


Formal
Methods
& Tools



University of Twente
The Netherlands


Introduction

Vertalerbouw HC1

VB HC1

Original design: Theo Ruys
University of Twente
Department of Computer Science
Formal Methods & Tools

Arend Rensink
kamer: Zilverling 5090
telefoon: 4862
email: rensink@cs.utwente.nl



Theo Ruys


Vertalerbouw - kernpunten

- ASTs staan centraal: geen one-pass compilatie
- nadruk op LL(k) compilatie (recursive descent)
- modern en populair practicumgereedschap ([ANTLR](#))
- Java als implementatietaal
- OO-achtige aanpak van het bouwen van een vertaler
- aandacht voor moderne taalaspecten
- minder aandacht voor theorie achter scanning en parsing
- eindopdracht: bouwen van eigen vertaler

VB HC 1

Ch. 1 - Introduction

2



Theo Ruys

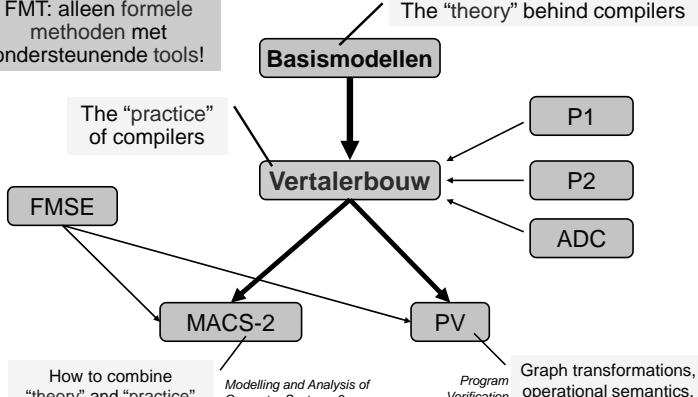
VB binnen INF/CS Curriculum (2)

FMT: alleen formele methoden met ondersteunende tools!

The “theory” behind compilers

The “practice” of compilers

How to combine “theory” and “practice” to build model checkers.




```

graph TD
    Basismodellen --> Vertalerbouw
    P1 --> Vertalerbouw
    P2 --> Vertalerbouw
    ADC --> Vertalerbouw
    Vertalerbouw --> MACS-2
    Vertalerbouw --> PV
    FMSE --> MACS-2
    GraphTransformations[Graph transformations, operational semantics, software verification.] --> PV
  
```

VB HC 1

Ch. 1 - Introduction

4



Theo Ruys

Organisatie (1)

- Hoorcolleges: 9x (ma 3+4 en di 6+7)
 - 7x theorie uit [Watt & Brown 2000]
 - 2x ANTLR (tool voor practicum)
- Practicum: 3 groepen
 - clusters: Zi 4054 A, B en C
 - indeling bij eerste practicum op woensdag 25 april 2012
 - verplicht: aanwezigheid wordt gecontroleerd
 - iedere groep heeft eigen studentassistent
- Onderwijsmateriaal:
 - [Watt & Brown 2000] en Triangle compiler bevatten fouten. Zie de website voor errata en bugfixes.
 - practicumhandleiding (via website)

VB HC 1

Ch. 1 - Introduction

5

© Theo Ruys

Organisatie (2)

- Beoordeling
 - twee opgavenseries (individueel) OS
 - eindopdracht (in tweetallen) P
 - eindcijfer = $(OS + P) / 2$
mits OS en P beiden ≥ 5.0 , anders 4
- Opgavenseries
 - eerste serie komt (uiterlijk) ma 02 mei 2011 beschikbaar
 - stricte deadlines
overschrijden van de deadline kost punten.
 - worden nagekeken door eigen studentassistent
 - strict individueel: fraude levert 1.0 voor het vak en uitsluiting voor dit studiejaar

Cijfers voor OS en P blijven dit jaar nog staan. Aparte delen van OS echter niet.

VB HC 1 Ch. 1 - Introduction 6

© Theo Ruys

Organisatie (3)

- Practicum deel 1 - introductie
 - 5 weken
 - oefenen met stof van hoorcollege
 - oefenen met ANTLR 3
 - succesvolle afronding noodzakelijk voor deel 2
 - advies: voorbereiden van de practica (met name wk 1 en 2)
- Practicum deel 2 - zelf een vertaler bouwen
 - 3 weken (+ 2 tentamenweken + 1 uitloopweek)
 - gebruikmaken van ANTLR
 - verschillende moeilijkheidsgraad
 - (maximum) cijfer hangt af van gekozen opdracht
 - deadline: woensdag 11 juli 2012 (= woensdag na tentamens)

VB HC 1 Ch. 1 - Introduction 7

© Theo Ruys

Overview of Lecture 1

- Ch 1 – Introduction
 - Levels of programming language
 - Programming Language Processors
 - Specification of Programming Languages
 - The Triangle Programming Language
- Ch 2 – Language Processors
 - Translators and Compilers
 - Interpreters
 - Real and abstract machines
 - Interpretive compilers
 - Portable compilers
 - Bootstrapping
 - Triangle language processors
- Organisatie van Vertalerbouw (in Dutch)

VB HC 1 Ch. 1 - Introduction 8

© Theo Ruys

Tijdsbesteding

5 ECTS = 140 uur

hoorcolleges: contacturen	9 x 2 =	18
zelfstudie naast hoorcolleges	7 x 3 =	21
practicum deel 1: contacturen	5 x 4 =	20
voorbereidingen practicum blok 1	5 x 2 =	10
opgavenseries	2 x 8 =	16
practicum deel 2: contacturen	3 x 6 =	18
eindopdracht (buiten practicumuren)	5 x 5 =	25
verslag		12
TOTAAL		140

VB HC 1 Ch. 1 - Introduction 9

© Theo Ruys

Ch 1 - Introduction

- 1.1 Levels of programming language
- 1.2 Programming Language Processors
- 1.3 Specification of Programming Languages
- 1.4 The Triangle Programming Language

VB HC 1 Ch. 1 - Introduction 10

© Theo Ruys

Why Compilers?

- Programming languages (Pascal, C, Java, Python,...)
 - “easy” to use and write by humans
 - “difficult” to be interpreted by computers
- Machine language (microcode/assembler)
 - “difficult” to use and write by humans
 - “easy” to execute by hardware/microprocessor
- Not just for programming languages
 - analysis/translation of natural language
 - analysis of generated data

Compilers: close interplay between theory and practice.

VB HC 1 Ch. 1 - Introduction 11

© Theo Ruys

Levels of Programming Languages (1)

- Programming language = notation for algorithms (which might be run on computers).
- Levels of programming language:
 - High-level (abstract)
Java, Pascal, Miranda


```
let
  var n: Integer
in
  n := n-1
```
 - Low-level (detailed)
assembler program


```
LOAD r1, n
LOAD r2, 1
SUB r1, r2
STORE r1, n
```
 - Machine code


```
00010001001001101
01000010010100011
11100011111...
```

VB HC 1 Ch. 1 - Introduction 12

© Theo Ruys

Levels of Programming Languages (2)

- Aspects typically found in high-level languages, but not in low-level languages:
 - expressions
 - control structures:
 - if-then-else, while, procedures
 - data types:
 - different types of data: boolean, integer, char, float
 - composite data types: arrays
 - user defined data types: records, classes
 - declarations:
 - type checking
 - **abstraction**
 - **encapsulation**

VB HC 1 Ch. 1 - Introduction 13

© Theo Ruys

Language Specification

Language of the end project has to be defined like Triangle!

- syntax (or *grammar*)
 - defines the structure of correct sentences
- contextual constraints (or *static semantics*)
 - well-formedness depends on context of an expression (e.g. scope rules, type rules)
- semantics
 - defines the meaning of correct sentences (denotational or operational semantics)
- [Watt & Brown 2000]
 - formal syntax using (E)BNF
 - informal contextual constraints
 - informal semantics

See Appendix B for a "complete" specification of Triangle.

VB HC 1 Ch. 1 - Introduction 15

© Theo Ruys

Syntax (1)

- Syntax is specified using "Context Free Grammars" (CFG, known from "Basismodellen").
A CFG G is defined by a 4-tuple (N, T, P, S)
 - N : a finite set of non-terminal symbols
 - T : a finite set of terminal symbols
 - P : a finite set of production rules
 - S : a start symbol
- CFGs are usually written using Backus-Naur Form (BNF) notation.
 - Two types of production rules
 - $N ::= \alpha$
 - $N ::= \alpha \mid \beta \mid \dots$

VB HC 1 Ch. 1 - Introduction 16

© Theo Ruys

Syntax (2)

- A CFG defines a set of strings, which is called the language of the CFG.
- Example:

id	::=	letter	
		id letter	non terminal
		id digit	
letter	::=	a b c ... z	
digit	::=	0 1 2 ... 9	terminals
- This grammar generates "identifiers": finite strings that start with a letter and are followed by zero or more letters or digits.
 - "a", "foo", "x123141243124124", "a1b2c3d4e5f6"

VB HC 1 Ch. 1 - Introduction 17

© Theo Ruys

(Mini) Triangle

- Triangle is a small, but realistic, Pascal-like language with **let-in** constructs for local declarations.
- Example:


```

let
  const MAX ~ 10;
  var n: Integer
in begin
  getint(var n);
  if (n>0) /\ (n<=MAX) then
    while n > 0 do begin
      putint(n); puteol();
      n := n-1
    end
  else
  end
end
      
```

Annotations for the example:

 - local declarations (pointing to the `let` block)
 - constant def ("~" instead "=") (pointing to `const MAX ~ 10;`)
 - variable may be changed by `getint` (pointing to `getint(var n);`)
 - sequence of commands can be grouped with `begin/end` (pointing to the `begin` block)
 - else is mandatory (but might be empty) (pointing to the `else` block)

VB HC 1 Ch. 1 - Introduction 18

© Theo Ruys

(Mini) Triangle – Syntax (1)

Program	::=	single-Command	
Command	::=	single-Command	
		Command ; single-Command	
single-Command	::=		<i>skip</i>
		V-name := Expression	
		Identifier (Expression)	
		if Expression	
		then single-Command	
		else single-Command	
		while Expression do single-Command	
		let Declaration in single-Command	
		begin Command end	
Expression	::=	primary-Expression	
		Expression Operator primary-Expression	
primary-Expression	::=	Integer-Literal	
		V-name	
		Operator primary-Expression	
		(Expression)	

VB HC 1 Ch. 1 - Introduction 19

© Theo Ruys

(Mini) Triangle – Syntax (2)

V-name	::=	Identifier
Declaration	::=	single-Declaration
		Declaration ; single-Declaration
single-Declaration	::=	const Identifier ~ Expression
		var Identifier : Type-denoter
Type-denoter	::=	Identifier
Operator	::=	+ - * / < > = \
Identifier	::=	Letter
		Identifier Letter
		Identifier Digit
Integer-Literal	::=	Digit
		Integer-Literal Digit
Comment	::=	! Graphic* eo1

VB HC 1 Ch. 1 - Introduction EBNF for: zero-or-more 20

© Theo Ruys

Syntax Trees (1)

- A Context-Free Grammar (CFG) is a specification of a rewrite system.
 - A CFG generates a language, which is the set of all strings of terminal symbols derivable from the start symbol.
 - Such a language is defined in terms of syntax trees and phrases (sentences).
- A syntax tree of a grammar G is an ordered labeled tree:
 - the leaves are labeled by terminal symbols
 - the interior nodes are labeled by nonterminal symbols
 - each nonterminal node labeled by N has children labeled by X_1, \dots, X_n , such that $N ::= X_1, \dots, X_n$ is a production rule.
- A phrase of G is a string of terminal symbols (taken from left to right) of a syntax tree of G.

VB HC 1 Ch. 1 - Introduction 21

© Theo Ruys

Syntax Trees (2)

Expr	::=	prim-Expr
		Expr Op prim-Expr
prim-Expr	::=	V-name
		Int-Literal
		...
V-name	::=	Identifier

- Example: $d + 10 * n$

Triangle's binary operators all have the same precedence.

VB HC 1 Ch. 1 - Introduction 22

Concrete vs. Abstract Syntax

- The defining grammar of a programming language specifies the concrete syntax of a language.
 - The concrete syntax is important for the programmer who needs to know exactly how to write syntactically well-formed programs.
 - But, the concrete syntax has no influence on the semantics of the programs.
- The abstract syntax omits irrelevant syntactic details and only specifies the essential structure of programs.
 - E.g. different concrete syntax for an assignment:

```

v := e
v <- e
assign e to v
set v=e
    
```

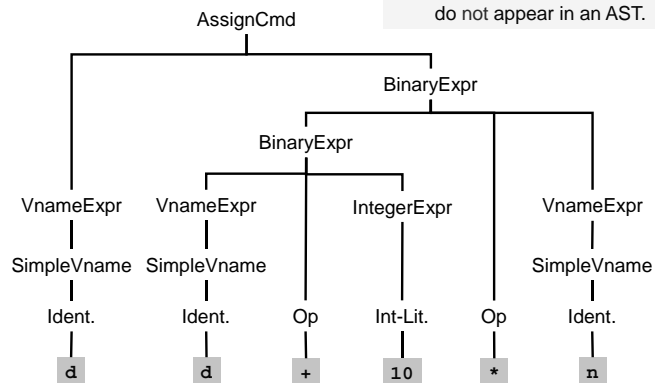
(Mini) Triangle – Abstract Syntax

Program	::=	Command	Program
Command	::=	Command ; Command	SequentialCmd
		V-name := Expression	AssignCmd
		Identifier (Expression)	CallCmd
		if Expression	IfCmd
		then Command	
		else Command	
		while Expression do Command	WhileCmd
		let Declaration in Command	LetCmd
Expression	::=	Integer-Literal	IntegerExpr
		V-name	VnameExpr
		Operator Expression	UnaryExpr
		Expression Operator Expression	BinaryExpr
V-name	::=	Identifier	SimpleVname
Declaration	::=	Declaration ; Declaration	SeqDecl
		const Identifier ~ Expression	ConstDecl
		var Identifier : Type-denoter	VarDecl
Type-denoter	::=	Identifier	SimpleTypeDen

Abstract Syntax Tree (1)

- Example: `d := d + 10 * n`

Structuring terminals (like `:=`, `if`, `while`, `let`, `begin`, `end`, etc.) do not appear in an AST.



Abstract Syntax Tree (2)

- In an AST, each node is labelled by a production rule. Consequently, an AST represents the phrase-structure of the program explicitly.
- The AST is a convenient structure for specifying
 - contextual constraints
 - semantics
 - code generation
- ASTs will be used extensively in [Watt & Brown 2000] and in the course of Vertalerbouw.
 - ANTLR, the tool used in the laboratory, works quite elegantly with ASTs.

© Theo Ruys

Context Constraints (1)

- Apart from the syntax rules, there may be rules which specify that a phrase is well-formed which depend on the context of the phrase: context constraints.
 - scope rules
 - Scope rules are concerned with the visibility of identifiers.
 - Every identifier which is used should first be declared.
 - In other words: every applied occurrence of an identifier is related to a binding occurrence of the same identifier.
 - type rules
 - Each value has a type.
 - Each operation has a type rule, which specifies the types of the operands and the type of the result of the operation.

VB HC 1 Ch. 1 - Introduction 27

© Theo Ruys

Context Constraints (2)

- Example: scope rule


```
let
  const m ~ 2;
  var n: Integer
in begin
  n := m*2;
end
```

declaration of n:
binding occurrence

use of n: applied occurrence

??

```
let
  var n: Integer
in begin
  ...
  n := m*2;
end
```

applied occurrence

If there is no enclosing scope (i.e. letCmd) where m is declared, the binding occurrence of m is missing: scope error!

VB HC 1 Ch. 1 - Introduction 28

© Theo Ruys

Context Constraints (3)

- Example: type rules


```
let
  var n: Integer
in begin
  ...
  while n > 0 do
    n := n-1;
  ...
end
```

Type rule for $E_1 > E_2$ (GreaterOp):
If E_1 and E_2 are both of type int, then the result is of type bool.

Type rule for **while** E do C (WhileCmd):
E must be a boolean.

Type rule for $V := E$ (AssignCmd):
The types of V and E must be equivalent.

Type rule for $E_1 - E_2$ (SubOp):
If E_1 and E_2 are both of type int, then the result is of type int.

VB HC 1 Ch. 1 - Introduction 29

© Theo Ruys

Semantics (1)

- Semantics is concerned with the meaning of programs, i.e., their behaviour when executed.
- Terminology:
 - Commands are executed and perform side effects.
 - Side effects:
 - changing the values of variables
 - perform I/O
 - Declarations are elaborated to produce bindings.
 - Expressions are evaluated and yield values (and may or may not perform side effects).
- The semantics of each specific form of command, expression, declaration, etc. has to be specified.

VB HC 1 Ch. 1 - Introduction 30

© Theo Ruys

(Mini) Triangle – Semantics (1)

- AssignCmd: $V := E$
 - The expression E is evaluated to yield a value v .
 - Then v is assigned to the variable name V .
- SequentialCmd: $C_1; C_2$
 - First C_1 is executed.
 - Then C_2 is executed.
- WhileCmd: **while** E **do** C
 - The expression E is evaluated to yield a truth-value t .
 - If t is true, C is executed, and then the WhileCmd is executed again.
 - If t is false, execution of the WhileCmd is completed.

VB HC 1 Ch. 1 - Introduction 31

© Theo Ruys

(Mini) Triangle – Semantics (2)

- Expression:
 - IntegerExpr: Integer-Literal
 - The expression yields the value of the Integer-Literal.
 - VnameExpr: V-name
 - The expression yields the value of the variable of V-name.
 - UnaryExpr: $op\ E$
 - The expression yields the value obtained by applying unary operator op on the value yielded by the expression E .
 - BinaryExpr: $E_1\ op\ E_2$
 - The expression yields the value obtained by applying binary operator op to the values yielded by the expression $E_1\ op\ E_2$.

Expressions in (Mini)Triangle do not have side effects.

VB HC 1 Ch. 1 - Introduction 32

© Theo Ruys

Triangle

We will practice with Triangle in the laboratory session of week 1.

- Triangle is a small, but realistic, Pascal-like language with **let-in** constructs for local declarations.
 - Triangle supports (user-defined) arrays, records, procedures and functions.
 - Triangle supports value- and reference parameter passing for procedures. Furthermore, procedures and functions can be passed to procedures.
 - Triangle is type complete: no operations are arbitrarily restricted in the types of the language.
 - Triangle expressions are free of side effects.
- See Appendix B for the (informal) specification of Triangle.
- See [Gosling et. al. 1996] for the language definition of Java.

VB HC 1 Ch. 1 - Introduction 33

© Theo Ruys

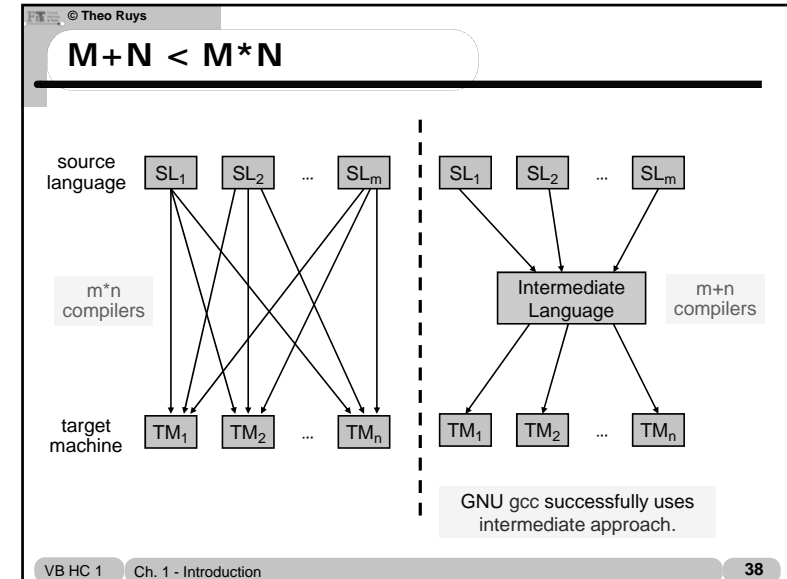
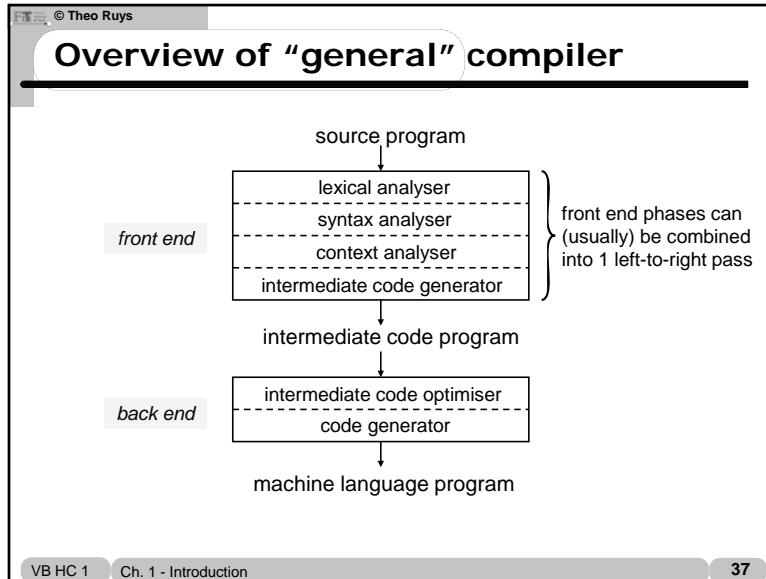
"General" compiler scheme

```

graph TD
    A[source program] --> B[Translator 1]
    B --> C[intermediate language program]
    C --> D[Translator 2]
    D --> E[machine language program]
  
```

The flowchart illustrates the general compiler scheme. It starts with a 'source program' (labeled *machine M_2*). This is processed by 'Translator 1' (labeled *front end*) to produce an 'intermediate language program' (labeled *machine M_1*). This intermediate program is then processed by 'Translator 2' (labeled *back end*) to produce the final 'machine language program' (labeled *machine M_0*).

VB HC 1 Ch. 1 - Introduction 36



© Theo Ruys

Ch 2 – Language Processors

- 2.1 Translators and Compilers
- 2.2 Interpreters
- 2.3 Real and abstract machines
- 2.4 Interpretive compilers
- 2.5 Portable compilers
- 2.6 Bootstrapping
- 2.7 Triangle language processors

VB HC 1 Ch. 1 - Introduction 42

- © Theo Ruys
- ## Translators (1)
- A translator accepts a text expressed in a source language and generates semantically-equivalent text expressed in a target language.
 - Some translators
 - C → x86 assembly
 - x86 assembly → x86 binary code
 - Java → JVM byte code
 - Java → C
 - JVM byte code → x86 assembly
 - JVM byte code → Java (java disassembler, e.g. jad)
 - Dutch → English
 - Natural-language translation/processing is AI-related (see HMI courses)
- VB HC 1 Ch. 1 - Introduction 43

© Theo Ruys

Translators (2)

- Terminology:
 - compiler: translates a high-level language into a low-level language.
several target instructions per source instruction
 - assembler: translates from an assembly language into machine code.
one machine instruction per source instruction
 - source program: source language (input) text.
 - object program: target language (output) text.
 - implementation language: programming language in which the translator itself is written.

VB HC 1 Ch. 1 - Introduction 44

© Theo Ruys

Tombstone diagrams (1)

- Tombstone diagrams
 - Set of “puzzle pieces” to reason about language processors and programs.
A complete diagram of a translator specifies how the source, target and implementation languages and the underlying machine are related.
 - four different kinds of pieces
 - combination rules to combine the pieces
 - not all pieces fit together

VB HC 1 Ch. 1 - Introduction 45

© Theo Ruys

Tombstone diagrams (2)

Program P expressed in language L.

Translator implemented in L, which translates programs from source language S to target language T.

Interpreter for language M, implemented in language L.

Machine M.

VB HC 1 Ch. 1 - Introduction 46

© Theo Ruys

Tombstone diagrams (3)

Examples:

WordCount program written in Java.

A compiler from Java to JVM byte code, written in Java.

Interpreter for JVM byte code, written in C.

The x86 processor family.

VB HC 1 Ch. 1 - Introduction 47

© Theo Ruys

Tombstone diagrams (4)

- Combination rule – like “domino”:
 - When combining pieces, the sides that touch each other should use the same implementation language.

VB HC 1 Ch. 1 - Introduction 48

© Theo Ruys

Tombstone diagrams (5)

- Compilation of a program:
 - When compiling a program P in source language S to a target language T , a new “tombstone piece” is obtained: a program P in language T .

To specify that the program is generated, we use a coloured background for these pieces.

This will only work, if we have a machine on which we can (perhaps indirectly) run programs written in L .

VB HC 1 Ch. 1 - Introduction 49

© Theo Ruys

Cross-compilation

- A cross-compiler runs on one machine (host machine) but generates code for a dissimilar machine (target machine).
 - Useful if the target machine
 - does not have enough memory to compile programs.
 - does not have tools to develop programs.
 - Examples: programs for PDAs, telephones, media players

download

VB HC 1 Ch. 1 - Introduction 50

© Theo Ruys

Two-stage compilation

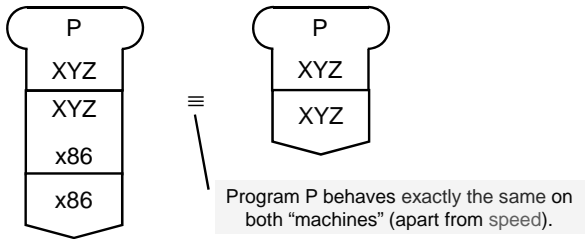
- A two-stage translator is a composition of two translators.
 - With compilers $S \rightarrow T$ and a $T \rightarrow U$, the source program in S is translated to target language U via T .
 - Easily generalized to an n -stage translator.

Compilation of a high-level programming language is usually an n -stage translation, as at least one intermediate language is used for the compilation to executable code.

VB HC 1 Ch. 1 - Introduction 51

© Theo Ruys

Real/Abstract Machines



- An interpreter is often called an abstract machine as opposed to its hardware counterpart, which is a real machine.
 - A well-known example of an abstract machine is the Java Virtual Machine (JVM).

VB HC 1 Ch. 1 - Introduction 56

© Theo Ruys

Interpretive Compilers

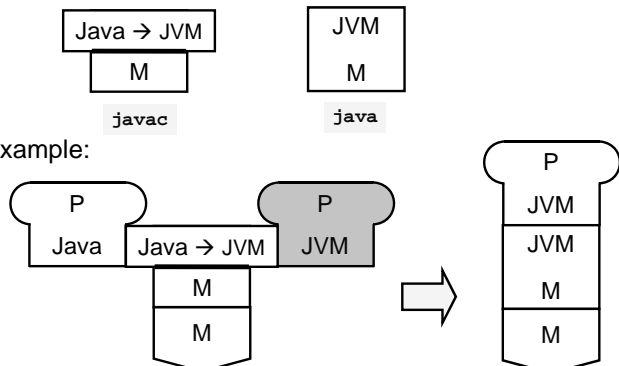
- Tradeoffs in “executing a program”:
 - compiler: long time to compile, but fast execution
 - interpreter: starts running immediately, but will be slow
- An interpretive compiler is a combination of a compiler and an interpreter. The key idea is to translate the source program into an intermediate language (IL).
 - the IL is in level between the source language and the ordinary machine code.
 - the instructions of the IL have simple formats and can therefore be analyzed easily and quickly
 - translation from the source language to IL is easy and fast.

VB HC 1 Ch. 1 - Introduction 57

© Theo Ruys

Java Development Kit (JDK)

- Sun’s JDK provides an implementation of an interpretive compiler for Java. Central to the JDK is the JVM.



- Example:

VB HC 1 Ch. 1 - Introduction 58

© Theo Ruys

Portable Compilers

- A program is portable to the extent that it can be (compiled and) run on any machine, without change.
 - Portability is measured in the proportion of code that remains unchanged when moving to a different machine.
 - Application programs in high-level languages should achieve a 95-99% portability.
 - In general, the portability for language processors is much lower, though (about 50%), because a compiler’s function is to generate machine code for a particular machine.
 - Unless you are able to parameterize the language processor in (a description of) the machine.
 - A compiler that generates intermediate code is potentially much more portable, though.

VB HC 1 Ch. 1 - Introduction 59

© Theo Ruys

Portable Compiler Kit for Java (1)

- Portable JDK:
- Solution: rewrite the interpreter for JVM for the target machine.

VB HC 1 Ch. 1 - Introduction 60

© Theo Ruys

Portable Compiler Kit for Java (2)

- Now we are able to compile Java programs using this JVM interpreter:

VB HC 1 Ch. 1 - Introduction 61

© Theo Ruys

What have you seen today?

- Language definition consists of
 - Syntax
 - Context constraints
 - Semantics
- Syntax definition
 - Based on Context-Free Grammars
- Compilation versus interpretation
 - Many different scenarios
 - Handy notation: Tombstone diagrams

VB HC 1 Ch. 1 - Introduction 62