

Principles of Programming Languages

<http://www.di.unipi.it/~andrea/Didattica/PLP-15/>

Prof. Andrea Corradini

Department of Computer Science, Pisa

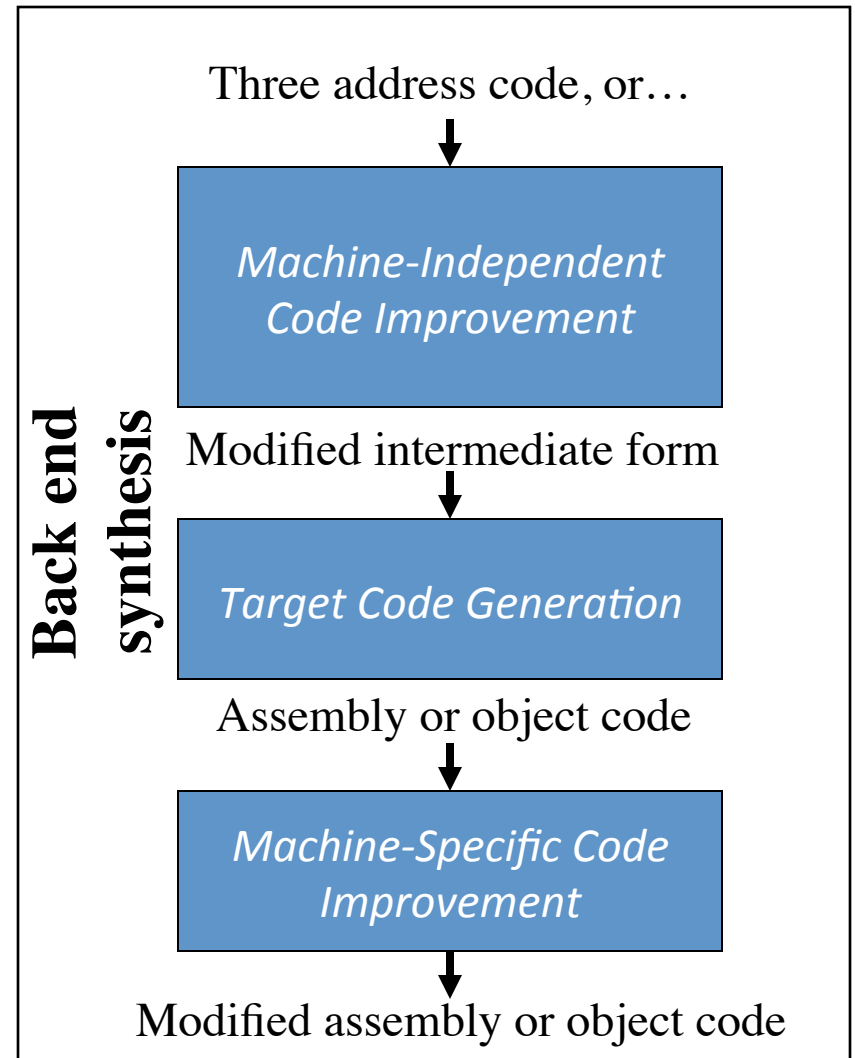
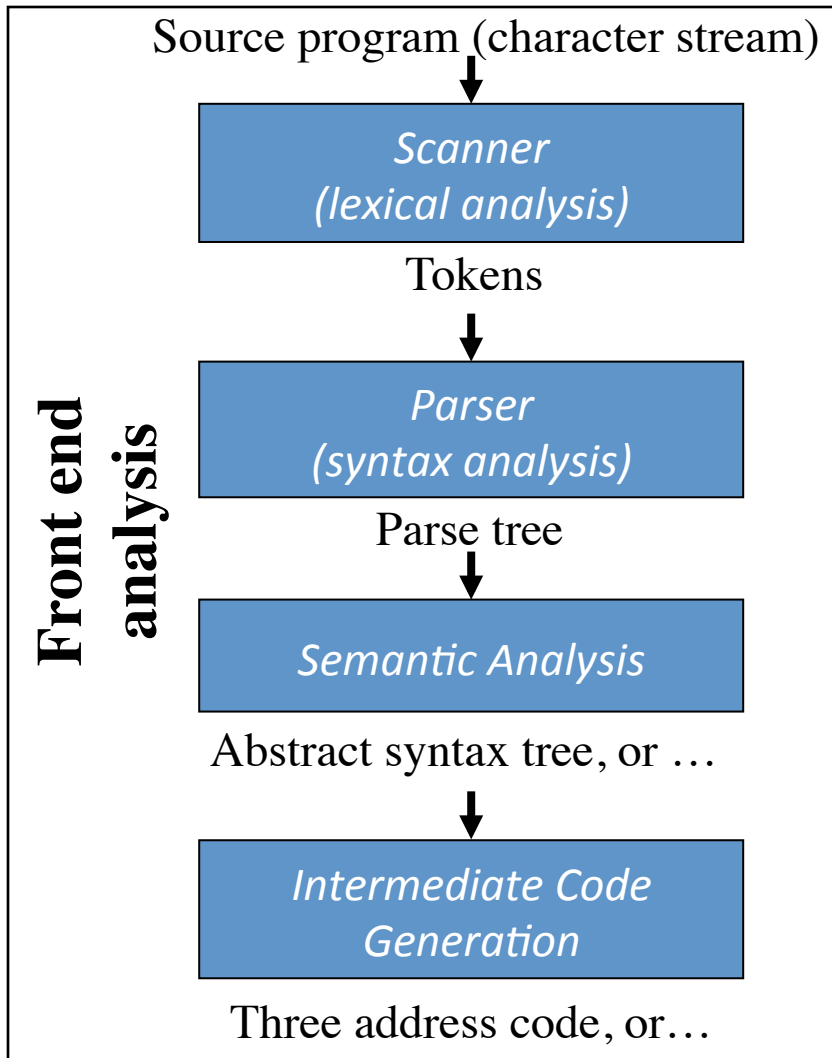
Lesson 3

- Overview of a syntax-directed compiler front-end

Overview of syntax-directed front-end

- (Context-Free) Grammars, Chomsky hierarchy
- Parse trees
- Ambiguity, associativity and precedence
- Syntax-directed translation
- Translation schemes
- Predictive recursive descent parsing
- Left factoring, elimination of left recursion

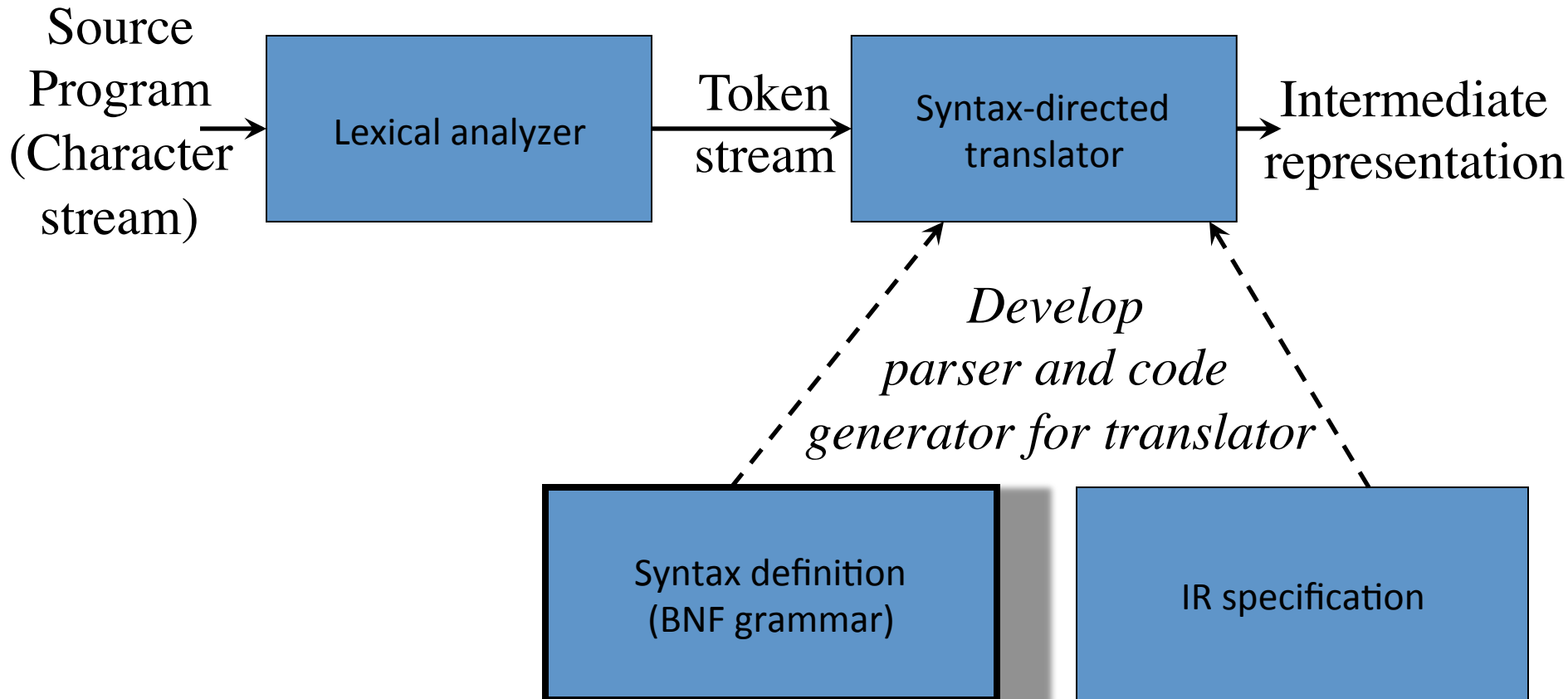
Compiler Front- and Back-end



A simple syntax-directed Compiler Front-end

- Overview of the front-end of a compiler with:
 - Definition of the context-free syntax of a programming language
 - Presentation of a source code parser: *top-down predictive parsing*
 - *Lexical analysis*
 - Implementing *syntax directed translation* to generate intermediate code

The Structure of the Front-End



Syntax Definition: Grammars

- A **grammar** is a 4-tuple $G = (N, T, P, S)$ where
 - T is a finite set of tokens (*terminal* symbols)
 - N is a finite set of *nonterminals*
 - P is a finite set of *productions* of the form
$$\alpha \rightarrow \beta$$
where $\alpha \in (NUT)^* N (NUT)^*$ and $\beta \in (NUT)^*$
 - $S \in N$ is a designated *start symbol*
- \mathbf{A}^* is the set of finite sequences of elements of \mathbf{A} . If $\mathbf{A} = \{a, b\}$, $\mathbf{A}^* = \{\epsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$
- $\mathbf{AB} = \{ab \mid a \in \mathbf{A}, b \in \mathbf{B}\}$

Notational Conventions Used

- Terminals
 $a, b, c, \dots \in T$
specific terminals: **0**, **1**, **id**, **+**
- Nonterminals
 $A, B, C, \dots \in N$
specific nonterminals: *expr*, *term*, *stmt*
- Grammar symbols
 $X, Y, Z \in (NUT)$
- Strings of terminals
 $u, v, w, x, y, z \in T^*$
- Strings of grammar symbols
 $\alpha, \beta, \gamma \in (NUT)^*$

Derivations

- A *one-step derivation* is defined by
$$\gamma \alpha \delta \Rightarrow \gamma \beta \delta$$
where $\alpha \rightarrow \beta$ is a production in the grammar
- In addition, we define
 - \Rightarrow is *leftmost* \Rightarrow_{lm} if γ does not contain a nonterminal
 - \Rightarrow is *rightmost* \Rightarrow_{rm} if δ does not contain a nonterminal
 - Transitive closure \Rightarrow^* (zero or more steps)
 - Positive closure \Rightarrow^+ (one or more steps)
- α is a *sentential form* if $S \Rightarrow^* \alpha$
- The *language generated by G* is defined by

$$L(G) = \{w \in T^* \mid S \Rightarrow^+ w\}$$

Derivation (Example)

Grammar $G = (\{E\}, \{+, *, (,), -, \text{id}\}, P, E)$ with productions

$$P = E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow (E)$$

$$E \rightarrow - E$$

$$E \rightarrow \text{id}$$

Example derivations:

$$E \Rightarrow - E \Rightarrow - \text{id}$$

$$E \Rightarrow_{rm} E + E \Rightarrow_{rm} E + \text{id} \Rightarrow_{rm} \text{id} + \text{id}$$

$$E \Rightarrow^* E$$

$$E \Rightarrow^* \text{id} + \text{id}$$

$$E \Rightarrow^+ \text{id} * \text{id} + \text{id}$$

Another grammar for expressions

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

Productions $P =$ $list \rightarrow list + digit$

$list \rightarrow list - digit$

$list \rightarrow digit$

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

A leftmost derivation:

list

$\Rightarrow_{lm} \underline{list} + digit$

$\Rightarrow_{lm} \underline{list} - digit + digit$

$\Rightarrow_{lm} \underline{digit} - digit + digit$

$\Rightarrow_{lm} 9 - \underline{digit} + digit$

$\Rightarrow_{lm} 9 - 5 + \underline{digit}$

$\Rightarrow_{lm} 9 - 5 + 2$

Chomsky Hierarchy: Language Classification

- A grammar G is said to be
 - *Regular* if it is *right linear* where each production is of the form
$$A \rightarrow w B \quad \text{or} \quad A \rightarrow w$$
or *left linear* where each production is of the form
$$A \rightarrow B w \quad \text{or} \quad A \rightarrow w \quad (w \in T^*)$$
 - *Context free* if each production is of the form
$$A \rightarrow \alpha$$
where $A \in N$ and $\alpha \in (N \cup T)^*$
 - *Context sensitive* if each production is of the form
$$\alpha A \beta \rightarrow \alpha \gamma \beta$$
where $A \in N$, $\alpha, \gamma, \beta \in (N \cup T)^*$, $|\gamma| > 0$
 - *Unrestricted*

Chomsky Hierarchy

$\mathcal{L}(\text{regular}) \subset \mathcal{L}(\text{context free}) \subset \mathcal{L}(\text{context sensitive}) \subset \mathcal{L}(\text{unrestricted})$

Where $\mathcal{L}(T) = \{ L(G) \mid G \text{ is of type } T \}$
That is: the set of all languages
generated by grammars G of type T

Examples:

Every *finite language* is regular! (construct a FSA for strings in $L(G)$)

$L_1 = \{ \mathbf{a}^n \mathbf{b}^n \mid n \geq 1 \}$ is context free

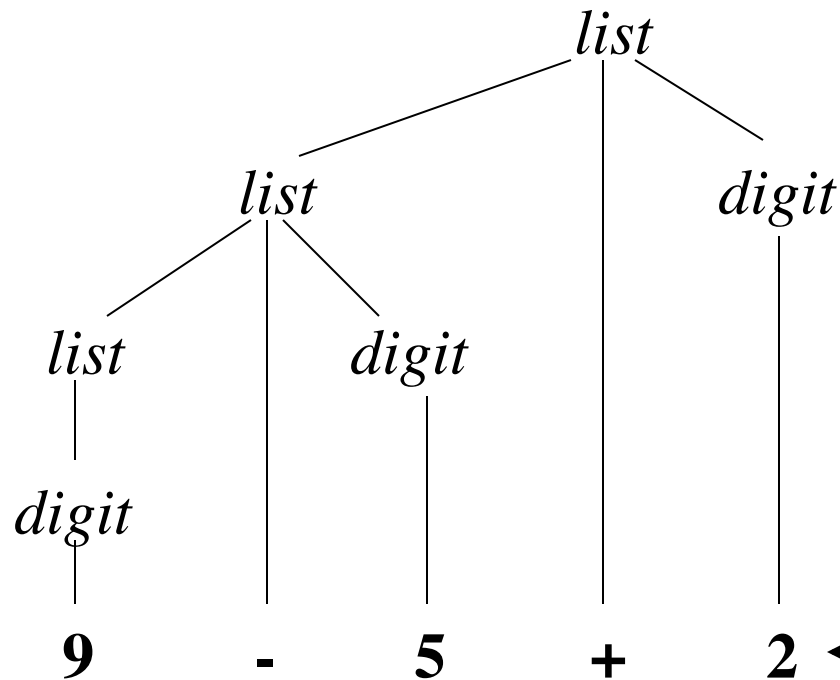
$L_2 = \{ \mathbf{a}^n \mathbf{b}^n \mathbf{c}^n \mid n \geq 1 \}$ is context sensitive

Parse Trees (context-free grammars)

- Tree-shaped representation of derivations
- 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 *internal 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:

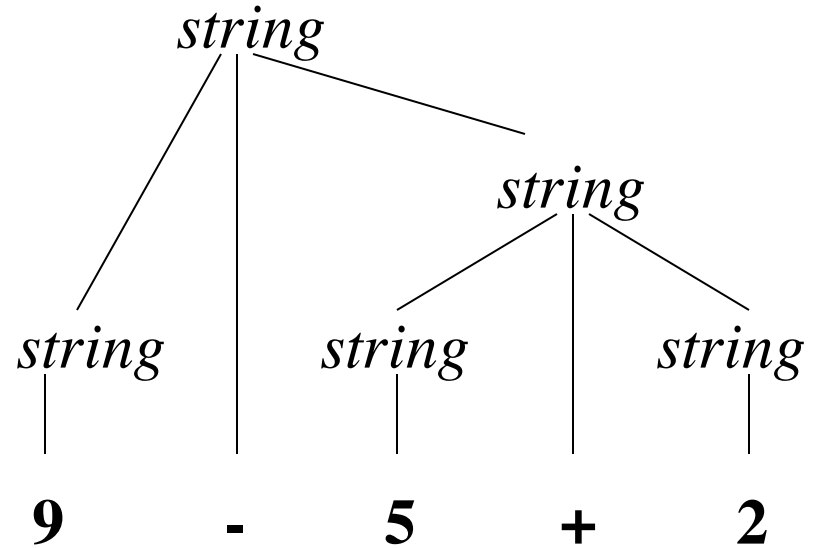
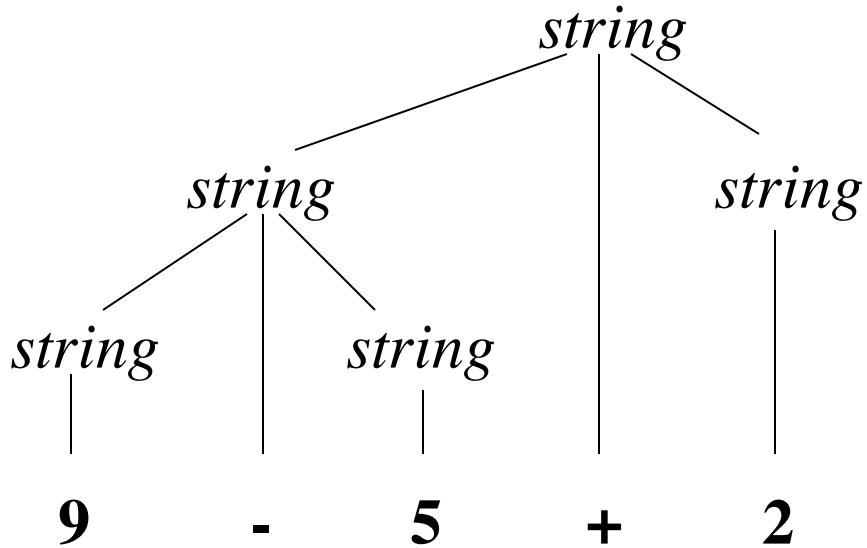
$$G = \langle \{string\}, \{+, -, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}, P, string \rangle$$

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

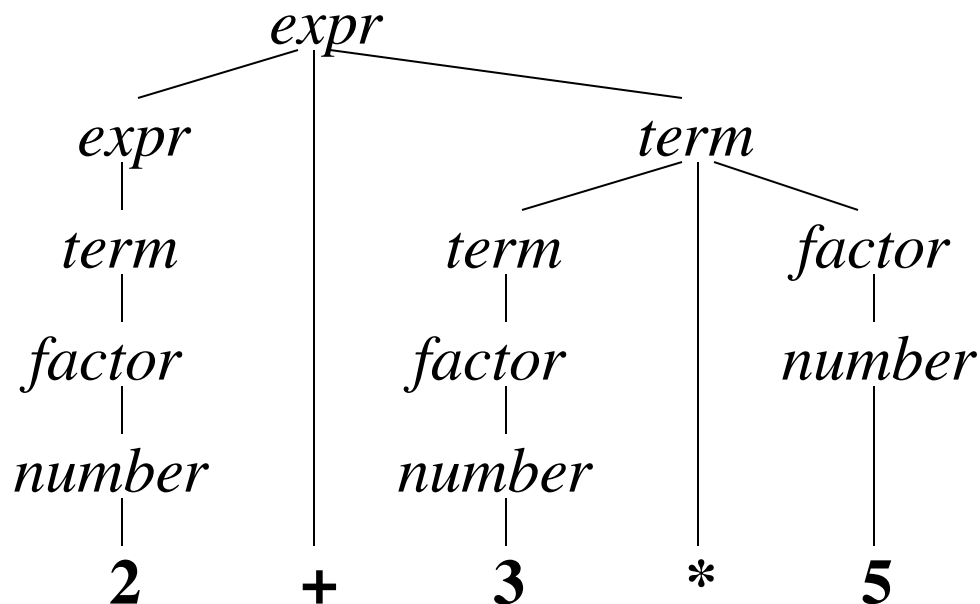
Operators with higher precedence “bind more tightly”

$expr \rightarrow expr + term \mid term$

$term \rightarrow term * factor \mid factor$

$factor \rightarrow number \mid (expr)$

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

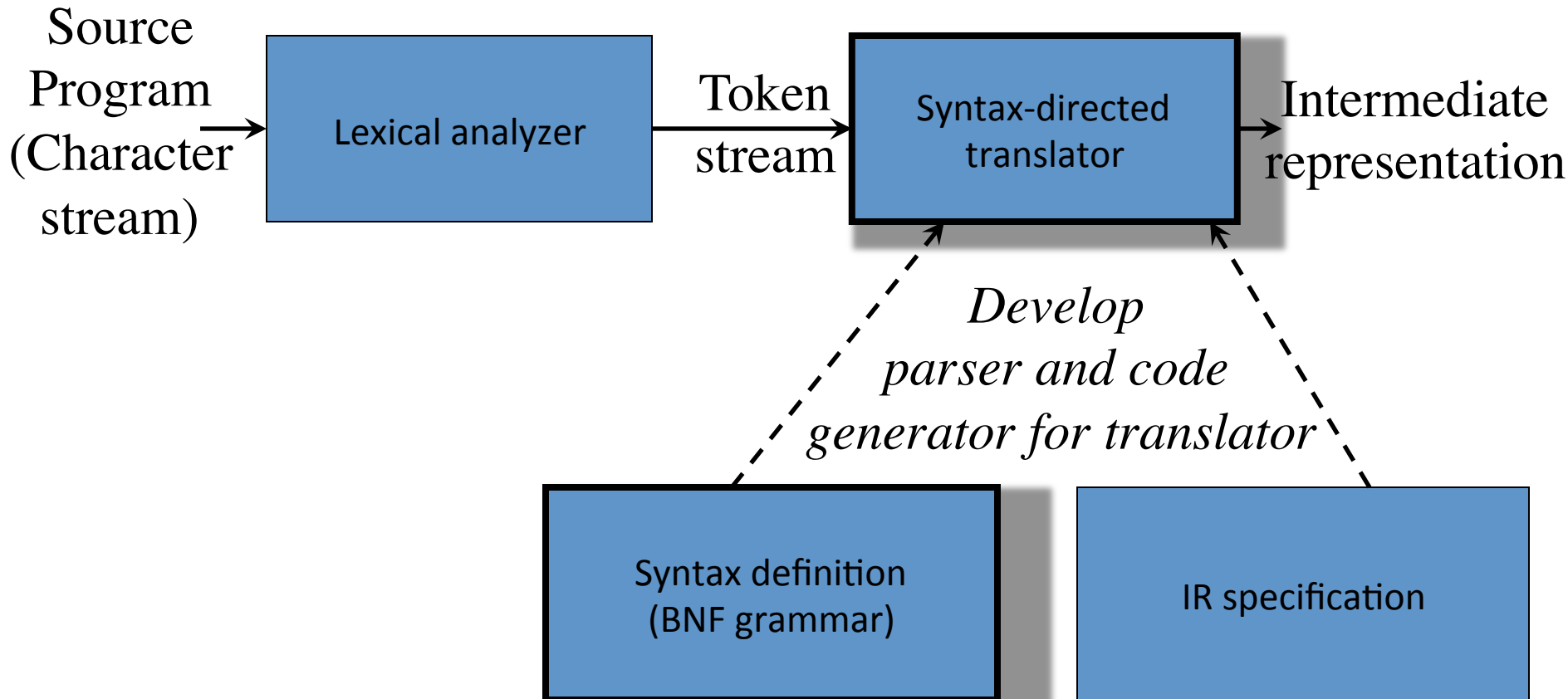


Syntax of Statements

$stmt \rightarrow id := expr$
| **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 Context Free 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

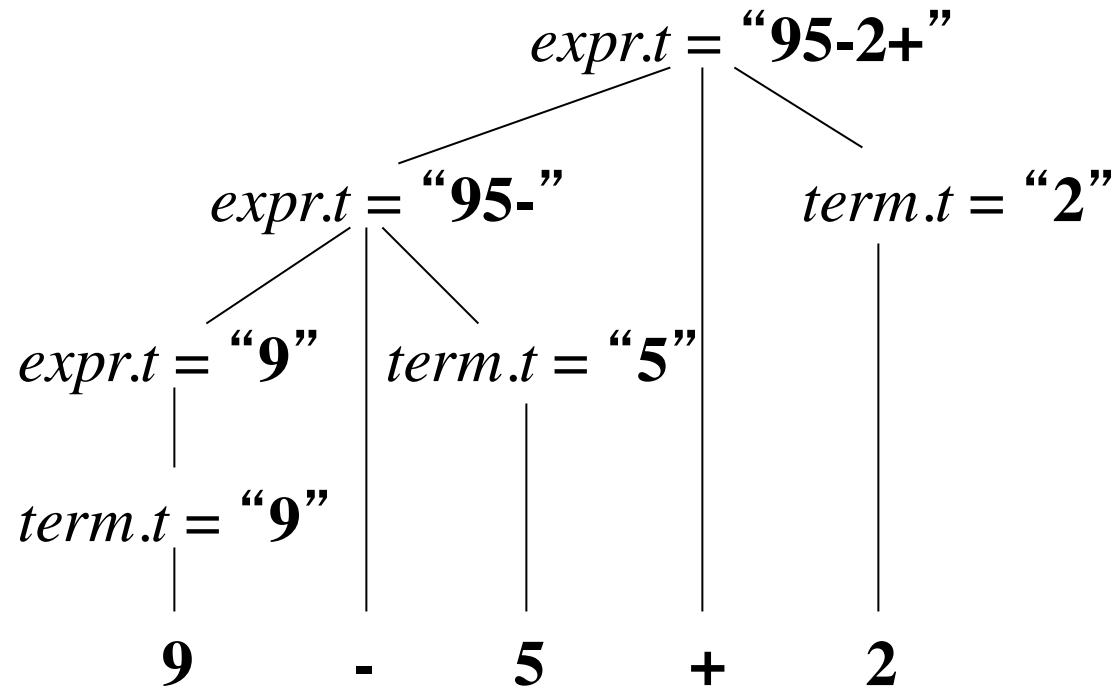
- An attribute is said to be ...
 - *synthesized* if its value at a parse-tree node is determined from the attribute values at the children of the node
 - *inherited* if its value at a parse-tree node is determined by the parent (by enforcing the parent's semantic rules)

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

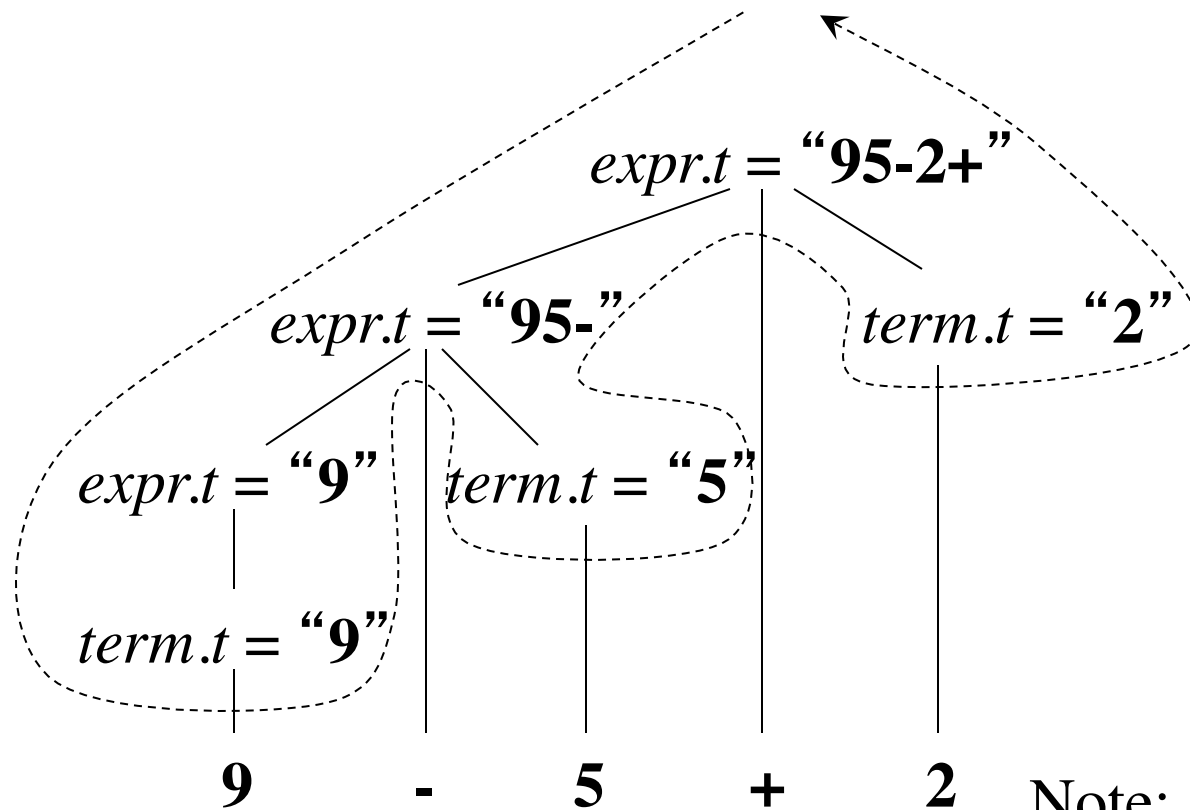
Example Annotated Parse Tree



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)




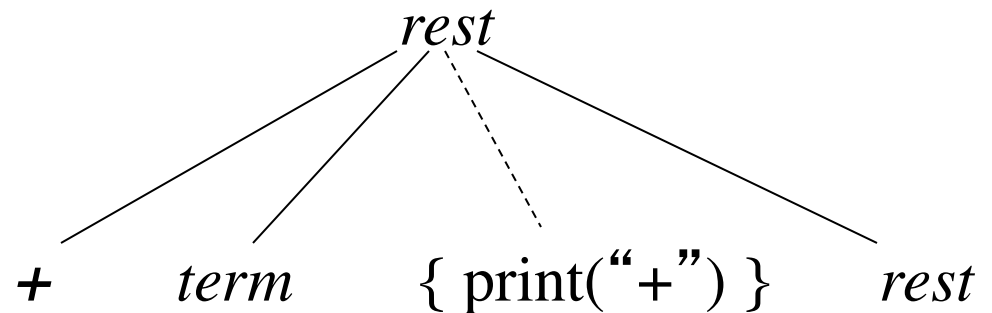
Note: all attributes are of the synthesized₂₆ type

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(“-”) }

$expr \rightarrow term$

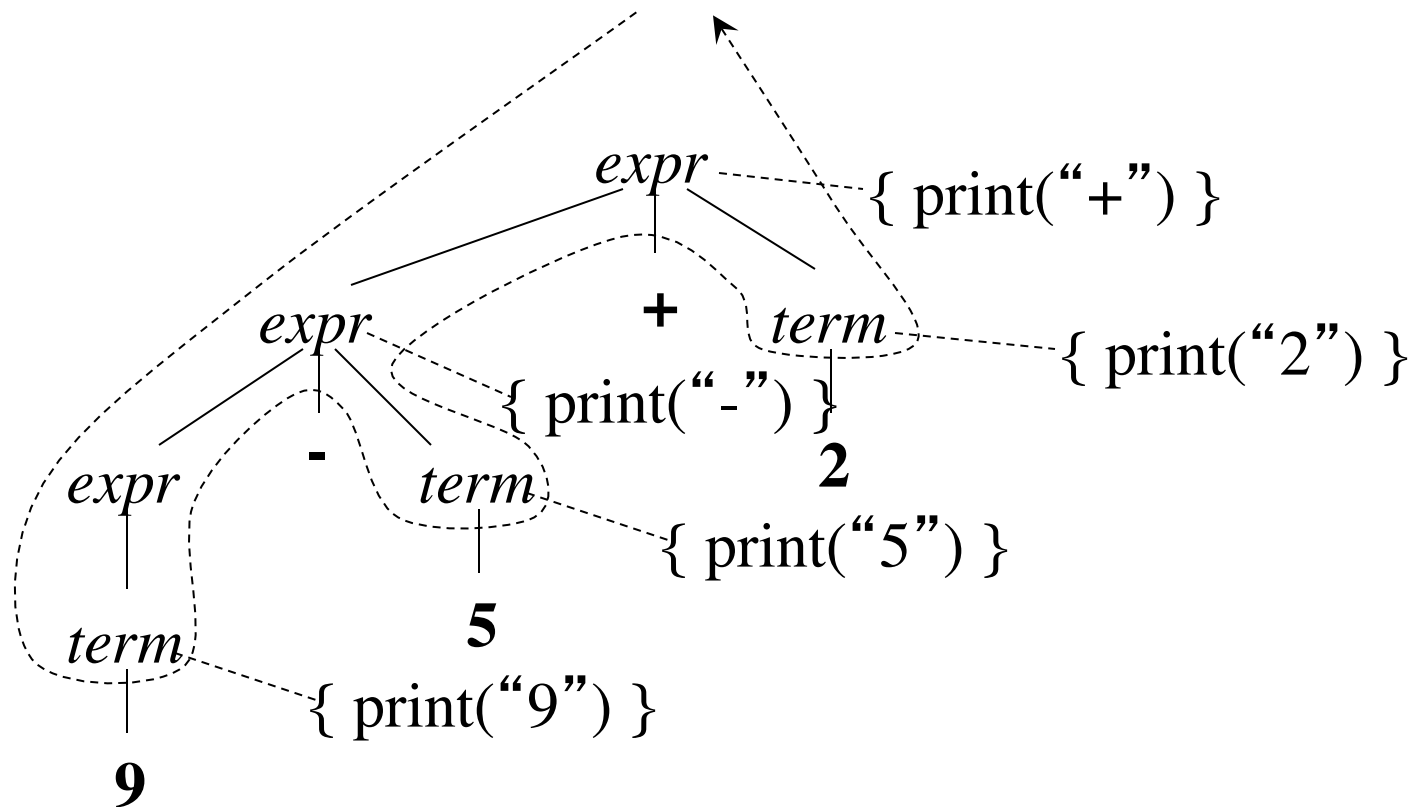
$term \rightarrow 0$ { print(“0”) }

$term \rightarrow 1$ { 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

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

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

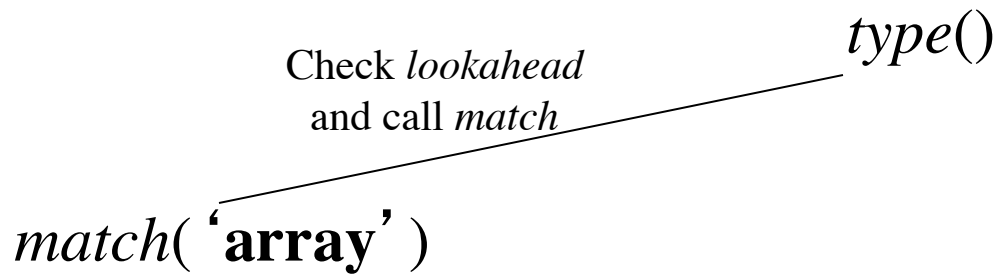
```
procedure simple();  
begin  
  if lookahead = 'integer' then  
    match('integer')  
  else if lookahead = 'char' then  
    match('char')  
  else if lookahead = 'num' then  
    match('num');  
    match('dotdot');  
    match('num')  
  else error()  
end;
```


Example Predictive Parser (Execution Step 1)

match(**'array'**)

Check *lookahead*
and call *match*

type()



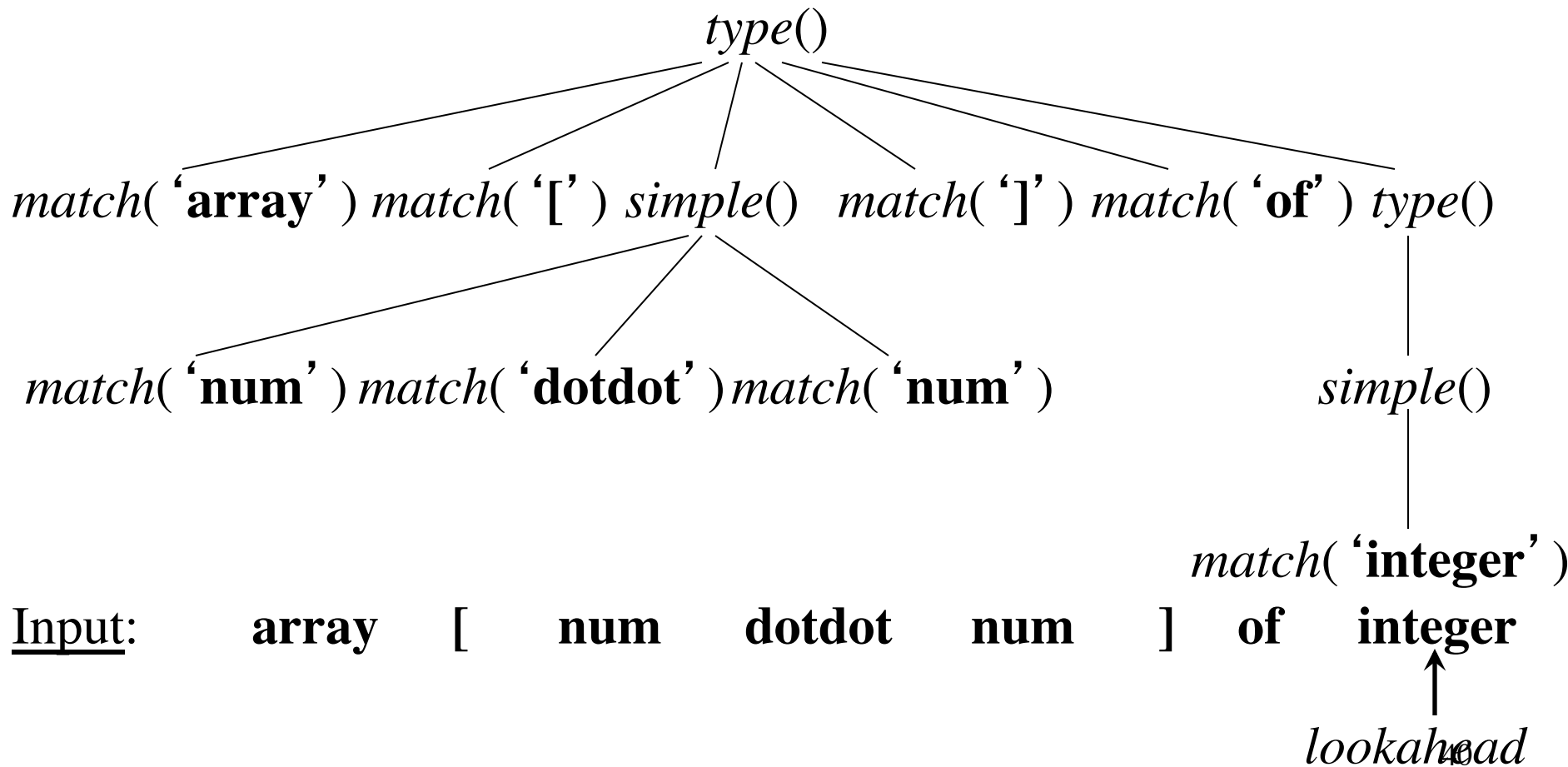
Input: **array** [**num** **dotdot** **num**] **of** **integer**

 ↑

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

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

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

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

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*

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