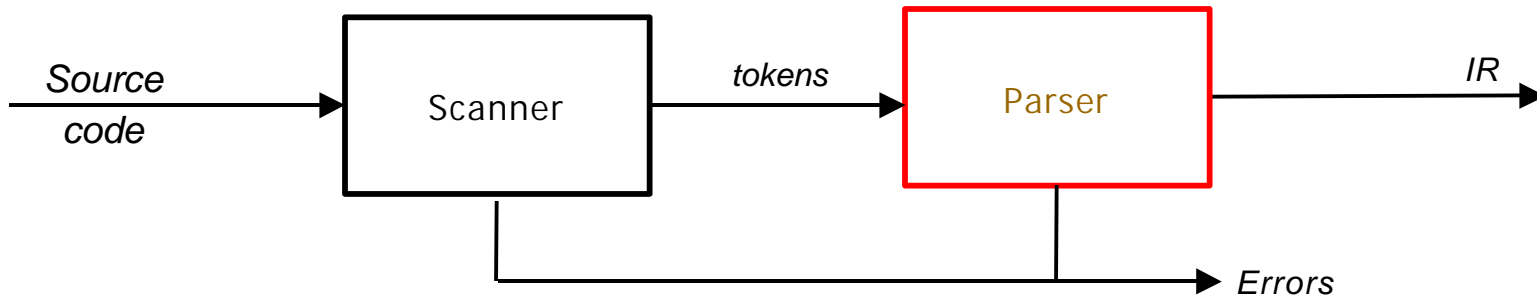




Introduction to Parsing

The Front End



Parser

- Checks the stream of words and their parts of speech (produced by the scanner) for grammatical correctness
- Determines if the input is syntactically well formed
- Guides checking at deeper levels than syntax
- May build an IR representation of the code

Think of this as the mathematics of diagramming sentences



The Study of Parsing

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 $L(G)$
- Need to keep in mind that our goal is building parsers, not studying the mathematics of arbitrary languages

Roadmap

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



Specifying Syntax with a Grammar

Context-free syntax is specified with a context-free grammar

$$\begin{aligned} \textit{SheepNoise} &\rightarrow \textit{SheepNoise} \underline{\textit{baa}} \\ &| \underline{\textit{baa}} \end{aligned}$$

This *CFG* defines the set of noises sheep normally make

It is written in a variant of Backus–Naur form

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

- S is the *start symbol* (*set of strings in $L(G)$*)
- N is a set of *non-terminal symbols* (*syntactic variables*)
- T is a set of *terminal symbols* (*words*)
- P is a set of *productions* or *rewrite rules* ($P : N \rightarrow (N \cup T)^+$)

Example due to Dr. Scott K. Warren



Deriving Syntax

We can use the *SheepNoise* grammar to create sentences

- > use the productions as *rewriting rules*

<i>Rule</i>	<i>Sentential Form</i>
—	<i>SheepNoise</i>
2	<u>baa</u>

<i>Rule</i>	<i>Sentential Form</i>
—	<i>SheepNoise</i>
1	<i>SheepNoise</i> <u>baa</u>
1	<i>SheepNoise</i> <u>baa</u> <u>baa</u>
2	<u>baa</u> <u>baa</u> <u>baa</u>

<i>Rule</i>	<i>Sentential Form</i>
—	<i>SheepNoise</i>
1	<i>SheepNoise</i> <u>baa</u>
2	<u>baa</u> <u>baa</u>

And so on ...

This example quickly runs out of intellectual steam ...



A More Useful Grammar

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

1	<i>Goal</i>	\rightarrow	<i>Expr</i>
2	<i>Expr</i>	\rightarrow	<i>Expr Op Expr</i>
3			<u>number</u>
4			<u>id</u>
5	<i>Op</i>	\rightarrow	+
6			-
7			*
8			/

<i>Rule</i>	<i>Sentential Form</i>
—	<i>Expr</i>
2	<i>Expr Op Expr</i>
4	<id, <u>x</u> > <i>Op Expr</i>
6	<id, <u>x</u> > - <i>Expr</i>
2	<id, <u>x</u> > - <i>Expr Op Expr</i>
3	<id, <u>x</u> > - <num, <u>2</u> > <i>Op Expr</i>
7	<id, <u>x</u> > - <num, <u>2</u> > * <i>Expr</i>
4	<id, <u>x</u> > - <num, <u>2</u> > * <id, <u>y</u> >

We denote this: $Expr \mathbf{P}^* \text{id} - \text{num}^* \text{id}$

- Such a sequence of rewrites is called a *derivation*
- Process of discovering a derivation is called *parsing*



Derivations

- 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

These are the two *systematic* derivations

(We don't care about randomly-ordered derivations!)

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

- Of course, there is a *rightmost* derivation
- Interestingly, it turns out to be different



The Two Derivations for $\underline{x} - \underline{2} * \underline{y}$

Rule	Sentential Form
—	<i>Expr</i>
2	<i>Expr Op Expr</i>
2	<i>Expr Op Expr Op Expr</i>
4	$\langle \text{id}, \underline{x} \rangle \text{ Op } \text{Expr Op } \text{Expr}$
6	$\langle \text{id}, \underline{x} \rangle - \text{Expr Op } \text{Expr}$
3	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle \text{ Op } \text{Expr}$
7	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \text{Expr}$
4	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, \underline{y} \rangle$

Leftmost derivation

Rule	Sentential Form
—	<i>Expr</i>
2	<i>Expr Op Expr</i>
4	<i>Expr Op</i> $\langle \text{id}, \underline{y} \rangle$
7	<i>Expr</i> * $\langle \text{id}, \underline{y} \rangle$
2	<i>Expr Op Expr</i> * $\langle \text{id}, \underline{y} \rangle$
3	<i>Expr Op</i> $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
6	<i>Expr</i> - $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
4	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, \underline{y} \rangle$

Rightmost derivation

In both cases, $\text{Expr} \stackrel{P}{\Rightarrow} \text{id} - \text{num} * \text{id}$

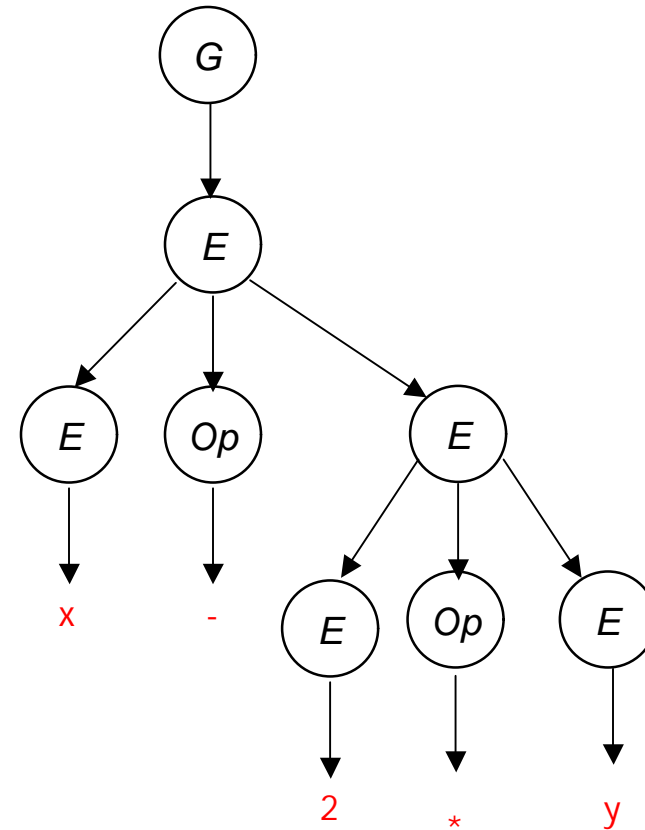
- The two derivations produce different parse trees
- The parse trees imply different evaluation orders!



Derivations and Parse Trees

Leftmost derivation

Rule	Sentential Form
—	<i>Expr</i>
2	<i>Expr Op Expr</i>
4	$\langle \text{id}, \underline{x} \rangle \text{ Op } \textit{Expr}$
6	$\langle \text{id}, \underline{x} \rangle - \textit{Expr}$
2	$\langle \text{id}, \underline{x} \rangle - \textit{Expr Op Expr}$
3	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle \text{ Op } \textit{Expr}$
7	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \textit{Expr}$
4	$\langle \text{id}, \underline{x} \rangle - \langle \text{num}, \underline{2} \rangle * \langle \text{id}, \underline{y} \rangle$



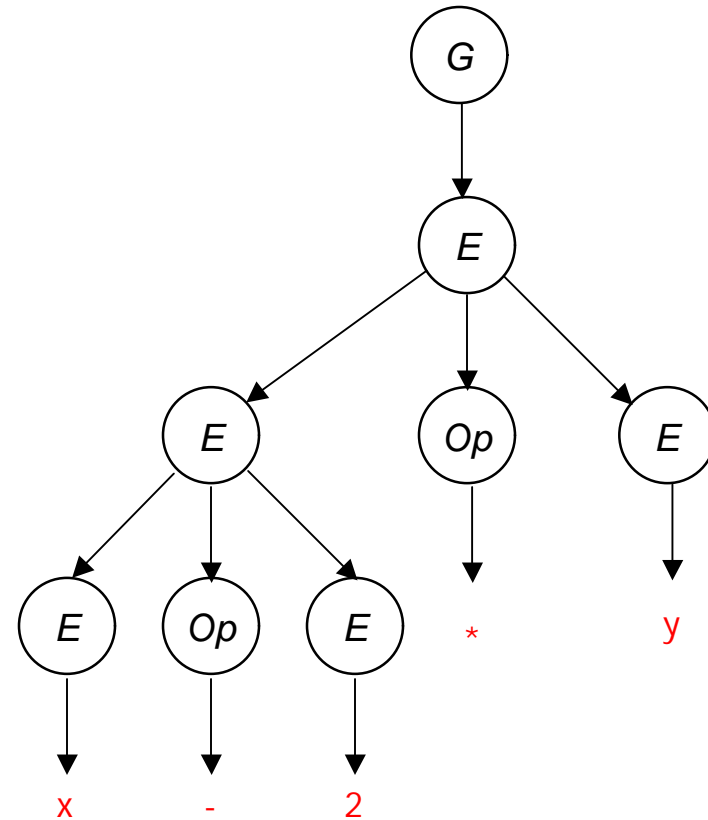
This evaluates as $\underline{x} - (\underline{2} * \underline{y})$



Derivations and Parse Trees

Rightmost derivation

Rule	Sentential Form
—	<i>Expr</i>
2	<i>Expr Op Expr</i>
4	<i>Expr Op</i> <id, <u>y</u> >
7	<i>Expr</i> * <id, <u>y</u> >
2	<i>Expr Op Expr</i> * <id, <u>y</u> >
3	<i>Expr Op</i> <num, <u>2</u> > * <id, <u>y</u> >
6	<i>Expr</i> - <num, <u>2</u> > * <id, <u>y</u> >
4	<id, <u>x</u> > - <num, <u>2</u> > * <id, <u>y</u> >



This evaluates as $(x - 2) * y$



Derivations and Precedence

*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 parser to recognize high precedence subexpressions first

For algebraic expressions

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



Derivations and Precedence

Adding the standard algebraic precedence produces:

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

This grammar is slightly larger

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

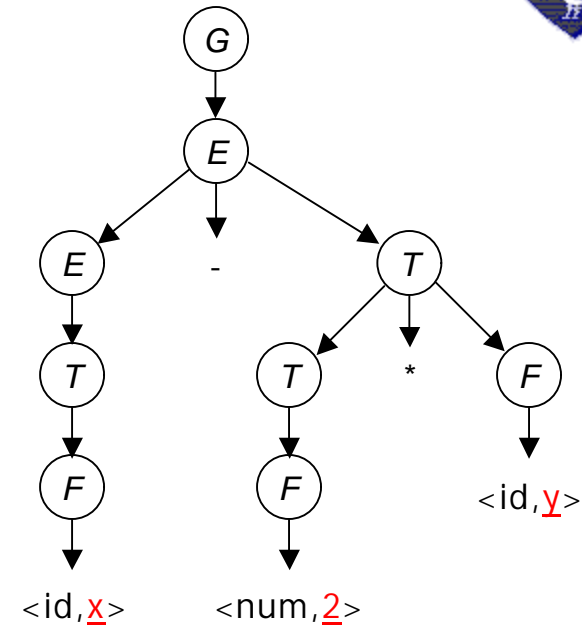
Let's see how it parses our example



Derivations and Precedence

Rule	Sentential Form
—	<i>Expr</i>
3	<i>Expr</i> - <i>Term</i>
5	<i>Expr</i> - <i>Term</i> * <i>Factor</i>
9	<i>Expr</i> - <i>Term</i> * $\langle \text{id}, \underline{y} \rangle$
7	<i>Expr</i> - <i>Factor</i> * $\langle \text{id}, \underline{y} \rangle$
8	<i>Expr</i> - $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
4	<i>Term</i> - $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
7	<i>Factor</i> - $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$
9	$\langle \text{id}, \underline{x} \rangle$ - $\langle \text{num}, \underline{2} \rangle$ * $\langle \text{id}, \underline{y} \rangle$

The rightmost derivation



Its parse tree

This produces $\underline{x} - (\underline{2} * \underline{y})$, along with an appropriate parse tree.

Both the leftmost and rightmost derivations give the same expression, because the grammar directly encodes the desired precedence.



Ambiguous Grammars

- Our original expression grammar had other problems

1	<i>Goal</i>	\rightarrow	<i>Expr</i>
2	<i>Expr</i>	\rightarrow	<i>Expr Op Expr</i>
3			<u>number</u>
4			<u>id</u>
5	<i>Op</i>	\rightarrow	+
6			-
7			*
8			/

<i>Rule</i>	<i>Sentential Form</i>
—	<i>Expr</i>
2	<i>Expr Op Expr</i>
2	<i>Expr Op Expr Op Expr</i>
4	$\langle \text{id, } \underline{x} \rangle \text{ Op Expr Op Expr}$
6	$\langle \text{id, } \underline{x} \rangle - \text{Expr Op Expr}$
3	$\langle \text{id, } \underline{x} \rangle - \langle \text{num, } \underline{2} \rangle \text{ Op Expr}$
7	$\langle \text{id, } \underline{x} \rangle - \langle \text{num, } \underline{2} \rangle * \text{Expr}$
4	$\langle \text{id, } \underline{x} \rangle - \langle \text{num, } \underline{2} \rangle * \langle \text{id, } \underline{y} \rangle$

- This grammar allows multiple leftmost derivations for $\underline{x} - \underline{2} * \underline{y}$
- Hard to automate derivation if > 1 choice
- The grammar is *ambiguous*

different choice than
the first time



Ambiguous Grammars

Definitions

- If a grammar has more than one leftmost derivation for a single *sentential form*, the grammar is ambiguous
- If a grammar has more than one rightmost derivation for a single sentential form, the grammar is ambiguous
- The leftmost and rightmost derivations for a sentential form may differ, even in an unambiguous grammar

Classic example — the *if-then-else* problem

```
Stmt @ if Expr then Stmt  
      | if Expr then Stmt else Stmt  
      | ... other stmts ...
```

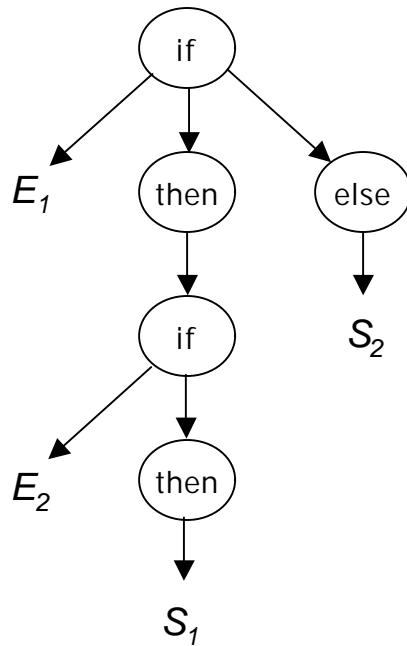
This ambiguity is entirely grammatical in nature



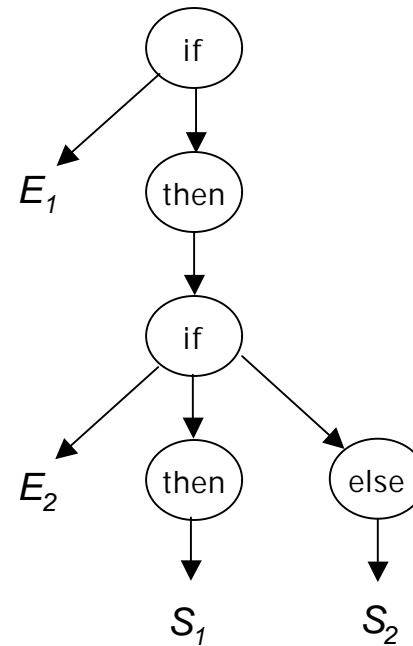
Ambiguity

This sentential form has two derivations

if Expr₁ then if Expr₂ then Stmt₁ else Stmt₂



*production 2, then
production 1*



*production 1, then
production 2*