

# Introduction to Parsing





Parser

- Checks the stream of <u>words</u> and their <u>parts of speech</u> (produced by the scanner) for grammatical correctness
- Builds an IR representation of the code

Think of this as the mathematics of diagramming sentences



The process of discovering a *derivation* for some sentence

- Need a mathematical model of syntax a grammar G
- Need an algorithm for testing membership in *L(G)*

Roadmap

- 1 Context-free grammars and derivations
- 2 Top-down parsing
  - → Hand-coded recursive descent parsers
- 3 Bottom-up parsing
  - $\rightarrow$  Generated LR(1) parsers

Specifying Syntax with a Grammar



Context-free syntax is specified with a context-free grammar (CFG)

# SheepNoise → SheepNoise baa | baa

This CFG defines the set of noises sheep normally make



It is written in a variant of Backus-Naur form, BNF notation

(words)

Formally, a grammar is a four tuple, G = (S, NT, T, P)

- *S* is the *start* (*or goal*) *symbol*
- NT is a set of non-terminal symbols (syntactic variables)
- T is a set of *terminal symbols*
- P is a set of productions or rewrite rules
  Production rules follow format NT→(NT∪T)<sup>+</sup>

Specifying Syntax with a Grammar



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

What are the:

S:

NT:

T:

**P**:

Specifying Syntax with a Grammar



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

What are the:

- S: SheepNoise
- NT: SheepNoise
  - T: baa
  - P: SheepNoise  $\rightarrow$  SheepNoise baa

SheepNoise  $\rightarrow$  baa



We can use the *SheepNoise* grammar to create sentences

 $\rightarrow$  use the productions as *rewriting rules* 

Rule	Sentential Form
-	SheepNoise
2	baa

Rule	Sentential Form
-	SheepNoise
1	SheepNoise <u>baa</u>
2	baa baa

And so on ...



To explore the uses of CFGs, we need a more complex grammar

1	Expr	$  \rightarrow$	Expr Op Expr
2			number
3			<u>id</u>
4	Ор	$\rightarrow$	+
5			-
6			*
7			1

What are the NT and T?



Rule	Sentential Form
_	Expr
1	Expr Op Expr
3	<id,<mark>x&gt; Op Expr</id,<mark>
5	<id,<u>x&gt; - Expr</id,<u>
1	<id,<u>x&gt; - Expr Op Expr</id,<u>
2	<id,<u>x&gt; - <num,<u>2&gt; Op Expr</num,<u></id,<u>
6	<id,<u>x&gt; - <num,<u>2&gt; * <i>Expr</i></num,<u></id,<u>
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>

- This <u>sequence of rewrites</u> is called a *derivation*
- Process of discovering a derivation is called *parsing*

We denote this derivation:  $Expr \Rightarrow^* \underline{id} - \underline{num} * \underline{id}$ 



- At each step, we choose a non-terminal to replace
- Different choices can lead to different derivations

Two derivations are of interest

- *Leftmost derivation* replace leftmost NT at each step
- *Rightmost derivation* replace rightmost NT at each step

The example on the preceding slide was a *leftmost* derivation



In both cases,  $Expr \Rightarrow * \underline{id} - \underline{num} * \underline{id}$ 

- The two derivations produce different parse trees
- The parse trees imply different evaluation orders!

Rule	Sentential Form	Rul
	Expr	
1	Expr Op Expr	1
3	<id x=""> Op Expr</id>	3
5	<id,<mark>x&gt; - Expr</id,<mark>	6
1	<id,<mark>x&gt; - Expr Op Expr</id,<mark>	1
2	<id,x> - <num,2> Op Expr</num,2></id,x>	2
6	<id,<u>x&gt; - <num,2> * <i>Expr</i></num,2></id,<u>	5
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>	3

Rule	Sentential Form
—	Expr
1	Expr Op Expr
3	<i>Expr Op</i> <id,<mark>y&gt;</id,<mark>
6	Expr * <id,<u>y&gt;</id,<u>
1	<i>Expr Op Expr</i> * <id,<u>y&gt;</id,<u>
2	<i>Expr Op&lt;</i> num, <mark>2</mark> > * <id,<u>y&gt;</id,<u>
5	<i>Expr-</i> <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u>
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>

Leftmost derivation

Rightmost derivation

### Derivations and Parse Trees



Leftmost derivation

Rule	Sentential Form
	Expr
1	Expr Op Expr
3	<id,<mark>x&gt; <i>Op Expr</i></id,<mark>
5	<id,<u>x&gt; - Expr</id,<u>
1	<id,<u>x&gt; - Expr Op Expr</id,<u>
2	<id,<u>x&gt; - <num,<u>2&gt; <i>Op Expr</i></num,<u></id,<u>
6	<id,<u>x&gt; - <num,<u>2&gt; * <i>Expr</i></num,<u></id,<u>
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>

Let's do the parse tree on the board

This evaluates as  $\underline{x} - (\underline{2} * \underline{y})$ 

### Derivations and Parse Trees

#### Leftmost derivation

Rule	Sentential Form
	Expr
1	Expr Op Expr
3	<id,<mark>x&gt; <i>Op Expr</i></id,<mark>
5	<id,<u>x&gt; - Expr</id,<u>
1	<id,<mark>x&gt; - Expr Op Expr</id,<mark>
2	<id,<u>x&gt; - <num,<u>2&gt; <i>Op Expr</i></num,<u></id,<u>
6	<id,<u>x&gt; - <num,<u>2&gt; * <i>Expr</i></num,<u></id,<u>
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>

This evaluates as  $\underline{x} - (\underline{2} * \underline{y})$ 







**Rightmost derivation** 

Rule	Sentential Form
	Expr
1	Expr Op Expr
3	<i>Expr Op</i> <id,<mark>y&gt;</id,<mark>
6	Expr * <id,<u>y&gt;</id,<u>
1	<i>Expr Op Expr</i> * <id,<u>y&gt;</id,<u>
2	<i>Expr Op</i> <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u>
5	<i>Expr-</i> <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u>
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>

Let's do the parse tree on the board

This evaluates as  $(\underline{x} - \underline{2}) * \underline{y}$ 

### Derivations and Parse Trees

#### Rightmost derivation

Rule	Sentential Form
_	Expr
1	Expr Op Expr
3	<i>Expr Op</i> <id,<mark>y&gt;</id,<mark>
6	Expr * <id,<u>y&gt;</id,<u>
1	<i>Expr Op Expr</i> * <id,<u>y&gt;</id,<u>
2	<i>Expr Op</i> <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u>
5	<i>Expr-</i> <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u>
3	<id,<u>x&gt; - <num,<u>2&gt; * <id,<u>y&gt;</id,<u></num,<u></id,<u>

This evaluates as  $(\underline{x} - \underline{2}) * \underline{y}$ 







These two derivations point out a problem with the grammar: It has no notion of precedence, or implied order of evaluation

To add precedence

- Create a non-terminal for each *level of precedence*
- Isolate the corresponding part of the grammar
- Force the parser to recognize high precedence subexpressions first

For algebraic expressions

- Multiplication and division, first
- Subtraction and addition, next

(level one) (level two)