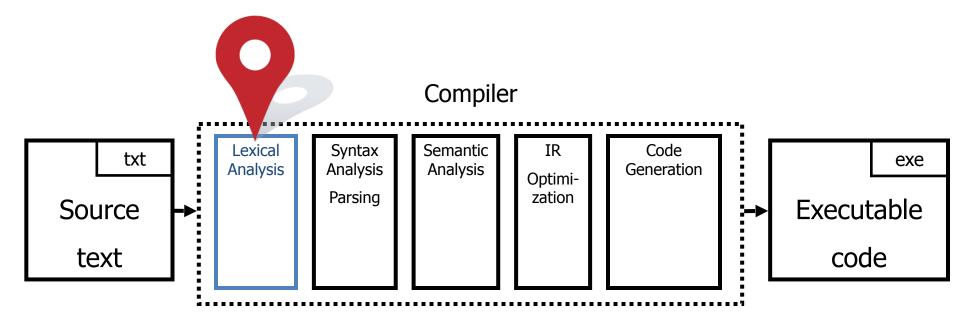
## THEORY OF GOMPILATION

#### LECTURE 01

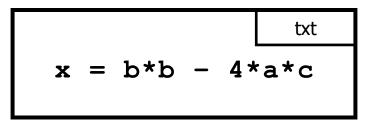


## You are here



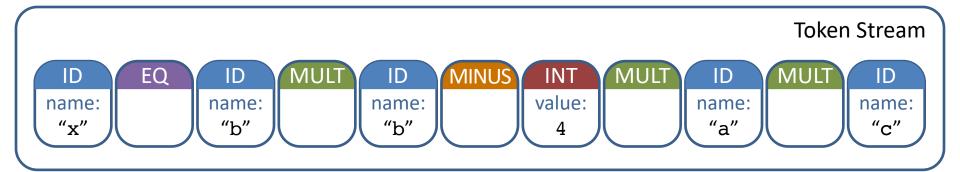
# From Characters to Tokens











# 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
- A token has:
  - a kind

lexical attributes

# **Example Tokens**

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



# Strings with Special Handling

Kind	Examples
Comments	<pre>/* Ceci n'est pas un commentaire */</pre>
Preprocessor directives	<pre>#include <foo.h></foo.h></pre>
Macros	#define THE_ANSWER 42
White spaces	\t \r\n
<b>→</b> I	C <sub>R</sub> L <sub>F</sub>
tab	space carriage line return feed



<sup>\*</sup> Note. these are *not* tokens, the are recognized, handled, and then discarded.

# 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, ...)

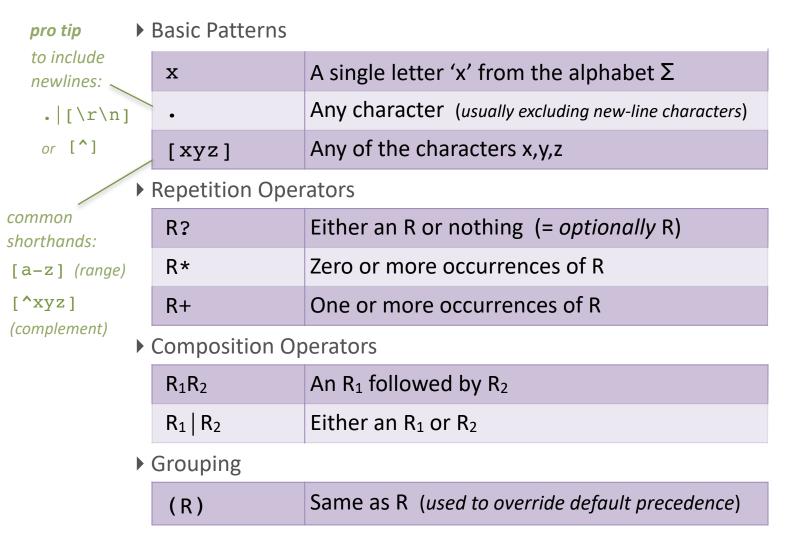
- Pattern a rule specifying a set of strings.
   Example: "a string that starts with a letter, followed by letters and digits"
- Token a pair of (kind, attributes)

#### How can we define tokens?

- Keywords easy!
  - if, then, else, for, while, ...
- Identifiers? x0, b2b, heIsNoOne, ...
- Numerical Values? 0, 1, 2, 3.4, 5.67, ...
- Strings? "there are infinitely many", ...
- Characterize unbounded sets of values using a bounded description?

# Regular Expressions

 $\sum$  = alphabet (set of characters)



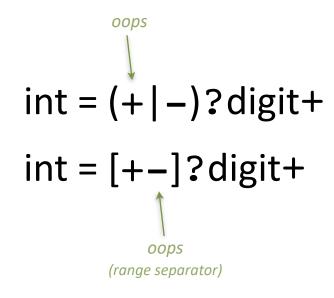


# More Examples

- if = if
- then = then
- word = [A-Za-z][a-z]\*

- digit = [0-9]
- digits = digit+

useful shorthands: once we have defined "digit", we can use it as a macro



# Nit: Escape characters

 What is the expression for one or more "+" symbols?

 backslash \ before any operator turns it into a standard character

# More Examples

- if = if
- then = then
- word = [A-Za-z][a-z]\*

- digit = [0-9]
- digits = digit+

• | ?? | = digits (\. digits)? (e (\+|-)? digits)?

# **Ambiguity**

```
• if = if
• id = letter_ (letter_ | digit)* | 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
     (⇒ tip: list rules for keywords before identifiers)

## Time for some action



https://compi.now.sh

# And now, let the real fun begin!

Generate a lexical analyzer automatically from token definitions

- Main idea
  - Use finite-state automata to match regular expressions

### Construction

x = b\*b - 4\*a\*c

#### **Token Definitions**

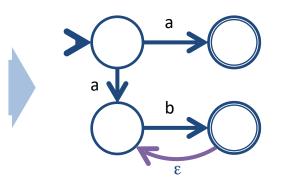
IF: if

TH: then

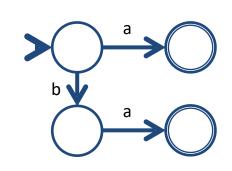
REL: [<>]=?

INT: [0-9]+

Non-deterministic Finite Automaton



Deterministic
Finite Automaton



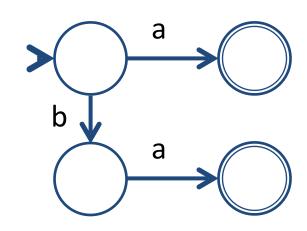
**Token Stream** 

 $\langle ID,"x"\rangle \langle EQ\rangle \langle ID,"b"\rangle \langle MULT\rangle \langle ID,"b"\rangle \langle MINUS\rangle \langle INT,4\rangle \langle MULT\rangle \langle ID,"a"\rangle \langle MULT\rangle \langle ID,"c"\rangle$ 

### Reminder: Finite-State Automata

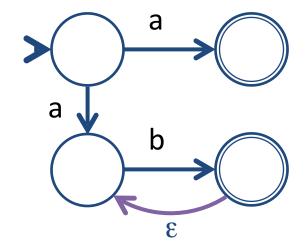
#### **Deterministic automaton**

- $M = (\Sigma, Q, \delta, q_0, F)$ 
  - $\Sigma$  alphabet
  - Q finite set of state
  - $ightharpoonup q_0 \in Q initial state$
  - ightharpoonup F  $\subseteq$  Q final states
  - $^{\bullet}$  δ : Q × Σ → Q − transition function



#### Non-Deterministic automaton

- $M = (\Sigma, Q, \delta, q_0, F)$ 
  - $\Sigma$  alphabet
  - Q finite set of state
  - $ightharpoonup q_0 \in Q$  initial state
  - ▶  $F \subseteq Q$  final states



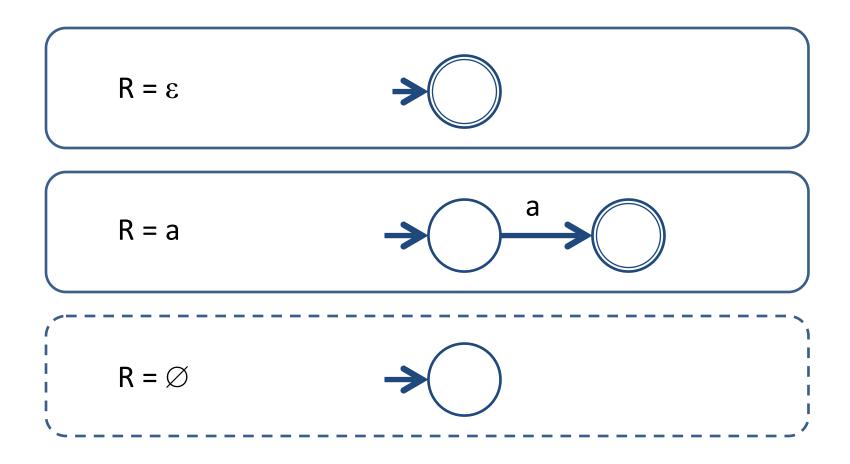
▶  $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$  – transition function

- Allows ε-transitions
- For a word w, M can reach a set of states or get stuck. If some state reached is final, M accepts w.

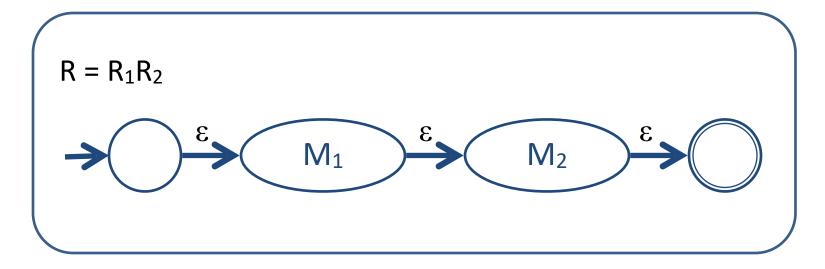
# From Regular Expressions to NFA

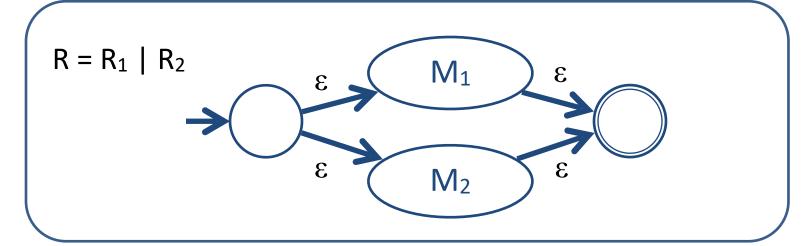
- Step 1: assign expression names and obtain pure regular expressions R<sub>1</sub>...R<sub>m</sub>
- Step 2: construct an NFA  $M_i$  for each regular expression  $R_i$
- Step 3: combine all M<sub>i</sub> into a single NFA
  - Ambiguity resolution: prefer longest accepting word

## Basic constructs

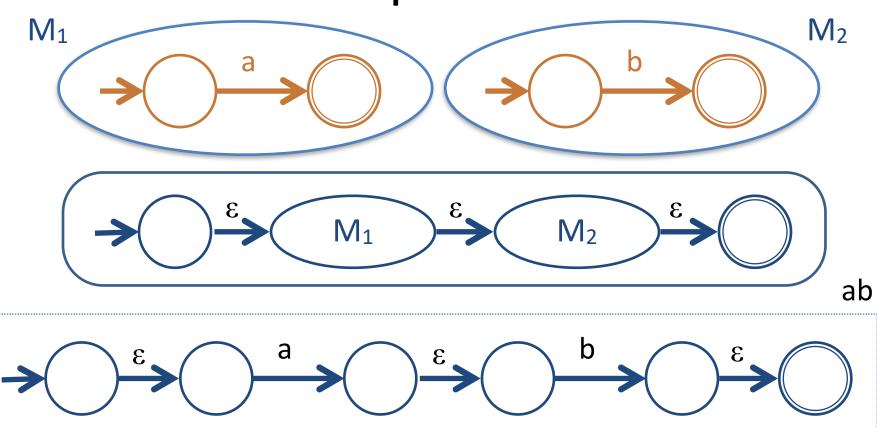


# Composition



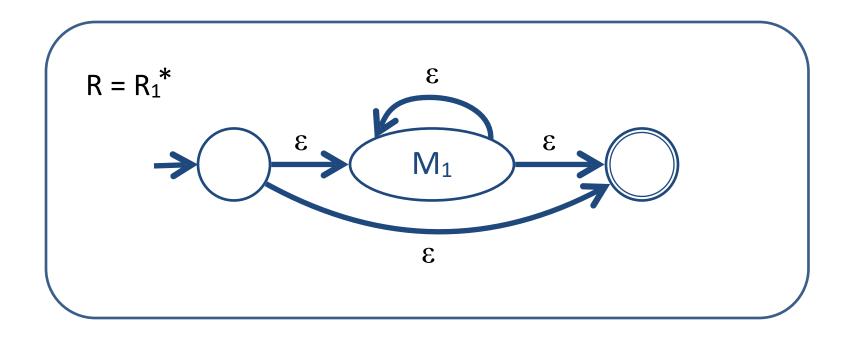


# Composition



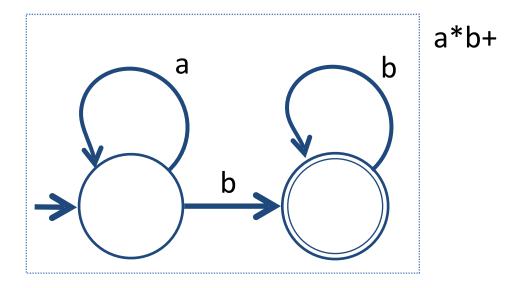
minimize:

# Repetition



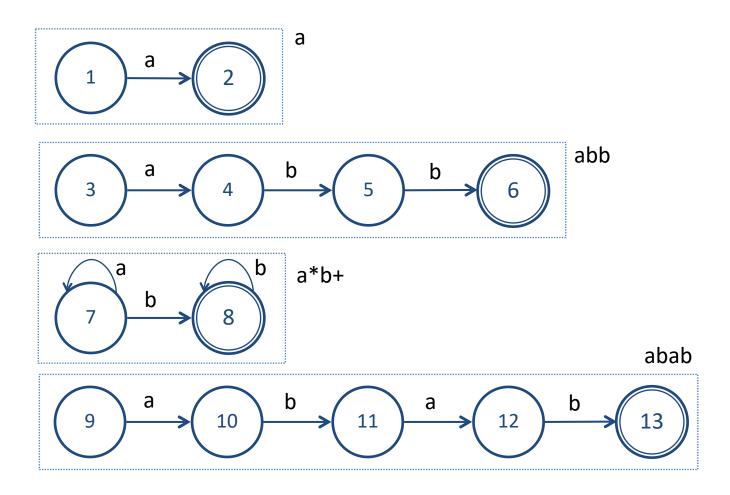
$$R_1 + = R_1 R_1^*$$

# Repetition



## One Automaton per Pattern

a abb a\*b+ abab

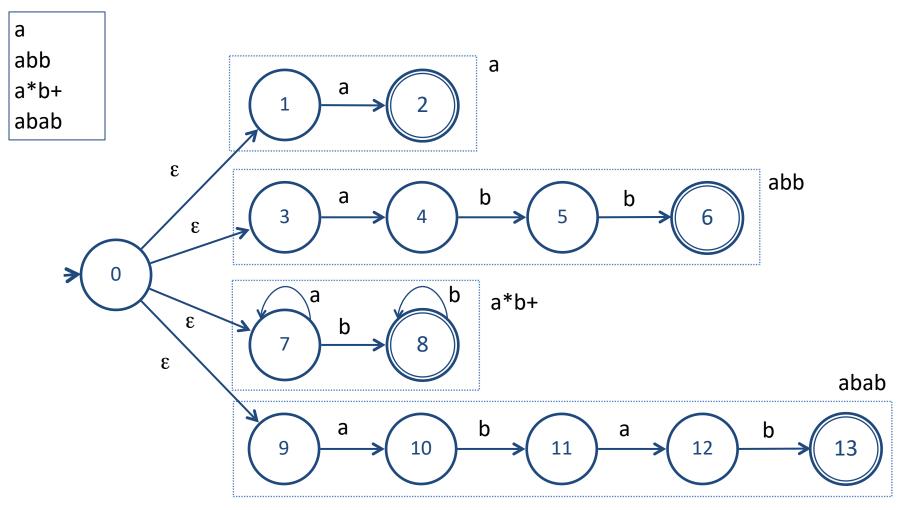


#### What now?

- We have one automaton per pattern.
   Naïve approach: try each automaton separately
- Given a word w:
  - ▶ Try M<sub>1</sub>(*w*)
  - ▶ Try M<sub>2</sub>(w)
  - **)** ...
  - ightharpoonup Try  $M_n(w)$
- Requires "rewinding" after every attempt

## **Combine Automata**

#### combines

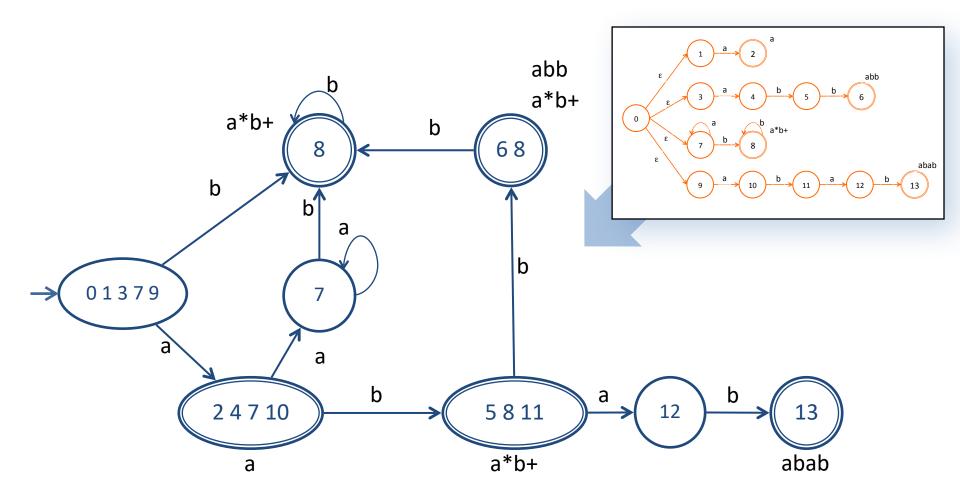


# Ambiguity resolution

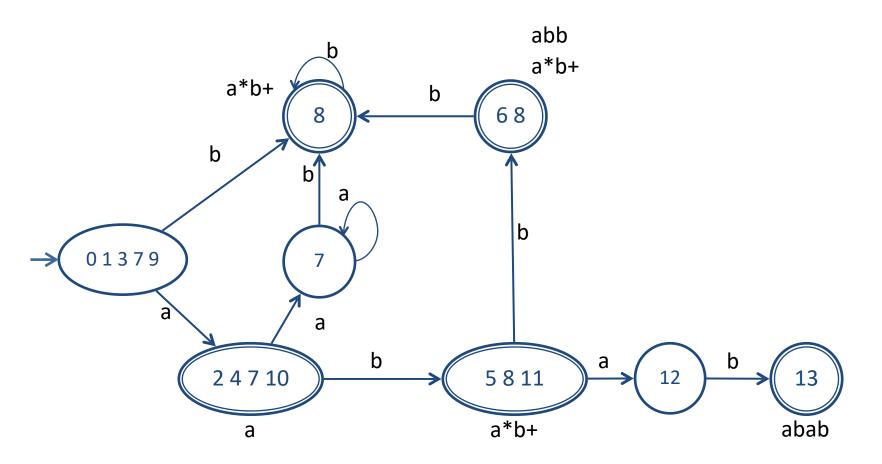
- Recall...
  - \* Longest word
  - \*Tie-breaker based on order of rules when words have same length

- Recipe
  - ▶ Turn NFA to DFA
  - ▶ Run until stuck, remember *last accepting state* this is the token to be returned

# **Corresponding DFA**



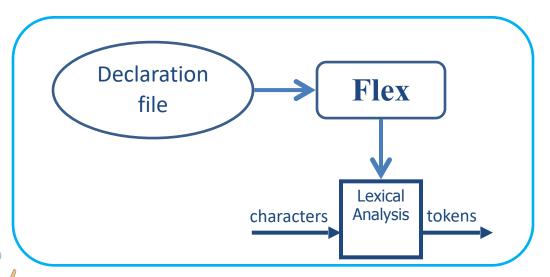
# **Example Inputs**



**abaa**: gets stuck after **aba** in state 12, backs up to state (5 8 11) pattern is a\*b+, token is **abba**: stops after second **b** in (6 8), pattern is abb because it comes first in spec

#### **Good News**

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



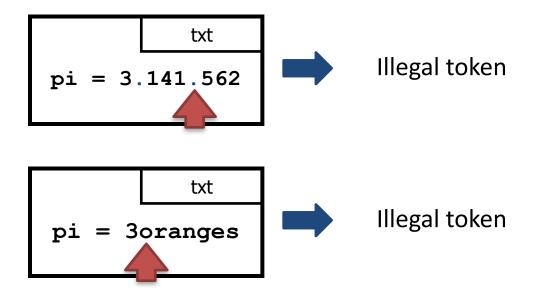
#### Intuitively:

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

### Flex declarations file

```
웅 {
#include <math.h>
int line number = 1;
용 }
WS [\t]
LETTER [a-zA-Z]
DIGIT [0-9]
ID {LETTER}({LETTER}|{DIGIT})*
응응
{DIGIT}+ { printf("number: %d\n", atoi(yytext)); return 1; }
{ID} { return 2;}
{WS} { /* ignore */ }
           { line_number++; /* and ignore */ }
\n
            { return -1; /* ERROR */ }
응응
int main() { return yylex(); }
```

# Errors in Lexical Analysis



Note: How can we make this an error?

# **Error Handling**

Sometimes the lexeme does not match any pattern

```
#@!?&
```

- Instead of aborting, we try to recover
- ► Easiest fix: skip characters until the beginning of a "legitimate" lexeme
- ▶ Alternatives: eliminate/add/replace one letter, replace order of two adjacent letters, etc.

- description in the compilation to continue
- " Con: errors that spread all over

# 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 Flex

