%
{NEWLINE}      { lineNumber++; }
[^\{NEWLINE}]  { }
```

102

## Catching errors

- What if input doesn't match any token definition?
- Trick: Add a "catch-all" rule that matches any character and reports an error
  - Add after all other rules

103

## A JFlex specification of C Scanner

```
import java_cup.runtime.*;
%%
%cup
%{
    private int lineNumber = 0;
}%
Letter= [a-zA-Z_]
Digit= [0-9]
%%
"\t"      { }
"\n"      { lineNumber++; }
";"       { return new Symbol(sym.SemiColumn); }
"++"      { return new Symbol(sym.PlusPlus); }
"+="      { return new Symbol(sym.PlusEq); }
"+"       { return new Symbol(sym.Plus); }
"while"   { return new Symbol(sym.While); }
{Letter}{Letter}|{Digit}*
          { return new Symbol(sym.Id, yytext()); }
"<="      { return new Symbol(sym.LessOrEqual); }
"<"       { return new Symbol(sym.LessThan); }
```

104

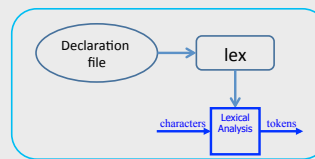
## A Simplified Scanner for C

```
Token nextToken()
{
    char c ;
    loop: c = getchar();
    switch (c){
        case '\n': goto loop ;
        case ';': return SemiColumn;
        case '+':
            c = getchar() ;
            switch (c) {
                case '+': return PlusPlus ;
                case '=': return PlusEqual;
                default: ungetc(c); return Plus;
            };
        case '<': ...
        case 'w': ...
    }
}
```

105

## Automatic Construction of Scanners

- Construction is done automatically by common tools
- lex is your friend
  - Automatically generates a lexical analyzer from declaration file
- Advantages: short declaration file, easily checked, easily modified and maintained



Intuitively:

- Lex builds DFA table
- Analyzer simulates (runs) the DFA on a given input

106

## Lecture Outline

- ✓ Role & place of lexical analysis
- ✓ What is a token?
- ✓ Regular languages
- ✓ Lexical analysis
- ✓ Error handling
- ✓ Automatic creation of lexical analyzers

107

## Start States

- It may be hard to specify regular expressions for certain constructs
  - Examples
    - Strings
    - Comments
- Writing automata may be easier
- Can combine both
- Specify partial automata with regular expressions on the edges
  - No need to specify all states
  - Different actions at different states

108

## Missing

- Creating a lexical analysis by hand
- Table compression
- Symbol Tables
- Nested Comments
- Handling Macros

109

## What does Lexical Analysis do?

- Input: program text (file)
- Output: sequence of tokens
  
- Read input file
- Identify language keywords and standard identifiers
- Handle include files and macros
- Count line numbers
- Remove whitespaces
- Report illegal symbols
  
- [Produce symbol table]

110

## Lexical Analysis – Impl.

- Lexical analyzer
  - Turns character stream into token stream
  - Tokens defined using regular expressions
  - Regular expressions -> NFA -> DFA construction for identifying tokens
  - Automated constructions of lexical analyzer using lex

111

## Automatic Construction of Scanners

- For most programming languages lexical analyzers can be easily constructed automatically
  - Exceptions:
    - Fortran
    - PL/1
- Lex/Flex/Jlex/JFlex are useful tools

112

## Next Week

- Syntax analysis (aka parsing)
- Lecture in **this room**
  - Trubowicz 101 (Law school)

113

114

## NFA vs. DFA

Automaton	SPACE	TIME
NFA	$O( r )$	$O( r ^* w )$
DFA	$O(2^{ r })$	$O( w )$

- $(a|b)^*a(a|b)(a|b)\dots(a|b)$   
n times

115

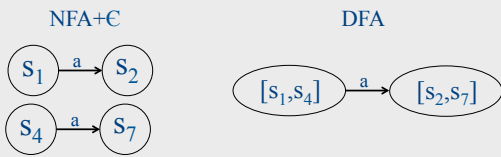
## From NFA+ $\epsilon$ to DFA

- Construction requires  $O(n)$  states for a reg-exp of length  $n$
- Running an NFA+ $\epsilon$  with  $n$  states on string of length  $m$  takes  $O(m \cdot n^2)$  time
- Solution: determinization via subset construction
  - Number of states worst-case exponential in  $n$
  - Running time  $O(m)$

116

## Subset construction

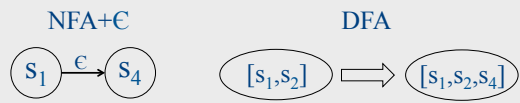
- For an NFA+ $\epsilon$  with states  $M=\{s_1, \dots, s_k\}$
- Construct a DFA with one state per set of states of the corresponding NFA
  - $M'=\{\emptyset, \{s_1\}, \{s_1, s_2\}, \{s_2, s_3\}, \{s_1, s_2, s_3\}, \dots\}$
- Simulate transitions between individual states for every letter



117

## Subset construction

- For an NFA+ $\epsilon$  with states  $M=\{s_1, \dots, s_k\}$
- Construct a DFA with one state per set of states of the corresponding NFA
  - $M'=\{\emptyset, \{s_1\}, \{s_1, s_2\}, \{s_2, s_3\}, \{s_1, s_2, s_3\}, \dots\}$
- Extend macro states by states reachable via  $\epsilon$  transitions



118