

Top-down Parsing Recursive Descent & LL(1)

Copyright 2010, Keith D. Cooper & Linda Torczon, all rights reserved.



Roadmap (Where are we?)

We set out to study parsing

- Specifying syntax
 - Context-free grammars \checkmark
- Top-down parsers
 - Algorithm & its problem with left recursion \checkmark
 - Ambiguity ✓
 - Left-recursion removal \checkmark
- Predictive top-down parsing today
 - The LL(1) Property
 - First and Follow sets
 - Simple recursive descent parsers
 - Table-driven LL(1) parsers

Predictive Parsing



Given $A \rightarrow \alpha \mid \beta$, the parser should be able to choose between $\alpha \& \beta$

Predictive Parser is a top-down parser free of backtracking



For some rhs $\alpha \in G$, FIRST(α) is set of tokens (terminals) that appear as first symbol in some string deriving from α

$$\mathbf{x} \in \mathsf{FIRST}(\alpha)$$
 iff $\alpha \Rightarrow^* \underline{\mathbf{x} \gamma}$, for some γ

Goal → SheepNoise SheepNoise → SheepNoise <u>baa</u> | <u>baa</u>

For SheepNoise: FIRST(Goal) = { <u>baa</u> } FIRST(SN) = { <u>baa</u> } FIRST(<u>baa</u>) = { <u>baa</u> }



LL(1) Property

If $A \rightarrow \alpha$ and $A \rightarrow \beta$ both appear in the grammar, we would like

 $\mathsf{FIRST}(\alpha) \cap \mathsf{FIRST}(\beta) = \emptyset$

This would allow the parser to make a correct choice with a lookahead of exactly one symbol !

Almost correct! See the next slide



Does not have LL(1) Property

What about ϵ -productions?



If $A \rightarrow \alpha$ and $A \rightarrow \beta$ and $\epsilon \in F_{IRST}(\alpha)$, then we need to ensure

 $FOLLOW(A) \cap FIRST(\beta) = \emptyset$

where,

Follow(A) = the set of terminal symbols that can immediately follow A in a sentential form

Formally,

Follow(A) = {t | (t is a terminal and $G \Rightarrow^* \alpha A \pm \beta$) or (t is eof and $G \Rightarrow^* \alpha A$)} Follow Sets







Definition of FIRST⁺($A \rightarrow \alpha$) if $\varepsilon \in FIRST(\alpha)$ then FIRST⁺($A \rightarrow \alpha$) = FIRST(α) \cup FOLLOW(A) else FIRST⁺($A \rightarrow \alpha$) = FIRST(α)

Grammar is *LL(1)* iff $A \rightarrow \alpha$ and $A \rightarrow \beta$ implies

$$\mathsf{FIRST}^{+}(A \rightarrow \alpha) \cap \mathsf{FIRST}^{+}(A \rightarrow \beta) = \emptyset$$



What If My Grammar Is Not LL(1)?

Can we transform a non-LL(1) grammar into an LL(1) grammar?

- In general, the answer is no
- In some cases, however, the answer is yes
- Perform:
 - -Eliminate left-recursion Friday

-Perform left factoring today



What If My Grammar Is Not LL(1)?

Given grammar G with productions $A \rightarrow \alpha \beta_1$ and $A \rightarrow \alpha \beta_2$

if α derives anything other than ε and FIRST⁺($A \rightarrow \alpha \beta_1$) \cap FIRST⁺($A \rightarrow \alpha \beta_2$) $\neq \emptyset$



This grammar is not LL(1)

Left Factoring



If we pull the common prefix, α , into a separate production, we may make the grammar LL(1).



Now, if FIRST⁺($A' \rightarrow \beta_1$) \cap FIRST⁺($A' \rightarrow \beta_2$) = \emptyset , G may be LL(1)

UNIVERSITY OF ELAWARE

Left Factoring



NT with common prefix

UNIVERSITY OF ELAWARE

Left Factoring

For each nonterminal A find the longest prefix α common to 2 or more alternatives for A if $\alpha \neq \varepsilon$ then replace all of the productions $A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \alpha \beta_3 | \dots | \alpha \beta_n | \gamma$ with Repeat until no NT has rhs' with a common prefix

Put common prefix α into a separate production rule

UNIVERSITY OF ELAWARE

Left Factoring

For each nonterminal A find the longest prefix α common to 2 or more alternatives for A if $\alpha \neq \varepsilon$ then replace all of the productions $A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \alpha \beta_3 | \dots | \alpha \beta_n | \gamma$ with $A \rightarrow \alpha A' | \gamma$ $A' \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n$ Repeat until no NT has rhs' with a common prefix

Create new Nonterminal (A') with all unique suffixes



Left Factoring

For each nonterminal A find the longest prefix α common to 2 or more alternatives for A if $\alpha \neq \varepsilon$ then replace all of the productions $A \rightarrow \alpha \beta_1 | \alpha \beta_2 | \alpha \beta_3 | \dots | \alpha \beta_n | \gamma$ with $A \rightarrow \alpha A' | \gamma$ $A' \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n$ Repeat until no NT has rhs' with a common prefix

Transformation makes some grammars into LL(1) grammars There are languages for which no LL(1) grammar exists 14



Left Factoring Example

Consider a simple right-recursive expression grammar

0	Goal	\rightarrow	Expr	To choose between 1, 2, & 3,
1	Expr	\rightarrow	Term + Expr	an LL(1) parser must look past
2			Term - Expr	the <u>number</u> or <u>id</u> to see the
3			Term	operator.
4	Term	\rightarrow	Factor * Term	FIRST ⁺ (1) = FIRST ⁺ (2) = FIRST ⁺ (3)
5			Factor / Term	and
6			Factor	FIRST ⁺ (4) = FIRST ⁺ (5) = FIRST ⁺ (6)
7	Factor	\rightarrow	<u>number</u>	Let's left factor this grammar.
8			<u>id</u>	

Left Factoring Example



After Left Factoring, we have

Goal	\rightarrow	Expr
Expr	\rightarrow	Term Expr'
Expr'	\rightarrow	+ Expr
		- Expr
		ε
Term	\rightarrow	Factor Term'
Term'	\rightarrow	* Term
		/ Term
		ε
Factor	\rightarrow	number
		<u>id</u>
	Goal Expr Expr' Term Term'	$\begin{array}{ccc} Goal & \rightarrow \\ Expr & \rightarrow \\ Expr' & \rightarrow \\ & & \\ & & \\ Term & \rightarrow \\ Term' & \rightarrow \\ & & \\ Factor & \rightarrow \\ & & \end{array}$

Clearly, FIRST+(2), FIRST+(3), & FIRST+(4) are disjoint, as are FIRST+(6), FIRST+(7), & FIRST+(8) The grammar now has the LL(1)

property



FIRST(α)

For some $\alpha \in (T \cup NT)^*$, define FIRST(α) as the set of tokens that appear as the first symbol in some string that derives from α

That is, $\underline{x} \in FIRST(\alpha)$ iff $\alpha \Rightarrow^* \underline{x} \gamma$, for some γ



for each $x \in T$, FIRST(x) $\leftarrow \{x\}$ for each $A \in NT$, $FIRST(A) \leftarrow \emptyset$ while (FIRST sets are still changing) do for each $p \in P$, of the form $A \rightarrow \beta$ do if β is $B_1B_2...B_k$ then begin; $rhs \leftarrow FIRST(B_1) - \{\varepsilon\}$ for $i \leftarrow 1$ to k-1 by 1 while $\varepsilon \in FIRST(B_i)$ do $rhs \leftarrow rhs \cup (FIRST(B_{i+1}) - \{\varepsilon\})$ end // for loop end //if-then if i = k and $\varepsilon \in FIRST(B_k)$ then $rhs \leftarrow rhs \cup \{\varepsilon\}$ $FIRST(A) \leftarrow FIRST(A) \cup rhs$ end // for loop // while loop end

Outer loop is monotone increasing for FIRST sets $\rightarrow | T \cup NT \cup \varepsilon |$ is bounded, so it terminates Inner loop is bounded by the length of the productions in the grammar

Set empty set for terminals



for each $x \in T$, FIRST(x) $\leftarrow \{x\}$ Outer loop is monotone for each $A \in NT$, FIRST(A) $\leftarrow \mathcal{O}^{4}$ increasing for FIRST while (FIRST sets are still changing) do sets for each $p \in P$, of the form $A \rightarrow \beta$ do \rightarrow $\top \cup \mathsf{NT} \cup \varepsilon \mid \mathsf{is}$ if β is $B_1B_2...B_k$ then begin; bounded, so it terminates $rhs \leftarrow FIRST(B_1) - \{\varepsilon\}$ for $i \leftarrow 1$ to k-1 by 1 while $\varepsilon \in FIRST(B_i)$ do Inner loop is bounded by the length of the $rhs \leftarrow rhs \cup (FIRST(B_{i+1}) - \{\varepsilon\})$ productions in the end // for loop grammar end //if-then if i = k and $\varepsilon \in FIRST(B_k)$ then $rhs \leftarrow rhs \cup \{\varepsilon\}$ $FIRST(A) \leftarrow FIRST(A) \cup rhs$ Set empty set for end // for loop First of nonterminals // while loop end



for each $x \in T$, FIRST(x) $\leftarrow \{x\}$ for each $A \in NT$, $FIRST(A) \leftarrow \emptyset$ while (FIRST sets are still changing) do for each $p \in P$, of the form $A \rightarrow \beta$ do if β is $B_1B_2...B_k$ then begin; $rhs \leftarrow FIRST(B_1) - \{\varepsilon\}$ for $i \leftarrow 1$ to k-1 by 1 while $\varepsilon \in FIRST(B_i)$ do $rhs \leftarrow rhs \cup (FIRST(B_{i+1}) - \{\varepsilon\})$ end // for loop end //if-then if i = k and $\varepsilon \in FIRST(B_k)$ then $rhs \leftarrow rhs \cup \{\varepsilon\}$ $FIRST(A) \leftarrow FIRST(A) \cup rhs$ end // for loop // while loop end

Outer loop is monotone increasing for FIRST sets → | T ∪ NT ∪ ε | is

bounded, so it terminates

Inner loop is bounded by the length of the productions in the grammar

Fixed point algorithm; Monotone because we always add to First sets



for each $x \in T$, FIRST(x) $\leftarrow \{x\}$ Outer loop is monotone for each $A \in NT$, $FIRST(A) \leftarrow \emptyset$ increasing for FIRST while (FIRST sets are still changing) do sets for each $p \in P$, of the form $A \rightarrow \beta$ do \rightarrow | T \cup NT $\cup \varepsilon$ | is if β is $B_1B_2...B_k$ then begin; bounded, so it terminates $rhs \leftarrow FIRST(B_1) - \{\varepsilon\}$ for $i \leftarrow 1$ to k-1 by 1 while $\varepsilon \in FIRST(B_i)$ do Inner loop is bounded by the length of the $rhs \leftarrow rhs \cup (FIRST(B_{i+1}) - \{\varepsilon\})$ productions in the end // for loop grammar end //if-then if i = k and $\varepsilon \in FIRST(B_k)$ then $rhs \leftarrow rhs \cup \{\varepsilon\}$ $FIRST(A) \leftarrow FIRST(A) \cup rhs$ Iterate through each end // for loop production // while loop end



for each $x \in T$, FIRST(x) $\leftarrow \{x\}$ Outer loop is monotone for each $A \in NT$, $FIRST(A) \leftarrow \emptyset$ increasing for FIRST while (FIRST sets are still changing) do sets for each $p \in P$, of the form $A \rightarrow \beta$ do \rightarrow | T \cup NT $\cup \varepsilon$ | is if β is $B_1B_2...B_k$ then begin;* bounded, so it terminates $rhs \leftarrow FIRST(B_1) - \{\varepsilon\}$ for $i \leftarrow 1$ to k-1 by 1 while $\varepsilon \in FIRST(B_i)$ do Inner loop is bounded by the length of the $rhs \leftarrow rhs \cup (FIRST(B_{i+1}) - \{\varepsilon\})$ productions in the end // for loop grammar end //if-then if i = k and $\varepsilon \in FIRST(B_k)$ then $rhs \leftarrow rhs \cup \{\varepsilon\}$ $FIRST(A) \leftarrow FIRST(A) \cup rhs$ end // for loop RHS is some set of T // while loop end and NT.



for each x ∈ T, FIRST(x) ← {x} for each A ∈ NT, FIRST(A) ← Ø	Outer loop is monotone increasing for FIRST		
while (FIRST sets are still changing) do for each p∈P. of the form A→B do	sets		
if β is $B_1B_2B_k$ then begin; $rhs \leftarrow FIRST(B_1) - \{\varepsilon\}$	\rightarrow 1 \bigcirc N1 $\bigcirc \varepsilon$ 15 bounded, so it terminates		
for $i \leftarrow 1$ to k-1 by 1 while $\varepsilon \in FIRST(B_i)$ do	Inner loop is bounded		
end // for loop	productions in the		
end // if-then if $i = k$ and $\varepsilon \in FIRST(B_k)$	granmar		
then rhs \leftarrow rhs $\cup \{\varepsilon\}^{n}$ ETDST(A) \leftarrow ETDST(A) \cup rhc			
end // for loop	Initialize rhs to First		
end // while loop	of first symbol minus epsilon		



for each $x \in T$, FIRST $(x) \leftarrow \{x\}$ for each $A \in NT$, FIRST $(A) \leftarrow \emptyset$ while (FIRST sets are still changing) do for each $p \in P$, of the form $A \rightarrow \beta$ do if β is $B_1B_2B_k$ then begin; rhs \leftarrow FIRST $(B_1) - \{\varepsilon\}$	Outer loop is monotone increasing for FIRST sets $\rightarrow T \cup NT \cup \varepsilon $ is bounded, so it terminates
for $i \leftarrow 1$ to $k-1$ by 1 while $\varepsilon \in FIRST(B_i)$ do $rhs \leftarrow rhs \cup (FIRST(B_{i+1}) - \{\varepsilon\})$ end // for loop end // if-then $if \ i = k \ and \ \varepsilon \in FIRST(B_k)$ $then \ rhs \leftarrow rhs \cup \{\varepsilon\}$	Inner loop is bounded by the length of the productions in the grammar
FIRST(A) ← FIRST(A) ∪ rhs end // for loop end // while loop	Iterate through symbols in production until have a symbol

until have a symbol that does not have epsilon in First set



FIRST

<u>num</u>

<u>id</u>

+

*

<u>eof</u>

ε

<u>num, id, (</u>

<u>num, id, (</u>

+, -, ε

<u>num, id, (</u>

*,/,ε

<u>num, id, (</u>

Term'

Factor

Expression Grammar

0	Goal	\rightarrow	Expr	Symbo
1	Expr	\rightarrow	Term Expr'	<u>num</u>
2	Expr'	\rightarrow	+ Term Expr'	<u>id</u>
3			- Term Expr'	+
4			ε	-
5	Term	\rightarrow	Factor Term'	*
6	Term'	\rightarrow	* Factor Term'	
7		Ι	/ Factor Term') T
8		I	ε	eof
9	Factor	\rightarrow	number	<u></u> З
10			id	Goal
11			(Expr)	Expr
				Expr'
				Term