

CS 314 Principles of Programming Languages

Lecture 5: Syntax Analysis (Parsing)

Zheng (Eddy) Zhang



Rutgers University

January 31, 2018

Class Information

- Homework 1 is being graded now.
The sample solution will be posted soon.
- Homework 2 posted. Due in one week.

Review: Context Free Grammars (CFGs)

- A formalism to for describing languages
- A CFG $G = \langle T, N, P, S \rangle$:
 1. A set T of terminal symbols (tokens).
 2. A set N of nonterminal symbols.
 3. A set P production (rewrite) rules.
 4. A special start symbol S .
- The language $L(G)$ is the set of sentences of terminal symbols in T^* that can be derived from the start symbol S :

$$L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$$

Review: Context Free Grammar

...

$\langle \text{if-stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle$

$\langle \text{expr} \rangle ::= \text{id } \langle = \rangle \text{ id}$

$\langle \text{stmt} \rangle ::= \text{id } := \text{num}$

...

Context free grammar

Rule 1 $\$1 \Rightarrow 1\&$

Rule 2 $\$0 \Rightarrow 0\$$

Rule 3 $\&1 \Rightarrow 1\$$

Rule 4 $\&0 \Rightarrow 0\&$

Rule 5 $\$\# \Rightarrow \rightarrow A$

Rule 6 $\&\# \Rightarrow \rightarrow B$

Not a context free grammar

CFGs are rewrite systems with restrictions on the form of rewrite (production) rules that can be used. The left hand side of a production rule can only be **one non-terminal symbol**.

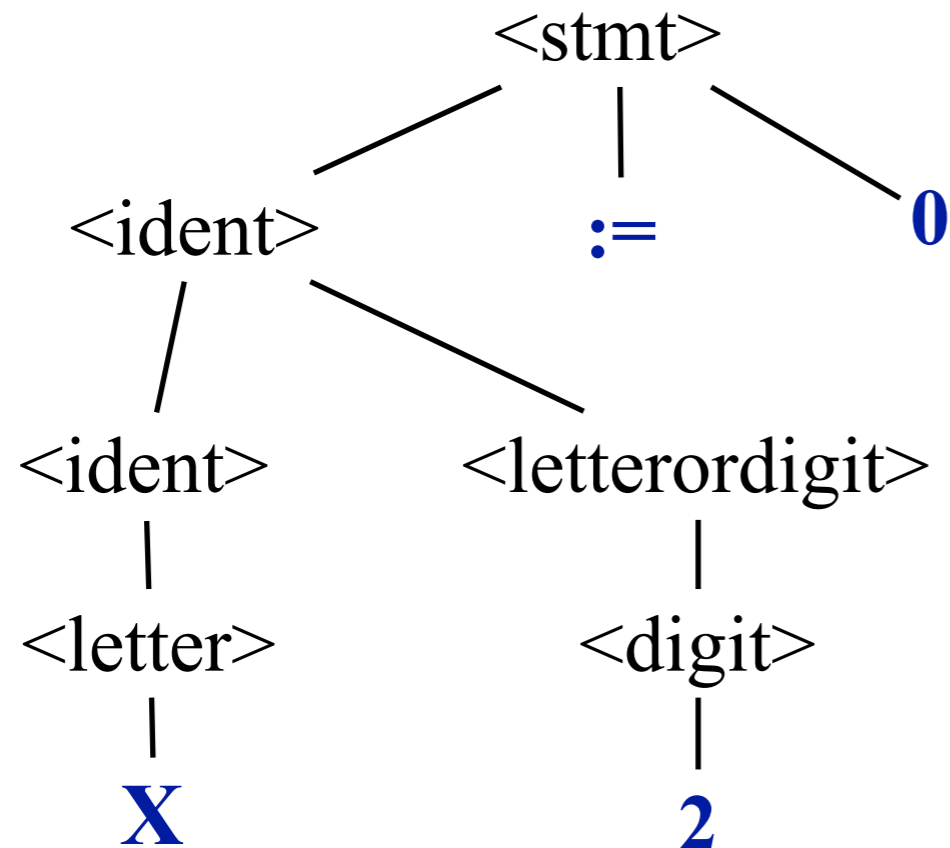
A Language May Have Many Grammars

Consider G' :

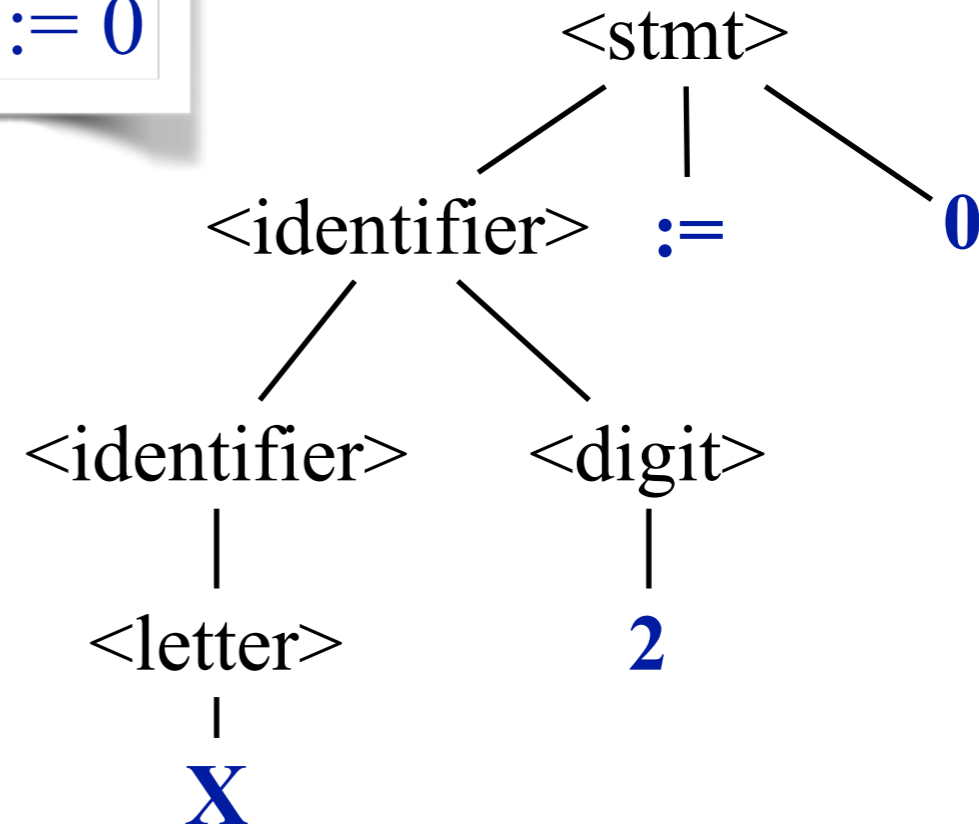
1. $\langle \text{letter} \rangle ::= \mathbf{A} \mid \mathbf{B} \mid \mathbf{C} \mid \dots \mid \mathbf{Z}$
2. $\langle \text{digit} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \dots \mid \mathbf{9}$
3. $\langle \text{ident} \rangle ::= \langle \text{letter} \rangle \mid$
4. $\langle \text{ident} \rangle \langle \text{letterordigit} \rangle$
5. $\langle \text{stmt} \rangle ::= \langle \text{ident} \rangle := \mathbf{0}$
6. $\langle \text{letterordigit} \rangle ::= \langle \text{letter} \rangle \mid \langle \text{digit} \rangle$

The Original Grammar G :

1. $\langle \text{letter} \rangle ::= \mathbf{A} \mid \mathbf{B} \mid \mathbf{C} \mid \dots \mid \mathbf{Z}$
2. $\langle \text{digit} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \dots \mid \mathbf{9}$
3. $\langle \text{identifier} \rangle ::= \langle \text{letter} \rangle \mid$
4. $\langle \text{identifier} \rangle \langle \text{letter} \rangle \mid$
5. $\langle \text{identifier} \rangle \langle \text{digit} \rangle$
6. $\langle \text{stmt} \rangle ::= \langle \text{identifier} \rangle := \mathbf{0}$



$X2 := 0$



Review: Grammars and Programming Languages

Many grammars may correspond to one programming language.

Good grammars:

- Captures the logic structure of the language
⇒ structure carries some semantic information
(example: expression grammar)
- Use meaningful names
- Are easy to read
- Are unambiguous
- ...

Review: Ambiguous Grammars

“Time flies like an arrow; fruit flies like a banana.”

A grammar \mathcal{G} is ambiguous iff there exists a $w \in L(\mathcal{G})$ such that there are:

- two distinct parse trees for w , or
- two distinct leftmost derivations for w , or
- two distinct rightmost derivations for w .

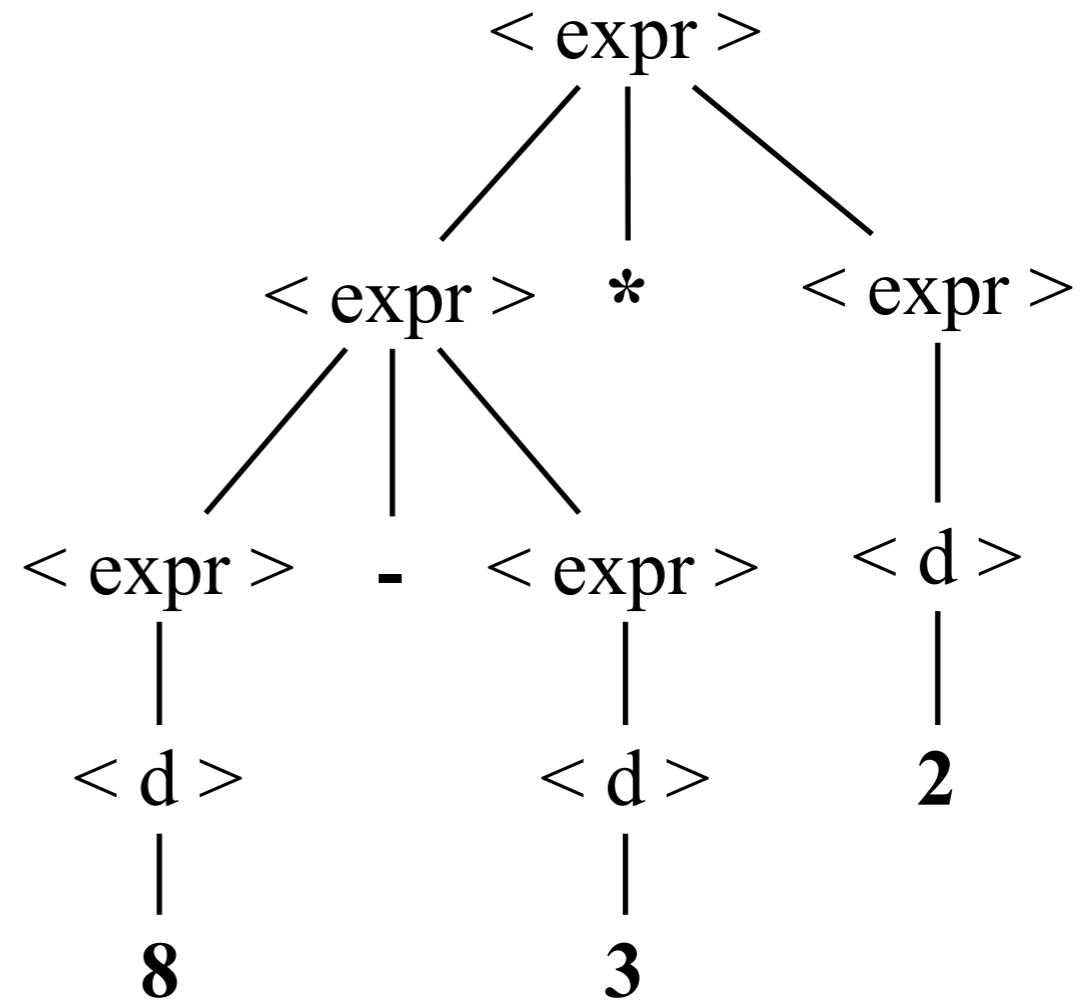
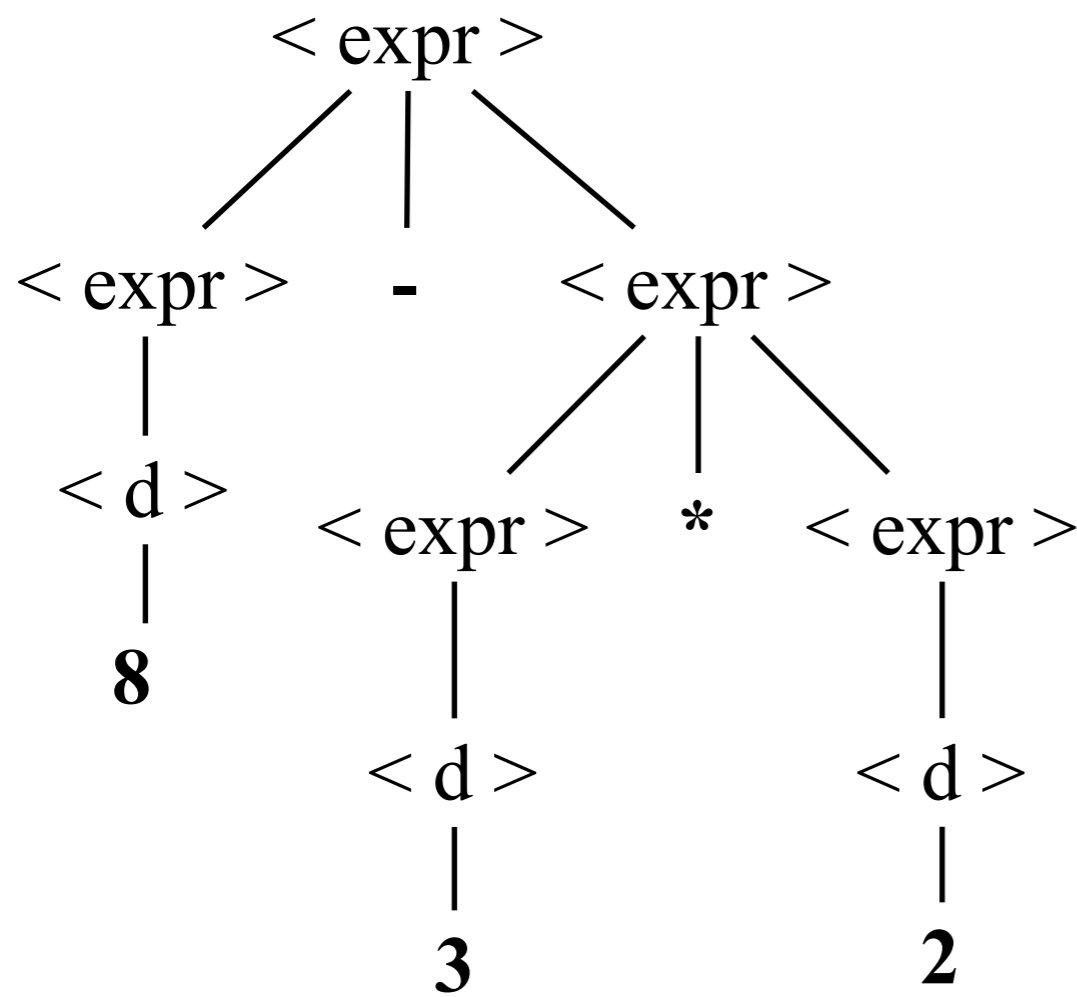
We want a unique semantics of our programs, which typically requires a unique syntactic structure.

Review: Arithmetic Expression Grammar

Parse “8 - 3 * 2”:

$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$
 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid$
 $\quad \langle \text{expr} \rangle * \langle \text{expr} \rangle \mid$
 $\quad \langle d \rangle \mid \langle 1 \rangle$
 $\langle d \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$
 $\langle 1 \rangle ::= a \mid b \mid c \mid \dots \mid z$

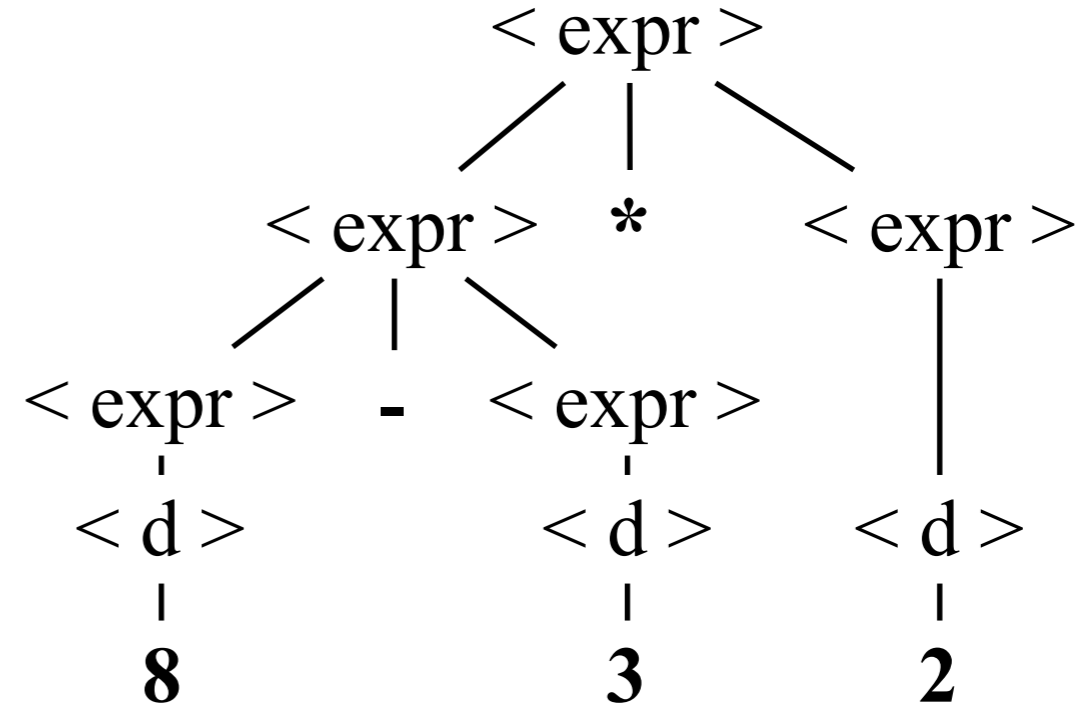
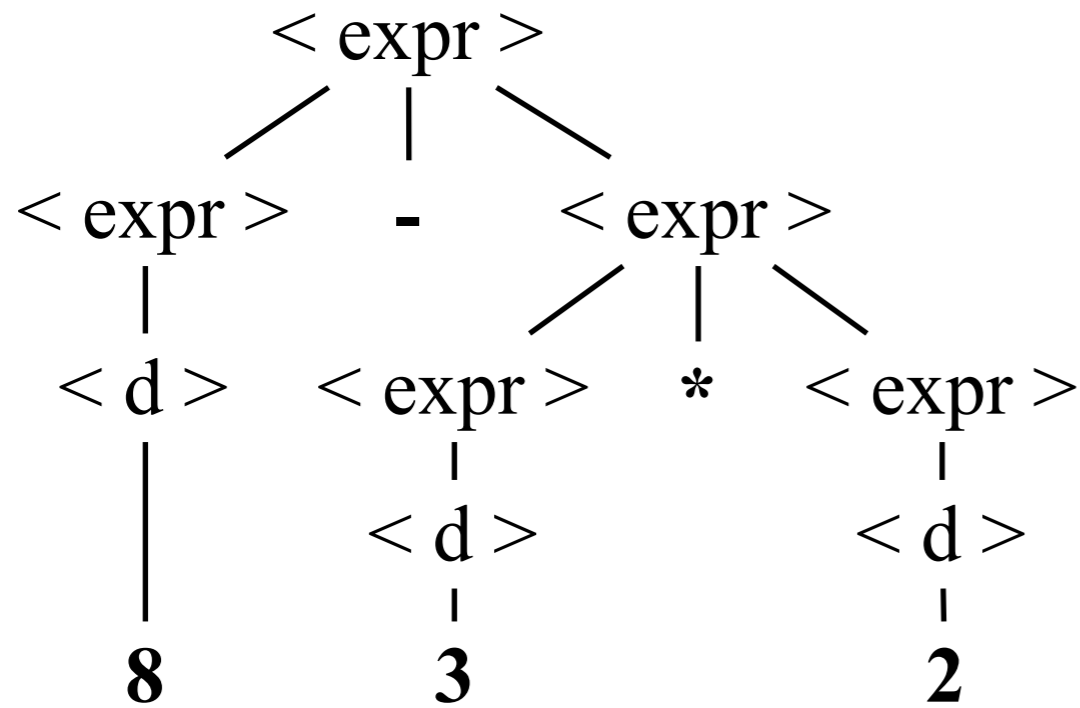
Two parse trees!



Review: Arithmetic Expression Grammar

Parse “8 - 3 * 2”:

Two leftmost derivations!



leftmost derivation	
$\langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{expr} \rangle - \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{expr} \rangle * \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{d} \rangle * \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{d} \rangle * \langle \text{d} \rangle$	

leftmost derivation	
$\langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{expr} \rangle * \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{expr} \rangle - \langle \text{expr} \rangle * \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{expr} \rangle * \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{d} \rangle * \langle \text{expr} \rangle$	\Rightarrow_L
$\langle \text{d} \rangle - \langle \text{d} \rangle * \langle \text{d} \rangle$	

Review: Ambiguity

How to deal with ambiguity?

- Change the language
Example: Adding new terminal symbols as delimiters.
Fix the *dangling else, expression* grammars.
- Change the grammar
Example: Impose **associativity** and **precedence** in an *arithmetic expression* grammar.

Changing the Grammar to Impose Precedence

```
< start > ::= < expr >
< expr > ::= < expr > + < expr > |
            < expr > - < expr > |
            < expr > * < expr > |
            < expr > / < expr > |
            < d > | < l >
< d >      ::= 0 | 1 | 2 | 3 | ... | 9
< l >      ::= a | b | c | ... | z
```

Original Grammar G

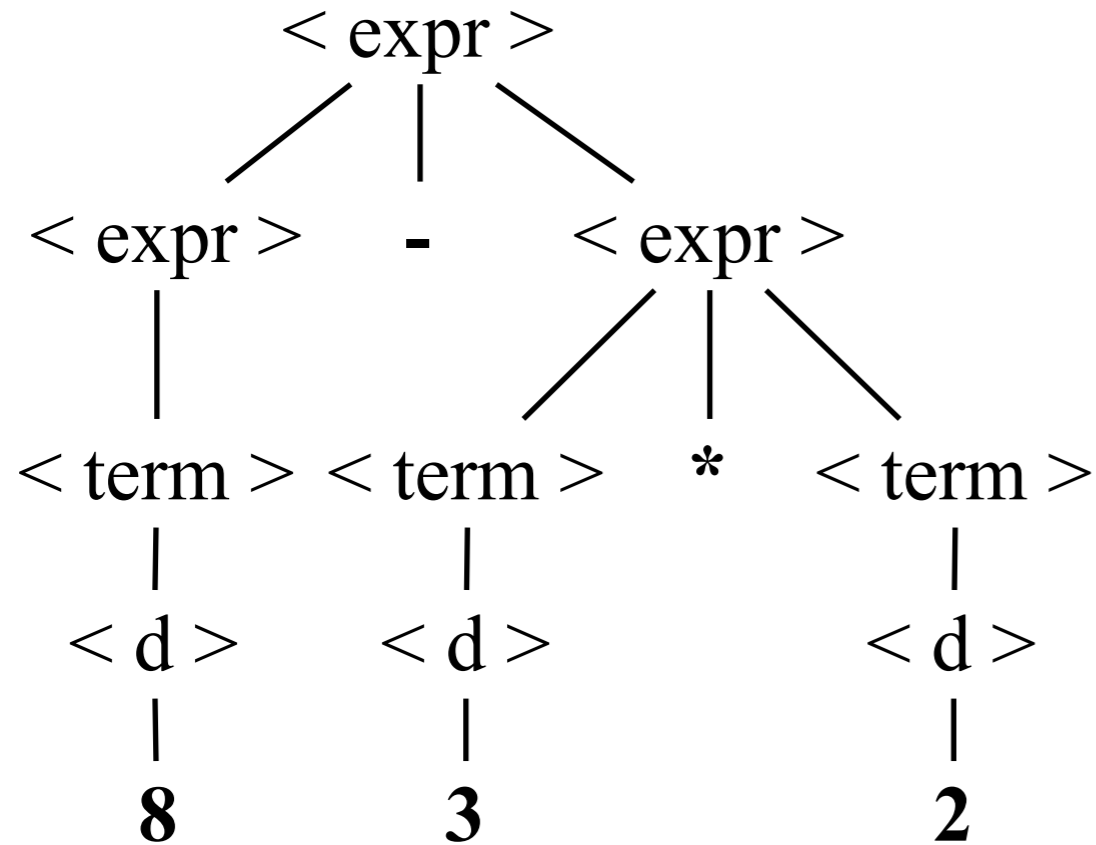


```
< start > ::= < expr >
< expr > ::= < expr > + < expr > |
            < expr > - < expr > |
            < term >
< term > ::= < term > * < term > |
            < term > / < term > |
            < d > | < l >
< d >      ::= 0 | 1 | 2 | 3 | ... | 9
< l >      ::= a | b | c | ... | z
```

Modified Grammar G'

Grouping in Parse Tree Now Reflects Precedence

Parse “8 - 3 * 2”:



Only One Possible Parse Tree

$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid$

$\langle \text{expr} \rangle - \langle \text{expr} \rangle \mid$

$\langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid$

$\langle \text{term} \rangle / \langle \text{term} \rangle \mid$

$\langle \text{d} \rangle \mid \langle 1 \rangle$

$\langle \text{d} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$

$\langle 1 \rangle ::= a \mid b \mid c \mid \dots \mid z$

Modified Grammar G'

Precedence

- *Low Precedence:*
Addition + and Subtraction -
- *Medium Precedence:*
Multiplication * and Division /
- *Highest Precedence:*
Exponentiation ^

$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid$

$\langle \text{expr} \rangle - \langle \text{expr} \rangle \mid$

$\langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid$

$\langle \text{term} \rangle / \langle \text{term} \rangle \mid$

$\langle \text{d} \rangle \mid \langle \text{l} \rangle$

$\langle \text{d} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \mathbf{3} \mid \dots \mid \mathbf{9}$

$\langle \text{l} \rangle ::= \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$

Still Have Ambiguity...

How about $3 - 2 - 1$?

$3 - 2 - 1$

OR?

$3 - 2 - 1$

$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \underline{\langle \text{expr} \rangle - \langle \text{expr} \rangle} \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid \langle \text{term} \rangle / \langle \text{term} \rangle \mid \langle \text{d} \rangle \mid \langle \text{l} \rangle$

$\langle \text{d} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \mathbf{3} \mid \dots \mid \mathbf{9}$

$\langle \text{l} \rangle ::= \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$

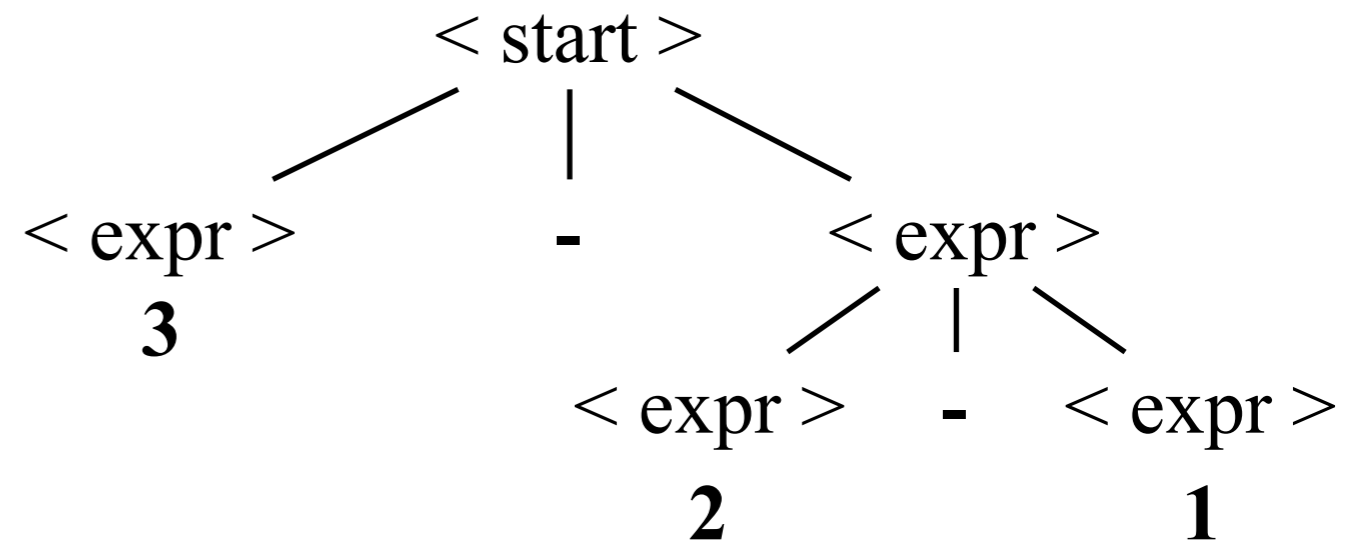
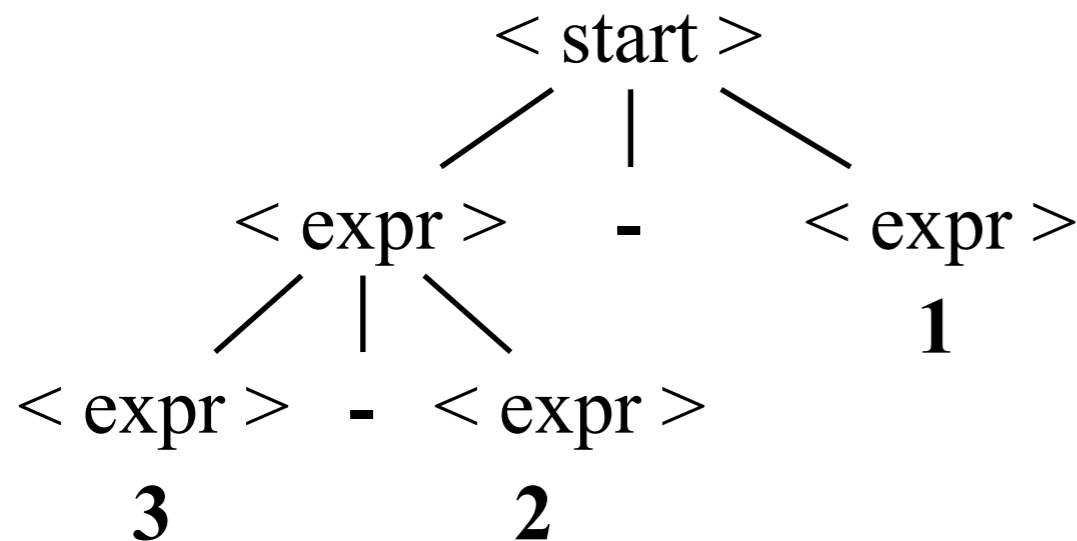
Still Have Ambiguity...

How about 3 - 2 - 1 ?

3 - 2 - 1

OR?

3 - 2 - 1



< start > ::= < expr >

< expr > ::= < expr > + < expr > | < expr > - < expr > | < term >

< term > ::= < term > * < term > | < term > / < term > | < d > | < 1 >

< d > ::= 0 | 1 | 2 | 3 | ... | 9

< 1 > ::= a | b | c | ... | z

Still Have Ambiguity...

- Grouping of operators of same precedence not disambiguated.
- Non-commutative operators: only one parse tree correct.

$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid \langle \text{term} \rangle / \langle \text{term} \rangle \mid \langle \text{d} \rangle \mid \langle \text{l} \rangle$

$\langle \text{d} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \mathbf{3} \mid \dots \mid \mathbf{9}$

$\langle \text{l} \rangle ::= \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$

Imposing Associativity

Same grammar with left / right recursion for - :

Our choices:

$\langle \text{expr} \rangle ::= \langle \text{d} \rangle - \langle \text{expr} \rangle \mid$
$\qquad \qquad \qquad \langle \text{d} \rangle$
$\langle \text{d} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \mathbf{3} \mid \dots \mid \mathbf{9}$

$\langle \text{expr} \rangle$	\Rightarrow
$\langle \text{d} \rangle - \langle \text{expr} \rangle$	\Rightarrow
$\langle \text{d} \rangle - \langle \text{d} \rangle - \langle \text{expr} \rangle$	\Rightarrow
$\langle \text{d} \rangle - \langle \text{d} \rangle - \langle \text{d} \rangle - \langle \text{expr} \rangle$	\Rightarrow
...	

Or:

$\langle \text{expr} \rangle ::= \langle \text{expr} \rangle - \langle \text{d} \rangle \mid$
$\qquad \qquad \qquad \langle \text{d} \rangle$
$\langle \text{d} \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \mathbf{3} \mid \dots \mid \mathbf{9}$

$\langle \text{expr} \rangle$	\Rightarrow
$\langle \text{expr} \rangle - \langle \text{d} \rangle$	\Rightarrow
$\langle \text{expr} \rangle - \langle \text{d} \rangle - \langle \text{d} \rangle$	\Rightarrow
$\langle \text{expr} \rangle - \langle \text{d} \rangle - \langle \text{d} \rangle - \langle \text{d} \rangle$	\Rightarrow
...	

Which one do we want for - in the calculator language?

Associativity

- Deals with operators of same precedence
- Implicit grouping or parenthesizing
- Left to right: $*$, $/$, $+$, $-$
- Right to left: $^$

Complete, Unambiguous Arithmetic Expression Grammar

$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$
 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{expr} \rangle \mid$
 $\quad \langle \text{expr} \rangle - \langle \text{expr} \rangle \mid$
 $\quad \langle \text{expr} \rangle * \langle \text{expr} \rangle \mid$
 $\quad \langle \text{expr} \rangle / \langle \text{expr} \rangle \mid$
 $\quad \langle \text{expr} \rangle ^ \langle \text{expr} \rangle \mid$
 $\quad \langle d \rangle \mid \langle l \rangle$
 $\langle d \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \mathbf{3} \mid \dots \mid \mathbf{9}$
 $\langle l \rangle ::= \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$

Original Ambiguous Grammar \mathcal{G}



$\langle \text{start} \rangle ::= \langle \text{expr} \rangle$
 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle \mid$
 $\quad \langle \text{expr} \rangle - \langle \text{term} \rangle \mid$
 $\quad \langle \text{term} \rangle$
 $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle \mid$
 $\quad \langle \text{term} \rangle / \langle \text{factor} \rangle \mid$
 $\quad \langle \text{factor} \rangle$
 $\langle \text{factor} \rangle ::= \langle g \rangle ^ \langle \text{factor} \rangle \mid$
 $\quad \langle g \rangle$
 $\langle g \rangle ::= (\langle \text{expr} \rangle) \mid \langle d \rangle \mid \langle l \rangle$
 $\langle d \rangle ::= \mathbf{0} \mid \mathbf{1} \mid \mathbf{2} \mid \dots \mid \mathbf{9}$
 $\langle l \rangle ::= \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \mid \dots \mid \mathbf{z}$

Unambiguous Grammar \mathcal{G}

Abstract versus Concrete Syntax

Concrete Syntax:

Representation of a construct in a particular language, including placement of keywords and delimiters.

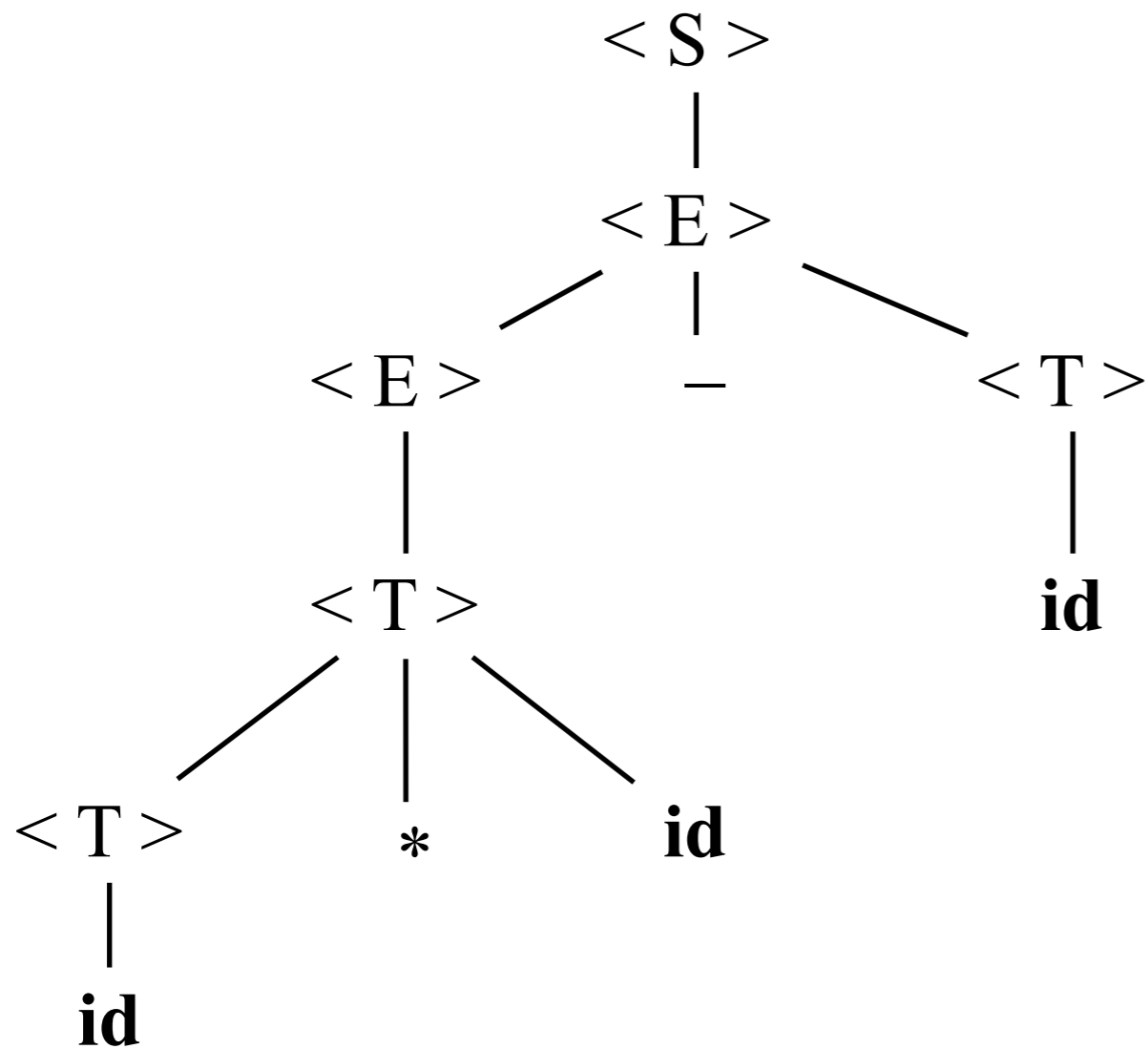
Abstract Syntax:

Structure of meaningful components of each language construct.

Example:

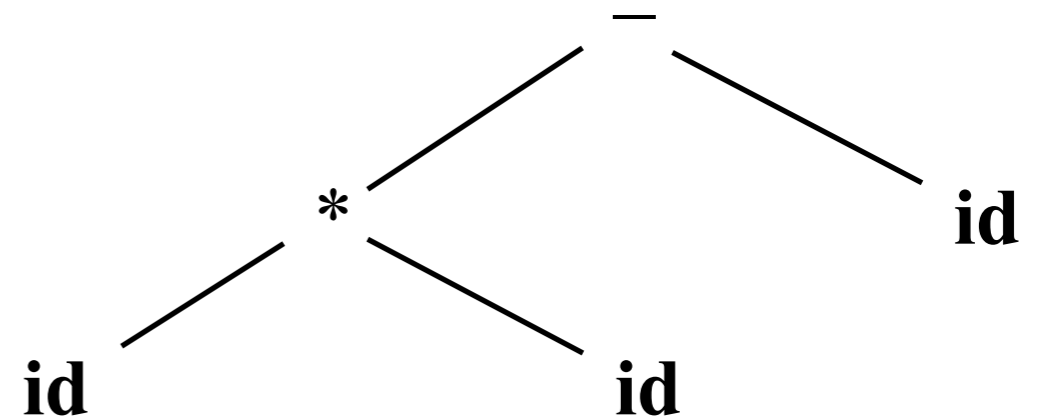
Consider $A * B - C$:

Concrete Syntax Tree



```
< S > ::= < E >
< E > ::= < E > - < T > | < T >
< T > ::= < T > * id | id
```

Abstract Syntax Tree (AST)



Abstract versus Concrete Syntax

Same abstract syntax, different concrete syntax:

Pascal **while** $x \neq A[i]$, **do**
 $i := i + 1$
end

C **while** ($x \neq A[i]$)
 $i = i + 1$;

Regular vs. Context Free

- All Regular languages are context free languages
- Not all context free languages are regular languages

Example:

$$N ::= X \mid Y$$
$$X ::= a \mid X b$$
$$Y ::= c \mid Y c$$

is equivalent to:

$$a b^* \mid c^+$$

Question:

Is $\{a^n b^n \mid n \geq 0\}$ a context free language?

Is $\{a^n b^n \mid n \geq 0\}$ a regular language?

Regular Grammars

CFGs with restrictions on the shape of production rules.

Left-linear:

$$X ::= a \mid X b$$
$$N ::= X a b$$

Right-linear:

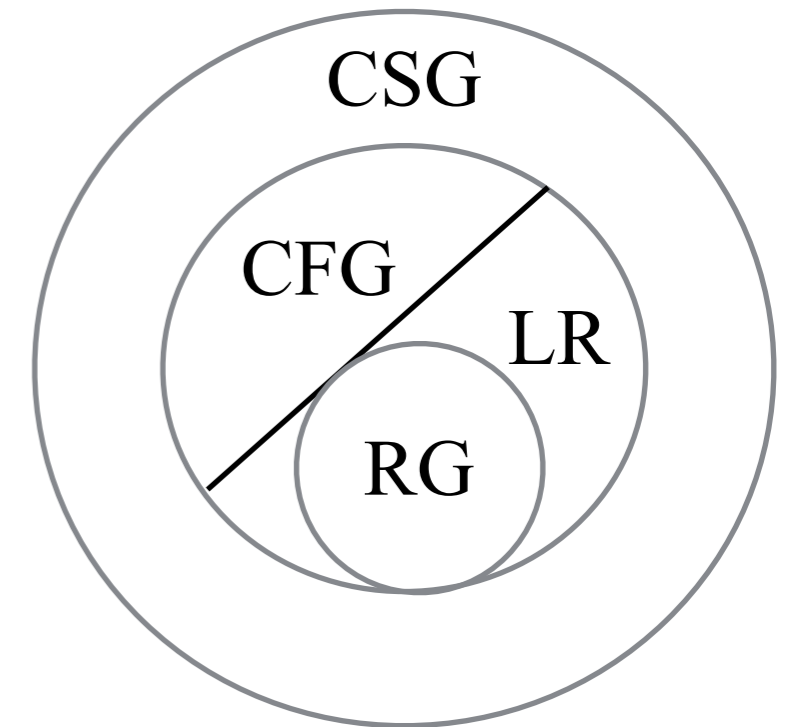
$$Y ::= a b \mid a b Y$$
$$N ::= b \mid b Y$$

Complexity of Parsing

Classification of languages that can be recognized by specific grammars.

Complexity:

Regular grammars	DFAs	$O(n)$
LR grammars	Kunth's algorithm	$O(n)$
Arbitrary CFGs	Earley's algorithm	$O(n^3)$
Arbitrary CSGs	LBAs	P-SPACE COMPLETE

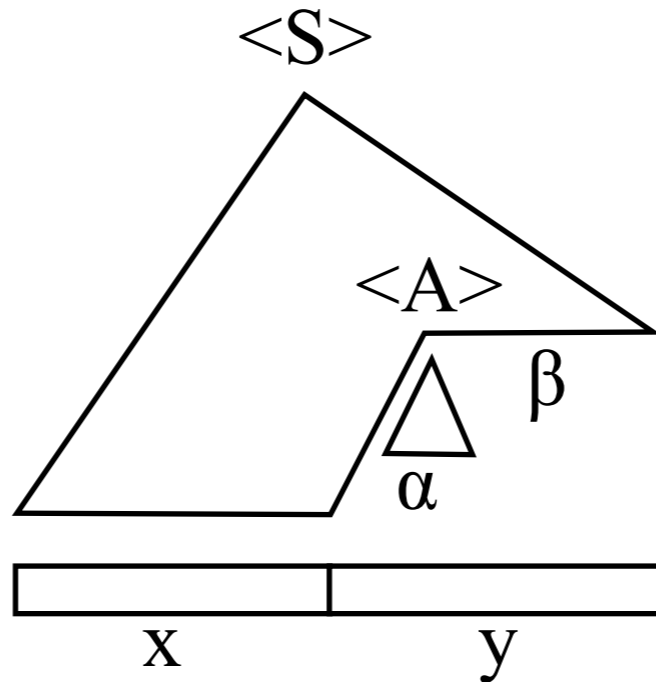


Reading:

Scott Chapter 2.3.4 (for LR parser) and Chapter 2.4.3 for language class.

Earley, Jay (1970), "An efficient context-free parsing algorithm", Communications of the ACM.

Top - Down Parsing - LL(1)



Basic Idea:

- The parse tree is constructed from the root, expanding non-terminal nodes on the tree's frontier following a **leftmost** derivation.
- The input program is read from **left** to right, and input tokens are read (consumed) as the program is parsed.
- The next non-terminal symbol is replaced using one of its rules. The particular choice has to be unique and uses parts of the input (partially parsed program), for instance the first **token** of the remaining input.

Top - Down Parsing - LL(1) (cont.)

Example:

$S ::= a S b \mid \varepsilon$

How can we parse (automatically construct a leftmost derivation) the input string **a a a b b b** using a PDA (push-down automaton) and only the first symbol of the remaining input?

INPUT:

a a a b b b eof

Predictive Parsing

Basic idea:

For any two productions $A ::= \alpha \mid \beta$, we would like a distinct way of choosing the correct production to expand.

For some rhs $\alpha \in G$, define **FIRST**(α) as the set of tokens that appear as the first symbol in some string derived from α .

That is

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

$\varepsilon \in \text{FIRST}(\alpha)$ iff $\alpha \Rightarrow^* \varepsilon$

Predictive Parsing

For a non-terminal A , define **FOLLOW**(A) as the set of terminals that can appear immediately to the right of A in some sentential form.

Thus, a non-terminal's **FOLLOW** set specifies the tokens that can legally appear after it. A terminal symbol has no **FOLLOW** set.

FIRST and **FOLLOW** sets can be constructed automatically

Predictive Parsing (cont.)

Key Property:

Whenever two productions $A ::= \alpha$ and $A ::= \beta$ both appear in the grammar, we would like

- $FIRST(\alpha) \cap FIRST(\beta) = \emptyset$, and
- if $\alpha \Rightarrow^* \varepsilon$, then $FIRST(\beta) \cap FOLLOW(A) = \emptyset$
- Analogue case for $\beta \Rightarrow^* \varepsilon$. Note: due to first condition, at most one of α and β can derive ε .

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

LL(1) Grammar

Define $FIRST^+(\delta)$ for rule $A ::= \delta$

- $FIRST(\delta) - \{ \varepsilon \} \cup \text{Follow}(A)$, if $\varepsilon \in FIRST(\delta)$
- $FIRST(\delta)$ otherwise

A Grammar is LL(1) iff

$(A ::= \alpha \text{ and } A ::= \beta)$ implies

$$FIRST^+(\alpha) \cap FIRST^+(\beta) = \emptyset$$

Next Lecture

Next Time:

- Review and more on LL(1) parsing
- Read Scott, Chapter 2.3.1 - 2.3.2 (and some chapter on companion site)