

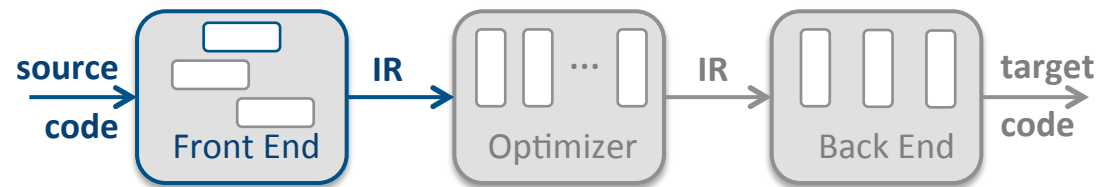


COMP 412
FALL 2017

Parsing II

Top-down parsing

Comp 412



Copyright 2017, Keith D. Cooper & Linda Torczon, all rights reserved.

Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.

Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved.

Chapter 3 in EaC2e



Definitions

- A context-free grammar G is **ambiguous** if there exists has more than one leftmost derivation for some sentence in $L(G)$
- A context-free grammar G is **ambiguous** if there exists has more than one rightmost derivation for some sentence in $L(G)$
- The leftmost and rightmost derivations for a sentential form may differ, even in an unambiguous grammar
 - *However, they must have the same parse tree*



We talked about syntactic ambiguity

- Ambiguity in the context-free syntax
 - Classic example is the **if-then-else** grammar

0	<i>Stmt</i>	→	<u>if</u> <i>Expr</i> <u>then</u> <i>Stmt</i>
1			<u>if</u> <i>Expr</i> <u>then</u> <i>Stmt</i> <u>else</u> <i>Stmt</i>
2			... other statements ...

- Fix ambiguity in context-free grammar by re-writing the grammar

0	<i>Stmt</i>	→	<u>if</u> <i>Expr</i> <u>then</u> <i>Stmt</i>
1			<u>if</u> <i>Expr</i> <u>then</u> <i>WithElse</i> <u>else</u> <i>Stmt</i>
2			... other statements ...
3	<i>WithElse</i>	→	<u>if</u> <i>Expr</i> <u>then</u> <i>WithElse</i> <u>else</u> <i>WithElse</i>
4			... other statements ...

Ambiguity



We must also deal with semantic ambiguity

- One syntax with two meanings
 - Classic example arose in Algol-like languages
 $A = f(17, 21)$
Is this a call to a function **f**? or a reference to an element of an array **f**?
- Disambiguating this kind of confusion requires context
 - Need the value of the declaration for **f**
 - An issue of type, not syntax
 - Requires either:
 1. An extra-grammatical solution
 - *Manage the ambiguity by accepting language and deferring disambiguation until the compiler has enough context (e.g., type information)*
 2. A different syntax
 - *Fix ambiguity in meaning by changing the language (e.g., C's `[]` or BCPL's `!`)*

Order of Operations or Precedence



Consider again the derivation of $x - 2 * y$

0	<i>Expr</i>	→	<i>Expr Op Value</i>
1			<i>Value</i>
2	<i>Value</i>	→	<u>number</u>
3			<u>identifier</u>
4	<i>Op</i>	→	<u>plus</u>
5			<u>minus</u>
6			<u>times</u>
7			<u>divide</u>

<i>Rule</i>	<i>Sentential Form</i>
—	<i>Expr</i>
0	<i>Expr Op Value</i>
0	<i>Expr Op Value Op Value</i>
1	<i>Value Op Value Op Value</i>
3	<id,x> <i>Op Value Op Value</i>
5	<id,x> — <i>Value Op Value</i>
2	<id,x> — <num,2> <i>Op Value</i>
6	<id,x> — <num,2> * <i>Value</i>
3	<id,x> — <num,2> * <id,y>

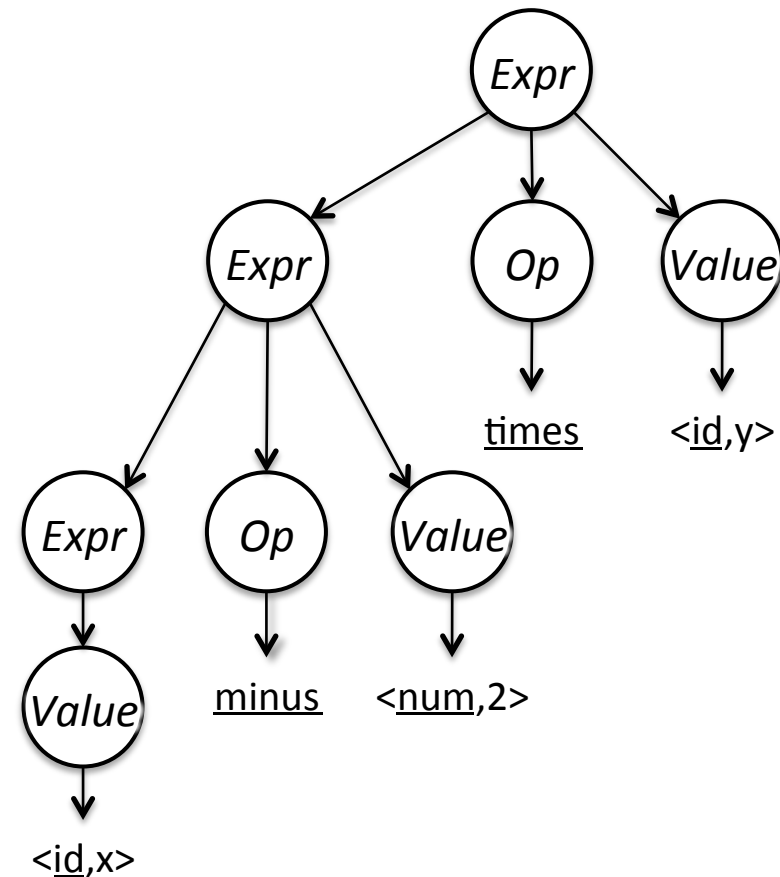
Leftmost derivation

Order of Operations



The leftmost derivation is unique, but it specifies the wrong precedence

Rule	Sentential Form
—	<i>Expr</i>
0	<i>Expr Op Value</i>
0	<i>Expr Op Value Op Value</i>
1	<i>Value Op Value Op Value</i>
3	<i><id,x> Op Value Op Value</i>
5	<i><id,x> — Value Op Value</i>
2	<i><id,x> — <num,2> Op Value</i>
6	<i><id,x> — <num,2> * Value</i>
3	<i><id,x> — <num,2> * <id,y></i>



Evaluates $(x - 2) * y$

Eliminating ambiguity does not necessarily produce the desired meaning. It produces a consistent meaning, but that meaning can be consistently wrong.

Order of Operations



How do you add precedence to a grammar?

To add precedence

- Decide how many *levels of precedence* the grammar needs
- Create a nonterminal for each level of precedence
- Isolate the corresponding part of the grammar
- Force the parser to recognize high precedence subexpressions first

For algebraic expressions, including $()$, $+$, $-$, $*$, and $/$

- Parentheses first (level 1)
- Multiplication and division, next (level 2)
- Subtraction and addition, last (level 3)

Derivations and Precedence



Adding the standard algebraic precedence produces:

	0	Goal	→	Expr
level 3	1	Expr	→	Expr + Term
	2			Expr - Term
	3			Term
level 2	4	Term	→	Term * Factor
	5			Term / Factor
	6			Factor
level 1	7	Factor	→	(Expr)
	8			<u>number</u>
	9			<u>id</u>

The new grammar is larger (7 vs. 9)

- Takes more rewriting to reach some of the terminal symbols
- Encodes expected precedence
- Produces same parse tree under leftmost & rightmost derivations
- Correctness trumps the speed of the parser

Let's see how it parses $x + 2 * y$

The “**classic expression grammar**”
See also Figure 7.7 on p. 351 of EaC2e

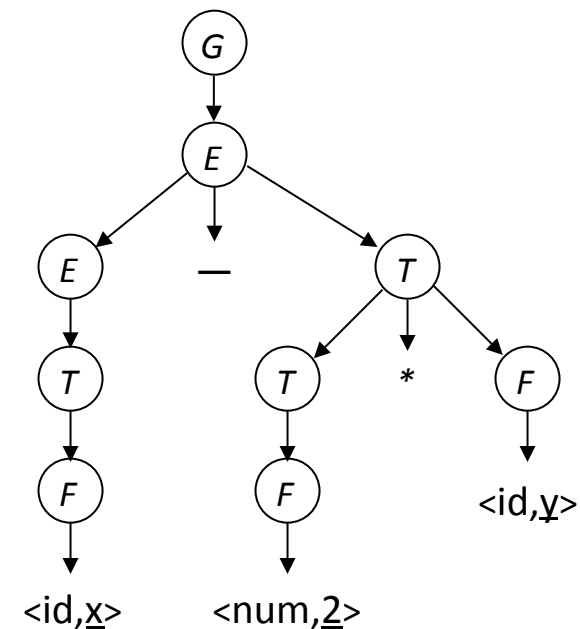
Both parentheses & precedence
are beyond the power of an RE

Derivations and Precedence



Rule	Sentential Form
—	Goal
0	Expr
2	Expr \rightarrow Term
3	Term \rightarrow Term
6	Factor \rightarrow Term
9	$\langle \text{id}, x \rangle \rightarrow$ Term
4	$\langle \text{id}, x \rangle \rightarrow$ Term $*$ Factor
6	$\langle \text{id}, x \rangle \rightarrow$ Factor $*$ Factor
8	$\langle \text{id}, x \rangle \rightarrow$ $\langle \text{num}, 2 \rangle * \text{Factor}$
9	$\langle \text{id}, x \rangle \rightarrow$ $\langle \text{num}, 2 \rangle * \langle \text{id}, y \rangle$

The leftmost derivation



Parse tree for $x - 2 * y$

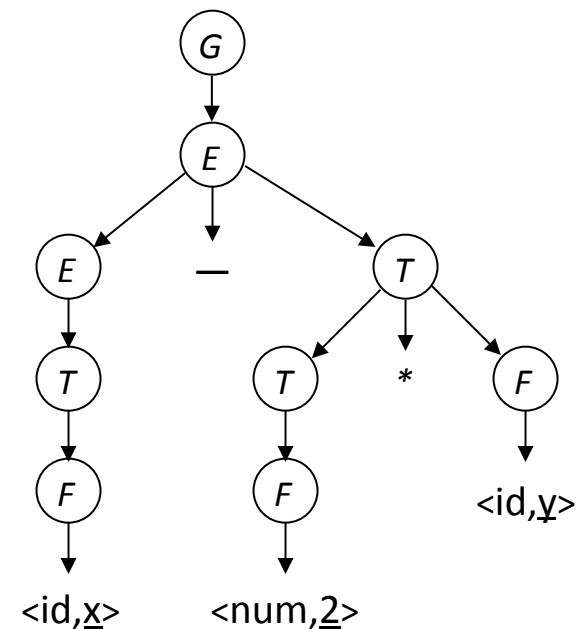
The classic expression grammar derives $\underline{x} - (\underline{2} * \underline{y})$ with the parse tree shown.. Both the leftmost and rightmost derivations give the same parse tree and value, because the grammar directly and explicitly encodes the desired precedence.

Derivations and Precedence



Rule	Sentential Form
—	Goal
0	Expr
2	Expr \rightarrow Term
4	Expr \rightarrow Term * Factor
9	Expr \rightarrow Term * <id,y>
6	Expr \rightarrow Factor * <id,y>
8	Expr \rightarrow <num,2> * <id,y>
3	Term \rightarrow <num,2> * <id,y>
6	Factor \rightarrow <num,2> * <id,y>
9	<id,x> \rightarrow <num,2> * <id,y>

The rightmost derivation



Parse tree for $x - 2 * y$

The classic expression grammar derives $\underline{x} - (\underline{2} * \underline{y})$ with the parse tree shown.. Both the leftmost and rightmost derivations give the same parse tree and value, because the grammar directly and explicitly encodes the desired precedence.

Parsing Techniques

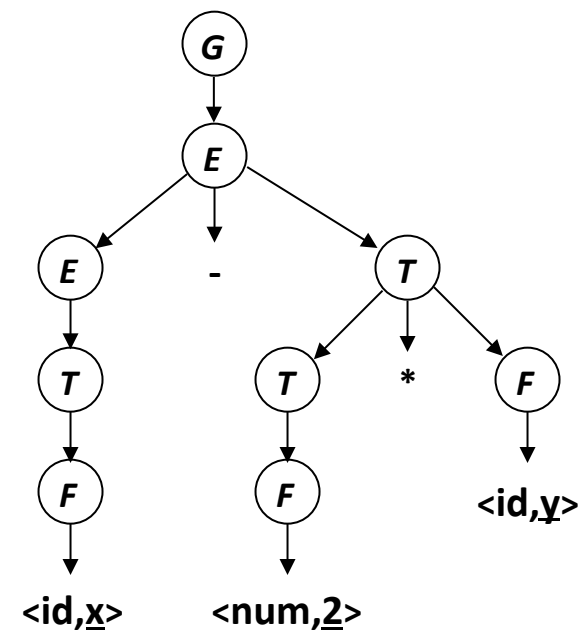


Top-down parsers (*LL(1), recursive descent*)

- Start at the root of the parse tree and grow toward leaves
- Pick a production & try to match the input
- Bad “pick” \Rightarrow may need to backtrack
- Large class of grammars are backtrack-free

Bottom-up parsers (*LR(1), operator precedence*)

- Start at the leaves and grow toward root
- As input is consumed, encode possibilities in an internal state
- Start in a state valid for legal first tokens
- We can make the process deterministic



*Parse tree for $x - 2 * y$*

Bottom-up parsers can recognize a larger class of *grammars* than can top-down parsers.

Top-Down Parsing



We will examine two ways of implementing top-down parsers



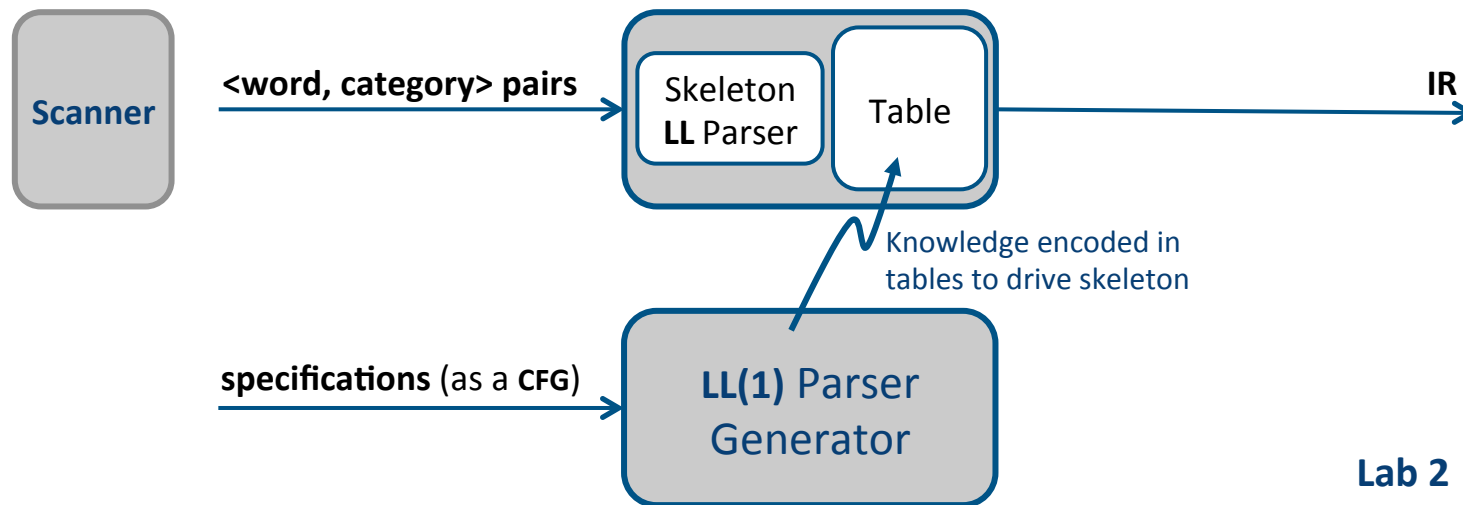
Recursive-Descent Parser

- Highly efficient, highly flexible form of parser
- Typically implemented as a hand-coded parser
 - Set of mutually-recursive routines
 - Works well for any “backtrack free” or “predictively parsable” grammar
- Easy to understand, easy to implement

Top-Down Parsing



We will examine two ways of implementing top-down parsers



Lab 2

Table-driven LL(1) Parser

- **LL(1)** Parser Generator takes as input a **CFG** that is backtrack free
- Skeleton Parser interprets the table produced by the generator
- In Lab 2, you will implement an **LL(1)** table generator
 - Your table generator will use a recursive-descent parser as its front end

Top-down Parsing



The Algorithm

- A top-down parser starts with the root of the parse tree
- The root node is labeled with the goal symbol of the grammar

Construct the root node of the parse tree

Repeat until lower fringe of the parse tree matches the input string

1. *At a node labeled with **NT** A , select a production with A on its **LHS** and, for each symbol on its **RHS**, construct the appropriate child*
2. *When a terminal symbol is added to the fringe and it doesn't match the fringe, backtrack*
3. *Find the next node to be expanded* *(label \in NT)*

The key is picking the right production in step 1

- *That choice should be guided by the input string*

The “Classic” Expression Grammar



Consider the Classic Expression Grammar

0	<i>Goal</i>	\rightarrow	<i>Expr</i>
1	<i>Expr</i>	\rightarrow	<i>Expr</i> + <i>Term</i>
2			<i>Expr</i> - <i>Term</i>
3			<i>Term</i>
4	<i>Term</i>	\rightarrow	<i>Term</i> * <i>Factor</i>
5			<i>Term</i> / <i>Factor</i>
6			<i>Factor</i>
7	<i>Factor</i>	\rightarrow	(<i>Expr</i>)
8			<u>number</u>
9			<u>id</u>

And the input $\underline{x} - \underline{2} * \underline{y}$

Example



Let's try to derive $\underline{x} - \underline{2} * \underline{y}$:

Goal

Rule	Sentential Form	Input
—	Goal	$\uparrow \underline{x} - \underline{2} * \underline{y}$

Lower fringe of the partially completed parse tree

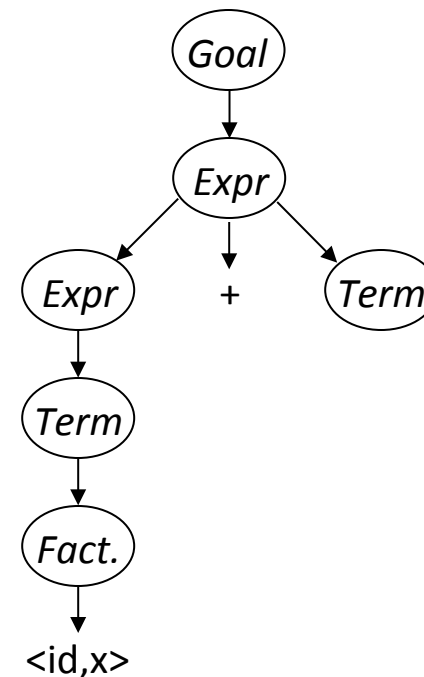
\uparrow is the position in the input buffer
Build a **leftmost** derivation, to work
with a left-to-right scanner.

Example



Let's try to derive $\underline{x} - \underline{2} * \underline{y}$:

Rule	Sentential Form	Input
—	Goal	$\uparrow \underline{x} - \underline{2} * \underline{y}$
0	Expr	$\uparrow \underline{x} - \underline{2} * \underline{y}$
1	Expr + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
3	Term + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
6	Factor + Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle \text{id}, \underline{x} \rangle + \text{Term}$	$\uparrow \underline{x} - \underline{2} * \underline{y}$
→	$\langle \text{id}, \underline{x} \rangle + \text{Term}$	$\underline{x} \uparrow - \underline{2} * \underline{y}$



This worked well, except that “−” doesn’t match “+”

The parser must backtrack to here

\uparrow is the position in the input buffer

Example

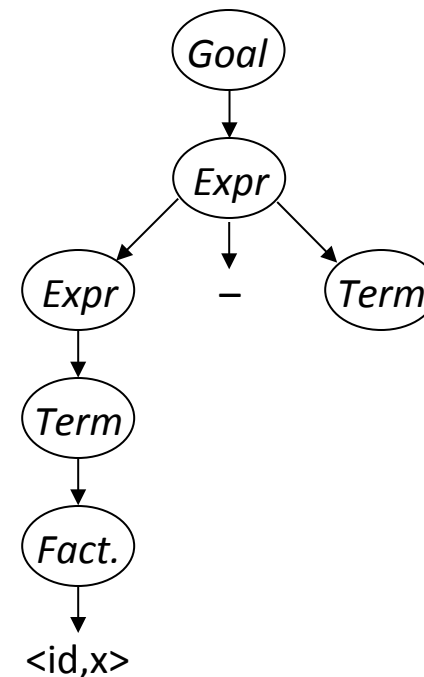


Continuing with $\underline{x} - \underline{2} * \underline{y}$:

Rule	Sentential Form	Input
—	Goal	$\uparrow \underline{x} - \underline{2} * \underline{y}$
0	Expr	$\uparrow \underline{x} - \underline{2} * \underline{y}$
<hr/>		
2	Expr - Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
3	Term - Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
6	Factor - Term	$\uparrow \underline{x} - \underline{2} * \underline{y}$
9	$\langle id, \underline{x} \rangle \ominus Term$	$\uparrow \underline{x} \ominus \underline{2} * \underline{y}$
→	$\langle id, \underline{x} \rangle - Term$	$\underline{x} \uparrow - \underline{2} * \underline{y}$
→	$\langle id, \underline{x} \rangle - Term$	$\underline{x} - \uparrow \underline{2} * \underline{y}$

Now, “-” and “-” match

Now we can expand Term to match “2”



⇒ Now, we need to expand *Term* - the last NT on the fringe



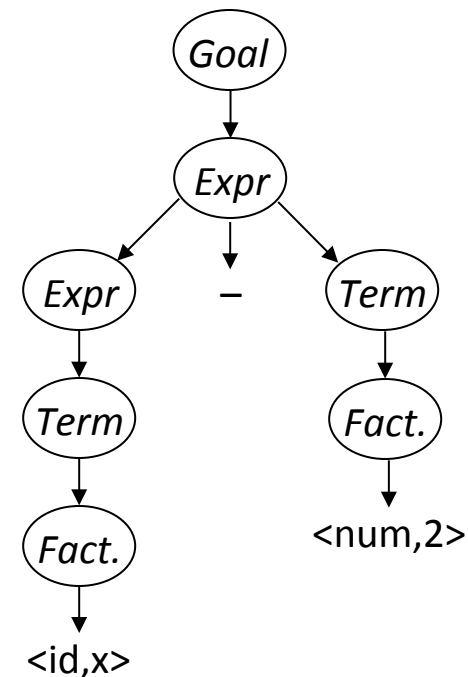
Example

Trying to match the “2” in $\underline{x} - \underline{2} * \underline{y}$:

Rule	Sentential Form	Input
\rightarrow	$\langle \text{id}, \underline{x} \rangle - \text{Term}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
6	$\langle \text{id}, \underline{x} \rangle - \text{Factor}$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
8	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle$	$\underline{x} - \uparrow \underline{2} * \underline{y}$
\rightarrow	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle$	$\underline{x} - \underline{2} \uparrow * \underline{y}$

Where are we?

- “2” matches “2”
 - We have more input, but no **NTs** left to expand
 - The expansion terminated too soon
- \Rightarrow Need to backtrack

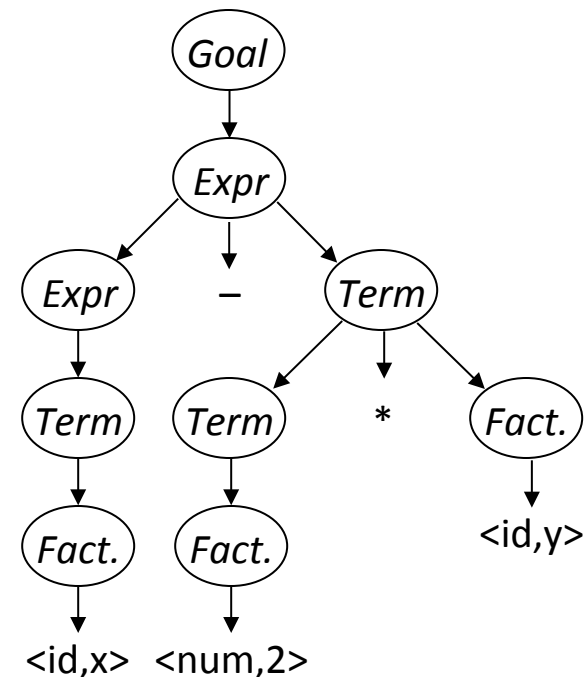


Example



Trying again with “2” in $x - 2 * y$:

Rule	Sentential Form	Input
→	$\langle id, x \rangle - Term$	$x - \uparrow 2 * y$
4	$\langle id, x \rangle - Term * Factor$	$x - \uparrow 2 * y$
6	$\langle id, x \rangle - Factor * Factor$	$x - \uparrow 2 * y$
8	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$x - \uparrow 2 * y$
→	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$x - 2 \uparrow * y$
→	$\langle id, x \rangle - \langle num, 2 \rangle * Factor$	$x - 2 * \uparrow y$
9	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$x - 2 * \uparrow y$
→	$\langle id, x \rangle - \langle num, 2 \rangle * \langle id, y \rangle$	$x - 2 * y \uparrow$



The Point:

For efficiency, the parser must make the correct choice when it expands a NT. Wrong choices lead to wasted effort.

Another possible parse



Other choices for expansion are possible

Rule	Sentential Form	Input
—	Goal	$\uparrow \underline{x} - \underline{2} * y$
0	Expr	$\uparrow \underline{x} - \underline{2} * y$
1	Expr + Term	$\uparrow \underline{x} - \underline{2} * y$
1	Expr + Term + Term	$\uparrow \underline{x} - \underline{2} * y$
1	Expr + Term + Term + Term	$\uparrow \underline{x} - \underline{2} * y$
1	... and so on	$\uparrow \underline{x} - \underline{2} * y$

Consumes no input!

This expansion doesn't terminate

- Wrong choice of expansion leads to non-termination
- **Non-termination** is a bad property for a parser to have
- Parser must make the right choice

The Classic Expression Grammar



0	<i>Goal</i>	\rightarrow	<i>Expr</i>
1	<i>Expr</i>	\rightarrow	<i>Expr</i> + <i>Term</i>
2			<i>Expr</i> - <i>Term</i>
3			<i>Term</i>
4	<i>Term</i>	\rightarrow	<i>Term</i> * <i>Factor</i>
5			<i>Term</i> / <i>Factor</i>
6			<i>Factor</i>
7	<i>Factor</i>	\rightarrow	(<i>Expr</i>)
8			<u>number</u>
9			<u>id</u>

Classic Expression Grammar

The possibility of an infinite sequence of expansions in a parser is ~~bad~~. **disastrous**

- The problem arises from *left recursion* in the grammar and a *leftmost* derivation¹
- **LHS** symbol cannot appear at start of the **RHS**
 - *Cannot derive from it in multiple steps, either*
- Top down parsers build leftmost derivations, so grammars with left recursion are not suitable for top-down parsing

¹ Similar problem arises with *right recursion* and a *rightmost* derivation.

Left Recursion



Top-down parsers cannot handle left-recursive grammars

Formally,

A grammar is **left recursive** if $\exists A \in NT$ such that a derivation $A \Rightarrow^+ A\alpha$ exists, for some string $\alpha \in (NT \cup T)^+$

Our classic expression grammar is left recursive

- This can lead to non-termination in a top-down parser
- In a top-down parser, any recursion must be right recursion
- We would like to convert the left recursion to right recursion

Non-termination is always a bad property in a compiler

Fortunately, we can eliminate left recursion in an algorithmic way.

Right recursion is defined in a symmetric way.

Eliminating Left Recursion



To remove left recursion, we can transform the grammar

Consider a grammar fragment of the form

$$\begin{aligned} Fee &\rightarrow Fee \ \alpha \\ &| \ \beta \end{aligned}$$

where neither α nor β start with Fee

Language is β followed
by 0 or more α 's

We can rewrite this fragment as

$$\begin{aligned} Fee &\rightarrow \beta Fie \\ Fie &\rightarrow \alpha Fie \\ &| \ \varepsilon \end{aligned}$$

where Fie is a new non-terminal

The new grammar defines the
same language as the old
grammar, using only right
recursion.

Added a reference to
the empty string

New Idea: the ε production

Eliminating Left Recursion



The expression grammar contains two cases of left recursion

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} + \text{Term} \\ &| \text{Expr} - \text{Term} \\ &| \text{Term} \end{aligned}$$
$$\begin{aligned} \text{Term} &\rightarrow \text{Term} * \text{Factor} \\ &| \text{Term} / \text{Factor} \\ &| \text{Factor} \end{aligned}$$

Applying the transformation yields

$$\begin{aligned} \text{Expr} &\rightarrow \text{Term Expr}' \\ \text{Expr}' &\rightarrow + \text{Term Expr}' \\ &| - \text{Term Expr}' \\ &| \epsilon \end{aligned}$$
$$\begin{aligned} \text{Term} &\rightarrow \text{Factor Term}' \\ \text{Term}' &\rightarrow * \text{Factor Term}' \\ &| / \text{Factor Term}' \\ &| \epsilon \end{aligned}$$

These fragments use only right recursion

Eliminating Left Recursion



Substituting them back into the grammar yields

0	<i>Goal</i>	\rightarrow	<i>Expr</i>
1	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>
2	<i>Expr'</i>	\rightarrow	$+ \textit{Term Expr'}$
3		$ $	$- \textit{Term Expr'}$
4		$ $	ϵ
5	<i>Term</i>	\rightarrow	<i>Factor Term'</i>
6	<i>Term'</i>	\rightarrow	$* \textit{Factor Term'}$
7		$ $	$/ \textit{Factor Term'}$
8		$ $	ϵ
9	<i>Factor</i>	\rightarrow	(\textit{Expr})
10		$ $	<u>number</u>
11		$ $	<u>id</u>

Right-recursive expression grammar

- This grammar is correct, if somewhat counter-intuitive.
 - A top-down parser will terminate using it.
 - A top-down parser may need to backtrack with it.
 - It is left associative, as was the original
- \Rightarrow Why didn't we just rewrite it so *Expr* was at the right end of the **RHS**, rather than the left end?

Eliminating Left Recursion



Substituting them back into the grammar yields

0	<i>Goal</i>	\rightarrow	<i>Expr</i>
1	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>
2	<i>Expr'</i>	\rightarrow	$+ \textit{Term Expr'}$
3		$ $	$- \textit{Term Expr'}$
4		$ $	ϵ
5	<i>Term</i>	\rightarrow	<i>Factor Term'</i>
6	<i>Term'</i>	\rightarrow	$* \textit{Factor Term'}$
7		$ $	$/ \textit{Factor Term'}$
8		$ $	ϵ
9	<i>Factor</i>	\rightarrow	(\textit{Expr})
10		$ $	<u>number</u>
11		$ $	<u>id</u>

Right-recursive expression grammar

- This grammar is correct, if

NOTE: This technique eliminates direct left recursion — when a production's RHS begins with its own LHS.

It does not address indirect left recursion. We will get there, in a couple of slides ...

original

\Rightarrow Why didn't we just rewrite it so *Expr* was at the right end of the RHS, rather than the left end?

Eliminating Left Recursion



Notice that we do not use the naïve right-recursive form

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

Transformed grammar fragment

Naïve right-recursive form

The naïve right-recursive form generates a different associativity (and parse tree) than did the original grammar.

The transformed grammar fragment generates the same parse tree as the original grammar did. (See § 3.5.3 in EaC2e.)

Eliminating Left Recursion

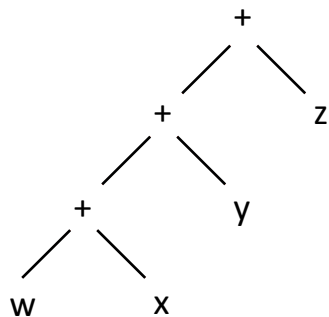


The naïve right-recursive form changes the associativity

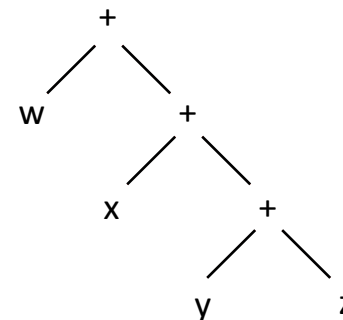
$$\begin{aligned} \text{Expr} &\rightarrow \text{Term Expr}' \\ \text{Expr}' &\rightarrow + \text{Term Expr}' \\ &\quad | - \text{Term Expr}' \\ &\quad | \epsilon \end{aligned}$$

$$\begin{aligned} \text{Expr} &\rightarrow \text{Term} + \text{Expr} \\ &\quad | \text{Term} - \text{Expr} \\ &\quad | \text{Term} \end{aligned}$$

Transformed grammar fragment



Naïve right-recursive form



ASTs for $w + x + y + z$

Parsing with the RR CEG



0	<i>Goal</i>	\rightarrow	<i>Expr</i>
1	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>
2	<i>Expr'</i>	\rightarrow	$+ \textit{Term Expr'}$
3		$ $	$- \textit{Term Expr'}$
4		$ $	ϵ
5	<i>Term</i>	\rightarrow	<i>Factor Term'</i>
6	<i>Term'</i>	\rightarrow	$* \textit{Factor Term'}$
7		$ $	$/ \textit{Factor Term'}$
8		$ $	ϵ
9	<i>Factor</i>	\rightarrow	(\textit{Expr})
10		$ $	<u>number</u>
11		$ $	<u>id</u>

Right Recursive CEG

<i>Rule</i>	<i>Sentential Form</i>	<i>x – 2 * y, again</i>
—	<i>Goal</i>	
0	<i>Expr</i>	
1	<i>Term Expr'</i>	
5	<i>Factor Term' Expr'</i>	
11	$\langle \text{id}, \underline{x} \rangle \textit{Term' Expr'}$	
8	$\langle \text{id}, \underline{x} \rangle \textit{Expr'}$	
3	$\langle \text{id}, \underline{x} \rangle - \textit{Term Expr'}$	
5	$\langle \text{id}, \underline{x} \rangle - \textit{Factor Term' Expr'}$	
10	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle \textit{Term' Expr'}$	
6	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \textit{Factor Term' Expr'}$	
11	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, y \rangle \textit{Term' Expr'}$	
8	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, y \rangle \textit{Expr'}$	
4	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, y \rangle$	

Parsing with RR CEG

Rule	Sentential Form
—	<i>Goal</i>
0	<i>Expr</i>
1	<i>Term Expr'</i>
5	<i>Factor Term' Expr'</i>
11	<i><id,x> Term' Expr'</i>
8	<i><id,x> Expr'</i>
3	<i><id,x> - Term Expr'</i>
5	<i><id,x> - Factor Term' Expr'</i>
10	<i><id,x> - <num,2> Term' Expr'</i>
6	<i><id,x> - <num,2> * Factor Term' Expr'</i>
11	<i><id,x> - <num,2> * <id,y> Term' Expr'</i>
8	<i><id,x> - <num,2> * <id,y> Expr'</i>
4	<i><id,x> - <num,2> * <id,y></i>

