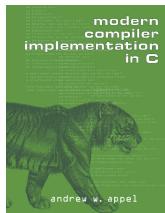


# Compilation

0368-3133

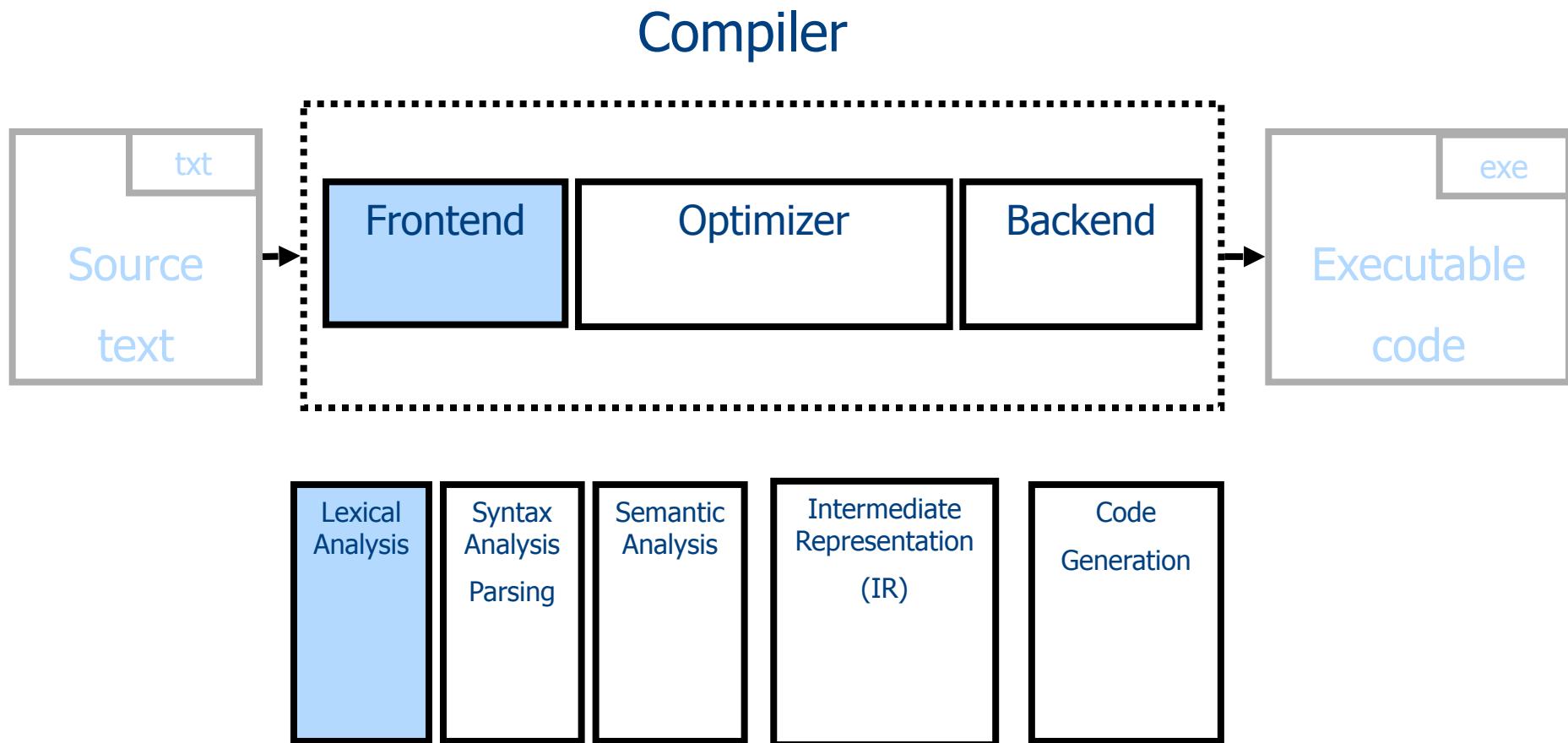
Lecture 1b:

# Lexical Analysis



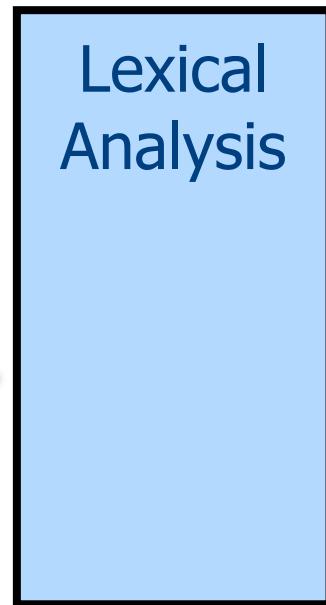
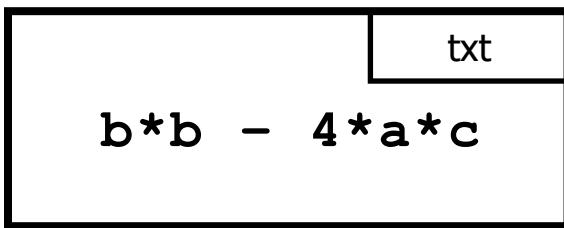
Modern Compiler Implementation in C: Chapter 2

# Conceptual Structure of a Compiler

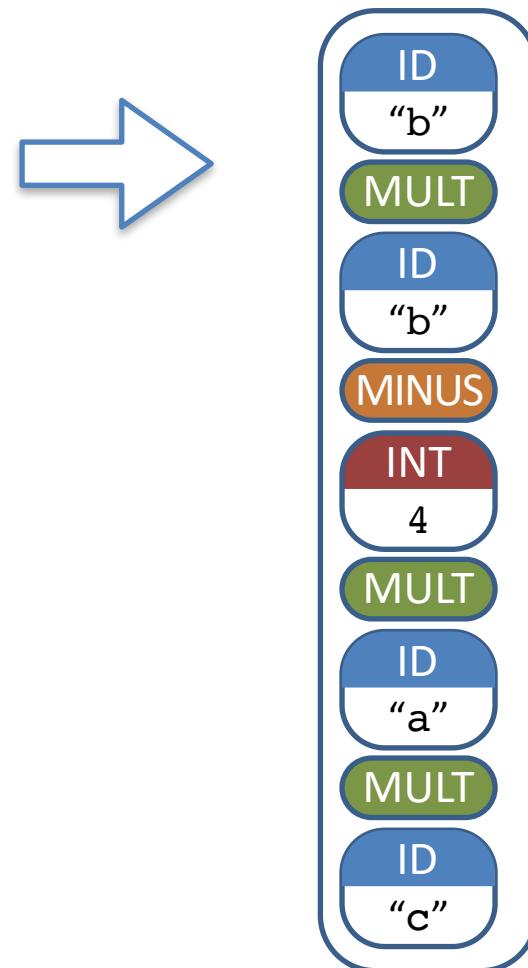


# What does Lexical Analysis do?

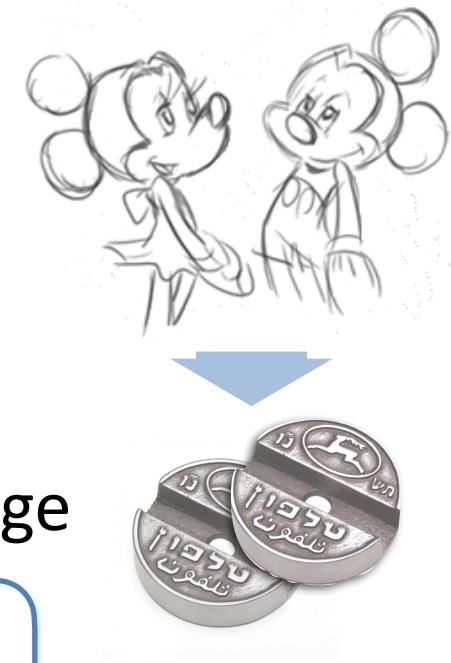
Stream of characters  
(program source code)



Stream of tokens  
("meaningful words")



# From Characters to Tokens



- What is a token?
  - Roughly – a “word” in the source language

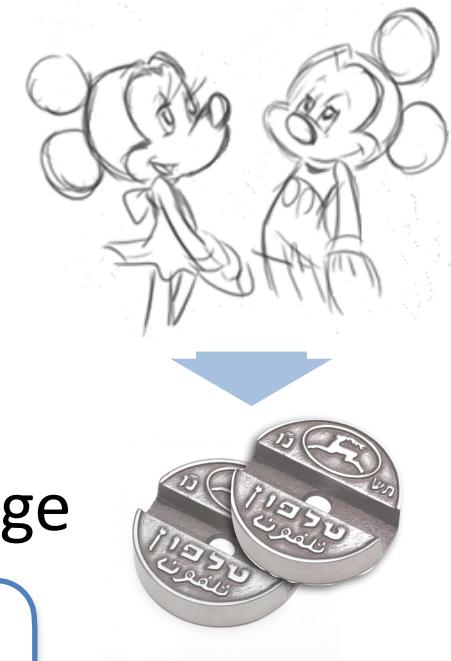
Identifiers

Values

Keywords

- Really – anything that should appear in the input to syntax analysis as a single unit
- Technically
  - Usually a pair of <kind, value>

# From Characters to Tokens



- What is a token?
  - Roughly – a “word” in the source language
    - Identifiers
    - Values
    - Keywords
  - Really – anything that should appear in the input to syntax analysis as a single unit
- Technically (in "real" compilers)
  - Usually a pair of <kind, value, line, column>

# Example Tokens

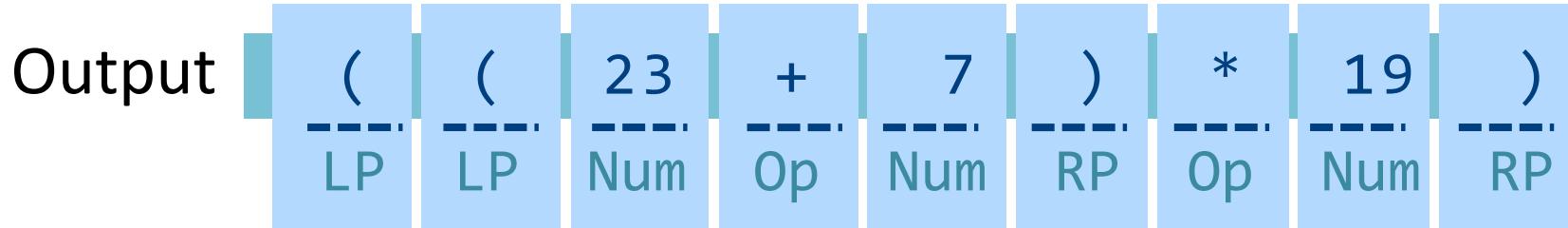
Kind	Examples
ID	<b>x, y, z0, foo, bar</b>
NUM	<b>42</b>
FLOATNUM	<b>3.141592654</b>
STRING	<b>"so long, and thanks for all the fish"</b>
IF	<b>if</b>
LPAREN	<b>(</b>
RPAREN	<b>)</b>
MINUS	<b>-</b>

# Example

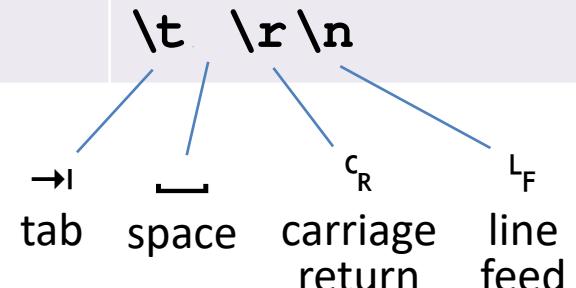
- Language: fully parenthesized numerical expressions
  - e.g.,  $((23 + 7) * 19)$
- What are the meaningful “words”?

# Example

- Language: fully parenthesized numerical expressions
  - e.g.,  $((23 + 7) * 19)$
- What are the meaningful “words”?
  - Numbers (e.g., 23, 7, 19)
  - Left (open) parantesis (i.e., ( )
  - Right (close) parenthesis (i.e. ) )
  - Operators (e.g., + , \*)



# Strings with Special Handling

Kind	Examples
Comments	<code>/* Ceci n'est pas un commentaire */</code>
Preprocessor directives	<code>#include &lt;foo.h&gt;</code>
Macros	<code>#define THE_ANSWER 42</code>
White spaces	<code>\t \r\n</code>  <p>The diagram illustrates the decomposition of white space characters. It shows four characters: a tab (\t), a space (\ ), a carriage return (\r), and a line feed (\n). Blue lines connect each character to its corresponding name below it: 'tab' connects to the tab character, 'space' connects to the space character, 'carriage return' connects to the carriage return character, and 'line feed' connects to the line feed character.</p>

\*In some languages, whitespaces do more than separate words

# Some Basic Terminology

- **Lexeme** (aka symbol) – a series of characters separated from the rest of the program according to a convention (space, semicolon, word boundary...)
- **Pattern** – a rule specifying a set of strings.  
Example: “*an identifier is a string that starts with a letter, followed by letters and digits*”
- **Token** – a pair of <pattern, attributes (of lexeme)>

“kind”  
“class”

“value”  
“representation”

# A More Realistic Example

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

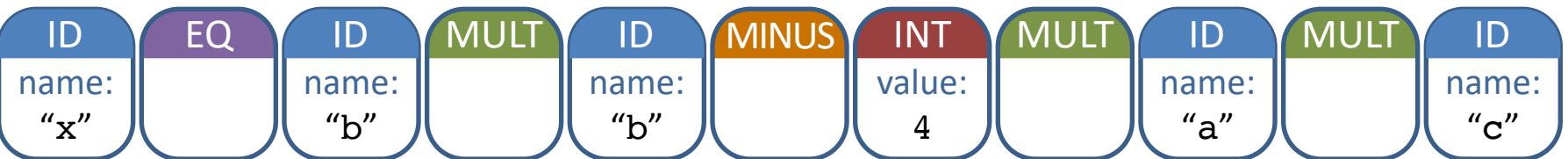
```
VOID ID("match0") LPAREN CHAR DEREF ID("s") RPAREN  
LBRACE  
IF LPAREN NOT ID("strcmp") LPAREN ID("s") COMMA  
STRING("0.0") COMMA NUM(3) RPAREN RPAREN  
RETURN REAL(0.0) SEMI  
RBRACE  
EOF
```

# Scanner\*: Characters to Patterns

txt  
 $x = b*b - 4*a*c$



Token Stream



# How can we define

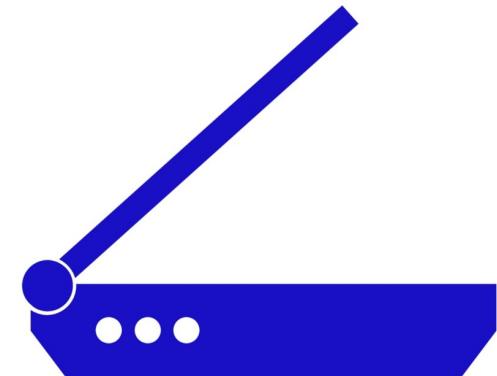
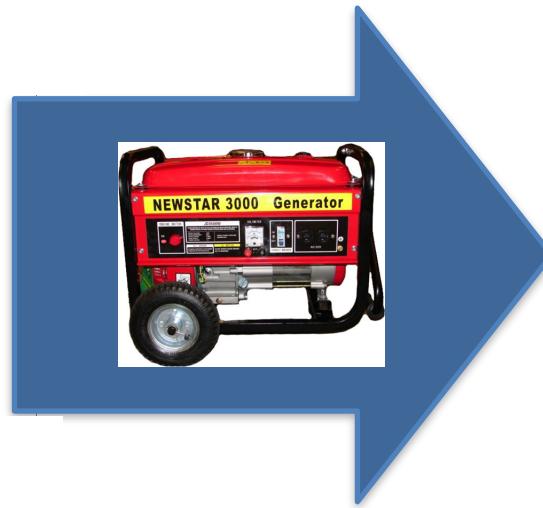
- Keywords — easy!
  - **if, then, else, for, while, ...**
- Identifiers?  
**x0, b2\_b, heIsNoOne, ...**
- Numerical Values?  
**0, 1, 2, 3.4, 5.67, ...**
- Strings?  
**"there are infinitely many", ...**
- Characterize **unbounded** sets of values using a **bounded** description?

An identifier is a sequence of letters and digits; the first character must be a letter.  
The underscore `_` counts as a letter.  
Upper- and lowercase letters are different.  
If the input stream has been parsed into tokens up to a given character, the next token is taken to include the longest string of characters that could possibly constitute a token.  
Blanks, tabs, newlines, and comments are ignored except as they serve to separate tokens.  
Some white space is required to separate otherwise adjacent identifiers, keywords, and constants.

# Regular languages & FSMs

- Formal languages
  - $\Sigma$  = finite set of letters
  - Word = sequence of letter
  - Language = set of words
- Regular languages defined equivalently by
  - Regular expressions (human-readable)
  - Finite-state machines (machine-executable)

# Goal: Regular Expressions $\Rightarrow$ Scanner



# How: Regular Expressions $\Rightarrow$ Scanner



Describe tokens as regular expressions



Create an automaton



Use automaton to scan the input

(compiler  
construction time)  
(compilation  
time)

# (Computer-Readable) Regular Expressions



# Regular Expressions over $\Sigma$

## ► Basic Patterns

x	A single letter 'x' from the alphabet $\Sigma$
.	Any character ( <i>usually excluding new-line characters</i> )
[xyz]	Any of the characters x,y,z

Highest precedence

## ► Repetition Operators

R?	Either an R or nothing (= <i>optionally</i> R)
R*	Zero or more occurrences of R
R+	One or more occurrences of R

Lowest precedence

## ► Composition Operators

$R_1R_2$	An $R_1$ followed by $R_2$
$R_1 R_2$	Either an $R_1$ or $R_2$

## ► Grouping

(R)	Same as R ( <i>used to override default precedence</i> )
-----	--

# Examples

- $(a|b)^* =$
- $(0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^+ =$
- $ab^* \mid cd? =$
- $[<>]=? \mid == \mid <> =$

# 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 in groups to denote a range
  - letter = [a-z] | [A-Z]  [a-zA-Z]
  - digit = [0-9]

# Shorthands

- Use names for expressions
  - $\text{letter} = \text{a} \mid \text{b} \mid \dots \mid \text{z} \mid \text{A} \mid \text{B} \mid \dots \mid \text{Z}$
  - $\text{letter\_} = \text{letter} \mid \underline{\phantom{x}}$
  - $\text{digit} = 0 \mid 1 \mid 2 \mid \dots \mid 9$
  - $\text{id} = \text{letter\_} (\text{letter\_} \mid \text{digit})^*$
- Use caret to denote negation
  - $\text{notabc} = [\wedge \text{abc}]$
  - $\text{notdigit} = [\wedge 0-9]$

# Escape characters

- What is the expression for one or more “+” symbols?
  - `(+)+` won't work
  - `(\+)+` will
- backslash \ before an operator turns it into a standard character
  - `\*, \?, \+, \[, ...`
- backslash \ before r / n turns it into carriage return (`\r`) / line feed (`\n`)

# More Examples

- `if` = `if`
- `then` = `then`
- `word` = [A-Za-z] [a-z]\*
- `digit` = [0-9]
- `digits` = `digit+`
- ?? = `digits ( \. digits ( e (\+|\-)?) digits )? )?`

E.g., ??

# Ambiguity

- if = if
  - id = letter\_ (letter\_ | digit)\*
- 
- letter\_ = [\_A-Za-z]  
digit = [0-9]

What's the problem?

# Ambiguity

- if = if
- id = letter\_ (letter\_ | digit)\*
  - letter\_ = [\_A-Za-z]
  - digit = [0-9]
- “**if**” matches both the pattern for reserved words and the pattern for identifiers... so what should it be?
- How about the identifier “**iffy**”?

# Ambiguity

- if = if
- id = letter\_ (letter\_ | digit)\*
  - letter\_ = [\_A-Za-z]
  - digit = [0-9]
- “**if**” matches both the pattern for reserved words and the pattern for 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  
( $\Rightarrow$  **tip: list rules for keywords before identifiers**)

# Building blocks of a lexical analyzer



# Goal: a lexical analyzer

- Given a list of token definitions, write a program such that
  - Conceptually,
    - Input: String of characters to be analyzed
    - Output: List of tokens
  - Practically, implement `get_next_token()`
- How?

# Building a Scanner – Take I

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

# There must be a better way!



# A better way

- Define tokens using regular expressions
- Convert regular expressions to a finite-state automaton
- Use finite-state automaton for detection

# Reg-exp vs. automata

- Regular expressions are declarative
  - Good for humans
  - Not “executable”
- Automata are operative
  - Define an *algorithm* for deciding whether a given word is in a regular language
  - Not a natural notation for humans

# Construction

Token Definitions



IF : if

TH : then

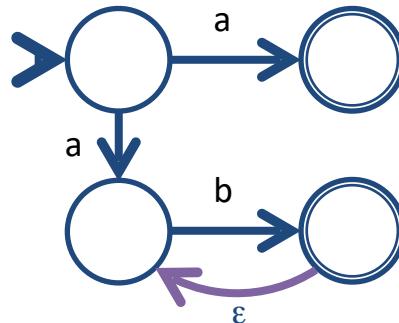
REL : [<>]=?

| == | <>

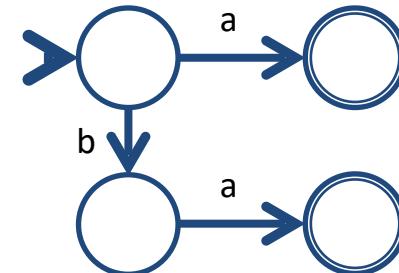
INT : [0-9]+



Non-deterministic  
“Finite Automaton”



Deterministic  
“Finite Automaton”



Token Stream

<ID,“x”> <EQ> <ID,“b”> <MULT> <ID,“b”> <MINUS>  
<INT,4> <MULT> <ID,“a”> <MULT> <ID,“c”>

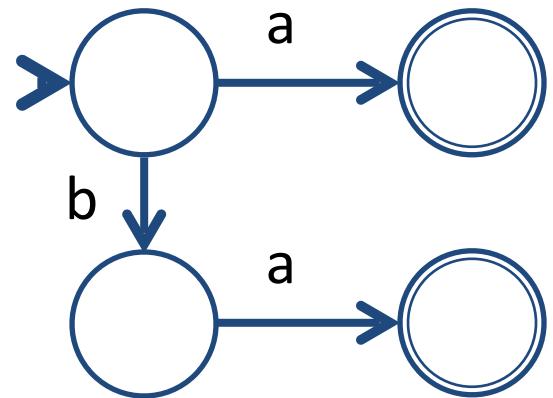
\* Actually, to handle ambiguity, a slight variant of finite-state automata is used



# Reminder: Finite-State Automata

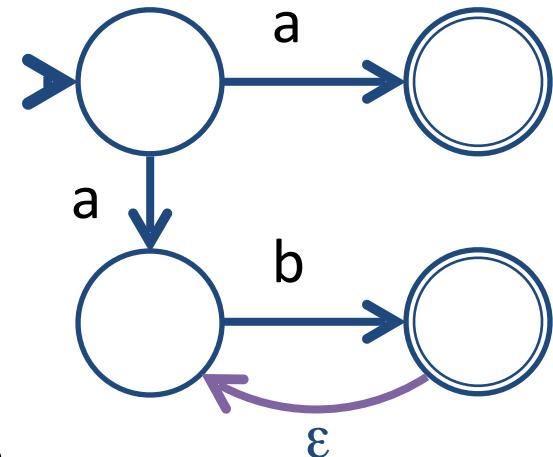
## Deterministic automaton

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

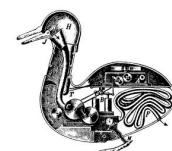


# Non-Deterministic automaton

- $M = (\Sigma, Q, \delta, q_0, F)$ 
  - $\Sigma$  - alphabet
  - $Q$  – finite set of states
  - $q_0 \in Q$  – initial state
  - $F \subseteq Q$  – final/accepting states
  - $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$  – transition function



- Allows  $\epsilon$ -transitions
- For a word  $w$ ,  $M$  can reach *a set of states* or *get stuck*
- If some state reached after reading  $w$  is final,  $M$  accepts  $w$

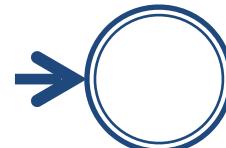


# From Regular Expressions to NFA

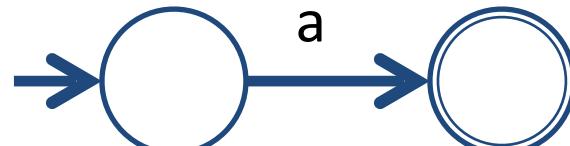


# Basic constructs

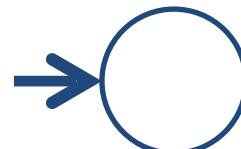
$R = \varepsilon$



$R = a$

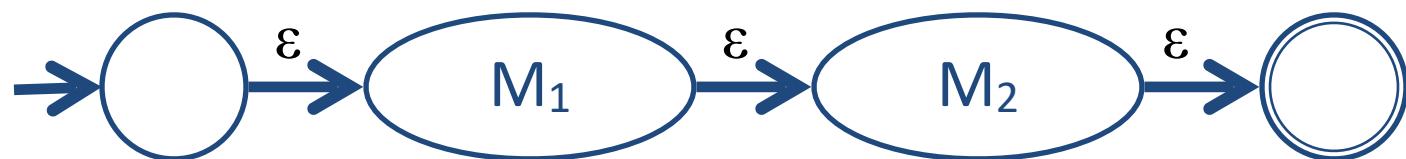


$R = \emptyset$

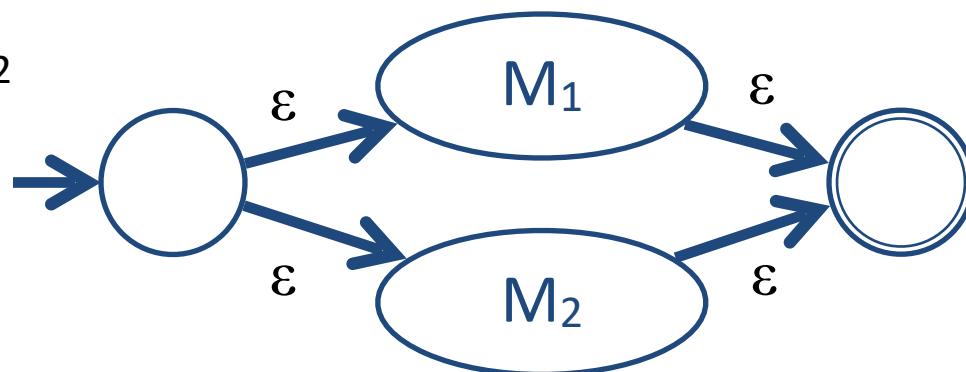


# Composition

$$R = R_1 R_2$$

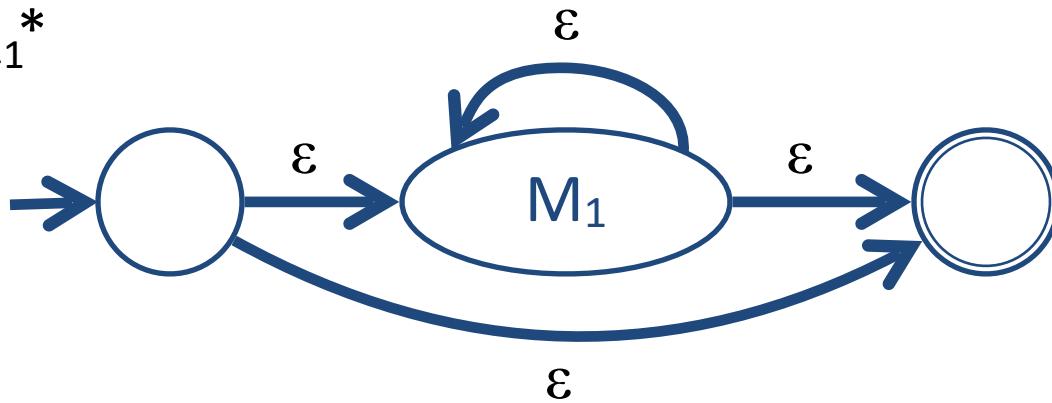


$$R = R_1 \mid R_2$$



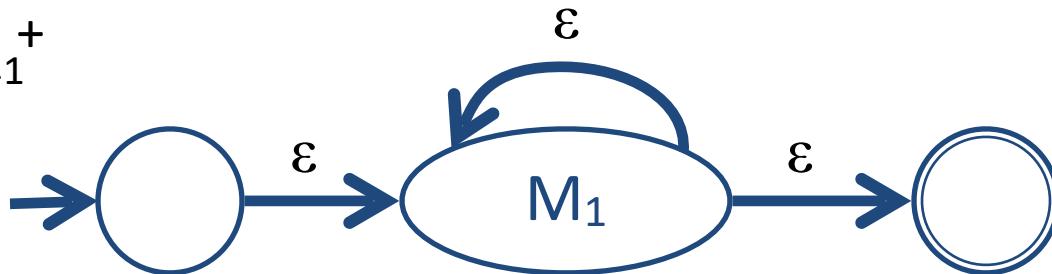
# Repetition

$$R = R_1^*$$



# Repetition

$$R = R_1^+$$

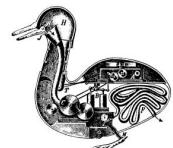


$$R_1^+ = R_1 R_1^*$$



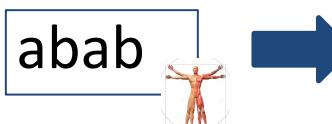
# From NFA to DFA

- Theorem: *there is an algorithm that transforms an NFA+ $\epsilon$  automaton into a DFA that accepts the same language*
- Recall “Modelim”

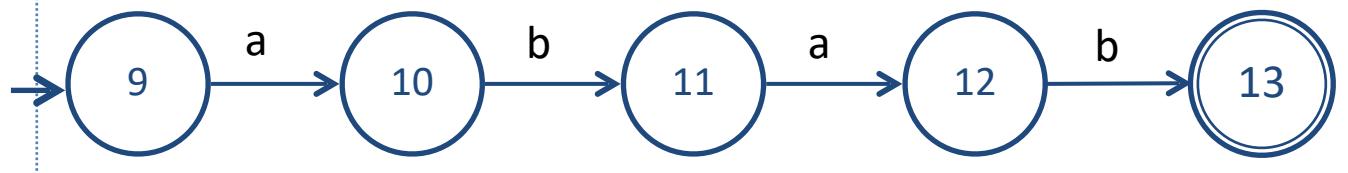
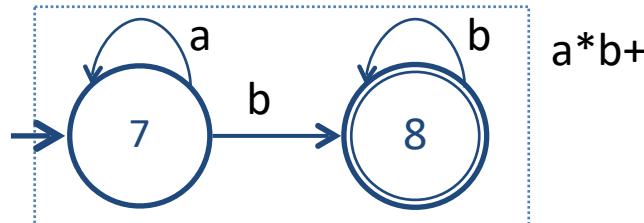
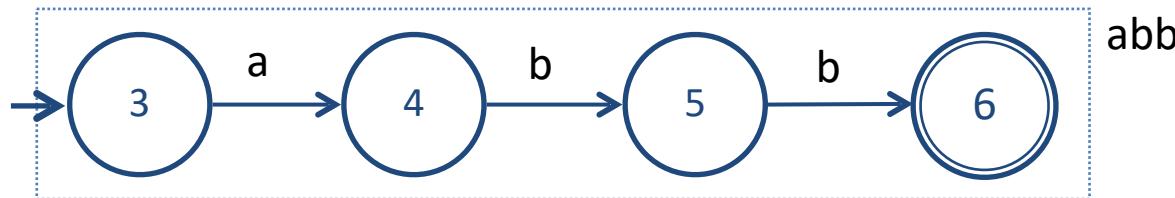
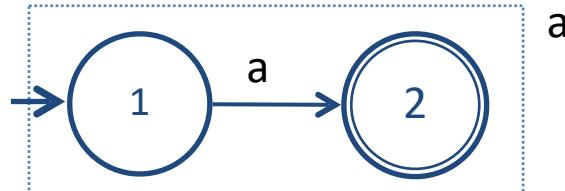


# Example: Regular Expressions to Automata

RE



DFA



# What now?

We have one DFM per pattern ( $RE_i \Rightarrow M_i$ ).

**Naïve approach:** try each automaton separately

- Recall *ambiguity resolution* ...
  - \* Longest word
  - \* Tie-breaker based on **order of rules** when words have same length

# What now?

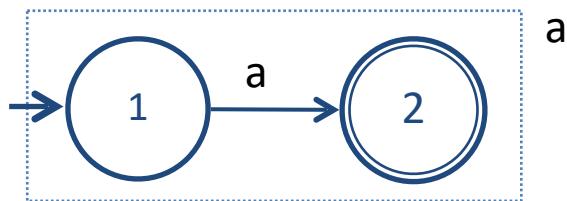
We have one DFM per pattern ( $RE_i \Rightarrow M_i$ ).

**Naïve approach:** try each automaton separately

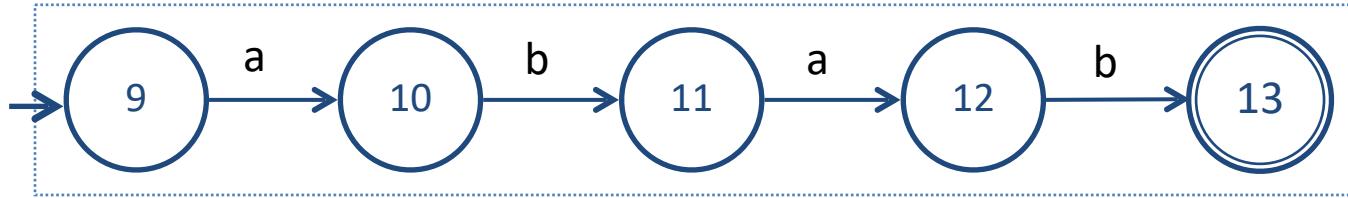
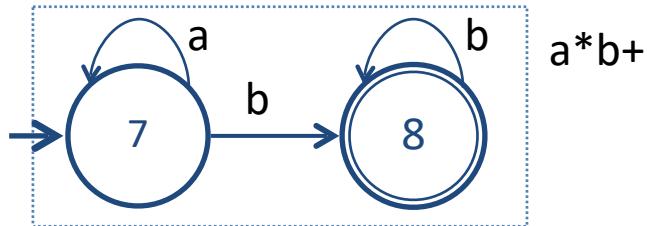
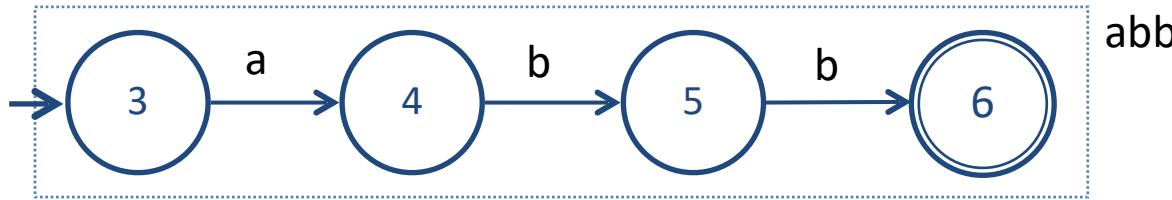
- Given a sequence of tokens as *input*\*:
    - Try  $M_1(\text{input})$
    - Try  $M_2(\text{input})$
    - ...
    - Try  $M_n(\text{input})$
- 
- Run  $M_i$  on *input* from start state until **stuck**
    - Remember **position<sub>i</sub>** after character read at the last final state and the accepted **prefix<sub>i</sub>**
  - Set  $m = \max\{position_1, \dots, position_n\}$ 
    - if  $m=0$  return ERROR
  - Set  $j = \min\{ i \mid position_i == m\}$
  - Advance input to  $position_j$
  - return  $(pattern_j, prefix_j)$

\*Requires “rewinding” input after every attempt

# Naïve Approach: Example Run



Requires “rewinding”  
after every attempt



**abaa:** Recognizes

**abba:** Recognizes

# A Better Solution

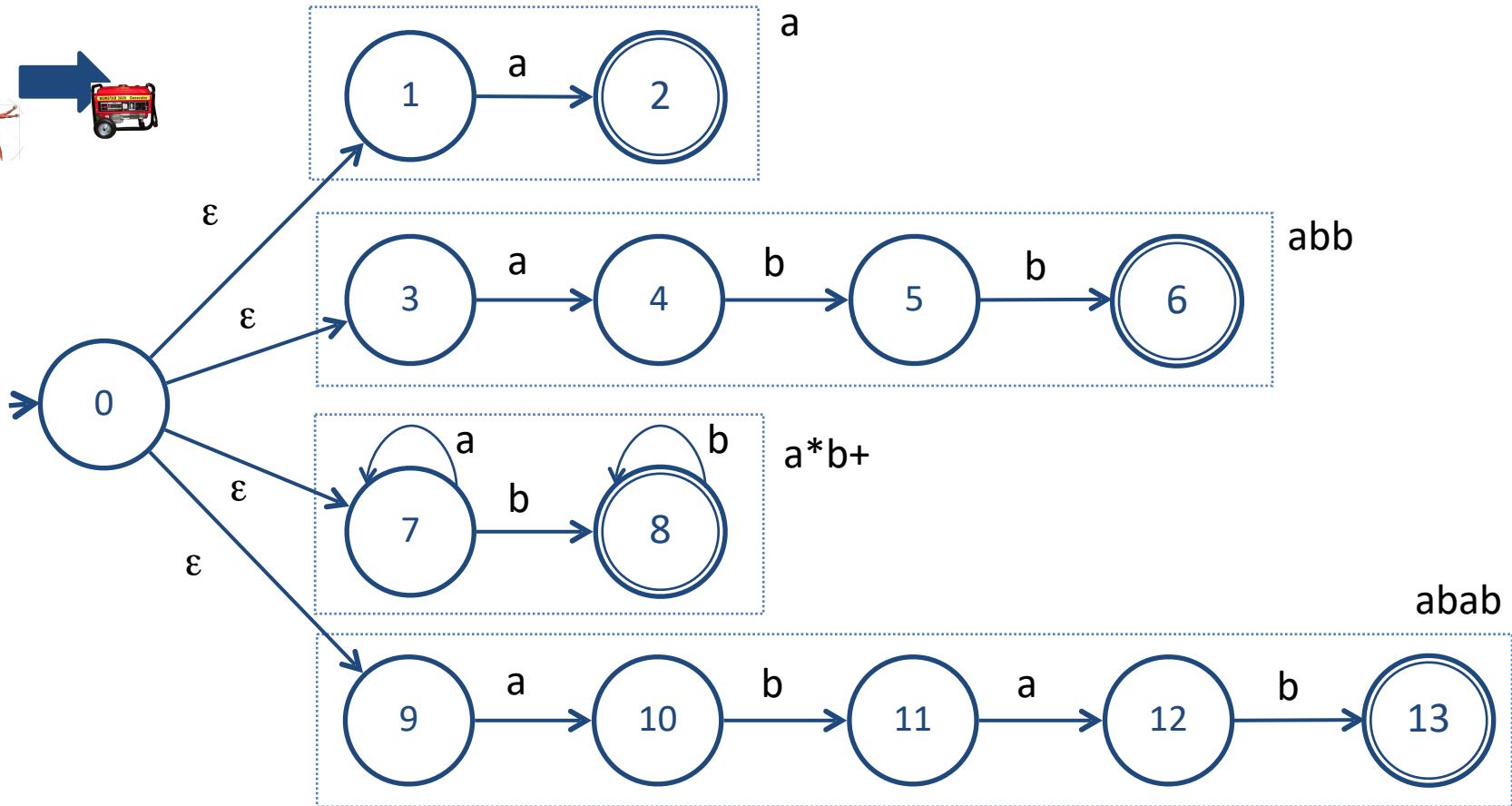
Construct a single DFM for all the patterns

- Record the origin patterns of the final states
- Run the DFM on input from start state until stuck as before
  - Again, ambiguity resolution based on length and order of rules
  - Again, requires rewinding, but only once!

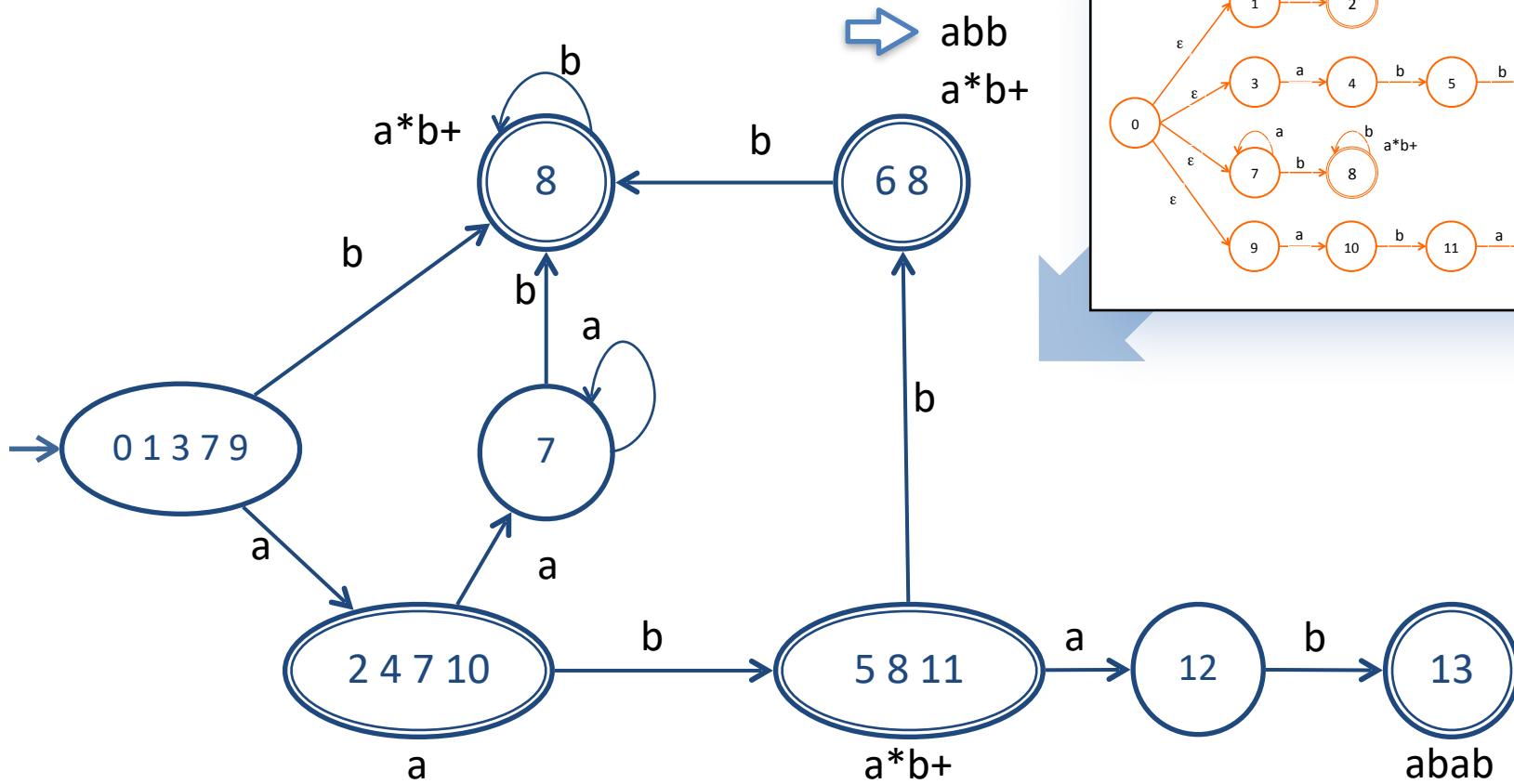
# Example: Combine Automata

combines

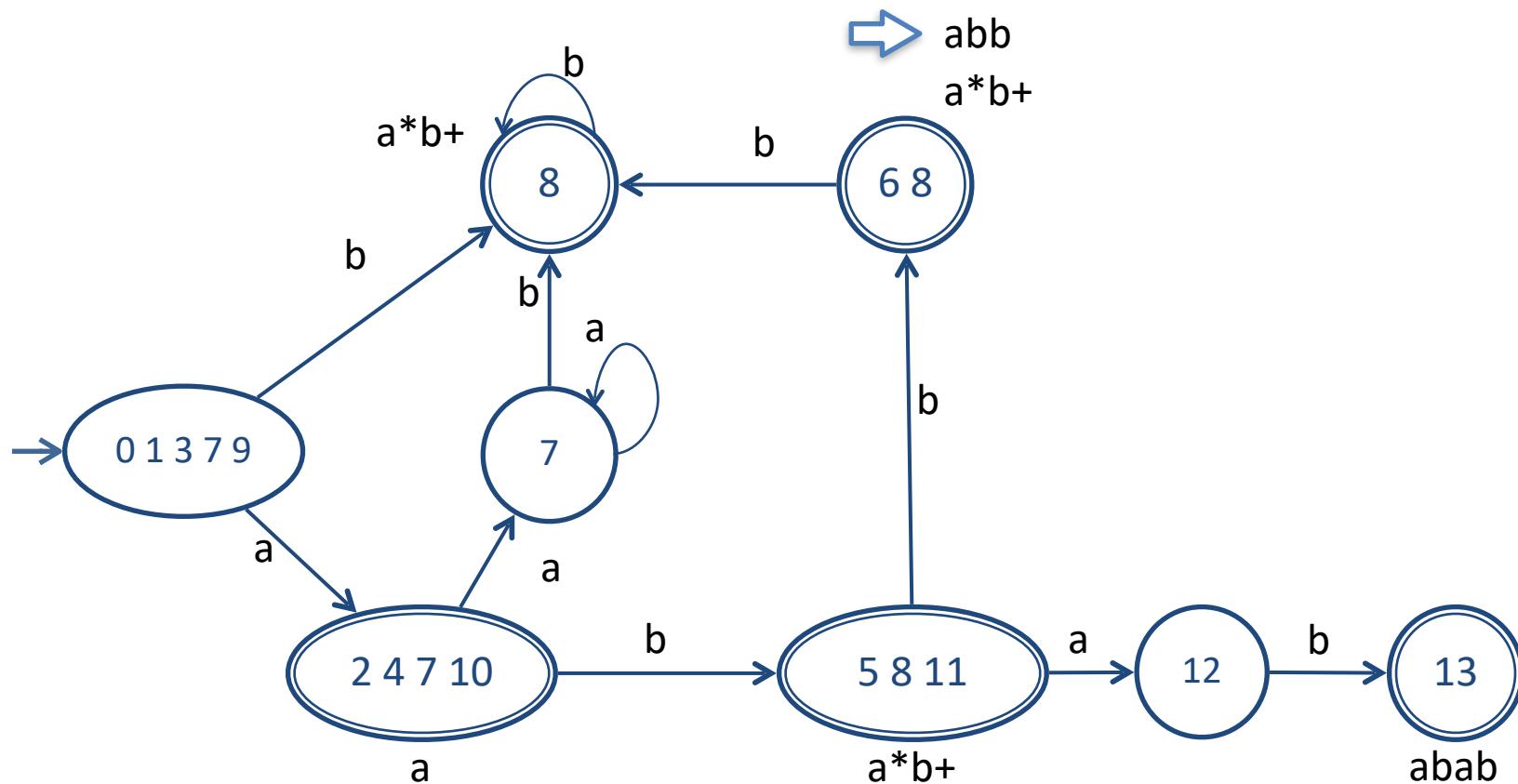
a  
abb  
 $a^*b^+$   
abab



# Example: Corresponding DFA

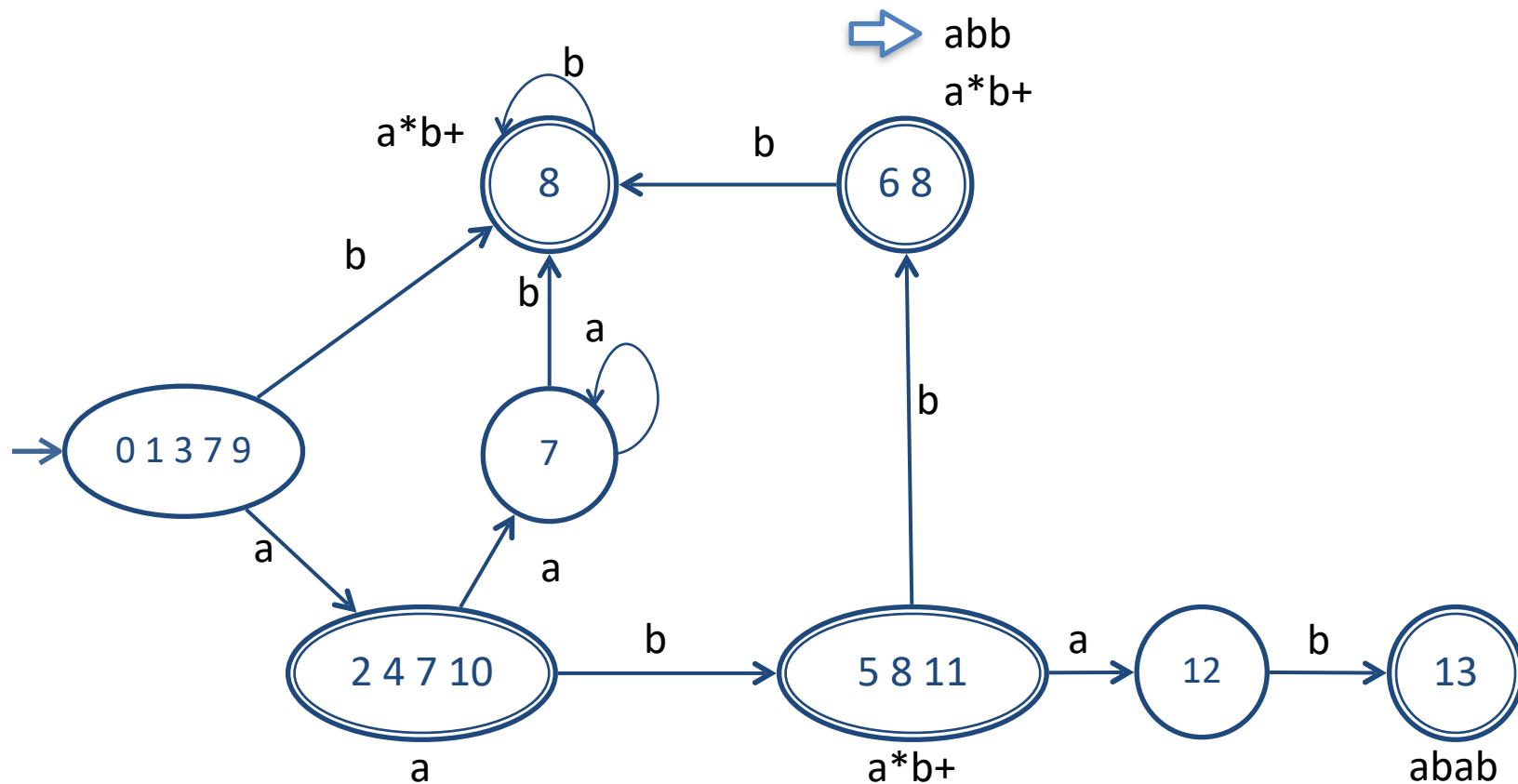


# Example Runs



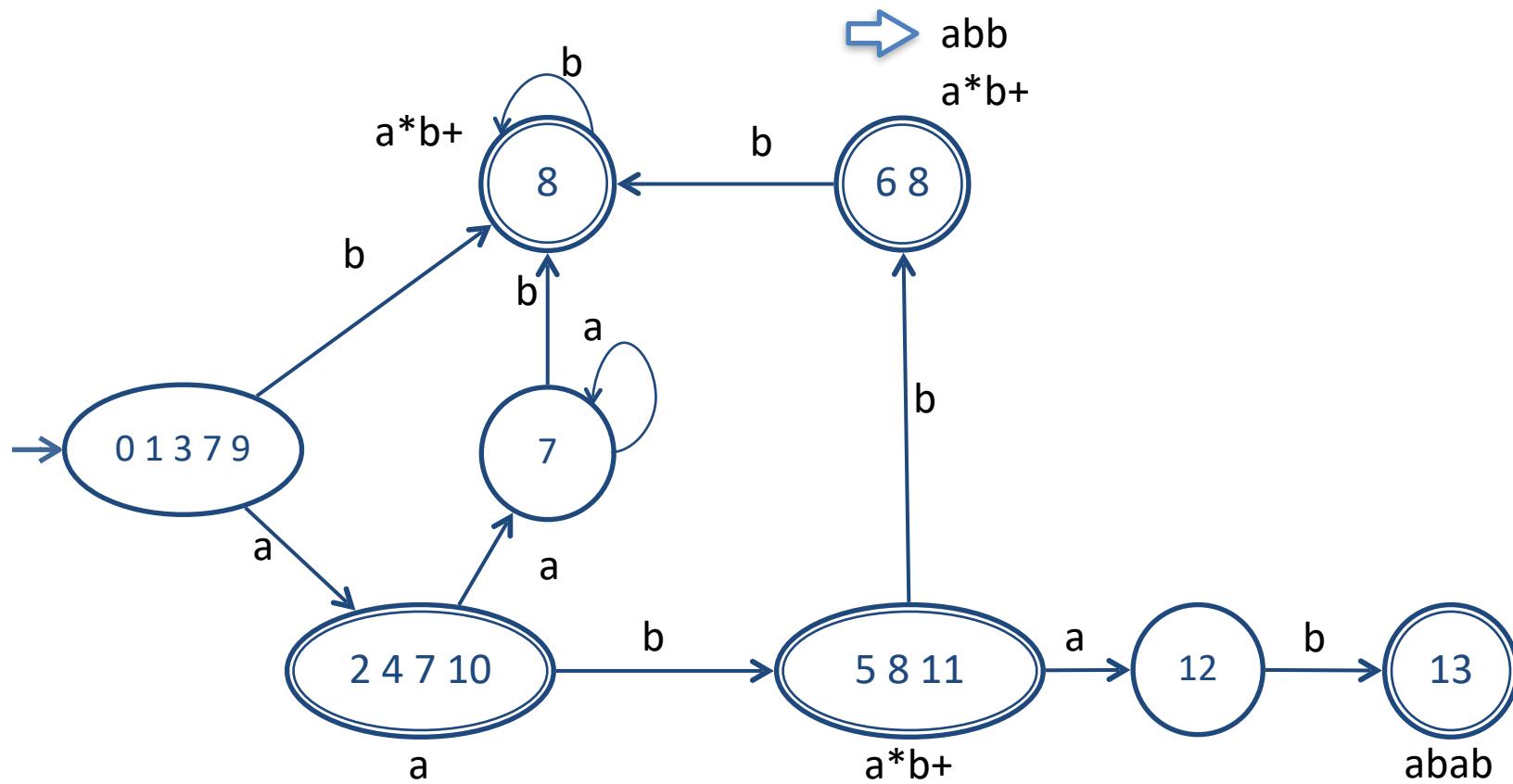
**abaa**: gets stuck after **aba** in state 12, backs up to state  $(5 \ 8 \ 11)$  pattern is  $a^*b^+$ , token is **ab**

# Example Runs



abaa →

# Example Runs



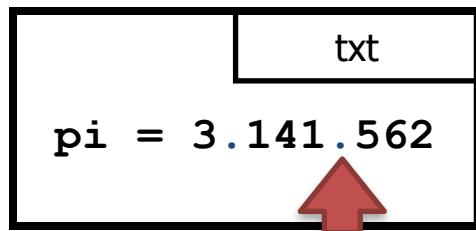
$abaa \rightarrow (a^*b+, ab) \quad (a, a) \quad (a, a)$

$abaa \rightarrow$

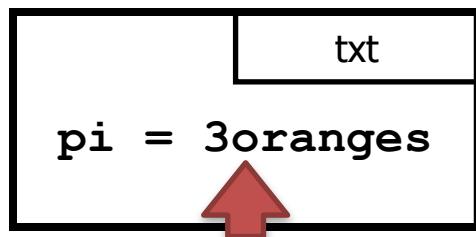
# Summary of Construction

- ✓ Describe tokens as regular expressions
  - ✓ Regular expressions “compiled” into a single DFA
  - ✓ Lexical analyzer simulates the run of an automaton with the given transition table on any input string
    - Use the following ambiguity resolution algorithm:
      - Prefer ***longest*** accepting word
      - Respect order of definitions
- (compilation time)  
(compiler construction time)

# Errors in Lexical Analysis



Illegal token



Illegal token

# Error Handling

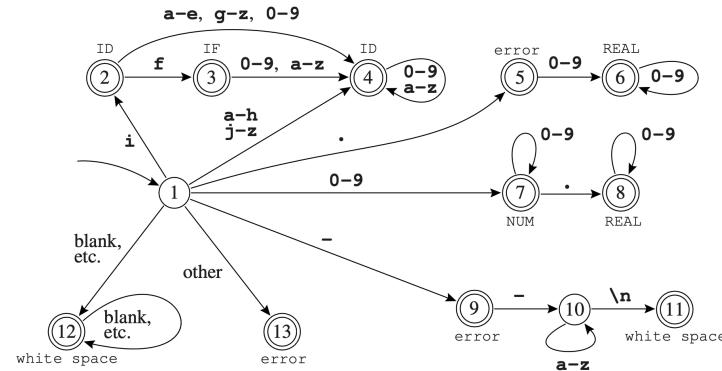
- Many errors cannot be identified at this stage
- Example: should “**fi**” be “**if**”? Or is it a routine name?
  - We will discover this later in the analysis
  - At this point, we just create an identifier token
- Sometimes the lexeme does not match any pattern
  - Easiest fix: skip characters until the beginning of a legitimate lexeme
  - Alternatives: eliminate/add/replace one letter, replace order of two adjacent letters, etc.
- Pro: allow the compilation to continue
- Con: errors that spread all over
- Record source location of token (line/column) for error reporting

**fi (a==f(x))**

# Implementing a Scanner



# Pseudocode: Encoding DFA



```

int edges[14][state] * ... , 0, 1, 2, 3, ..., -, e, f, g, h, i, j, ... */
/* state 0 */ {0, ..., 0, 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 */ , ..., 4, 4, 4, 4, ..., 0, 4, 3, 4, 4, 4, 4, ..., 0, 0},
/* state 3 */ , ..., 4, 4, 4, 4, ..., 0, 4, 4, 4, 4, 4, 4, ..., 0, 0},
/* state 4 */ , ..., 4, 4, 4, 4, ..., 0, 4, 4, 4, 4, 4, 4, ..., 0, 0},
/* state 5 */ {0, ..., 6, 6, 6, 6, ..., 0, 0, 0, 0, 0, 0, 0, ..., 0, 0},
/* state 6 */ {0, ..., 6, 6, 6, 6, ..., 0, 0, 0, 0, 0, 0, 0, ..., 0, 0},
/* state 7 */ ...
/* state ... */
/* state 13 */ {0, ..., 0, 0, 0, 0, ..., 0, 0, 0, 0, 0, 0, ..., 0, 0}
};

bool isFinal[14] = {false, false, true, ... , true, ... };

Kind pattern[14] = { Error, null, ID, IF, ... }; // Final state → Token Kind

```

# Pseudocode: Scanning routine

```
input = “ ... ”; startPos = 0;

Token get_next_token() {
    lastFinal = 0;      // Stuck state
    currentState = 1;  // Initial state
    currentPos = posAtLastFinal = startPos;

    while (currentState != 0)  {
        currentChar = input[currentPos]
        nextState = edges[currentState][currentChar];
        if (isFinal[nextState]) {
            lastFinal = nextState ;
            posAtLastFinal = currentPos;
        }
        currentState = nextState;
        currentPos++;
    }

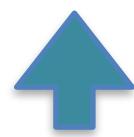
    lexeme = input[startPos .. posAtLastFinal];
    startPos = posAtLastFinal + 1;
    return (pattern[LastFinal], lexeme);
}
```

# Pseudocode: Scanning routine

```
input = " ... "; startPos = 0;

Token get_next_token() {
    lastFinal = 0;      // Stuck state
    currentState = 1;  // Initial state
    currentPos = posAtLastFinal = startPos;

    while (currentState != 0)  {
        currentChar = input[currentPos]
        nextState = edges[currentState][currentChar];
        if (isFinal[nextState]) {
            lastFinal = nextState ;
            posAtLastFinal = currentPos;
        }
        currentState = nextState;
        currentPos++;
    }
    lexeme = input[startPos .. posAtLastFinal];
    startPos = posAtLastFinal + 1;
    return (pattern[LastFinal], attr[lastFinal](lexeme), initPos);
}
```

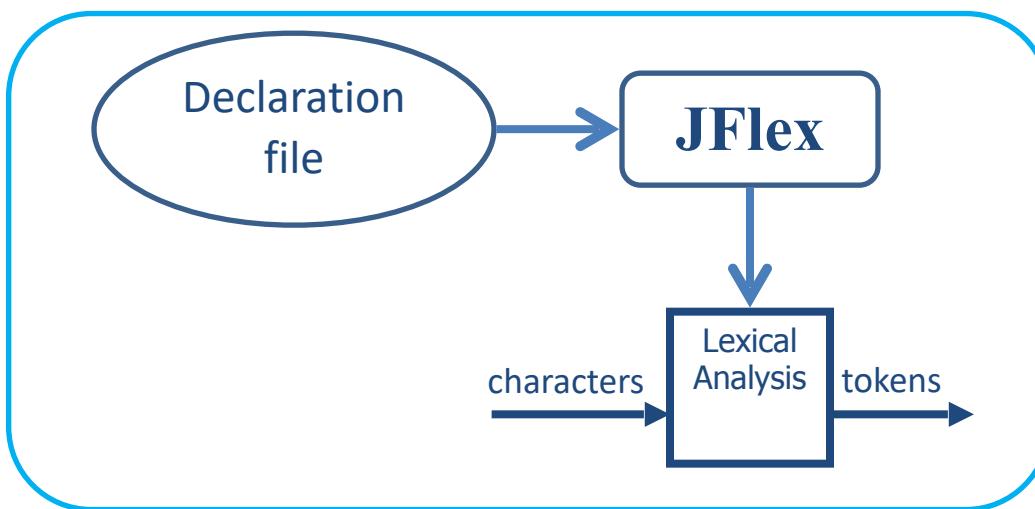


NUM: "18" → 0x12  
ID: "a18" → "a18"



# Good News: Compiler Compiler

- Construction is done **automatically** by common tools
- **JFLex/Flex/Lex** are your friends
  - Automatically generate a lexical analyzer from a declaration file
- Advantages: short declaration file, easily checked, easily modified and maintained



Intuitively:

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

# JFlex Spec File

**Java code:** Copied verbatim to generated Java file  
(before the lexer class)

%%

%{

**Java code:** Copied verbatim into the lexer class

%}

%line, %column,  
%jcup

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

**JFlex directives:** options, shorthand macros, state names

%%

“if”  
{LETTER}({LETTER}|{DIGIT})\*

{return new Symbol(Symbol.IF); }  
{return new Symbol(Symbol.ID, yytext()); }

**Lexical analysis rules:** Optional state, pattern { action }

# Creating a ~~Scanner~~ using JFlex

```
import java_cup.runtime.*;  
  
%%  
  
%cup  
%{ // Copied into the generated Lexer class source code  
    private int wsCounter = 0;  
    public int getWSCounter() { return wsCounter; }  
}  
  
LineTerminator  = \r | \n | \r\n  
WhiteSpace      = {LineTerminator} | [\s\t]  
%%  
  
{WhiteSpace} { wsCounter++; }  
.           { }
```



# 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

*rules ...*

.

{ `error("Illegal Token");` }

*to include any  
character + newlines:*

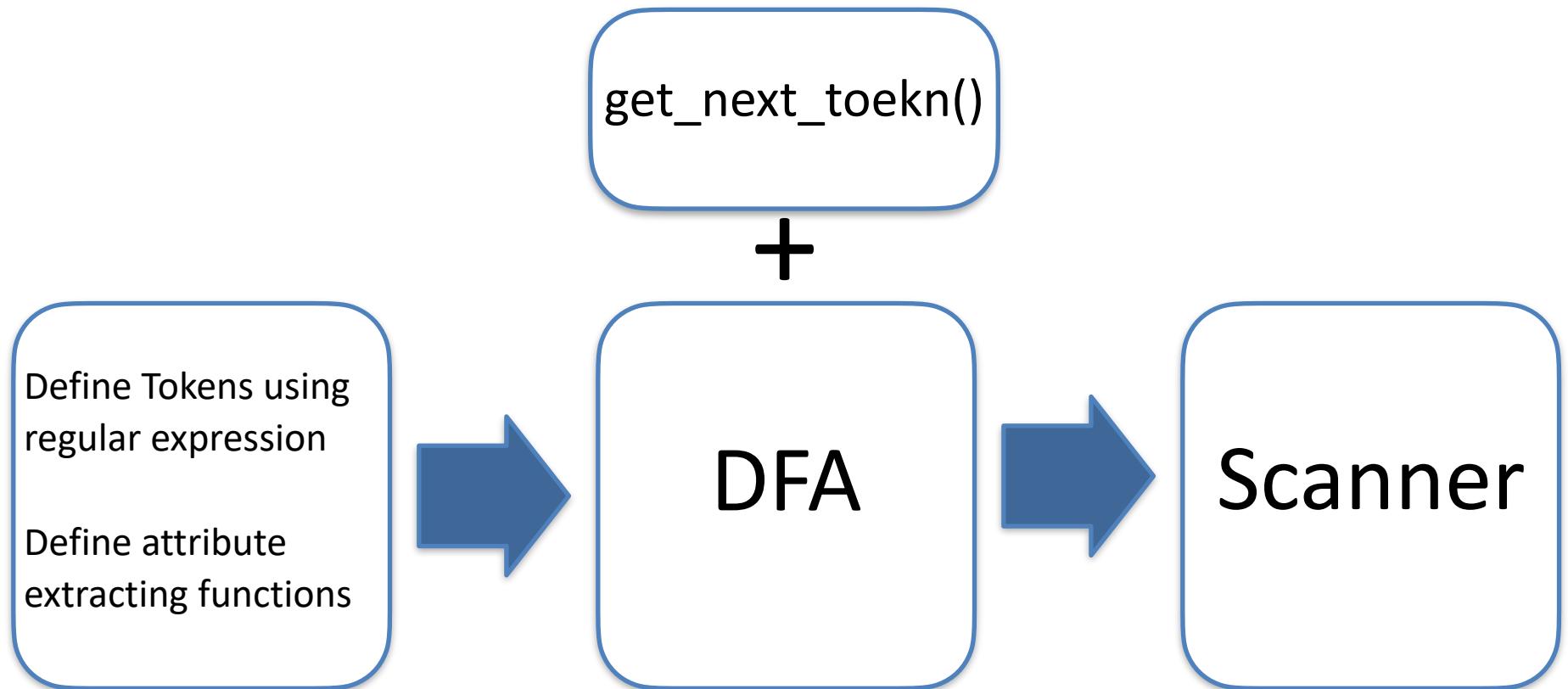
.|[\r\n]

*or*

[^]

conclusion

# Constructing A Scanner



# From Regular Expressions to DFA

(compiler construction time)

- 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,  
determinize
  - Associate every accepting state with a single pattern,  
respecting the order of definitions



# Using DFM to Recognize Tokens

(compilation time)

- Step 1: Read input (program) character by character while simulating the run of the DFM until getting stuck
  - Record last accepting state + corresponding prefix
- Step 2: Rewind input to the longest accepted prefix
- Step 3: Return the token associated with the last accepting state



# DFA Construction

- Turning an NFA into a DFA is expensive, but
  - Exponential in the worst-case
  - In practice, it works fine
- The construction is done **once per language**
  - At Compiler construction time
  - Not at compilation time
- Minimize the automaton
  - Separate accepting states by accepted token kinds

# Missing

- Efficient Scanners
  - Transition table compression
  - Input buffering
  - Using switch and gotos instead of tables
- Symbol Tables
- Nested Comments
- Handling Macros
- Creating a lexical analysis by hand

# Summary

- **Lexical analyzer**
  - ✓ Turns character stream into token stream
  - ✓ Tokens defined using regular expressions
  - ✓ Regex → NFA → DFA construction for identifying tokens
  - ✓ Automated constructions of lexical analyzer using **JFlex**

# Exercise (From Exam)

## שאלה 1 (20 נקודות)

נתונה רשימה מאורעות. לכל אחד מהמאורעות הסבירו בקצרה (3-1 שורות) האם הוא מתרחש במהלך בניית הקומpileר / הקומפילציה / loading / linking / assembly / ריצת התוכנית. אם בחרתם בזמן קומPILEציה, הסבירו מהו השלב המסויים בקומPILEר בו מתרחש המאורע ומהם מבני הנתונים הרלוונטיים. אם מאורע יכול להתרחש במספר שלבים, ציינו את כולם.

Possible answers (at this point):

- Compiler construction time
- Lexical Analysis
- Later stages

# Exercise (From Exam)

## שאלה 1 (20 נקודות)

נתונה רשימה מאורעות. לכל אחד מהאירועים הסבירו בקצרה (3-1 שורות) האם הוא מתרחש במהלך בניית הקומpileר / הקומPILEציה / loading / linking / assembly / ריצת התוכנית. אם בחרתם בזמן קומPILEציה, הסבירו מהו השלב המסויים בקומPILEר בו מתרחש המאורע ומהם מבני הנתונים הרלוונטיים. אם מאורע יכול להתרחש במספר שלבים, ציינו את כולם.

- א. (4 נקודות) נמצא כי נקראת מトודה המוגדרת במחלקה D.
- ב. (4 נקודות) נמצא כי אין ריצה בה המשתנה x מכיל ערכים אי זוגיים.
- ג. (4 נקודות) נמצא כי אין אף ריצה בה יש לכתוב את הערך של המשתנה הЛОКАלי x במחסנית.
- ד. (4 נקודות) הוחלט כי הארגומנט הראשון בפעולות כפל יהיה ברגיסטר %eax.
- ה. (4 נקודות) נמצא כי יש identifier בשם myid בתוכנית.



# Exercise (From Exam)

## שאלה 1 (20 נקודות)

נתונה רשימה מאורעות. לכל אחד מהמאורעות הסבירו בקצרה (1-3 שורות) האם הוא מתרחש במהלך בנית הקומpileר / הקומPILEציה / loading / linking / assembly / ריצת התוכנית. אם בחרתם בזמן קומPILEציה, הסבירו מהו השלב המשמעותי בקומPILEר בו מתרחש המאורע ומהם מבני הנתונים הרלוונטיים. אם מאורע יכול להתתרחש במספר שלבים, ציינו את כולם.

- א. (4 נקודות) נקבע כי השדה `f` שניגשים אליו במתודה של מחלקה `D` נורש מה `superclass` של `D`.
- ב. (4 נקודות) נמצא כי אין ריצה של התוכנית בה אין גישה מחוץ לגבולות מערך.
- ג. (4 נקודות) הוחלט כי יש לקבל את רצף התווים `else` כ `token` מסווג `ELSE` ולא `.identifier` ←
- ד. (4 נקודות) הוחלט לקודד גישה לזכרון מסוימת באמצעות הפקודה `movl 4(%ebx,%ecx,8), %eax`
- ה. (4 נקודות) נקבעה ה `loading address` של `executable`

# Exercise (From Exam)

## שאלה 1 (20 נקודות)

נתונה רשימה מאורעות. לכל אחד מהמאורעות הסבירו בקצרא (1-3 שורות) האם הוא מתרחש במהלך בניית הקומפיאר / הקומפיאלציה / ריצת התוכנית. אם בחרתם בזמן בזמן קומפיאלציה, הסבירו מהו השלב המסוים בקומפיאר בו מתרחש המאורע ומהם מבני הנתונים הרלבנטיים. אם מאורע יכול להתרחש במספר שלבים, ציינו את כולם.

א. (4 נקודות) ערך ה register callee-save משוחזר מהמחסנית.

ב. (4 נקודות) יש בדקוק רקורסיה שמאלית.

ג. (4 נקודות) אין רегистר פנוי עבור המשתנה הזמן  $t170$ .

ד. (4 נקודות) השם  $$1r$  אינו שם חוקי.

ה. (4 נקודות) יש פונקציה בשם  $r1$ .



# Exercise (From Exam)

נתונה רשימה מאורעות. לכל אחד מהאירועים הסבירו בקצרה (3-1 שורות)

1. אם הוא מתרחש במהלך בניית הקומפיילר/הקומPILEZA / Linking / Loading / ריצת התוכנית.
  2. אם בחרת בזמן קומPILEZA, הסבר מהו שלב המסורים בקומפיילר בו מתרחש המאורע ומהם מבני הנ吐ונים הרלבנטיים.
  3. אם מאורע יכול להתרחש במספר שלבים, ציינו את כולם.
- א. (4 נקודות) נקבע כי הפרמטרים יוכנסו למחסנית בסדר הפוך מסדרם בראשימת הארגומנטים..
- ב. (4 נקודות) נמצאה כי יש בתוכנית מחרוזת קבועה. 
- ג. (4 נקודות) נתגלה כי הקוד של הפרוצדורה `printf` מהספריה הסטנדרטית אינו נמצא.
- ד. (4 נקודות) נתגלה כי **בכל** פעם שמתבצע הפקודה  $y + 42 = x$  ערכו של  $y$  הינו 0.
- ה. (4 נקודות) נמצא כי **שמות כל המשתנים** בתוכנית מתחלים באות  $x$ . 

identifiers

הוחלט

# Exercise (From Exam)

## שאלה 2 (20 נקודות)

נתונה רשימה מאורעות. לכל אחד מהמאורעות הסבירו בקצרה (3-1 שורות)

1. אם הוא מתרחש בזמן בניית הקומפִילר, בזמן קומפִילציה או בזמן ריצה
2. אם בחרת בזמן קומפִילציה, הסבר מהו שלב המסורם בקומפִילר בו מתרחש המאורע ומהם מבני הנתונים הרלוונטיים.
3. אם בחרת בזמן ריצה, הסבר האם מתרחש לפני ביצוע התוכנית, בזמן הביצוע, או לאחריו.
  - א. (4 נקודות) התגלה כי המתכנת השתמש במילת מפתח (keyword) כשם משתנה. 
  - ב. (4 נקודות) התגלתה בקוד פעלות חלוקה באפס.
  - ג. (4 נקודות) התגלתה מחלוקת שיש בה שדה מטיפוס לא מוגדר.
  - ד. (4 נקודות) התגלה כי משתנה מסוים ערכו תמיד חיובי.
  - ה. (4 נקודות) התגלה כי הדקדוק של שפת התכנות הינו רב משמעות.

# Exercise (From Exam)

## שאלה 1 (20 נקודות)

נתונה רשימה מאורעות. לכל אחד מהאירועים הסבירו בקצרה (3-1 שורות) loading/linking/קומפילציה/loading/linking/kompilezja/ריצת התוכנית.

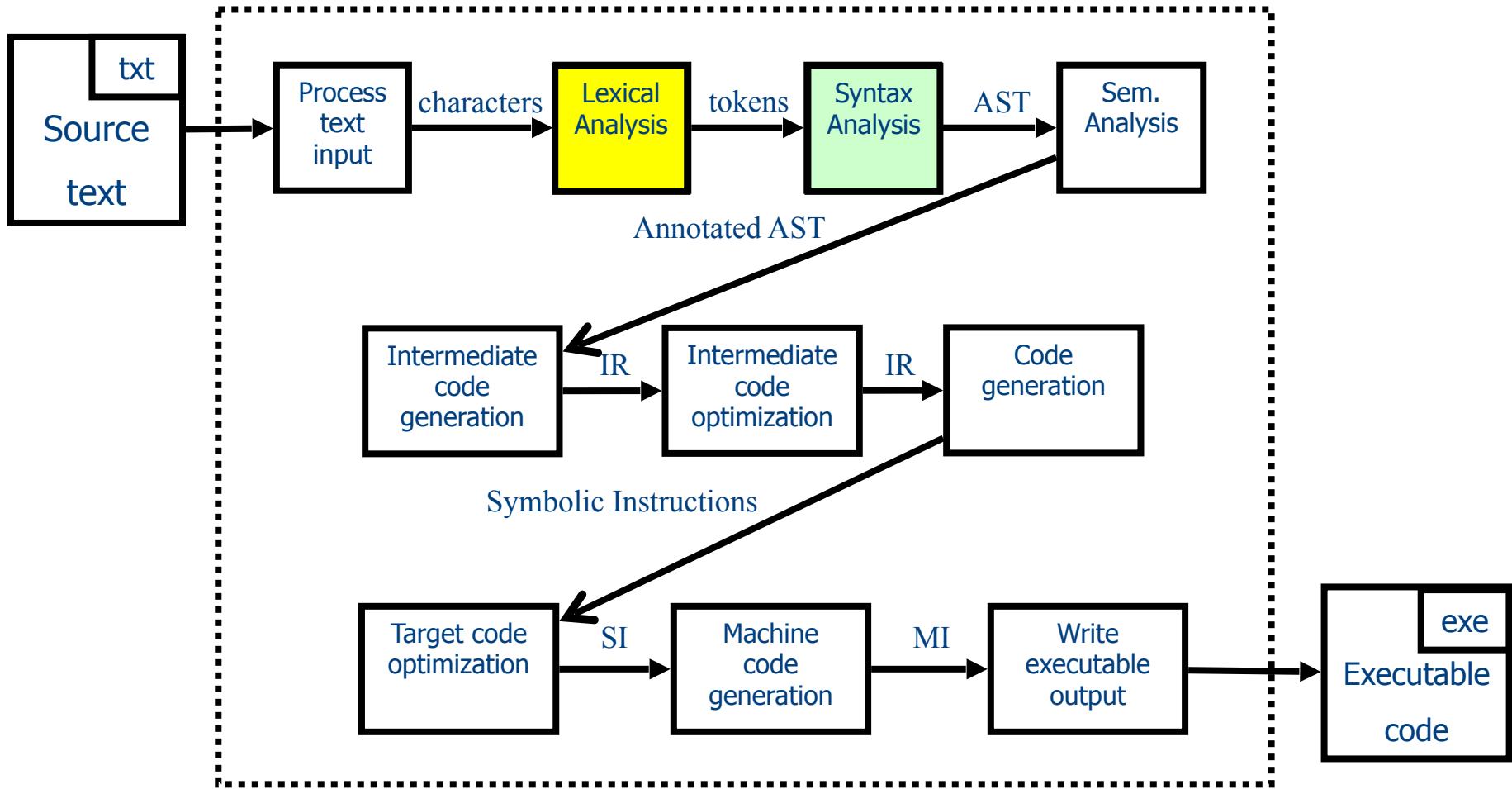
1. אם הוא מתרחש במהלך בנית הקומpileר/הקומPILEזיה/loading/linking/ריצת התוכנית.
2. אם בחרת בזמן קומPILEזיה, הסביר מהו שלב המסויים בקומPILEר בו מתרחש המאורע ומהם מבני הנתונים הרלבנטיים.

- א. (4 נקודות) התגלה כי ישנה שגיאה בקוד הבא: X+++++Y  
- הניתן כי אופרטורים ++ + מותרים, וכי X ו Y הם מסוג int. (אפשר להפעיל ++ רק ישירות על משתנים)
- ב. (4 נקודות) התגלה כי אין די רегистרים במכונה עבור קוד הבינאים הנוכחי.
- ג. (4 נקודות) התגלה כי המתודה הנקראת לא מוגדרת.
- ד. (4 נקודות) התגלה כי אין ספירה סטנדרטיבית נדרשת.
- ה. (4 נקודות) התגלה כי תמיד המתודה foo הנקראת שייכת למחלקה האב.

```
If (...) // ... = random choice
    x = new MyClass()
else
    x = new MySuperClass()
x.foo()
```

מה יקרה אם הקוד הוא Y++ + ++X? (רווחים מפרידים בין tokens)

# The Real Anatomy of a Compiler



# The end