

LL(1) Parsing, Refactoring and Recursive Descent

CSCI 3136: Principles of Programming Languages

Agenda

- Announcements
 - Assignment 4 is out and due June 14 (hopefully)
 - Midterm on Wednesday, June 19, in CHEB 170, 10 – 11:30am.
 - Final Exam on Friday, August 2, in CHEB 170, 1:00 – 4:00pm
- Readings:
 - Today: 2.3.0, 2.3.1
 - Note: I recommend using alternative texts for this part of the course:
 - E..g, Hopcorft et al, “Introduction to Automata Theory”
- Lecture Contents
 - Building an LL(1) Parser
 - The PREDICT Table
 - Constructing FIRST, FOLLOW, and PREDICT
 - Is a Grammar LL(1)?
 - Refactoring
 - Recursive Descent

Building an LL(1) Parser

- Basic Challenge: Given current token, which production does the parser select if next item in sentential form is a nonterminal

E.g., if S is on the stack and input is $+$, then parser must select production $S \rightarrow +SS$

- In general: for input \mathbf{a} , sentential form $A \dots$, either
 - $A \Rightarrow \alpha \Rightarrow^* \mathbf{a}\beta$
 - $A \Rightarrow \alpha \Rightarrow^* \epsilon$ and derivation of A is succeeded by \mathbf{a} .
- Intuitively, \mathbf{a} is in the *predictor set* of $A \rightarrow \alpha$
if $A\beta \Rightarrow \alpha\beta \Rightarrow^* \mathbf{a}\gamma$, for $\beta, \gamma \in \Sigma^*$
I.e., the parser selects $A \rightarrow \alpha$ if \mathbf{a} is the input and in the *predictor set* of $A \rightarrow \alpha$

LL(1) Grammars

- **Definition:** A grammar is LL(1) if the predictor sets of all productions with the same LHS are disjoint.
- E.g. S-Grammars are LL(1)
Grammar
 1. $S \rightarrow + SS$
 2. $S \rightarrow - SS$
 3. $S \rightarrow * SS$
 4. $S \rightarrow / SS$
 5. $S \rightarrow \text{neg } S$
 6. $S \rightarrow \text{integer}$

PREDICT Table

Production	Predictor Set
$S \rightarrow + SS$	{+}
$S \rightarrow - SS$	{-}
$S \rightarrow * SS$	{*}
$S \rightarrow / SS$	{/}
$S \rightarrow \text{neg } S$	{neg}
$S \rightarrow \text{integer}$	{integer}

Constructing PREDICT: The 3 Tables

- **FIRST(α)**: the set of leftmost terminals **a** to be derived from $\alpha \in (V \cup \Sigma)^*$
 $\alpha \Rightarrow^* a\beta$
- **FOLLOW(X)**: the set of first terminals **a** that follows variable X in a derivation
 $S \Rightarrow^* \alpha X a \beta$
- **PREDICT($A \rightarrow \alpha$)**: the set of terminals that predict this production given **A**



The FIRST Table: Example

Grammar

- $T \rightarrow AB$
- $A \rightarrow PQ$
- $A \rightarrow BC$
- $P \rightarrow pP$
- $P \rightarrow \epsilon$
- $Q \rightarrow qQ$
- $Q \rightarrow \epsilon$
- $B \rightarrow bB$
- $B \rightarrow e$
- $C \rightarrow cC$
- $C \rightarrow f$

		Sym: σ	FIRST(σ)
Σ		p	{p}
		q	{q}
		b	{b}
		e	{e}
		c	{c}
		f	{f}
		T	{p,q,b,e}
V		A	{p,q,b,e, ϵ }
		P	{p, ϵ }
		Q	{q, ϵ }
		B	{b,e}
		C	{c,f}


The FIRST Table

- **Definition:** $\text{FIRST}(\sigma)$, $\forall \sigma \in (V \cup \Sigma)^*$:

- For $a \in \Sigma$, $a \in \text{FIRST}(\sigma)$ if $\sigma \Rightarrow^* a\beta$
- $\varepsilon \in \text{FIRST}(\sigma)$ if $\sigma \Rightarrow^* \varepsilon$

- **Notes:**

- For $a \in \Sigma$, $\text{FIRST}(a) = \{a\}$
- Precompute $\text{FIRST}(X)$ only for $X \in V$
- Generate $\text{FIRST}(\sigma)$, $\sigma \in (V \cup \Sigma)^*$ as needed



Precompute FIRST sets for all terminals and variables

- **To compute FIRST for terminals and variables**

- For $a \in \Sigma$, $\text{FIRST}(a) = \{a\}$
- For $X \in V$, $\text{FIRST}(X) = \emptyset$
- Repeat until no new additions to $\text{FIRST}(X)$, $X \in V$ are possible:
 $\forall (X \rightarrow \alpha) \in P$,
 $\text{FIRST}(X) = \text{FIRST}(X) \cup \text{FIRST}(\alpha)$

The FIRST Table (Part 2)

- **To compute FIRST(α) (in general)**

- $\alpha = \alpha_1 \alpha_2 \dots \alpha_k, \alpha_i \in V \cup \Sigma$
- $\text{FIRST}(\alpha) = \emptyset$
- For $i = 1, 2, \dots, k$:
 - $\text{FIRST}(\alpha) = \text{FIRST}(\alpha) \cup (\text{FIRST}(\alpha_i) \setminus \{\epsilon\})$
 - if $\epsilon \notin \text{FIRST}(\alpha_i)$ then return
- $\text{FIRST}(\alpha) = \text{FIRST}(\alpha) \cup \epsilon$

The FIRST Table: Example

Loop until no more changes:

$$\forall X \rightarrow \alpha \in P,$$

$$\text{FIRST}(X) = \text{FIRST}(X) \cup \text{FIRST}(\alpha)$$

Grammar

- $T \rightarrow AB$
- $A \rightarrow PQ$
- $A \rightarrow BC$
- $P \rightarrow pP$
- $P \rightarrow \epsilon$
- $Q \rightarrow qQ$
- $Q \rightarrow \epsilon$
- $B \rightarrow bB$
- $B \rightarrow e$
- $C \rightarrow cC$
- $C \rightarrow f$

Symbol	Iter. 0	Iter. 1	Iter.2	Iter. 3
p	{p}	{p}	{p}	{p}
q	{q}	{q}	{q}	{q}
b	{b}	{b}	{b}	{b}
e	{e}	{e}	{e}	{e}
c	{c}	{c}	{c}	{c}
f	{f}	{f}	{f}	{f}
T	\emptyset	\emptyset	\emptyset	{p,q,b,e}
A	\emptyset	\emptyset	{p,q,b,e, ϵ }	{p,q,b,e, ϵ }
P	\emptyset	{p, ϵ }	{p, ϵ }	{p, ϵ }
Q	\emptyset	{q, ϵ }	{q, ϵ }	{q, ϵ }
B	\emptyset	{b,e}	{b,e}	{b,e}
C	\emptyset	{c,f}	{c,f}	{c,f}

The FOLLOW Table

- **Definition:** $\text{FOLLOW}(X), \forall X \in V$:
 - For $a \in \Sigma$, $a \in \text{FOLLOW}(X)$ if $\exists S \Rightarrow^* \alpha X a \beta$
 - $\epsilon \in \text{FOLLOW}(X)$ if $S \Rightarrow^* \alpha X$
- **To Compute FOLLOW**
 - $\text{FOLLOW}(S) = \{\epsilon\}$
 - For $X \in V$, $\text{FOLLOW}(X) = \emptyset$
 - Repeat until no new additions to $\text{FOLLOW}(X)$, $X \in V$ are possible:
 - For each $B \in V$
 - $\forall (X \rightarrow \alpha B \beta) \in P$,
 - $\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup (\text{FIRST}(\beta) \setminus \{\epsilon\})$
 - if $\epsilon \in \text{FIRST}(\beta)$ then $\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup \text{FOLLOW}(X)$

The FOLLOW Table: Example

Grammar

- $T \rightarrow AB$
- $A \rightarrow PQ$
- $A \rightarrow BC$
- $P \rightarrow pP$
- $P \rightarrow \epsilon$
- $Q \rightarrow qQ$
- $Q \rightarrow \epsilon$
- $B \rightarrow bB$
- $B \rightarrow e$
- $C \rightarrow cC$
- $C \rightarrow f$

Symbol	Iter. 0	Iter. 1	Iter.2
T	$\{\epsilon\}$	$\{\epsilon\}$	$\{\epsilon\}$
A	\emptyset	$\{b,e\}$	$\{b,e\}$
P	\emptyset	$\{q\}$	$\{q,b,e\}$
Q	\emptyset	\emptyset	$\{b,e\}$
B	\emptyset	$\{\epsilon,c,f\}$	$\{\epsilon,c,f\}$
C	\emptyset	\emptyset	$\{b,e\}$

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

For $X \in V$, $\text{FOLLOW}(X) = \emptyset$

Repeat until no new additions to $\text{FOLLOW}(X)$, $X \in V$ are possible:

For each $B \in V$

$\forall X \rightarrow \alpha B \beta \in P$,

$\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup (\text{FIRST}(\beta) \setminus \{\epsilon\})$

if $\epsilon \in \text{FIRST}(\beta)$ then $\text{FOLLOW}(B) = \text{FOLLOW}(B) \cup \text{FOLLOW}(X)$

The PREDICT Table

- **Definition:** For $a \in \Sigma \cup \{\epsilon\}$, $a \in \text{PREDICT}(A \rightarrow \alpha)$ if
 - $a \in \text{FIRST}(\alpha) \setminus \{\epsilon\}$ or
 - $\epsilon \in \text{FIRST}(\alpha)$ and $a \in \text{FOLLOW}(A)$
- **To Compute PREDICT**
 - For each $(A \rightarrow \alpha) \in P$, $\text{PREDICT}(A \rightarrow \alpha) = \emptyset$
 - For each $(A \rightarrow \alpha) \in P$
 - $\text{PREDICT}(A \rightarrow \alpha) = \text{FIRST}(\alpha) \setminus \{\epsilon\}$
 - if $\epsilon \in \text{FIRST}(\alpha)$ then
$$\text{PREDICT}(A \rightarrow \alpha) = \text{PREDICT}(A \rightarrow \alpha) \cup \text{FOLLOW}(A)$$

The PREDICT Table: Example

Symbol	FIRST	FOLLOW
T	{p,q,b,e}	{ ϵ }
A	{p,q,b,e, ϵ }	{b,e}
P	{p, ϵ }	{q,b,e}
Q	{q, ϵ }	{b,e}
B	{b,e}	{ ϵ ,c,f}
C	{c,f}	{b,e}

For each $(A \rightarrow \alpha) \in P$
 $PREDICT(A \rightarrow \alpha) = FIRST(\alpha) \setminus \{\epsilon\}$
 if $\epsilon \in FIRST(\alpha)$ then
 $PREDICT(A \rightarrow \alpha) = PREDICT(A \rightarrow \alpha) \cup FOLLOW(A)$

Since the predictor sets overlap for A productions, this is not an LL(1) grammar

Production	Predictor Set
$T \rightarrow AB$	{p,q,b,e}
$A \rightarrow PQ$	{p,q,b,e}
$A \rightarrow BC$	{b,e}
$P \rightarrow pP$	{p}
$P \rightarrow \epsilon$	{q,b,e}
$Q \rightarrow qQ$	{q}
$Q \rightarrow \epsilon$	{b,e}
$B \rightarrow bB$	{b}
$B \rightarrow e$	{e}
$C \rightarrow cC$	{c}
$C \rightarrow f$	{f}

How to Prove a Grammar is LL(1)

- Construct PREDICT Table for the grammar
- This grammar is not LL(1) if and only if two productions with the same left hand side have non-disjoint predictor sets.
- Note: It's actually possible to build the FIRST, FOLLOW, and PREDICT tables by simply looking at the grammar.
- What happens if our grammar is not LL(1)?

Limitations and Problems with LL(1)

- There exist context free languages that do not have LL(1) grammars
- There is no known algorithm to determine whether a language is LL(1)
- There is an algorithm to decide whether a grammar is LL(1) (we just saw it)
- Most obvious grammars for most programming languages are usually not LL(1)
- In many cases a non-LL(1) grammar can be refactored into an LL(1) grammar

Refactoring Grammars

- Two common problems:
Which production do you use? (Both have α in FIRST)
- Left recursion
 - $A \rightarrow A\beta$
 - $A \rightarrow \alpha$
- Common Prefix
 - $A \rightarrow \alpha\beta$
 - $A \rightarrow \alpha\gamma$

Dealing with Left Recursion

- Idea: Replace Left Recursion with Right Recursion

$$\begin{array}{l} A \rightarrow A\beta \\ A \rightarrow \alpha \end{array}$$



$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta A' \\ A' \rightarrow \varepsilon \end{array}$$

- Note: As a side-effect the grammar may cease to capture some properties such as left-associativity

Example of Eliminating Left Recursion

- Consider the grammar fragment:

Block \rightarrow '{' Statements '}'

Statements \rightarrow Statements Statement

Statements $\rightarrow \epsilon$

- Replace this with:

Block \rightarrow '{' Statements '}'

Statements \rightarrow Statement Statements

Statements $\rightarrow \epsilon$

Dealing with Common Prefix

- Idea: Remove common prefix by left factoring

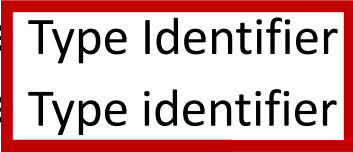



$$\begin{array}{l} A \rightarrow \alpha\beta \\ A \rightarrow \alpha\gamma \end{array}$$



$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta \\ A' \rightarrow \gamma \end{array}$$

Example of Eliminating Common Prefix

Bad Grammar

Field \rightarrow Type Identifier ';' 
Field \rightarrow Type identifier '(' Args ')' ';' 
Type \rightarrow Identifier 
Type \rightarrow Identifier Array 
Array \rightarrow '[' ']' Array
Array $\rightarrow \epsilon$

Better grammar

Field \rightarrow Type Identifier FieldBody ';'
FieldBody \rightarrow '(' Args ')'
FieldBody $\rightarrow \epsilon$
Type \rightarrow Identifier Array
Array \rightarrow '[' ']' Array
Array $\rightarrow \epsilon$

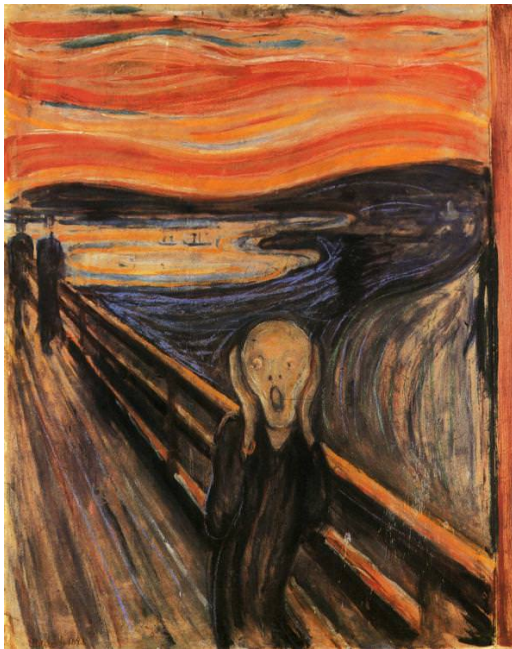
LL(1) Parser Implementation

- Two efficient approaches:
 - Recursive Descent
 - Deterministic Pushdown Automata (DPDA)
- Recursive Descent is easier to understand and implement.

Recursive Descent

Use PREDICT
table

Idea: For each variable
 X , write a procedure:
`parse_X()`



```
parse_X:
    t = peek_next_token()
    select  $X \rightarrow \alpha$  based on t
    for each  $\alpha_i \in \alpha$  :
        if  $\alpha_i == Y_1 \in V$ :
            parse_Y1()
        elseif  $\alpha_i == Y_2 \in V$ :
            parse_Y2()

    ...
    elseif  $\alpha_i == Y_k \in V$ :
        parse_Yk()
    elseif  $\alpha_i == t$ :
        remove_token()
        t = peek_next_token()
    else:
        syntax error
```

Example

Grammar

- $S \rightarrow \text{Add} \mid \text{Sub} \mid \text{Mul}$
- $S \rightarrow \text{Div} \mid \text{Neg} \mid \text{Val}$
- $\text{Add} \rightarrow + S S$
- $\text{Sub} \rightarrow - S S$
- $\text{Mul} \rightarrow * S S$
- $\text{Div} \rightarrow / S S$
- $\text{Neg} \rightarrow \text{neg } S$
- $\text{Val} \rightarrow \text{integer}$

```
parse_S:
    t = peek_at_token()
    select S  $\rightarrow$   $\alpha$  based on t
    for each  $\alpha_i \in \alpha$  :
        if  $\alpha_i == \mathbf{Add} \in V$ :
            parse_Add()
        elseif  $\alpha_i == \mathbf{Sub} \in V$ :
            parse_Sub()

        ...
        elseif  $\alpha_i == \mathbf{Val} \in V$ :
            parse_Val()
elseif  $\alpha_i == t$ :
    remove_token()
    t = peek_at_token()
    else:
        syntax error
```

Example

More Concrete

Grammar

- **$S \rightarrow \text{Add} \mid \text{Sub} \mid \text{Mul}$**
- **$S \rightarrow \text{Div} \mid \text{Neg} \mid \text{Val}$**
- $\text{Add} \rightarrow + S S$
- $\text{Sub} \rightarrow - S S$
- $\text{Mul} \rightarrow * S S$
- $\text{Div} \rightarrow / S S$
- $\text{Neg} \rightarrow \text{neg } S$
- $\text{Val} \rightarrow \text{integer}$

```
parse_S:
    t = peek_at_token()
    if t == '+' :
        n = parse_Add()
    elseif t == '-':
        n = parse_Sub()
    elseif t == '*':
        n = parse_Mul()
    elseif t == '/':
        n = parse_Div()
    elseif t == 'neg':
        n = parse_Neg()
    elseif t is integer:
        n = parse_Val()
    else:
        syntax error

    return TreeNode(n)
```


Example

More Concrete

Grammar

- $S \rightarrow \text{Add} \mid \text{Sub} \mid \text{Mul}$
- $S \rightarrow \text{Div} \mid \text{Neg} \mid \text{Val}$
- **$\text{Add} \rightarrow + S S$**
- $\text{Sub} \rightarrow - S S$
- $\text{Mul} \rightarrow * S S$
- $\text{Div} \rightarrow / S S$
- $\text{Neg} \rightarrow \text{neg } S$
- $\text{Val} \rightarrow \text{integer}$

`parse_Add:`

```
t = peek_at_token()
if t != '+':
    # never happens
    syntax error
remove_token()
s1 = parse_S()
s2 = parse_S()
return TreeNode(t, s1, s2)
```