

A Rule of While Language Syntax

// Where things work very nicely for recursive descent!

statmt ::=

println (stringConst , ident)

| ident = expr

| if (expr) statmt (else statmt)?

| while (expr) statmt

| { statmt }*

Parser for the `statmt` (rule \rightarrow code)

```
def skip(t : Token) = if (lexer.token == t) lexer.next
  else error("Expected"+ t)
// statmt ::=
def statmt = {
  // println ( stringConst , ident )
  if (lexer.token == Println) { lexer.next;
    skip(openParen); skip(stringConst); skip(comma);
    skip(identifier); skip(closedParen)
  // | ident = expr
  } else if (lexer.token == Ident) { lexer.next;
    skip(equality); expr
  // | if ( expr ) statmt (else statmt)?
  } else if (lexer.token == ifKeyword) { lexer.next;
    skip(openParen); expr; skip(closedParen); statmt;
    if (lexer.token == elseKeyword) { lexer.next; statmt }
  // | while ( expr ) statmt
```

Continuing Parser for the Rule

```
// | while ( expr ) statmt
```

```
} else if (lexer.token == whileKeyword) { lexer.next;  
    skip(openParen); expr; skip(closedParen); statmt
```

```
// | { statmt* }
```

```
} else if (lexer.token == openBrace) { lexer.next;  
    while (isFirstOfStatmt) { statmt }  
    skip(closedBrace)
```

```
} else { error("Unknown statement, found token " +  
    lexer.token) }
```

How the parser decides which alternative to follow?

```
statmt ::= println ( stringConst , ident )
        | ident = expr
        | if ( expr ) statmt (else statmt)?
        | while ( expr ) statmt
        | { statmt* }
```

- Look what each alternative starts with to decide what to parse
- Here: we have terminals at the beginning of each alternative!
- More generally, we have ‘first’ computation, as for regular expressions

- Consider a grammar G and non-terminal N

$L_G(N) = \{ \text{set of strings that } N \text{ can derive} \}$

e.g. $L(\text{statmt})$ – all statements of while language

$\text{first}(N) = \{ a \mid aw \text{ in } L_G(N), a - \text{terminal}, w - \text{string of terminals} \}$

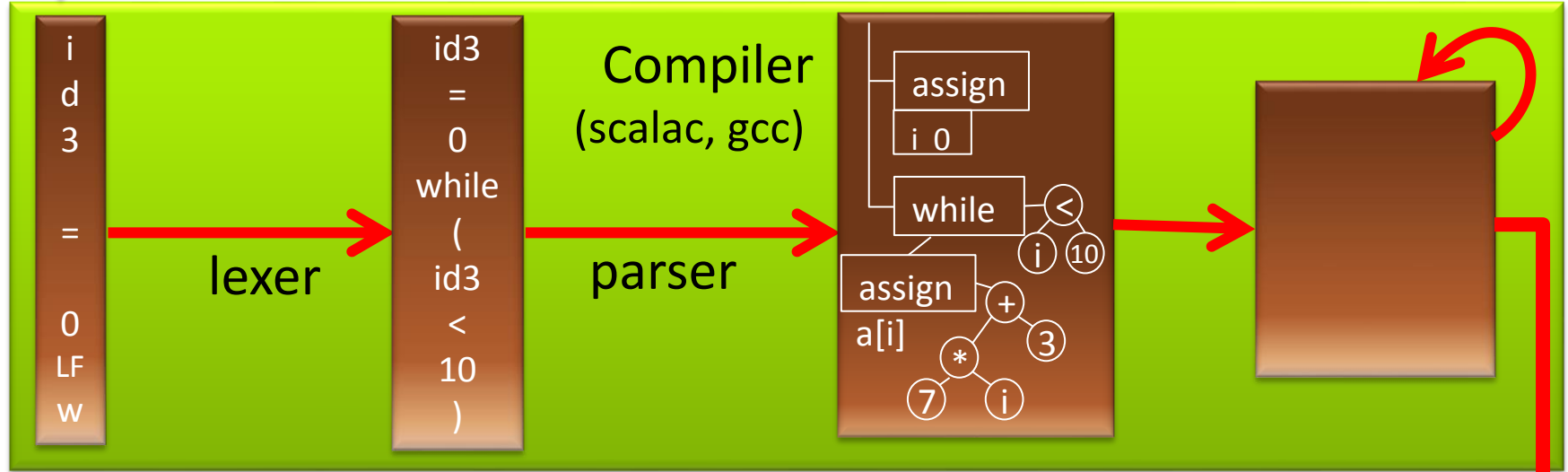
$\text{first}(\text{statmt}) = \{ \text{println, ident, if, while, } \{ \} \}$

$\text{first}(\text{while (expr) statmt}) = \{ \text{while} \}$

Compiler Construction

source code

```
id3 = 0  
while (id3 < 10) {  
  println("",id3);  
  id3 = id3 + 1 }  
}
```



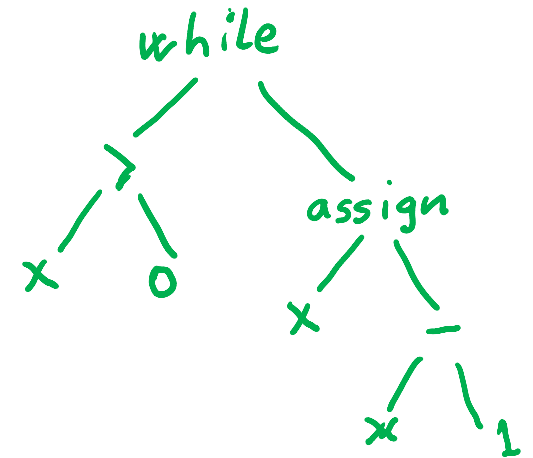
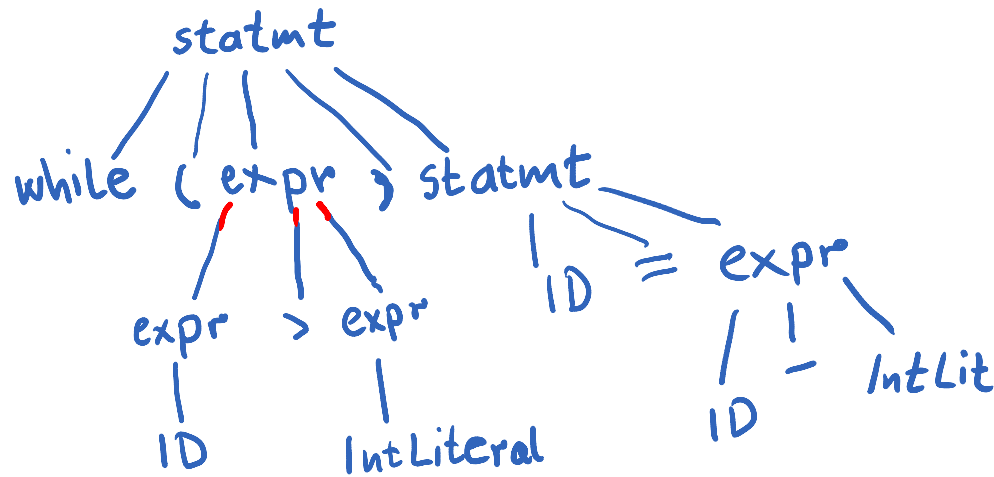
characters

words
(tokens)

trees

Parse Tree vs Abstract Syntax Tree (AST)

while (x > 0) x = x - 1



Pretty printer: takes abstract syntax tree (AST) and outputs the leaves of one possible (concrete) parse tree.

$\text{parse}(\text{prettyPrint}(\text{ast})) \approx \text{ast}$

Parse Tree vs Abstract Syntax Tree (AST)

- Each node in **parse tree** has children corresponding **precisely to right-hand side of grammar rules**. The definition of parse trees is fixed given the grammar
 - **Often compiler never actually builds parse trees in memory**
- Nodes in **abstract syntax tree (AST)** contain only useful information and usually omit the punctuation signs. We can choose our own syntax trees, to make it convenient for both construction in parsing and for later stages of compiler or interpreter
 - **A compiler typically directly builds AST**

Abstract Syntax Trees for Statements

`statmt ::= println (stringConst , ident)`

`| ident = expr`

`| if (expr) statmt (else statmt)?`

`| while (expr) statmt`

`| { statmt* }`

abstract class `Statmt`

case class `PrintlnS(msg : String, var : Identifier)` **extends** `Statmt`

case class `Assignment(left : Identifier, right : Expr)` **extends** `Statmt`

case class `If(cond : Expr, trueBr : Statmt,`
`falseBr : Option[Statmt])` **extends** `Statmt`

case class `While(cond : Expr, body : Expr)` **extends** `Statmt`

case class `Block(sts : List[Statmt])` **extends** `Statmt`

Abstract Syntax Trees for Statements

```
statmt ::= println ( stringConst , ident )
        | ident = expr
        | if ( expr ) statmt (else statmt)?
        | while ( expr ) statmt
        | { statmt* }
```

abstract class Statmt

case class PrintlnS(msg : String, var : Identifier) **extends** Statmt

case class Assignment(left : Identifier, right : Expr) **extends** Statmt

case class If(cond : Expr, trueBr : Statmt,
falseBr : Option[Statmt]) **extends** Statmt

case class While(cond : Expr, body : Statmt) **extends** Statmt

case class Block(sts : List[Statmt]) **extends** Statmt

Our Parser Produced Nothing 😞

```
def skip(t : Token) : unit = if (lexer.token == t) lexer.next  
  else error("Expected"+ t)
```

```
// statmt ::=
```

```
def statmt : Unit = {
```

```
  // println ( stringConst , ident )
```

```
  if (lexer.token == Println) { lexer.next;
```

```
    skip(openParen); skip(stringConst); skip(comma);
```

```
    skip(identifier); skip(closedParen)
```

```
  // | ident = expr
```

```
  } else if (lexer.token == Ident) { lexer.next;
```

```
    skip(equality); expr
```

New Parser: Returning an AST 😊

```
def expect(t : Token) : Token = if (lexer.token == t) { lexer.next;t}
  else error("Expected"+ t)
// statmt ::=
def statmt : Statmt = {
  // println ( stringConst , ident )
  if (lexer.token == Println) { lexer.next;
    skip(openParen); val s = getString(expect(stringConst));
    skip(comma);
    val id = getIdent(expect(identifier)); skip(closedParen)
    PrintlnS(s, id)
  // | ident = expr
  } else if (lexer.token.class == Ident) { val lhs = getIdent(lexer.token)
    lexer.next;
    skip(equality); val e = expr
    Assignment(lhs, e)
```

Constructing Tree for 'if'

```
def expr : Expr = { ... }
```

```
// statmt ::=
```

```
def statmt : Statmt = {
```

```
...
```

```
// if ( expr ) statmt (else statmt)?
```

```
// case class If(cond : Expr, trueBr: Statmt, falseBr: Option[Statmt])
```

```
  } else if (lexer.token == ifKeyword) { lexer.next;  
    skip(openParen); val c = expr; skip(closedParen);
```

```
    val trueBr = statmt
```

```
    val elseBr = if (lexer.token == elseKeyword) {  
      lexer.next; Some(statmt) } else None
```

```
    If(c, trueBr, elseBr) // made a tree node 😊
```

```
  }
```

Task: Constructing AST for 'while'

```
def expr : Expr = { ... }
```

```
// statmt ::=
```

```
def statmt : Statmt = {
```

```
  ...
```

```
  // while ( expr ) statmt
```

```
  // case class While(cond : Expr, body : Expr) extends Statmt
```

```
  } else if (lexer.token == WhileKeyword) {
```

```
  } else
```

Here each alternative started with
different token

statmt ::=

println (stringConst , ident)
| ident = expr
| if (expr) statmt (else statmt)?
| while (expr) statmt
| { statmt* }

What if this is not the case?

Left Factoring Example: Function Calls

statmt ::=

println (stringConst , ident)

foo = 42 + x

foo (u , v)

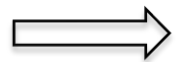


| ident = expr

| if (expr) statmt (else statmt)?

| while (expr) statmt

| { statmt* }



| ident (expr (, expr)*)

code to parse the grammar:

```
} else if (lexer.token.class == Ident) {
```

```
    ???
```

```
}
```

Left Factoring Example: Function Calls

statmt ::=

println (stringConst , ident)



| ident assignmentOrCall

| if (expr) statmt (else statmt)?

| while (expr) statmt

| { statmt* }

assignmentOrCall ::= "=" expr | (expr (, expr)*)

code to parse the grammar:

```
} else if (lexer.token.class == Ident) {
```

```
    val id = getIdentifier(lexer.token); lexer.next
```

```
    assignmentOrCall(id)
```

```
} // Factoring pulls common parts from alternatives
```


Beyond Statements: Parsing Expressions

While Language with Simple Expressions

`statmt ::=`

`println (stringConst , ident)`

`| ident = expr`

`| if (expr) statmt (else statmt)?`

`| while (expr) statmt`

`| { statmt* }`

`expr ::= intLiteral | ident`

`| expr (+ | /) expr`

Abstract Syntax Trees for Expressions

```
expr ::= intLiteral | ident  
      | expr + expr | expr / expr
```

abstract class Expr

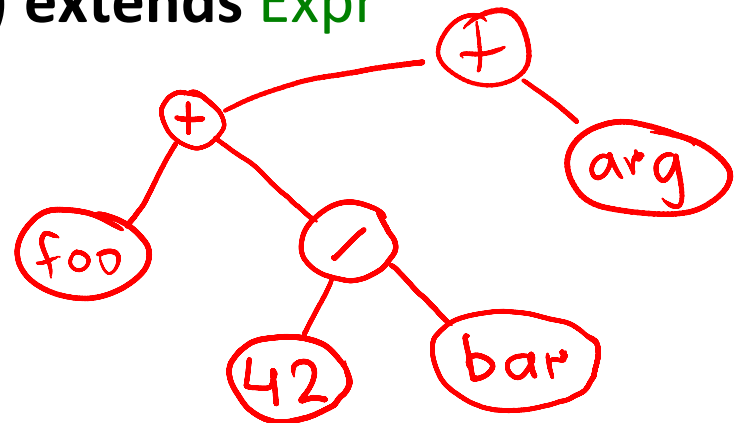
case class IntLiteral(x : Int) **extends** Expr

case class Variable(id : Identifier) **extends** Expr

case class Plus(e1 : Expr, e2 : Expr) **extends** Expr

case class Divide(e1 : Expr, e2 : Expr) **extends** Expr

foo + 42 / bar + arg



Parser That Follows the Grammar?

```
expr ::= intLiteral | ident  
      | expr + expr | expr / expr
```

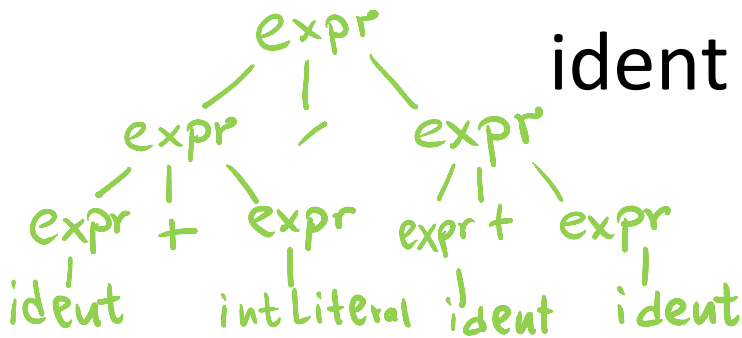
input:
foo + 42 / bar + arg

```
def expr : Expr = {  
  if (??) IntLiteral(getInt(lexer.token))  
  else if (??) Variable(getIdent(lexer.token))  
  else if (??) {  
    val e1 = expr; val op = lexer.token; val e2 = expr  
    op match Plus {  
      case PlusToken => Plus(e1, e2)  
      case DividesToken => Divides(e1, e2)  
    }  
  }  
}
```

When should parser enter the recursive case?!

Ambiguous Grammars

```
expr ::= intLiteral | ident  
      | expr + expr | expr / expr
```



ident + intLiteral / ident + ident

Each node in parse tree is given by one grammar alternative.

Ambiguous grammar: if some token sequence has multiple parse trees (then it is has multiple abstract trees).

Ambiguous grammar: if some token sequence
has multiple parse trees
(then it is usually has multiple abstract trees)

Two parse trees, each following the grammar,
their leaves both give the same token
sequence.

Ambiguous Expression Grammar

```
expr ::= intLiteral | ident  
      | expr + expr | expr / expr
```

ident + intLiteral / ident + ident

Each node in parse tree is given by one grammar alternative.

Show that the input above has two parse trees!

Exercise: Balanced Parentheses I

Show that the following balanced parentheses grammar is ambiguous (by finding two parse trees for some input sequence).

$$B ::= \varepsilon \mid (B) \mid B B$$

Remark

- The same parse tree can be derived using two different derivations, e.g.

$B \rightarrow (B) \rightarrow (BB) \rightarrow ((B)B) \rightarrow ((B)) \rightarrow (())$

$B \rightarrow (B) \rightarrow (BB) \rightarrow ((B)B) \rightarrow (())B \rightarrow (())$

this correspond to different orders in which nodes in the tree are expanded.

- Ambiguity refers to the fact that there are actually multiple *parse trees*, not just multiple derivations.

Exercise: Balanced Parentheses

Show that the following balanced parentheses grammar is ambiguous (by finding two parse trees for some input sequence) **and find unambiguous grammar for the same language.**

$$B ::= \varepsilon \mid (B) \mid B B$$

Not Quite Solution

- This grammar:

$$B ::= \varepsilon \mid A$$
$$A ::= () \mid A A \mid (A)$$

solves the problem with multiple ε symbols generating different trees.

Does string $()()()$ have a unique parse tree?

Solution for Unambiguous Parenthesis Grammar

- Proposed solution:

$$B ::= \varepsilon \mid B (B)$$

- How to come up with it?
- Clearly, rule $B ::= B B$ generates any sequence of B's. We can also encode it like this:

$$B ::= C^*$$

$$C ::= (B)$$

- Now we express sequence using recursive rule that does not create ambiguity:

$$B ::= \varepsilon \mid C B$$

$$C ::= (B)$$

- but now, look, we "inline" C back into the rules for so we get exactly the rule

$$B ::= \varepsilon \mid B (B)$$

This grammar is not ambiguous and is the solution. We did not prove unambiguity (we only tried to find ambiguous trees but did not find any).

Exercise:

Left Recursive and Right Recursive

We call a production rule “left recursive” if it is of the form

$$A ::= A p$$

for some sequence of symbols p . Similarly, a “right-recursive” rule is of a form

$$A ::= q A$$

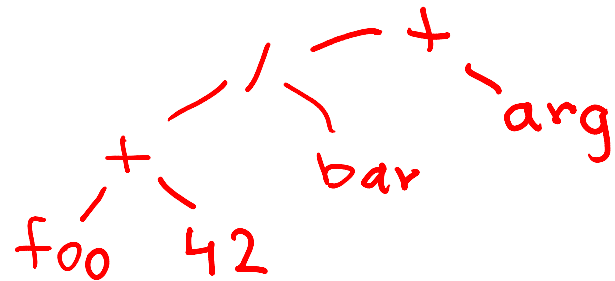
Is every context free grammar that contains both left and right recursive rule for a some nonterminal A ambiguous?

An attempt to rewrite the grammar

```
expr ::= simpleExpr (( + | / ) simpleExpr)*  
simpleExpr ::= intLiteral | ident
```

```
def simpleExpr : Expr = { ... }  
def expr : Expr = {  
  var e = simpleExpr  
  while (lexer.token == PlusToken ||  
         lexer.token == DividesToken) {  
    val op = lexer.token  
    val eNew = simpleExpr  
    op match {  
      case TokenPlus => { e = Plus(e, eNew) }  
      case TokenDiv => { e = Divide(e, eNew) }  
    }  
  }  
  e }
```

foo + 42 / bar + arg



Not ambiguous, but gives wrong tree.

Making Grammars Unambiguous

- some useful recipes -

Ensure that there is always only one parse tree

Construct the correct abstract syntax tree

Goal: Build Expression Trees

abstract class Expr

case class Variable(id : Identifier) **extends** Expr

case class Minus(e1 : Expr, e2 : Expr) **extends** Expr

case class Exp(e1 : Expr, e2 : Expr) **extends** Expr

$e_1 - e_2 - e_3$

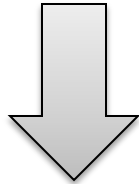
different parse trees give ASTs:

Minus(e1, Minus(e2,e3)) $e_1 - (e_2 - e_3)$

Minus(Minus(e1,e2),e3) $(e_1 - e_2) - e_3$

1) Layer the grammar by priorities

$\text{expr} ::= \text{ident} \mid \text{expr} - \text{expr} \mid \text{expr} \wedge \text{expr} \mid (\text{expr})$



$\text{expr} ::= \text{term} (- \text{term})^*$
 $\text{term} ::= \text{factor} (\wedge \text{factor})^*$
 $\text{factor} ::= \text{id} \mid (\text{expr})$

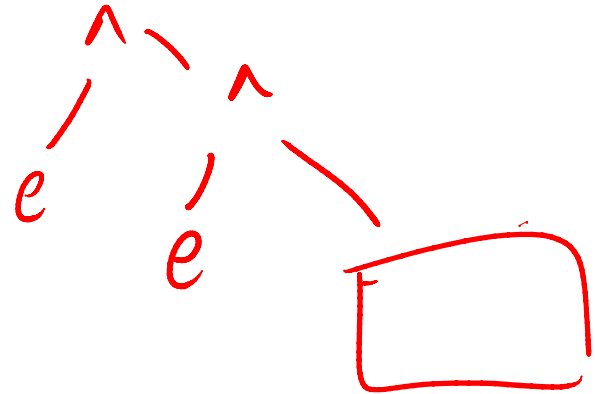
lower priority binds weaker,
so it goes outside

3) Building trees: right-associative "^"

RIGHT-associative operator – using recursion
(or also loop and then reverse a list)

$x \wedge y \wedge z \rightarrow x \wedge (y \wedge z)$
`Exp(Var("x"), Exp(Var("y"), Var("z"))))`

```
def expr : Expr = {  
  val e = factor  
  if (lexer.token == ExpToken) {  
    lexer.next  
    Exp(e, expr)  
  } else e  
}
```



Manual Construction of Parsers

- Typically one applies previous transformations to get a nice grammar
- Then we write recursive descent parser as set of mutually recursive procedures that check if input is well formed
- Then enhance such procedures to construct trees, paying attention to the associativity and priority of operators

Exercise: Unary Minus

1) Show that the grammar

$$A ::= - A$$
$$A ::= A - \text{id}$$
$$A ::= \text{id}$$

is ambiguous by finding a string that has two different syntax trees.

2) Make two different unambiguous grammars for the same language:

a) One where prefix minus binds stronger than infix minus.

b) One where infix minus binds stronger than prefix minus.

3) Show the syntax trees using the new grammars for the string you used to prove the original grammar ambiguous.

Exercise: Dangling Else I

The dangling-else problem happens when the conditional statements are parsed using the following grammar.

$S ::= \text{id} := E$

$S ::= \text{if } E \text{ then } S$

$S ::= \text{if } E \text{ then } S \text{ else } S$

Find an unambiguous grammar that accepts the same conditional statements and matches the else statement with the nearest unmatched if.

Discussion of Dangling Else

```
if (x > 0) then
  if (y > 0) then
    z = x + y
else x = - x
```

- This is a real problem languages like C, Java
 - resolved by saying **else** binds to innermost **if**
- Can we design grammar that allows all programs as before, but only allows parse trees where else binds to innermost if?

Solution?

Exercise: Dangling Else II

Suppose that in addition assignments and if statements we have statement sequencing:

$S ::= S ; S$

$S ::= \text{id} := E$

$S ::= \text{if } E \text{ then } S$

$S ::= \text{if } E \text{ then } S \text{ else } S$

Find an unambiguous grammar that accepts the same conditional statements, matches the else statement with the nearest unmatched if, and behaves similarly to statements in Java.

Sources of Ambiguity in this Example

- Ambiguity arises in this grammar here due to:
 - dangling **else**
 - binary rule for sequence (;) as for parentheses
 - priority between if-then-else and semicolon (;)

```
if (x > 0)
```

```
    if (y > 0)
```

```
        z = x + y;
```

```
        u = z + 1    // last assignment is not inside if
```

Wrong parse tree -> wrong generated code

How we Solved It

We identified a wrong tree and tried to refine the grammar to prevent it, by making a copy of the rules. Also, we changed some rules to disallow sequences inside if-then-else and make sequence rule non-ambiguous. The end result is something like this:

```
S ::= ε | A S // a way to write S ::= A*
A ::= id := E
A ::= if E then A
A ::= if E then A' else A
A' ::= id := E
A' ::= if E then A' else A'
```

At some point we had a useless rule, so we deleted it.

We also looked at what a practical grammar would have to allow sequences inside if-then-else. It would add a case for blocks, like this:

```
A ::= { S }
A' ::= { S }
```

We could factor out some common definitions (e.g. define A in terms of A'), but that is not important for this problem.

Solution?

Formalizing and Automating Recursive Descent: LL(1) Parsers

Recursive Descent - LL(1)

- See wiki for
 - computing first, nullable, follow for non-terminals of the grammar
 - construction of parse table using this information
 - LL(1) as an interpreter for the parse table

Grammar vs Recursive Descent Parser

```
expr ::= term termList
termList ::= + term termList
           | - term termList
           | ε
term ::= factor factorList
factorList ::= * factor factorList
            | / factor factorList
            | ε
factor ::= name | ( expr )
name ::= ident
```

```
def expr = { term; termList }
def termList =
  if (token==PLUS) {
    skip(PLUS); term; termList
  } else if (token==MINUS)
    skip(MINUS); term; termList
  }
def term = { factor; factorList }
...
def factor =
  if (token==IDENT) name
  else if (token==OPAR) {
    skip(OPAR); expr; skip(CPAR)
  } else error("expected ident or ")
```

Rough General Idea

$A ::= B_1 \dots B_p$
| $C_1 \dots C_q$
| $D_1 \dots D_r$



```
def A =  
  if (token ∈ T1) {  
    B1 ... Bp  
  } else if (token ∈ T2) {  
    C1 ... Cq  
  } else if (token ∈ T3) {  
    D1 ... Dr  
  } else error("expected T1,T2,T3")
```

where:

$T1 = \mathbf{first}(B_1 \dots B_p)$

$T2 = \mathbf{first}(C_1 \dots C_q)$

$T3 = \mathbf{first}(D_1 \dots D_r)$

$\mathbf{first}(B_1 \dots B_p) = \{a \in \Sigma \mid B_1 \dots B_p \Rightarrow \dots \Rightarrow aw\}$

$T1, T2, T3$ should be **disjoint** sets of tokens.

Computing **first** in the example

```
expr ::= term termList
termList ::= + term termList
           | - term termList
           | ε
term ::= factor factorList
factorList ::= * factor factorList
            | / factor factorList
            | ε
factor ::= name | ( expr )
name ::= ident
```

$\text{first}(\text{name}) = \{\mathbf{ident}\}$

$\text{first}(\text{(expr)}) = \{ (\ }$

$\text{first}(\text{factor}) = \text{first}(\text{name})$

$\cup \text{first}(\text{(expr)})$

$= \{\mathbf{ident}\} \cup \{ (\ }$

$= \{\mathbf{ident}, (\ }$

$\text{first}(* \text{ factor factorList}) = \{ * \}$

$\text{first}(/ \text{ factor factorList}) = \{ / \}$

$\text{first}(\text{factorList}) = \{ *, / \}$

$\text{first}(\text{term}) = \text{first}(\text{factor}) = \{\mathbf{ident}, (\ }$

$\text{first}(\text{termList}) = \{ +, - \}$

$\text{first}(\text{expr}) = \text{first}(\text{term}) = \{\mathbf{ident}, (\ }$

Algorithm for **first**

Given an arbitrary context-free grammar with a set of rules of the form $X ::= Y_1 \dots Y_n$ compute first for each right-hand side and for each symbol.

How to handle

- alternatives for one non-terminal
- sequences of symbols
- nullable non-terminals
- recursion

Rules with Multiple Alternatives

$$A ::= B_1 \dots B_p \\ | C_1 \dots C_q \\ | D_1 \dots D_r$$

$$\text{first}(A) = \text{first}(B_1 \dots B_p) \\ \cup \text{first}(C_1 \dots C_q) \\ \cup \text{first}(D_1 \dots D_r)$$

Sequences

$$\text{first}(B_1 \dots B_p) = \text{first}(B_1)$$

if not nullable(B_1)

$$\text{first}(B_1 \dots B_p) = \text{first}(B_1) \cup \dots \cup \text{first}(B_k)$$

if nullable(B_1), ..., nullable(B_{k-1}) and
not nullable(B_k) or $k=p$

Abstracting into Constraints

recursive grammar: constraints over finite sets: expr' is $\text{first}(\text{expr})$

```
expr ::= term termList
termList ::= + term termList
           | - term termList
           | ε
term ::= factor factorList
factorList ::= * factor factorList
            | / factor factorList
            | ε
factor ::= name | ( expr )
name ::= ident
```

nullable: termList, factorList

```
expr' = term'
termList' = {+}
           ∪ {-}

term' = factor'
factorList' = {*}
            ∪ {/}

factor' = name' ∪ { ( }
name' = { ident }
```

For this nice grammar, there is no recursion in constraints. Solve by substitution.

Example to Generate Constraints

$$\begin{aligned} S &::= X \mid Y \\ X &::= \mathbf{b} \mid S Y \\ Y &::= Z X \mathbf{b} \mid Y \mathbf{b} \\ Z &::= \varepsilon \mid \mathbf{a} \end{aligned}$$

$$\begin{aligned} S' &= X' \cup Y' \\ X' &= \end{aligned}$$

terminals: \mathbf{a}, \mathbf{b}

non-terminals: S, X, Y, Z

reachable (from S):

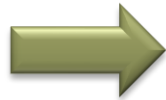
productive:

nullable:

First sets of terminals:

$$S', X', Y', Z' \subseteq \{\mathbf{a}, \mathbf{b}\}$$

Example to Generate Constraints

$$\begin{aligned} S &::= X \mid Y \\ X &::= \mathbf{b} \mid S Y \\ Y &::= Z X \mathbf{b} \mid Y \mathbf{b} \\ Z &::= \varepsilon \mid \mathbf{a} \end{aligned}$$

$$\begin{aligned} S' &= X' \cup Y' \\ X' &= \{\mathbf{b}\} \cup S' \\ Y' &= Z' \cup X' \cup Y' \\ Z' &= \{\mathbf{a}\} \end{aligned}$$

terminals: **a, b**

non-terminals: **S, X, Y, Z**

reachable (from S): **S, X, Y, Z**

productive: **X, Z, S, Y**

nullable: **Z**

These constraints are recursive.
How to solve them?

$$S', X', Y', Z' \subseteq \{\mathbf{a}, \mathbf{b}\}$$

How many candidate solutions

- in this case?
- for k tokens, n nonterminals?

Iterative Solution of **first** Constraints

	S'	X'	Y'	Z'
1.	$\{\}$	$\{\}$	$\{\}$	$\{\}$
2.	$\{\}$	$\{b\}$	$\{b\}$	$\{a\}$
3.	$\{b\}$	$\{b\}$	$\{a,b\}$	$\{a\}$
4.	$\{a,b\}$	$\{a,b\}$	$\{a,b\}$	$\{a\}$
5.	$\{a,b\}$	$\{a,b\}$	$\{a,b\}$	$\{a\}$

$$S' = X' \cup Y'$$

$$X' = \{b\} \cup S'$$

$$Y' = Z' \cup X' \cup Y'$$

$$Z' = \{a\}$$

- Start from all sets empty.
- Evaluate right-hand side and assign it to left-hand side.
- Repeat until it stabilizes.

Sets grow in each step

- initially they are empty, so they can only grow
- if sets grow, the RHS grows (U is monotonic), and so does LHS
- they cannot grow forever: in the worst case contain all tokens

Constraints for Computing Nullable

- Non-terminal is nullable if it can derive ε

$S ::= X \mid Y$
 $X ::= \mathbf{b} \mid S Y$
 $Y ::= Z X \mathbf{b} \mid Y \mathbf{b}$
 $Z ::= \varepsilon \mid \mathbf{a}$



$S' = X' \mid Y'$
 $X' = 0 \mid (S' \& Y')$
 $Y' = (Z' \& X' \& 0) \mid (Y' \& 0)$
 $Z' = 1 \mid 0$

$S', X', Y', Z' \in \{0,1\}$

0 - not nullable

1 - nullable

| - disjunction

& - conjunction

	S'	X'	Y'	Z'
1.	0	0	0	0
2.	0	0	0	1
3.	0	0	0	1

again monotonically growing

Computing first and nullable

- Given any grammar we can compute
 - for each non-terminal X whether $\text{nullable}(X)$
 - using this, the set $\text{first}(X)$ for each non-terminal X
- General approach:
 - generate constraints over finite domains, following the structure of each rule
 - solve the constraints iteratively
 - start from least elements
 - keep evaluating RHS and re-assigning the value to LHS
 - stop when there is no more change

Rough General Idea

$A ::= B_1 \dots B_p$
| $C_1 \dots C_q$
| $D_1 \dots D_r$



```
def A =  
  if (token ∈ T1) {  
    B1 ... Bp  
  } else if (token ∈ T2) {  
    C1 ... Cq  
  } else if (token ∈ T3) {  
    D1 ... Dr  
  } else error("expected T1,T2,T3")
```

where:

$T1 = \text{first}(B_1 \dots B_p)$

$T2 = \text{first}(C_1 \dots C_q)$

$T3 = \text{first}(D_1 \dots D_r)$

$T1, T2, T3$ should be **disjoint** sets of tokens.

Exercise 1

$A ::= B \text{ EOF}$

$B ::= \varepsilon \mid B B \mid (B)$

- Tokens: **EOF**, (,)
- Generate constraints and compute nullable and first for this grammar.
- Check whether first sets for different alternatives are disjoint.

Exercise 2

$S ::= B \text{ EOF}$

$B ::= \varepsilon \mid (B) B$

- Tokens: **EOF**, (,)
- Generate constraints and compute nullable and first for this grammar.
- Check whether first sets for different alternatives are disjoint.

Important Exercise 3

Compute nullable, first for this grammar:

$\text{stmtList} ::= \varepsilon \mid \text{stmt stmtList}$

$\text{stmt} ::= \text{assign} \mid \text{block}$

$\text{assign} ::= \text{ID} = \text{ID} ;$

$\text{block} ::= \text{beginof ID stmtList ID ends}$

Describe a parser for this grammar and explain how it behaves on this input:

beginof myPrettyCode

x = u;

y = v;

myPrettyCode **ends**

Problem Identified

$\text{stmtList} ::= \varepsilon \mid \text{stmt stmtList}$

$\text{stmt} ::= \text{assign} \mid \text{block}$

$\text{assign} ::= \mathbf{ID} = \mathbf{ID} ;$

$\text{block} ::= \mathbf{beginof ID stmtList ID ends}$

Problem parsing stmtList :

- \mathbf{ID} could start alternative stmt stmtList
- \mathbf{ID} could **follow** stmt , so we may wish to parse ε that is, do nothing and return
- For nullable non-terminals, we must also compute what follows them

General Idea when nullable(A)

$A ::= B_1 \dots B_p$
| $C_1 \dots C_q$
| $D_1 \dots D_r$



```
def A =  
  if (token ∈ T1) {  
    B1 ... Bp  
  } else if (token ∈ (T2 U TF)) {  
    C1 ... Cq  
  } else if (token ∈ T3) {  
    D1 ... Dr  
  } // no else error, just return
```

where:

$T_1 = \mathbf{first}(B_1 \dots B_p)$

$T_2 = \mathbf{first}(C_1 \dots C_q)$

$T_3 = \mathbf{first}(D_1 \dots D_r)$

$T_F = \mathbf{follow}(A)$

Only one of the alternatives can be nullable (here: 2nd)
 T_1, T_2, T_3, T_F should be pairwise **disjoint** sets of tokens.

LL(1) Grammar - good for building recursive descent parsers

- Grammar is LL(1) if for each nonterminal X
 - first sets of different alternatives of X are disjoint
 - if nullable(X), first(X) must be disjoint from follow(X)
- For each LL(1) grammar we can build recursive-descent parser
- Each LL(1) grammar is unambiguous
- If a grammar is not LL(1), we can sometimes transform it into equivalent LL(1) grammar

Computing if a token can follow

first($B_1 \dots B_p$) = $\{a \in \Sigma \mid B_1 \dots B_p \Rightarrow \dots \Rightarrow aw\}$

follow(X) = $\{a \in \Sigma \mid S \Rightarrow \dots \Rightarrow \dots Xa \dots\}$

There exists a derivation from the start symbol that produces a sequence of terminals and nonterminals of the form $\dots Xa \dots$
(the token a follows the non-terminal X)

Rule for Computing Follow

Given $X ::= YZ$ (for reachable X)

then $\mathbf{first}(Z) \subseteq \mathbf{follow}(Y)$

and $\mathbf{follow}(X) \subseteq \mathbf{follow}(Z)$

now take care of nullable ones as well:

For each rule $X ::= Y_1 \dots Y_p \dots Y_q \dots Y_r$

$\mathbf{follow}(Y_p)$ should contain:

- $\mathbf{first}(Y_{p+1} Y_{p+2} \dots Y_r)$
- also $\mathbf{follow}(X)$ if $\mathbf{nullable}(Y_{p+1} Y_{p+2} Y_r)$

Compute nullable, first, follow

stmtList ::= ϵ | stmt stmtList

stmt ::= assign | block

assign ::= **ID = ID ;**

block ::= **beginof ID stmtList ID ends**

Is this grammar LL(1)?

Conclusion of the Solution

The grammar is not LL(1) because we have

- nullable(stmtList)
- $\text{first}(\text{stmt}) \cap \text{follow}(\text{stmtList}) = \{\mathbf{ID}\}$
- If a recursive-descent parser sees **ID**, it does not know if it should
 - finish parsing stmtList or
 - parse another stmt

Table for LL(1) Parser: Example

$S ::= B \text{ EOF}$
(1)

$B ::= \varepsilon \mid B (B)$
(1) (2)

nullable: B

$\text{first}(S) = \{ (\}$

$\text{follow}(S) = \{ \}$

$\text{first}(B) = \{ (\}$

$\text{follow}(B) = \{), (, \text{EOF} \}$

empty entry:
when parsing S,
if we see),
report error

Parsing table:

	EOF	()
S	{1}	{1}	{ }
B	{1}	{1,2}	{1}

**parse conflict - choice ambiguity:
grammar not LL(1)**

1 is in entry because (is in follow(B)

2 is in entry because (is in first(B(B))

Table for LL(1) Parsing

Tells which alternative to take, given current token:

choice : Nonterminal x Token \rightarrow Set[Int]

$$\begin{array}{l} A ::= (1) B_1 \dots B_p \\ \quad | (2) C_1 \dots C_q \\ \quad | (3) D_1 \dots D_r \end{array}$$

if $t \in \text{first}(C_1 \dots C_q)$ add 2
to choice(A,t)
if $t \in \text{follow}(A)$ add K to choice(A,t)
where K is nullable alternative

For example, when parsing A and seeing token t

choice(A,t) = {2} means: parse alternative 2 ($C_1 \dots C_q$)

choice(A,t) = {1} means: parse alternative 3 ($D_1 \dots D_r$)

choice(A,t) = {} means: report syntax error

choice(A,t) = {2,3} : not LL(1) grammar

Transform Grammar for LL(1)

$S ::= B \text{ EOF}$

$B ::= \varepsilon \mid B (B)$
(1) (2)

Transform the grammar so that parsing table has no conflicts.

$S ::= B \text{ EOF}$

$B ::= \varepsilon \mid (B) B$
(1) (2)

Left recursion is bad for LL(1)

Old parsing table:

	EOF	()
S	{1}	{1}	{}
B	{1}	{1,2}	{1}

**conflict - choice ambiguity:
grammar not LL(1)**

- 1 is in entry because (is in follow(B)
- 2 is in entry because (is in first(B(B))

	EOF	()
S			
B			

choice(A,t)

Parse Table is Code for Generic Parser

```
var stack : Stack[GrammarSymbol] // terminal or non-terminal
stack.push(EOF);
stack.push(StartNonterminal);
var lex = new Lexer(inputFile)
while (true) {
  X = stack.pop
  t = lex.curent
  if (isTerminal(X))
    if (t==X) if (X==EOF) return success
    else lex.next // eat token t
  else parseError("Expected " + X)
else { // non-terminal
  cs = choice(X)(t) // look up parsing table
  cs match { // result is a set
  case {i} => { // exactly one choice
    rhs = p(X,i) // choose correct right-hand side
    stack.pushRev(rhs) } // pushes symbols in rhs so leftmost becomes top of stack
  case {} => parseError("Parser expected an element of " + unionOfAll(choice(X)))
  case _ => crash("parse table with conflicts - grammar was not LL(1)")
  }
}
```


What if we cannot transform the grammar into LL(1)?

1) Redesign your language

2) Use a more powerful parsing technique