

# CS406: Compilers

Spring 2021

Week 6: Parsers (LR(k)) and Semantic Processing

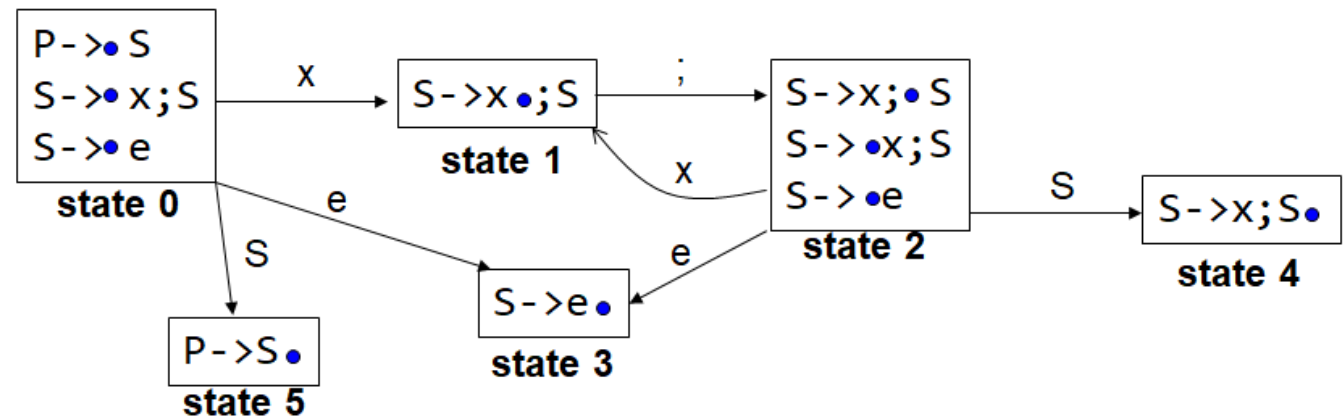
# LR(0) Parsing

- Previous Example of LR Parsing was LR(0)
  - No (0) lookahead involved
  - Operate based on the parse stack state and with goto and action tables (How?)

# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$  == saying that  $\alpha$  e.g. prefix of **x;x** is seen in the input string

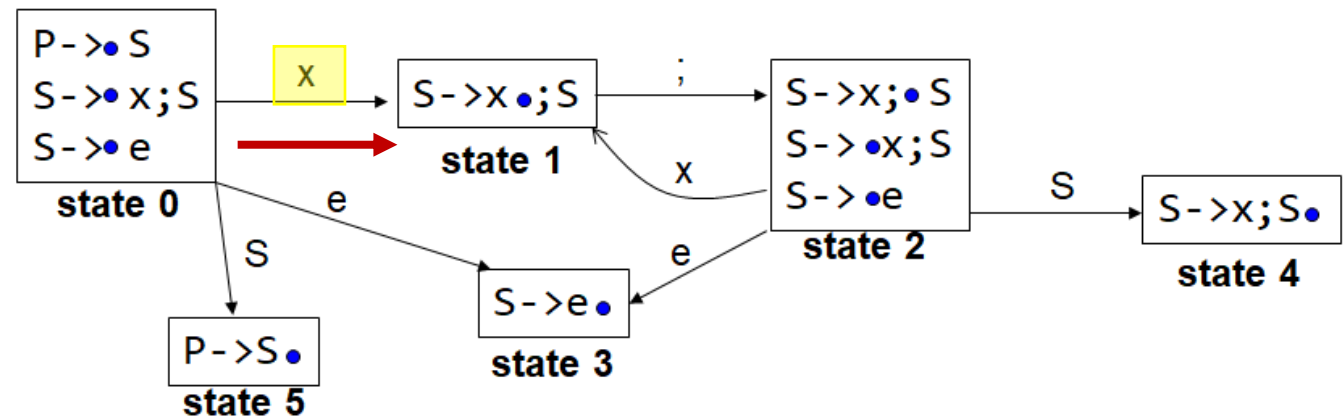
Parse Stack
0
0 1
0 1 2
0 1 2 <b>1</b>



# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$  == saying that a prefix of **x;x** is seen in the input string

Parse Stack
0
0 1
0 1 2
0 1 2 <b>1</b>

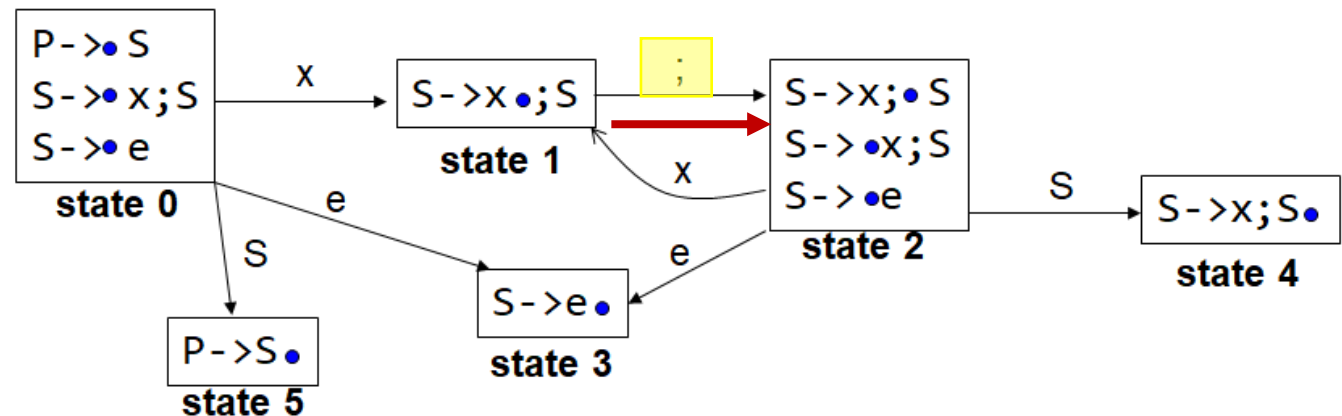


Go from state 0 to state 1 consuming x

# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$  == saying that a prefix of **x;x** is seen in the input string

Parse Stack
0
0 1
0 1 2
0 1 2 <b>1</b>

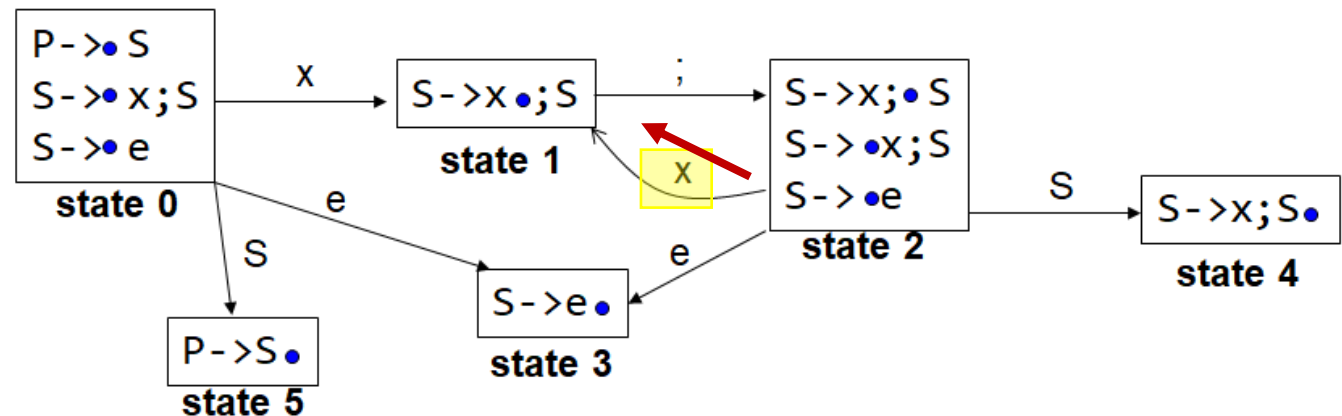


Go from state 1 to state 2 consuming ;

# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$  == saying that a prefix of **x;x** is seen in the input string

Parse Stack
0
0 1
0 1 2
0 1 2 <b>1</b>



Go from state 2 to state 1 consuming x

# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$ .  
=> we are in some state  $s$

# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$ .  
=> we are in some state  $s$ .  
We reduce by  $X \rightarrow \beta$  if state  $s$  contains  $X \rightarrow \beta \bullet$
- Note: reduction is done based solely on the current state.



# LR(0) Parsing

- Assume: Parse stack contains  $\alpha$ .

=> we are in some state  $s$ .

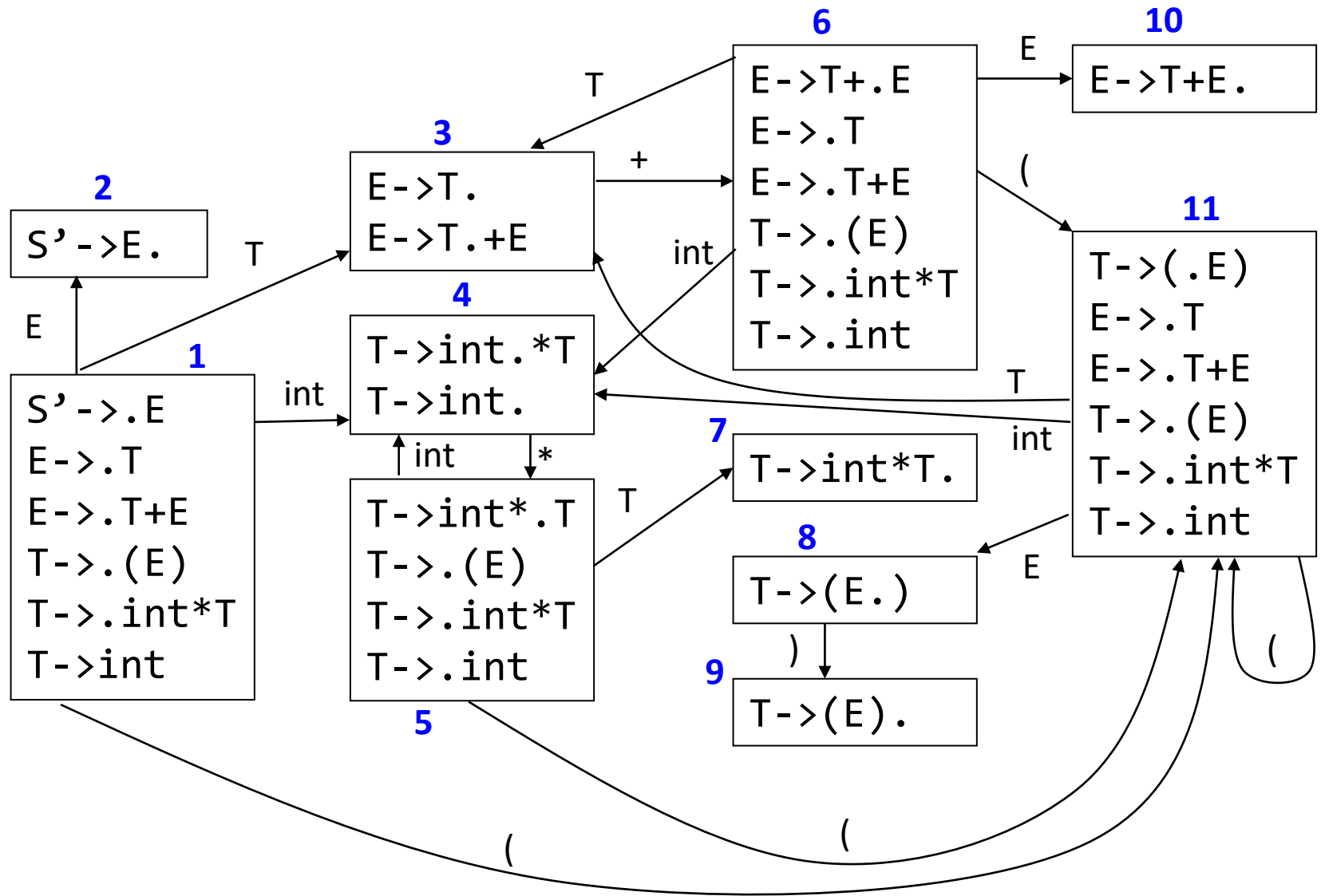
- Assume: Next input is  $t$

We **shift** if  $s$  contains  $X \rightarrow \beta \bullet t$

==  $s$  has a transition labelled  $t$

# LR(0) Parsing

- What if  $s$  contains  $X \rightarrow \beta \bullet t\omega$  and  $X \rightarrow \beta \bullet$  ?



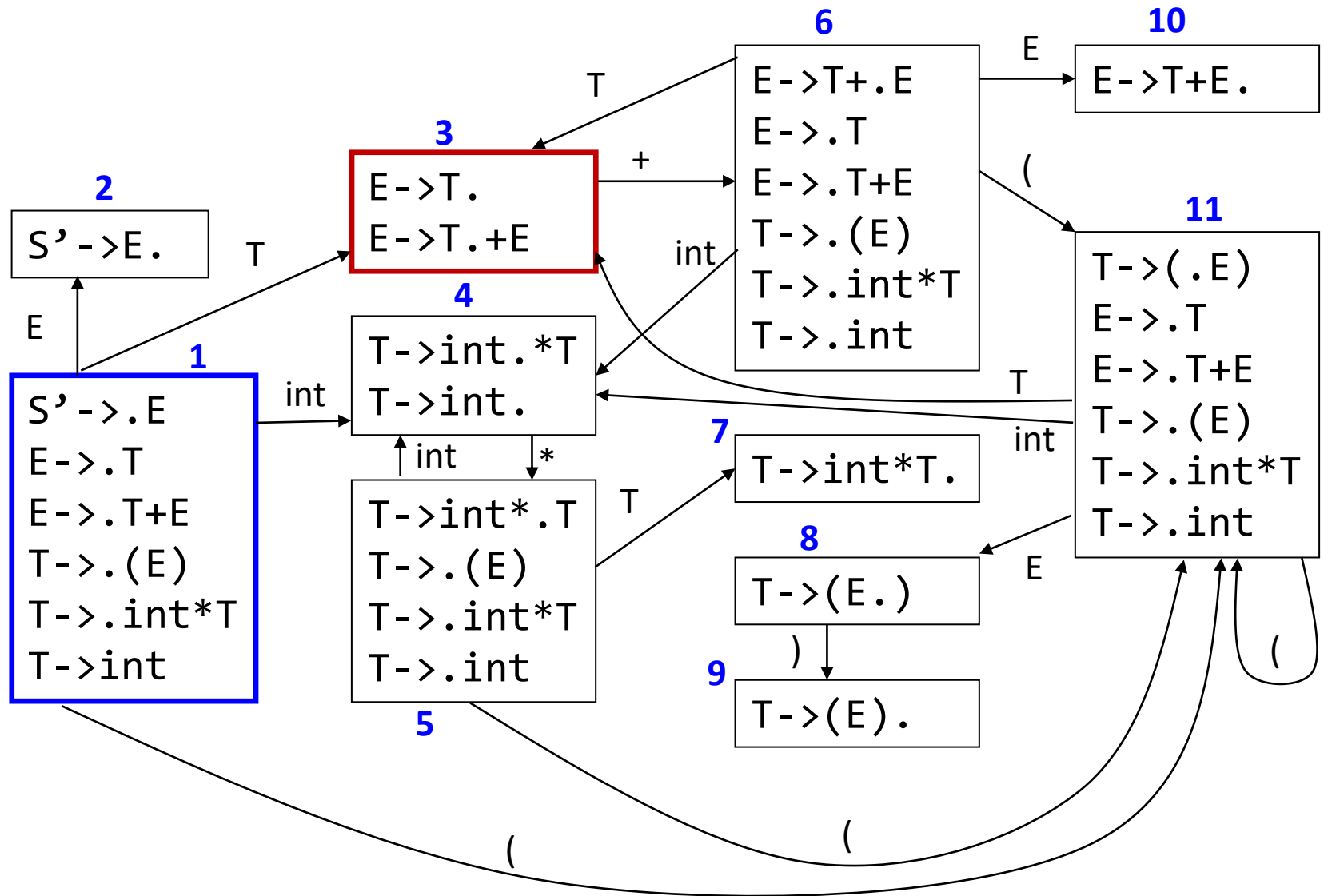
*Conflicts or not?*

# SLR Parsing

- SLR Parsing improves the shift-reduce conflict states of LR(0):

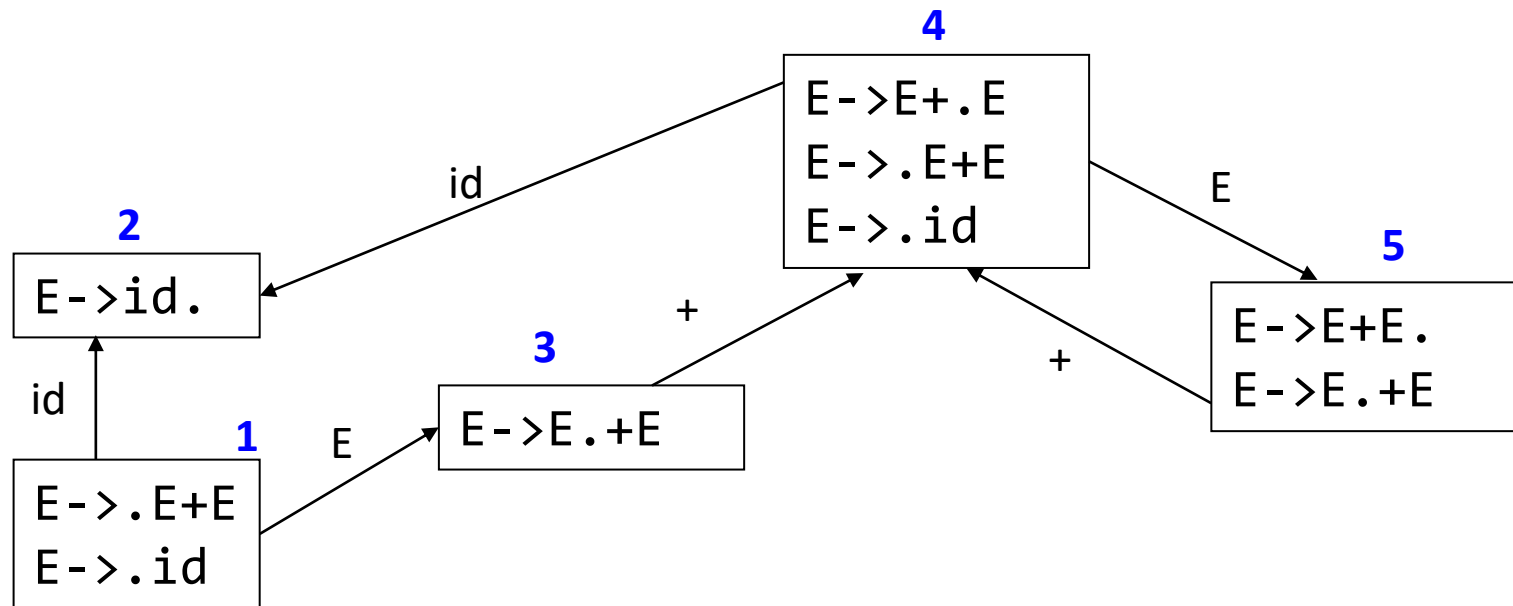
Reduce  $X \rightarrow \beta \bullet$  only if

$t \in \text{Follow}(X)$



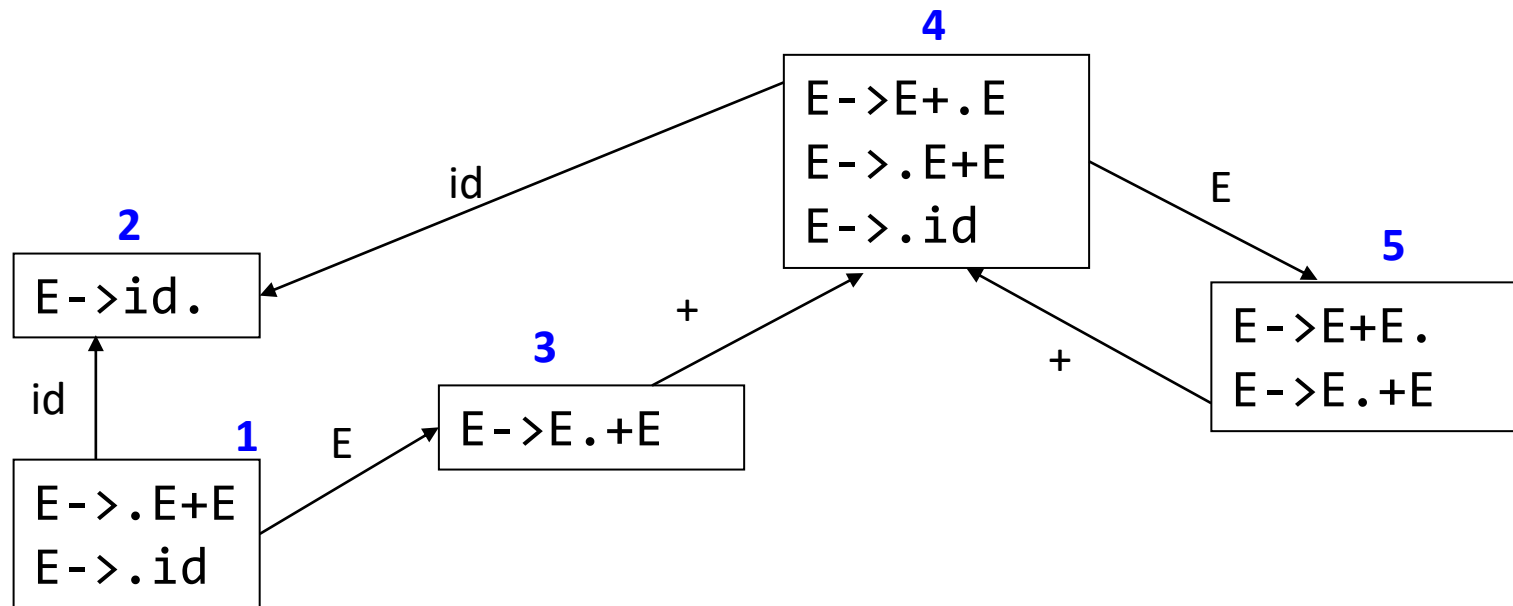
$\text{Follow}(E) = \{ \$, ) \} \Rightarrow \text{reduce by } E \rightarrow T \text{ only if } \underline{\text{next input}} \text{ is } \$ \text{ or } )$

*lookahead 1*



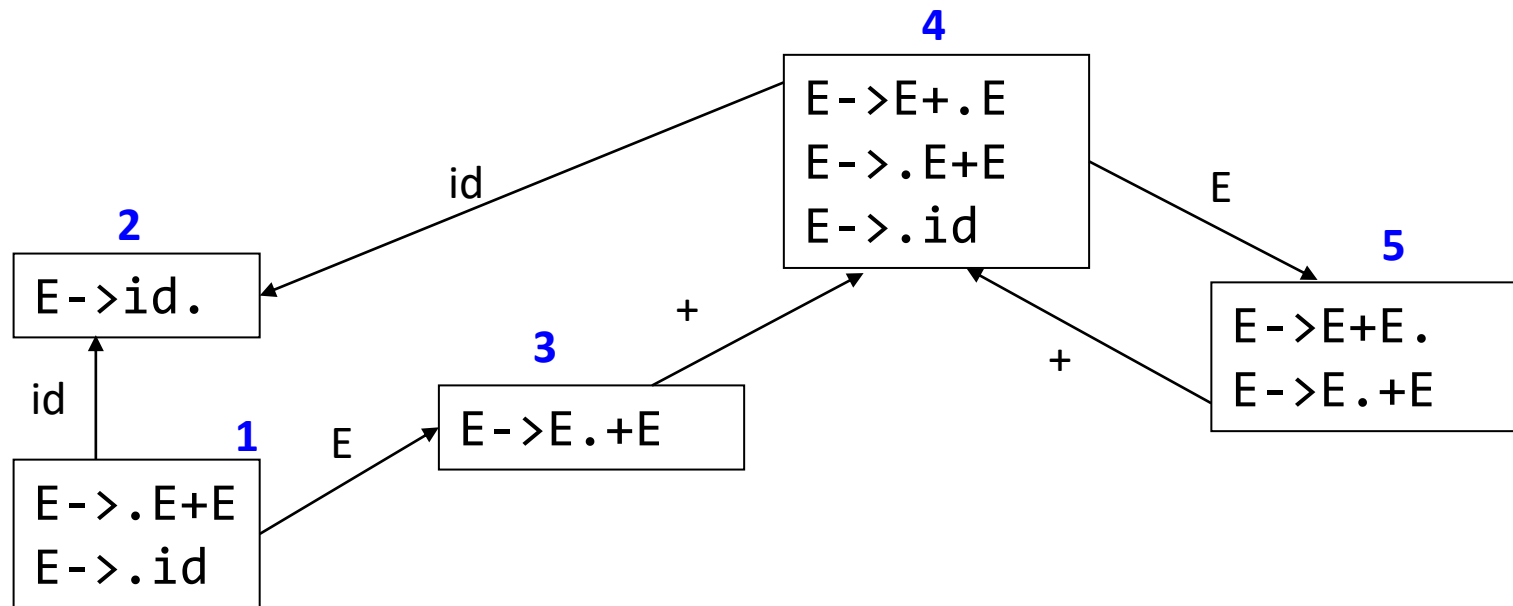
What about the grammar  $E \rightarrow E + E \mid id$  ?

LR(0)?



What about the grammar  $E \rightarrow E + E \mid id$  ?

LR(0)? SLR(1)?

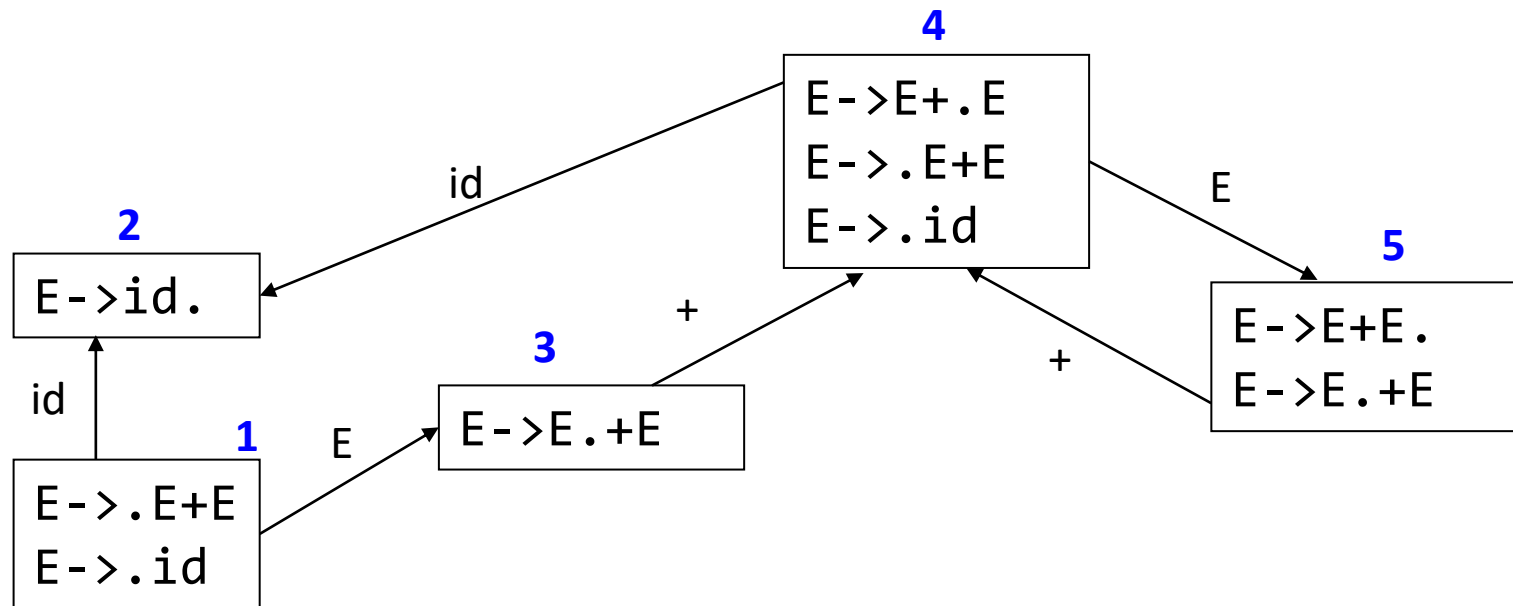


What about the grammar  $E \rightarrow E + E \mid id$  ?

LR(0)? SLR(1)?

$\text{Follow}(E) = \{+, \$\} \Rightarrow$  in state 5, reduce by  $E \rightarrow T$ . only if next input is \$ or +





What about the grammar  $E \rightarrow E + E \mid id$  ?

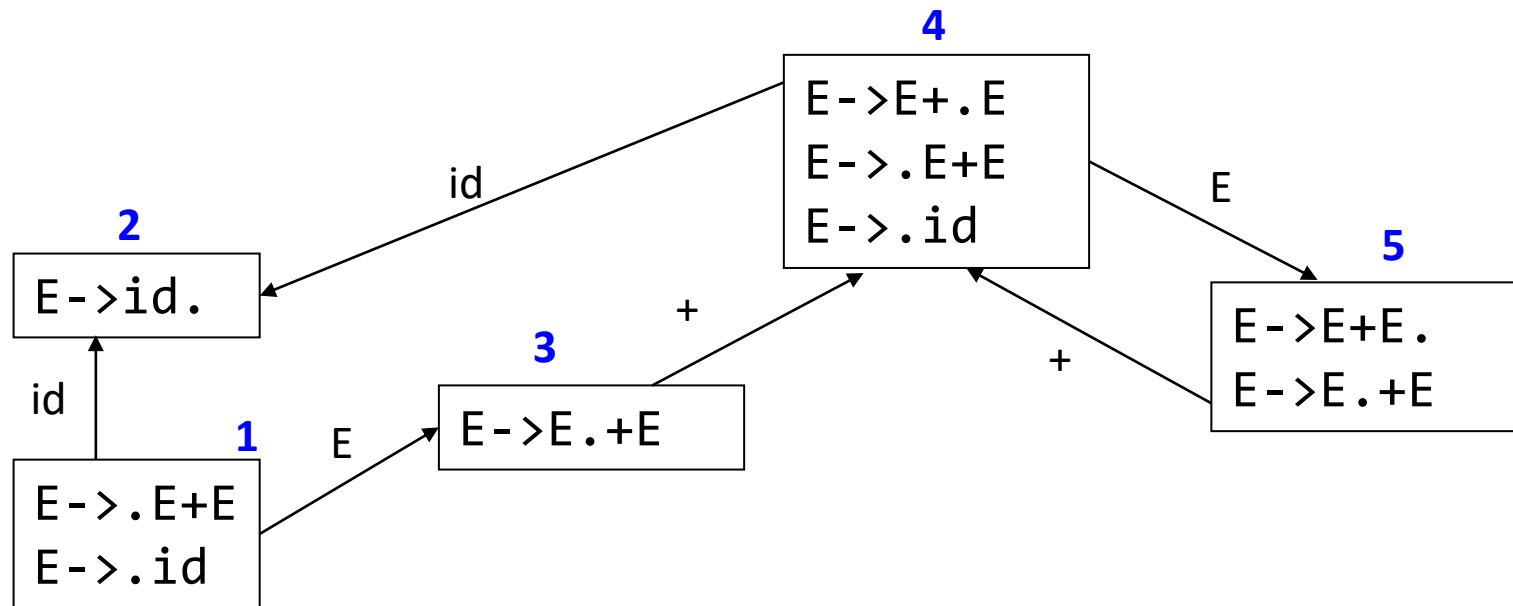
LR(0)? SLR(1)?

$\text{Follow}(E) = \{+, \$\} \Rightarrow$  in state 5, reduce by  $E \rightarrow T$ . only if next input is  $\$$  or  $+$

But state 5 has  $E \rightarrow E.+E$  (shift if next input is +)  
**Shift-reduce conflict!**

# LR( $k$ ) parsers

- LR(0) parsers
  - No lookahead
  - Predict which action to take by looking only at the symbols currently on the stack
- LR( $k$ ) parsers
  - Can look ahead  $k$  symbols
  - Most powerful class of deterministic bottom-up parsers
  - LR(1) and variants are the most common parsers



What about the grammar  $E \rightarrow E + E \mid id$  ?

LR(0)? SLR(1)?

$\text{Follow}(E) = \{+, \$\} \Rightarrow$  in state 5, reduce by  $E \rightarrow T$ . only if next input is  $\$$  or  $+$

But state 5 has  $E \rightarrow E.+E$  (shift if next input is  $+$ )

Shift-reduce conflict!

%left +

*says reduce if the next input symbol is + i.e. prioritize rule  $E+E$ . over  $E.+E$*

# Top-down vs. Bottom-up parsers

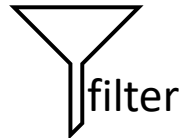
- Top-down parsers expand the parse tree in *pre-order*
  - Identify parent nodes before the children
- Bottom-up parsers expand the parse tree in *post-order*
  - Identify children before the parents
- Notation:
  - LL(1): Top-down derivation with 1 symbol lookahead
  - LL(k): Top-down derivation with k symbols lookahead
  - LR(1): Bottom-up derivation with 1 symbol lookahead

# Semantic Processing

**Lexical Analysis** — Detects programs with illegal tokens



**Parsing** — Detects programs with ill-formed programming constructs i.e. invalid parse tree structure



**Semantic Routines** — Detects all remaining errors



“Front-end”

# Semantic Processing

- Syntax-directed / syntax-driven
  - Routines (called as semantic routines) interpret the meaning of programming constructs based on the syntactic structure
- Routines play a dual role
  - Analysis – Semantic analysis
    - undefined vars, undefined types, uninitialized variables, type errors that can be caught at compile time, unreachable code, etc.
  - Synthesis – Generation of intermediate code
    - 3 address code
- Routines create semantic records to aid the analysis and synthesis

# Semantic Processing

- **Syntax-directed translation:** notation for attaching program fragments to grammar productions.
  - Program fragments are executed when productions are matched
  - The combined execution of all program fragments produces the translation of the program

e.g.  $E \rightarrow E + T$  { print( '+' ) }

Output: program fragments may create AST and 3 Address Codes

- **Attributes:** any 'quality' associated with a terminal and non-terminal e.g. type, number of lines of a code, first line of the code block etc.

# Why Semantic Analysis?

- Context-free grammars cannot specify all requirements of a language
  - Identifiers declared before their use (scope)
  - Types in an expression must be consistent  

```
STRING str:= "Hello";  
str := str + 2;
```
  - Number of formal and actual parameters of a function must match
  - Reserved keywords cannot be used as identifiers
  - A Class is declared only once in a OO language, a method can be overridden.
  - ...



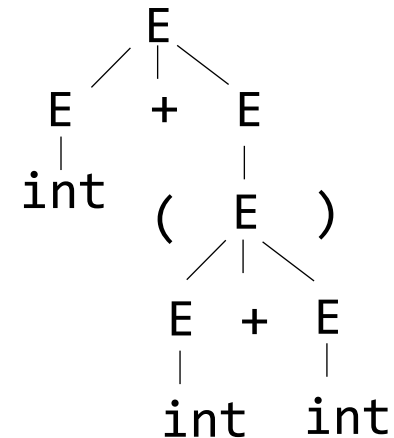
# Abstract Syntax Tree

- Abstract Syntax Tree (AST) or Syntax Tree can be the input for semantic analysis.
  - What is Concrete Syntax Tree? – the parse tree
- ASTs are like parse trees but ignore certain details:

E.g. Consider the grammar:

$E \rightarrow E + E$   
 $E \rightarrow ( E )$   
 $E \rightarrow \text{int}$

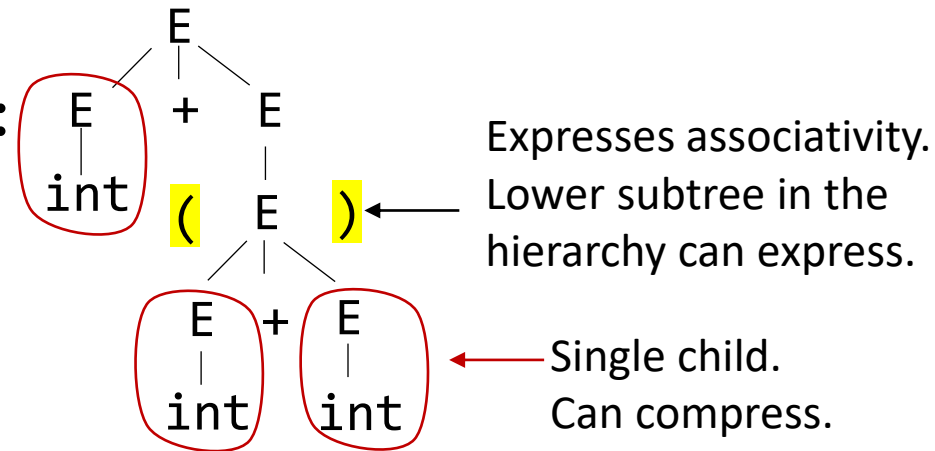
The parse tree for  $1+(2+3)$



# Abstract Syntax Tree - Example

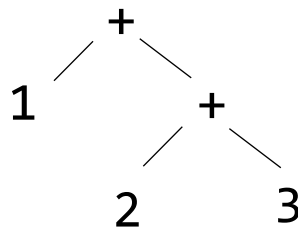
- Not all details (nodes) of the parse tree are helpful for semantic analysis

The parse tree for  $1+(2+3)$  :

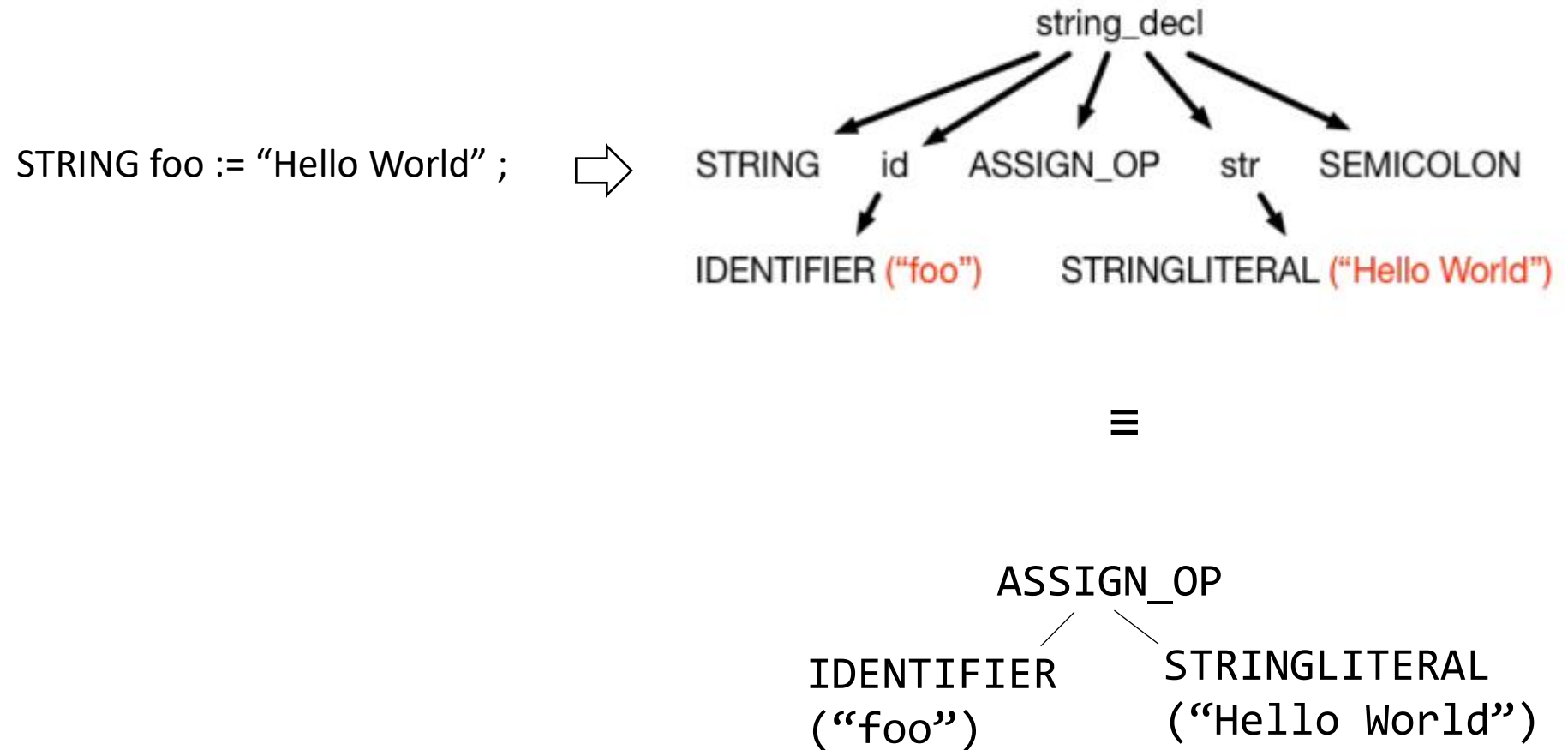


We need to compute the result of the expression. So, a simpler structure is sufficient:

AST for  $1+(2+3)$  :



# AST - Example

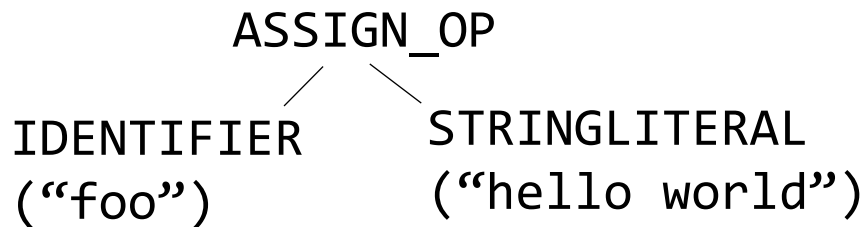


# Semantic Analysis – Example

- Context-free grammars cannot specify all requirements of a language
  - Identifiers declared before their use (scope)
  - Types in an expression must be consistent
    - Type checks
    - STRING str:= “Hello”;
    - str := str + 2;
  - Number of formal and actual parameters of a function must match
  - Reserved keywords cannot be used as identifiers
  - A Class is declared only once in a OO language, a method can be overridden.
  - ...

# Scope

- **Goal:** matching identifier declarations with uses
- Most languages require this!
- Scope confines the activity of an identifier



What if `foo` is declared as a `STRING` in an enclosing scope but is an `INT` in the current scope?

in different parts of the program:

- Same identifier may refer to different things
- Same identifier may not be accessible

# Static vs. Dynamic Scope

- Most languages are statically scoped
  - Scope depends only the program text (not runtime behavior)
  - A variable refers to the closest defined instance

```
INT w, x;
```

```
{
```

```
    FLOAT x, z;
```

```
    f(x, w, z);
```

```
}
```

```
g(x)
```

x is a FLOAT here

x is an INT here

```
f(){
```

```
    a=4; g();
```

```
}
```

```
g() { print(a); }
```

value of a is 4 here

- In dynamically scoped languages
  - Scope depends on the execution context
  - A variable refers to the closest enclosing binding in the execution of the program

# Exercise: Static vs. Dynamic Scope

```
#define a (x+1)
```

```
int x = 2;
```

statically scoped or dynamically scoped?



```
void b() { int x = 1; printf("%d\n",a);}
```

```
void c() { printf("%d\n",a);}
```

```
int main() { b(); c(); }
```

# Symbol Table

- Data structure that tracks the bindings of identifiers. Specifically, returns the current binding.
  - E.g., stores a mapping of names to types
- Should provide for efficient retrieval and frequent insertion and deletion of names. Should consider scopes.

```
{  
    int x = 0;  
    //accessing y here should be illegal  
    {  
        int y = 1;  
    }  
}
```

- Can use stacks, binary trees, hash maps for implementation



# Symbol Table and Classes in OO Language

- Class names may be used before their definition
- Can't use symbol table (to check class definition)
  - Gather all class names
  - Do multiple passes.
- Semantic analysis is done in multiple passes
- One of the goals of semantic analysis is to create/update data structures that help the next round of analysis

# Semantic Analysis – How?

- Recursive descent of AST
    - Process a node, n
    - Recurse into children of n and process them
    - Finish processing the node, n
- ⇒ Do a postorder processing of the AST
- As you visit a node, you will add information depending upon the analysis performed
    - The information is referred to as attributes of the node

# Case study - Semantic Analysis of Expressions

- Fully parenthesized expression (FPE)
  - Expressions (algebraic notation) are the normal way we are used to seeing them. E.g.  $2 + 3$
  - *Fully-parenthesized* expressions are simpler versions: every binary operation is enclosed in parenthesis
    - E.g.  $2 + 3$  is written as  $(2+3)$
    - E.g.  $(2 + (3 * 7))$
    - We can ignore order-of-operations (PEMDAS rule) in FPEs.

# FPE – definition

- Either a:
  1. A number (integer in our example) OR
  2. *Open parenthesis* ‘(’ followed by  
*FPE* followed by  
*an operator* (‘+’, ‘-’, ‘\*’, ‘/’) followed by  
*FPE* followed by  
*closed parenthesis* ‘)’

# FPE – Notation

1.  $E \rightarrow \text{INTLITERAL}$

2.  $E \rightarrow (E \text{ op } E)$

3.  $\text{op} \rightarrow \text{ADD} \mid \text{SUB} \mid \text{MUL} \mid \text{DIV}$

# Implementing a parser for FPE

1. One function defined for every non-terminal

- E, op

2. One function defined for every production

- E1, E2

3. One function defined for all terminals

- IsTerm

1.E  $\rightarrow$  INTLITERAL

2.E  $\rightarrow$  (E op E)

3.op  $\rightarrow$  ADD | SUB | MUL | DIV

# Implementing a parser for FPE

1. One function defined for every non-terminal

- E, op

2. One function defined for every production

- E1, E2

3. One function defined for all terminals

- IsTerm

1.E  $\rightarrow$  INTLITERAL

2.E  $\rightarrow$  (E op E)


3.op  $\rightarrow$  ADD | SUB | MUL | DIV

# Implementing a parser for FPE

*This function checks if the next token returned by the scanner matches the expected token. Returns true if match. false if no match.*

Assume that a scanner module has been provided.  
The scanner has one function, `GetNextToken`, that  
returns the next token in the sequence.

Can be any one of: `INTLITERAL`, `LPAREN`,  
`RPAREN`, `ADD`, `SUB`, `MUL`, `DIV`



```
bool IsTerm(Scanner* s, TOKEN tok) {  
    return s->GetNextToken() == tok;  
}
```

The diagram shows two arrows. One arrow originates from the text 'returns the next token in the sequence.' and points to the parameter `s` in the function signature. The other arrow originates from the text 'Can be any one of: INTLITERAL, LPAREN, RPAREN, ADD, SUB, MUL, DIV' and points to the parameter `tok` in the function signature.



# Implementing a parser for FPE

1. One function defined for every non-terminal

- E, op

2. One function defined for every production

- E1, E2

3. One function defined for all terminals

- IsTerm

1.E  $\rightarrow$  INTLITERAL

2.E  $\rightarrow$  (E op E)

3.op  $\rightarrow$  ADD | SUB | MUL | DIV

# Implementing a parser for FPE

*This function implements production #1:  $E \rightarrow \text{INTLITERAL}$*

*Returns true if the next token returned by the scanner is an INTLITERAL. false otherwise.*

```
bool E1(Scanner* s) {  
    return IsTerm(s, INTLITERAL);  
}
```

# Implementing a parser for FPE

1. One function defined for every non-terminal

- E, op

2. One function defined for every production

- E1, E2

3. One function defined for all terminals

- IsTerm

1.E  $\rightarrow$  INTLITERAL

2.E  $\rightarrow$  (E op E)

3.op  $\rightarrow$  ADD | SUB | MUL | DIV

# Implementing a parser for FPE

*This function implements production #2:  $E \rightarrow (E \text{ op } E)$*

*Returns true if the Boolean expression on line 2 returns true. false otherwise.*

```
1: bool E2(Scanner* s) {  
  
2:   return IsTerm(s, LPAREN) &&  
           E(s) &&  
           OP(s) &&  
           E(s) &&  
           IsTerm(s, RPAREN);  
  
3: }
```

# Implementing a parser for FPE

1. One function defined for every non-terminal

- E, op

2. One function defined for every production

- E1, E2

3. One function defined for all terminals

- IsTerm

1.E  $\rightarrow$  INTLITERAL

2.E  $\rightarrow$  (E op E)

3.op  $\rightarrow$  ADD | SUB | MUL | DIV

# Implementing a parser for FPE

*This function implements production #3:  $op \rightarrow ADD \mid SUB \mid MUL \mid DIV$*

*Returns true if the next token returned by the scanner is any one from ADD, SUB, MUL, DIV. false otherwise.*

```
bool OP(Scanner* s) {  
  
    TOKEN tok = s->GetNextToken();  
  
    if((tok == ADD) || (tok == SUB) || (tok ==  
        MUL) || (tok == DIV))  
        return true;  
  
    return false;  
  
}
```

# Implementing a parser for FPE

1. One function defined for every non-terminal

- **E**, op

2. One function defined for every production

- E1, E2

3. One function defined for all terminals

- IsTerm

1.E  $\rightarrow$  INTLITERAL

2.E  $\rightarrow$  (E op E)


3.op  $\rightarrow$  ADD | SUB | MUL | DIV

# Implementing a parser for FPE

*This function implements the routine for matching non-terminal E*

```
bool E(Scanner* s) {  
    TOKEN* prevToken = s->GetCurTokenSequence();  
    if(!E1(s)) {  
        s->SetCurTokenSequence(prevToken);  
        return E2(s);  
    }  
    return true;  
}
```

Assume that GetCurTokenSequence  
returns a reference to the first token in  
a sequence of tokens maintained by  
the scanner





# Implementing a parser for FPE

*This function implements the routine for matching non-terminal E*

```
bool E(Scanner* s) {  
  
    TOKEN* prevToken = s->GetCurTokenSequence();  
    if(!E1(s)) {  
        s->SetCurTokenSequence(prevToken);  
        return E2(s);  
    }  
    return true;  
}
```

```
//This line implements the check to see if the sequence of tokens match production #1:  
E->INTLITERAL.
```

# Implementing a parser for FPE

*This function implements the routine for matching non-terminal E*

```
bool E(Scanner* s) {  
  
    TOKEN* prevToken = s->GetCurTokenSequence();  
    if(!E1(s)) {  
        s->SetCurTokenSequence(prevToken);  
        return E2(s);  
    }  
    return true;  
}
```

//because E1(s) calls s->GetNextToken() internally, the reference to the sequence of tokens would have moved forward. This line restores the reference back to the first node in the sequence so that the scanner provides the correct sequence to the call E2 in next line

# Implementing a parser for FPE

*This function implements the routine for matching non-terminal E*

```
bool E(Scanner* s) {  
  
    TOKEN* prevToken = s->GetCurTokenSequence();  
    if(!E1(s)) {  
        s->SetCurTokenSequence(prevToken);  
        return E2(s);  
    }  
    return true;  
}
```

```
//This line implements the check to see if the sequence of tokens match production #2:  
E->(E op E)
```

# Implementing a parser for FPE

```
IsTerm(Scanner* s, TOKEN tok) { return s->GetNextToken() == tok;}

bool E1(Scanner* s) {
    return IsTerm(s, INTLITERAL);
}

bool E2(Scanner* s) { return IsTerm(s, LPAREN) && E(s) && OP(s) && E(s) && IsTerm(s, RPAREN); }

bool OP(Scanner* s) {
    TOKEN tok = s->GetNextToken();
    if((tok == ADD) || (tok == SUB) || (tok == MUL) || (tok == DIV))
        return true;
    return false;
}

bool E(Scanner* s) {
    TOKEN* prevToken = s->GetCurTokenSequence();
    if(!E1(s)) {
        s->SetCurTokenSequence(prevToken);
        return E2(s);
    }
    return true;
}
```

***Start the parser by invoking E().***

***Value returned tells if the expression is FPE or not.***

# Exercise

- What parsing technique does this parser use?

# Building Abstract Syntax Trees

- Can build while parsing a FPE
  - Via bottom-up building of the tree
- Create subtrees, make those subtrees left- and right-children of a newly created root.
  - Modify recursive parser:
    - If:  
token == INTLITERAL, return a reference to newly created node containing a number
    - Else:  
store references to nodes that are left- and right- expression subtrees  
Create a new node with value = 'OP'

# Building Abstract Syntax Trees

- Can build while parsing a FPE
  - Via bottom-up building of the tree
- Create subtrees, make those subtrees left- and right-children of a newly created root.
  - Modify recursive parser:
    - If:  
token == INTLITERAL, return a reference to newly created node containing a number
    - Else:  
store references to nodes that are left- and right- expression subtrees  
Create a new node with value = 'OP'

# Building Abstract Syntax Trees

*This function creates an AST node and adds information that stores the value of an INTLITERAL in the node. A reference to the AST node is returned.*

```
TreeNode* IsTerm(Scanner* s, TOKEN tok) {  
    TreeNode* ret = NULL;  
    TOKEN nxtToken = s->GetNextToken();  
    if(nxtToken == tok)  
        ret = CreateTreeNode(nxtToken.val);  
    return ret;  
}
```



# Building Abstract Syntax Trees

*E1 needs to change because IsTerm returns a TreeNode\*.  
E1 returns a TreeNode\* now.*

```
TreeNode* E1(Scanner* s) {  
    return IsTerm(s, INTLITERAL);  
}
```

# Building Abstract Syntax Trees

- Can build while parsing a FPE
  - Via bottom-up building of the tree
- Create subtrees, make those subtrees left- and right-children of a newly created root.
  - Modify recursive parser:
    - If:
      - token == INTLITERAL, return a reference to newly created node containing a number
    - Else:
      - store references to nodes that are left- and right- expression subtrees
      - Create a new node with value = 'OP'

# Building Abstract Syntax Trees

*This function creates an AST node and adds information that stores the value of an op in the node. A reference to the AST node is returned.*

```
TreeNode* OP(Scanner* s, TOKEN tok) {  
    TreeNode* ret = NULL;  
    TOKEN tok = s->GetNextToken();  
    if((tok == ADD) || (tok == SUB) || (tok ==  
    MUL) || (tok == DIV))  
        ret = CreateTreeNode(tok.val);  
    return ret;  
}
```

# Building Abstract Syntax Trees

*This function sets the references to left- and right- expression subtrees if those subtrees are valid FPEs. Returns reference to the AST node corresponding to the op value, NULL otherwise.*

```
TreeNode* E2(Scanner* s, TOKEN tok) {  
    TOKEN nxtTok = s->GetNextToken();  
    if(nxtTok == LPAREN) {  
        TreeNode* left = E(s); if(!left) return NULL;  
        TreeNode* root  = OP(s); if(!root) return NULL;  
        TreeNode* right = E(s); if(!right) return NULL;  
        nxtTok = s->GetNextToken();  
        if(nxtTok != RPAREN); return NULL;  
        //set left and right as children of root.  
        return root;  
    }  
}
```

# Building Abstract Syntax Trees

*E needs to change because E1, E2, and OP return a TreeNode\*.  
E returns a TreeNode\* now.*

```
TreeNode* E(Scanner* s) {  
  
    TOKEN* prevToken = s->GetCurTokenSequence();  
    TreeNode* ret = E1(s);  
    if(!ret) {  
        s->SetCurTokenSequence(prevToken);  
        ret = E2(s);  
    }  
    return ret;  
}
```

# Building Abstract Syntax Trees

```
TreeNode* IsTerm(Scanner* s, TOKEN tok) {
    TreeNode* ret = NULL;
    TOKEN nxtToken = s->GetNextToken();
    if(nxtToken == tok)
        ret = CreateTreeNode(nxtToken.val);
    return ret;
}

TreeNode* E1(Scanner* s) {
    return IsTerm(s, INTLITERAL);
}

TreeNode* E2(Scanner* s) {
    TOKEN nxtTok = s->GetNextToken();
    if(nxtTok == LPAREN) {
        TreeNode* left = E(s);
        if(!left) return NULL;
        TreeNode* root = OP(s);
        if(!root) return NULL;
        TreeNode* right = E(s);
        if(!right) return NULL;
        nxtTok = s->GetNextToken();
        if(nxtTok != RPAREN); return NULL;
        //set left and right as children of root.
        return root;
    }
}
```

# Building Abstract Syntax Trees

```
TreeNode* OP(Scanner* s) {
    TreeNode* ret = NULL;
    TOKEN tok = s->GetNextToken();
    if((tok == ADD) || (tok == SUB) || (tok == MUL) || (tok == DIV))
        ret = CreateTreeNode(tok.val);
    return ret;
}

TreeNode* E(Scanner* s) {
    TOKEN* prevToken = s->GetCurTokenSequence();
    TreeNode* ret = E1(s);
    if(!ret) {
        s->SetCurTokenSequence(prevToken);
        ret = E2(s);
    }
    return ret;
}
```

***Start the parser by invoking E().  
Value returned is the root of the AST.***

# Exercise

- Did we build the AST bottom-up or top-down?



# Identifying Semantic Actions for FPE Grammar

- What did we do when we saw an INTLITERAL?
  - Create a TreeNode
  - Initialize it with a value (string equivalent of INTLITERAL in this case)
  - Return a pointer to TreeNode

# Identifying Semantic Actions for FPE Grammar

- What did we do when we saw an E (parenthesized expression)?
  - Create an AST node with two children. The node contains the binary operator OP stored as a string. Children point to roots of subtrees representing E.