



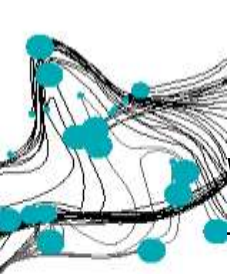
COMPILER CONSTRUCTION:

CC 2-1: PREDICTIVE PARSING, LL(1), FIRST

MODULE 8: PROGRAMMING PARADIGMS

30 APRIL 2019





WHERE ARE WE?

- **Overall structure:** Compiler front end
 - Scanning
 - Parsing ← To be continued today
 - Type checking
- **Seen last time:** Parse trees
 - Correspondence to derivations
 - Grammar ambiguity
 - Operator associativity
- **Today:** Bottom-up versus top-down parsing
 - LL(1) criterion
 - FIRST function for lookahead

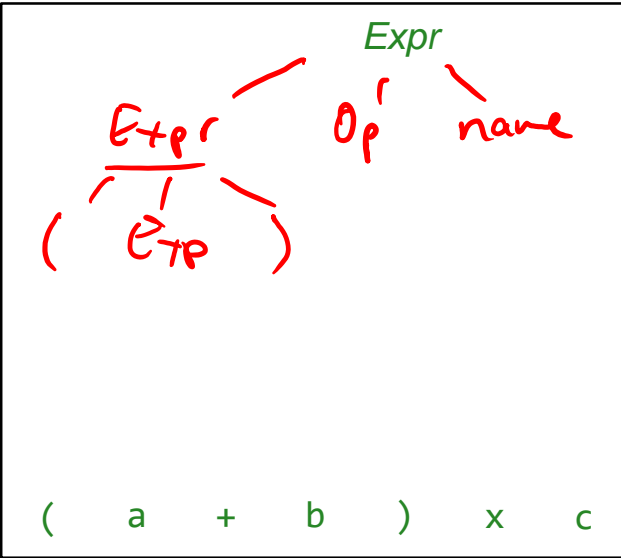


BOTTOM-UP VS TOP-DOWN PARSING

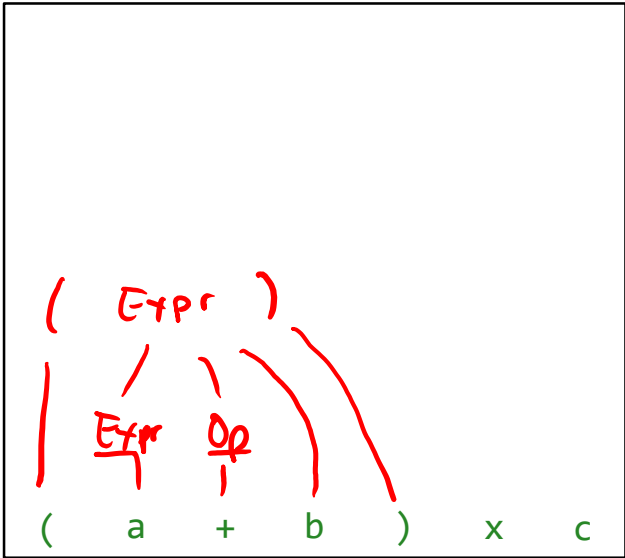
1	<i>Expr</i>	→	(<i>Expr</i>)
2			<i>Expr Op</i> name
3			name
4	<i>Op</i>	→	+
5			-
6			×
7			÷

- The task of finding a parse tree starts with
 - Start symbol
 - Text to be parsed
- Two fundamental ways to detect the structure in between

Top-down



Bottom-up



PREDICTIVE PARSING: LOOK-AHEAD

- To be avoided: backtracking
 - Introduces complexity and overhead
- Wanted: the ability to “guess” correctly which rule to apply
 - Called *predictive parsing*
 - Achieved by looking at first symbol(s) in token stream: *look-ahead*
- Consider top-down parsing of $1 - 3 + 4 - 2$

- Right-recursive grammar

$Expr \rightarrow Operand + Expr \mid Operand - Expr \mid Operand$

- No good: all rules apply on first symbol
- Refactor to

$Expr \rightarrow Operand Next \textcircled{1}$

$Next \rightarrow + Expr \textcircled{2} \mid - Expr \textcircled{3} \mid \epsilon \textcircled{4}$

- This works: apply rules $\textcircled{1} \textcircled{3} \textcircled{1} \textcircled{2} \textcircled{1} \textcircled{3} \textcircled{1} \textcircled{4}$

$Expr$	\rightarrow	Operand + $Expr$
		Operand - $Expr$
		Operand

PREDICTIVE PARSING: LOOK-AHEAD

- To be avoided: backtracking
 - Introduces complexity and overhead
- Wanted: the ability to “guess” correctly which rule to apply
 - Called *predictive parsing*
 - Achieved by looking at first symbol(s) in token stream: *look-ahead*
- Consider top-down parsing of 1 - 3 + 4 - 2
 - Left-recursive grammar

$Expr \rightarrow Expr + Operand \mid Expr - Operand \mid Operand$

- No good: how deep is the tree?

$Expr \rightarrow Operand \quad Expr'$

$Expr$	\rightarrow	$Expr + Operand$
	$ $	$Expr - Operand$
	$ $	$Operand$

$Expr' \rightarrow + Expr$
 $Expr' \rightarrow - Expr$
 $Expr' \rightarrow \epsilon$

LL(1): LEFT-TO-RIGHT, LEFTMOST, LOOKAHEAD 1

- At every state, the parser has
 - A stack of instantiated rules
 - Each element of the stack is the remaining part of a RHS
 - ▪ Initially, the stack consists only of the start symbol
 - A partially built parse tree
- Parser algorithm
 1. Remove the first symbol α from the top element of the stack
 - If $\alpha \in T$, it should be the next token on the input, otherwise report error
 - If $\alpha \in NT$, we need to know which RHS of α to choose
 - Based on the next token in the input (lookahead!)
 - Let's say $\alpha \rightarrow \beta$ with $\beta = \beta_1\beta_2 \cdots \beta_n$
 - Put $\beta_1\beta_2 \cdots \beta_n$ on the stack and in the parse tree (as children of α)
 2. Pop the stack as long as the top element is the empty string ϵ
 3. If the stack is empty, we're done

$A \rightarrow \underline{\underline{BcdE}}$

β
|
c
|
d
|
E

LL(1): LEFT-TO-RIGHT, LEFTMOST, LOOKAHEAD 1

- Parser algorithm
 - Remove the first symbol α from the top element of the stack
 - If $\alpha \in T$, it should be the next token on the input, otherwise report error
 - If $\alpha \in NT$, we need to know which RHS of α to choose
 - Based on the next token in the input (lookahead!)
 - Let's say $\alpha \rightarrow \beta$ with $\beta = \beta_1\beta_2 \cdots \beta_n$
 - Put $\beta_1\beta_2 \cdots \beta_n$ on the stack and in the parse tree (as children of α)
 - Pop the stack as long as the top element is the empty string ϵ
 - If the stack is empty, we're done
- Consider parsing of $1 - 3 + 4$

$Expr \rightarrow Operand\ Next$ ①
 $Next \rightarrow + Expr$ ②
| $- Expr$ ③
| ϵ ④

LL(1): LEFT-TO-RIGHT, LEFTMOST, LOOKAHEAD 1

- Top-down parsing with lookahead 1
 - When is a grammar suitable for this?
 - Recall: In a CFG, rules are of the form $A \rightarrow \beta$ with $A \in NT$ and $\beta \in (NT \cup T)^*$
- Algorithm: compute *FIRST*, *FOLLOW*, *FIRST+*
 - *FIRST*: $(NT \cup T)^* \rightarrow 2^T$ (function from sentential forms to sets of terminals)
 - *FIRST*(β) yields *initial* terminals with which a β -derivation may start
 - In particular applied to nonterminals (*FIRST*(A))
 - *FOLLOW*: $NT \rightarrow 2^T$ (function from non-terminals to sets of terminals)
 - *FOLLOW*(A) yields terminals that may *follow* A in a derivable sentential form
 - *FIRST+*: $Rule \rightarrow 2^T$ (function from grammar rules to sets of terminals)
 - *FIRST+*($A \rightarrow \beta$) yields *initial* terminals of A -derivations starting with this rule
- LL(1) criterion → use for "based on the next token in the input (lookahead!)"
 - *FIRST+*($A \rightarrow \beta_1$) and *FIRST+*($A \rightarrow \beta_2$) should be disjoint for distinct β_1 and β_2
 - If they overlap on terminal a , we wouldn't know what to do when finding a

EXAMPLE COMPUTATION OF LL(1) CRITERION

0	<i>Goal</i>	\rightarrow	<i>Expr</i>	6	<i>Term'</i>	\rightarrow	\times <i>Factor</i> <i>Term'</i>
1	<i>Expr</i>	\rightarrow	<i>Term</i> <i>Expr'</i>	7			\div <i>Factor</i> <i>Term'</i>
2	<i>Expr'</i>	\rightarrow	$+$ <i>Term</i> <i>Expr'</i>	8			ϵ
3			$-$ <i>Term</i> <i>Expr'</i>	9	<i>Factor</i>	\rightarrow	$($ <i>Expr</i> $)$
4			ϵ	10			num
5	<i>Term</i>	\rightarrow	<i>Factor</i> <i>Term'</i>	11			name

FIRST for terminal:
terminal itself

	num	name	+	-	\times	\div	()	eof	ϵ
FIRST	num	name	+	-	\times	\div	()	eof	ϵ

Special "fake" symbols

FIRST for non-terminal:

	<i>Expr</i>	<i>Expr'</i>	<i>Term</i>	<i>Term'</i>	<i>Factor</i>
FIRST	(, num, name	+, -, ϵ	(, num, name	\times, \div, ϵ	(, num, name

FOLLOW:
never ϵ

	<i>Expr</i>	<i>Expr'</i>	<i>Term</i>	<i>Term'</i>	<i>Factor</i>
FOLLOW					

	Production	FIRST set	FIRST ⁺ set
4	<i>Expr'</i> $\rightarrow \epsilon$		
8	<i>Term'</i> $\rightarrow \epsilon$		

FIRST⁺ equals **FIRST** except if ϵ is in **FIRST**; then take union with **FOLLOW**

EXAMPLE COMPUTATION OF LL(1) CRITERION

0	<i>Goal</i>	\rightarrow	<i>Expr</i>	6	<i>Term'</i>	\rightarrow	\times <i>Factor Term'</i>
1	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>	7			\div <i>Factor Term'</i>
2	<i>Expr'</i>	\rightarrow	$+$ <i>Term Expr'</i>	8			ϵ
3			$-$ <i>Term Expr'</i>	9	<i>Factor</i>	\rightarrow	$($ <i>Expr</i> $)$
4			ϵ	10			num
5	<i>Term</i>	\rightarrow	<i>Factor Term'</i>	11			name

FIRST for terminal:
terminal itself

	num	name	+	-	\times	\div	$($	$)$	eof	ϵ
FIRST										

FIRST for non-terminal:

	<i>Expr</i>	<i>Expr'</i>	<i>Term</i>	<i>Term'</i>	<i>Factor</i>
FIRST					

FOLLOW:
never ϵ

	<i>Expr</i>	<i>Expr'</i>	<i>Term</i>	<i>Term'</i>	<i>Factor</i>
FOLLOW					

	Production	FIRST set	FIRST ⁺ set
4	<i>Expr'</i> \rightarrow ϵ		
8	<i>Term'</i> \rightarrow ϵ		

FIRST⁺ equals **FIRST** except if ϵ is in **FIRST**; then take union with **FOLLOW**