# Compilers

Arthur Hoskey, Ph.D. Farmingdale State College Computer Systems Department

#### Parsing (syntactic analysis)



 Parsing - The process of constructing a derivation from a specific input sentence.

 Taken from Engineering a Compiler 2<sup>nd</sup> edition by Cooper and Torczon, 2012.



- Syntatic Analysis Check if the source code conforms to the rules of the language.
- The rules of the language are defined using a context-free grammar.



- Parser The parser takes a stream of tokens as input and produces an abstract syntax tree.
- The parser's input stream of tokens is generated by the scanner during the lexical analysis phase.



• Parse tree examples...

# **Parsing Tree Examples**

• Productions  $A \rightarrow id = Exp$   $Exp \rightarrow Exp Op id \mid id$  $Op \rightarrow + \mid - \mid * \mid /$  Terminals: id, =, +, -, \*, / Nonterminals: A, Exp, Op Start: A

Parse the following: w = x + y - z



• Productions  $A \rightarrow id = Exp$   $Exp \rightarrow Exp Op id | id$   $Op \rightarrow + | - | * | /$ Parse the following: w = x + y - z Terminals: id, =, +, -, \*, / Nonterminals: A, Exp, Op Start: A

**Start with A** 













Terminals: id, =, +, -, \*, / Nonterminals: A, Exp, Op Start: A

**Exp**→id

Op→-



Op→+



Terminals: id, =, +, -, \*, / Nonterminals: A, Exp, Op Start: A

Variables w, x, y, z will respectively be associated with the id nodes



x+y will be evaluated first.
z will be subtracted from the result of x+y because of the tree structure.

#### **Parsing Example 1 - Observations**



# **Parsing Example 2**

Productions
A → id = Exp
Exp → Exp Op id | id
Op → + | - | \* | /
Parse the following: w = x + y \* z

**Start with A** 











• Productions  $A \rightarrow id = Exp$   $Exp \rightarrow Exp Op id | id$  $Op \rightarrow + | - | * | /$ 

Having + and \* defined at the same level in the grammar means they have the same precedence

Parse the following:  $\mathbf{w} = \mathbf{x} + \mathbf{y} * \mathbf{z}$ 

Variables w, x, y, z will respectively be associated with the id nodes

#### **ERROR**

x+y will STILL be evaluated first. This does not follow the normal order of operations.
This is incorrect because z will be multiplied with the result of x+y (because of the tree structure). Need a grammar that expresses the correct precedence.



#### **Parsing Example 2 - Observations**

Productions
A → id = Exp
Exp → Exp + Term | Exp - Term | Term
Term → Term \* Fact | Term / Fact | Fact
Fact → id | num

• Parse the following: w = x + y \* z

This grammar gives higher precedence to \* and /.

Operations that are deeper in the parse tree have higher precedence.

#### **Grammar for Correct Expression Precedence**

Productions
A → id = Exp
Exp → Exp + Term | Exp - Term | Term
Term → Term \* Fact | Term / Fact | Fact
Fact → id | num
Parse the following: w = x + y \* z





# Productions A → id = Exp Exp → Exp + Term | Exp - Term | Term Term → Term \* Fact | Term / Fact | Fact Fact → id | num Parse the following: w = x + y \* z











# Productions A → id = Exp Exp → Exp + Term | Exp - Term | Term Term → Term \* Fact | Term / Fact | Fact Fact → id | num | (Exp)

Parse the following: w = (x + y) \* z

Add parenthesis to the grammar to create higher precedence for + and – when necessary.

Adding the parenthesis at a deeper level in the grammar compared to Exp and Term gives the parenthesis higher precedence.

#### **Add Parenthesis to the Expression Grammar**

```
Productions
A → id = Exp
Exp → Exp + Term | Exp - Term | Term
Term → Term * Fact | Term / Fact | Fact
Fact → id | num | (Exp )
Parse the following: w = (x + y) * z
```

Start with A



Productions
A → id = Exp
Exp → Exp + Term | Exp - Term | Term
Term → Term \* Fact | Term / Fact | Fact
Fact → id | num | (Exp )
Parse the following: w = (x + y) \* z



Fact→(Exp)



• Productions  $A \rightarrow id = Exp$   $Exp \rightarrow Exp + Term | Exp - Term | Term$   $Term \rightarrow Term * Fact | Term / Fact | Fact$   $Fact \rightarrow id | num | (Exp )$ Parse the following: w = (x + y) \* z



Productions •  $A \rightarrow id = Exp$  $Exp \rightarrow Exp + Term | Exp - Term | Term$ Term → Term \* Fact | Term / Fact | Fact Fact  $\rightarrow$  id | num | (Exp) Parse the following: w = (x + y) \* z



Productions •  $A \rightarrow id = Exp$  $Exp \rightarrow Exp + Term | Exp - Term | Term$ Term → Term \* Fact | Term / Fact | Fact Fact  $\rightarrow$  id | num | (Exp) Parse the following: w = (x + y) \* z



Productions •  $A \rightarrow id = Exp$  $Exp \rightarrow Exp + Term | Exp - Term | Term$ Term → Term \* Fact | Term / Fact | Fact Fact  $\rightarrow$  id | num | (Exp) Parse the following: w = (x + y) \* z



• Productions  $A \rightarrow id = Exp$   $Exp \rightarrow Exp + Term | Exp - Term | Term$   $Term \rightarrow Term * Fact | Term / Fact | Fact$   $Fact \rightarrow id | num | (Exp )$ Parse the following: w = (x + y) \* z



• Productions  $A \rightarrow id = Exp$   $Exp \rightarrow Exp + Term | Exp - Term | Term$   $Term \rightarrow Term * Fact | Term / Fact | Fact$   $Fact \rightarrow id | num | (Exp )$ Parse the following: w = (x + y) \* z

Variables w, x, y, z will respectively be associated with the id nodes

 (x+y) will be evaluated first. z will be multiplied by the result of (x+y) because of the tree structure. The \*
 operator must wait for its lhs value to be computed.

### **Parsing Example 4 - Observations**



- Left recursion The first symbol of the rhs is the same as the lhs nonterminal.
- For example:
- A → Ay | z Left recursion because the lhs symbol A is the first symbol on the rhs



- Productions
   A → Ay | z
- Parse the following: zyyy

#### **Left Recursive Grammar**
Productions
 A → Ay | z

#### • Parse the following: zyyy



- Right recursion The last symbol of the rhs is the same as the lhs nonterminal.
- For example:

A → zB B → yB | ε N Right recursion because the lhs symbol B is the last symbol on the rhs

# **Right Recursion**

- Productions A  $\rightarrow$  zB B  $\rightarrow$  yB |  $\epsilon$
- Parse the following: zyyy



• Now on to LL(1) grammars...



## <u>LL(1)</u>

- LL(1) is a top-down parsing technique. We begin with the start symbol and apply substitutions.
- First L in LL(1). The first L means to process the input from left to right.
- Second L in LL(1). The second L stands for a leftmost derivation (this means always expand the leftmost nonterminal).
- The 1 in LL(1). The 1 means there is a 1 symbol lookahead.
- LL(1) grammars are only able to parse a subset of all context-free grammars.



- LL(1) parsing requires that when we are doing a substitution(expanding a nonterminal), the choice of production to use is unambiguous.
- When applying a substitution for a nonterminal, there should only be one possible substitution given one token of lookahead.
- For example, assume we are expanding nonterminal B in the following grammar and the lookahead is 'c'. The choice of which production to use is unambiguous.



- Assume the following grammar.
- $A \rightarrow Bw$
- $B \rightarrow sC$
- $B \rightarrow sz$
- $C \rightarrow y$

### Question

If we are expanding nonterminal B and the lookahead is 's', is this valid when using LL(1) parsing on the grammar above?

# **Requirements for Parsing LL(1)**

- Assume the following grammar.
- $A \rightarrow Bw$
- $B \rightarrow sC$  $B \rightarrow sz$

 $C \rightarrow y$ 

Ambiguous. If we are expanding B and the lookahead is 's', we do NOT know which production to use. It is ambiguous (both production start with 's'). This grammar is not valid for LL(1) parsing.

## Question

If we are expanding nonterminal B and the lookahead is 's', is this valid when using LL(1) parsing on the grammar above?

## <u>Answer</u>

#### No

**Requirements for Parsing LL(1)** 

- The first+ sets of productions for a given nonterminal must be disjoint.
- In the following grammar the first+ sets for B's productions are not disjoint (they both include s).



There is no way to know which substitution to apply since s is in both first+ sets!!!

First+( $B \rightarrow sC$ ) = { s } First+( $B \rightarrow sD$ ) = { s }

## **Requirements for Parsing LL(1)**

- Assume the following grammar.
- $A \rightarrow Bw$
- $B \rightarrow Cs$
- $B \rightarrow sz$
- $C \rightarrow y$
- $C \rightarrow \epsilon$

#### <u>Question</u> Are the First+ sets of the productions for B disjoint?

# **Requirements for Parsing LL(1)**



<u>Question</u> Are the First+ sets of the productions for B disjoint?

#### <u>Answer</u>

No. The First+ sets of  $B \rightarrow Cs$  and  $B \rightarrow sD$  intersect (they are not disjoint). This grammar is NOT LL(1).

**Requirements for Parsing LL** 

- Assume the following grammar.
- $A \rightarrow Bs$
- $B \rightarrow C$
- $B \rightarrow sz$
- $C \rightarrow y$
- $C \rightarrow \epsilon$

#### <u>Question</u> Are the First+ sets of the productions for B disjoint?

# **Requirements for Parsing LL(1)**

- Assume the following grammar.
- $A \rightarrow Bs$ C can disappear in this production because of  $C \rightarrow \epsilon$ . $B \rightarrow C$ Since the right side of this production can<br/>disappear, we must include members of Follow(B)<br/>in First+( $B \rightarrow C$ ).<br/>Note:First(C)={y,  $\epsilon$ } and Follow(B)={s}.

First+( $B \rightarrow C$ ) = First(C) + Follow(B) = {y,  $\epsilon$ , s} First+( $B \rightarrow sz$ ) = {s}

Question

Are the First+ sets of the productions for B disjoint?

#### **Answer**

No. The First+ sets of  $B \rightarrow C$  and  $B \rightarrow sz$  intersect (they are not disjoint). This grammar is NOT LL(1).

**Requirements for Parsing LL** 

- Assume the following grammar.
- $A \rightarrow Bw$
- $B \rightarrow C$
- $B \rightarrow sz$
- $C \rightarrow y$
- $C \rightarrow \epsilon$

## <u>Question</u> Are the First+ sets of the productions for B disjoint?

# **Requirements for Parsing LL(1)**



#### **Answer**

Yes. The First+ sets of  $B \rightarrow C$  and  $B \rightarrow sz$  are disjoint. This grammar is LL(1) because all productions for each nonterminal have disjoint First+ sets.

**Requirements for Parsing LL(1)** 

### **Transforming Grammars to LL(1)**

- You can try and transform a grammar to LL(1) using the following techniques:
  - Change to right recursive.
  - Use left factoring
- Every context free grammar is not LL(1), so these techniques are not guaranteed to work.



#### **Predictive Parser**

- A predictive parser is able to determine the correct production to use given a small number of lookahead symbols.
- The right recursive form of a grammar might work for LL(1).
- Converting to right recursive form will allow you to use a predictive parser.
- If a grammar is LL(1) then you can build a recursive descent parser for it. Recursive descent parsing means parsing from the start symbol down (top-down parsing).

## **Predictive Parser**

#### **Convert the Expression Grammar to Right Recursive**

Original Expression Grammar (without the assignment)
 Exp → Exp + Term | Exp - Term | Term
 Term → Term \* Fact | Term / Fact | Fact
 Fact → id | num | (Exp )

Right Recursive Expression Grammar
 Exp → Term ExpEnd
 ExpEnd → + Term ExpEnd | - Term ExpEnd | ε
 Term → Fact TermEnd
 TermEnd → \* Fact TermEnd | / Fact TermEnd | ε
 Fact → id | num | (Exp )

## **Right Recursive Expression Grammar**

## Now on to right recursive grammars...

## **Right Recursive Grammar**

```
Productions (Exp is start symbol)
Exp → Term ExpEnd
ExpEnd → + Term ExpEnd | - Term ExpEnd | ε
Term → Fact TermEnd
TermEnd → * Fact TermEnd | / Fact TermEnd | ε
Fact → id | num | (Exp )
Parse the following: a+2*b
```























#### **Requirement for Parsing LL(1) - Revisited**

- Assume the following grammar.
- $\mathsf{A} \rightarrow \mathsf{cd}$
- $A \rightarrow ce$
- $A \rightarrow cf$

<u>Question</u> Are the First+ sets of the productions for nonterminal A disjoint?

## **Requirements for Parsing LL(1) -Revisited**

#### **Requirement for Parsing LL(1) - Revisited**

• Assume the following grammar.

 $A \rightarrow cd$ First+ sets are NOT disjoint (grammar is not LL(1)) $A \rightarrow ce$ First+( $A \rightarrow cd$ ) = {c} $A \rightarrow cf$ First+( $A \rightarrow ce$ ) = {c}First+( $A \rightarrow cf$ ) = {c}

Question

Are the First+ sets of the productions for nonterminal A disjoint?

#### <u>Answer</u> No. They all contain c. This grammar is not LL(1).

We can transform these productions so that they are disjoint though! Requirements for Parsing LL(1) -Revisited

#### Left Factoring Rule

• Assume the following grammar. A  $\rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid ... \mid \alpha\beta_n \mid \gamma_1 \mid \gamma_2 \mid ... \mid \gamma_n$  Greek symbols: α is alpha β is beta γ is gamma

- To left factor A move the symbols following a into new productions.
- Grammar after left factoring.
- $A \rightarrow \mathbf{aB} | \gamma_1 | \gamma_2 | ... | \gamma_n$  $B \rightarrow \beta_1 | \beta_2 | ... | \beta_n$

α originally appeared in multiple A productions but it has now been
 factored out into one A production. The γ productions in A remain unchanged because they do not contain α.

 Taken from Engineering a Compiler 2<sup>nd</sup> edition by Cooper and Torczon, 2012.

**Left Factoring Rule** 

#### Left Factoring

• Assume the following grammar.

- $A \rightarrow cd$
- $A \rightarrow ce$
- $A \rightarrow cf$

Question

What does the grammar look like after left factoring nonterminal A's productions?



# Left Factoring• Assume the following grammar. $A \rightarrow cd$ $A \rightarrow cd$ $A \rightarrow ce$ $A \rightarrow cf$ Leave c in production A and<br/>move d, e, and f into the new B<br/>productions

#### Question

What does the grammar look like after left factoring nonterminal A's productions?

#### **Answer**

- $A \rightarrow cB$ The First+ sets are $B \rightarrow d$ now disjoint
- $B \rightarrow d$  $B \rightarrow e$  $B \rightarrow f$

# **Left Factoring**
### Left Factoring

• Assume the following grammar.

- $A \rightarrow cd$
- $A \rightarrow ce$
- $\mathsf{A} \to \mathsf{x}$

 $\mathsf{A} \to \mathsf{y}$ 

Question

What does the grammar look like after left factoring nonterminal A's productions?



© 2023 Arthur Hoskey. All rights reserved.

### Left Factoring

- Assume the following grammar.
- $A \rightarrow cd$  $A \rightarrow ce$

 $A \rightarrow x$ 

 $A \rightarrow y$ 

# Leave c in production A and move d, e, and f into the new B productions

#### Question

What does the grammar look like after left factoring nonterminal A's productions?

## When applying the left factoring rule:

Answer	c is a
$A \rightarrow cB$	
$A \rightarrow x$	a is p <sub>1</sub>
$\Delta \rightarrow \gamma$	e is $\beta_2$
	x is y <sub>1</sub>
	v is v
$P \rightarrow G$	<b>y</b> 10 <b>y</b> 2







© 2023 Arthur Hoskey. All rights reserved.