

Drawing Hands  
M.C. Escher, 1948

<http://lara.epfl.ch/cc>

# Computer Language Processing (Compiler Construction)

## Staff:

- Viktor Kuncak – Lectures
- Etienne Kneuss – Labs
- Ravichandhran Kandhadai Madhavan – Exercises

# Computer Language Processing

## Language can be

- Natural language (French, English, ...)
- **computer language** (Java, Scala, C, SQL, ...)
- Mathematical language: a set of strings
  - Can represent both of the above

## We can **process** the language: do something with strings in the language

- manually: mathematical proof, literary criticism
- using **computers**: algorithms that work on strings

## Computer language processing in this course: **processing computer languages using computers**

# Compilers

A typical compiler processes a general-purpose, Turing-complete, language and translates it into the form where it can be efficiently executed

- gcc and clang: map C into machine instructions
- Java compiler: Java source code into bytecodes (.class files)
- Just-in-time (JIT) compiler inside the Java Virtual Machine (JVM): translate .class files into machine instructions (while running the program)

This is the focus of this class and the class project.

# Example: javac

## - from Java to Bytecode

```
...  
while (i < 10) {  
    System.out.println(j);  
    i = i + 1;  
    j = j + 2*i+1;  
}  
...
```

```
javac Test.java  
    → Test.class  
javap -c Test
```

```
4: iload_1  
5: bipush 10  
7: if_icmpge 32  
10: getstatic #2; //System.out  
13: iload_2  
14: invokevirtual #3; //println  
17: iload_1  
18: iconst_1  
19: iadd  
20: istore_1  
21: iload_2  
22: iconst_2  
23: iload_1  
24: imul  
25: iadd  
26: iconst_1  
27: iadd  
28: istore_2  
29: goto 4  
32: return
```

You will build  
a compiler that  
generates such  
code

# Example: gcc

- from C to Intel x86

```
#include <stdio.h>
int main(void) {
    int i = 0;
    int j = 0;
    while (i < 10) {
        printf("%d\n", j);
        i = i + 1;
        j = j + 2*i+1;
    }
}
```

gcc test.c -S  
→ test.s

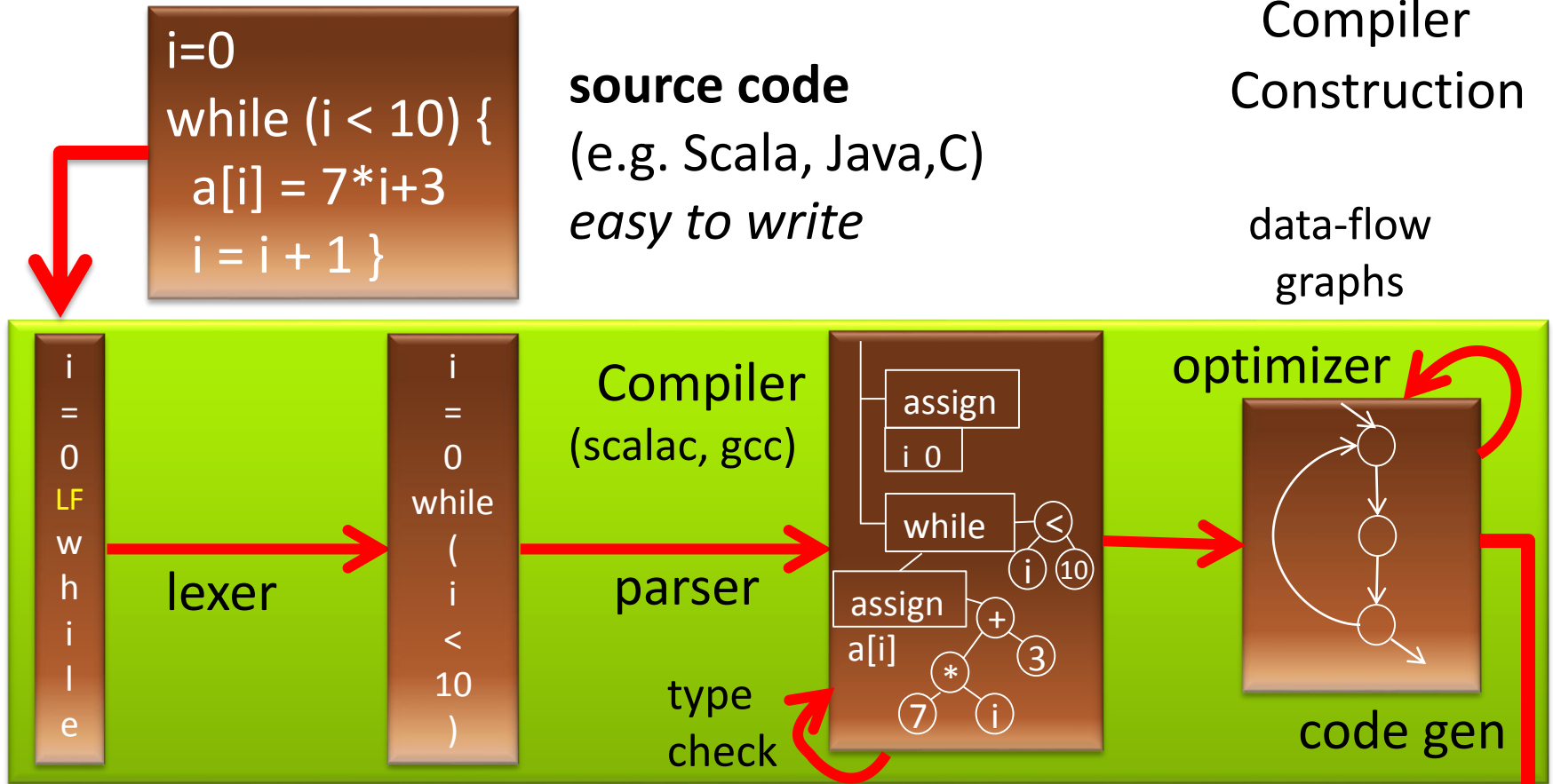
```
jmp .L2
.L3:    movl -8(%ebp), %eax
        movl %eax, 4(%esp)
        movl $.LC0, (%esp)
        call printf
        addl $1, -12(%ebp)
        movl -12(%ebp), %eax
        addl %eax, %eax
        addl -8(%ebp), %eax
        addl $1, %eax
        movl %eax, -8(%ebp)

.L2:    cmpl $9, -12(%ebp)
        jle .L3
```

# Compiler Construction

**source code**  
(e.g. Scala, Java, C)  
*easy to write*

data-flow graphs



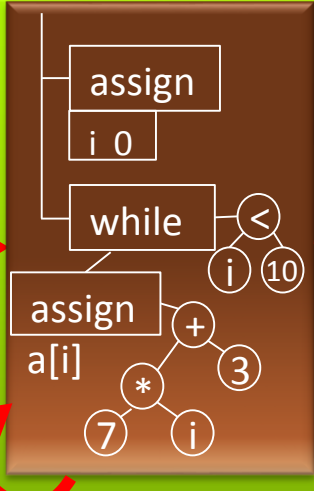
```
i=0
while (i < 10) {
  a[i] = 7*i+3
  i = i + 1 }
```

```
i
=
0
LF
w
h
i
l
e
```

lexer

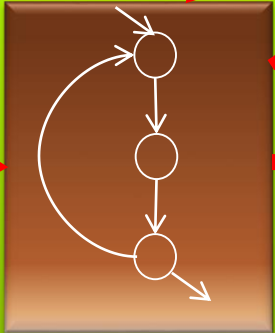
```
i
=
0
while
(
i
<
10
)
```

parser



type check

optimizer



code gen

characters

words

trees

**machine code**  
(e.g. x86, ARM, JVM)  
*efficient to execute*

```
mov R1,#0
mov R2,#40
mov R3,#3
jmp +12
mov (a+R1),R3
add R1, R1, #4
add R3, R3, #7
cmp R1, R2
blt -16
```

# Compilers are Important

**Source code** (e.g. Scala, Java, C, C++, Python) – designed to be easy for programmers to use

- should correspond to way programmers think
- help them be productive: avoid errors, write at a higher level, use abstractions, interfaces

**Target code** (e.g. x86, arm, JVM, .NET) – designed to efficiently run on hardware / VM

- fast, low-power, compact, low-level

**Compilers bridge these two worlds**, they are essential for building complex software

NATIONAL PHYSICAL LABORATORY

TEDDINGTON, MIDDLESEX, ENGLAND

PAPER 2-3

---

AUTOMATIC PROGRAMMING  
PROPERTIES AND PERFORMANCE OF  
FORTRAN SYSTEMS I AND II

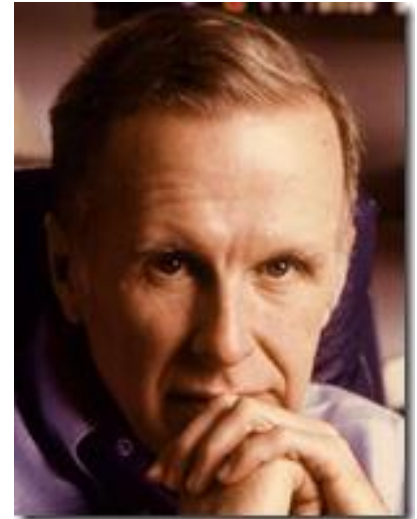
---

by

J. W. BACKUS

To be presented at a Symposium on  
The Mechanization of Thought Processes,  
which will be held at the National Physical  
Laboratory, Teddington, Middlesex, from 24th-  
27th November 1958. The papers and the discussions  
are to be published by H.M.S.O. in the Proceedings  
of the Symposium. This paper should not be repro-  
duced without the permission of the author and of  
the Secretary, National Physical Laboratory.

A pioneering  
compiler:  
**FORTRAN**  
(FORmula  
TRANslator)



Backus-Naur  
Form - BNF

Turing Award  
1977

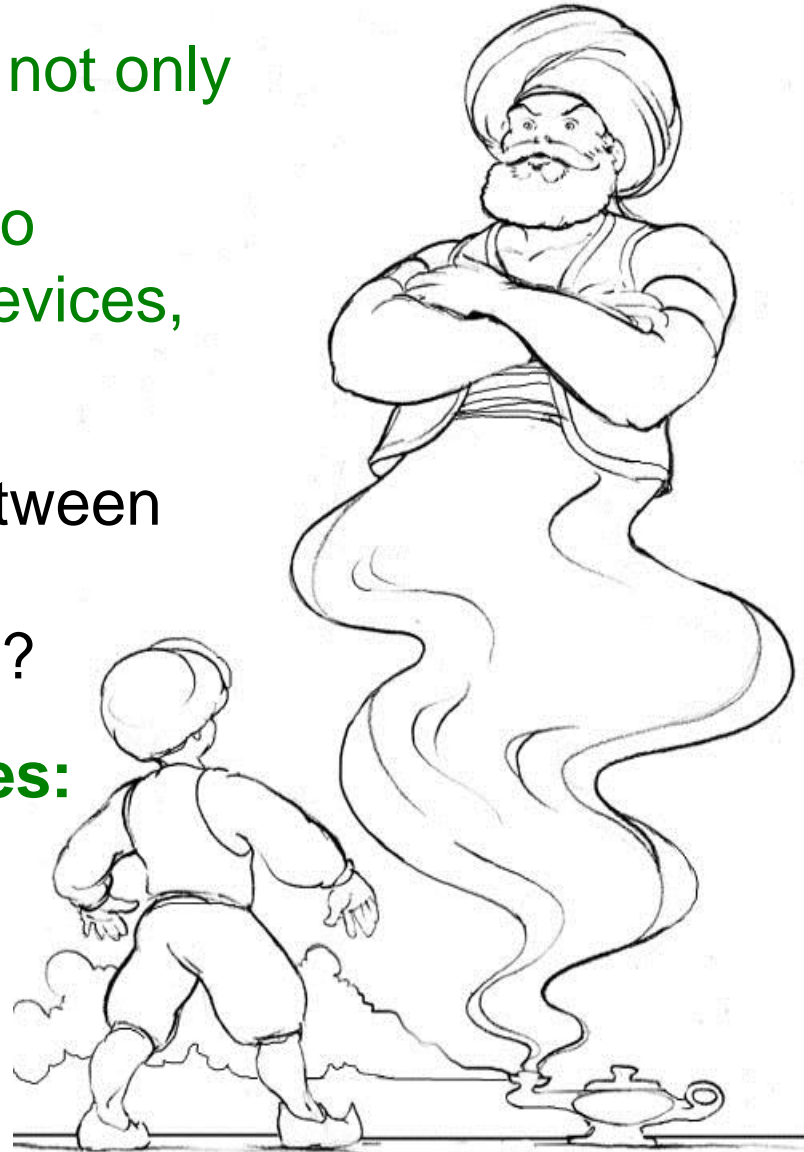


# Challenges for Future

Can target code **commands** include not only execution of commands on standard microprocessors processors, but also automatic design of new hardware devices, and control of physical devices?

Can compilers bridge the gap between wishes and commands, and help humans make the right decisions?

Can source code programs be **wishes**: specification languages, math, natural language phrases, diagrams, other forms of communication closer to engineers and users?



# Some of Topics You Learn in Course

- **Develop a compiler for a Java-like language**
  - Write a compiler from start to end
  - Generates Java Virtual Machine (JVM) code  
(We provide you code stubs, libraries **in Scala**)
- **Compiler generators – using and making them**
- **Analyze complex text**
  - Automata, regular expressions, grammars, parsing
- **Automatically detecting errors in code**
  - name resolution, type checking, data-flow analysis
- **Machine-like code generation**

# Potential Uses of Knowledge Gained

- understand how compilers work, use them better
- gain experience with building complex software
- build compiler for your next great language
- extend language with a new construct you need
- adapt existing compiler to new target platform (e.g. embedded CPU or graphics processor)
- regular expression handling in editors, grep
- build an XML parsing library
- process complex input box in an application (e.g. expression evaluator)
- parse simple natural language fragments

# Schedule and Activities (6 credits)

- All activities take place in INM 202
  - Mondays **10:15**-12:00,
  - Wednesday **8:15**-10:00 and continuing to:
  - Wednesday **10:15**-12:00
- Lectures, Labs, Exercises
- At home
  - Continue with programming the compiler
  - Practice solving problems to prepare for quizzes
- If you need more help, email us:
  - we will arrange additional meetings

# How We Compute Your Grade

The grade is based on a weighted average of:

- 50% : project (submit, explain if requested)
  - submit through our wonderful online system
  - do them in groups of 2, exceptionally 1 or 3
- 25% : quiz for the first part of the course
- 25% : quiz for the second part of the course
  - will be on the last Wednesday of classes

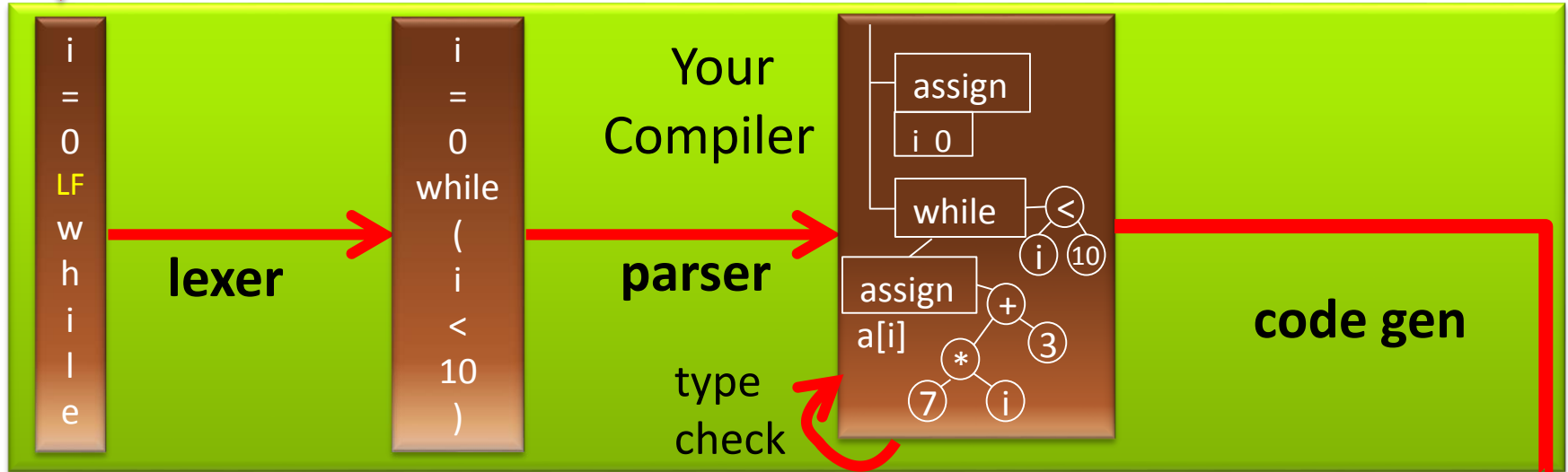
# Collaboration and Its Boundaries

- For clarification questions, discuss them in the mailing list, which we monitor
- Work in groups of **2** for project
  - everyone should know every part of code
  - we may ask you to explain specific parts of code
- Do not copy lab solutions from other groups!
  - we use code plagiarism detection tools
  - we will check if you fully understand your code
- Do the quizzes *individually*
  - You wouldn't steal a handbag.
  - You wouldn't steal a car.
  - You wouldn't steal a compiler!

# Your Compiler Construction

```
i=0  
while (i < 10) {  
  a[i] = 7*i+3  
  i = i + 1 }  
}
```

source code  
simplified Java-like  
language



characters

words

trees

```
21: iload_2  
22: iconst_2  
23: iload_1  
24: imul  
25: iadd  
26: iconst_1  
27: iadd  
28: istore_2
```

Each two weeks you will add next phase

- keep same groups
- it is essential to not get behind schedule
- final addition to compiler is **your choice!**

# EPFL Course Dependencies

- Theoretical Computer Science (CS-251)
  - If have not taken it, check the book “Introduction to the Theory of Computation” by Michael Sipser
- Knowledge of the Scala language (see web)
- Helpful general background
  - Discrete structures (CS-150), Algorithms (CS-250)
- This course provides background for MSc:
  - Advanced Compilers
  - Synthesis Analysis & Verification
  - Foundations of Software



# Course Materials

## Official Textbook:

Andrew W. Appel, Jens Palsberg:

Modern Compiler Implementation in Java  
(2nd Edition). Cambridge University Press, 2002

We do not strictly follow it—understand, not copy

- program in Scala instead of Java
- use pattern matching instead of visitors
- hand-written parsers in the project  
(instead of using a parser generator)

Lectures in course wiki: <http://lara.epfl.ch/w/cc>

# Additional Materials

- Compilers: Principles, Techniques, and Tools (2nd Edition) by Alfred V. Aho, Monica S. Lam, Ravi Sethi, Jeffrey D. Ullman
  - comprehensive
- Compiler Construction by Niklaus Wirth
  - concise, has main ideas

“**Niklaus Emil Wirth** (born February 15, 1934) is a Swiss computer scientist, best known for designing several programming languages, including Pascal, and for pioneering several classic topics in software engineering. In 1984 he won the Turing Award for developing a sequence of innovative computer languages.”

- Additional recent books (2011-2012):
  - **Aarne Ranta**: Implementing Programming Languages
  - **H.Seidl, R.Wilhelm, S.Haack**: Compiler Design (3 vols, Springer)

# Describing the **Syntax** of Languages

# Syntax (from Wikipedia)

...In linguistics, **syntax** (from Ancient Greek σύνταξις "arrangement" from σύν - *syn*, "together", and τάξις - *táxis*, "an ordering") is the study of the principles and rules for constructing phrases and sentences in natural languages.

...In computer science, the **syntax** of a programming language is the set of rules that define the combinations of symbols that are considered to be correctly structured programs in that language.

# Describing Syntax: Why

- Goal: document precisely (a superset of) meaningful programs (for users, implementors)
  - Programs outside the superset: meaningless
  - We say programs inside make *syntactic* sense (They may still be ‘wrong’ in a deeper sense)
- Describing syntactically valid programs
  - There exist arbitrarily long valid programs, we cannot list all of them explicitly!
  - Informal English descriptions are imprecise, cannot use them as language reference

# Describing Syntax: How

- Use theory of formal languages (from TCS)
  - regular expressions & finite automata
  - context-free grammars
- We can use such precise descriptions to
  - document what each compiler should support
  - manually derive compiler phases (lexer, parser)
  - automatically construct these phases using compiler generating tools
- We illustrate this through an example

# *While* Language – Idea

- Small language used to illustrate key concepts
- Simpler than the language for which you implement your compiler
- ‘while’ and ‘if’ are the control statements
  - no procedures, no exceptions
- the only variables are of ‘int’ type
  - no variable declarations, they are initially zero
  - no objects, pointers, arrays

# While Language – Example Programs

```
while (i < 100) {  
  j = i + 1;  
  while (j < 100) {  
    println(" ",i);  
    println(", ",j);  
    j = j + 1;  
  }  
  i = i + 1;  
}
```

Nested loop

```
x = 13;  
while (x > 1) {  
  println("x=", x);  
  if (x % 2 == 0) {  
    x = x / 2;  
  } else {  
    x = 3 * x + 1;  
  }  
}
```

Does the program terminate  
for every initial value of x?  
(Collatz conjecture - open)

Even though it is simple, while is Turing-complete.



# Reasons for Unbounded Program Length

constants of any length

variable names of any length

String constants of any length

(words - tokens)

```
while (i < 100) {  
    j = i + 5*(j + 2*(k + 7*(j+k) + i));  
    while (293847329 > j) {  
        while (k < 100) {  
            someName42a = someName42a + k;  
            k = k + i + j;  
            println("Nice number", k)  
        }  
    }  
}
```

nesting of expressions

nesting of statements

# Tokens (Words) of the *While* Language

Ident ::=

letter (letter | digit)\*

regular  
expressions



integerConst ::=

digit digit\*

stringConst ::=

" AnySymbolExceptQuote\* "

keywords

if else while println

special symbols

( ) && < == + - \* / % ! - { } ; ,

letter ::= a | b | c | ... | z | A | B | C | ... | Z

digit ::= 0 | 1 | ... | 8 | 9

# Double Floating Point Constants

Different rules in different languages

1)  $\text{digit digit}^* [ . ] [ \text{digit digit}^* ]$

2)  $\text{digit digit}^* [ . \text{digit digit}^* ]$

3)  $\text{digit}^* . \text{digit digit}^*$

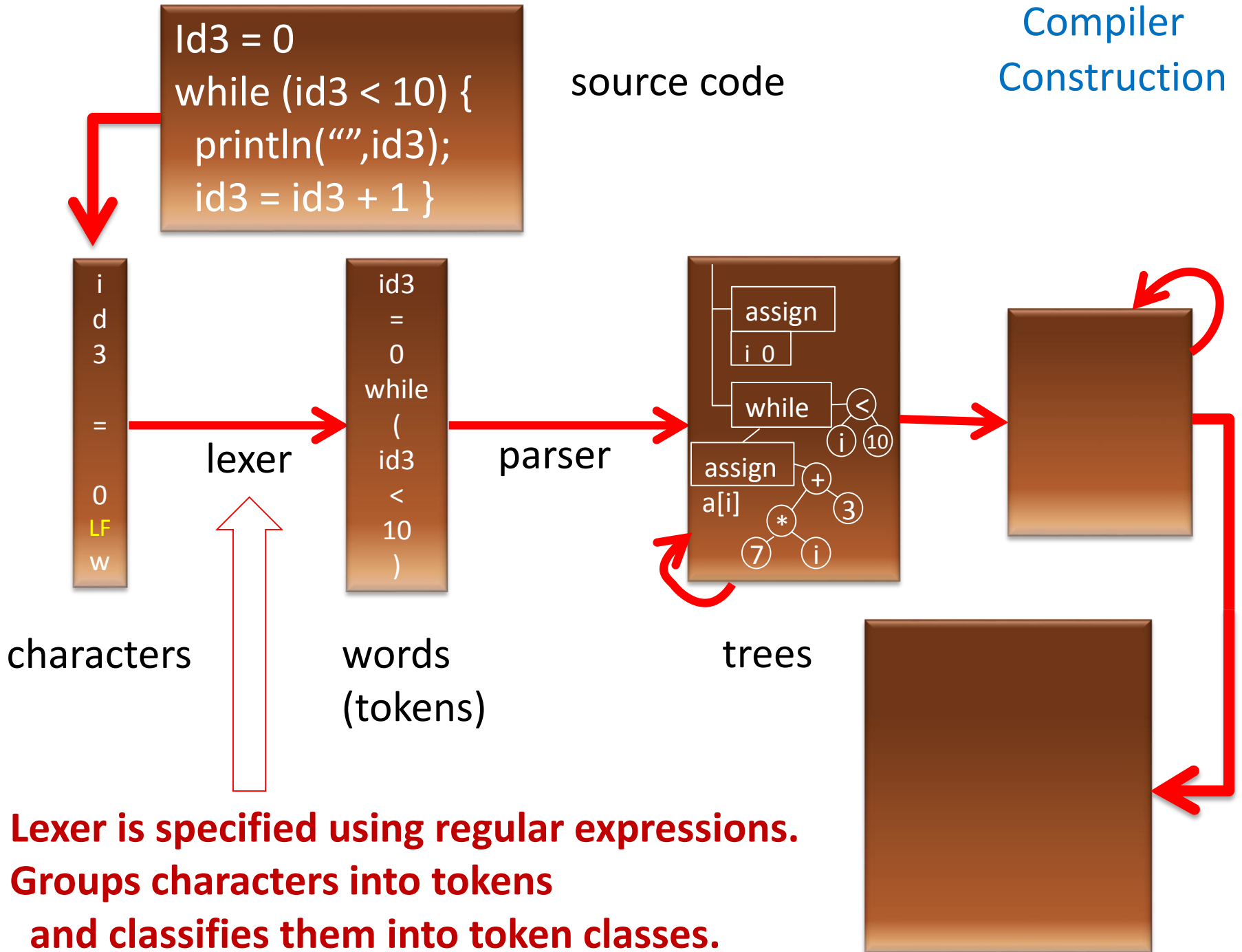
4)  $\text{digit digit}^* . \text{digit digit}^*$

# Identifiers

```
while (i < 100) {  
  j = i + 5*(j + 2*(k + 7*(j+k) + i));  
  while (293847329 > j) {  
    while (k < 100) {  
      someName42a = someName42a + k;  
      k = k + i + j;  
      println("Nice number", k)  
    }  
  }  
}
```

letter (letter | digit)\*

# Compiler Construction



**Lexer is specified using regular expressions.  
Groups characters into tokens  
and classifies them into token classes.**

# More Reasons for Unbounded Length

constants of  
any length

```
while (i < 100) {  
    j = i + 5*(j + 2*(k + 7*(j+k) + i));
```

variable names  
of any length

```
    while (293847329847 > j) {
```

```
        while (k < 100) {
```

```
            someName42a = someName42a + k;
```

```
            k = k + i + j;
```

```
            println("Nice number", k)
```

String constants  
of any length

```
        }  
    }  
}
```

(words - tokens)

nesting of  
expressions

nesting of  
statements

(sentences)

# Sentences of the *While* Language

We describe sentences using **context-free grammar (Backus-Naur form)**. Terminal symbols are tokens (words)

program ::= statmt\*

statmt ::= println( stringConst , ident )

| ident = expr

| if ( expr ) statmt (else statmt)?

| while ( expr ) statmt

| { statmt\* }

nesting of  
statements



expr ::= intLiteral | ident

| expr ( && | < | == | + | - | \* | / | % ) expr

| ! expr | - expr

nesting of  
expressions



# *While* Language without Nested Loops

statmt ::= println( stringConst , ident )

| ident = expr

| if ( expr ) statmt (else statmt)?

| while ( expr ) statmtww

| { statmt\* }

statmtww ::= println( stringConst , ident )

| ident = expr

| if ( expr ) statmtww (else statmtww)?

| { statmtww\* }



# Compiler Construction

```
Id3 = 0
while (id3 < 10) {
  println("",id3);
  id3 = id3 + 1 }
}
```

source code



i  
d  
3  
  
=  
  
0  
LF  
w

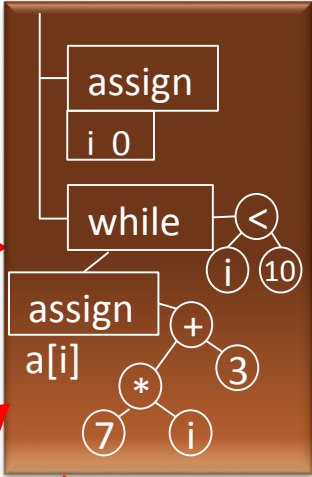
characters

lexer

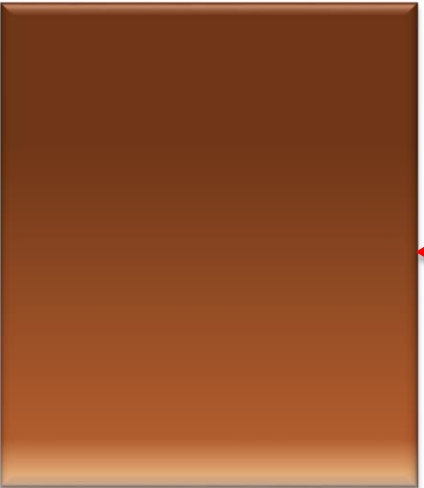
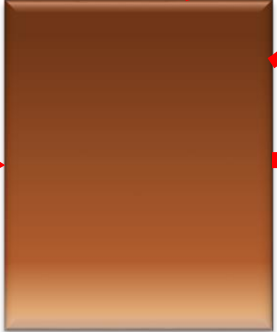
id3  
=  
0  
while  
(  
id3  
<  
10  
)

words  
(tokens)

parser

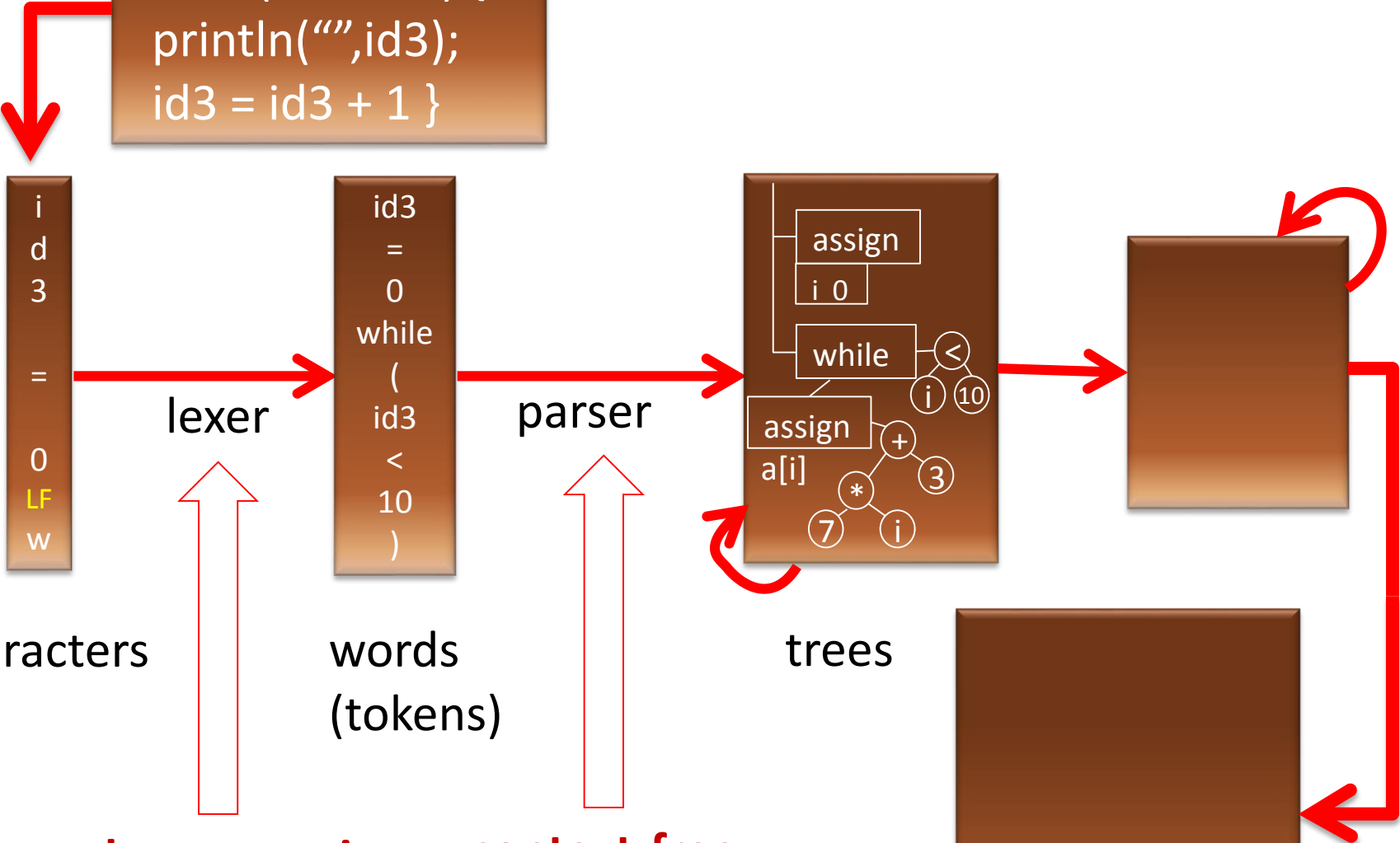


trees



regular expressions  
for tokens

context-free  
grammar



# Abstract Syntax - Trees

To get abstract syntax (trees, cases classes),  
start from context-free grammar for tokens, then

- remove punctuation characters
- interpret rules as **tree descriptions**, not string descriptions

statmt ::= println( stringConst , ident )	PRINT(String,ident)
ident = expr	ASSIGN(ident,expr)
if ( expr ) statmt (else statmt)?	IF(expr,stmt,Option[statmt])
while ( expr ) statmt	WHILE(expr,statmt)
{ statmt* }	BLOCK(List[statmt])

concrete syntax

Scala trees for this  
abstract syntax

**abstract class** statmt

**case class** PRINT(id:ident) **extends** statmt

**case class** ASSIGN(id:ident, e:expr) **extends** statmt

**case class** IF(e:expr, s1:statmt, s2:Option[statmt])

**extends** statmt ...

**abstract syntax**

# ocaml vs Haskell vs Scala

ocaml:

```
type tree = Leaf | Node of tree * int * tree
```

Haskell:

```
data Tree = Leaf | Branch Tree Int Tree
```

Scala:

```
abstract class Tree  
case object Leaf extends Tree  
case class Node(left:Tree,x:Int,right:Tree)  
  extends Tree
```

# Example of Parsing

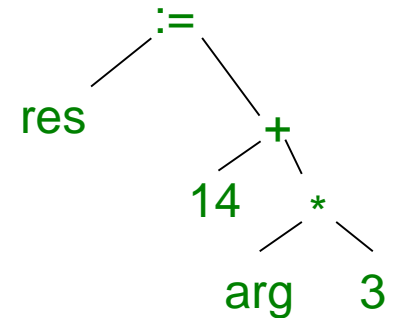
res = 14 + arg \* 3

Lexer:

res = 14 + arg \* 3

Parser:

```
ASSIGN(res,  
  PLUS(CONST(14),  
    TIMES(VAR(arg),CONST(3))))
```

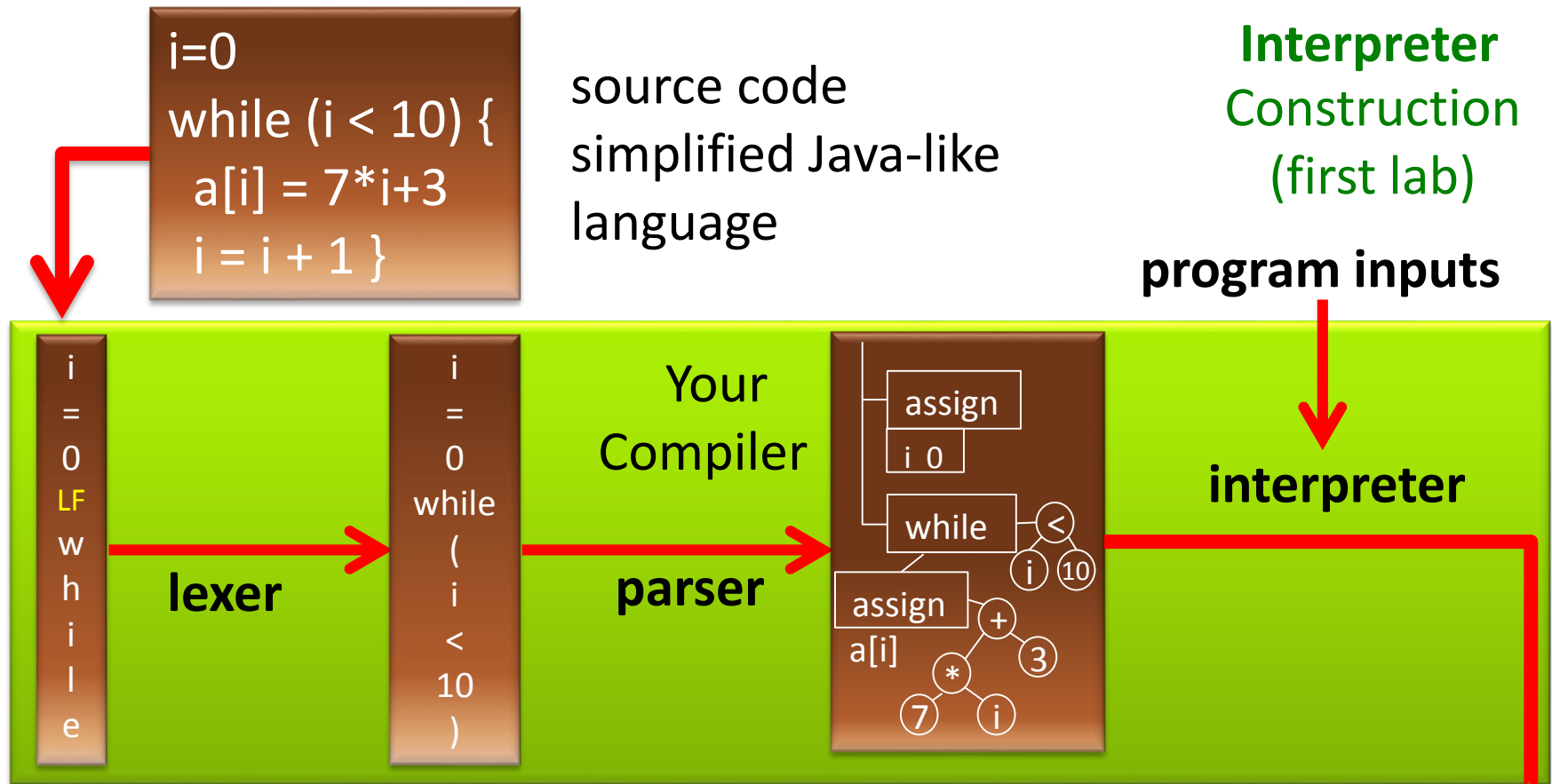


Code generator then “prints” this tree into instructions.

# Interpreters

# What is an interpreter?

- Interpreter is a simpler way to implement a language
- Usually it is easier to build than a compiler
- It can be used as one way to define the meaning of programs: programs should compute whatever the interpreter returns
- Your first lab: build an interpreter
  - the front end (lexer and parser) will be given to you as a class file



characters

words

trees

### Differences with the compiler

- does not generate code
- waits for program input, evaluates program tree, computes the result

# Reminder about Formal Languages



# Languages Formally

- A *word* is a finite, possibly empty, sequence of elements from some set  $\Sigma$

$\Sigma$  – *alphabet*,  $\Sigma^*$  – set of all words over  $\Sigma$

- For *lexer*: characters; for *parser*: token classes
- $uv$  denotes concatenation of words  $u$  and  $v$
- By a *language* we mean a subset of  $\Sigma^*$

– union, intersection, complement wrt.  $\Sigma^*$

$$L_1 L_2 = \{ u_1 u_2 \mid u_1 \text{ in } L_1, u_2 \text{ in } L_2 \}$$

$$L^0 = \{\varepsilon\}$$

$$L^{k+1} = L L^k$$

$$L^* = \bigcup_k L^k \quad (\text{Kleene star})$$

# Are there finitely many tokens?

- There are finitely many **token classes**
  - identifier
  - string
  - {
  - }
  - (
  - ... (many, but finitely many)

There is unbounded number of *instances* of token classes identifier and string

When we discuss grammars, we work with token classes.

# Examples of Languages

$$\Sigma = \{a, b\}$$

$$\Sigma^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$$

Examples of two languages, subsets of  $\Sigma^*$  :

$$L_1 = \{a, bb, ab\} \quad (\text{finite language, three words})$$

$$L_2 = \{ab, abab, ababab, \dots\}$$

$$= \{(ab)^n \mid n \geq 0\} \quad (\text{infinite language})$$

# Examples of Operations

$L = \{a, ab\}$

$LL = \{aa, aab, aba, abab\}$

$L^* = \{a, ab, aa, aab, aba, abab, aaa, \dots\}$

(is  $bb$  inside  $L^*$  ?)

$= \{w \mid \text{immediately before each } b \text{ there is a } a\}$