
Syntax Analysis – Part III

(Top-Down Parsers)

Yongjun Park
Hanyang University

From Last Time: Predictive Parsing

- **LL(1) grammar:**
 - For a given non-terminal, the lookahead symbol uniquely determines the production to apply
 - Top-down parsing = predictive parsing
 - Driven by predictive parsing table of
 - non-terminals x terminals \rightarrow productions



From Last Time: Parsing with Table

$$S \rightarrow ES'$$

$$S' \rightarrow \varepsilon \mid +S$$

$$E \rightarrow \text{num} \mid (S)$$

Partly-derived String

$\rightarrow ES'$

$\rightarrow (S)S'$

$\rightarrow (ES')S'$

$\rightarrow (1S')S'$

$\rightarrow (1+ES')S'$

$\rightarrow (1+2S')S'$

Lookahead

(

1

1

+

2

+

parsed part **unparsed part**

$(1+2+(3+4))+5$

$(1+2+(3+4))+5$

$(1+2+(3+4))+5$

$(1+2+(3+4))+5$

$(1+2+(3+4))+5$

$(1+2+(3+4))+5$

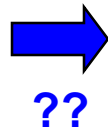
| | num | + | (|) | \$ |
|----|--------------------------|------------------|-------------------|---------------------------|---------------------------|
| S | $\rightarrow ES'$ | | $\rightarrow ES'$ | | |
| S' | | $\rightarrow +S$ | | $\rightarrow \varepsilon$ | $\rightarrow \varepsilon$ |
| E | $\rightarrow \text{num}$ | | $\rightarrow (S)$ | | |



How to Construct Parsing Tables?

Needed: Algorithm for automatically generating a predictive parse table from a grammar

$S \rightarrow ES'$
 $S' \rightarrow \varepsilon \mid +S$
 $E \rightarrow \text{number} \mid (S)$



| | num | + | (|) | \$ |
|----|-----|----|-----|---------------|---------------|
| S | ES' | | ES' | | |
| S' | | +S | | ε | ε |
| E | num | | (S) | | |

Constructing Parse Tables

- **Can construct predictive parser if:**
 - For every non-terminal, every lookahead symbol can be handled by at most 1 production
- **FIRST(β) for an arbitrary string of terminals and non-terminals β is:**
 - Set of symbols that might begin the fully expanded version of β

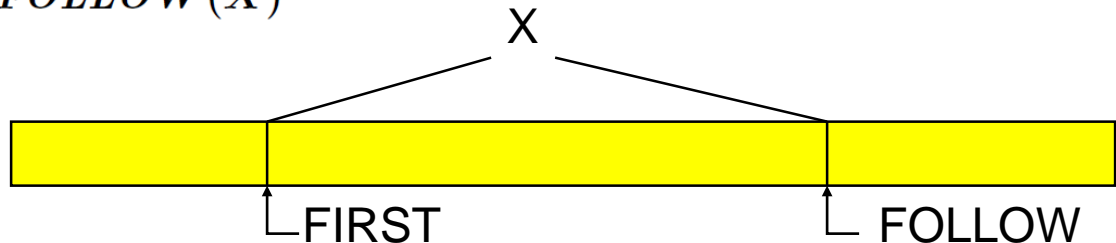
If $\alpha \Rightarrow^* c\beta$, then $c \in FIRST(\alpha)$.

If $\alpha \Rightarrow^* \epsilon$, $\epsilon \in FIRST(\alpha)$.

- **FOLLOW(X) for a non-terminal X is:**
 - Set of symbols that might follow the derivation of X in the input stream

If $S \Rightarrow^* \alpha X a \beta$, then $a \in FOLLOW(X)$.

If $S \Rightarrow^* \alpha X$, $\$ \in FOLLOW(X)$



Parse Table Entries

- Consider a production $X \rightarrow \beta$
- Add $\rightarrow \beta$ to the X row for each symbol in $\text{FIRST}(X)$
- If β can derive ε (β is nullable), add $\rightarrow \beta$ for each symbol in $\text{FOLLOW}(X)$
- Grammar is LL(1) if no conflicting entries

$S \rightarrow ES'$
 $S' \rightarrow \varepsilon \mid +S$
 $E \rightarrow \text{number} \mid (S)$

| | num | + | (|) | \$ |
|----|-----|----|-----|---------------|---------------|
| S | ES' | | ES' | | |
| S' | | +S | | ε | ε |
| E | num | | (S) | | |

Computing Nullable

- **X is nullable if it can derive the empty string:**
 - If it derives ε directly ($X \rightarrow \varepsilon$)
 - If it has a production $X \rightarrow YZ \dots$ where all RHS symbols (Y,Z) are nullable
- **Algorithm: assume all non-terminals are non-nullable, apply rules repeatedly until no change**

$S \rightarrow ES'$

$S' \rightarrow \varepsilon \mid +S$

$E \rightarrow \text{number} \mid (S)$

Only S' is nullable

Computing FIRST

- **Determining FIRST(X)**
 1. if X is a terminal, then add X to FIRST(X)
 2. if $X \rightarrow \varepsilon$ then add ε to FIRST(X)
 3. if X is a nonterminal and $X \rightarrow Y_1Y_2...Y_k$ then a is in FIRST(X) if a is in FIRST(Y_i) and ε is in FIRST(Y_j) for $j = 1...i-1$ (i.e., its possible to have an empty prefix $Y_1 ... Y_{i-1}$)
 4. if ε is in FIRST($Y_1Y_2...Y_k$) then ε is in FIRST(X)

FIRST Example

$$S \rightarrow ES'$$

$$S' \rightarrow \varepsilon \mid +S$$

$$E \rightarrow \text{number} \mid (S)$$

Apply rule 1: $\text{FIRST}(\text{num}) = \{\text{num}\}$, $\text{FIRST}(+) = \{+\}$, etc.

Apply rule 2: $\text{FIRST}(S') = \{\varepsilon\}$

Apply rule 3: $\text{FIRST}(S) = \text{FIRST}(E) = \{\}$

$$\text{FIRST}(S') = \text{FIRST}('+') + \{\varepsilon\} = \{\varepsilon, +\}$$

$$\text{FIRST}(E) = \text{FIRST}(\text{num}) + \text{FIRST}('(') = \{\text{num}, (\}$$

Rule 3 again: $\text{FIRST}(S) = \text{FIRST}(E) = \{\text{num}, (\}$

$$\text{FIRST}(S') = \{\varepsilon, +\}$$

$$\text{FIRST}(E) = \{\text{num}, (\}$$



Computing FOLLOW

- **Determining FOLLOW(X)**
 1. if S is the start symbol then \$ is in FOLLOW(S)
 2. if $A \rightarrow \alpha B \beta$ then add all $\text{FIRST}(\beta) \neq \varepsilon$ to FOLLOW(B)
 3. if $A \rightarrow \alpha B$ or ($\alpha B \beta$ and ε is in $\text{FIRST}(\beta)$) then add FOLLOW(A) to FOLLOW(B)



FOLLOW Example

$S \rightarrow ES'$

$S' \rightarrow \varepsilon \mid +S$

$E \rightarrow \text{number} \mid (S)$

$\text{FIRST}(S) = \{\text{num}, (\}$

$\text{FIRST}(S') = \{\varepsilon, + \}$

$\text{FIRST}(E) = \{\text{num}, (\}$

Apply rule 1: $\text{FOL}(S) = \{\$ \}$

Apply rule 2: $S \rightarrow ES' \text{ FOL}(E) += \{\text{FIRST}(S') - \varepsilon\} = \{+\}$

$S' \rightarrow \varepsilon \mid +S$ -

$E \rightarrow \text{num} \mid (S) \text{ FOL}(S) += \{\text{FIRST}('(') - \varepsilon\} = \{ \$,) \}$

Apply rule 3: $S \rightarrow ES' \text{ FOL}(E) += \text{FOL}(S) = \{+, \$,)\}$

(because S' is nullable)

$\text{FOL}(S') += \text{FOL}(S) = \{ \$,) \}$

Putting it all Together

$\text{FIRST}(S) = \{\text{num}, (\}$

$\text{FIRST}(S') = \{\epsilon, + \}$

$\text{FIRST}(E) = \{\text{num}, (\}$

$\text{FOLLOW}(S) = \{ \$,) \}$

$\text{FOLLOW}(S') = \{ \$,) \}$

$\text{FOLLOW}(E) = \{ +,), \$ \}$

❖ Consider a production $X \rightarrow \beta$

❖ Add $\rightarrow \beta$ to the X row for each symbol in $\text{FIRST}(X)$

❖ If β can derive ϵ (β is nullable), add $\rightarrow \beta$ for each symbol in $\text{FOLLOW}(X)$

$S \rightarrow ES'$

$S' \rightarrow \epsilon \mid +S$

$E \rightarrow \text{number} \mid (S)$

| | num | + | (|) | \$ |
|----|-----|----|-----|------------|------------|
| S | ES' | | ES' | | |
| S' | | +S | | ϵ | ϵ |
| E | num | | (S) | | |

Ambiguous Grammars

Construction of predictive parse table for ambiguous grammar results in conflicts in the table (ie 2 or more productions to apply in same cell)

$$S \rightarrow S + S \mid S * S \mid \text{num}$$

$$\text{FIRST}(S+S) = \text{FIRST}(S*S) = \text{FIRST}(\text{num}) = \{ \text{num} \}$$



Class Problem

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid \text{num} \mid \varepsilon$

1. Compute FIRST and FOLLOW sets for this G
2. Compute parse table entries

Top-Down Parsing Up to This Point

- **Now we know**
 - How to build parsing table for an LL(1) grammar (ie FIRST/FOLLOW)
 - How to construct recursive-descent parser from parsing table
 - Call tree = parse tree
- **Open question – Can we generate the AST?**

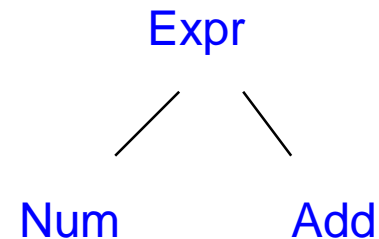


Creating the Abstract Syntax Tree

- Some class definitions to assist with AST construction

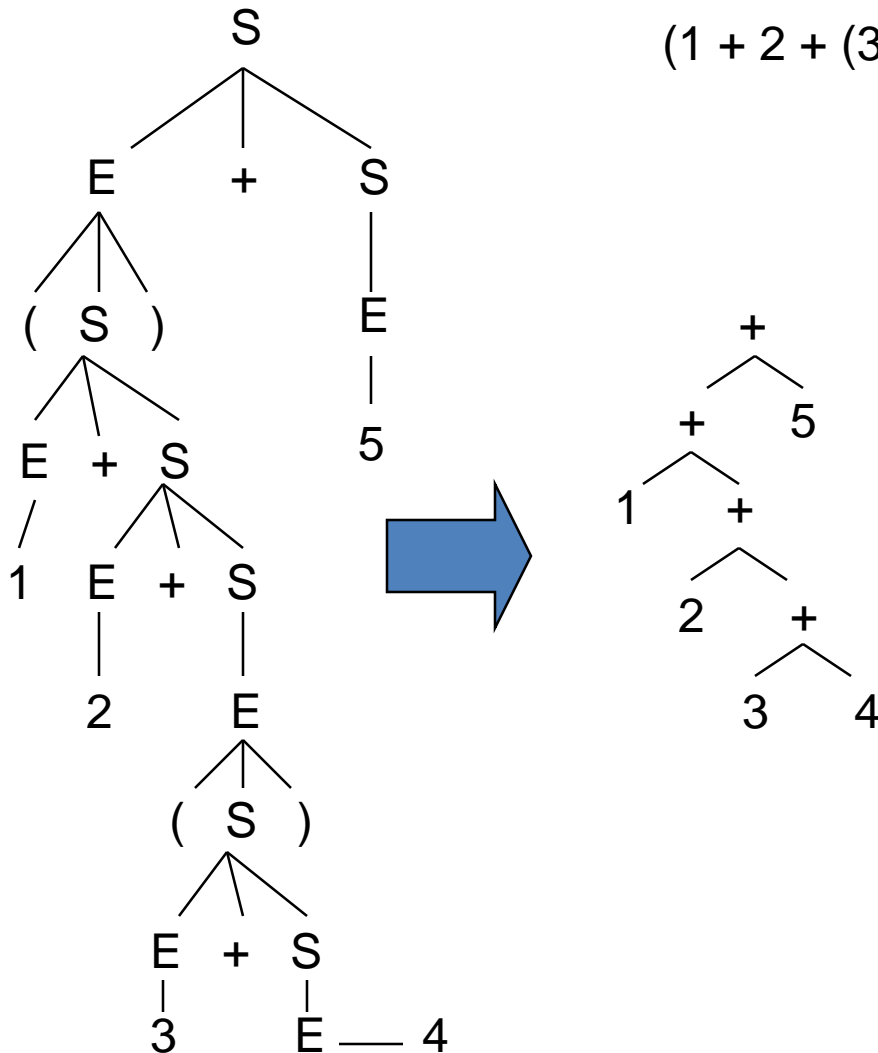
- **class Expr {}**
- **class Add extends Expr {**
 - Expr left, right;
 - Add(Expr L, Expr R) {
 - left = L; right = R;
 - }
- **}**
- **class Num extends Expr {**
 - int value;
 - Num(int v) {value = v;}
- **}**

Class Hierarchy



Creating the AST

$(1 + 2 + (3 + 4)) + 5$



- We got the parse tree from the call tree

- Just add code to each parsing routine to create the appropriate nodes

- Works because parse tree and call tree are the same shape, and AST is just a compressed form of the parse tree

AST Creation: parse_E

- **Expr parse_E() {**
 - switch (token) {
 - case num: // E → number
 - Expr result = Num(token.value);
 - token = input.read(); return result;
 - case '(': // E → (S)
 - token = input.read();
 - Expr result = parse_S();
 - if (token != ')') ParseError();
 - token = input.read(); return result;
 - default: ParseError();
 - }
- }

Remember, this is lookahead token

$S \rightarrow ES'$
 $S' \rightarrow \varepsilon \mid +S$
 $E \rightarrow \text{number} \mid (S)$

AST Creation: parse_S

- **Expr parse_S() {**

- switch (token) {

- case num:

- case '(': // S → ES'

- Expr left = parse_E();

- Expr right = parse_S'();

- if (right == NULL) return left;

- else return new Add(left,right);

- default: ParseError();

- }

- }

$S \rightarrow ES'$
 $S' \rightarrow \varepsilon \mid +S$
 $E \rightarrow \text{number} \mid (S)$

Grammars

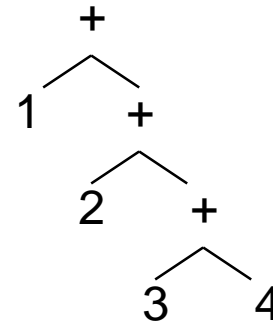
- Have been using grammar for language “sums with parentheses” $(1+2+(3+4))+5$
- Started with simple, right-associative grammar
 - $S \rightarrow E + S \mid E$
 - $E \rightarrow \text{num} \mid (S)$
- Transformed it to an LL(1) by left factoring:
 - $S \rightarrow ES'$
 - $S' \rightarrow \varepsilon \mid +S$
 - $E \rightarrow \text{num} \mid (S)$
- What if we start with a left-associative grammar?
 - $S \rightarrow S + E \mid E$
 - $E \rightarrow \text{num} \mid (S)$

Reminder: Left vs Right Associativity

Consider a simpler string on a simpler grammar: “1 + 2 + 3 + 4”

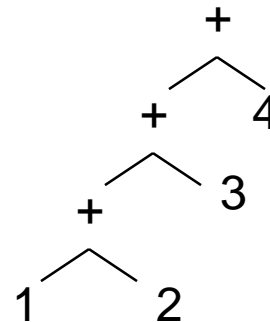
Right recursion : right associative

$S \rightarrow E + S$
 $S \rightarrow E$
 $E \rightarrow \text{num}$



Left recursion : left associative

$S \rightarrow S + E$
 $S \rightarrow E$
 $E \rightarrow \text{num}$



Left Recursion

$S \rightarrow S + E$

$S \rightarrow E$

$E \rightarrow \text{num}$

“1 + 2 + 3 + 4”

| derived string | lookahead | read/unread |
|----------------|-----------|-------------|
| S | 1 | 1+2+3+4 |
| S+E | 1 | 1+2+3+4 |
| S+E+E | 1 | 1+2+3+4 |
| S+E+E+E | 1 | 1+2+3+4 |
| E+E+E+E | 1 | 1+2+3+4 |
| 1+E+E+E | 2 | 1+2+3+4 |
| 1+2+E+E | 3 | 1+2+3+4 |
| 1+2+3+E | 4 | 1+2+3+4 |
| 1+2+3+4 | \$ | 1+2+3+4 |

Is this right? If not, what's the problem?

Left-Recursive Grammars

- Left-recursive grammars don't work with top-down parsers: we don't know when to stop the recursion
- **Left-recursive grammars are NOT LL(1)!**
 - $S \rightarrow S\alpha$
 - $S \rightarrow \beta$
- **In parse table**
 - Both productions will appear in the predictive table at row S in all the columns corresponding to $\text{FIRST}(\beta)$

Eliminate Left Recursion

- **Replace**

- $X \rightarrow X\alpha_1 \mid \dots \mid X\alpha_m$
- $X \rightarrow \beta_1 \mid \dots \mid \beta_n$

- **With**

- $X \rightarrow \beta_1X' \mid \dots \mid \beta_nX'$
- $X' \rightarrow \alpha_1X' \mid \dots \mid \alpha_mX' \mid \varepsilon$

- **See complete algorithm in Dragon book**

Class Problem

Transform the following grammar to eliminate left recursion:

$$E \rightarrow E + T \mid T$$
$$T \rightarrow T * F \mid F$$
$$F \rightarrow (E) \mid \text{id}$$

Class Problem

Transform the following grammar to eliminate left recursion:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

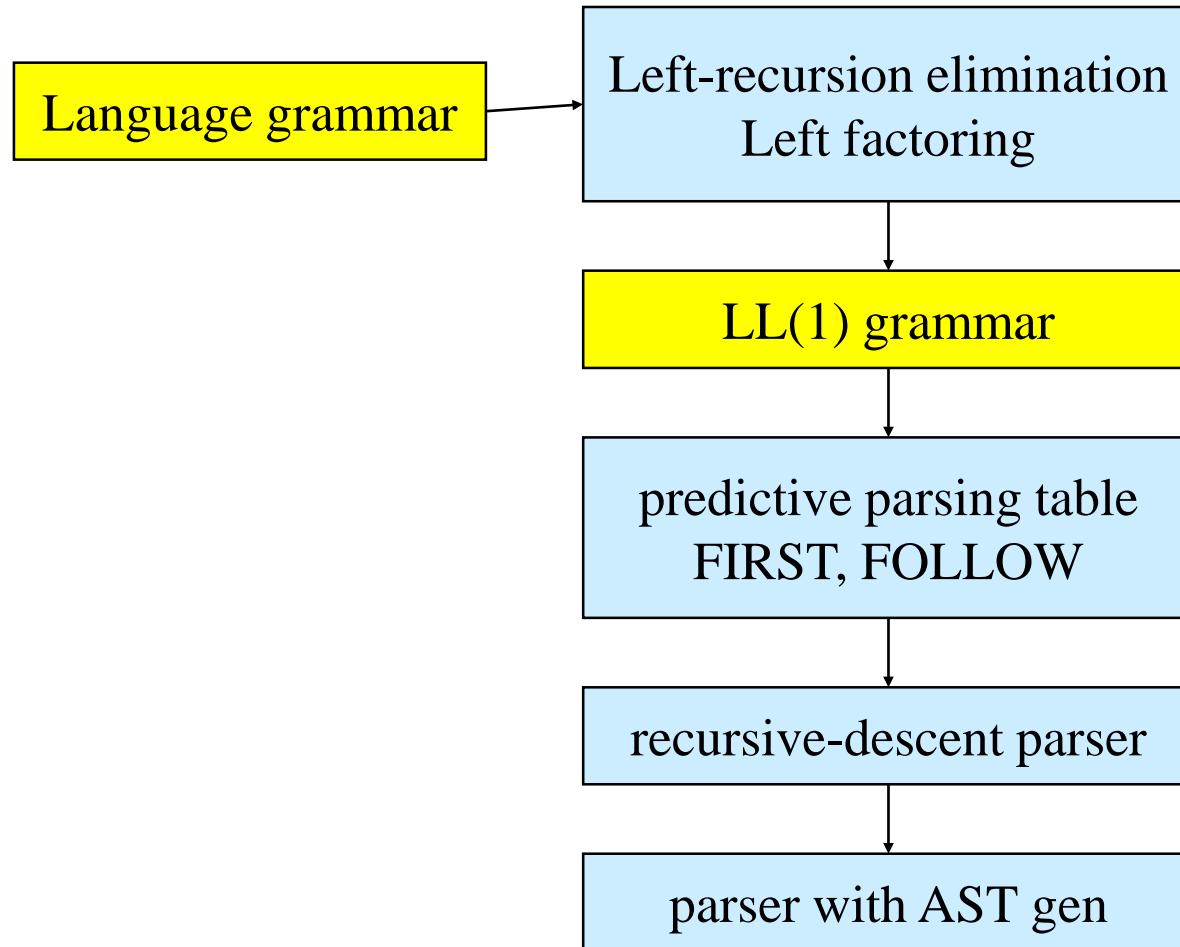
$$F \rightarrow (E) \mid \text{id}$$

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

Creating an LL(1) Grammar

- **Start with a left-recursive grammar**
 - $S \rightarrow S + E$
 - $S \rightarrow E$
 - and apply left-recursion elimination algorithm
 - $S \rightarrow ES'$
 - $S' \rightarrow +ES' \mid \varepsilon$
- **Start with a right-recursive grammar**
 - $S \rightarrow E + S$
 - $S \rightarrow E$
 - and apply left-factoring to eliminate common prefixes
 - $S \rightarrow ES'$
 - $S' \rightarrow +S \mid \varepsilon$

Top-Down Parsing Summary



Next Topic: Bottom-Up Parsing

- **A more power parsing technology**
- **LR grammars – more expressive than LL**
 - Construct right-most derivation of program
 - Left-recursive grammars, virtually all programming languages are left-recursive
 - Easier to express syntax
- **Shift-reduce parsers**
 - Parsers for LR grammars
 - Automatic parser generators (yacc, bison)