

Principles of Programming Languages

<http://www.di.unipi.it/~andrea/Didattica/PLP-14/>

Prof. Andrea Corradini

Department of Computer Science, Pisa

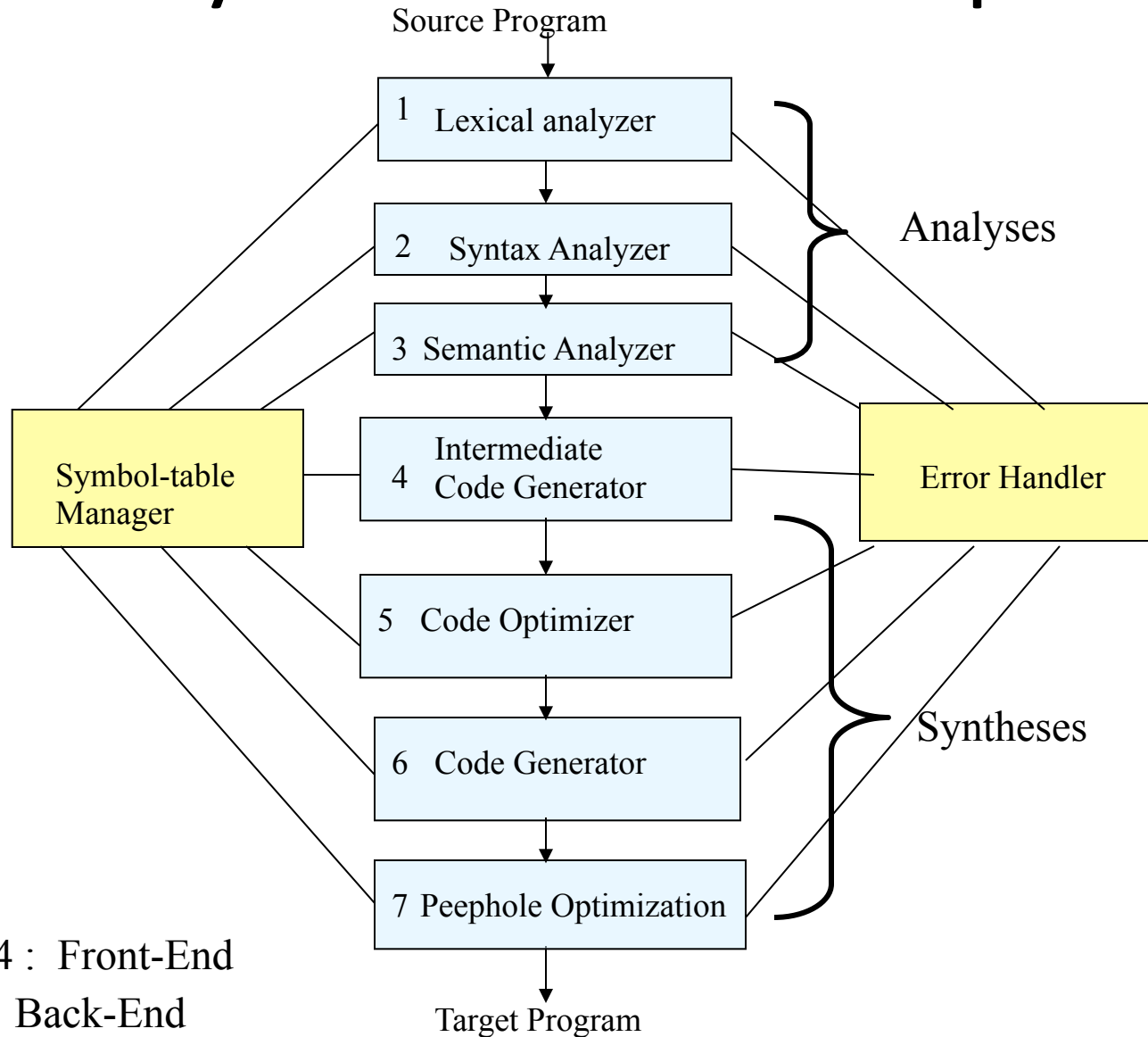
Lesson 2

- The structure of a compiler
- Overview of a Simple Compiler Front-end
 - Predictive top-down parsing
 - Syntax directed translation
 - Lexical analysis

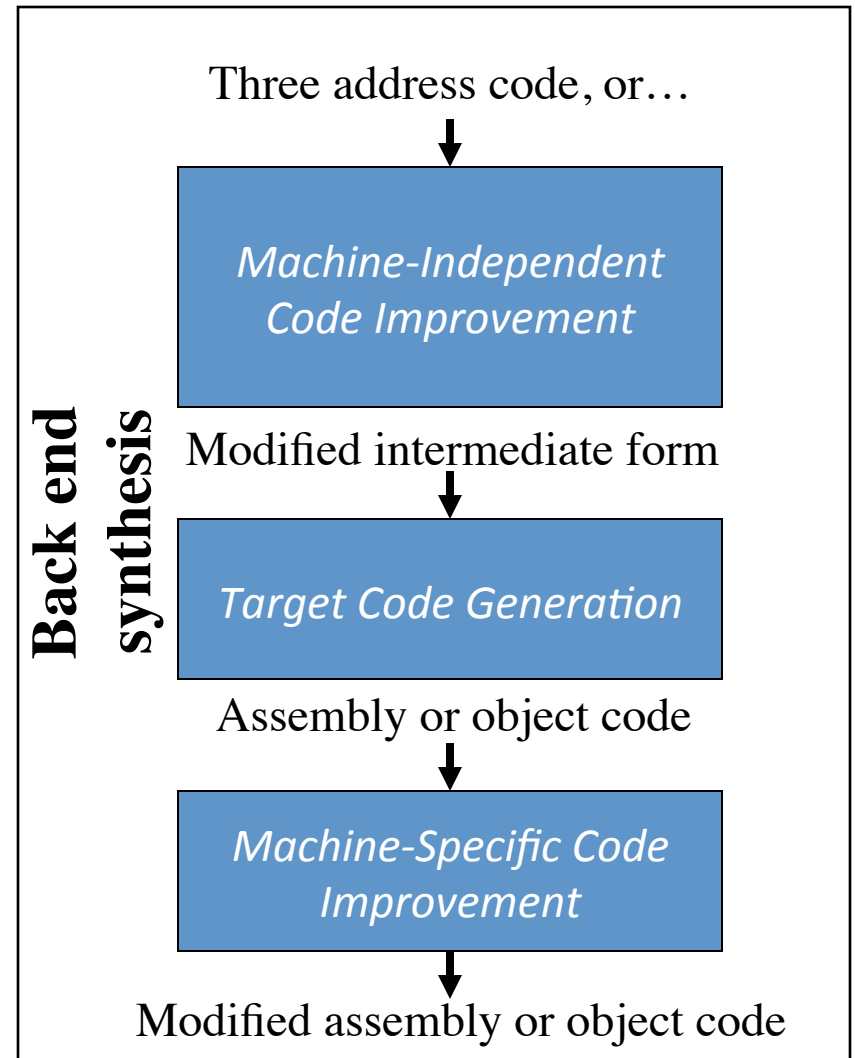
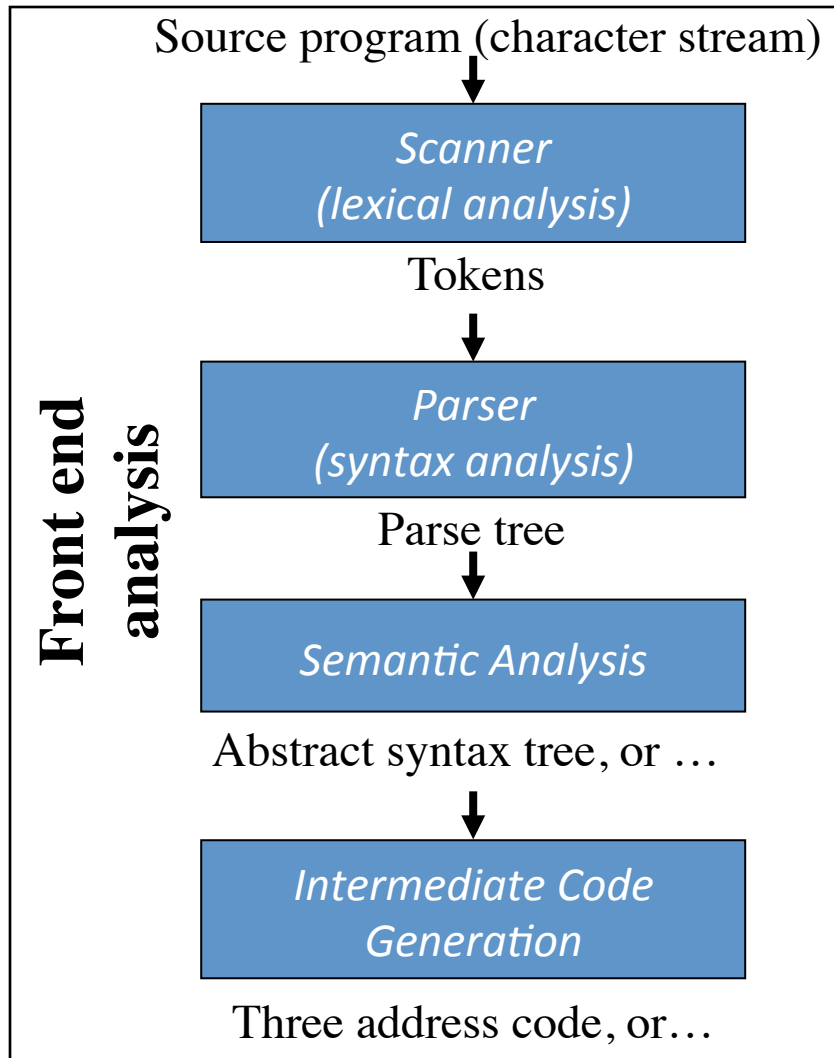
Admins

- Office Hours:
 - Wednesday, 9 - 11 ← my proposal
 - *Monday, 18 - 19:30*
 - *Friday, 9-11*
- Check your data and add the University ID (matricola) in the sheet

The Many Phases of a Compiler



Compiler Front- and Back-end



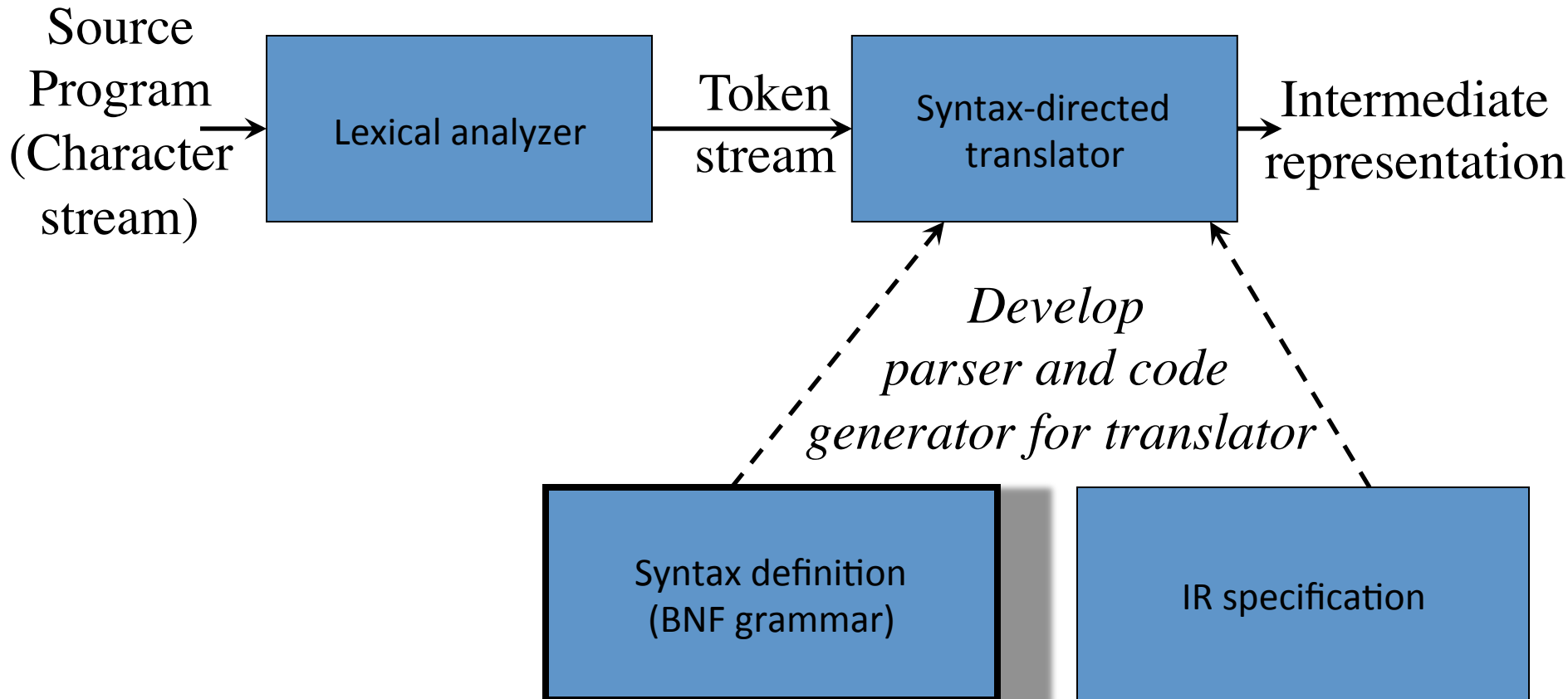
Single-pass vs. Multi-pass Compilers

- A collection of compilation phases is done only once (*single pass*) or multiple times (*multi pass*)
- **Single pass:** more efficient and uses less memory
 - requires everything to be defined before being used
 - standard for languages like Pascal, FORTRAN, C
 - Influenced the design of early programming languages
- **Multi pass:** needs more memory (to keep entire program), usually slower
 - needed for languages where declarations e.g. of variables may follow their use (Java, ADA, ...)
 - allows better optimization of target code

Overview of a Simple Compiler Front-end

- Building a compiler involves:
 - Defining the *syntax* of a programming language
 - Develop a source code parser: we consider here *predictive parsing*
 - Implementing *syntax directed translation* to generate intermediate code

The Structure of the Front-End



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*

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 $P =$

$$list \rightarrow list + digit$$

$$list \rightarrow list - digit$$

$$list \rightarrow digit$$

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

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

Derivation for the Example Grammar

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

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

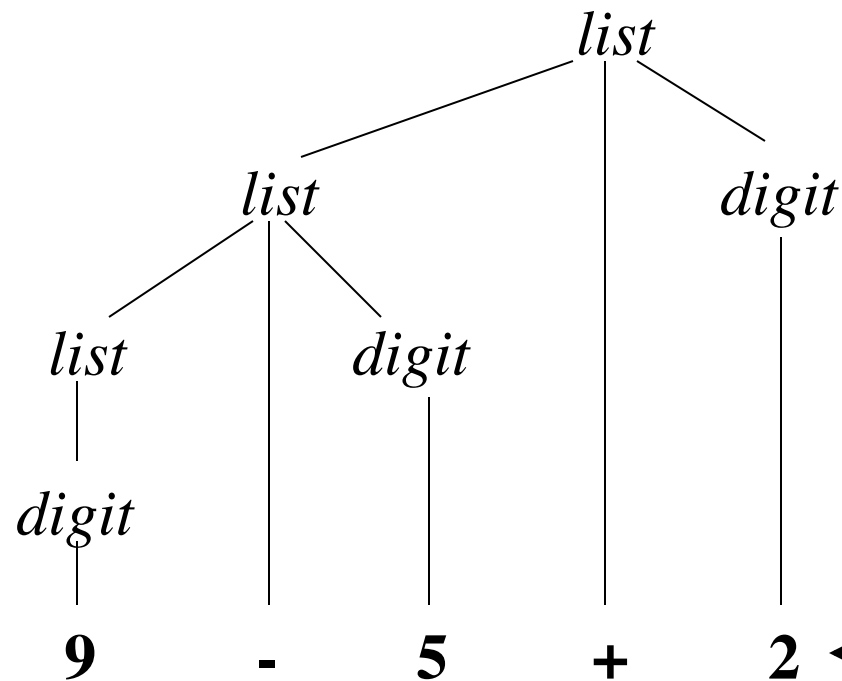
Likewise, a *rightmost derivation* replaces the rightmost nonterminal in each step

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
- 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 immediate *children* X_1, X_2, \dots, X_n where X_i is a (non)terminal or ε (ε denotes the *empty string*)

Parse Tree for the Example Grammar

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



The sequence of
leaves is called the
yield of the parse tree

Ambiguity

Consider the following context-free grammar:

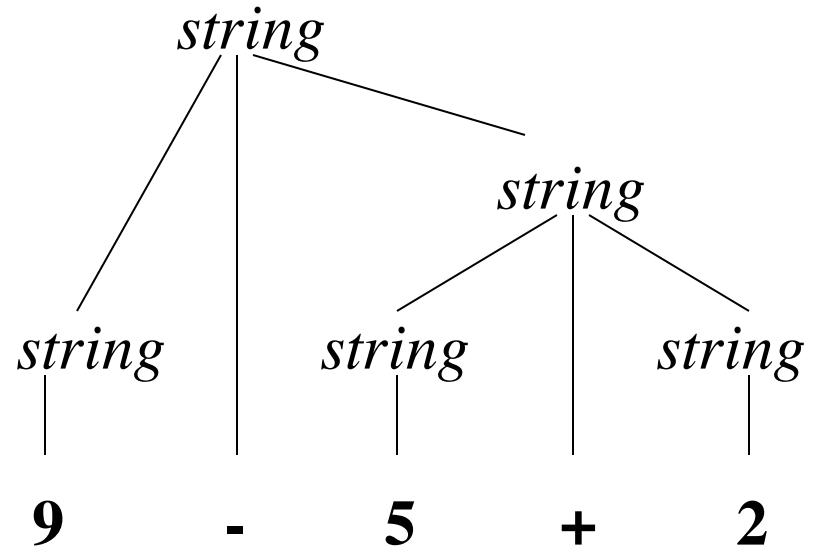
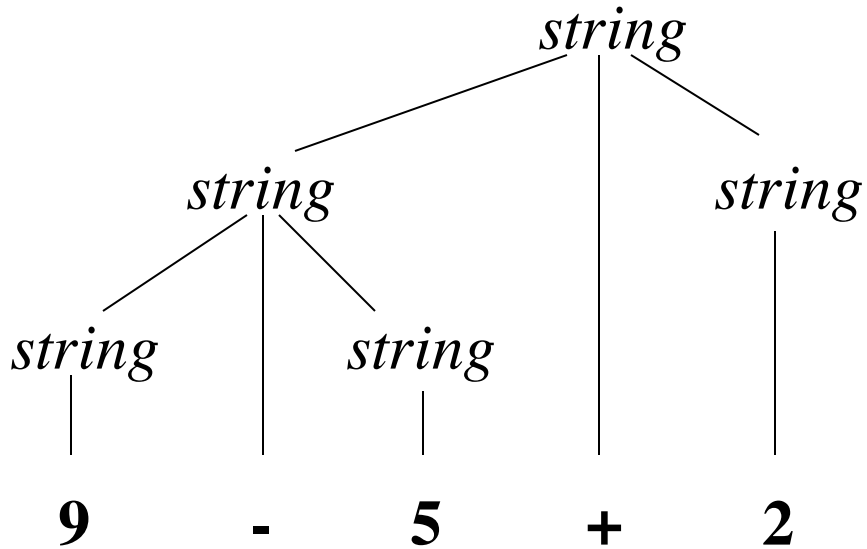
$$G = \langle \{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string \rangle$$

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 represents the string **9-5+2**

Ambiguity (cont'd)



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)**

Precedence of Operators

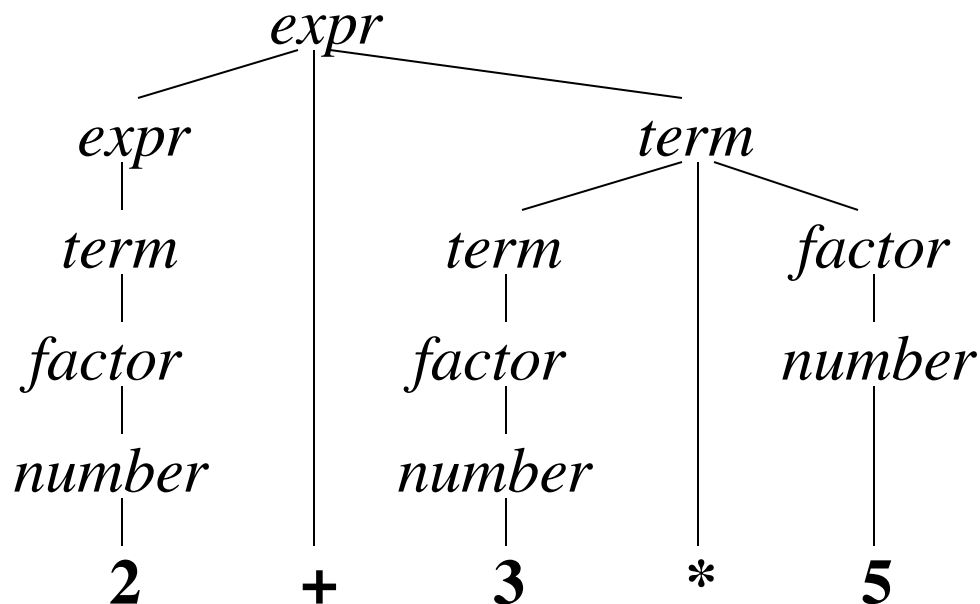
Operators with higher precedence “bind more tightly”

$expr \rightarrow expr + term \mid term$

$term \rightarrow term * factor \mid factor$

$factor \rightarrow number \mid (expr)$

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

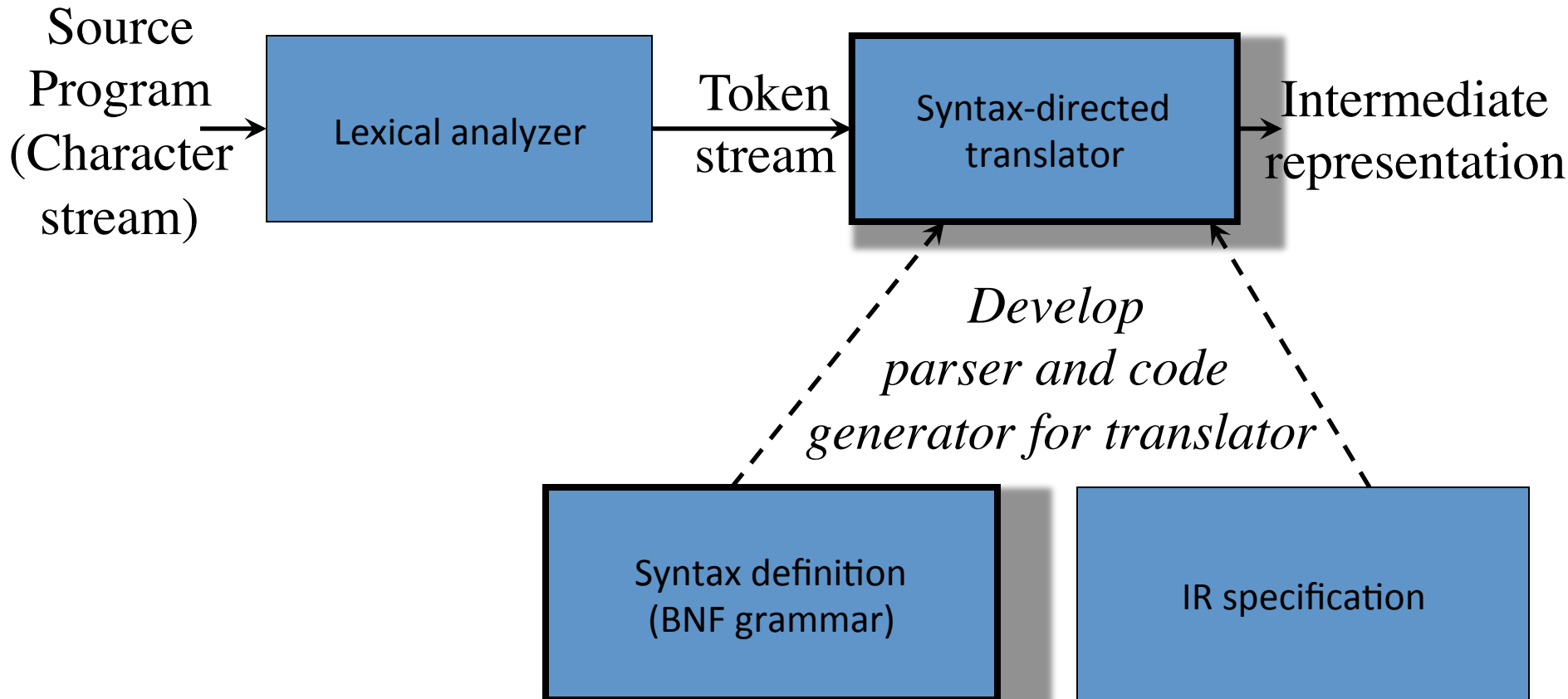


Syntax of Statements

$stmt \rightarrow$ **id** := *expr*
| **if** *expr* **then** *stmt*
| **if** *expr* **then** *stmt* **else** *stmt*
| **while** *expr* **do** *stmt*
| **begin** *opt_stmts* **end**

$opt_stmts \rightarrow$ *stmt* ; *opt_stmts*
| ϵ

The Structure of the Front-End



Syntax-Directed Translation


- Uses a CF grammar to specify the syntactic structure of the language
- AND associates a set of *attributes* with the terminals and nonterminals of the grammar
- AND associates with each production a set of *semantic rules* to compute values of attributes
- A parse tree is traversed and semantic rules applied: after the tree traversal(s) are completed, the attribute values on the nonterminals contain the translated form of the input

Synthesized and Inherited Attributes

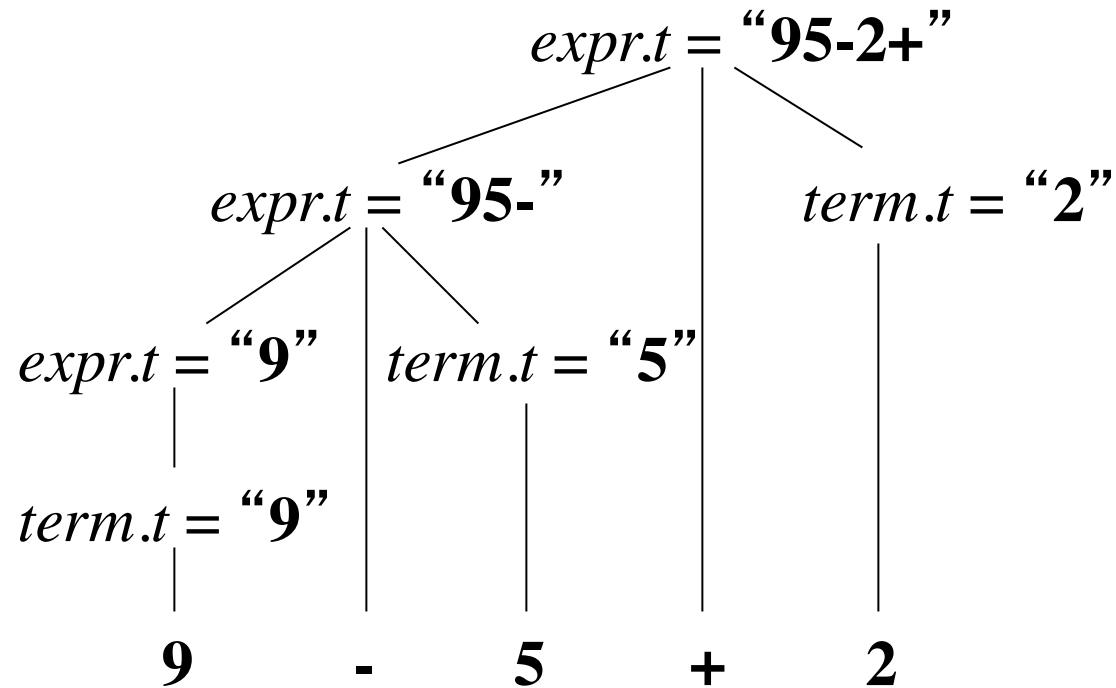
- An attribute is said to be ...
 - *synthesized* if its value at a parse-tree node is determined from the attribute values at the children of the node
 - *inherited* if its value at a parse-tree node is determined by the parent (by enforcing the parent's semantic rules)

Example Attribute Grammar (Postfix Form)

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

String concat operator


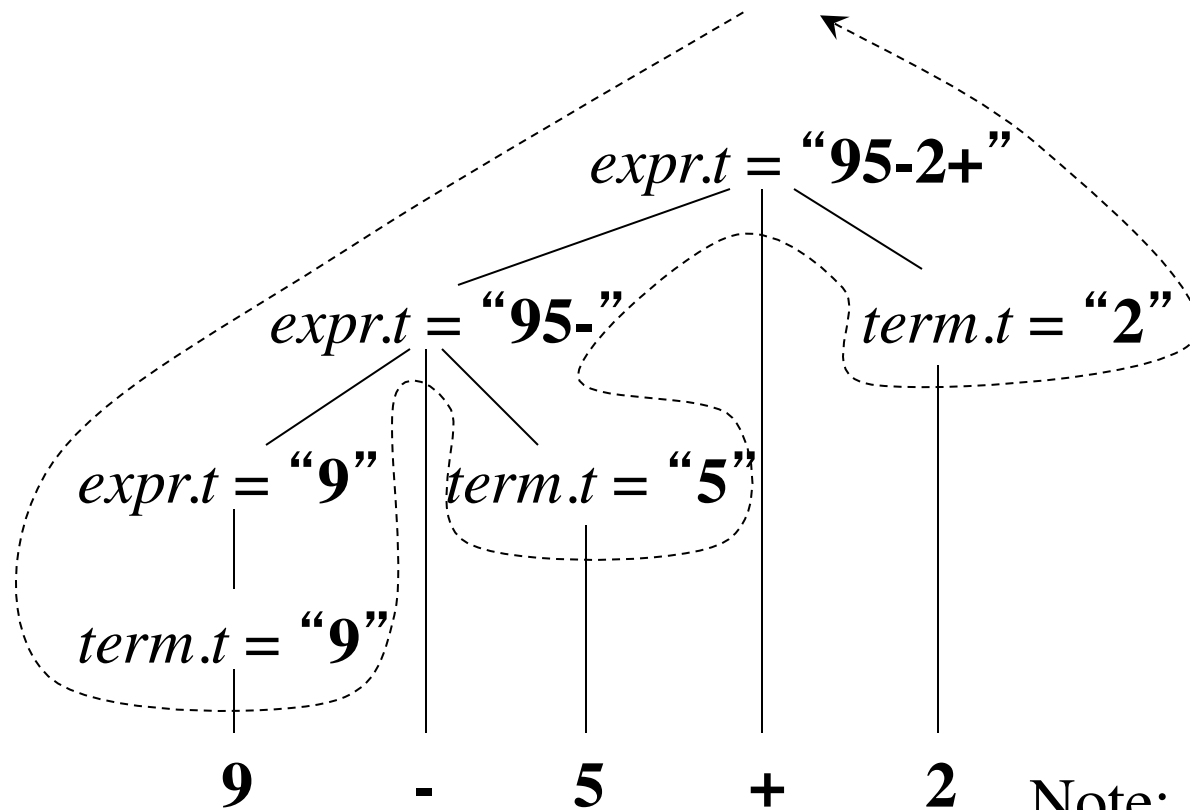
Example Annotated Parse Tree



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
```


Depth-First Traversals (Example)




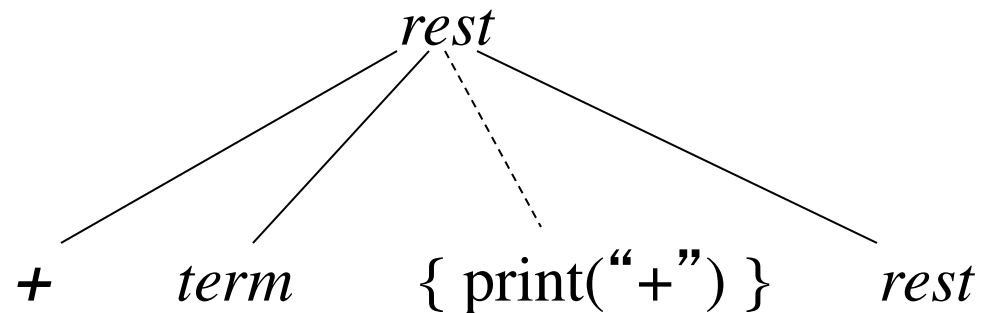
Note: all attributes are of the synthesized type

Translation Schemes

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

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

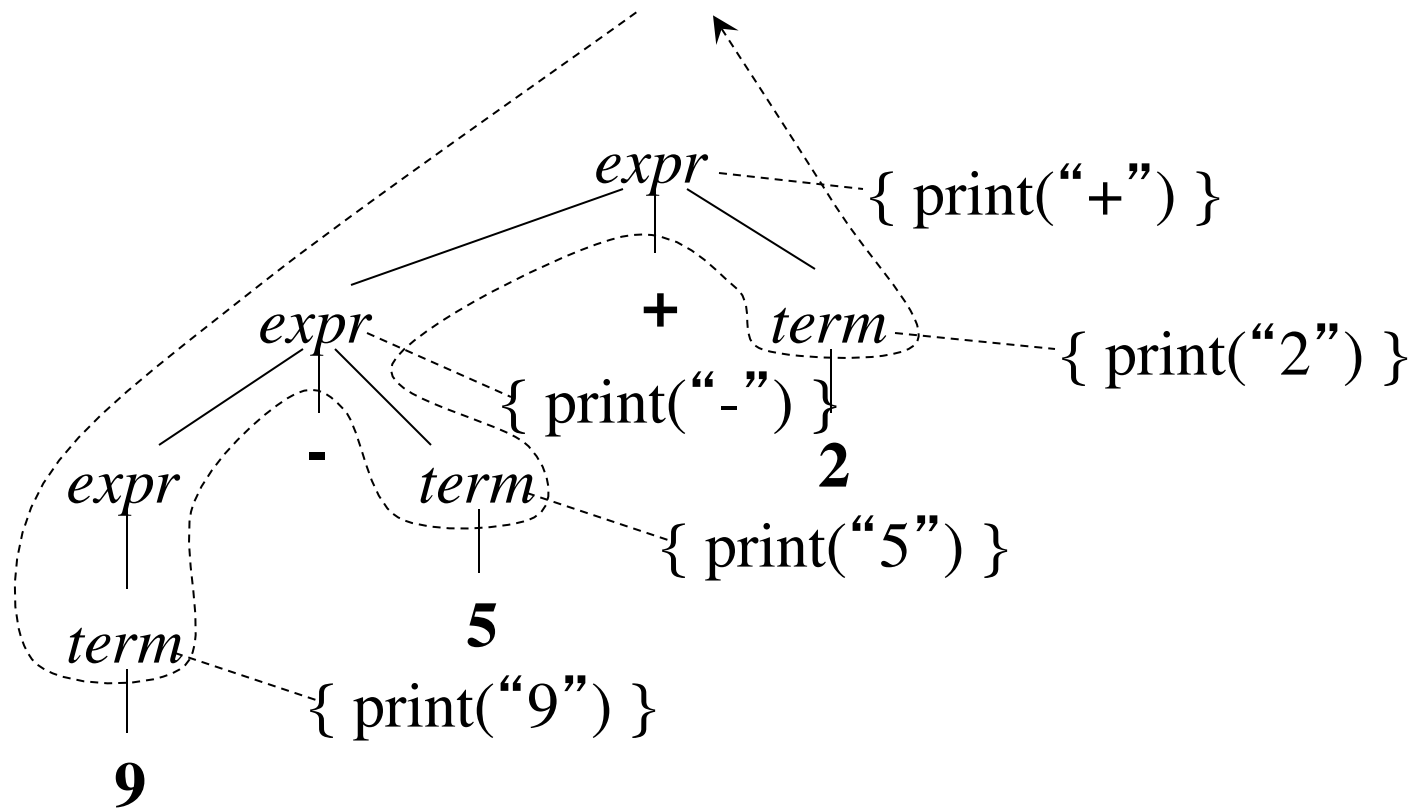

Embedded
semantic action



Example Translation Scheme for Postfix Notation

$expr \rightarrow expr + term$ { print(“+”) }
 $expr \rightarrow expr - term$ { print(“-”) }
 $expr \rightarrow term$
 $term \rightarrow 0$ { print(“0”) }
 $term \rightarrow 1$ { print(“1”) }
...
 $term \rightarrow 9$ { print(“9”) }

Example Translation Scheme (cont'd)



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

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 source code
- *Top-down parsing* “constructs” a parse tree from root to leaves
- *Bottom-up parsing* “constructs” a parse tree from leaves to root

Predictive Parsing

- *Recursive descent parsing* is a top-down parsing method
 - Each nonterminal has one (recursive) procedure that is 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

Example Predictive Parser (Grammar)

type → *simple*
| **^ id**
| **array [*simple*] of *type***

simple → **integer**
| **char**
| **num dotdot num**

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;
```

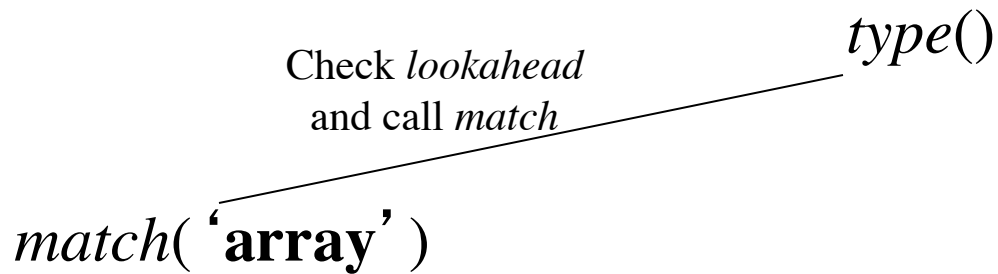
```
procedure simple();  
begin  
  if lookahead = 'integer' then  
    match('integer')  
  else if lookahead = 'char' then  
    match('char')  
  else if lookahead = 'num' then  
    match('num');  
    match('dotdot');  
    match('num')  
  else error()  
end;
```


Example Predictive Parser (Execution Step 1)

match('array')

Check *lookahead*
and call *match*

type()



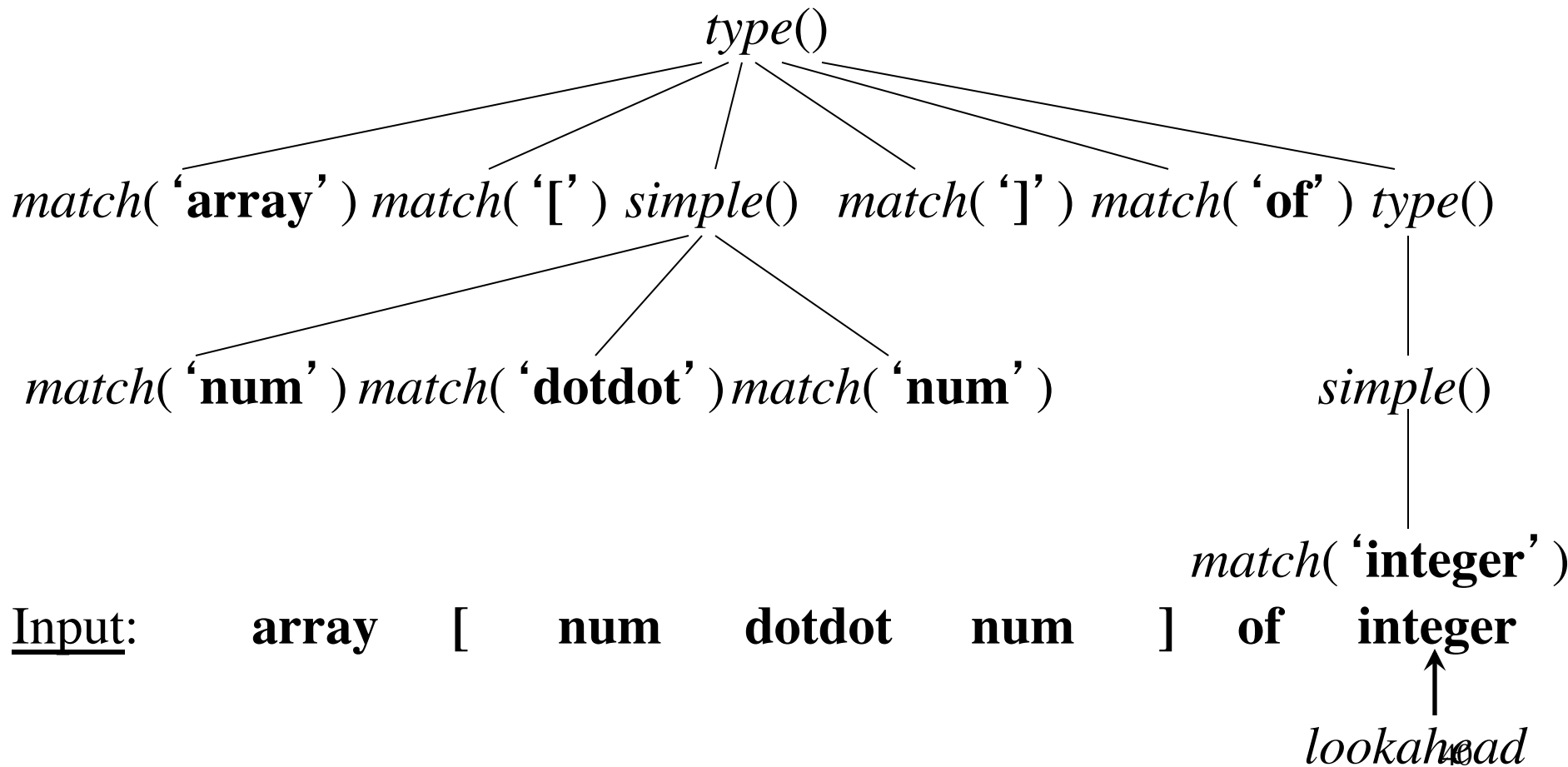
Input: **array** [**num** **dotdot** **num**] **of** **integer**

 ↑

lookahead



Example Predictive Parser (Execution Step 8)



FIRST

$\text{FIRST}(\alpha)$ is the set of terminals that appear as the first symbols of one or more strings generated from α

type \rightarrow *simple*
| \wedge **id**
| **array** [*simple*] **of type**
simple \rightarrow **integer**
| **char**
| **num dotdot num**

$\text{FIRST}(\textit{simple}) = \{ \mathbf{integer}, \mathbf{char}, \mathbf{num} \}$

$\text{FIRST}(\wedge \mathbf{id}) = \{ \wedge \}$

$\text{FIRST}(\textit{type}) = \{ \mathbf{integer}, \mathbf{char}, \mathbf{num}, \wedge, \mathbf{array} \}$

How to use FIRST

We use FIRST to write a predictive parser as follows

| | | |
|----------------------------------|--|--|
| $expr \rightarrow term\ rest$ | | procedure <i>rest</i> (); |
| $rest \rightarrow +\ term\ rest$ | | begin |
| $- term\ rest$ | | if <i>lookahead</i> in <u>FIRST(+ term rest)</u> then |
| ϵ | | <i>match</i> ('+'); <i>term</i> (); <i>rest</i> () |
| | | else if <i>lookahead</i> in <u>FIRST(- term rest)</u> then |
| | | <i>match</i> (' - '); <i>term</i> (); <i>rest</i> () |
| | | 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

Left Factoring

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

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

We can use *left factoring* to fix the problem

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

Left Recursion

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

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

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

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

A Translator for Simple Expressions

$expr \rightarrow expr + term \quad \{ \text{print}("+") \}$

$expr \rightarrow expr - term \quad \{ \text{print}("-") \}$

$expr \rightarrow term$

$term \rightarrow 0 \quad \{ \text{print}("0") \}$

$term \rightarrow 1 \quad \{ \text{print}("1") \}$

...

$term \rightarrow 9 \quad \{ \text{print}("9") \}$

After left recursion elimination:

$expr \rightarrow term \text{ rest}$

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

$rest \rightarrow - term \quad \{ \text{print}("-") \} \text{ rest}$

$rest \rightarrow \epsilon$

$term \rightarrow 0 \quad \{ \text{print}("0") \}$

$term \rightarrow 1 \quad \{ \text{print}("1") \}$

...

$term \rightarrow 9 \quad \{ \text{print}("9") \}$

Code of the translator

$expr \rightarrow term\ rest$

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

$rest \rightarrow -\ term\ \{ \text{print}("-") \}\ rest$

$rest \rightarrow \epsilon$

$term \rightarrow 0\ \{ \text{print}("0") \}$

$term \rightarrow 1\ \{ \text{print}("1") \}$

...

$term \rightarrow 9\ \{ \text{print}("9") \}$

```
main()
{   lookahead = getchar();
    expr();
}

expr()
{   term(); rest(); }

rest ()
{   if (lookahead == '+')
    {match('+'); term(); putchar('+'); rest();
    }
    else if (lookahead == '-')
    {match('-'); term(); putchar('-'); rest();
    }
    else {};
}

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);
}
```

Optimized code of the translator

$expr \rightarrow term\ rest$

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

$rest \rightarrow -\ term\ \{ \text{print}(\text{"-"}) \}\ rest$

$rest \rightarrow \epsilon$

$term \rightarrow 0\ \{ \text{print}(\text{"0"}) \}$

$term \rightarrow 1\ \{ \text{print}(\text{"1"}) \}$

...

$term \rightarrow 9\ \{ \text{print}(\text{"9"}) \}$

```
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);
}
```