

# Introduction to Grammar

# Introduction to Grammar

*Every intelligible sentence must follow a grammar*<sup>1</sup>

Grammar describes how the basic building blocks of a language can be used to construct a “valid” sentence in that language

A sentence is said to be **syntactical correct** if it follows the grammar of the language in which the sentence is written. Every programming language has a corresponding grammar/syntax.

---

<sup>1</sup>For Yoda always exceptions are there

# Grammar: Informal vs. Formal

Informal Grammar: collection of rules.

- Subject and verb must agree on the number (singular or plural).
- You may not start a sentence with because because because is a conjunction.

Formal Grammar: specifies the atoms in the language and the rules of composition of these atoms to form a valid sentence in the language.

- Typically does not have any exceptions.
- Typically avoids ambiguity in the application of rules.

# Formal Grammars

- Regular grammar for regular languages: Typical pattern-based searching
- Context-free grammar for context-free languages: Typical programs
- Context-sensitive grammar for context-sensitive languages: Typical counting patterns, scoping in programs
- Unrestricted grammar for recursively enumerable languages

Grammars and expressive power: COM S 331 Theory of Computing.  
*Chomsky Hierarchy.*

# Context-Free Grammar (CFG)

A CFG contains:

- A set of terminals or atoms in the language
- A set of non-terminals
- A set of (production rules) which describes how a non-terminal can be expanded/rewritten to a sequence of terminals and non-terminals

# Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1

## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1, 0..0



## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1, 0..0, 0.0.1

## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1, 0..0, 0.0.1, +-23.1

## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1, 0..0, 0.0.1, +-23.1, 23.+1

## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1, 0..0, 0.0.1, +-23.1, 23.+1, 12.

## Example CFG

- Terminals: 0 1 2 3 4 5 6 7 8 9 . + -
- Non-terminals: real-number, part, digit, sign
- Production rules:
  - A digit is a single terminal except . + and -
  - A part is a sequence of digits
  - A sign is either + or -
  - A real-number is a sign followed by a part optionally followed by a . and another part

### Example

+23.1, 0..0, 0.0.1, +-23.1, 23.+1, 12., 12

# Formal Definition of CFG

A CFG is a tuple  $G = (\Sigma, V, S, P)$ , where

- $\Sigma$  is a set of terminals
- $V$  is a set of non-terminals such that  $\Sigma \cap V = \emptyset$
- $S \in V$  is a start non-terminal
- $P$  is a set of product rules, each of the form:  $X \longrightarrow \omega$ , such that  $X \in V$  and  $\omega \in (\Sigma \cup V)^+$

# Formal Definition of CFG: Real number example

A CFG for real number is a tuple  $G = (\Sigma, V, S, P)$ , where

- $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, .\}$
- $V = \{\text{real-number, part, digit, sign}\}$
- $S = \text{real-number}$
- $P$

# Formal Definition of CFG: Real number example

A CFG for real number is a tuple  $G = (\Sigma, V, S, P)$ , where

- $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, .\}$
- $V = \{\text{real-number}, \text{part}, \text{digit}, \text{sign}\}$
- $S = \text{real-number}$
- $P$  (Backus-Naur form/BNF)

$\text{real-number} \longrightarrow \text{sign part} . \text{ part}$   
 $\qquad\qquad\qquad | \text{ sign part}$

$\text{sign} \longrightarrow + \mid -$

$\text{part} \longrightarrow \text{digit} \mid \text{digit part}$

$\text{digit} \longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$



## CFG Cont'd.

A grammar  $G$  generates a string over terminals if there exists a sequence of application of production rules starting from the the start symbol.

$$\begin{aligned} \text{real-number} &\longrightarrow \text{sign part} . \text{ part} \\ &\quad \quad \quad | \text{sign part} \\ \text{sign} &\longrightarrow + \mid - \\ \text{part} &\longrightarrow \text{digit} \mid \text{digit part} \\ \text{digit} &\longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

real-number

## CFG Cont'd.

A grammar  $G$  generates a string over terminals if there exists a sequence of application of production rules starting from the the start symbol.

$$\begin{aligned} \text{real-number} &\longrightarrow \text{sign part} \mid \text{sign part} . \text{ part} \\ \text{sign} &\longrightarrow + \mid - \\ \text{part} &\longrightarrow \text{digit} \mid \text{digit part} \\ \text{digit} &\longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

$\text{real-number} \rightarrow \text{sign part}$

## CFG Cont'd.

A grammar  $G$  generates a string over terminals if there exists a sequence of application of production rules starting from the the start symbol.

$$\begin{aligned} \text{real-number} &\longrightarrow \text{sign part} \mid \text{sign part} . \text{ part} \\ \text{sign} &\longrightarrow + \mid - \\ \text{part} &\longrightarrow \text{digit} \mid \text{digit part} \\ \text{digit} &\longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

$\text{real-number} \rightarrow \text{sign part} \rightarrow +\text{part}$

## CFG Cont'd.

A grammar  $G$  generates a string over terminals if there exists a sequence of application of production rules starting from the the start symbol.

$$\begin{aligned} \text{real-number} &\longrightarrow \text{sign part} \mid \text{sign part} . \text{ part} \\ \text{sign} &\longrightarrow + \mid - \\ \text{part} &\longrightarrow \text{digit} \mid \text{digit part} \\ \text{digit} &\longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

$\text{real-number} \rightarrow \text{sign part} \rightarrow +\text{part} \rightarrow + \text{digit}$

## CFG Cont'd.

A grammar  $G$  generates a string over terminals if there exists a sequence of application of production rules starting from the the start symbol.

$$\text{real-number} \longrightarrow \text{sign part} \mid \text{sign part} . \text{ part}$$

$$\text{sign} \longrightarrow + \mid -$$

$$\text{part} \longrightarrow \text{digit} \mid \text{digit part}$$

$$\text{digit} \longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

$$\text{real-number} \rightarrow \text{sign part} \rightarrow + \text{part} \rightarrow + \text{digit} \rightarrow + 1$$

## CFG Cont'd.

A grammar  $G$  generates a string over terminals if there exists a sequence of application of production rules starting from the the start symbol.

$$\begin{aligned} \text{real-number} &\longrightarrow \text{sign part} . \text{ part} \\ &\quad \quad \quad | \text{sign part} \\ \text{sign} &\longrightarrow + \mid - \\ \text{part} &\longrightarrow \text{digit} \mid \text{digit part} \\ \text{digit} &\longrightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \end{aligned}$$

$\text{real-number} \rightarrow \text{sign part} \rightarrow +\text{part} \rightarrow + \text{digit} \rightarrow + 1$

How many strings can you generate using the above grammar?

# Language

A language of a grammar  $G$  is denoted by  $L(G)$ , which is the set of all strings generated by  $G$ .

# Exercise

Find the pattern of strings generated by the following grammars.

$$① \quad S \longrightarrow aSb \mid ab$$

$$② \quad S \longrightarrow aSa \mid bSb \mid \epsilon$$

$$③ \quad S \longrightarrow SS \mid (S) \mid ()$$

$$④ \quad S \longrightarrow i \ c \ S \ t \ S \mid a$$

$$⑤ \quad S \longrightarrow \text{part} \mid S + S \mid S - S$$



# Derivation

Every programming language has a grammar describing the valid syntax.

Compilers or interpreters use the given grammar to validate the grammatical correctness of a given program as a sequence of strings.

Typical compiler errors are syntactic/grammatical errors.

# Compilation Technique

A program  $P$  is a string  $s$  over term (keywords, symbols, numbers, operators, etc).  $P$  is said to be syntactically correct if and only if  $s$  can be generated/derived from the grammar  $G$  for the language in which the program is written (i.e.,  $s \in L(G)$ ).

# Compilation Technique

A program  $P$  is a string  $s$  over term (keywords, symbols, numbers, operators, etc).  $P$  is said to be syntactically correct if and only if  $s$  can be generated/derived from the grammar  $G$  for the language in which the program is written (i.e.,  $s \in L(G)$ ).

The derivation of  $s$  from  $G$  is a sequence of application of production rules of the grammar.

# Compilation Technique

A program  $P$  is a string  $s$  over term (keywords, symbols, numbers, operators, etc).  $P$  is said to be syntactically correct if and only if  $s$  can be generated/derived from the grammar  $G$  for the language in which the program is written (i.e.,  $s \in L(G)$ ).

The derivation of  $s$  from  $G$  is a sequence of application of production rules of the grammar.

The derivation of  $+1$  as per the grammar for real numbers is  
 $\text{real-number} \rightarrow \text{sign part} \rightarrow +\text{part} \rightarrow + \text{digit} \rightarrow + 1$

Grammars and Compilers: COM S 440 Principles and Practice of Compiling

# Example

$\text{digit} \longrightarrow 0 \mid 1 \mid \dots \mid 9$

$\text{part} \longrightarrow \text{digit} \mid \text{digit part}$

$S \longrightarrow \text{part} \mid S + S \mid S - S \mid S * S \mid S / S$

- $1 + 2$
- $1 + 2 + 3$
- $1 + 2 - 3$
- $1 + 2 * 3$

# Leftmost vs. Rightmost Derivation

- Leftmost derivation: At each derivation point, the leftmost non-terminal is expanded
- Rightmost derivation: At each derivation point, the rightmost non-terminal is expanded

# Leftmost vs. Rightmost Derivation

Derivation of  $1 + 2 + 3$

## Leftmost

---


$$\begin{aligned}
 S &\rightarrow S + S \\
 &\rightarrow \text{part} + S \\
 &\rightarrow \text{digit} + S \\
 &\rightarrow 1 + S \\
 &\rightarrow 1 + S + S \\
 &\rightarrow 1 + \text{part} + S \\
 &\rightarrow 1 + \text{digit} + S \\
 &\rightarrow 1 + 2 + S \\
 &\rightarrow 1 + 2 + \text{part} \\
 &\rightarrow 1 + 2 + \text{digit} \\
 &\rightarrow 1 + 2 + 3
 \end{aligned}$$

## Rightmost

---


$$\begin{aligned}
 S &\rightarrow S + S \\
 &\rightarrow S + \text{part} \\
 &\rightarrow S + \text{digit} \\
 &\rightarrow S + 3 \\
 &\rightarrow S + S + 3 \\
 &\rightarrow S + \text{part} + 3 \\
 &\rightarrow S + \text{digit} + 3 \\
 &\rightarrow S + 2 + 3 \\
 &\rightarrow \text{part} + 2 + 3 \\
 &\rightarrow \text{digit} + 2 + 3 \\
 &\rightarrow 1 + 2 + 3
 \end{aligned}$$

# Leftmost vs. Rightmost Derivation

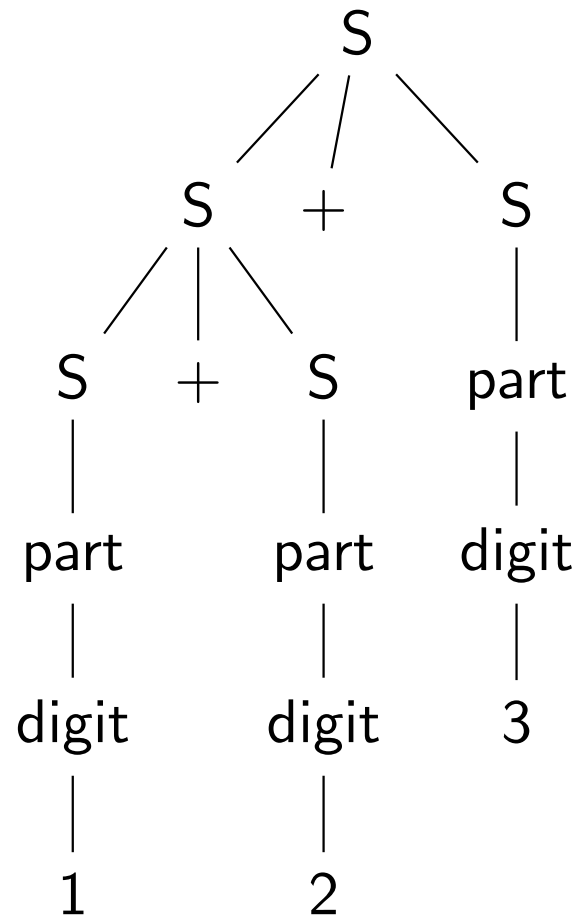


# Parse Tree

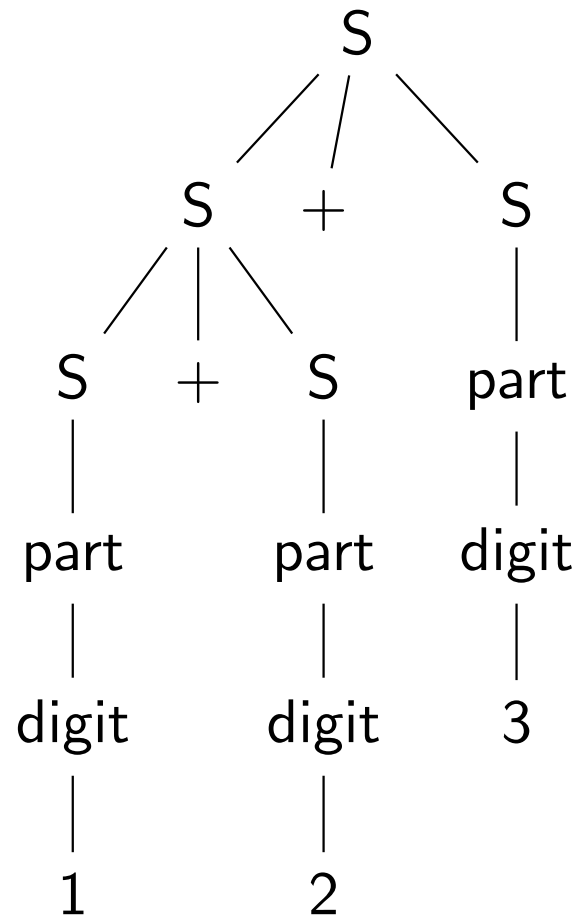
A parse tree results from the derivation sequence.

- Each node in the tree is a terminal or non-terminal in the production rule.
- Each edge in the tree from a non-terminal results from the application of production rule on the non-terminal.
- Application of production rule always result in new nodes in the tree.
- A terminal is a leaf node

# Parse Tree for $1 + 2 + 3$



# Parse Tree for $1 + 2 + 3$



The same parse tree can be generated by left-most and right-most derivation.

# Parse Tree & Semantics

## Semantics

Meaning of a syntactically correct sentence.

- 1 Classify the terminals into atom and operator classes.
- 2 Associate meaning with each atom.
- 3 Associate meaning with the application of operator(s) on atom(s).

Evaluate the semantics using parse tree.

# Example Semantics

$\text{digit} \longrightarrow 0 \mid 1 \mid \dots \mid 9$

$\text{part} \longrightarrow \text{digit} \mid \text{digitpart}$

$S \longrightarrow \text{part} \mid S + S \mid S - S \mid S * S \mid S / S$

Atoms: 0 ... 9 and their sequences (numbers).

- Meaning of a number  $n$  can be  $n$  points.

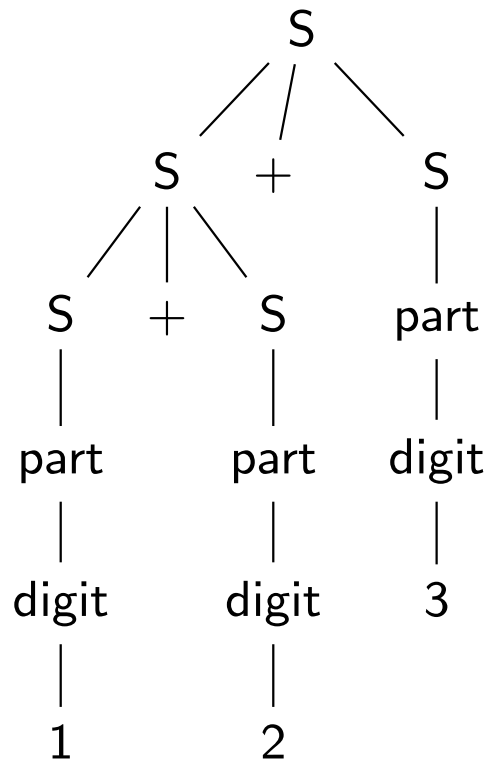
Operators:  $+$ ,  $-$ ,  $*$ ,

- Meaning of application of  $+$  on two numbers  $m$  and  $n$  can be placing  $m$  points besides  $n$  points.

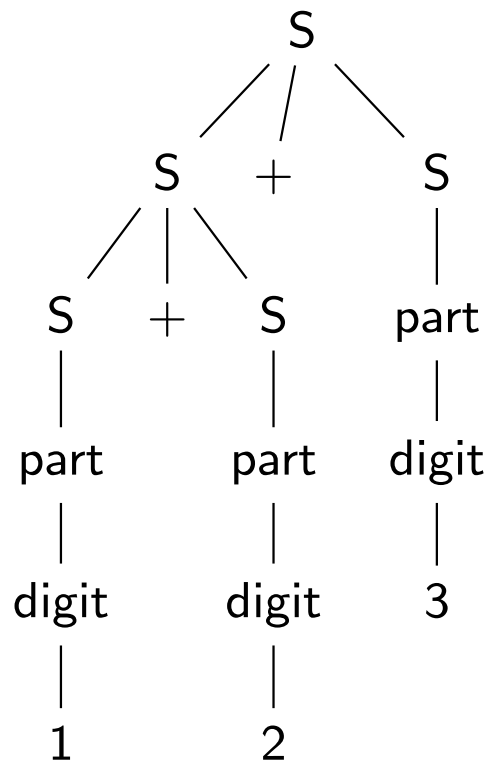
# Parse Tree and Semantics

Evaluate the semantics using parse tree.

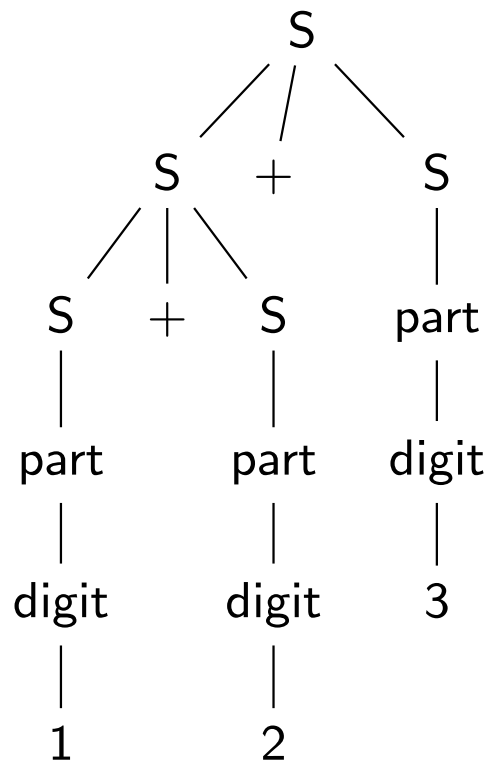
- Start from the leaf-nodes to create the atoms and find their meanings.
- Apply the operators on the generated atoms to obtain the meaning of the application of the operators.
- Continue until you reach the root-node.



# Evaluating Semantics



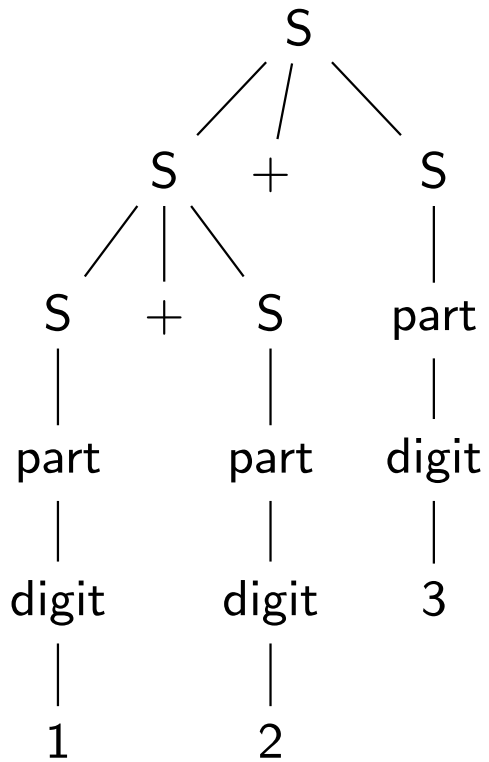
# Evaluating Semantics



sem(1)

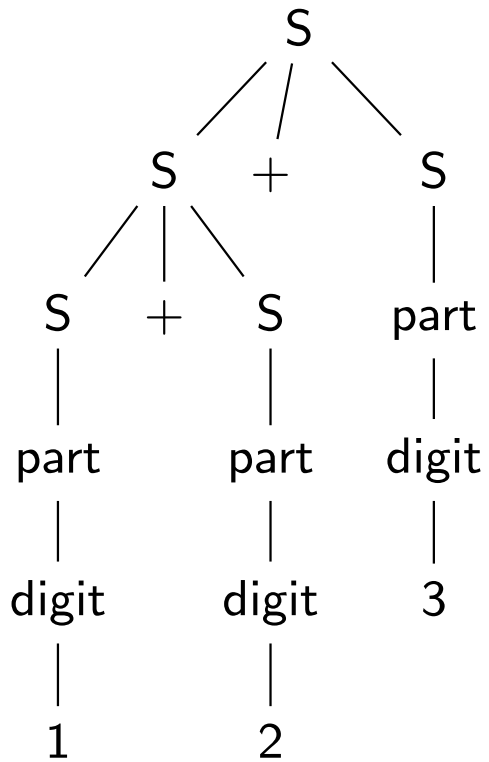


# Evaluating Semantics



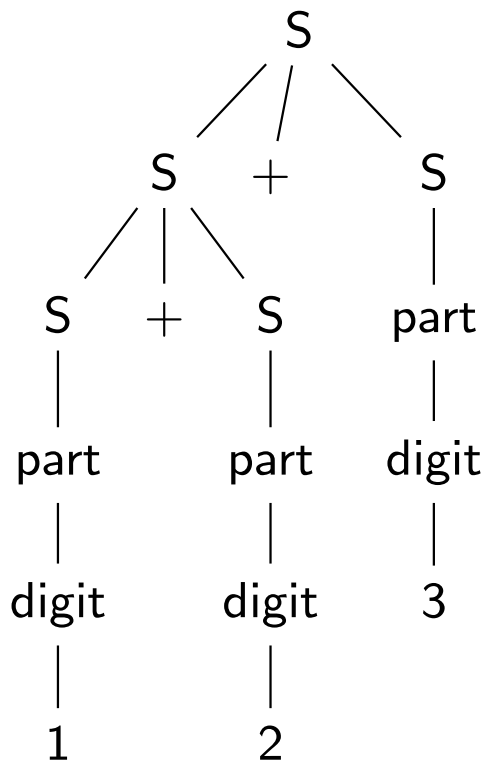
sem(1) sem(2)

# Evaluating Semantics



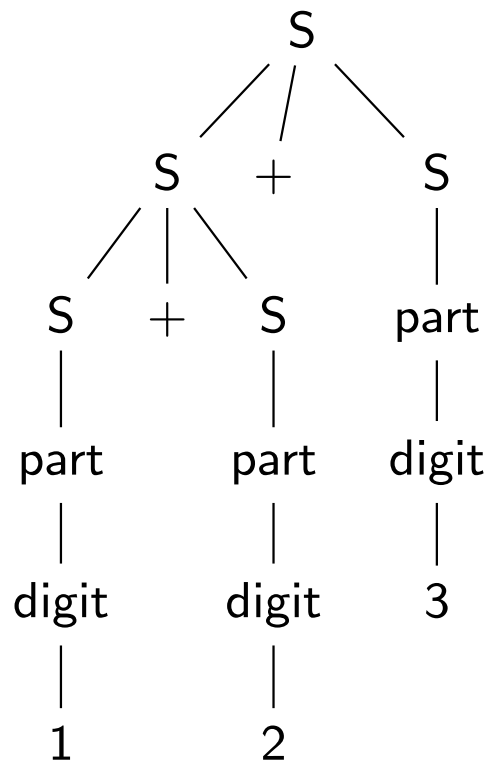
sem(1)   sem(2)   sem(3)

# Evaluating Semantics



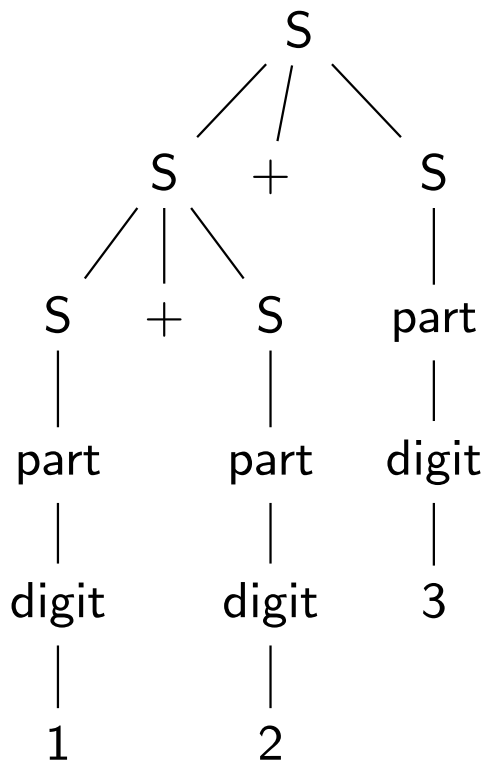
$\text{sem}(+, \text{sem}(1), \text{sem}(2)) \quad \text{sem}(3)$

# Evaluating Semantics



$\text{sem}(+, \text{sem}(+, \text{sem}(1), \text{sem}(2)), \text{sem}(3))$

# Evaluating Semantics



$\text{sem}(+, \text{sem}(+, \text{sem}(1), \text{sem}(2)), \text{sem}(3))$