

Compilerconstructie

najaar 2012

<http://www.liacs.nl/home/rvw11et/coco/>

Rudy van Vliet

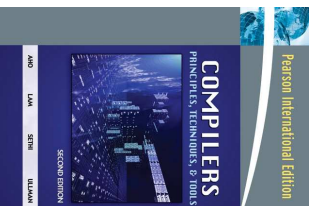
kamer 124 Snellius, tel. 071-527 5777
rvvliet(at)liacs.nl

college 1, dinsdag 4 september 2012

Overview

1

Course Outline

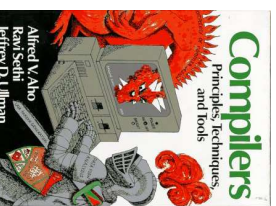


- In class, we discuss the theory using the 'dragon book' by Aho et al.
- The theory is applied in the practicum to build a compiler that converts Pascal code to MIPS instructions.

A.V. Aho, M.S. Lam, R. Sethi, and J.D. Ullman, Compilers: Principles, Techniques, & Tools, Addison-Wesley, 2007, ISBN: 978-0-321-49169-5 (international edition).

3

Earlier edition



- Dragon book has been revised in 2006
- In Second edition good improvements are made
 - Parallelism
 - *
 - * Array data-dependence analysis
- First edition may also be used

A.V. Aho, R. Sethi, and J.D. Ullman, Compilers: Principles, Techniques, and Tools, Addison-Wesley, 1986, ISBN-10: 0-201-10085-6 / 0-201-10194-7 (international edition).

5

Course Outline

- Grading: Combination of the grades from the written exam and the practicum
- You need to pass all 4 assignments to obtain a grade
- Final grade is only accepted if all grades are ≥ 5.5
- Then, you obtain 6 EC

Studying only from the lecture slides may not be sufficient. Relevant book chapters will be given.

7

Why this course

It's part of the general background of a software engineer

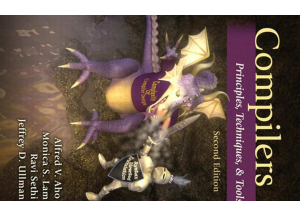
- How do compilers work?
- How do computers work?
- What machine code is generated for certain language constructs?
- Working on a non-trivial programming project

After the course

- Know how to build a compiler for a simplified progr. language
- Know how to use compiler construction tools, such as generators for scanners and parsers
- Be familiar with compiler analysis and optimization techniques

2

Course Outline



- In class, we discuss the theory using the 'dragon book' by Aho et al.
- The theory is applied in the practicum to build a compiler that converts Pascal code to MIPS instructions.

A.V. Aho, M.S. Lam, R. Sethi, and J.D. Ullman, Compilers: Principles, Techniques, & Tools, Addison-Wesley, 2006, ISBN: 978-0-321-54798-9.

4

Course Outline

- Contact
 - Room 124, tel. 071-5275777, rvvliet(at)liacs.nl
 - Course website: <http://www.liacs.nl/home/rvvliet/coco>
- Practicum
 - 4 self-contained assignments
 - These assignments are done by groups of two persons
 - Assignments are submitted by e-mail
 - Assistants: Teddy Zhai, Sven van Haastregt

6

Course Outline

- (tentative)
1. Overview
 2. Lexical Analysis
 3. Syntax Analysis Part 1
 4. Syntax Analysis Part 2
 5. **Assignment 1**
 6. Static Type Checking
 7. **Assignment 2**
 8. Intermediate Code Generation
 9. **Assignment 3**
 10. Code Generation
 11. Code optimization
 12. **Assignment 4**
 13. Daedalus

8

Practicum

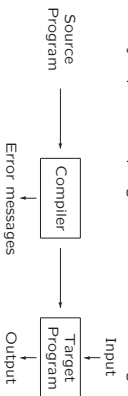
- Assignment 1: Calculator
- Assignment 2: Parsing & Syntax tree
- Assignment 3: Intermediate code
- Assignment 4: Assembly generation

2 academic hours of Lab session + 3 weeks to complete (except assignment 1)

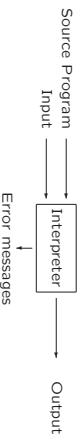
9

Compilers and Interpreters

- **Compilation:**
Translation of a program written in a source language into a semantically equivalent program written in a target language



- **Interpretation:**
Performing the operations implied by the source program.

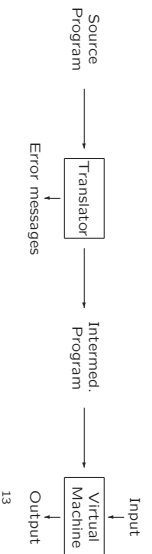


11

Compilers and Interpreters

- Hybrid compiler (Java):

- Translation of a program written in a source language into a semantically equivalent program written in an intermediate language (bytecode)
- Interpretation of intermediate program by virtual machine, which simulates physical machine



13

Other tools that use A-S Model

- Editors (syntax highlighting, text auto completion)
- Text formatters (L^AT_EX, MS Word)

15

Short History of Compiler Construction

- Formerly 'a mystery', today one of the best known areas of computing
- 1957 Fortran first compilers
- (arithmetic expressions, statements, procedures)
- 1960 Algol first formal language definition (grammars in Backus-Naur form, block structure, recursion, ...)
- 1970 Pascal user-defined types, virtual machines (P-code)
- 1985 C++ object-orientation, exceptions, templates
- 1995 Java just-in-time compilation

We only consider **imperative languages**
Functional languages (e.g., Lisp) and logical languages (e.g., Prolog) require different techniques.

10

Compilers and Interpreters

- **Compiler:** Translates source code into machine code, with scanner, parser, ... , code generator
- **Interpreter:** Executes source code 'directly', with scanner, parser
Statements in, e.g., a loop are scanned and parsed again and again

12

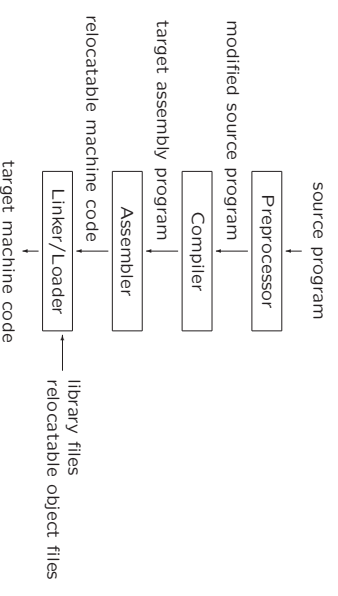
Analysis-Synthesis Model of Compilation

There are two parts to compilation:

- **Analysis**
 - Determines the operations implied by the source program which are recorded in an intermediate representation, e.g., a tree structure
- **Synthesis**
 - Takes the intermediate representation and translates the operations therein into the target program

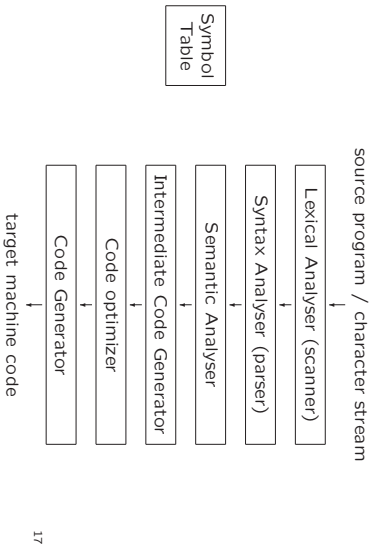
14

Compilation flow



16

The Phases of a Compiler



Symbol Table

Lexical Analyser (scanner)

Syntax Analyser (parser)

Semantic Analyser

Intermediate Code Generator

Code optimizer

Code Generator

target machine code

17

The Phases of a Compiler

Character stream:

position = initial + rate * 60

Lexical Analyser (scanner)

Token stream:

(id, 1) (=) (id, 2) (+) (id, 3) (*) (60)

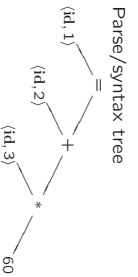
18

The Phases of a Compiler

Token stream:

(id, 1) (=) (id, 2) (+) (id, 3) (*) (60)

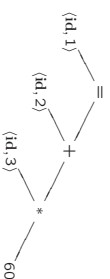
Syntax Analyser (parser)



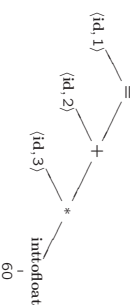
19

The Phases of a Compiler

Parse/syntax tree



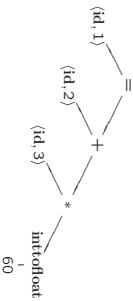
Semantic Analyser (parser)



20

The Phases of a Compiler

Parse/syntax tree



Intermediate Code Generator

Intermediate code (three-address code):

```

t1 = intfloat(60)
t2 = id3 * t1
t3 = id2 + t2
idt = t3
  
```

21

The Phases of a Compiler

Intermediate code (three-address code):

```

t1 = intfloat(60)
t2 = id3 * t1
t3 = id2 + t2
idt = t3
  
```

Code Optimizer

```

t1 = id3 * 60.0
idt = id2 + t1
  
```

22

The Phases of a Compiler

```

t1 = id3 * 60.0
idt = id2 + t1
  
```

Code Generator

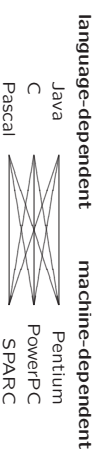
```

LDF R2, id3
MULF R2, R2, #60.0
LDF R1, id2
ADDF R1, R1, R2
STF id1, R1
  
```

23

The Grouping of Phases

- Front End: scanning, parsing, semantical analysis (source code → intermediate representation)
- Back End: code optimizing, code generation (intermediate representation → target machine code)



24

Passes: Single-Pass Compilers

Phases work in an interleaved way

```
do
  scan token
  parse token
  check token
  generate code for token
while (not eof)
```

Portion of code is generated while reading portion of source program

25

Compiler-Construction Tools

Software development tools are available to implement one or more compiler phases

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

27

What is a grammar?

Context-free grammar is a 4-tuple with

- A set of tokens (*terminal* symbols)
- A set of *nonterminals* (syntactic variables)
- A set of *productions*: rules how to decompose nonterminals
- A designated *start/ symbol* (nonterminal)

Example: Context-free grammar for simple expressions:

$$G = (\{list, digit\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, list)$$

with productions P:

```
list → list + digit
list → list - digit
list → digit
digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

29

Derivation (Example)

$$G = (\{list, digit\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, list)$$

$$list \rightarrow list + digit \mid list - digit \mid digit$$

$$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

Example: 9-5+2

```
list ⇒ list + digit
      ⇒ list - digit + digit
      ⇒ digit - digit + digit
      ⇒ 9 - digit + digit
      ⇒ 9 - 5 + digit
      ⇒ 9 - 5 + 2
```

This is an example of **leftmost derivation**, because we replaced the **leftmost nonterminal** (underlined) in each step

31

Passes: Multi-Pass Compilers

Phases are separate 'programs', which run sequentially

```
characters → Scanner → tokens → Parser → tree
                                     → Semantical analyser → ... → code
```

Each phase reads from a file and writes to a new file.

Time vs memory

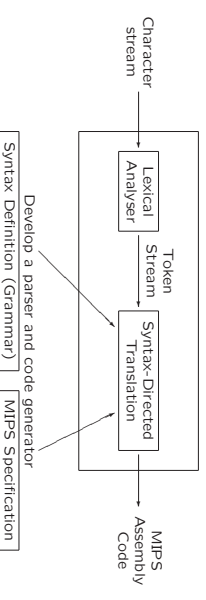
Why multi-pass?

- If the language is complex
- If portability is important

Today: often two-pass compiler

26

The Structure of our compiler



Syntax directed translation:

The compiler uses the syntactic structure of the language to generate output

28

Derivation

Given a context-free grammar, we can determine the set of all strings (sequences of tokens) generated by the grammar using derivations:

- We begin with the start symbol
- In each step, we replace one nonterminal in the current form with one of the right-hand sides of a production for that nonterminal

30

Parse Tree

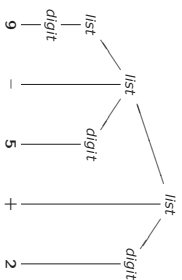
(derivation tree in F12)

- The root of the tree is labelled by the start symbol
- Each leaf of the tree is labelled by a terminal (=token) or ϵ (=empty)
- Each interior node is labelled by a nonterminal
- If node A has children X_1, X_2, \dots, X_n , then there must be a production $A \rightarrow X_1 X_2 \dots X_n$

32

Parse Tree (Example)

Parse tree of the string $9 - 5 + 2$ using grammar G'



Yield of the parse tree: the sequence of leafs (left to right)

Parsing: the process of finding a parse tree for a given string

Language: the set of strings that can be generated by some parse tree

33

Ambiguity

Consider the following context-free grammar:

$$G' = (\{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string)$$

with productions P

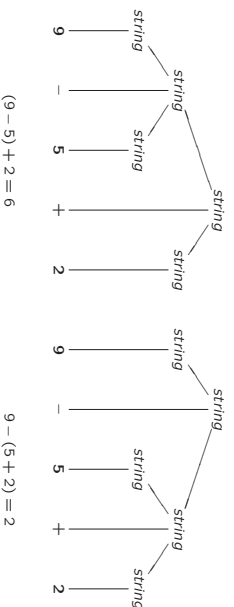
$$string \rightarrow string + string \mid string - string \mid 0 \mid 1 \mid \dots \mid 9$$

This grammar is ambiguous, because more than one parse tree generates the string $9 - 5 + 2$

34

Ambiguity (Example)

Parse trees of the string $9 - 5 + 2$ using grammar G'



35

Associativity of Operators

By convention

$$\left. \begin{aligned} 9 + 5 + 2 &= (9 + 5) + 2 \\ 9 - 5 - 2 &= (9 - 5) - 2 \end{aligned} \right\} \text{left associative}$$

In most programming languages:

$+$, $-$, $*$, $/$ are left associative

$**$, $=$ are right associative:

$$\begin{aligned} a ** b ** c &= a ** (b ** c) \\ a = b = c &= a = (b = c) \end{aligned}$$

36

Precedence of Operators

Consider: $9 + 5 * 2$

Is this $(9 + 5) * 2$ or $9 + (5 * 2)$?

Associativity does not resolve this

Precedence of operators: $\begin{matrix} + - \\ * / \end{matrix} \mid \begin{matrix} \text{increasing} \\ \text{precedence} \end{matrix}$

A grammar for arithmetic expressions:

$$\begin{aligned} expr &\rightarrow expr + term \mid expr - term \mid term \\ term &\rightarrow term * factor \mid term / factor \mid factor \\ factor &\rightarrow digit \mid (expr) \\ digit &\rightarrow 0 \mid 1 \mid \dots \mid 9 \end{aligned}$$

Parse tree for $9 + 5 * 2$. . .

37

Syntax-Directed Translation

Using the syntactic structure of the language to generate output corresponding to some input

Two techniques:

- Syntax-directed definition
- Translation scheme

Example: translation of infix notation to postfix notation

Infix	Postfix
$(9 - 5) + 2$	$95 - 2 +$
$9 - (5 + 2)$	$952 + -$

What is $952 + -3*$?

38

Syntax-Directed Translation

Using the syntactic structure of the language to generate output corresponding to some input

Two variants:

- Syntax-directed definition
- Translation scheme

Example: translation of infix notation to postfix notation

Simple infix expressions generated by

$$\begin{aligned} expr &\rightarrow expr_1 + term \mid expr_1 - term \mid term \\ term &\rightarrow 0 \mid 1 \mid \dots \mid 9 \end{aligned}$$

39

Syntax-Directed Definition

- Uses a context-free grammar to specify the syntactic structure of the language

- Associates a set of attributes with (non)terminals

- Associates with each production a set of semantic rules for computing values for the attributes

- The attributes contain the translated form of the input after the computations are completed (in example: postfix notation corresponding to subtree)

40

Syntax-Directed Definition (Example)

Production	Semantic rule
$expr \rightarrow expr_1 + term$	$expr.t = expr_1.t \parallel term.t \parallel '+'$
$expr \rightarrow expr_1 - term$	$expr.t = expr_1.t \parallel term.t \parallel '-'$
$expr \rightarrow term$	$expr.t = term.t$
$term \rightarrow 0$	$term.t = '0'$
$term \rightarrow 1$	$term.t = '1'$
...	...
$term \rightarrow 9$	$term.t = '9'$

Result: annotated parse tree

Example: $9 - 5 + 2$

41

Synthesized and Inherited Attributes

An attribute is said to be ...

- **synthesized** if its value at a parse tree node N is determined from attribute values at the children of N (and at N itself)
- **inherited** if its value at a parse tree node N is determined from attribute values at the parent of N (and at N itself and its siblings)

We consider synthesized attributes

42

Depth-First Traversal

- A syntax-directed definition does not impose an evaluation order of the attributes on a parse tree
- Different orders might be suitable
- **Tree traversal**: a specific order to visit the nodes of a tree (always starting from the root node)
- **Depth-first traversal**
 - Start from root
 - Recursively visit children (in any order)
- Hence, visit nodes far away from the root as quickly as it can (DF)

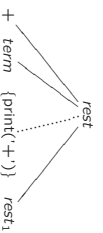
43

Translation Scheme

A translation scheme is a context-free grammar with semantic actions embedded in the bodies of the productions

Example

$rest \rightarrow + term \{ \text{print}('+') \} rest_1$



45

Parsing

- Process of determining if a string of tokens can be generated by a grammar
- For any context-free grammar, there is a parser that takes at most $O(n^3)$ time to parse a string of n tokens
- Linear algorithms sufficient for parsing programming languages
- Two methods of parsing:
 - **Top-down** constructs parse tree from root to leaves
 - **Bottom-up** constructs parse tree from leaves to root

Cf. top-down PDA and bottom-up PDA in FI2

47

A Possible DF Traversal

Postorder traversal

```

procedure visit (node N)
{
  for (each child C of N, from left to right)
  { visit (C);
  }
  evaluate semantic rules at node N;
}
  
```

Can be used to determine synthesized attributes / annotated parse tree

44

Translation Scheme (Example)

```

expr -> expr1 + term {print('+')}
expr -> expr1 - term {print('-')}
expr -> term
term -> 0 {print('0')}
term -> 1 {print('1')}
...
term -> 9 {print('9')}
  
```

Example: parse tree for $9 - 5 + 2$

Implementation requires postorder traversal

46

Parsing (Top-Down Example)

```

stmt -> expr ;
      | if (expr) stmt
      | for (optexpr ; optexpr ; optexpr) stmt
      | other
optexpr -> ε
        | expr
  
```

How to determine parse tree for

for (:expr ; expr)other

Use lookahead: current terminal in input

48

Predictive Parsing

- Recursive-descent parsing is a top-down parsing method:
 - Executes a set of recursive procedures to process the input
 - Every nonterminal has one (recursive) procedure parsing the nonterminal's syntactic category of input tokens
- Predictive parsing is a special form of recursive-descent parsing:
 - The lookahead symbol unambiguously determines the production for each nonterminal

49

Predictive Parsing (Example)

```
void stmt()
{
    switch (lookahead)
    {
        case expr:
            match(expr); match('('); break;
        case if:
            match(if); match('('); match(expr); match(')'); stmt();
            break;
        case for:
            match(for); match('(');
                stmt(); match('('); match(expr); match(')'); optexpr();
            match(')'); stmt(); break;
        case other:
            match(other); break;
        default:
            report("syntax error");
    }
}

void match(terminal t)
{
    if (lookahead==t) lookahead = nextTerminal;
    else report("syntax error");
}
```

50

Using FIRST

- Let α be string of grammar symbols
- $FIRST(\alpha)$ is the set of terminals that appear as first symbols of strings generated from α

Simple example:

```
stmt → expr ;
      | if (expr ) stmt
      | for (optexpr ; optexpr ) stmt
      | other
```

Right-hand side may start with nonterminal...

51

Using FIRST

- Let α be string of grammar symbols
- $FIRST(\alpha)$ is the set of terminals that appear as first symbols of strings generated from α
- When a nonterminal has multiple productions, e.g.,
 $A \rightarrow \alpha \mid \beta$
then $FIRST(\alpha)$ and $FIRST(\beta)$ must be disjoint in order for predictive parsing to work

52

Compiler constructie

college 1
Overview

Chapters for reading: 1.1, 1.2, 2.1-2.5

53