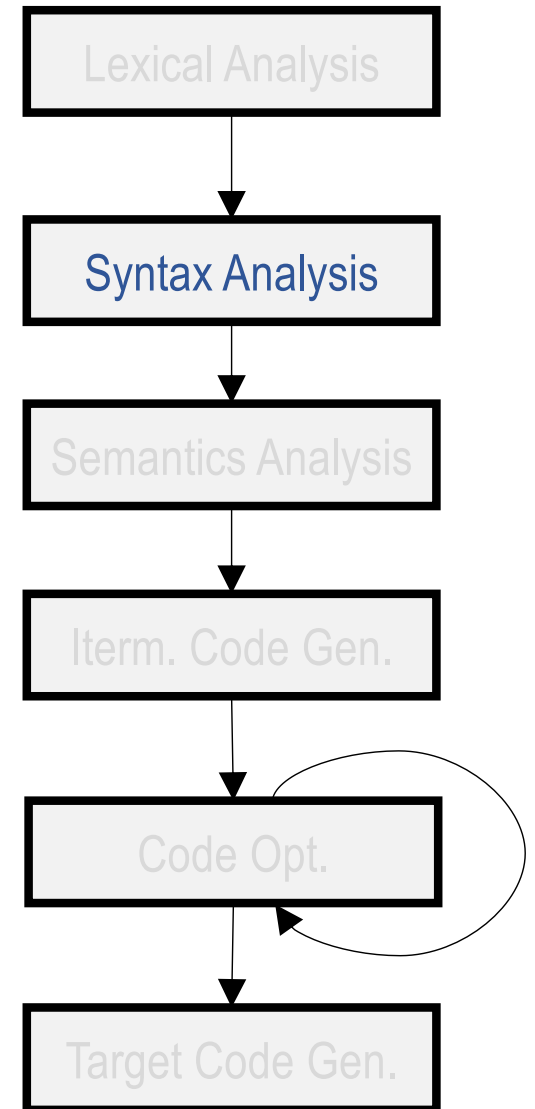


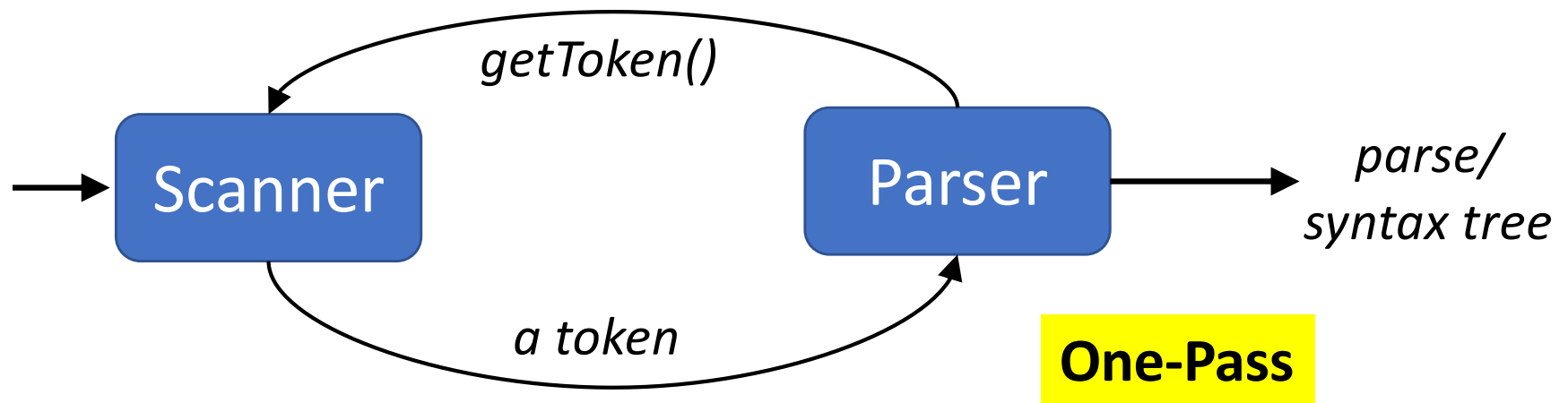
# Syntax Analysis

(Chapters 3-5)



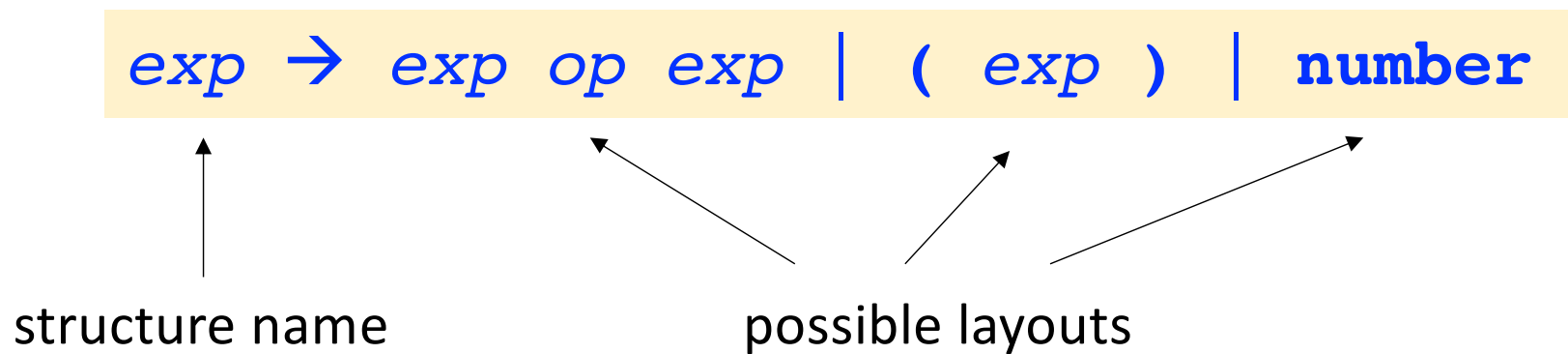
# Parsing

- Determine the syntax (structure) of a program based on the token sequence; Typically, parser drives scanner



# Context-Free Grammars

- A Rule in BNF
  - left-hand side (**LHS**) defines the name of a structure (**nonterminal**)
  - right-hand side (**RHS**) gives its possible layouts (sequence of **terminals** & **nonterminals**)



# Context-Free Grammars

- Derivation

- a sequence of replacements of structure names by layouts on RHS, starting from the **start nonterminal**

Start  
Nonterminal

$exp \Rightarrow exp\ op\ exp$   
 $\Rightarrow exp\ op\ number$   
 $\Rightarrow number\ op\ number$   
 $\Rightarrow number\ * \ number$

(1)  $exp \rightarrow exp\ op\ exp$   
(2)  $exp \rightarrow ( exp )$   
(3)  $exp \rightarrow number$   
(4)  $op \rightarrow +$   
(5)  $op \rightarrow -$   
(6)  $op \rightarrow *$

a legal token string  
with terminals only

CFG

# Context-Free Grammars

- Derivation

- **Leftmost derivation**: the leftmost nonterminal is replaced in each step

$exp \Rightarrow exp\ op\ exp$   
 $\Rightarrow ( exp )\ op\ exp$   
 $\Rightarrow ( exp\ op\ exp )\ op\ exp$   
 $\Rightarrow ( number\ op\ exp )\ op\ exp$   
 $\Rightarrow ( number - exp )\ op\ exp$   
 $\Rightarrow ( number - number )\ op\ exp$   
 $\Rightarrow ( number - number ) * exp$   
 $\Rightarrow ( number - number ) * number$

(1)  $exp \rightarrow exp\ op\ exp$   
(2)  $exp \rightarrow ( exp )$   
(3)  $exp \rightarrow number$   
(4)  $op \rightarrow +$   
(5)  $op \rightarrow -$   
(6)  $op \rightarrow *$

CFG

# Context-Free Grammars

- Derivation

- **Rightmost derivation**: the rightmost nonterminal is replaced in each step

$exp \Rightarrow exp\ op\ exp$   
 $\Rightarrow exp\ op\ \mathbf{number}$   
 $\Rightarrow exp\ * \ \mathbf{number}$   
 $\Rightarrow (\ exp )\ * \ \mathbf{number}$   
 $\Rightarrow (\ exp\ op\ exp )\ * \ \mathbf{number}$   
 $\Rightarrow (\ exp\ op\ \mathbf{number} )\ * \ \mathbf{number}$   
 $\Rightarrow (\ exp\ -\ \mathbf{number} )\ * \ \mathbf{number}$   
 $\Rightarrow (\ \mathbf{number}\ -\ \mathbf{number} )\ * \ \mathbf{number}$

(1)  $exp \rightarrow exp\ op\ exp$   
(2)  $exp \rightarrow (\ exp )$   
(3)  $exp \rightarrow \mathbf{number}$   
(4)  $op \rightarrow +$   
(5)  $op \rightarrow -$   
(6)  $op \rightarrow *$

**CFG**

# Context-Free Grammars

- Language

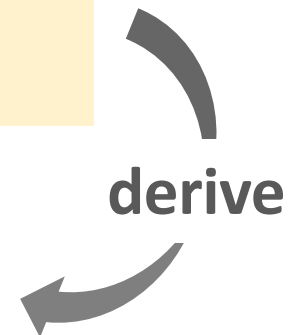
- the set of *token strings* obtained from ***all possible derivations*** is the ***language*** defined by the CFG

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

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

$( \ number \ - \ number \ ) \ * \ number$

a legal token string



# Context-Free Grammars

- Comparison to Regex
  - both use alternation, concatenation, and name
  - no repetition “\*”, but recursion (“→” instead of “=”)
  - more powerful (balanced parentheses  $S \rightarrow (S)S \mid \epsilon$ )
  - terminals are tokens

CFG

```
exp → exp op exp | ( exp ) | number
op  → + | - | *
```

recursion

regex

```
number = digit digit*
digit  = 0|1|2|3|4|5|6|7|8|9
```

tokens



# Context-Free Grammars

- Example
  - Nested if-statments in a C-like form

```
statement → if-stmt | other  
if-stmt → if ( exp ) statement  
          | if ( exp ) statement else statement  
exp → 0 | 1
```

```
other  
if (0) other  
if (1) other  
if (0) other else other  
if (0) if (0) other  
...
```

# Context-Free Grammars

- Example
  - A sequence of statements

*i*  
*s ;*  
*s ; s ;*  
*s ; s ; s ;*  
*...*

*stmt\_sequence*  $\rightarrow$  *stmt ; stmt\_sequence* | *stmt ;*  
*stmt*  $\rightarrow$  *s* |  $\epsilon$

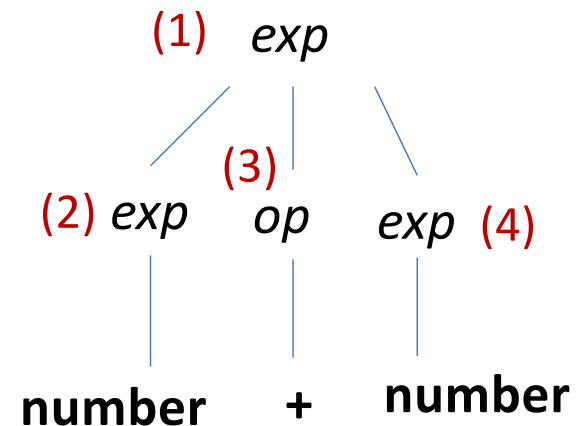
# Parse/Syntax Tree

# Parse Tree

- A labeled tree corresponding to a derivation
  - **Internal nodes:** nonterminals
  - **Leaf nodes:** terminals
  - **Children – parent relation:** a derivation step

(1)  $exp \Rightarrow exp \ op \ exp$   
(2)  $\Rightarrow \mathbf{number} \ op \ exp$   
(3)  $\Rightarrow \mathbf{number} \ + \ exp$   
(4)  $\Rightarrow \mathbf{number} \ + \ \mathbf{number}$

Leftmost Derivation

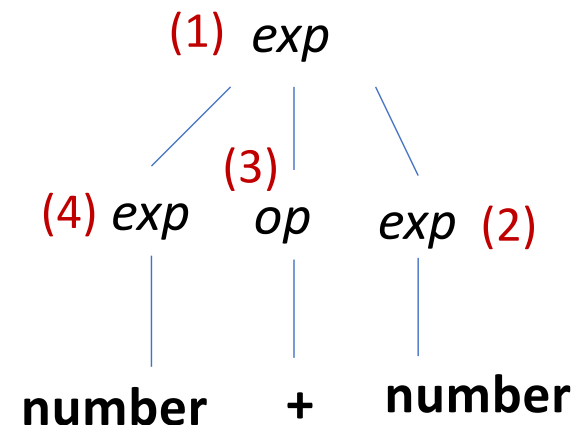


# Parse Tree

- A labeled tree corresponding to a derivation
  - **Internal nodes:** nonterminals
  - **Leaf nodes:** terminals
  - **Children – parent relation:** a derivation step

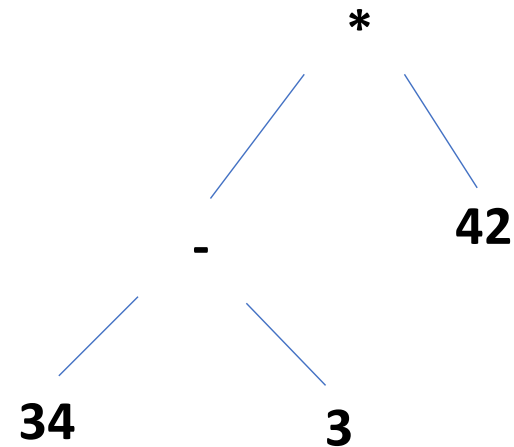
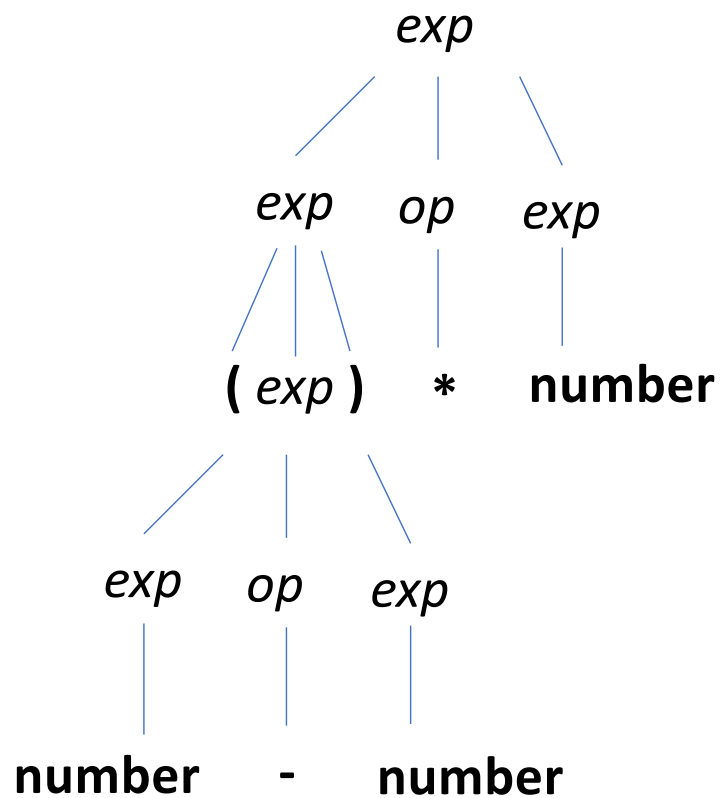
(1)  $exp \Rightarrow exp\ op\ exp$   
(2)  $\Rightarrow exp\ op\ number$   
(3)  $\Rightarrow exp\ +\ number$   
(4)  $\Rightarrow number\ +\ number$

Rightmost Derivation



# Parse Tree

- **Issue:** match derivation well, but not in a concise form



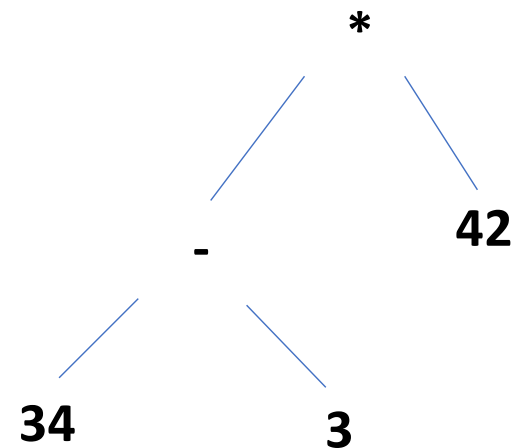
# Abstract Syntax Tree (AST)

- a.k.a. Syntax Tree
  - abstract syntactic structure (may miss some syntax details)
  - token sequence may NOT be recoverable from the AST
  - still contains all necessary info for translation

(34 - 3) \* 42

((34 - 3)) \* 42

((34 - 3) \* 42)



# Abstract Syntax Tree (AST)

- 2nd Example

**if (0) other else other**

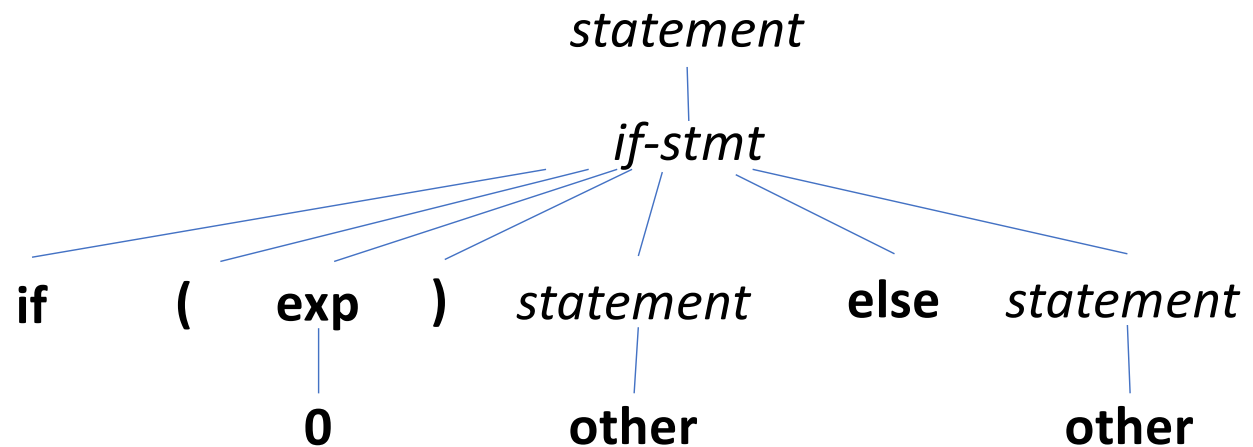
*statement* → *if-stmt* | **other**

*if-stmt* → **if** ( *exp* ) *statement*

          | **if** ( *exp* ) *statement* **else** *statement*

*exp* → **0** | **1**

## Parse Tree





# Abstract Syntax Tree (AST)

- 2nd Example

```
if (0) other else other
```

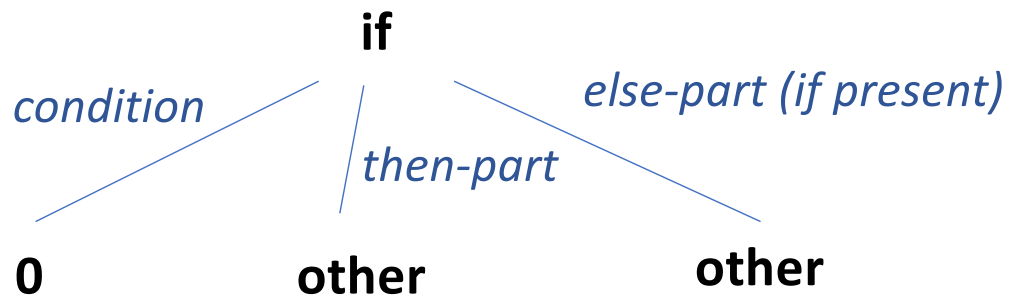
```
statement → if-stmt | other
```

```
if-stmt → if ( exp ) statement
```

```
          | if ( exp ) statement else statement
```

```
exp → 0 | 1
```

**AST**



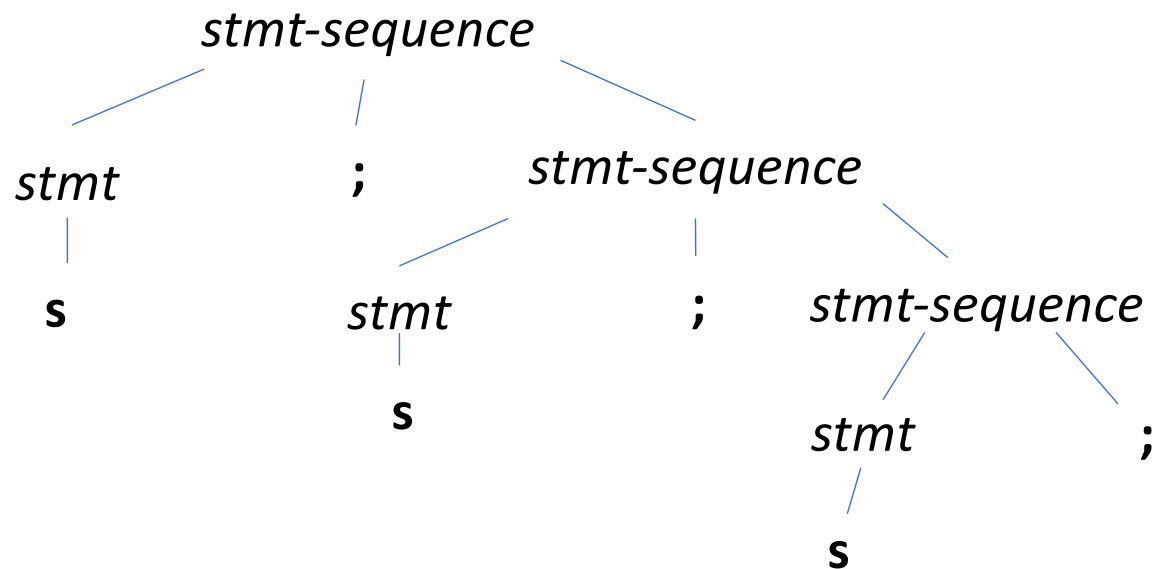
# Abstract Syntax Tree (AST)

- 3rd Example

**s ; s ; s ;**

$stmt\_sequence \rightarrow stmt ; stmt\_sequence \mid stmt ;$   
 $stmt \rightarrow s \mid \epsilon$

**Parse Tree**

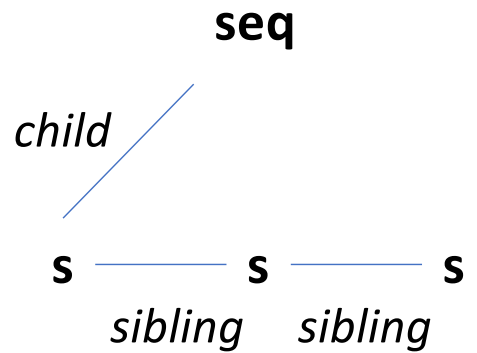
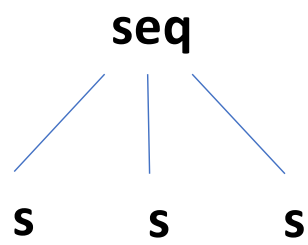
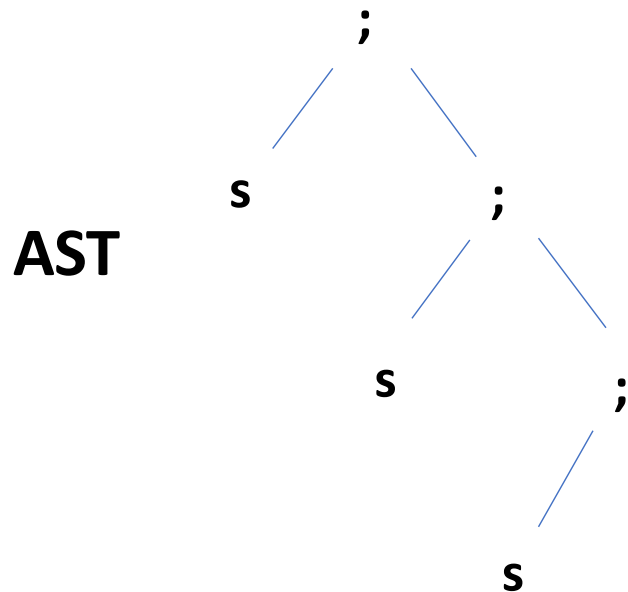


# Abstract Syntax Tree (AST)

- 3rd Example

**s ; s ; s ;**

*stmt\_sequence* → *stmt*; *stmt\_sequence* | *stmt*;  
*stmt* → **s** | ε



**(leftmost-child right-sibling)**

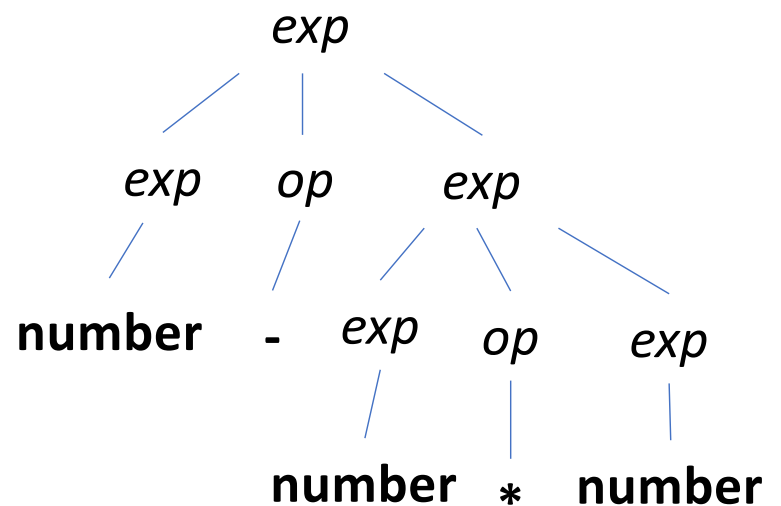
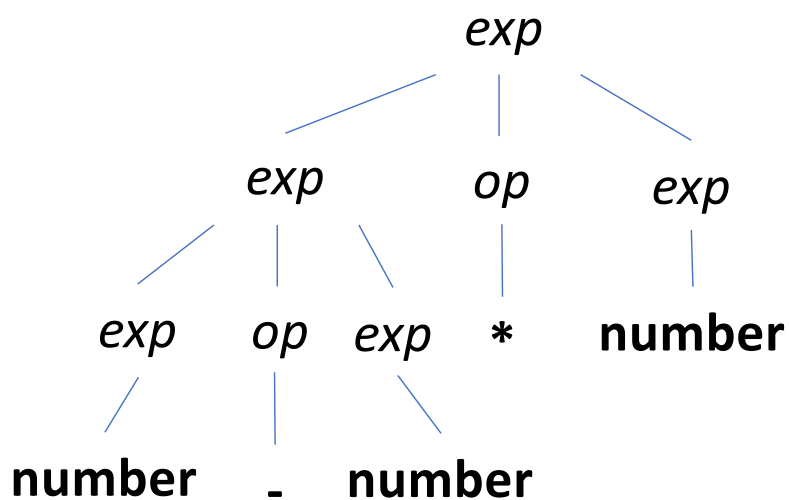
# Ambiguous Grammar

# Ambiguous Grammar

- A grammar allowing a string to have **more than one parse tree**

$exp \rightarrow exp\ op\ exp \mid ( exp ) \mid number$   
 $op \rightarrow + \mid - \mid * \mid /$

34 - 3 \* 42

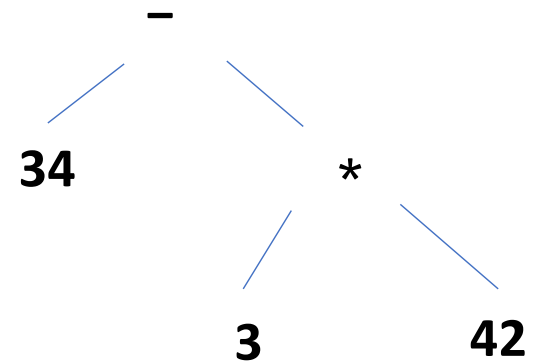
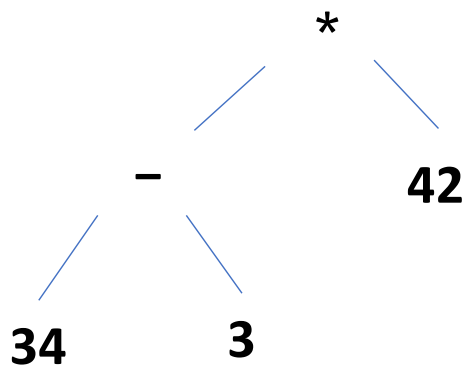


# Ambiguous Grammar

- A grammar allowing a string to have **more than one parse tree**

$exp \rightarrow exp\ op\ exp \mid ( exp ) \mid number$   
 $op \rightarrow + \mid - \mid * \mid /$

34 - 3 \* 42



# Ambiguous Grammar

- Problem – **multiple leftmost / multiple rightmost derivations**

$exp \rightarrow exp\ op\ exp \mid ( exp ) \mid number$   
 $op \rightarrow + \mid - \mid * \mid /$

34 – 3 \* 42

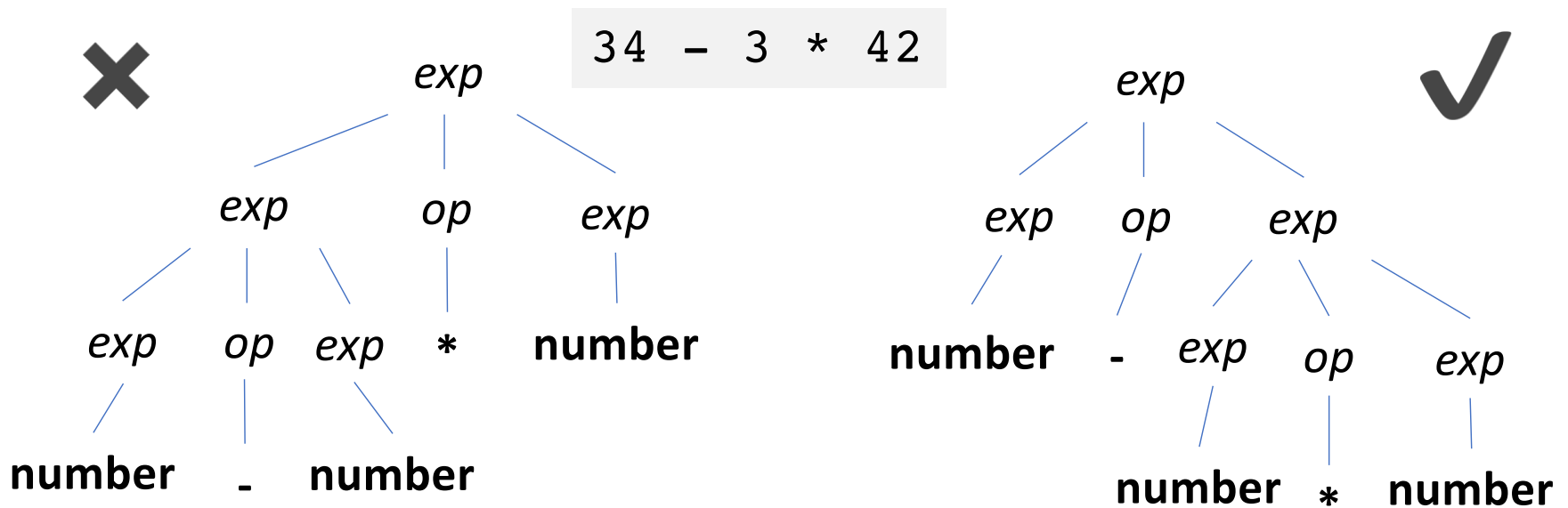
$exp \Rightarrow exp\ op\ exp$   
 $\Rightarrow exp\ op\ exp\ op\ exp$   
 $\Rightarrow number\ op\ exp\ op\ exp$   
 $\Rightarrow number - exp\ op\ exp$   
 $\Rightarrow number - number\ op\ exp$   
 $\Rightarrow number - number * exp$   
 $\Rightarrow number - number * number$

$exp \Rightarrow exp\ op\ exp$   
 $\Rightarrow number\ op\ exp$   
 $\Rightarrow number - exp$   
 $\Rightarrow number - exp\ op\ exp$   
 $\Rightarrow number - number\ op\ exp$   
 $\Rightarrow number - number * exp$   
 $\Rightarrow number - number * number$

**Two leftmost derivations**

# Ambiguous Grammar

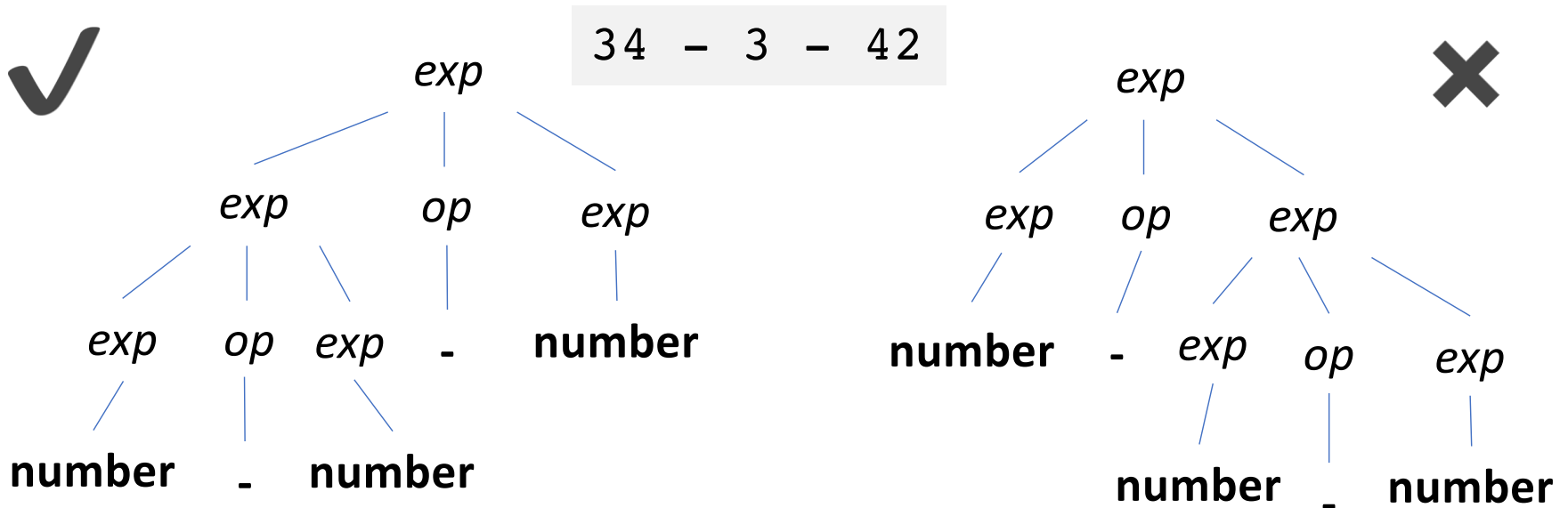
- Disambiguation rule needed → \* has precedence over -





# Ambiguous Grammar

- **Disambiguation rule needed** → - is **left associative**:  
left to right calculation



# Ambiguous Grammar

- Add disambiguation rules to grammars

\* **and** / have **precedence** over - and +

- and / are **left associative**

$exp \rightarrow exp\ op\ exp \mid ( exp ) \mid number$

$op \rightarrow + \mid - \mid * \mid /$

# Ambiguous Grammar

## Precedence Rules

### Rewrite the grammar

group operators of same precedence;

write a different production rule for each precedence group

```
exp → exp op exp | (exp) | number  
op → + | - | * | /
```



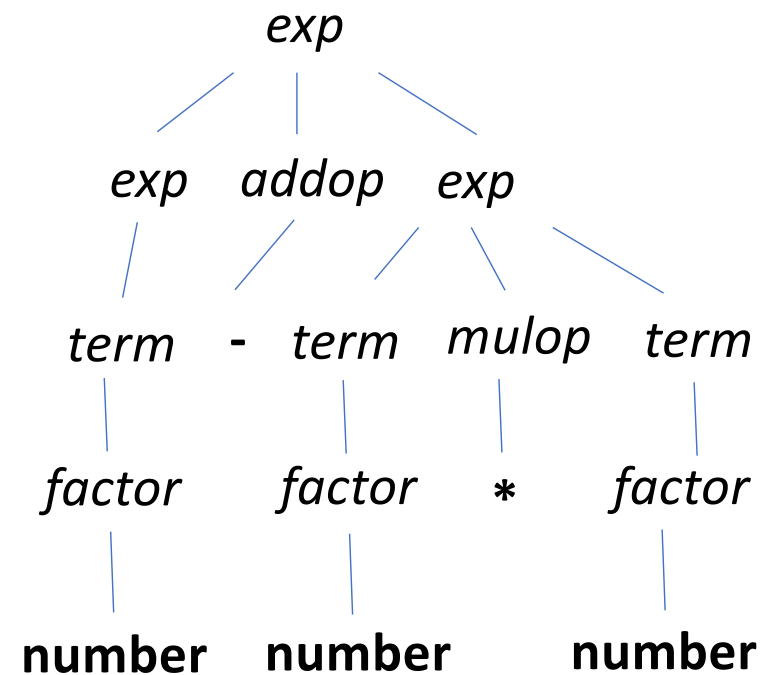
```
exp → exp addop exp | term  
addop → + | -  
term → term mulop term | factor  
mulop → * | /  
factor → (exp) | number
```

# Ambiguous Grammar

- What is the parse tree with the new grammar?

34 - 3 \* 42

$exp \rightarrow exp \text{ addop } exp \mid term$   
 $addop \rightarrow + \mid -$   
 $term \rightarrow term \text{ mulop } term \mid factor$   
 $mulop \rightarrow * \mid /$   
 $factor \rightarrow ( exp ) \mid \mathbf{number}$



# Ambiguous Grammar

- Some ambiguity still exists due to associativity rules

34 - 3 - 42

```
exp → exp addop exp | term
addop → + | -
term → term mulop term | factor
mulop → * | /
factor → ( exp ) | number
```

**recursion occurring on both sides**

# Ambiguous Grammar

- Associativity Constraints

- **Rewrite the grammar**

Replace one of the recursions with the base case

```
exp → exp op exp | ( exp ) | number
op → + | - | * | /
```



```
exp → exp addop term | term
addop → + | -
term → term mulop factor | factor
mulop → * | /
factor → ( exp ) | number
```

# Ambiguous Grammar

- Example parse tree for the new grammar

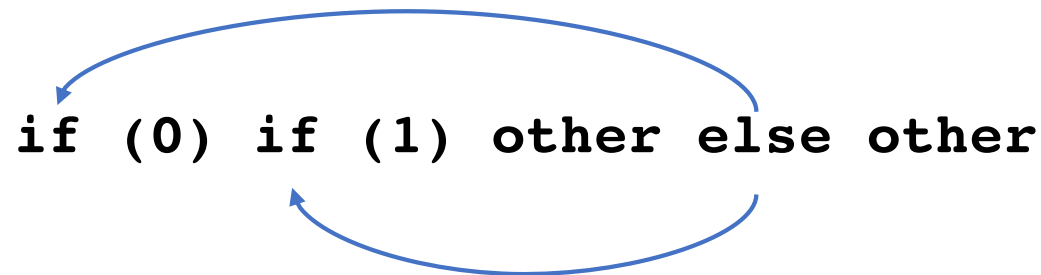
34 - 3 - 42

```
exp → exp addop term | term  
addop → + | -  
term → term mulop factor | factor  
mulop → * | /  
factor → ( exp ) | number
```

# Ambiguous Grammar

- Revisit the if-statments

```
statement → if-stmt | other
if-stmt  → if ( exp ) statement
          | if ( exp ) statement else statement
exp      → 0 | 1
```



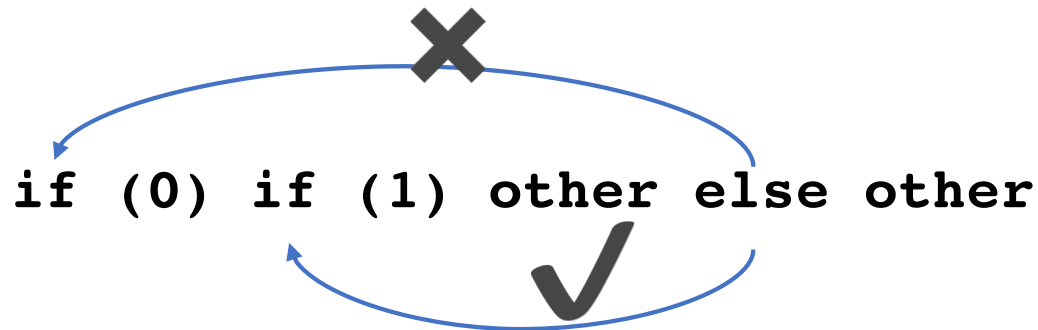
“Dangling Else” Problem



# Ambiguous Grammar

- Revisit the if-statments

```
statement → if-stmt | other
if-stmt → if ( exp ) statement
          | if ( exp ) statement else statement
exp → 0 | 1
```



disambiguation rule: most closely nested rule

# Ambiguous Grammar

- Revisit the if-statements

```
statement → if-stmt | other
if-stmt  → if ( exp ) statement
          | if ( exp ) statement else statement
exp      → 0 | 1
```

```
statement → matched-stmt | unmatched-stmt
matched-stmt → if (exp) matched-stmt else matched-stmt
              | other
unmatched-stmt → if (exp) statement
                | if (exp) matched-stmt else unmatched-stmt
exp           → 0 | 1
```

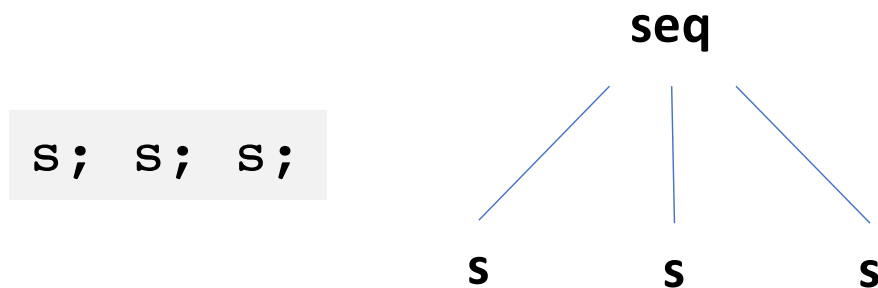
Harder to write (often not used in practice)

# Ambiguous Grammar

- Harmless Ambiguity

- ambiguous grammar may always produce unique AST  
*(but may produce different parse trees)*

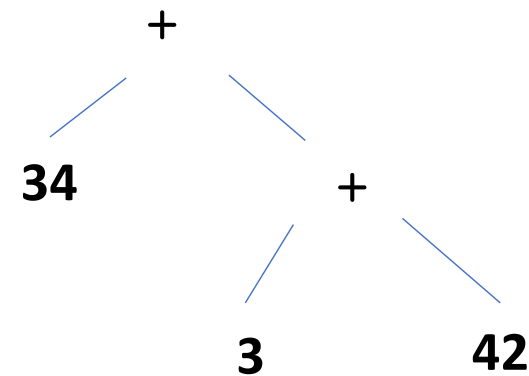
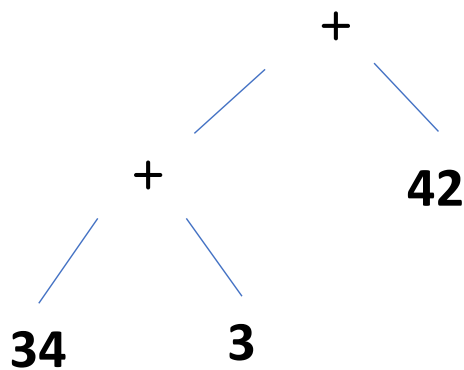
```
stmt_sequence → stmt-sequence ; stmt-sequence | stmt ;  
stmt → s | ε
```



# Ambiguous Grammar

- Harmless Ambiguity
  - even for ambiguous grammar producing different ASTs

34 + 3 + 42



Both carry the same semantics