**CS 4120**
**Introduction to Compilers**

Ross Tate
Cornell University

Lecture 5: Top-down parsing

---

# Parsing Top-down

$$S \rightarrow E + S \mid E$$
$$E \rightarrow \mathbf{num} \mid (\,S\,)$$

- **Goal:** construct a leftmost derivation of string while reading in token stream
- Partly-derived String  Lookahead  parsed part unparsed part

| | | |
|---|---|---|
| $S$ | ( | (1+2+(3+4))+5 |
| $\rightarrow E$+$S$ | ( | (1+2+(3+4))+5 |
| $\rightarrow$ ($S$) +$S$ | 1 | (1+2+(3+4))+5 |
| $\rightarrow$ ($E$+$S$)+$S$ | 1 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+$S$)+$S$ | 2 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+$E$+$S$)+$S$ | 2 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2+$S$)+$S$ | ( | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2+$E$)+$S$ | ( | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2+($S$))+$S$ | 3 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2+($E$+$S$))+$S$ | 3 | (1+2+(3+4))+5 |

2

---

# Problem

$$S \rightarrow E + S \mid E$$
$$E \rightarrow \mathbf{num} \mid (\,S\,)$$

- Want to decide which production to apply based on next symbol

- (1)    $S \rightarrow E \rightarrow (S) \rightarrow (E) \rightarrow (1)$
- (1)+2  $S \rightarrow E + S \rightarrow (S) + S \rightarrow (E) + S$
  $\rightarrow (1)+E \rightarrow (1)+2$

- *Why is this hard?*

CS 4120 Introduction to Compilers

3

---

# Grammar is Problematic

- This grammar cannot be parsed top-down with only a single look-ahead symbol
- Not **LL(1)**
- **L**eft-to-right-scanning, **L**eft-most derivation, **1** look-ahead symbol
- Is it LL(k) for some k?
- Can rewrite grammar to allow top-down parsing: create LL(1) grammar for same language

CS 4120 Introduction to Compilers

4

---

# Making a grammar LL(1)

$S \rightarrow E + S$
$S \rightarrow E$
$E \rightarrow \mathbf{num}$
$E \rightarrow (\,S\,)$

⬇

$S \rightarrow ES'$
$S' \rightarrow \varepsilon$
$S' \rightarrow + S$
$E \rightarrow \mathbf{num}$
$E \rightarrow (\,S\,)$

- **Problem**: can't decide which $S$ production to apply until we see symbol after first expression

- **Left factoring**: Factor common $S$ prefix, add new non-terminal $S'$ at decision point. $S'$ derives $(+E)*$

CS 4120 Introduction to Compilers

5

---

# Parsing with new grammar

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

| | | |
|---|---|---|
| $S$ | ( | (1+2+(3+4))+5 |
| $\rightarrow E$ $S'$ | ( | (1+2+(3+4))+5 |
| $\rightarrow$ ($S$) $S'$ | 1 | (1+2+(3+4))+5 |
| $\rightarrow$ ($E$ $S'$) $S'$ | 1 | (1+2+(3+4))+5 |
| $\rightarrow$ (1 $S'$) $S'$ | + | (1+2+(3+4))+5 |
| $\rightarrow$ (1+$E$ $S'$) $S'$ | 2 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 $S'$) $S'$ | + | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 + $S$) $S'$ | ( | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 + $E$ $S'$) $S'$ | ( | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 + ($S$) $S'$) $S'$ | 3 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 + ($E$ $S'$) $S'$) $S'$ | 3 | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 + (3 $S'$) $S'$) $S'$ | + | (1+2+(3+4))+5 |
| $\rightarrow$ (1+2 + (3 + $E$) $S'$) $S'$ | 4 | (1+2+(3+4))+5 |

6

# Predictive Parsing

- **LL(1)** grammar:
  - for a given non-terminal, the look-ahead symbol uniquely determines the production to apply
  - uses predictive parsing
  - driven by *predictive parsing table* of
    non-terminals **x** input symbols → productions

CS 4120 Introduction to Compilers

7

# Predictive Parse Table

$$
\begin{aligned}
S &\to E\,S' \\
S' &\to \varepsilon \mid +\,S \\
E &\to \mathbf{num} \mid (\,S\,)
\end{aligned}
$$

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

8

# How to Implement?

- Table can be converted easily into a
  ***recursive-descent parser***

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

- Three procedures: parse_S, parse_S', parse_E

CS 4120 Introduction to Compilers

9

# Recursive-Descent Parser

```
void parse_S () {          lookahead token
    switch (token) {
        case num: parse_E(); parse_S'(); return;
        case '(': parse_E(); parse_S'(); return;
        default: throw new ParseError();
    }
}
```

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

CS 4120 Introduction to Compilers

10

# Recursive-Descent Parser

```
void parse_S'() {
    switch (token) {
        case '+': token = input.read(); parse_S(); return;
        case ')': return;
        case EOF: return;
        default: throw new ParseError();
    }
}
```

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

CS 4120 Introduction to Compilers

11

# Recursive-Descent Parser
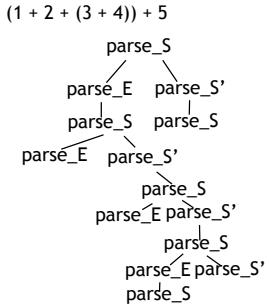
```
void parse_E() {
    switch (token) {
        case number: token = input.read(); return;
        case '(': token = input.read(); parse_S();
            if (token != ')') throw new ParseError();
            token = input.read(); return;
        default: throw new ParseError(); }
}
```

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

CS 4120 Introduction to Compilers

12

2

## Call Tree = Parse Tree

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

parse_S

parse_E  parse_S'

parse_S  parse_S

parse_E  parse_S'

parse_S

parse_E parse_S'

parse_S

parse_E parse_S'

parse_S

$S$
$E$   $S'$
$(\ S\ )$  $+\ S$
$E$  $S'$    5
1   $+\ S$
$E$  $S'$
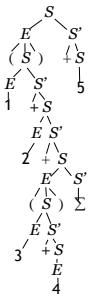2   $+\ S$
$E$   $S'$
$(\ S\ )$  $\Sigma$
$E$  $S'$
3   $+\ S$
$E$
4

CS 4120 Introduction to Compilers

13

---

## How to Construct Parsing Tables

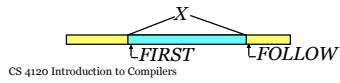- Needed: algorithm for automatically generating a predictive parse table from a grammar

$S \rightarrow E\ S'$
$S' \rightarrow \varepsilon \mid +\ S$
$E \rightarrow \textbf{num} \mid (\ S\ )$

?

| | num | + | ( | ) | EOF |
|---|---|---|---|---|---|
| **S** | ES' | | ES' | | |
| **S'** | | +S | | ε | ε |
| **E** | num | | (S) | | |

CS 4120 Introduction to Compilers

14

---

## Constructing Parse Tables

- Can construct predictive parser if:
  – For every non-terminal, every look-ahead symbol can be handled by at most one production
- FIRST($\gamma$) for arbitrary string of terminals and non-terminals $\gamma$ is:
  - set of symbols that might begin a fully expanded version of $\gamma$
- FOLLOW(X) for a non-terminal $X$ is:
  - set of symbols that might follow the derivation of X in the input stream

$X$

$\llcorner$FIRST  $\llcorner$FOLLOW

CS 4120 Introduction to Compilers

15

---

## Computing nullable and FIRST

- $X$ is nullable if it can derive the empty string:
  - if it derives ε directly ($X \rightarrow \varepsilon$)
  - if it has a production $X \rightarrow YZ...$ where all RHS symbols ($Y$, $Z$) are nullable
  - **Algorithm:** Assume all non-terminals non-nullable, apply rules repeatedly until no change in status
- Determining FIRST($\gamma$)
  – $FIRST(X) \supseteq FIRST(\gamma)$   if $X \rightarrow \gamma$
  – $FIRST(\textbf{a}\ \beta) = \{\ \textbf{a}\ \}$
  – $FIRST(X\ \beta) \supseteq FIRST(X)$
  – $FIRST(X\ \beta) \supseteq FIRST(\beta)$ if $X$ is nullable
  – **Algorithm:** Assume $FIRST(\gamma) = \{\ \}$ for all $\gamma$, apply rules repeatedly to build $FIRST$ sets.

CS 4120 Introduction to Compilers

16

---

## Computing FOLLOW

- $FOLLOW(S) \supseteq \{\ \$\ \}$
- If $X \rightarrow \alpha Y \beta$,
      $FOLLOW(Y) \supseteq FIRST(\beta)$
- If $X \rightarrow \alpha Y \beta$ and $\beta$ is nullable (or non-existent),
      $FOLLOW(Y) \supseteq FOLLOW(X)$
- **Algorithm:** Assume $FOLLOW(X) = \{\ \}$ for all $X$, apply rules repeatedly to build $FOLLOW$ sets

- Common theme: iterative analysis. Start with initial assignment, apply rules until no change

CS 4120 Introduction to Compilers

17

---

## Example

$S \rightarrow E\ S'$
$S' \rightarrow \varepsilon \mid +\ S$
$E \rightarrow \textbf{num} \mid (\ S\ )$

- nullable
  - only $S'$ is nullable
- FIRST
  – $FIRST(E\ S') = \{\ \textbf{num},\ (\ \}$
  – $FIRST(+S) = \{\ +\ \}$
  – $FIRST(\textbf{num}) = \{\ \textbf{num}\ \}$
  – $FIRST(\ (S)\ ) = \{\ (\ \},$    $FIRST(S') = \{\ +\ \}$
- FOLLOW
  – $FOLLOW(S) = \{\ \$,\ )\ \}$
  – $FOLLOW(S') = \{\ \$,\ )\ \}$
  – $FOLLOW(E) = \{\ +,\ ),\ \$\}$

18

3

## Parse Table Entries

- Consider a production $X \rightarrow \gamma$

$$
\begin{aligned}
S &\rightarrow E\,S' \\
S' &\rightarrow \varepsilon \mid + S \\
E &\rightarrow \mathbf{num} \mid (\,S\,)
\end{aligned}
$$

- Add $\gamma$ to the $X$ row for each symbol in FIRST($\gamma$)

|    | num  | + | ( | ) | EOF |
|----|------|---|---|---|-----|
| S  | ES'  |   | ES' |   |     |
| S' |      | +S |   | ε | ε   |
| E  | num  |   | (S) |   |     |

- If $\gamma$ can derive $\varepsilon$ ($\gamma$ is *nullable*), add $\gamma$ for each symbol in FOLLOW($X$)
- Grammar is LL(1) if no conflicting entries

CS 4120 Introduction to Compilers

19

## Ambiguous grammars

- Construction of predictive parse table for ambiguous grammar results in *conflicts* (*but converse does not hold*)

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

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

|   | **num** | + | * |
|---|---------|---|---|
| $S$ | **num**, $S + S$, $S * S$ |   |   |

20

## Completing the Parser

- Now we know how to construct a recursive-descent parser for an LL(1) grammar.

- LL(k) generalizes this to k lookahead tokens.

- LL(k) parser generators can be used to automate the process (e.g. ANTLR)

- Can we use recursive descent to build an abstract syntax tree too?

CS 4120 Introduction to Compilers

21

## Creating the AST

```
abstract 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);
}
```
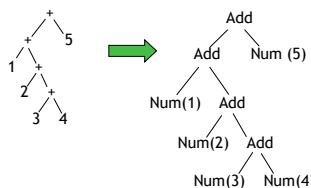
Expr
Num   Add

CS 4120 Introduction to Compilers

22

## AST Representation

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



How to generate this structure during recursive-descent parsing?

CS 4120 Introduction to Compilers

23

## Creating the AST

- Just add code to each parsing routine to create the appropriate nodes!
- Works because parse tree and call tree have same shape
- parse_S, parse_S', parse_E all return an Expr:

```
void parse_E()  ⟹ Expr parse_E()
void parse_S()  ⟹ Expr parse_S()
void parse_S'() ⟹ Expr parse_S'()
```

CS 4120 Introduction to Compilers

24

## AST creation code

```
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 != ')') throw new ParseError();
        token = input.read(); return result;
   default: throw new ParseError();
 }
}
```

## parse_S

```
Expr parse_S() {
    switch (token) {
        case num:
        case '(':
            Expr left = parse_E();
            Expr right = parse_S'();
            if (right == null) return left;
            else return new Add(left, right);
        default: throw new ParseError();
    }
}
```

$$S \rightarrow E\,S'$$
$$S' \rightarrow \varepsilon \mid + S$$
$$E \rightarrow \mathbf{num} \mid (\,S\,)$$

## Or…an Interpreter!

```
int parse_E() {
    switch(token) {
        case number:
            int result = token.value;
            token = input(); return result;
        case '(':
            token = input.read();
            int result = parse_S();
            if (token != ')') throw new ParseError();
            token = input.read(); return result;
        default: throw new ParseError(); }}

int parse_S() {
    switch (token) {
        case number:
        case '(':
            int left = parse_E();
            int right = parse_S'();
            if (right == 0) return left;
            else return left + right;
        default: throw new ParseError(); } }
```

## Summary

- We can build a recursive-descent parser for LL(1) grammars
  - Make parsing table from *FIRST*, *FOLLOW* sets
  - Translate to recursive-descent code
  - Instrument with abstract syntax tree creation
- Systematic approach avoids errors, detects ambiguities
- Next time: converting a grammar to LL(1) form, bottom-up parsing