

### ΕΠΛ323 - Θεωρία και Πρακτική Μεταγλωττιστών

#### Lecture 1 Introduction

Elias Athanasopoulos eliasathan@cs.ucy.ac.cy



- Office hours
  - Every Tuesday, 10:00 12:00, B105 (ΘΕΕ01)
- Credits: 7.5 ECTS
- Lecture Timetable
  - Monday, Thursday, 12:00 13:30, 006 (ΧΩΔ01)
  - Wednesday, 9:00-10:00, 110 (XΩΔ01), only when announced!



- Labs (Dr. Petros Panayi)
  - Every Wednesday, B103 (ΘΕΕ01)
  - Group 1: 15:00 17:00
  - Group 2: 17:00 19:00



- Score
  - Programming Assignments (Quizzes), 15%
  - Project (multiple steps), 15%
  - Midterm, 30%
  - Final, 40%
- Requirements
  - Average grade of written exams should be at least 4.5
  - Final grade should be at least 5



- Required courses
  - CS132 (C programming),
  - CS211 (Complexity),
  - CS231 (Data Structures)

Students with conflicts should let me know



- Reading material
  - Dragon book, Compilers Principles, Techniques, and Tools (2nd Edition), A.V. Aho, M.S. Lam, R.
     Sethi and J.D. Ullman, Prentice Hall, 2006.

### Communication



- Labs
  - Blackboard
- Lectures

– http://www.cs.ucy.ac.cy/courses/EPL323



(Chapter 1 from Dragon Book)

### INTRODUCTION

# Why this course is important?



- Many core concepts of CS in a real-life setting
  - CS132: Behind the scenes of C programming
  - CS211, CS231: DFAs, parsing algorithms
  - CS372, CS221: Assembly, Computer Architecture
- Compilers are everywhere
  - From exotic devices (IoT) to web browsers
     (JavaScript engines, DOM parsers, etc.)
- Tweaking compilers for Security

### Compiler is a tough one



- Significant effort to build a compiler from scratch
  - First Fortran compiler took 18 staff-years
- Nowadays, many tools exist to help
  - Our understanding is much better
  - Better techniques, better programming languages

### What is a compiler?



- A compiler is a program that
  - reads a program written in one language (source)
  - and translates it to an equivalent program in another language (target)
  - **important**: error reporting during translation



### Examples

![](_page_11_Picture_1.jpeg)

- GCC (Gnu Compiler Collection)
   gcc, g++, javac, etc.
- LLVM (Low Level Virtual Machine)
  - clang, clang++
- "Compilers" are everywhere,
  - Pretty printers (i.e., color syntax in editors), static checkers, interpreters (i.e., scripting languages), etc.

# Analysis-Synthesis Model

![](_page_12_Picture_1.jpeg)

- There are two parts in compilation:
  - Analysis
  - Synthesis
- Analysis
  - Breaks up the *source program* to subparts and creates intermediate representation(s)
- Synthesis
  - Constructs the *target program* from intermediate representation(s)

### Example 1 (LaTeX)

![](_page_13_Picture_1.jpeg)

```
\begin{table}[tb]
    \centering
    \caption{We name gadgets based on their type (prefix), payload (body),
    and exit instruction (suffix). In total, we name 2$\times$3$\times$3=18
    different gadget types.}
    begin{tabular}{|c|c|c|}
        \hline
        \textbf{Gadget type} & \textbf{Payload instructions} &
            \textbf{Exit instruction} \\
        \hline
        {Prefix} & {Body} & {Suffix} \\
        \hline
        \begin{tabular}{1}
        CS - Call site\\
        EP - Entry point\\
                                           TABLE II: We name gadgets based on their type (prefix),
        \end{tabular} &
                                           payload (body), and exit instruction (suffix). In total, we name
        . . .
                                           2 \times 3 \times 3 = 18 different gadget types.
```

\$ pdflatex main.tex

Gadget type	<b>Payload instructions</b>	Exit instruction
Prefix	Body	Suffix
CS - Call site EP - Entry point	IC - Indirect call F - Fixed function call none - Other instructions	R - Return IC - Indirect call IJ - Indirect jump

## Example 2 (Database)

![](_page_14_Picture_1.jpeg)

SELECT AVG(grade)
FROM student, class
WHERE
class.name = "epl223" AND

class.id = student.id;

![](_page_14_Figure_4.jpeg)

### Requirements

![](_page_15_Picture_1.jpeg)

- Compiler
  - Reliability
  - Fast execution
  - Low memory overhead
  - Good error reporting
  - Error recovery
  - Portability
  - Maintainability
- Target program
  - Fast execution
  - Low memory overhead

# Source code/program

![](_page_16_Picture_1.jpeg)

• Easy to read/write by human

```
int expr(int n) {
    int d;
```

d = 4 \* n \* n \* (n + 1) \* (n + 1);

```
return d;
}
```

# Assembly and Machine Code

![](_page_17_Picture_1.jpeg)

- Optimized for execution by a machine (CPU)
- Less descriptive
- Hard to be processed by a human

lda \$30,-32(\$30)	
stq \$26,0(\$30)	
stq \$15,8(\$30)	
bis \$30,\$30,\$15	
bis \$16,\$16,\$1	
stl \$1,16(\$15)	
lds \$f1,16(\$15)	
sts \$f1,24(\$15)	
ldl \$5,24(\$15)	
bis \$5,\$5,\$2	
s4addq \$2,0,\$3	
ldl \$4,16(\$15)	
mull \$4,\$3,\$2	
ldl \$3,16(\$15)	

### Optimizations

![](_page_18_Picture_1.jpeg)

Compilers have several layers of optimizations

}

#### No optimizations

· 3	
.expr:	
stw 31,-4(1)	lwz 11,64(31)
stwu 1,-40(1)	addi 9,11,1
mr 31,1	mullw 0,0,9
stw 3,64(31)	stw 0,24(31)
lwz 0,64(31)	lwz 0,24(31)
mr 9,0	mr 3,0
slwi 0,9,2	b L2
lwz 9,64(31)	L2:
mullw 0,0,9	lwz 1,0(1)
lwz 11,64(31)	lwz 31,-4(1)
addi 9,11,1	blr
mullw 0,0,9	

```
int expr(int n) {
    int d;
    d = 4 * n * n * (n + 1) * (n + 1);
    return d;
```

```
Optimizations
$ gcc -03
```

.expr: addi 9,3,1 slwi 0,3,2 mullw 3,3,0 mullw 3,3,9 mullw 3,3,9 blr

### Cross-compiler

![](_page_19_Picture_1.jpeg)

 Compilers can generate code for different machines (targets) int expr(int n) {

}

For x86

\$ gcc -03 -b i586

expr:	
pushl	%ebp
movl	%esp, %ebp
movl	8(%ebp), %eax
leal	1(%eax), %edx
imull	%eax, %eax
imull	%edx, %eax
imull	%edx, %eax
sall	\$2, %eax
popl	%ebp
ret	

```
nt expr(int n) {
    int d;
    d = 4 * n * n * (n + 1) * (n + 1);
    return d;
```

```
For PowerPC
$ gcc -03 -b powerpc
```

```
.expr:
addi 9,3,1
slwi 0,3,2
mullw 3,3,0
mullw 3,3,9
mullw 3,3,9
blr
```

### **Compilation life cycle**

![](_page_20_Picture_1.jpeg)

- Phases
  - Source code is transformed to intermediate representations
  - Each intermediate representation is suitable for a particular processing (lexical, syntax, optimization, etc.)
- In each phase the program is translated to a form closer to the machine representation and less similar to the (human-oriented) source representation

### **Compiler Phases**

![](_page_21_Picture_1.jpeg)

![](_page_21_Figure_2.jpeg)

# Analysis of the source program

- Linear analysis
  - Source is treated as a stream of characters (left-toright) and is grouped into tokens
- Hierarchical analysis
  - Tokens are further grouped in larger grammatical structures (e.g., nested parentheses and blocks)
- Semantic analysis
  - Certain checks are performed to ensure the validity of the identified grammatical structures

### Lexical Analysis

Λεξιλογική Ανάλυση

- Linear scanning
- Consider the expression position := initial + rate \*60
- Lexical analysis produces
   id(1) op(:=) id(2) op(+) id(3) op(\*) cons(60)
   id: identifier, op: operator, cons: constant
- Symbol Table

1	position	
2	initial	
3	rate	
4		

![](_page_23_Picture_7.jpeg)

## Syntax Analysis

Συντακτική Ανάλυση

- Hierarchical
- Involves grouping the tokens into grammatical phases
- Constructs the structure with the token relationship id(1)
   position := initial + rate \* 60
   id(3)

### Simple Grammar

![](_page_25_Figure_1.jpeg)

- The hierarchical structure of the program is usually expressed by recursive rules
  - 1. Any *identifier* is an expression
  - 2. Any *number* is an expression
  - 3. If *expression*<sub>1</sub> and *expression*<sub>2</sub> are expressions, then so are:

 $expression_1 + expression_2$  $expression_1 * expression_2$ (  $expression_1$  )

# Applying the grammar

![](_page_26_Picture_1.jpeg)

![](_page_26_Figure_2.jpeg)

### Semantic Analysis

Σημασιολογική Ανάλυση

• Checks the program for semantic errors

60

- Gathers type information
- Operands and operators
- Type-checking

position := initial + rate \*

![](_page_27_Figure_6.jpeg)

![](_page_27_Picture_7.jpeg)

### Error detection and reporting

![](_page_28_Picture_1.jpeg)

- All phases can issue errors
- A compiler that stops at the first error is not helpful
- Most of the errors are handled in the syntax/semantic analysis phases
  - Lexical analysis detects errors where a stream of characters does not form a valid token
  - Syntax analysis detects errors where the stream of valid tokens violate the structure rules (syntax)
  - Semantic analysis detects errors where the syntax is valid by the operation not (adding an array with a real number)

### Intermediate Code and Optimization

Ενδιάμεσος Κώδικας και Βελτιστοποίηση

![](_page_29_Picture_2.jpeg)

• Each phase produces intermediate code

temp1 := id(3) \* 60.0
id(1) := id(2) + temp1

three-address code: a simple assembly-like language, which consists of instructions, each of which has at most three operands

![](_page_29_Figure_8.jpeg)

### **Code Generation**

Παραγωγή Κώδικα

![](_page_30_Picture_2.jpeg)

- The last phase of the compiler is the generation of the target code
- Register allocation
  - Each expression should use registers that are available
- Relocation information
  - Variables are stored in relocatable addresses

MOVF	id3, R2
MULF	#60.0, R2
MOVF	id2, R1
ADDF	R2, R1
MOVF	R1, id1

### Generic Picture Compiler and Friends

![](_page_31_Picture_1.jpeg)

![](_page_31_Figure_2.jpeg)

### Front and Back ends

- Separation of common tasks
- Makes design and implementation easier
- K compilers for N machines
  - N back ends, K front ends
  - Instead of K\*N compilers

![](_page_32_Figure_6.jpeg)

### Passes

![](_page_33_Picture_1.jpeg)

- A *pass* is when the compiler reads the source code (or intermediate files)
- The number of passes depends on the source and target language and the running environment
- Different phases that cooperate can be grouped to a single pass (not always possible)
- When grouping is not possible
  - Backpatching: leave empty information that is going to be filled by a later phase/pass

# **Compiler-construction Tools**

![](_page_34_Picture_1.jpeg)

- Parser generators
- Scanner generators
- Syntax-directed translation engines
- Automatic code generators
- Data-flow engines