

University of Waterloo

# Parsers

Jon Eyolfson

June 3, 2013

# Goals

Formulate a grammar for integer arithmetic

Create an AST and interpreter for our grammar

# Our Grammar

Here is our attempt at the grammar without multiplication and division for now:

$$T = \{+, -, \text{INT}\}$$

$$N = \{e\}$$

$$S = e$$

$$P = \begin{cases} e \rightarrow e + e \\ e \rightarrow e - e \\ e \rightarrow \text{INT} \end{cases}$$

**Start AST**

# Abstract Syntax Trees

An AST is a tree representation of the input language

This tree should only represents the important details of the input

# Unimportant Details in Our Grammar

Our parse trees have the following unimportant details:

- Rule names
- Brackets (or grouping)

# Our AST Nodes

**BinaryExpr** need to know what operator to use, the left and right side

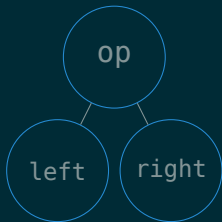
**IntegerExpr** to store the value of an integer

Both of which are **Expr** (which is an abstract base class)

# Our AST Node Representations

We will represent our AST nodes as follows:

## BinaryExpr



op is either +, -, \* or /  
left and right are **Exprs**

## IntegerExpr

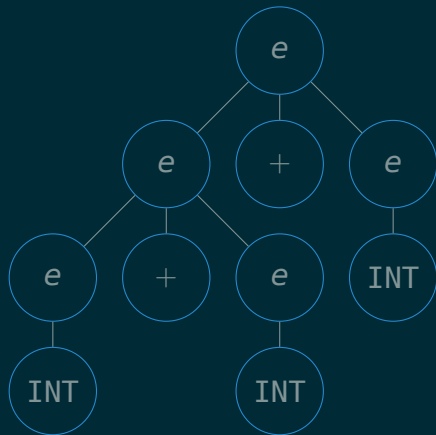


value is an int (in Java we can  
pass the regular expression  
matched by the INT token to  
Integer.parseInt)



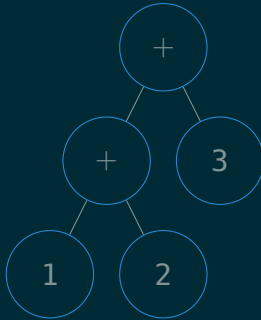
## AST Question (1)

If our input is "1 + 2 + 3" and our parse tree is the following:



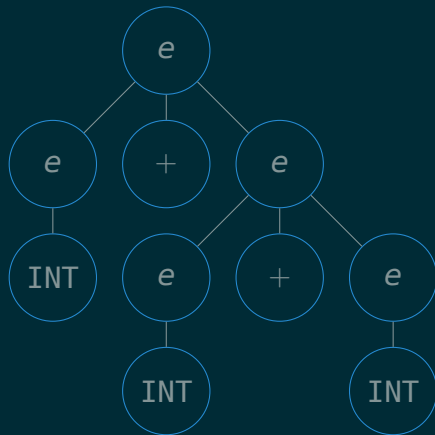
What is our AST?

## AST Answer (1)



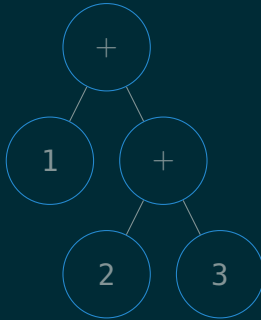
## AST Question (2)

Now consider the same input and the following parse tree:



What is our AST now?

## AST Answer (2)



**End AST**

**Start Interpreter**

# Our Own Interpreter

An interpreter executes the AST directly

Our result should be equivalent to following all the mathematical rules for the input

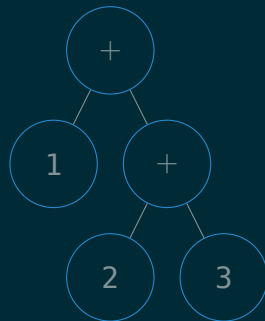
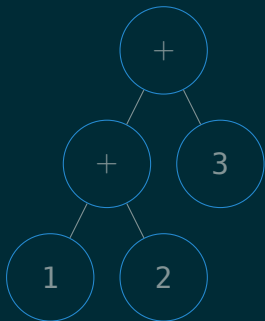
The rules in this case is the DMAS part of BEDMAS

# Our Interpreter Implementation

```
public int eval(BinaryExpr e) {  
    switch (e.op) {  
        case '*':  
            return e.left.eval() * e.right.eval();  
        case '/':  
            return e.left.eval() / e.right.eval();  
        case '+':  
            return e.left.eval() + e.right.eval();  
        case '-':  
            return e.left.eval() - e.right.eval();  
    }  
}  
  
public int eval(IntegerExpr i) {  
    return i.value;  
}
```

# AST Evaluation (1)

We call eval with the root node, what is the result?



The result in both cases is 6



# Ambiguity

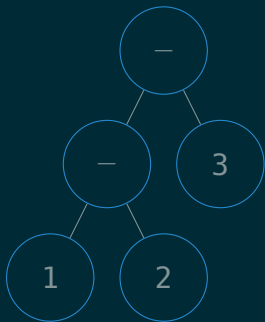
With ambiguity either parse tree/AST is correct

Therefore, the implementer can decide which one they want to use

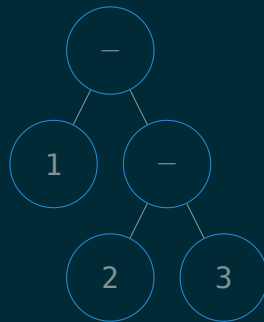
What happens if we use  $-$  instead of  $+$ ?

## AST Evaluation (2)

Again, calling `eval` with the root node, what is the result?



The result is -4



The result is 2

Now, depending on the implementation you get different results

**End Interpreter**

**Start Unambiguity**

# Removing Ambiguity

We have to enforce:

- Precedence
- Associativity

# Precedence Rules

- ① Each level of precedence should be it's own rule
- ② The lowest level of precedence should be the top level rule (highest level of precedence is the last rule)
- ③ Each rule should use the next highest level of precedence
- ④ The highest level of precedence includes the operands

## Our Grammar (Attempt 3)

Following the rules, our grammar is now:

$$T = \{+, -, \text{INT}\}$$

$$N = \{e, p\}$$

$$S = e$$

$$P = \begin{cases} e \rightarrow e + e \\ e \rightarrow e - e \\ e \rightarrow p \\ p \rightarrow \text{INT} \end{cases}$$

But the grammar is still ambiguous

# Associativity Rules

Consider  $e$  as our current level of precedence and  $p$  as the next highest

$e \rightarrow e \circ p$ , would make it left associative

$e \rightarrow p \circ e$ , would make it right associative

## Our Grammar (Attempt 4)

We expect left associativity, our grammar is now:

$$T = \{+, -, \text{INT}\}$$

$$N = \{e, p\}$$

$$S = e$$

$$P = \begin{cases} e \rightarrow e + p \\ e \rightarrow e - p \\ e \rightarrow p \\ p \rightarrow \text{INT} \end{cases}$$



# Parsing

Let's use our previous grammar attempt

Parse the following input: "1 - 2 - 3"

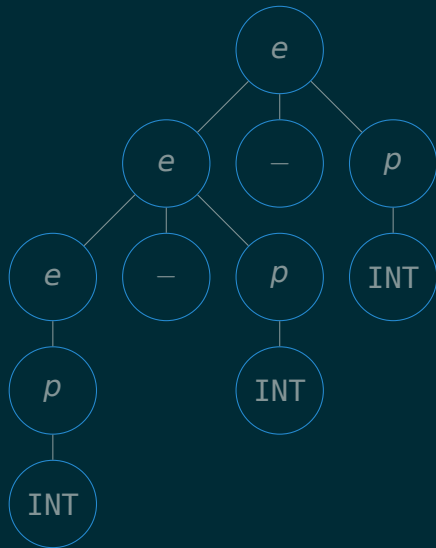
We have the following terminals:  $\text{INT} - \text{INT} - \text{INT}$

## Previous Derivation

$$\begin{aligned} e &\Rightarrow e - p \\ &\Rightarrow e - p - p \\ &\Rightarrow p - p - p \\ &\Rightarrow \text{INT} - p - p \\ &\Rightarrow \text{INT} - \text{INT} - p \\ &\Rightarrow \text{INT} - \text{INT} - \text{INT} \end{aligned}$$

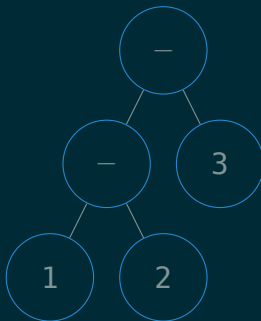
This is the only derivation

## Previous Parse Tree



What about the corresponding AST?

## Previous AST



Our interpreter would return -4 as expected

# Our Grammar (Attempt 5)

Let's add multiplication and division:

$$T = \{+, -, *, /, \text{INT}\}$$

$$N = \{e, p\}$$

$$S = e$$

$$P = \left\{ \begin{array}{l} e \rightarrow e + p \\ e \rightarrow e - p \\ e \rightarrow e * p \\ e \rightarrow e / p \\ e \rightarrow p \\ p \rightarrow \text{INT} \end{array} \right.$$

# Parsing

Let's use our previous grammar attempt

Parse the following input: "1 - 2 \* 3"

We have the following terminals:  $\text{INT} \mid \text{INT} * \text{INT}$

# Previous Derivation

$$e \Rightarrow e * p$$

$$\Rightarrow e - p * p$$

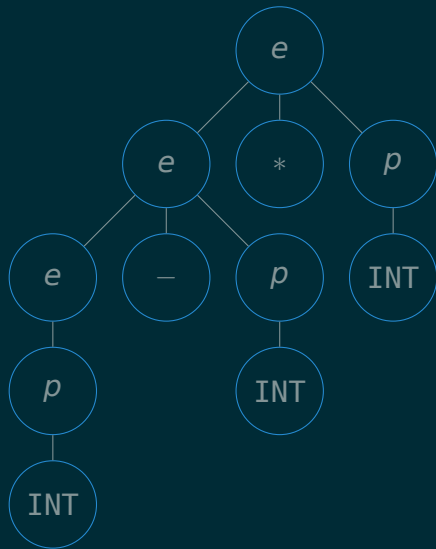
$$\Rightarrow p - p * p$$

$$\Rightarrow \text{INT} - p * p$$

$$\Rightarrow \text{INT} - \text{INT} * p$$

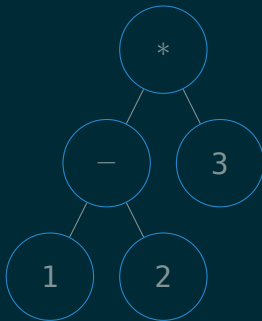
$$\Rightarrow \text{INT} - \text{INT} * \text{INT}$$

## Previous Parse Tree





## Previous AST



Our interpreter would return  $-3$  which is wrong

How would we fix this?

# Precedence Rules

- ① Each level of precedence should be it's own rule
- ② The lowest level of precedence should be the top level rule (highest level of precedence is the last rule)
- ③ Each rule should use the next highest level of precedence
- ④ The highest level of precedence includes the operands

## Our Grammar (Attempt 6)

We have to add another level of precedence ( $t$ ) for multiplication and division:

$$T = \{+, -, *, /, \text{INT}\}$$

$$N = \{e, t, p\}$$

$$S = e$$

$$P = \left\{ \begin{array}{l} e \rightarrow e + t \\ e \rightarrow e - t \\ e \rightarrow t \\ t \rightarrow t * p \\ t \rightarrow t / p \\ t \rightarrow p \\ p \rightarrow \text{INT} \end{array} \right.$$

# Parsing

Let's use our previous grammar attempt

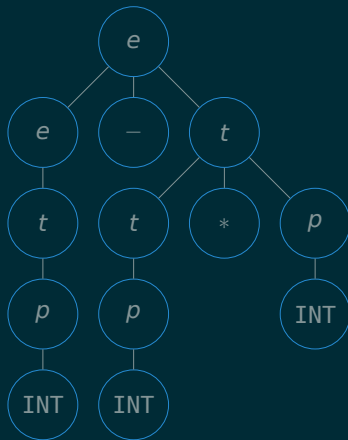
Parse the following input: "1 - 2 \* 3"

We have the following terminals:  $\text{INT} \mid \text{INT} * \text{INT}$

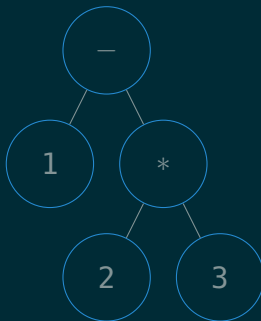
## Previous Derivation

$$\begin{aligned} e &\Rightarrow e - t \\ &\Rightarrow t - t \\ &\Rightarrow p - t \\ &\Rightarrow \text{INT} - t \\ &\Rightarrow \text{INT} - t * p \\ &\Rightarrow \text{INT} - p * p \\ &\Rightarrow \text{INT} - \text{INT} * p \\ &\Rightarrow \text{INT} - \text{INT} * \text{INT} \end{aligned}$$

## Previous Parse Tree



## Previous AST



Our interpreter would return -5 which is correct

# Grouping and Unary Precedence

Let's add grouping `()` and a unary `-`

Grouping is always done at the highest level of precedence and calls the lowest level precedence rule

Unary operators follow the same rules as others, except there is no associativity (so it should call itself)



# Our Completed Grammar

$$T = \{+, -, *, /, (, ), \text{INT}\}$$

$$N = \{e, t, f, p\}$$

$$S = e$$

$$P = \left\{ \begin{array}{l} e \rightarrow e + t \\ e \rightarrow e - t \\ e \rightarrow t \\ t \rightarrow t * f \\ t \rightarrow t / f \\ t \rightarrow f \\ f \rightarrow -f \\ f \rightarrow p \\ p \rightarrow (e) \\ p \rightarrow \text{INT} \end{array} \right.$$

# Grammar Question

Consider the grammar:

$$T = \{\text{IF}, \text{ELSE}\}$$

$$N = \{\text{ifstmt}, \text{expr}, \text{stmt}\}$$

$$S = \text{stmt}$$

$$P = \begin{cases} \text{stmt} \rightarrow \text{ifstmt} \\ \text{ifstmt} \rightarrow \text{IF expr stmt} \\ \text{ifstmt} \rightarrow \text{IF expr stmt ELSE stmt} \end{cases}$$

Is this ambiguous?

# Grammar Answer

Yes, it is ambiguous

Consider the **sentential form**:

IF *expr* IF *expr stmt* ELSE *stmt*

We have two leftmost derivations of this form

# Previous Derivations

$stmt \Rightarrow ifstmt$   
 $\Rightarrow IF\ expr\ stmt$   
 $\Rightarrow IF\ expr\ ifstmt$   
 $\Rightarrow IF\ expr\ IF\ expr\ stmt\ ELSE\ stmt$

$stmt \Rightarrow ifstmt$   
 $\Rightarrow IF\ expr\ stmt\ ELSE\ stmt$   
 $\Rightarrow IF\ expr\ ifstmt\ ELSE\ stmt$   
 $\Rightarrow IF\ expr\ IF\ expr\ stmt\ ELSE\ stmt$

This is called the dangling else problem

# Dangling Else Example

Consider the following Java code:

```
if (x > 0) if (x > 3) x = 3 else x = 0
```

When  $x = 1$  the final value of  $x$  depends on our parse

# Removing the Dangling Else

There are 4 main approaches:

- Add an `endif` statement (like in VHDL)
- Specify the meaning in the language specification and leave it up to the implementations (like Java, which always does inner)
- Introduce indenting into the grammar (like Python)
- Require every `if` to have an `else` (like Haskell)

**End Unambiguity**

**Start Techniques**

# Top-Down Parsing

Consider the grammar:

$$T = \{0, 1\}$$

$$N = \{s, b\}$$

$$S = s$$

$$P = \begin{cases} s \rightarrow 1b0 \\ s \rightarrow 0b1 \\ b \rightarrow 10 \\ b \rightarrow 11 \end{cases}$$

Let's do our leftmost derivation one token at a time

This is called **top-down parsing**



# Top-Down Parse (1)

Current token: 1

$$s \Rightarrow 1b0$$

## Top-Down Parse (2)

Current token: 1

$$s \Rightarrow 1b0$$

$$s \Rightarrow 1100$$

This is a guess since we don't know the next token

## Top-Down Parse (3)

Current token: 1

Our previous guess was wrong, so we backtrack to

$$s \Rightarrow 1b0$$

And guess again

$$s \Rightarrow 1b0$$

$$s \Rightarrow 1110$$

## Top-Down Parse (4)

Current token: 0

$$s \Rightarrow 1b0$$

$$s \Rightarrow 1110$$

This is the last token

Therefore 1110 is valid in our grammar

# Recursive Descent Parsing

An example of a top-down parser is a **recursive descent parser**

Translating the grammar follows these rules:

- Each non-terminal turns into a function, the body turns into the production
- Each non-terminal in the production is a function call to the corresponding non-terminal
- Each terminal in the production is a call to the lexer's consume with the terminal as the argument

You can use `inspect` with a terminal as an argument to check the current token and decide which rule to follow

There could be backtracking as well if you need it

# Implementing Our Grammar

Let's implement our first rule  $e \rightarrow e + t$

One implementation is the following:

```
public void E() {  
    E();  
    consume("+");  
    T();  
}
```

Do you see the major problem with this?

# Removing Left Recursion

Left recursion leads to infinite recursion in our function

We can remove left recursion in grammars of the form:

$$P = \begin{cases} a \rightarrow ax \\ a \rightarrow y \end{cases}$$

where  $a$  is a non-terminal,  $x$  and  $y$  are any sequence of terminals/non-terminals

By replacing it with:

$$P = \begin{cases} a \rightarrow ya' \\ a' \rightarrow xa' \\ a' \rightarrow \varepsilon \end{cases}$$

Note that  $\varepsilon$  represents no symbols

# Removing Left Recursion Question

Remove left recursion from our rule:

$$P = \begin{cases} e \rightarrow e + t \\ e \rightarrow t \end{cases}$$



# Removing Left Recursion Answer

Following the rules we would replace it with:

$$P = \begin{cases} e \rightarrow t e' \\ e' \rightarrow + t e' \\ e' \rightarrow \varepsilon \end{cases}$$

**End Techniques**

**Start EBNF**

# Grammar Equivalence

Let's demonstrate that the grammars are equivalent:

$$P = \begin{cases} a \rightarrow ax \\ a \rightarrow y \end{cases}$$

$$P = \begin{cases} a \rightarrow ya' \\ a' \rightarrow xa' \\ a' \rightarrow \varepsilon \end{cases}$$

What are some example derivations?

# Extended Backus-Naur Form

$y, yx, yxx, yxxx$

If  $x$  and  $y$  were just characters, this would be the same as the regular expression  $yx^*$

Extended Backus-Naur Form (EBNF) is just that, it's BNF with regular expression operations (+, \*, ( ), ?)

# EBNF Example

Assume the grammar:

$$T = \{A, B, C\}$$

$$N = \{a, x\}$$

$$S = a$$

$$P = \left\{ a \rightarrow A x C \right.$$

Let's consider some variations for the rule  $x$  in BNF and EBNF

# EBNF Variation (1)

The following productions are equivalent:

$$P = \begin{cases} a \rightarrow A x C \\ x \rightarrow x B \\ x \rightarrow \varepsilon \end{cases}$$

$$P = \begin{cases} a \rightarrow A x C \\ x \rightarrow B^* \end{cases}$$

$$P = \{ a \rightarrow A B^* C \}$$

## EBNF Variation (2)

The following productions are equivalent:

$$P = \begin{cases} a \rightarrow A x C \\ x \rightarrow B \\ x \rightarrow \varepsilon \end{cases}$$

$$P = \begin{cases} a \rightarrow A x C \\ x \rightarrow B? \end{cases}$$

## EBNF Variation (3)

The following productions are equivalent:

$$P = \begin{cases} a \rightarrow A x C \\ x \rightarrow B y \\ y \rightarrow B y \\ y \rightarrow \varepsilon \end{cases}$$

$$P = \begin{cases} a \rightarrow A x C \\ x \rightarrow B^+ \end{cases}$$



# EBNF Implementation Example

For recursive descent parsers \* translates to a while loop, ? translates to an if statement

$$P = \{x \rightarrow B^*\}$$

```
public void x() {  
    while (inspect("B")) { consume("B"); }  
}
```

$$P = \{x \rightarrow B^?\}$$

```
public void x() {  
    if (inspect("B")) { consume("B"); }  
}
```

# EBNF Question

Consider our rewritten productions:

$$P = \begin{cases} e \rightarrow t e' \\ e' \rightarrow + t e' \\ e' \rightarrow \varepsilon \end{cases}$$

Let's convert this grammar to EBNF

# EBNF Answer

Recognize that  $e'$  can be transformed:

$$P = \{ e \rightarrow t (+ t)^* \}$$

# Converting Our Grammar to EBNF

We also have a  $-$  at that level of precedence:

$$P = \begin{cases} e \rightarrow t (+ t)^* \\ e \rightarrow t (- t)^* \end{cases}$$

We can factor out the common parts of both productions

$$P = \left\{ e \rightarrow t ((+ t) \mid (- t))^* \right.$$

$$P = \left\{ e \rightarrow t ((+|-) t)^* \right.$$

The production is now in a form that's easily translated into a recursive descent parser

# EBNF for Associativity

Let's specify our previous associativity rules in EBNF

Consider  $e$  as the current level of precedence and  $p$  is the next highest level of precedence

$e \rightarrow p (0P p)^*$  is left associative

$e \rightarrow p (0P e)?$  is right associative

**End EBNF**

**Start Language Theory**

# Deterministic Top-Down Parsing

We can have a deterministic recursive descent parser if we have no backtracking

To avoid backtracking, we make only correct guesses

To make only correct guesses, we restrict the grammar so we can always choose the correct next step based on the next token

# LL(1)

*LL* is the class of grammars you've been parsing

- L for left-to-right scan of the input

- L for leftmost derivation

An  $LL(k)$  grammar can be parsed top-down with no backtracking looking only at the next  $k$  symbols in the input

$LL(1)$  can be parsed with a recursive descent parser with a one token (or terminal) look-ahead



# End Language Theory