

Chapter 3

Context-Free Grammars and Parsing

2022 Spring&Summer

Outline

- Context-Free Grammars (CFGs)
- Derivations
- Parse trees and abstract syntax
- Ambiguous grammars

Where We Are

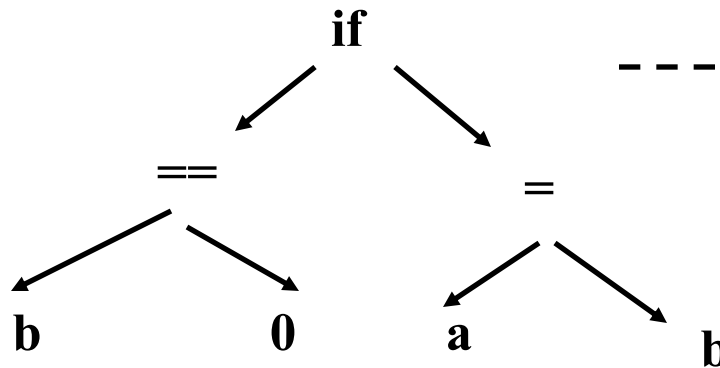
Source code
(character stream)

```
if (b == 0) a = b;
```

**Token
stream**

if	(b	==	0)	a	=	b	;
----	---	---	----	---	---	---	---	---	---

**Abstract Syntax
Tree (AST)**



Lexical Analysis

**Syntax Analysis
(parsing)**

Semantic Analysis

3.1 the parsing process

The parser may be viewed as a *function* that takes as its input the sequence of tokens produced by the scanner and produces as its output the syntax tree.

- 1、 The parser calls a scanner procedure such as *getToken* to fetch the next token from the input
- 2、 In a single-pass compiler the parser will incorporate all the other phases of a compiler,
parse() ;
- 3、 More commonly, a compiler will be multi-pass the further passes will use the syntax tree as their input.

3.1 the parsing process

- The structure of the syntax tree : dependent on the particular syntactic structure of the language.
- This tree is usually defined as a dynamic data structure
each node consists of a record whose fields include the attributes needed for the remainder of the compilation process.

3.1 The parsing process

One problem that is more difficult: the treatment of errors.

1、 Error in the scanner

Generate an error token and consume the offending character.

2、 Error in the parser

Report an error message

Recover from the error and continue parsing (to find as many errors as possible).

Sometimes, a parser may perform **error repair**.

3.2 Context-free grammars

- A context-free grammar is a specification for the syntactic structure of a programming language.
- A context-free grammar involves recursive rules.

For example:

$$exp \rightarrow exp\ op\ exp \mid (exp) \mid \textit{number}$$
$$op \rightarrow + \mid - \mid *$$

3.2.1 Comparison to regular expression notation

The context-free grammar:

$$\textit{exp} \rightarrow \textit{exp op exp} \mid (\textit{exp}) \mid \textit{number}$$
$$\textit{op} \rightarrow + \mid - \mid *$$

The regular expression:

$$\textit{number} = \textit{digit digit}^*$$
$$\textit{digit} = 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

3.2.2 Specification of context-free grammar rules

- A Context-Free Grammar is a 4-tuple $(V, \Sigma, S, \rightarrow)$, where
 - V is a finite set of nonterminal symbols
 - Σ is a finite set of terminal symbols
 - $S \in V$ is a distinguished nonterminal, the start symbol
 - $\rightarrow \subseteq V \times (V \cup \Sigma)^*$ is a finite relation, the productions
- Context Free Grammar is abbreviated CFG

3.2.2 Specification of context-free grammar rules

- The choice operation given by the vertical bar metasymbol is also not necessary in grammar rules.

For example, simple expression grammar could be written as follows:

$$\textit{exp} \rightarrow \textit{exp} \textit{op} \textit{exp}$$
$$\textit{exp} \rightarrow (\textit{exp})$$
$$\textit{exp} \rightarrow \textit{number}$$
$$\textit{op} \rightarrow +$$
$$\textit{op} \rightarrow -$$
$$\textit{op} \rightarrow *$$

3.2.2 Specification of context-free grammar rules

- Use **uppercase** letters for structure names and **lowercase** letters for individual token symbols (which often are just single characters).

$$E \rightarrow E O E \mid (E) \mid n$$

$$O \rightarrow + \mid - \mid *$$

3.2.3 Derivations and the language defined by a grammar

• Let $G = (V, \Sigma, S, \rightarrow)$ be a CFG. The “directly derives” relation (\Rightarrow) is defined as $\{ (\alpha A \gamma, \alpha \beta \gamma) \mid A \rightarrow \beta \}$.

• Example

Let G be the grammar with productions $S \rightarrow aSbS \mid \varepsilon$

$S \Rightarrow aSbS$

$aSbS \Rightarrow aaSbSbS$

$aaSbSbS \Rightarrow aabSbS$

$aabSbS \Rightarrow aabbS$

$aabbS \Rightarrow aabbaSbS$

$aabbaSbS \Rightarrow aabbabS$

$aabbabS \Rightarrow aabbab$

3.2.3 Derivations and the language defined by a grammar

- **Derivation:** Grammar rules determine the legal strings of token symbols by means of derivations.
- A **derivation** is a sequence of replacements of structure names by choices on the right-hand sides(RHS) of grammar rules.
- A **derivation** begins with a single structure name and ends with a string of token symbols.
- At each step in a derivation, a single replacement is made using one choice from a grammar rule.

nonterminal is LHS of production
string is RHS of production

3.2.3 Derivations and the language defined by a grammar

Figure 3.1 : a derivation

• $exp \rightarrow exp\ op\ exp \mid (exp) \mid \textit{number}$

• $op \rightarrow + \mid - \mid *$

(1) $exp \Rightarrow exp\ op\ exp$	$[exp \rightarrow exp\ op\ exp]$
(2) $\Rightarrow exp\ op\ \textit{number}$	$[exp \rightarrow \textit{number}]$
(3) $\Rightarrow exp\ *\ \textit{number}$	$[op \rightarrow *]$
(4) $\Rightarrow (exp)\ *\ \textit{number}$	$[exp \rightarrow (exp)]$
(5) $\Rightarrow \{ exp\ op\ exp \}\ *\ \textit{number}$	$[exp \rightarrow exp\ op\ exp]$
(6) $\Rightarrow (exp\ op\ \textit{number})\ *\ \textit{number}$	$[exp \rightarrow \textit{number}]$
(7) $\Rightarrow (exp - \textit{number})\ *\ \textit{number}$	$[op \rightarrow -]$
(8) $\Rightarrow (\textit{number} - \textit{number})\ *\ \textit{number}$	$[exp \rightarrow \textit{number}]$

3.2.3 Derivations and the language defined by a grammar

- Derivation steps use a different arrow from the arrow metasymbol in the grammar rules.

$$L(G) = \{ s \mid \text{exp} \Rightarrow^* s \}$$

The set of all strings of token symbols obtained by derivations from the *exp* symbol is the language defined by the grammar of expressions.

3.2.3 Derivations and the language defined by a grammar

1. G represents the expression grammar
2. s represents an arbitrary string of token symbols (sometimes called a **sentence**)
3. The symbol $s \Rightarrow^*$ stand for a derivation consisting of a sequence of replacements as described earlier.
4. Grammar rules are sometimes called productions because they "produce" the strings in $L(G)$ via derivations.

3.2.3 Derivations and the language defined by a grammar

- **the start symbol:** the most general structure is listed first in the grammar rules.
- **nonterminals:** structure names are also called nonterminals, since they always must be replaced further on in a derivation .
- **terminals:** symbols in the alphabet are called terminals, since they terminate a derivation. Terminals are usually tokens in compiler applications.

3.2.3 Derivations and the language defined by a grammar

- The grammar for a programming language often defines a structure called *program*
- Language of this structure is the set of all syntactically legal programs of the programming language.
- For example : a BNF for *pascal* language

program \rightarrow *program-heading; program-block*

program-heading \rightarrow

program-block \rightarrow

3.2.3 Derivations and the language defined by a grammar

- Example : Consider the following extremely simplified grammar of statements:

Statement \rightarrow *if-stmt* | **other**

if-stmt \rightarrow **if** (*exp*) *statement*
 | **if** (*exp*) *statement* **else** *statement*

exp \rightarrow 0 | 1

- Examples of strings in this language are

if (0) *other*

if (1) *other*

if (0) *other else other*

if (1) *other else other*

if (0) *if* (0) *other*

if (0) *if* (1) *other else other*

if (1) *other else if* (0) *other else other*

3.2.3 Derivations and the language defined by a grammar

Recursion: the grammar rule $A \rightarrow A a \mid a$

or the grammar rule $A \rightarrow a A \mid a$

- Generates the language $\{a^n \mid n \text{ an integer } \geq 1\}$ (the set of all strings of one or more a's)
- (the regular expression a^+)
- The string *aaaa* can be generated by the first grammar rule with the derivation

$A \Rightarrow Aa \Rightarrow Aaa \Rightarrow Aaaa \Rightarrow aaaa$

3.2.3 Derivations and the language defined by a grammar

- **Left Recursive:** the nonterminal A appears as the first symbol on the right-hand side of the rule defining A .
- **Right Recursive:** the nonterminal A appears as the last symbol on the right-hand side of the rule defining A .

3.2.3 Derivations and the language defined by a grammar

- Consider a rule of the form: $A \rightarrow A\alpha \mid \beta$
where α and β represent arbitrary strings and β does not begin with A .
- This rule generates all strings of the form $\beta, \beta\alpha, \beta\alpha\alpha, \beta\alpha\alpha\alpha, \dots$
- This grammar rule is equivalent in its effect to the regular expression $\beta\alpha^*$.
- $A \rightarrow \alpha A \mid \beta$ generates all strings $\beta, \alpha\beta, \alpha\alpha\beta, \alpha\alpha\alpha\beta, \dots$

3.2.3 Derivations and the language defined by a grammar

- To generate the same language as the regular expression a^* we must have a notation for a grammar rule that generates the empty string (since the regular expression a^* matches the empty string).

$empty \rightarrow$

- Use the *epsilon* metasymbol for the empty string (similar to its use in regular expressions):

$empty \rightarrow \epsilon$

3.2.3 Derivations and the language defined by a grammar

- Such a grammar rule is called an ε -production (an "epsilon production").
- A grammar that generates a language containing the empty string must have at least one ε -production.
- A grammar that is equivalent to the regular expression a^* as
$$A \rightarrow A a \mid \varepsilon \quad \text{or} \quad A \rightarrow a A \mid \varepsilon$$
- Both grammars generate the language
$$\{ a^n \mid n \text{ an integer } \geq 0 \} = L(a^*).$$

3.2.3 Derivations and the language defined by a grammar

- Example :

$statement \rightarrow if\text{-}stmt \mid \textit{other}$

$if\text{-}stmt \rightarrow \mathbf{if} (exp) statement \textit{else-part}$

$else\text{-}part \rightarrow \mathbf{else} statement \mid \varepsilon$

$exp \rightarrow \mathbf{0} \mid \mathbf{1}$

the ε -production indicates that the structure *else-part* is optional.

3.2.3 Derivations and the language defined by a grammar

- Example : $A \rightarrow (A) A \mid \varepsilon$

This grammar generates the strings of all "balanced parentheses."

For example, the string $((() (())) ()$

$$\begin{aligned} A &\Rightarrow (A) A \Rightarrow (A)(A)A \Rightarrow (A)(A) \Rightarrow (A)() \Rightarrow \\ &((A)A)() \Rightarrow (() A) () \Rightarrow (() (A)A) () \Rightarrow (() (A)) () \\ &\Rightarrow (() ((A)A)) () \Rightarrow (() (() A)) () \Rightarrow (() (())) () \end{aligned}$$

3.2.3 Derivations and the language defined by a grammar

- Example : Consider the following grammar G for a sequence of statements:

$stmt\text{-}sequence \rightarrow stmt ; stmt\text{-}sequence \mid stmt$

$stmt \rightarrow s$

- This grammar generates sequences of one or more statements separated by semicolons

$L(G) = \{ s, s;s, s;s;s... \}$

3.2.3 Derivations and the language defined by a grammar

If we want to allow statement sequences to also be empty, we could write the following grammar G' :

$stmt\text{-}sequence \rightarrow stmt ; stmt\text{-}sequence \mid \epsilon$

$stmt \rightarrow s$

$L(G') = \{ \epsilon, s;, s;s;, s;s;s;, \dots \}$

3.2.3 Derivations and the language defined by a grammar

- This turns the semicolon into a **statement terminator** rather than a statement **separator**
- Allow statement sequences to be empty, but retain the semicolon as a statement separator.

stmt-sequence \rightarrow *nonempty-stmt-sequence* $|\epsilon$

nonempty-stmt-sequence \rightarrow *stmt*; *nonempty-stmt-sequence* $|$
stmt

stmt \rightarrow **s**

- Care must be taken in the placement of the ϵ -production when constructing optional structures.

3.3 Parse trees and abstract syntax trees

- Derivations do not uniquely represent the structure of the strings they construct.
- There are many derivations for the same string.
 - the string of tokens *(number - number) * number* from our simple expression grammar using the derivation in Figure 3.1.
 - A second derivation for this string is given in Figure 3.2. (Page 107)

3.3.1 Parse trees

- | | |
|--|-------------------------------------|
| (1) $exp \Rightarrow exp\ op\ exp$ | $[exp \rightarrow exp\ op\ exp]$ |
| (2) $\Rightarrow exp\ op\ \mathbf{number}$ | $[exp \rightarrow \mathbf{number}]$ |
| (3) $\Rightarrow exp\ *\ \mathbf{number}$ | $[op \rightarrow *]$ |
| (4) $\Rightarrow (exp) *\ \mathbf{number}$ | $[exp \rightarrow (exp)]$ |
| (5) $\Rightarrow (exp\ op\ exp) *\ \mathbf{number}$ | $[exp \rightarrow exp\ op\ exp]$ |
| (6) $\Rightarrow (exp\ op\ number) *\ \mathbf{number}$ | $[exp \rightarrow \mathbf{number}]$ |
| (7) $\Rightarrow (exp - number) *\ \mathbf{number}$ | $[op \rightarrow -]$ |
| (8) $\Rightarrow (number - number) *\ \mathbf{number}$ | $[exp \rightarrow \mathbf{number}]$ |

3.3.1 Parse trees

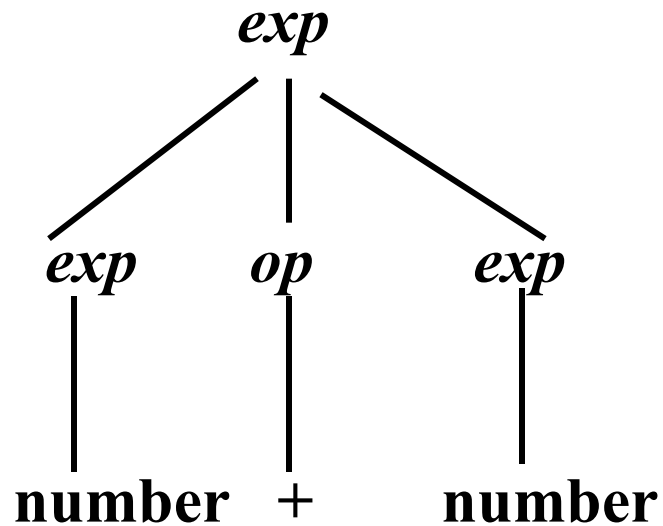
A **parse tree** corresponding to a derivation is a labeled tree

- the interior nodes are labeled by **nonterminals**,
- the leaf nodes are labeled by **terminals**
- the children of each internal node represent the replacement of the associated nonterminal in one step of the derivation.

.

3.3.1 Parse trees

corresponds to the parse tree



3.3.1 Parse trees

The above parse tree corresponds to the three derivations:

1 (1) $exp \Rightarrow exp\ op\ exp$
 (2) \Rightarrow ***number*** $op\ exp$
 (3) \Rightarrow ***number*** $+ exp$
 (4) \Rightarrow ***number*** $+ number$

leftmost derivation

2 (1) $exp \Rightarrow exp\ op\ exp$
 (2) $\Rightarrow exp\ op$ ***number***
 (3) $\Rightarrow exp +$ ***number***
 (4) \Rightarrow ***number*** $+ number$

rightmost derivation

3 (1) $exp \Rightarrow exp\ op\ exp$
 (2) $\Rightarrow exp + exp$
 (3) \Rightarrow ***number*** $+ exp$
 (4) \Rightarrow ***number*** $+ number$

neither leftmost nor
rightmost derivation

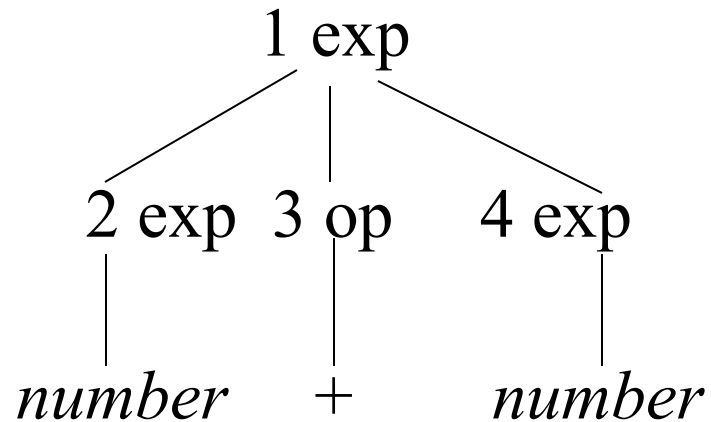
3.3.1 Parse trees

It is possible to distinguish particular derivations that are uniquely associated with the parse tree.

- **A leftmost derivation:** a derivation in which the leftmost nonterminal is replaced at each step in the derivation.
 - Corresponds to the preorder numbering of the internal nodes of its associated parse tree.
- **A rightmost derivation:** a derivation in which the rightmost nonterminal is replaced at each step in the derivation.
 - Corresponds to the postorder numbering of the internal nodes of its associated parse tree.

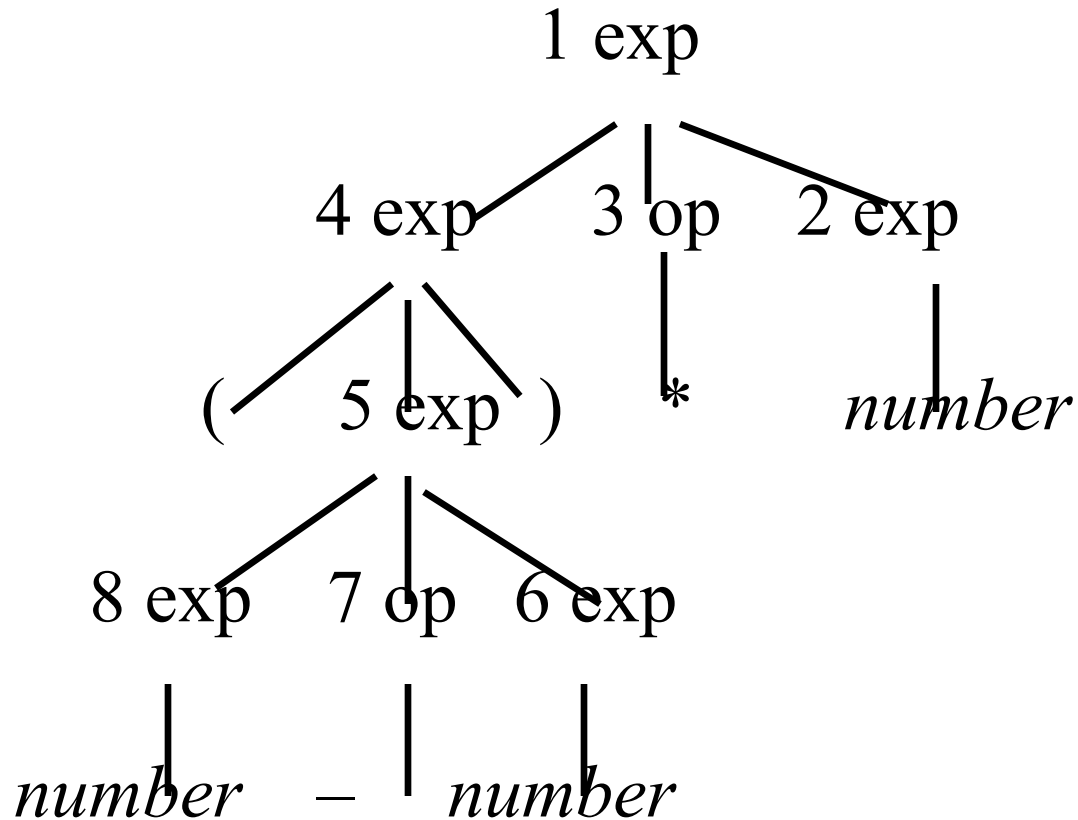
3.3.1 Parse trees

The parse tree corresponds the first derivation



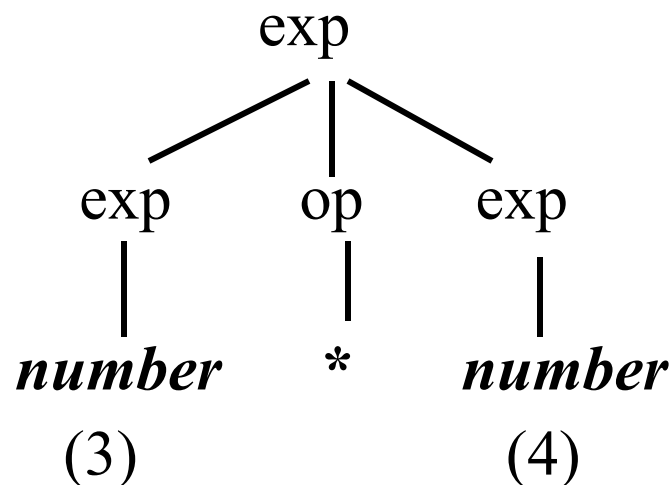
3.3.1 Parse trees

The parse tree for the arithmetic expression $(34-3)*42$



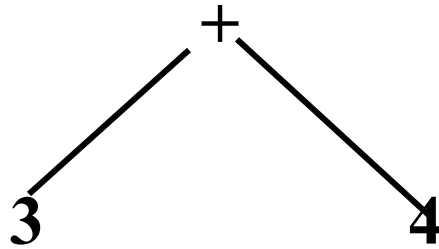
3.3.2 Abstract syntax trees

- The principle of syntax-directed translation states that the meaning, or semantics, of the string $3+4$ should be directly related to its syntactic structure as represented by the parse tree.
- The principle of syntax-directed translation means that the parse tree should imply that the value 3 and the value 4 are to be added.



3.3.2 Abstract syntax trees

A much simpler way to represent this same information, namely, as the tree



3.3.2 Abstract syntax trees

- Abstract Syntax Trees, or Syntax Trees:
 1. Such trees represent abstractions of the actual source code token sequences, and *the token sequences cannot be recovered from them (unlike parse trees)*.
 2. A parse tree is a representation for the structure of ordinary called **concrete syntax** when comparing it to **abstract syntax**.
 3. **Abstract syntax** can be given a formal definition using a BNF-like notation, just like concrete syntax.

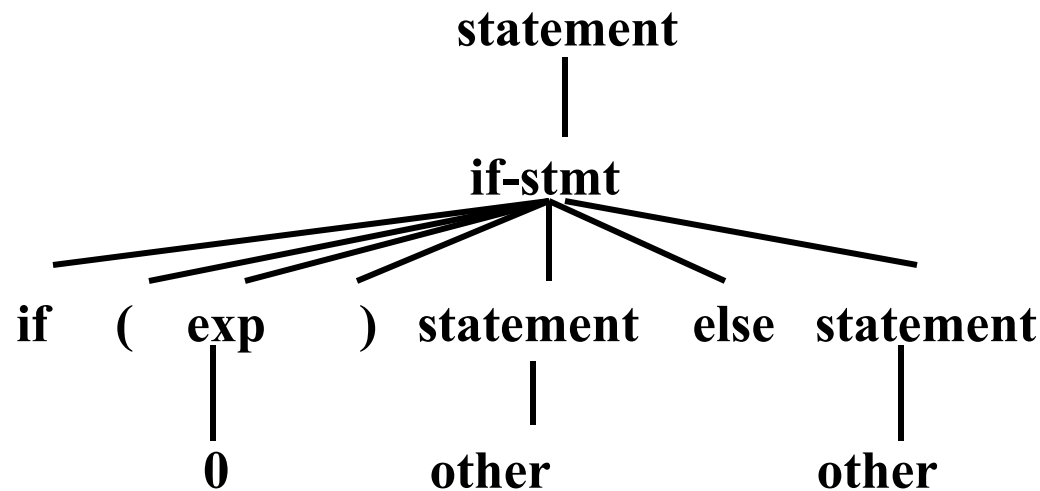
3.3.2 Abstract syntax trees

Example 3.8: Consider the grammar for simplified if-statements of Example 3.4:

$statement \rightarrow if\text{-}stmt \mid \textit{other}$

$if\text{-}stmt \rightarrow \textit{if} (exp) statement$
 $\quad \quad \quad \mid \textit{if} (exp) statement \textit{else} statement$

$exp \rightarrow 0 \mid 1$



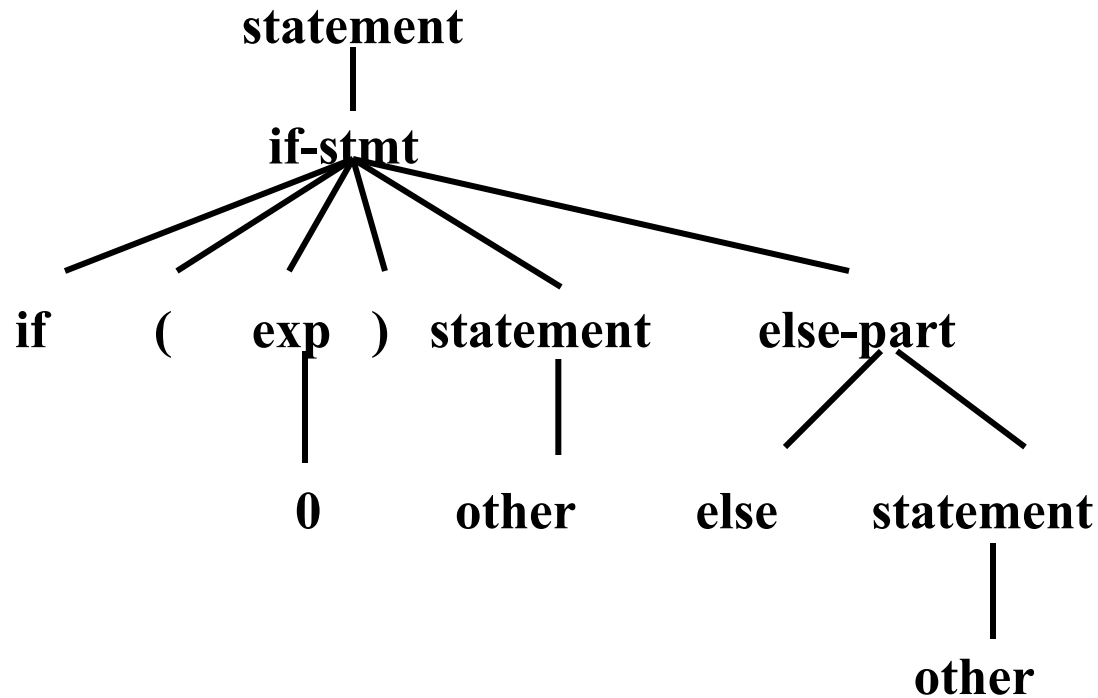
3.3.2 Abstract syntax trees

$statement \rightarrow if\text{-}stmt \mid \textit{other}$

$if\text{-}stmt \rightarrow \textit{if} (exp) statement \textit{else-part}$

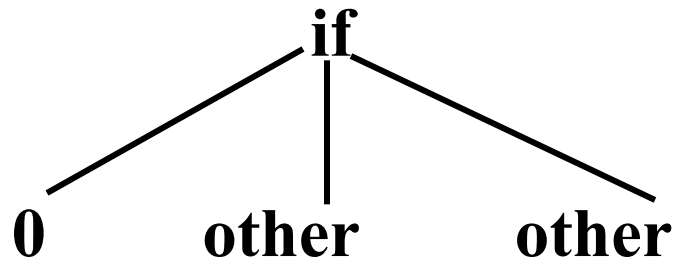
$\textit{else-part} \rightarrow \textit{else} statement \mid \varepsilon$

$exp \rightarrow 0 \mid 1$



3.3.2 Abstract syntax trees

A syntax tree for the previous string:



3.3.2 Abstract syntax trees

A set of C declarations that would be appropriate for the structure of the statements and expressions in this example is as follows:

```
typedef enum {ExpK, StmtK} NodeKind;  
typedef enum {Zero, One} ExpKind;  
typedef enum {IfK, OtherK} StmtKind;  
typedef struct streenode  
    { NodeKind kind;  
    ExpKind ekind;      .  
    StmtKind skind;  
    struct streenode  
        *test, *thenpart, *elsepart;  
    } STreeNode;  
typedef STreeNode * SyntaxTree;
```

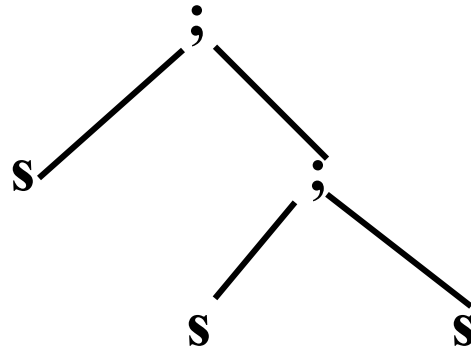
3.3.2 Abstract syntax trees

Example 3.9: Consider the grammar of a sequence of statements separated by semicolons from Example 3.7:

stmt-sequence \rightarrow *stmt* ; *stmt-sequence* | *stmt*

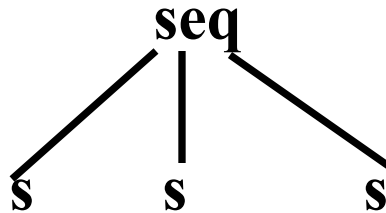
stmt \rightarrow *s*

A possible syntax tree for this same string is:



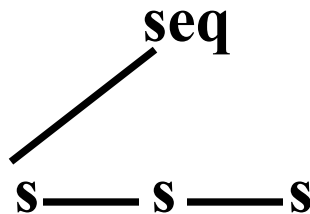
3.3.2 Abstract syntax trees

To bind all the statement nodes in a sequence together with just one node



The problem: a *seq* node may have an arbitrary number of children

The solution: use the standard *leftmost-child right-sibling* representation for a tree (presented in most data structures texts).



3.3.2 Abstract syntax trees

With this arrangement, we can also do away with the connecting *seq* node, and the syntax tree then becomes simply:

s — s — s

Homework of Chapter 3

- 3.2 Given the grammar $A \rightarrow AA \mid (A) \mid \varepsilon$,
- 1). Describe the language it generates,
- 2). Show that it is ambiguous.
- 3.3 Given the grammar
- $\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$
- $\text{addop} \rightarrow + \mid -$
- $\text{term} \rightarrow \text{term mulop factor} \mid \text{factor}$
- $\text{mulop} \rightarrow *$
- $\text{factor} \rightarrow (\text{exp}) \mid \text{number}$
- write down leftmost derivations, parse trees, and abstract syntax trees for the following expressions:
- a. $3+4*5-6$ b. $3*(4-5+6)$ c. $3-(4+5*6)$

Homework of Chapter 3

- 3.4 The following grammar generates all regular expressions over the alphabet of letters (we have used quotes to surround operators, since the vertical bar is an operator as well as a metasyMBOL):
 - $\text{rexp} \rightarrow \text{rexp} \text{ " | " rexp}$
 - | rexp rexp
 - $\text{ | rexp " * "$
 - $\text{ | " (" rexp ") "$
 - | letter
- a. Give a derivation for the regular expression $(ab|b)^*$ using this grammar.
- b. Show that this grammar is ambiguous.
- c. Rewrite this grammar to establish the correct precedences for the operators.
- d. What associativity does your answer in part (c) give to the binary operators? Why?