

Compiler Construction

Dr. Ir. Bart Kienhuis
Computer Systems Group
LIACS

1

Why This Course

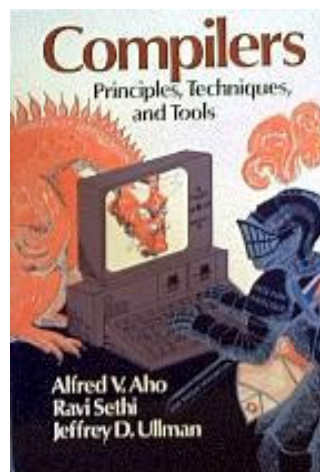
- ⌘ Know how to build a compiler for a (simplified) (programming) language
- ⌘ Know how to use compiler construction tools, such as generators for scanners and parsers
- ⌘ Be able to write LL(1), LR(1) grammars (for new languages)
- ⌘ Be familiar with compiler analysis and optimization techniques
- ⌘ ... learn how to work on a larger software project!

2

Course Outline

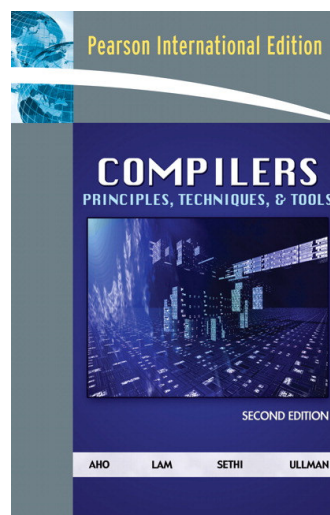
- ⌘ In class, we discuss the theory using the 'dragon' book by Aho et al.
- ⌘ In the practicum, the theory is applied when building a compiler that converts Pascal code to MIPS instructions.

A.V. Aho, R. Sethi, en J.D. Ullman, Compilers: Principles, Techniques, and Tools, Addison-Wesley, 1986, ISBN: 0-201-10088-6.



New edition

- ⌘ Dragon book has been revised in 2006
- ⌘ In Second edition good improvements are made
- ⌘ **Publisher:** Addison Wesley; 2 edition (August 31, 2006)
- ⌘ **Language:** English
- ⌘ **ISBN-10:** 0321486811



Course Outline

⌘ Contact hours

- ☒ Official communication medium is email
- ☒ Blackboard (<http://blackboard.leidenuniv.nl>)
- ☒ All material needed is available here

⌘ Practicum

- ☒ Different from previous years, we now offer 5 self contained assignments
- ☒ These assignments are done by groups of two persons
- ☒ Assignments are handed in into CVS

5

Course Outline

⌘ Grading

- ☒ 2 ECTS Written Exam
- ☒ 5 ECTS Practicum

⌘ You need to pass all 5 assignments

⌘ No 'late' submissions will be accepted!

- ☒ If you miss these assignments, you have to wait until next year

6

Course Outline (Tentative)

- ⌘ 04/09/08 Introduction
- ⌘ 11/09/08 Lexical and Syntax Analysis
- ⌘ 18/09/08 Syntax Analysis assignment 1
- ⌘ 25/09/08 NO CLASS
- ⌘ 02/10/08 Type Checking assignment 2
- ⌘ 09/10/08 Intermediate Code Generation 1
- ⌘ 16/10/08 NO CLASS
- ⌘ 23/10/08 Intermediate Code Generation 2 assignment 3
- ⌘ 30/10/08 Code Generation 1
- ⌘ 06/11/08 Code Generation 2 assignment 4
- ⌘ 13/11/08 Run-Time Organization
- ⌘ 20/11/08 Code Optimizations assignment 5
- ⌘ 27/11/08 ELECTIVE
- ⌘ 04/12/08 backup date

7

Practicum

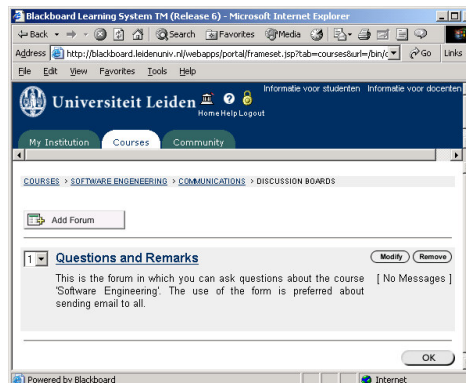
- ⌘ 28/09-02/10 assignment 1, Calculator
- ⌘ 02/10-23/10 assignment 2, Parsing & Syntax tree
- ⌘ 23/10-06/11 assignment 3, Intermediate code
- ⌘ 06/11-20/11 assignment 4, Assembly generation
- ⌘ 20/11-04/12 assignment 5, Optimizations

- ⌘ All deadlines are at 17.00h (5 pm).
- ⌘ The deadlines are strict.
- ⌘ Submission takes place in CVS

8

Blackboard Coco

- ⌘ Please enroll on-line to Coco in Blackboard
- ⌘ Communication about Coco is shared between everyone.
- ⌘ Use the 'Forum' option to ask me questions.
- ⌘ If you ask me directly, I will submit also to the forum.



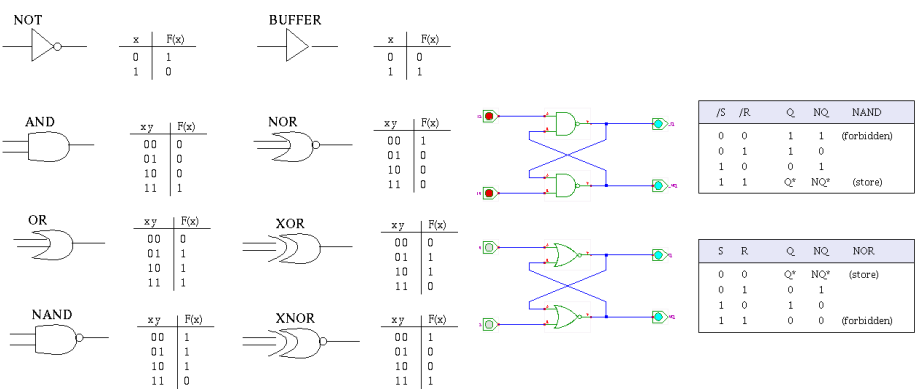
9

Introduction

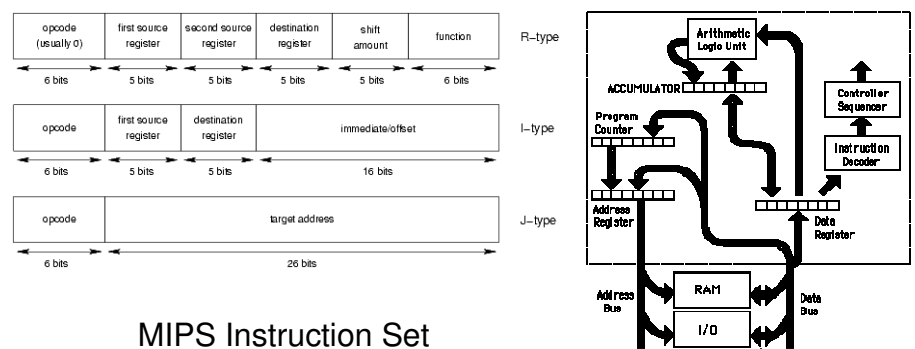
- ⌘ Compiler Construction
- ⌘ Missing Link between
 - ⊞ Digital Technique
 - ⊞ Boolean Logic
 - ⊞ Flip-Flops
 - ⊞ Computer Architectures
 - ⊞ Memory
 - ⊞ Instructions

10

Digital Technique



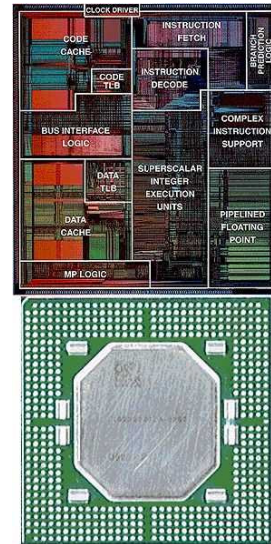
Computer Architecture



MIPS Instruction Set

opcode field opcode instruction format

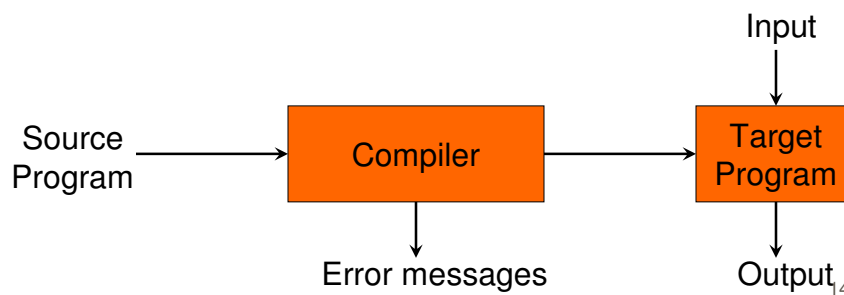
000010	j	J-type
000011	jal	J-type
000100	beq	I-type
000101	bne	I-type
001000	Addi	I-type
001001	Addiu	I-type
001010	Slti	I-type
001011	Sltiu	I-type
001100	Andi	I-type



Compilers and Interpreters

⌘ "Compilation"

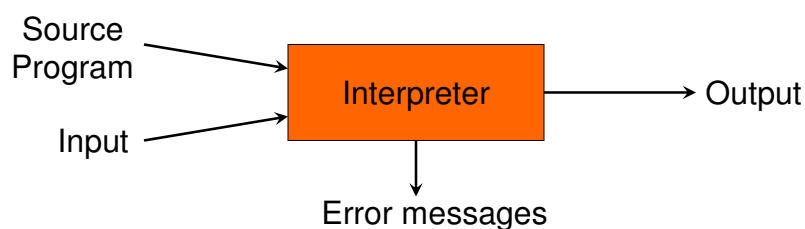
☒ Translation of a program written in a source language into a semantically equivalent program written in a target language



Compilers and Interpreters (cont'd)

⌘ "Interpretation"

- ☒ Performing the operations implied by the source program



15

The Analysis-Synthesis Model of Compilation

⌘ There are two parts to compilation:

- ☒ *Analysis* determines the operations implied by the source program which are recorded in a tree structure
- ☒ *Synthesis* takes the tree structure and translates the operations therein into the target program

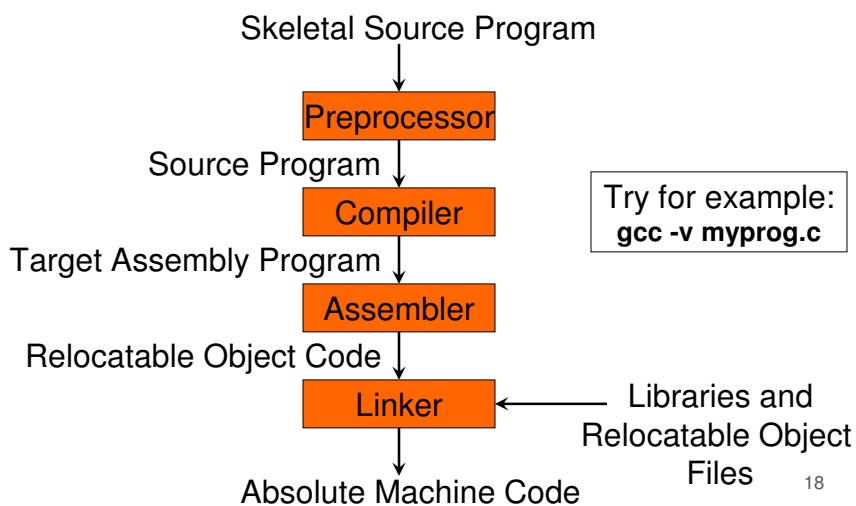
16

Other Tools that Use the Analysis-Synthesis Model

- ⌘ *Editors* (syntax highlighting)
- ⌘ *Pretty printers* (e.g. doxygen)
- ⌘ *Static checkers* (e.g. lint and splint)
- ⌘ *Interpreters*
- ⌘ *Text formatters* (e.g. TeX and LaTeX)
- ⌘ *Silicon compilers* (e.g. VHDL)
- ⌘ *Query interpreters/compilers* (Databases)

17

Preprocessors, Compilers, Assemblers, and Linkers



The Phases of a Compiler

Phase	Output	Sample
<i>Programmer</i>	Source string	A=B+C;
<i>Scanner</i> (performs <i>lexical analysis</i>)	Token string	'A', '=', 'B', '+', 'C', '\n'; And <i>symbol table</i> for
<i>Parser</i> (performs <i>syntax analysis</i> based on the grammar of the programming language)	Parse tree or abstract syntax tree	identifiers = / \ A + / \ B C
<i>Semantic analyzer</i> (type checking, etc)	Parse tree or abstract syntax tree	
<i>Intermediate code generator</i>	Three-address code, quads, or RTL	int2fp B t1 + t1 C t2 := t2 A
<i>Optimizer</i>	Three-address code, quads, or RTL	int2fp B t1 + t1 #2.3 A
<i>Code generator</i>	Assembly code	MOVF #2.3, r1 ADDF2 r1, r2 MOVF r2, A
<i>Peephole optimizer</i>	Assembly code	ADDF2 #2.3, r2 MOVF r2, A

The Grouping of Phases

⌘ Compiler front and back ends:

- ☒ Analysis (*machine independent* front end)
- ☒ Synthesis (*machine dependent* back end)

⌘ Passes

- ☒ A collection of phases may be repeated only once (*single pass*) or multiple times (*multi pass*)
- ☒ Single pass: usually requires everything to be defined before being used in source program
- ☒ Multi pass: compiler may have to keep entire program representation in memory

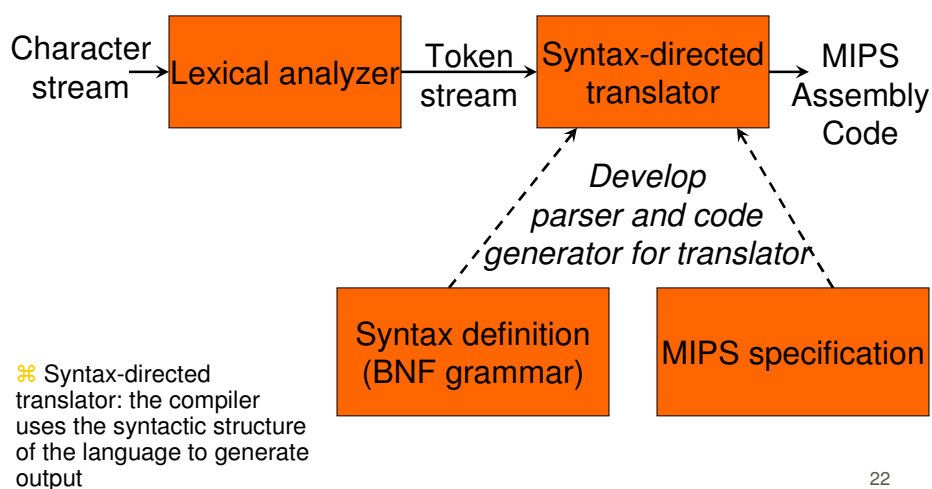
Compiler-Construction Tools

⌘ 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*

21

The Structure of our Compiler



22

Syntax Definition

- ⌘ Context-free grammar is a 4-tuple with
 - ☒ A set of tokens (*terminal* symbols)
 - ☒ A set of *nonterminals*
 - ☒ A set of *productions*
 - ☒ A designated *start symbol*

23

Example Grammar

Context-free grammar for simple expressions:

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

with productions $P =$

$$list \rightarrow list + digit$$

$$list \rightarrow list - digit$$

$$list \rightarrow digit$$

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

24

Derivation

⌘ Given a CF grammar we can determine the set of all *strings* (sequences of tokens) generated by the grammar using *derivation*

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

25

Derivation for the Example Grammar

$$\begin{aligned}
 & \textit{list} \\
 \Rightarrow & \underline{\textit{list}} + \textit{digit} \\
 \Rightarrow & \underline{\textit{list}} - \textit{digit} + \textit{digit} \\
 \Rightarrow & \underline{\textit{digit}} - \textit{digit} + \textit{digit} \\
 \Rightarrow & \mathbf{9} - \underline{\textit{digit}} + \textit{digit} \\
 \Rightarrow & \mathbf{9} - \mathbf{5} + \underline{\textit{digit}} \\
 \Rightarrow & \mathbf{9} - \mathbf{5} + \mathbf{2}
 \end{aligned}$$

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

26

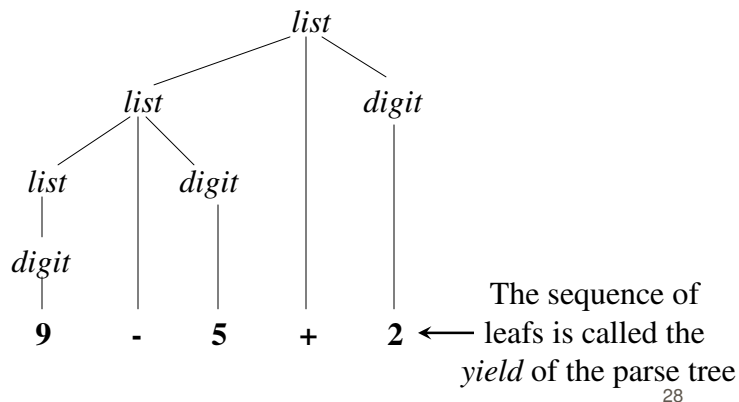
Parse Trees

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

27

Parse Tree for the Example Grammar

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



28

Ambiguity

Consider the following context-free grammar:

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

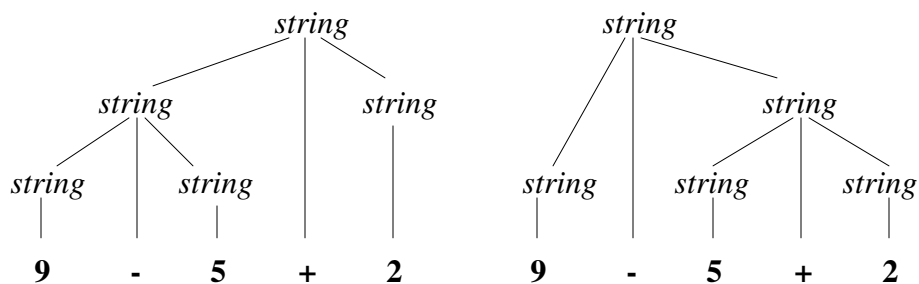
with production $P =$

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

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

29

Ambiguity (cont'd)



30

Associativity of Operators

Left-associative operators have left-recursive productions

$$\textit{left} \rightarrow \textit{left} + \textit{term} \mid \textit{term}$$

String **a+b+c** has the same meaning as **(a+b)+c**

Right-associative operators have right-recursive productions

$$\textit{right} \rightarrow \textit{term} = \textit{right} \mid \textit{term}$$

String **a=b=c** has the same meaning as **a=(b=c)**

31

Precedence of Operators

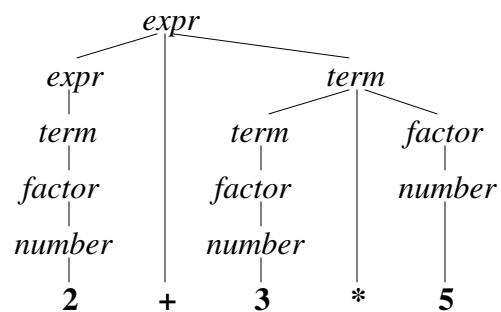
Operators with higher precedence “bind more tightly”

$$\textit{expr} \rightarrow \textit{expr} + \textit{term} \mid \textit{term}$$

$$\textit{term} \rightarrow \textit{term} * \textit{factor} \mid \textit{factor}$$

$$\textit{factor} \rightarrow \textit{number} \mid (\textit{expr})$$

String **2+3*5** has the same meaning as **2+(3*5)**



32

Syntax of Statements

$$\begin{aligned}
 \text{stmt} \rightarrow & \text{ id := expr} \\
 & | \text{ if expr then stmt} \\
 & | \text{ if expr then stmt else stmt} \\
 & | \text{ while expr do stmt} \\
 & | \text{ begin opt_stmts end} \\
 \text{opt_stmts} \rightarrow & \text{ stmt ; opt_stmts} \\
 & | \epsilon
 \end{aligned}$$

33

Syntax-Directed Translation

- ⌘ Uses a CF grammar to specify the syntactic structure of the language
- ⌘ AND associates a set of *attributes* with (non)terminals
- ⌘ AND 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

34

Synthesized and Inherited Attributes

- ⌘ An attribute is said to be ...
 - ☒ *synthesized* if its value at a parse-tree node is determined from the attribute values at the children of the node
 - ☒ *inherited* if its value at a parse-tree node is determined by the parent (by enforcing the parent's semantic rules)

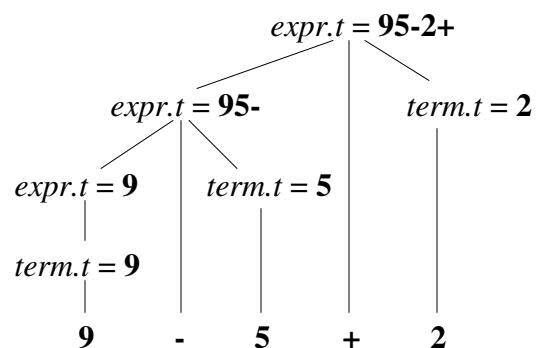
35

Example Attribute Grammar

Production	Semantic Rule
$expr \rightarrow expr_1 + term$	$expr.t := expr_1.t \parallel term.t \parallel "+"$
$expr \rightarrow expr_1 - term$	$expr.t := expr_1.t \parallel term.t \parallel "-"$
$expr \rightarrow term$	$expr.t := term.t$
$term \rightarrow 0$	$term.t := "0"$
$term \rightarrow 1$	$term.t := "1"$
...	...
$term \rightarrow 9$	$term.t := "9"$

36

Example Annotated Parse Tree



37

Depth-First Traversals

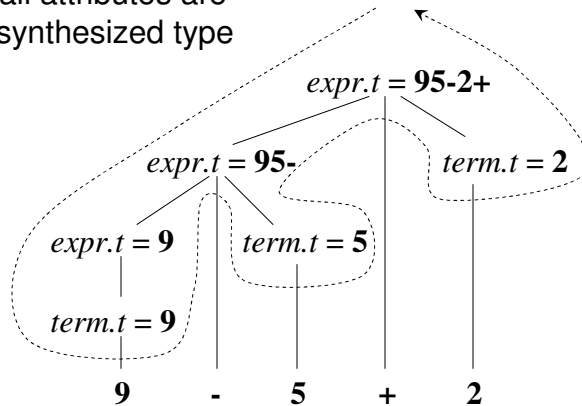
```

procedure visit(n : node);
begin
    for each child m of n, from left to right do
        visit(m);
    evaluate semantic rules at node n
end
    
```

38

Depth-First Traversals (Example)

Note: all attributes are of the synthesized type



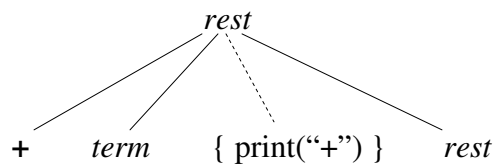
39

Translation Schemes

⌘ A *translation scheme* is a CF grammar embedded with *semantic actions*

$$rest \rightarrow + term \{ print("+") \} rest$$

Embedded
semantic action



40

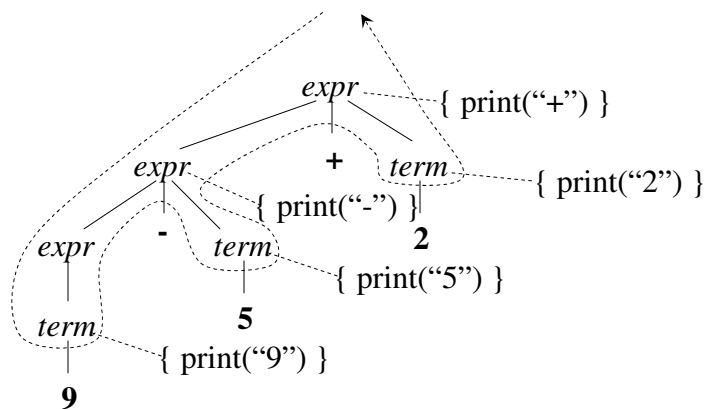
Example Translation Scheme

```

expr → expr + term  { print(“+”) }
expr → expr - term  { print(“-”) }
expr → term
term → 0              { print(“0”) }
term → 1              { print(“1”) }
...
term → 9              { print(“9”) }
    
```

41

Example Translation Scheme (cont'd)



Translates **9-5+2** into postfix **95-2+**

42

Parsing

- ⌘ Parsing = *process of determining if a string of tokens can be generated by a grammar*
- ⌘ For any CF grammar there is a parser that takes at most $O(n^3)$ time to parse a string of n tokens
- ⌘ Linear algorithms suffice for parsing programming language
- ⌘ *Top-down parsing* “constructs” parse tree from root to leaves
- ⌘ *Bottom-up parsing* “constructs” parse tree from leaves to root

43

Predictive Parsing

- ⌘ *Recursive descent parsing* is a top-down parsing method
 - ☒ Every nonterminal has one (recursive) procedure responsible for parsing the nonterminal’s syntactic category of input tokens
 - ☒ When a nonterminal has multiple productions, each production is implemented in a branch of a selection statement based on input look-ahead information
- ⌘ *Predictive parsing* is a special form of recursive descent parsing where we use one lookahead token to unambiguously determine the parse operations

44

Example Predictive Parser (Grammar)

```

type → simple
      | ^ id
      | array [ simple ] of type
simple → integer
      | char
      | num dotdot num
    
```

45

Example Predictive Parser (Program Code)

```

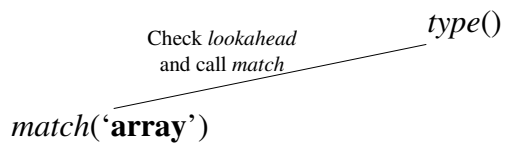
procedure match(t : token);
begin
  if lookahead = t then
    lookahead := nexttoken()
  else error()
end;

procedure type();
begin
  if lookahead in { 'integer', 'char', 'num' } then
    simple()
  else if lookahead = '^' then
    match('^'); match(id)
  else if lookahead = 'array' then
    match('array'); match('['); simple();
    match(']'); match('of'); type()
  else error()
end;

procedure simple();
begin
  if lookahead = 'integer' then
    match('integer')
  else if lookahead = 'char' then
    match('char')
  else if lookahead = 'num' then
    match('num');
    match('dotdot');
    match('num')
  else error()
end;
    
```

46

Example Predictive Parser (Execution Step 1)

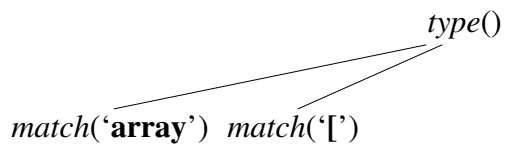


Input: array [num dotdot num] of integer

 ↑

 lookahead

Example Predictive Parser (Execution Step 2)

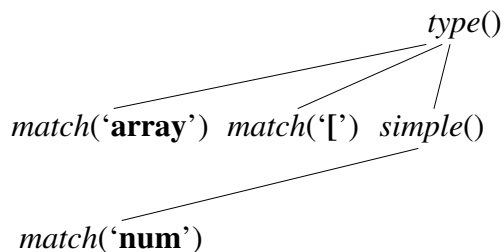


Input: array [num dotdot num] of integer

 ↑

 lookahead

Example Predictive Parser (Execution Step 3)



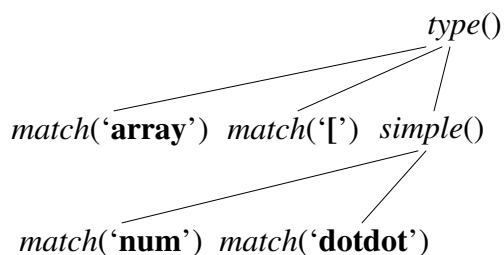
Input: array [num dotdot num] of integer

↑

lookahead

49

Example Predictive Parser (Execution Step 4)



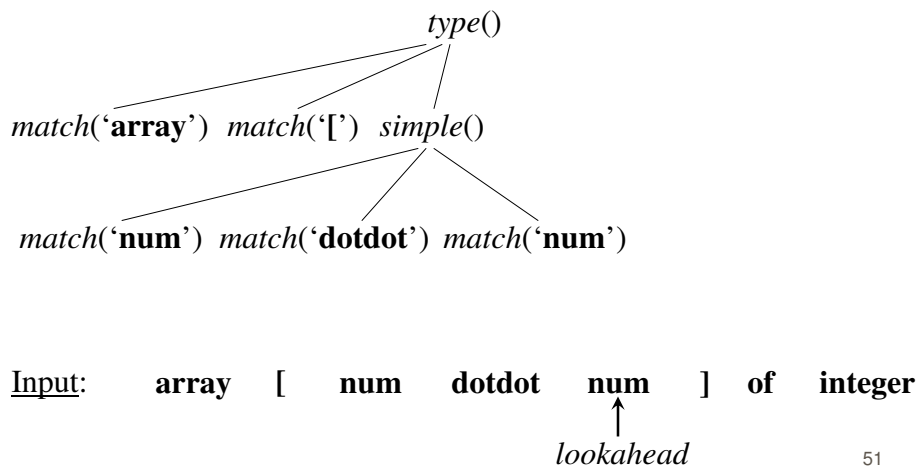
Input: array [num dotdot num] of integer

↑

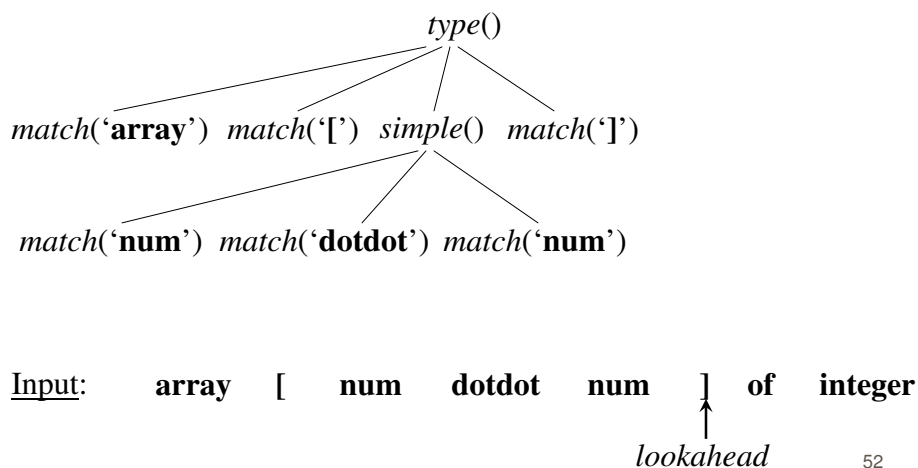
lookahead

50

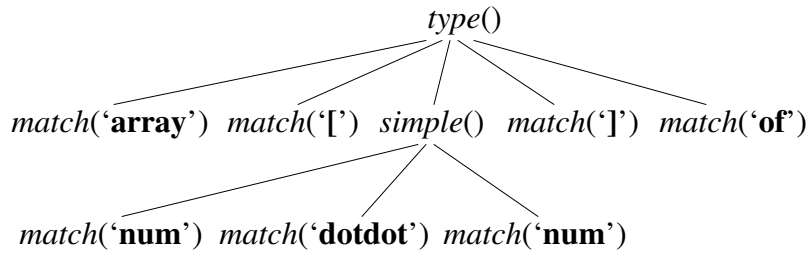
Example Predictive Parser (Execution Step 5)



Example Predictive Parser (Execution Step 6)

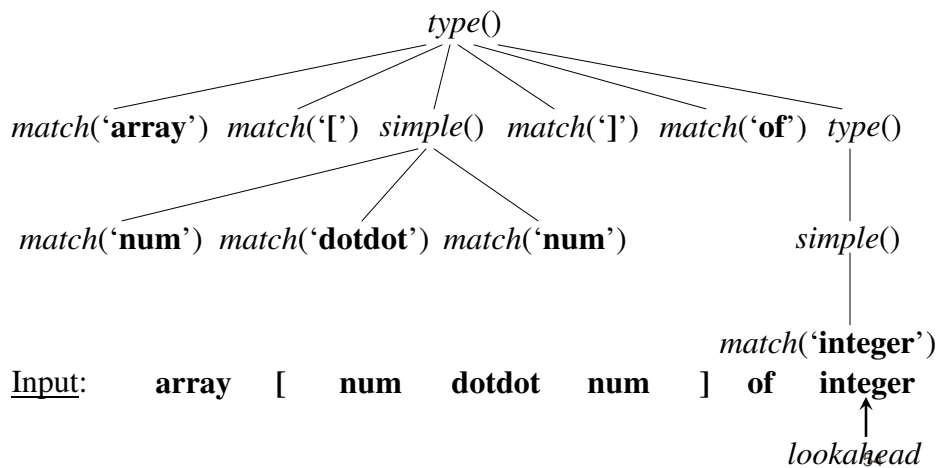


Example Predictive Parser (Execution Step 7)



Input: array [num dotdot num] of integer
↑
lookahead 53

Example Predictive Parser (Execution Step 8)



Input: array [num dotdot num] of integer
↑
lookahead

FIRST

FIRST(α) is the set of terminals that appear as the first symbols of one or more strings generated from α

```

type → simple
      | ^ id
      | array [ simple ] of type
simple → integer
      | char
      | num dotdot num
    
```

```

FIRST(simple) = { integer, char, num }
FIRST(^ id)   = { ^ }
FIRST(type)   = { integer, char, num, ^, array }
    
```

55

Using FIRST

We use FIRST to write a predictive parser as follows

```

expr → term rest
rest → + term rest
      | - term rest
      | ε
    
```

```

procedure rest();
begin
  if lookahead in FIRST(+ term rest) then
    match('+'); term(); rest()
  else if lookahead in FIRST(- term rest) then
    match('-'); term(); rest()
  else return
end;
    
```

When a nonterminal A has two (or more) productions as in

$$A \rightarrow \alpha$$

$$| \beta$$

Then FIRST(α) and FIRST(β) must be disjoint for predictive parsing to work

56

Left Factoring

When more than one production for nonterminal A starts with the same symbols, the FIRST sets are not disjoint

$$\begin{aligned} stmt &\rightarrow \mathbf{if\ expr\ then\ stmt} \\ &\quad | \mathbf{if\ expr\ then\ stmt\ else\ stmt} \end{aligned}$$

We can use *left factoring* to fix the problem

$$\begin{aligned} stmt &\rightarrow \mathbf{if\ expr\ then\ stmt\ opt_else} \\ opt_else &\rightarrow \mathbf{else\ stmt} \\ &\quad | \epsilon \end{aligned}$$

Left factoring: if not clear what to choose, rewrite the production until we have seen enough to make a decision.

57

Left Recursion

When a production for nonterminal A starts with a *self reference* then a predictive parser loops forever

$$\begin{aligned} A &\rightarrow A\ \alpha \\ &\quad | \beta \\ &\quad | \gamma \end{aligned}$$

We can eliminate *left recursive productions* by systematically rewriting the grammar using *right recursive productions*

$$\begin{aligned} A &\rightarrow \beta R \\ &\quad | \gamma R \\ R &\rightarrow \alpha R \\ &\quad | \epsilon \end{aligned}$$

58

A Translator for Simple Expressions

```

expr → expr + term { print(“+”) }
expr → expr - term { print(“-”) }
expr → term
term → 0 { print(“0”) }
term → 1 { print(“1”) }
...
term → 9 { print(“9”) }
    
```

After left recursion elimination:

```

expr → term rest
rest → + term { print(“+”) } rest | - term { print(“-”) } rest | ε
term → 0 { print(“0”) }
term → 1 { print(“1”) }
...
term → 9 { print(“9”) }
    
```

59

	<pre> main() { lookahead = getchar(); expr(); } expr() { term(); while (1) /* optimized by inlining rest() and removing recursive calls */ { if (lookahead == '+') { match('+'); term(); putchar('+'); } else if (lookahead == '-') { match('-'); term(); putchar('-'); } else break; } } term() { if (isdigit(lookahead)) { putchar(lookahead); match(lookahead); } else error(); } match(int t) { if (lookahead == t) { lookahead = getchar(); } else error(); } error() { printf("Syntax error\n"); exit(1); } </pre>
<i>expr</i> → <i>term rest</i>	}
<pre> <i>rest</i> → + <i>term</i> { print(“+”) } <i>rest</i> - <i>term</i> { print(“-”) } <i>rest</i> ε </pre>	
<pre> <i>term</i> → 0 { print(“0”) } <i>term</i> → 1 { print(“1”) } ... <i>term</i> → 9 { print(“9”) } </pre>	}

60