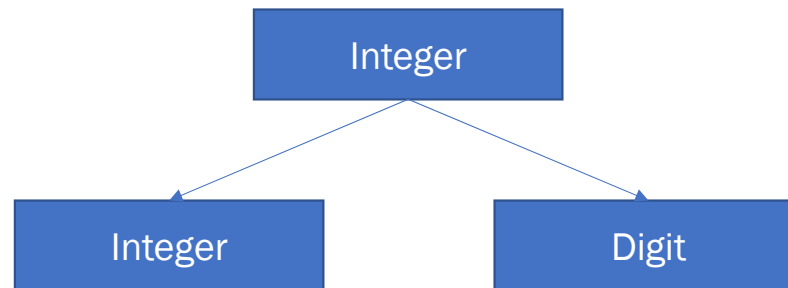Parsing

Professor: Suman Saha

# Parse Trees

- Derivation in graphical form
- The root node always contains the start symbol
- Each internal node has as its direct descendants the elements that appear on the right-hand side grammar
- The leaves of the parse tree are always terminal symbol
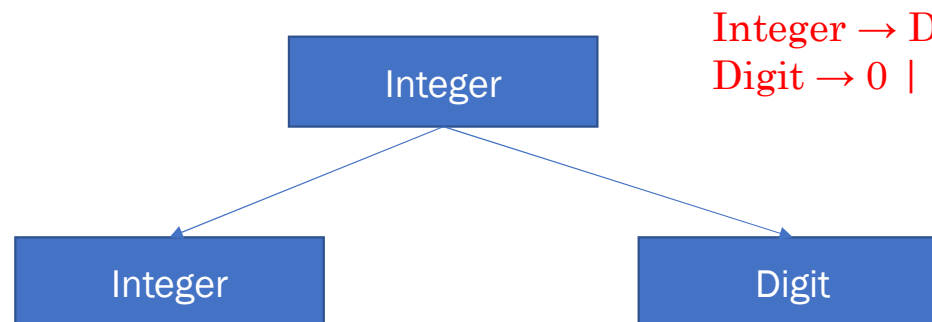
# Parse Trees

- Derivation in graphical form
- The root node always contains the start symbol
- Each internal node has as its direct descendants the elements that appear on the right-hand side grammar
- The leaves of the parse tree are always terminal symbol
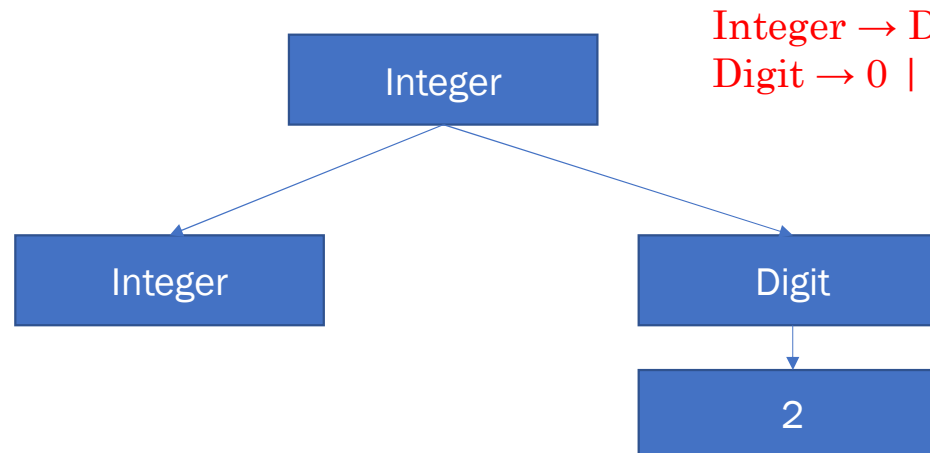
$$Integer \rightarrow Integer\ Digit$$

# Parse Trees

## Parse Tree of 352

Integer → Digit | Integer Digit
Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

## Parse Tree of 352

$$Integer \rightarrow Digit \mid Integer\ Digit$$
$$Digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

```
                    Integer
                   /       \
            Integer         Digit
                              |
                              2
```

# Parse Trees

## Parse Tree of 352

Integer → Digit | Integer Digit
Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

# Parse Trees

## Parse Tree of 352
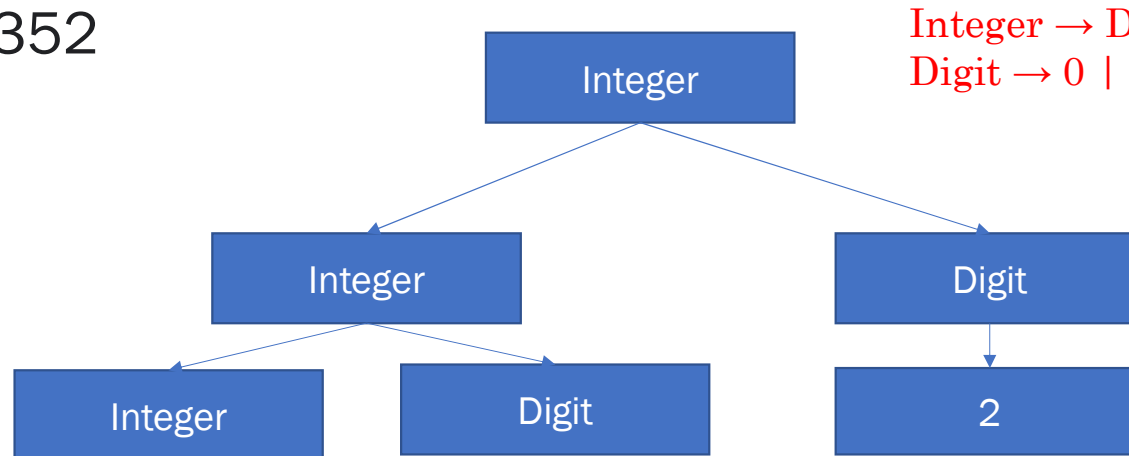
Integer → Digit | Integer Digit
Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
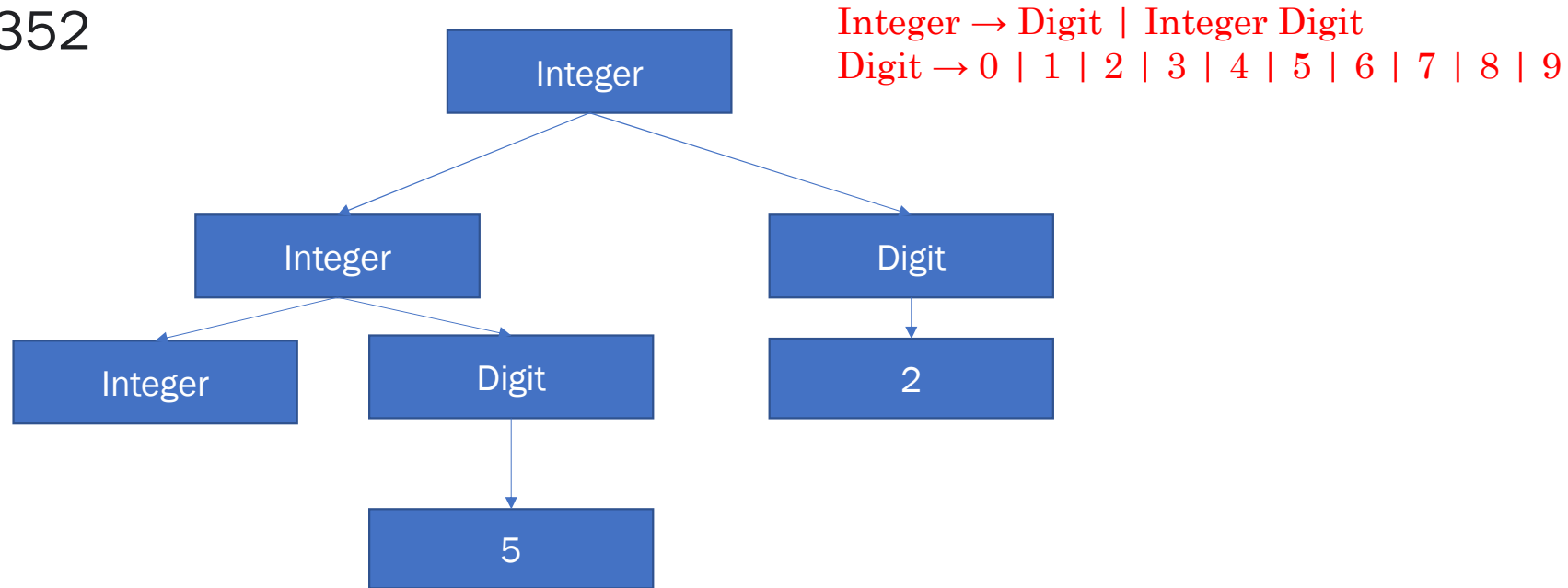
# Parse Trees

## Parse Tree of 352

Integer → Digit | Integer Digit
Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

# Parse Trees

Parse Tree of 352
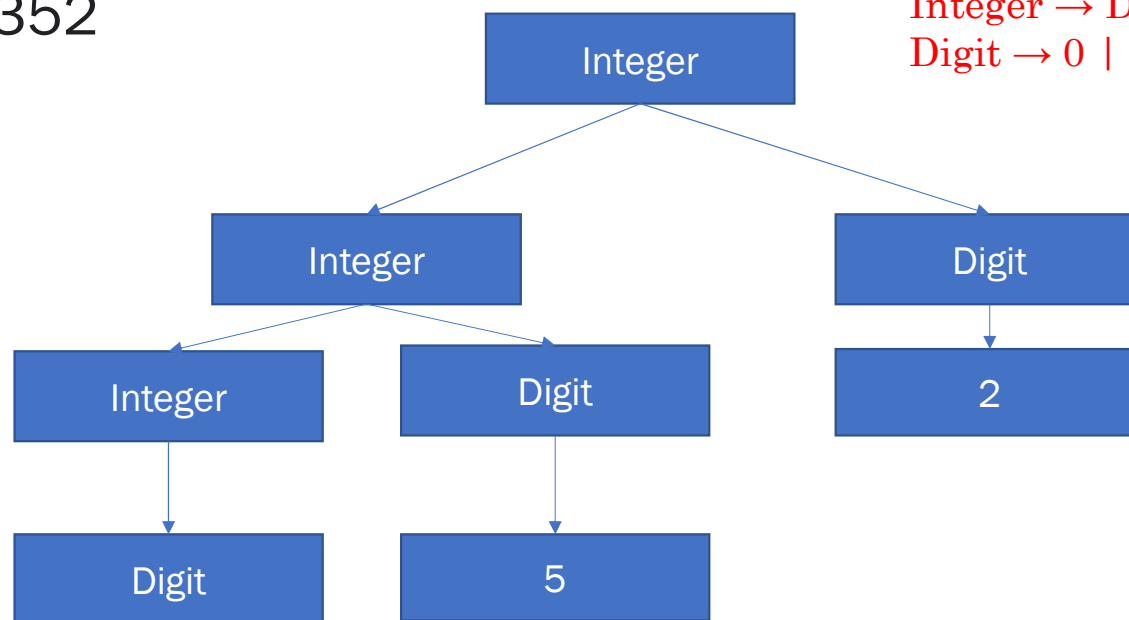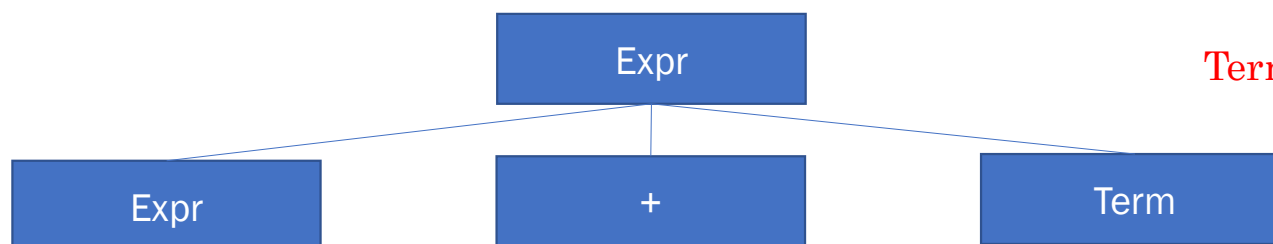
Integer → Digit | Integer Digit
Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- Given grammar

$$
\begin{aligned}
\mathrm{Expr} \rightarrow \quad & \mathrm{Expr} + \mathrm{Term} \\
& | \ \mathrm{Expr} - \mathrm{Term} \\
& | \ \mathrm{Term}
\end{aligned}
$$

$$
\mathrm{Term} \rightarrow \ 0 \ | \ \ldots \ | \ 9 \ | \ (\mathrm{Expr})
$$

# Parse Trees

Parse of the string $5 - 4 + 3$

$$Expr \rightarrow \quad Expr + Term$$
$$| \ Expr - Term$$
$$| \ Term$$
$$Term \rightarrow \quad 0 \ | \ ... \ | \ 9 \ | \ (Expr)$$

Parse of the string $5 - 4 + 3$



$Expr \rightarrow$      $Expr + Term$
         $| \ Expr - Term$
         $| \ Term$

$Term \rightarrow$      $0 \ | \ ... \ | \ 9 \ | \ (Expr)$

# Parse Trees

Parse of the string $5 - 4 + 3$

$$Expr \rightarrow \quad Expr + Term$$
$$| \ Expr - Term$$
$$| \ Term$$
$$Term \rightarrow \quad 0 \ | \ \dots \ | \ 9 \ | \ (Expr)$$

# Parse Trees

Parse of the string $5 - 4 + 3$

$$Expr \rightarrow \quad Expr + Term$$
$$| \ Expr - Term$$
$$| \ Term$$
$$Term \rightarrow \quad 0 \ | \ \ldots \ | \ 9 \ | \ (Expr)$$

# Parse Trees

Parse of the string $5 - 4 + 3$

Expr → Expr + Term
| Expr – Term
| Term

Term → 0 | ... | 9 | (Expr)

# Parse Trees

Parse of the string $5 - 4 + 3$

$Expr \rightarrow \quad Expr + Term$
$| \; Expr - Term$
$| \; Term$

$Term \rightarrow \quad 0 \; | \; \dots \; | \; 9 \; | \; (Expr)$

# Parse Trees

Parse of the string $5 - 4 + 3$

$$Expr \rightarrow Expr + Term$$
$$| \; Expr - Term$$
$$| \; Term$$
$$Term \rightarrow 0 \; | \; \dots \; | \; 9 \; | \; (Expr)$$



$$(5\text{-}4)+3$$
$$\neq$$
$$5 - (4 + 3)$$

# Ambiguous Grammar

- A grammar is ambiguous if, for any string
    - it has more than one parse tree, or
    - there is more than one right-most derivation, or
    - there is more than one left-most derivation

    (the three conditions are equivalent)

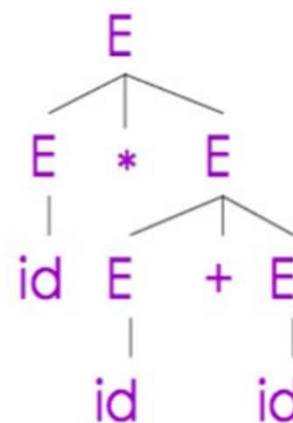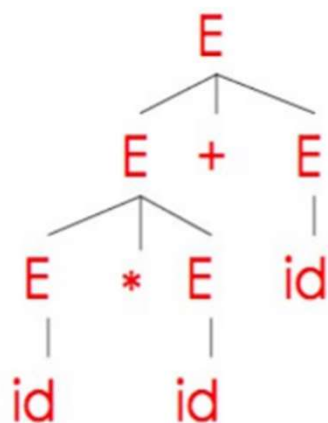# Ambiguous Grammar

- A grammar is ambiguous if, for any string
  - it has more than one parse tree, or
  - there is more than one right-most derivation, or
  - there is more than one left-most derivation

  (the three conditions are equivalent)

- Unambiguous grammar is preferred
- However, ambiguity may tolerable
- Tradeoff between the size of the grammar and the information it is trying to convey

# Ambiguity

- Ambiguity = Program structure is not uniquely defined

- $E \rightarrow E + E \mid E * E \mid (E) \mid id$

- String id * id + id has two parse trees:

Given grammar

$$Stmt \rightarrow \text{if Expr then Stmt}$$
$$| \text{ if Expr then Stmt else Stmt}$$
$$| \text{ other}$$

Example:

if E1 then

if E2 then S1 else S2

**PennState**

Stmt

if · Expr · then · Stmt

if E1 then
  if E2 then S1 else S2

# Dangling Else Problem

```
                    Stmt
       ┌──────┬──────┴──────┐
      if     Expr   then    Stmt
              │
             E1
```
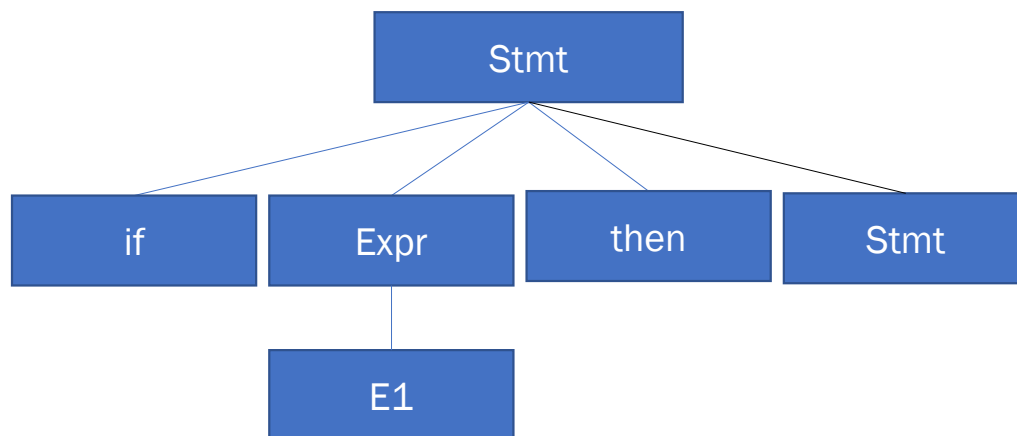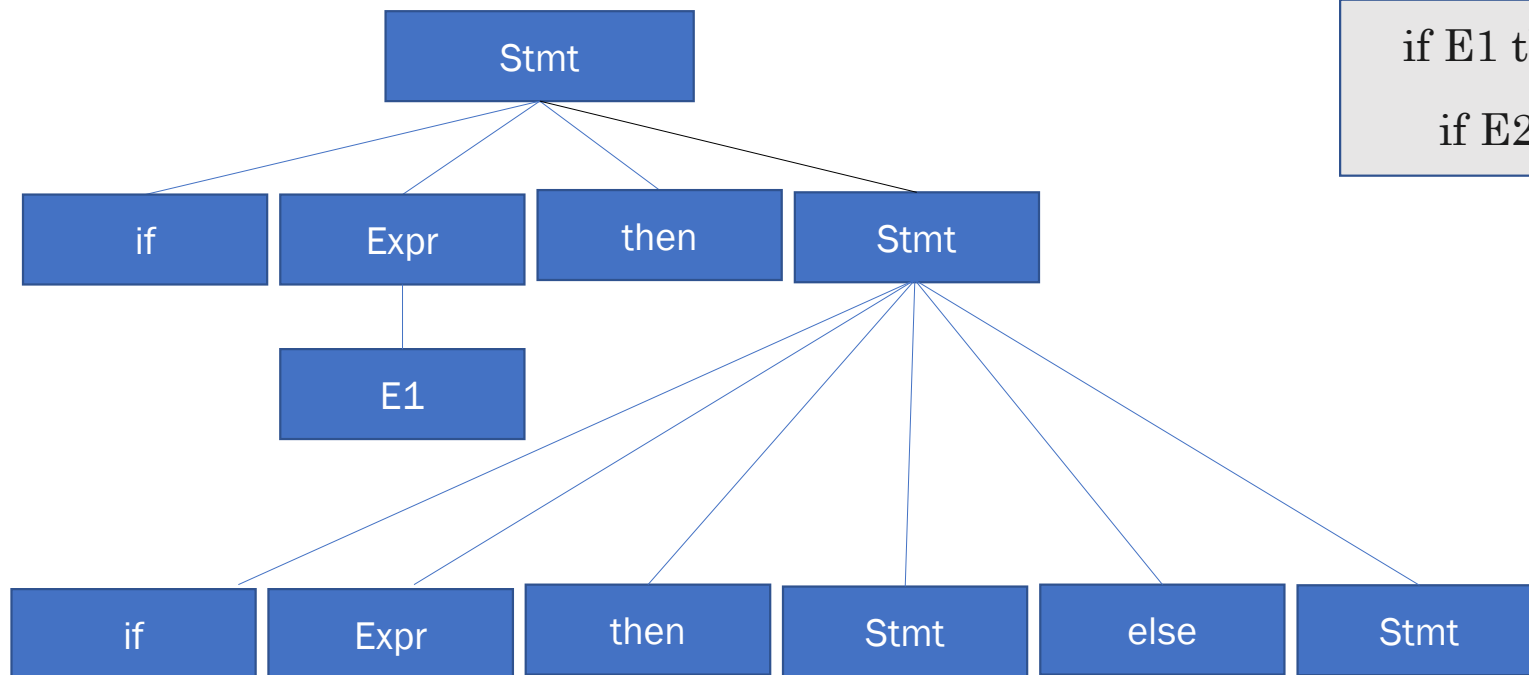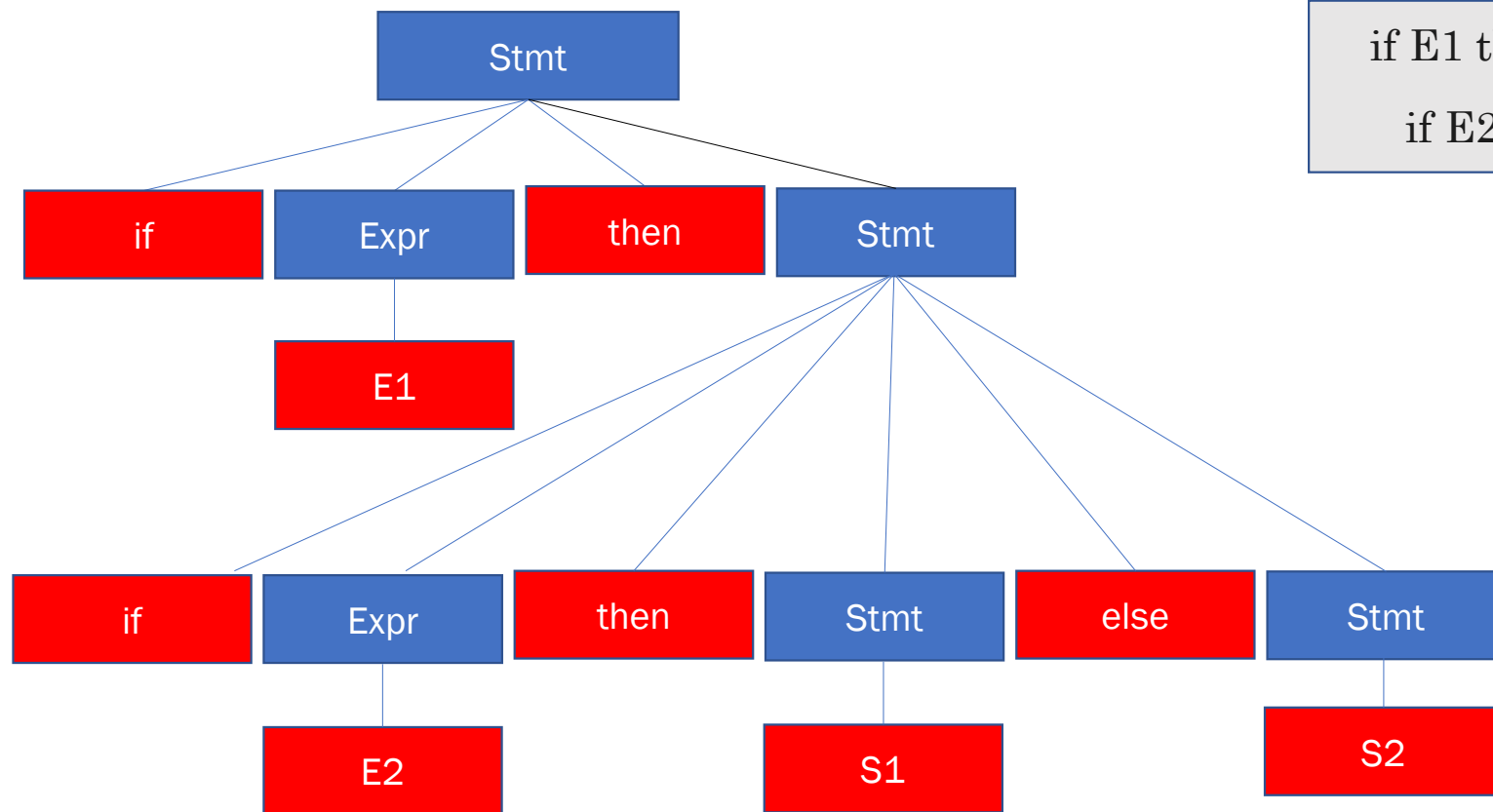
if E1 then

    if E2 then S1 else S2

# Dangling Else Problem



if E1 then
  if E2 then S1 else S2

# Dangling Else Problem

# Dangling Else Problem

if E1 then

    if E2 then S1 else S2

```
                         Stmt
        ┌──────┬──────┬──────┬──────┬──────┐
       if     Expr   then   Stmt   else   Stmt
               │             │             │
               E1     ┌───┬──┴──┬────┐     S2
                     if  Expr  then  Stmt
                          │           │
                          E2          S1
```

# Dealing with Ambiguity

- There are several ways to handle ambiguity
- We will discuss one of them

- Rewriting the grammar

## Rewriting the grammar
- use a different nonterminal for each precedence
- Start with the lowest precedence

## Rewriting the grammar

- use a different nonterminal for each precedence
- Start with the lowest precedence
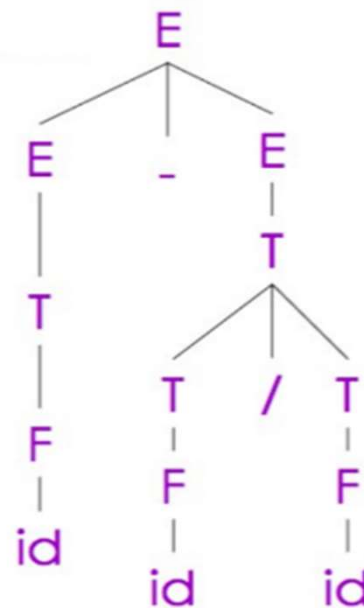
$$E \rightarrow E - E \mid E / E \mid (E) \mid id$$

rewrite to

$$E \rightarrow E - E \mid T$$
$$T \rightarrow T / T \mid F$$
$$F \rightarrow id \mid (E)$$

Parse tree for id – id / id

$$E \rightarrow E - E \mid T$$
$$T \rightarrow T / T \mid F$$
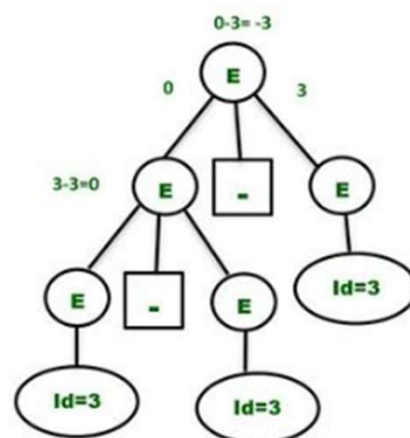$$F \rightarrow id \mid (E)$$

# Associativity

- The grammar captures operator precedence, but it is still ambiguous!
  - fails to express that both subtraction and division are left associative;

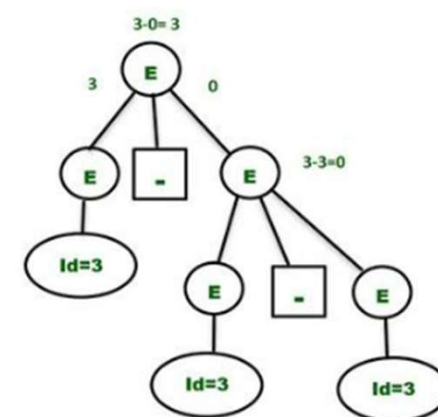$$E \rightarrow E - E \mid id$$

Parse tree for $id - id - id$.

Let's consider a single value of id = 3

- Grammar is recursive in nonterminal $X$ if:

$X \rightarrow + \ldots X \ldots$

$\rightarrow +$ means "in one or more steps, $X$ derives a sequence of symbols that includes an $X$"

# Recursion

- Grammar is recursive in nonterminal $X$ if:

  $$X \rightarrow + \ldots X \ldots$$

  $\rightarrow +$ means "in one or more steps, $X$ derives a sequence of symbols that includes an $X$"

- Grammar is <span style="color:#2196C9">left recursive</span> in $X$ if:

  $$X \rightarrow + X \ldots$$

  In one or more steps, X derives a sequence of symbols that starts with an $X$

# Recursion

- Grammar is recursive in nonterminal $X$ if:

  $X \rightarrow + \ldots X \ldots$

  $\rightarrow +$ means "in one or more steps, $X$ derives a sequence of symbols that includes an $X$"

- Grammar is left recursive in $X$ if:

  $X \rightarrow + X \ldots$

  In one or more steps, X derives a sequence of symbols that starts with an $X$

- A grammar is right recursive in $X$ if:

  $X \rightarrow + \ldots X$

  In one or more steps, $X$ derives a sequence of symbols that ends with and X

- The grammar given above is both left and right recursive in non-terminals exp and term

- To correctly expresses operator associativity:
  - For left associativity, use left recursion
  - For right associativity, use right recursion

  Here's the correct grammar

  $$E \rightarrow E - T \mid T$$
  $$T \rightarrow T/ \, F \mid F$$
  $$F \rightarrow id \mid (E)$$

# Abstract Syntax Tree

- An abstract syntax tree (AST) is a simplified version of a parse tree. An AST only contains information related to analyzing the source text and ignores extra syntactic information used for parsing text.

- Why do we need alternative?
  - Fewer intermediate nodes and subtrees
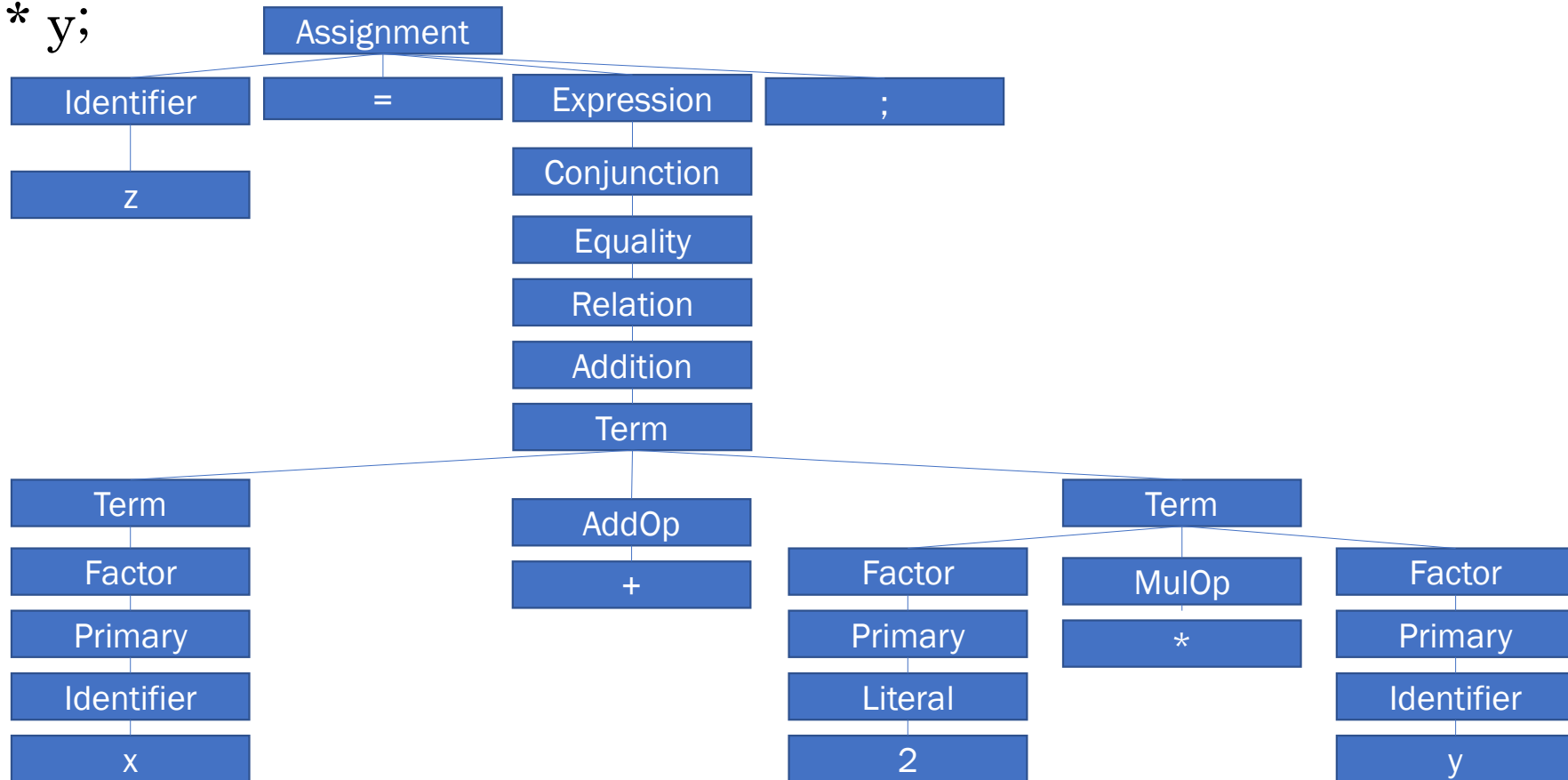  - All information like parse tree but smaller

# Transform Parse tree into AST

- Discard all the punctuation, such as semicolon

- Discard all nonterminal which are trivial roots, ones with only a single subtree

- Finally, replace the remaining non-terminals with the operators which are a leaf of one of their immediate subtrees.
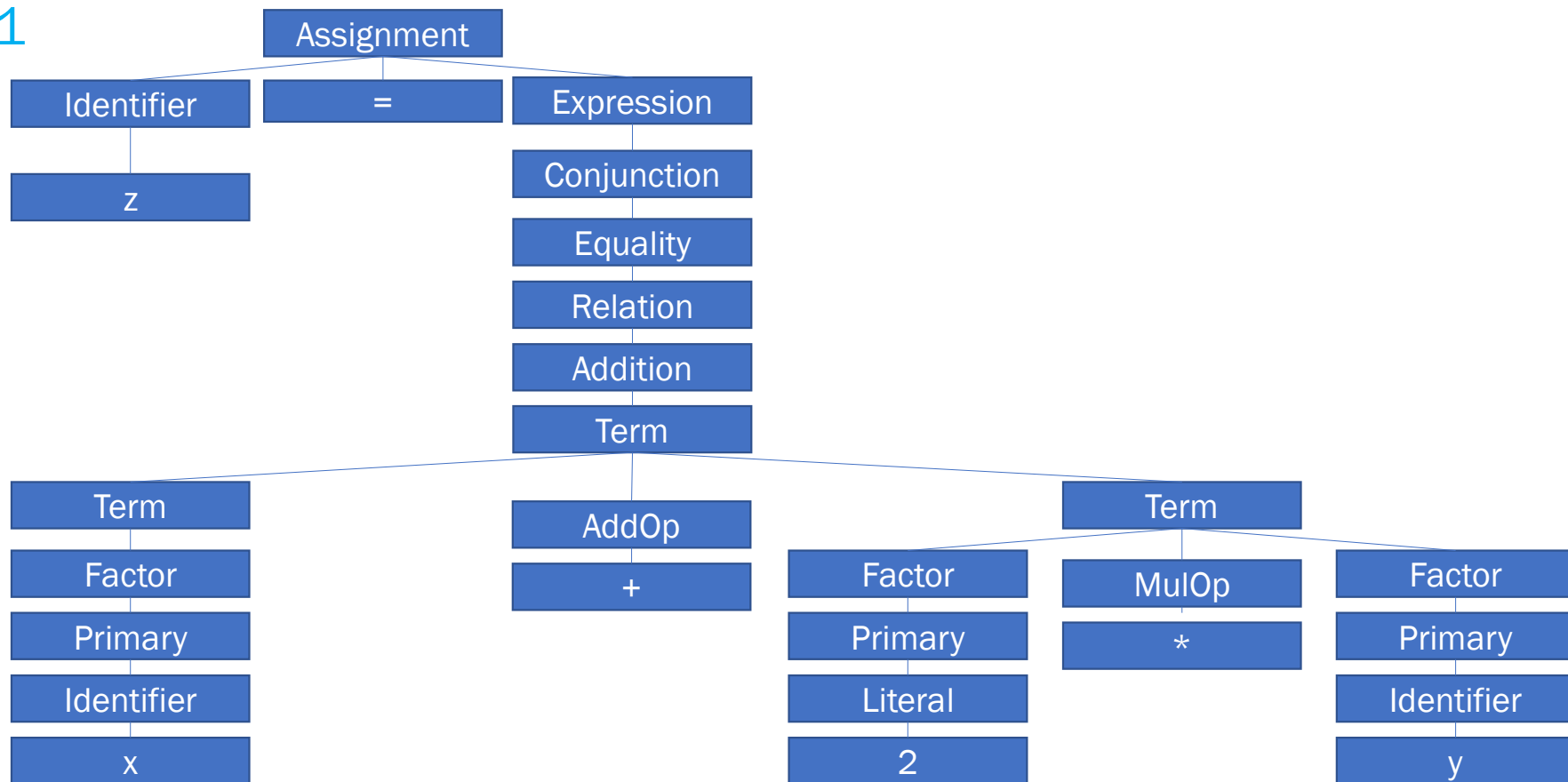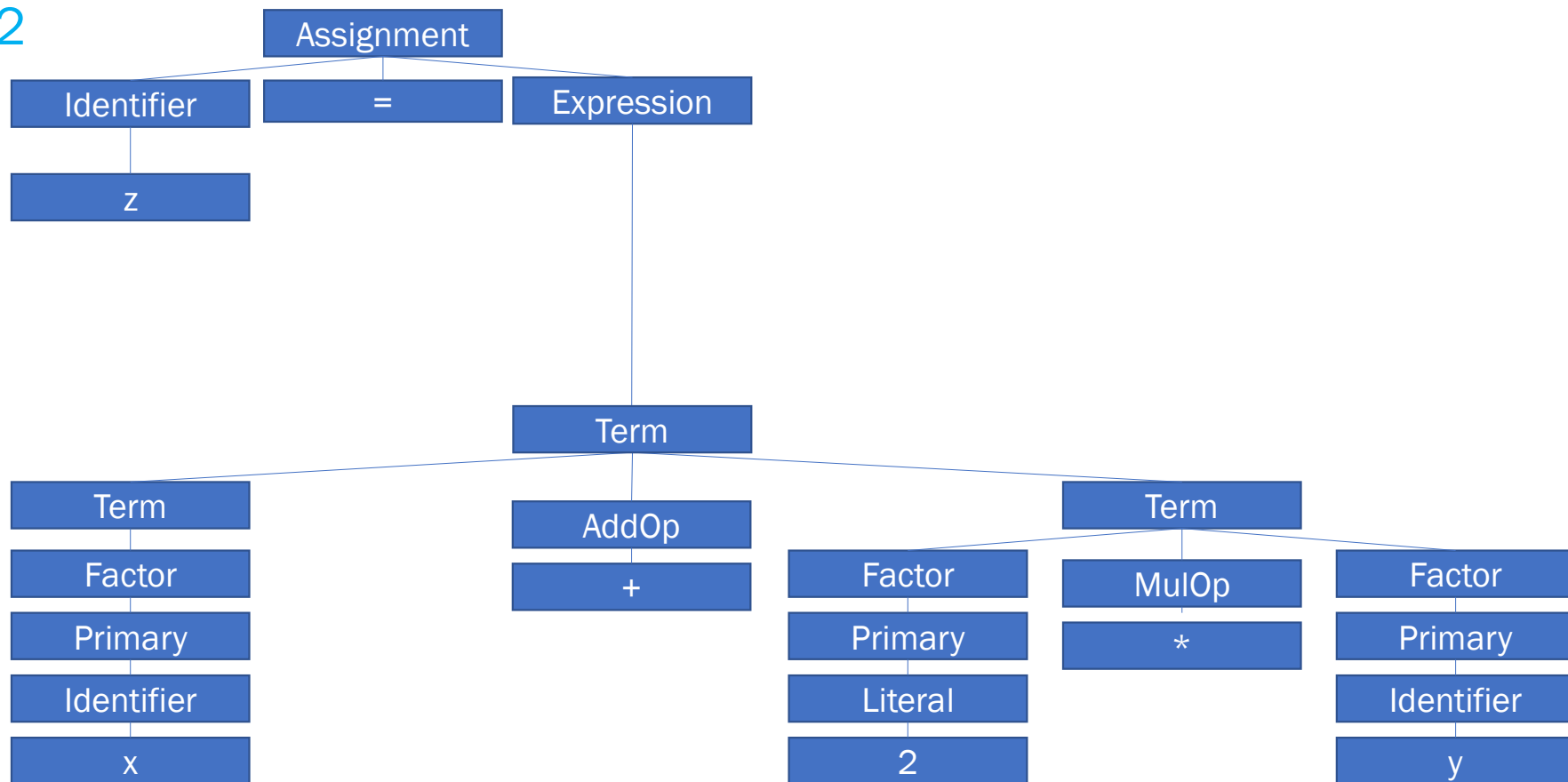
# Example (Parse Tree)

- z = x + 2 * y;

# Example (Parse Tree)

Apply Rule 1

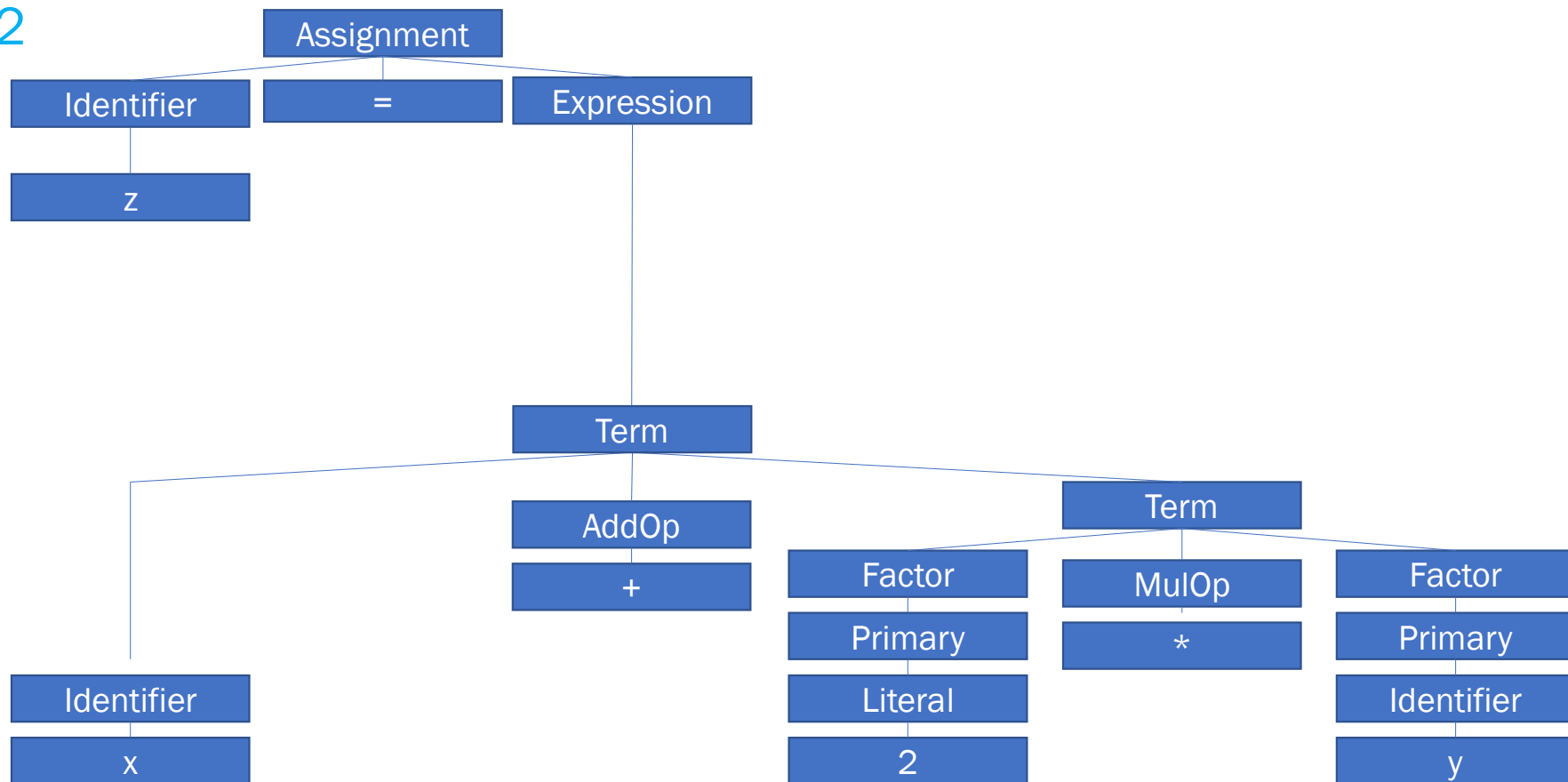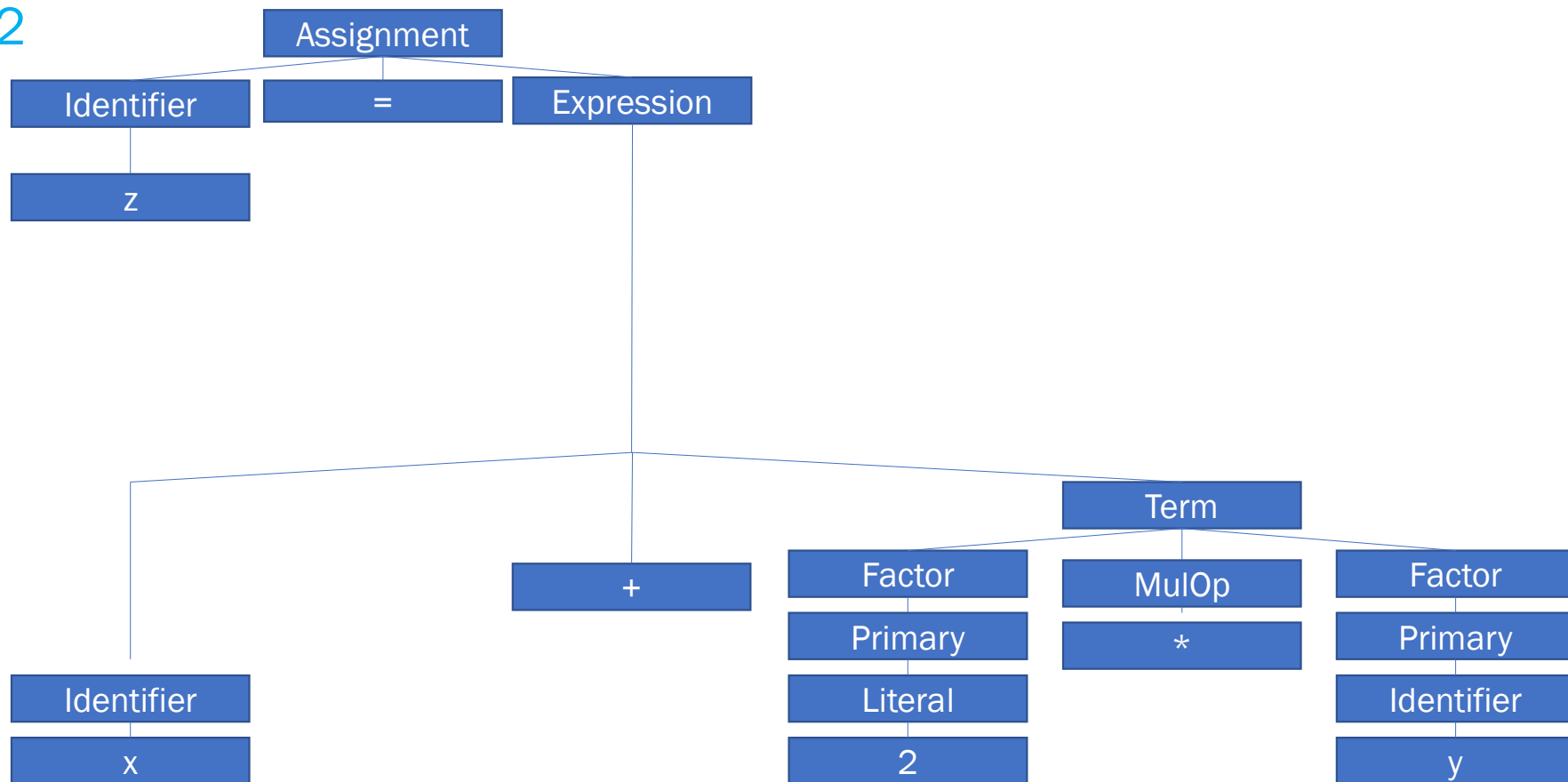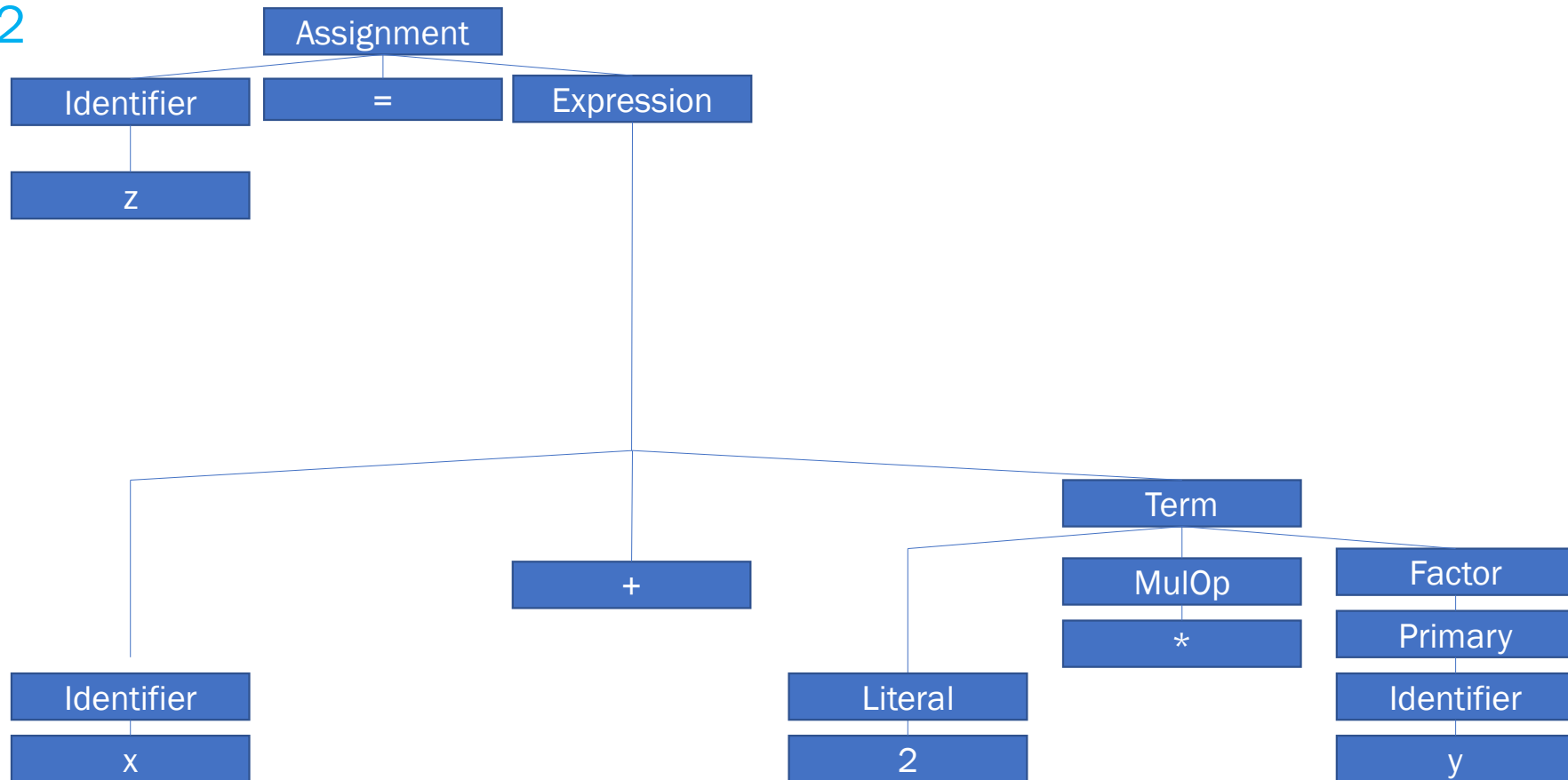# Example (Parse Tree)

Apply Rule 2

# Example (Parse Tree)

Apply Rule 2
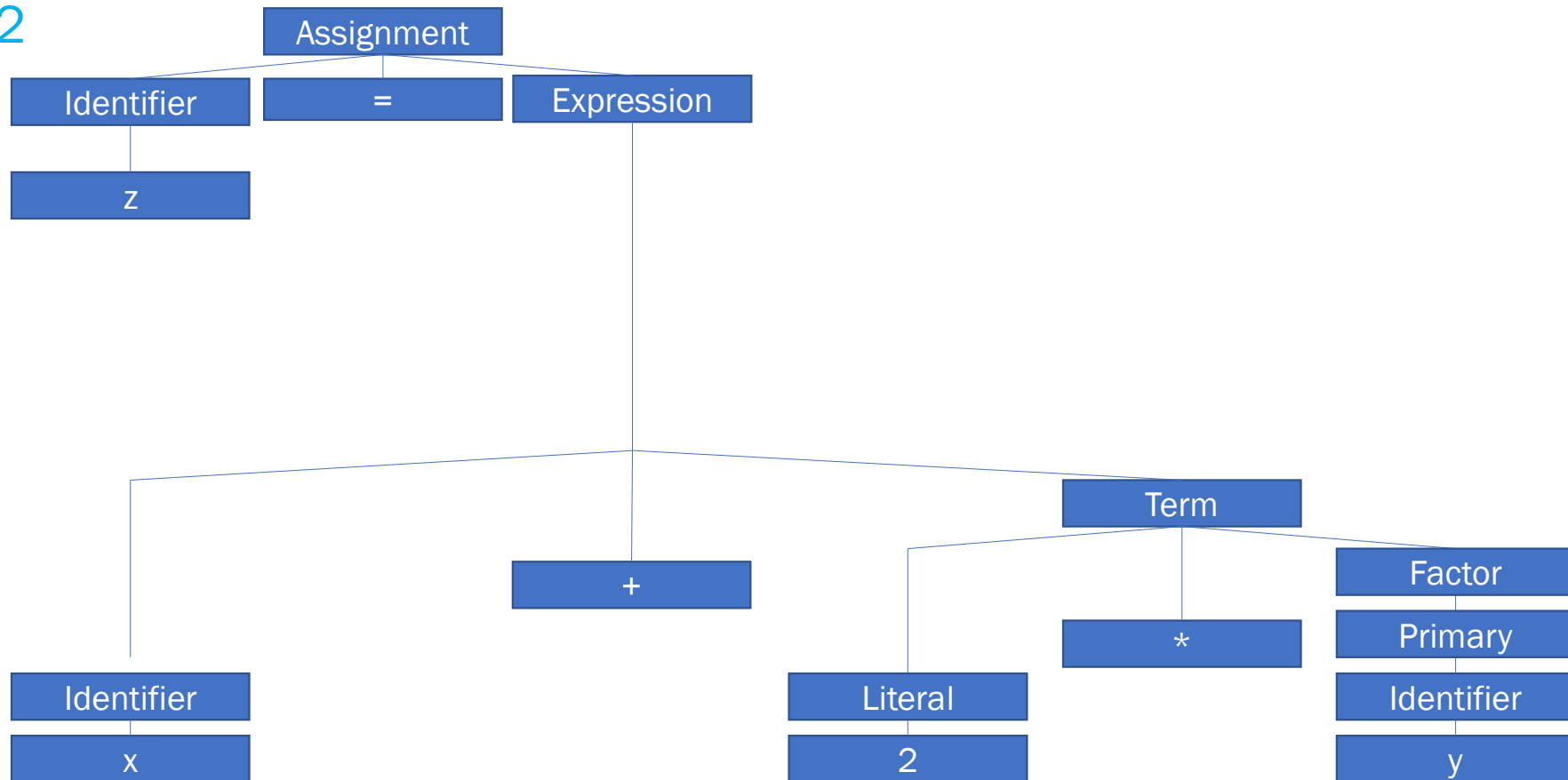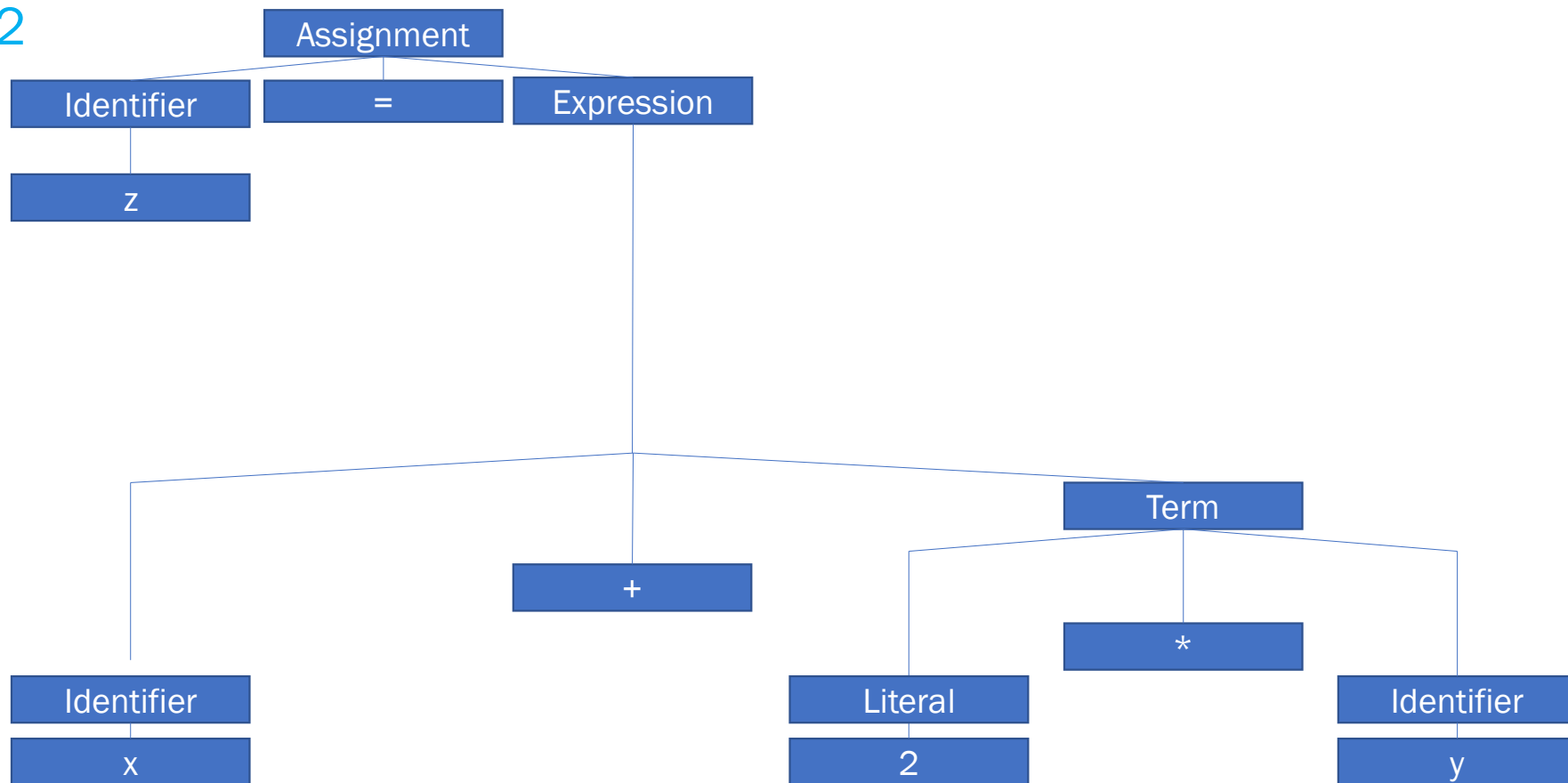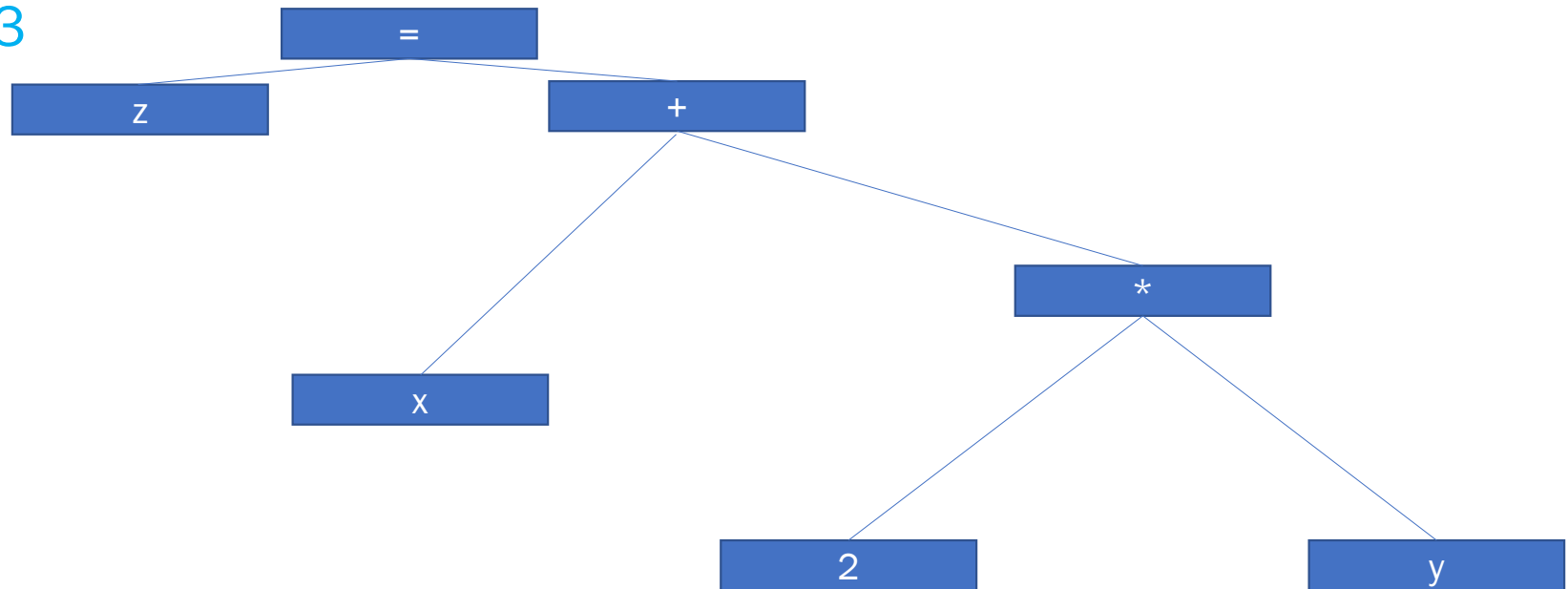
# Example (Parse Tree)
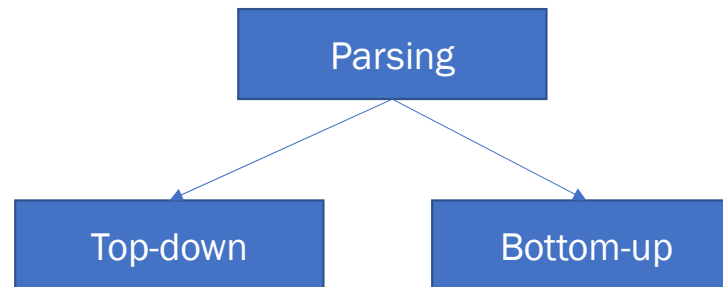
Apply Rule 2

# Example (Parse Tree)

Apply Rule 2

Apply Rule 3

# Parsing

- Top-down: the parser starts constructing the parse tree from the start symbol and then tries to transform the start symbol to the input

- Bottom-up: parsing with the input symbols and tries to construct the parse tree up to the start symbols

```
              Parsing
             /       \
       Top-down    Bottom-up
```

# Bottom-up Parsing (Example)

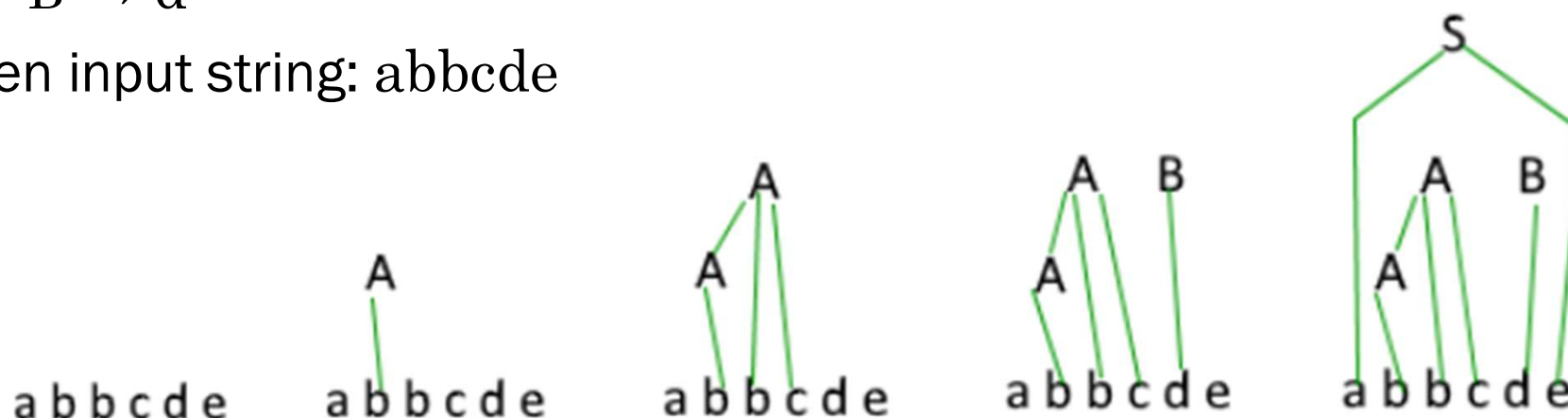- Say, we have grammar

$$S \rightarrow aABe$$

$$A \rightarrow Abc \,|\, b$$

$$B \rightarrow d$$

- Given input string: abbcde

- Say, we have grammar

$$E \to T + E \mid T$$

$$T \to int * T \mid int \mid (E)$$

- Given input string: int * int + int

- Say, we have grammar

    E → T + E | T

    T → int * T | int | (E)

- Given input string: int * int + int

    int * int + int

    int * T + int

- Say, we have grammar

    $E \rightarrow T + E \mid T$

    $T \rightarrow int * T \mid int \mid (E)$

- Given input string: int * int + int

    int * int + int

    int * T + int

    T + int

- Say, we have grammar

    E → T + E | T

    T → int * T | int | (E)


- Given input string: int * int + int

    int * int + int

    int * T + int

    T + int

    T + T

- Say, we have grammar

    $E \rightarrow T + E \mid T$

    $T \rightarrow int * T \mid int \mid (E)$

- Given input string: int * int + int

    int * int + int

    int * T + int

    T + int

    T + T

    T + E

- Say, we have grammar

    $E \rightarrow T + E \ | \ T$

    $T \rightarrow int * T \ | \ int \ | \ (E)$

- Given input string: int * int + int

    int * int + int

    int * T + int

    T + int

    T + T

    T + E

    E

- Say, we have grammar

    Integer → Digit | Integer Digit

    Digit → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

- Given input string: 24567

# Reading and Exercises

Reading

- Chapter: 2.2 (Michael Scott Book)

Exercises

- Exercises: 2.9 and 2.10 (Michael Scott Book)