

# CSE211: Compiler Design

Oct. 13, 2020

- **Topic:** PLY and parsing with derivatives

- **Questions:**

*How would you implement what we've talked about so far?*

*Have you ever implemented a parser?*

```
import ply.lex as lex
import ply.yacc as yacc

# All token names must be declared in this global variable
tokens = ['NUM', 'PLUS', 'MINUS', 'TIMES', 'DIVIDE', 'EX

t_NUM = r'[0-9]+'
t_PLUS = r'\+'
t_MINUS = r'\-'
t_TIMES = r'\*'
t_DIVIDE = r'\/'
t_EXP = r'\^'
t_LPAR = r'\('
t_RPAR = r'\)'

# Ignore spaces and tabs
t_ignore = ' '

# Required in case of a lexing error
def t_error(t):
    print("lexing error!")
    exit(0)

# Build the lexer
lexer = lex.lex()

def p_expression(p):
    'expr : term'
    p[0] = p[1]
```

# Announcements

- End of module 1!
  - Module 2 starts next week
- Homework is released: Due in 2 weeks
  - Questions?
- References for today is an academic paper and code docs

# CSE211: Compiler Design

Oct. 13, 2020

- **Topic:** PLY and parsing with derivatives

- **Questions:**

*How would you implement what we've talked about so far?*

*Have you ever implemented a parser?*

```
import ply.lex as lex
import ply.yacc as yacc

# All token names must be declared in this global variable
tokens = ['NUM', 'PLUS', 'MINUS', 'TIMES', 'DIVIDE', 'EX

t_NUM = r'[0-9]+'
t_PLUS = r'\+'
t_MINUS = r'\-'
t_TIMES = r'\*'
t_DIVIDE = r'\/'
t_EXP = r'\^'
t_LPAR = r'\('
t_RPAR = r'\)'

# Ignore spaces and tabs
t_ignore = ' '

# Required in case of a lexing error
def t_error(t):
    print("lexing error!")
    exit(0)

# Build the lexer
lexer = lex.lex()

def p_expression(p):
    'expr : term'
    p[0] = p[1]
```

# Today

- Building a simple interpreter
- Parsing regular expressions with derivatives
- Homework overview

# Reminder

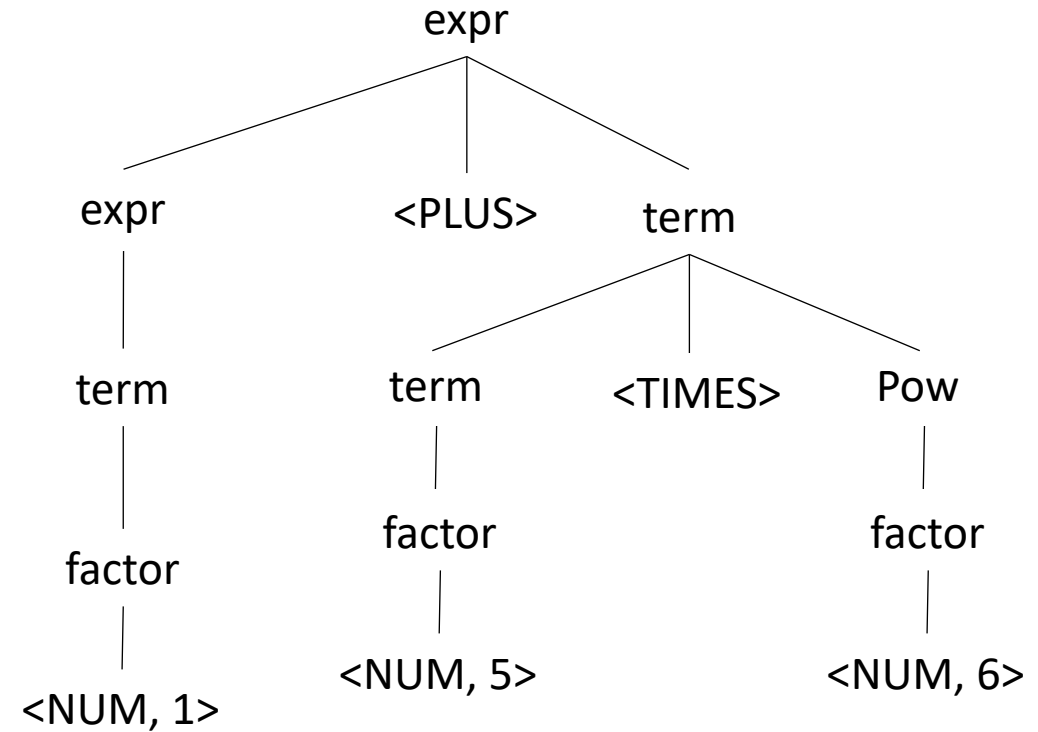
What are the levels for?

Operator	Name	Productions
+,-	Expr	: Expr + Term   Expr - Term   Term
*,/	Term	: Term * Pow : Term / Pow   Pow
^	Pow	: Factor ^ Pow   Factor
()	Factor	: ( Expr )   NUM

# Reminder

input: 1+5\*6

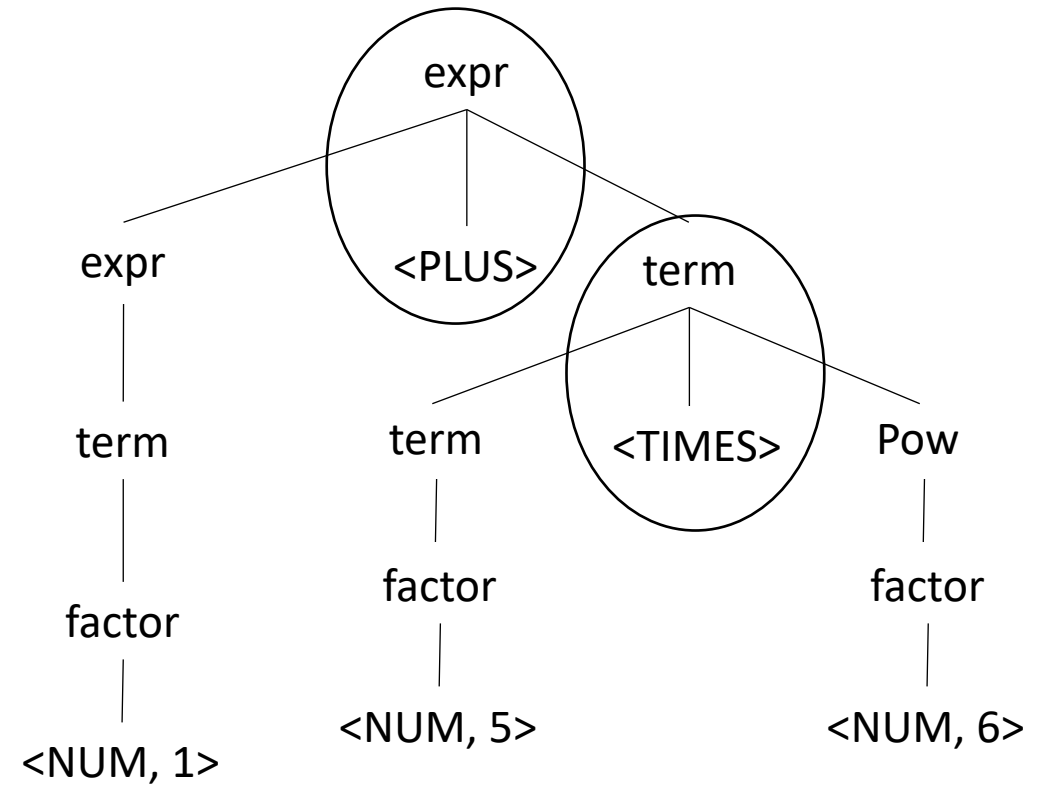
Operator	Name	Productions
+,-	Expr	: Expr + Term   Expr - Term   Term
*,/	Term	: Term * Pow : Term / Pow   Pow
^	Pow	: Factor ^ Pow   Factor
()	Factor	: ( Expr )   NUM



# Reminder

input: 1+5\*6

Operator	Name	Productions
+,-	Expr	: Expr + Term   Expr - Term   Term
*,/	Term	: Term * Pow : Term / Pow   Pow
^	Pow	: Factor ^ Pow   Factor
()	Factor	: ( Expr )   NUM



Higher precedence is lower level

# Reminder

What are the recursion choices?

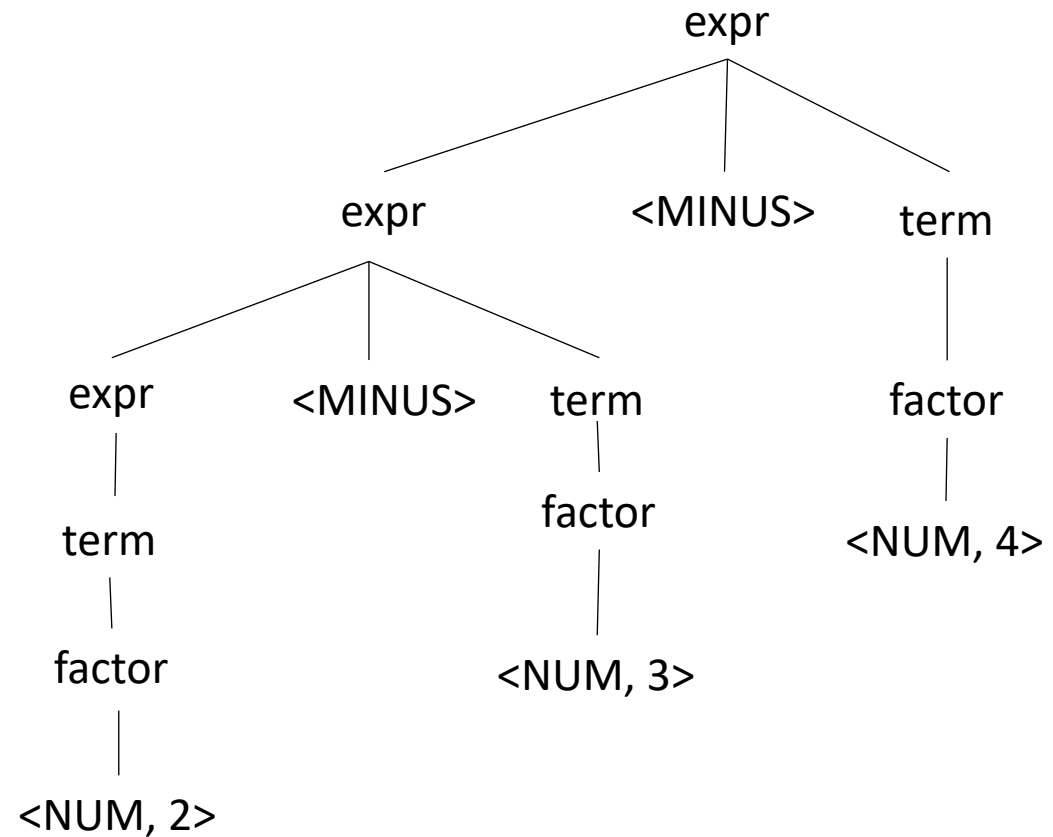
Operator	Name	Productions
+,-	Expr	: Expr + Term   Expr - Term   Term
*,/	Term	: Term * Pow : Term / Pow   Pow
^	Pow	: Factor ^ Pow   Factor
()	Factor	: ( Expr )   NUM



# Reminder

input: 2-3-4

Operator	Name	Productions
+,-	Expr	: Expr + Term   Expr - Term   Term
*,/	Term	: Term * Pow : Term / Pow   Pow
^	Pow	: Factor ^ Pow   Factor
()	Factor	: ( Expr )   NUM



# Production actions

Operator	Name	Productions	Action
+,-	Expr	: Expr + Term   Expr - Term   Term	
*,/	Term	: Term * Pow : Term / Pow   Pow	
^	Pow	: Pow ^ Factor   Factor	
()	Factor	: ( Expr )   NUM	

One action per row.

a struct p with values for each production element, with p[0] as the return value

For example:

“(5)” will eventually execute the production rule:

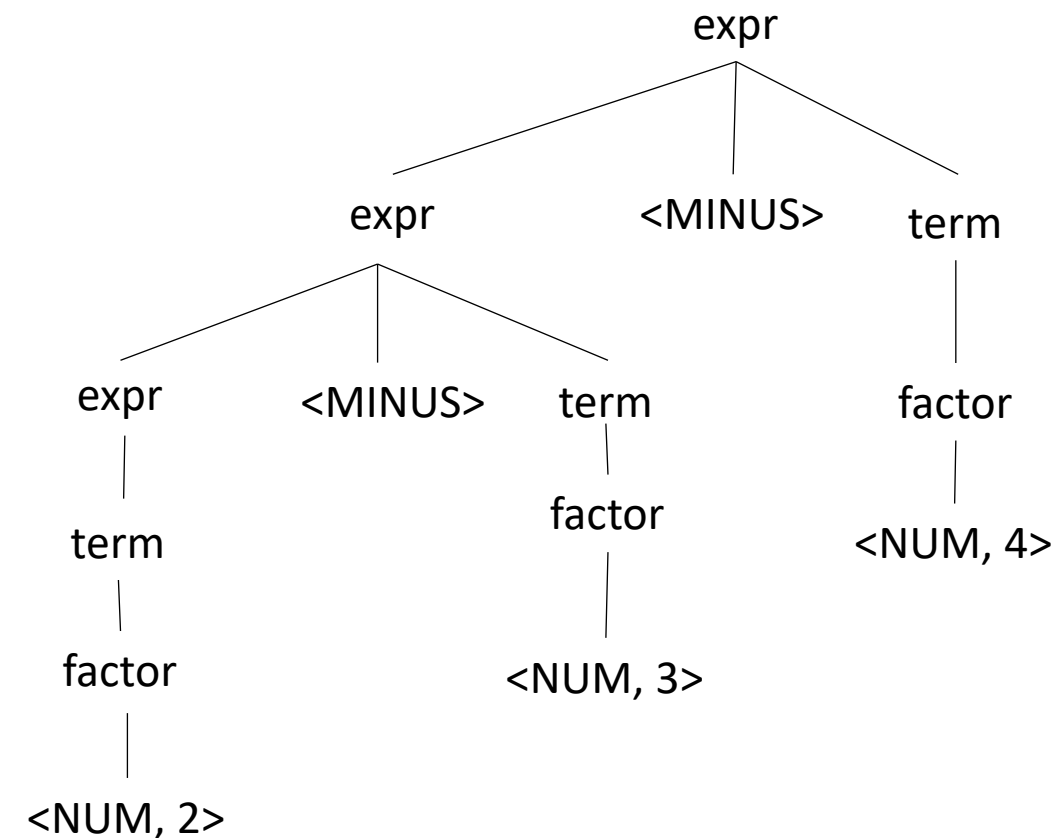
Factor: ( Expr )

at that point:

p = {'(', 5, ')'}

# Production actions

Operator	Name	Productions	Action
+,-	Expr	: Expr + Term   Expr - Term   Term	{p[0] = p[1] + p[3]} {p[0] = p[1] - p[3]} {p[0] = p[1]}
*,/	Term	: Term * Pow : Term / Pow   Pow	{p[0] = p[1] * p[3]} {p[0] = p[1] / p[3]} {p[0] = p[1]}
^	Pow	: Pow ^ Factor   Factor	{p[0] = p[1] ^ p[3]} {p[0] = p[1]}
()	Factor	: ( Expr )   NUM	{p[0] = p[2]} {p[0] = int(p[1])}



# Let's implement it!

- Related to the PLY calc example:
  - <https://github.com/dabeaz/ply/blob/master/example/calc/calc.py>
- You can also reference the PLY documentation:
  - <https://www.dabeaz.com/ply/ply.html>
- I will provide the code examples on Canvas

# Lecture Part 2: Parsing RE's with Derivatives

- A simple, functional style regular expression parser implementation
- Uses an regular expression AST

# Language Derivatives

- A language is a (potentially infinite) set of strings  $\{s_1, s_2, s_3, s_4, \dots\}$
- A language is regular if it can be captured using a regular expression
- Examples of regular languages:
  - $\{a\}, \{+\}, \{+,-,*,\}$
  - $\{1, 1+1, 1+1+1\}$
  - $\{\text{""}\}$ , also called  $\{\varepsilon\}$
  - $\{\}$

***Subtle distinction between  $\{\}$  and  $\{\varepsilon\}$***

# Language Derivatives

- The Derivative of language  $L$  with respect to character  $c$  (noted  $\delta_c(L)$  ) is:

for all  $s$  in  $L$ , if  $s$  begins with  $c$ , then  $s[1:]$  is in  $\delta_c(L)$

- Examples:

# Language Derivatives Examples

- $L = \{“a”\}$
- $\delta_a(L) = \{\varepsilon\}$
- $\delta_b(L) = \{\}$



# Language Derivatives Examples

- $L = \{ "+", "-", "*", "/" \}$

- $\delta_+(L) = \{\varepsilon\}$

- $\delta_\wedge(L) = \{\}$

- $\delta_*(L) = \{\varepsilon\}$

# Language Derivatives Examples

- $L = \{“aaa”, “ab”, “ba”, “bba”\}$

- $\delta_a(L) = \{“aa”, “b”\}$

- $\delta_{aa}(L) = \{“a”\}$

- $\delta_b(L) = \{“a”, “ba”\}$

- $\delta_{ba}(L) = \{\epsilon\}$

# Language Derivatives Examples

- $L = \{“1”, “1+1”, “1+1+1”, “1+1+1+1”, \dots\}$

- $\delta_+(L) = \{\}$

- $\delta_1(L) = \{\varepsilon, “+1”, “+1+1”, “+1+1+1”, \dots\}$

- $\delta_{1+}(L) = \{“1”, “1+1”, “1+1”, \dots\} = L$

# Regular expressions are closed under derivatives

- Given a regular expression  $re$ , any derivative of  $re$  is also a regular expression

- *regular expression =*

|  $\{\}$

|  $\{\epsilon\}$

| "a" (single character)

|  $re_{lhs} \mid re_{rhs}$

|  $re_{lhs} \cdot re_{rhs}$

|  $re_{starred}^*$

# Derivative Base Cases

- $\delta_c(re)$ , where  $re$  is:

- $\{\}$

- $\{\}$

- $\{\varepsilon\}$

- $\{\}$

- "a" (single character)

# Derivative Base Cases

- $\delta_c(re)$ , where  $re$  is:

- $\{\}$

return  $\{\}$

- $\{\varepsilon\}$

return  $\{\}$

- $"a"$  (single character)

*return if  $"a" == c$  then  $\{\varepsilon\}$ , else  $\{\}$*

# Derivative Recursive Cases

- $\delta_c(re)$ , where  $re$  is:

- $re_{lhs} / re_{rhs}$

return  $\delta_c(re_{lhs}) \mid \delta_c(re_{rhs})$

- $re_{starred}^*$

return  $\delta_c(re_{starred}) \cdot re_{starred}^*$

- $re_{lhs} \cdot re_{rhs}$

return  $\delta_c(re_{lhs}) \cdot re_{rhs} \mid$

$if \ \varepsilon \text{ in } re_{lhs} \text{ then } \delta_c(re_{rhs}) \text{ else } \{\}$

# Nullable operator

- $\text{NULL}(re) =$

*if  $\epsilon \in re$  then:  $\{\epsilon\}$   
else:  $\{\}$*



# Nullable operator

- $\text{NULL}(re) =$

*if  $\epsilon \in re$  then:  $\{\epsilon\}$   
else:  $\{\}$*

- $\delta_c(re)$ , where  $re$  is:

- $re_{rhs} \cdot re_{lhs}$

$\delta_c(re_{rhs}) \cdot re_{lhs} \mid$

*if  $\epsilon$  in  $re_{rhs}$  then  $\delta_c(re_{lhs})$  else  $\{\}$*

# Nullable operator

- $\text{NULL}(re) =$

*if  $\epsilon \in re$  then:  $\{\epsilon\}$   
else:  $\{\}$*

- $\delta_c(re)$ , where  $re$  is:

- $re_{rhs} \cdot re_{lhs}$

$\delta_c(re_{rhs}) \cdot re_{lhs} \mid$

$\text{NULL}(re_{rhs}) \cdot \delta_c(re_{lhs})$

# Nullable Base Case

- $NULL(re)$ , where  $re$  is:

- $\{\}$

return  $\{\}$

- $\{\epsilon\}$

return  $\{\epsilon\}$

- "a" (single character)

return  $\{\}$

# Nullable Recursive Cases

- $NULL(re)$ , where  $re$  is:

- $re_{lhs} \mid re_{rhs}$

$$NULL(re_{lhs}) \mid NULL(re_{rhs})$$

- $re_{starred}^*$

$$\{\varepsilon\}$$

- $re_{lhs} \cdot re_{rhs}$

$$NULL(re_{lhs}) \cdot NULL(re_{rhs})$$

# Parsing REs with derivative

given a function  $\delta_c$  to compute the derivative of an RE, the NULL function, an RE  $re$ , and a string  $s = c_1 \cdot c_2 \cdot c_3 \dots$  (concat of characters)

Can we check if  $re$  matches  $s$ ?

# Parsing REs with derivative

given a function  $\delta_c$  to compute the derivative of an RE, the NULL function, an RE  $re$ , and a string  $s = c_1 . c_2 . c_3 \dots$  (concat of characters)

Can we check if  $re$  matches  $s$ ?

$$L(re) = \{.. s ..\}$$

# Parsing REs with derivative

given a function  $\delta_c$  to compute the derivative of an RE, the NULL function, an RE  $re$ , and a string  $s = c_1 \cdot c_2 \cdot c_3 \dots$  (concat of characters)

Can we check if  $re$  matches  $s$ ?

$$L(re) = \{.. s ..\} \left| \begin{array}{l} \delta_{c_1}(re) \\ \\ L(\delta_{c_1}(re)) = \{.. s[1:] ..\} \end{array} \right.$$

# Parsing REs with derivative

given a function  $\delta_c$  to compute the derivative of an RE, the NULL function, an RE  $re$ , and a string  $s = c_1 \cdot c_2 \cdot c_3 \dots$  (concat of characters)

Can we check if  $re$  matches  $s$ ?

$L(re) = \{.. s ..\}$	$\delta_{c_1}(re)$	$\delta_{c_2}(\delta_{c_1}(re)) = \delta_{c_1, c_2}(re)$
	$L(\delta_{c_1}(re)) = \{.. s[1:] ..\}$	$L(\delta_{c_1, c_2}(re)) = \{.. s[2:] ..\}$



# Parsing REs with derivative

given a function  $\delta_c$  to compute the derivative of an RE, the NULL function, an RE  $re$ , and a string  $s = c_1 \cdot c_2 \cdot c_3 \dots$  (concat of characters)

Can we check if  $re$  matches  $s$ ?

	$\delta_{c_1}(re)$	$\delta_{c_2}(\delta_{c_1}(re)) = \delta_{c_1, c_2}(re)$	$\delta_s(re)$
$L(re) = \{.. s ..\}$			
	$L(\delta_{c_1}(re)) = \{.. s[1:] ..\}$	$L(\delta_{c_1, c_2}(re)) = \{.. s[2:] ..\}$	$L(\delta_s(re)) = \{.. \epsilon ..\}$

# Parsing REs with derivative

given a function  $\delta_c$  to compute the derivative of an RE, the NULL function, an RE  $re$ , and a string  $s = c_1 \cdot c_2 \cdot c_3 \dots$  (concat of characters)

Can we check if  $re$  matches  $s$ ?

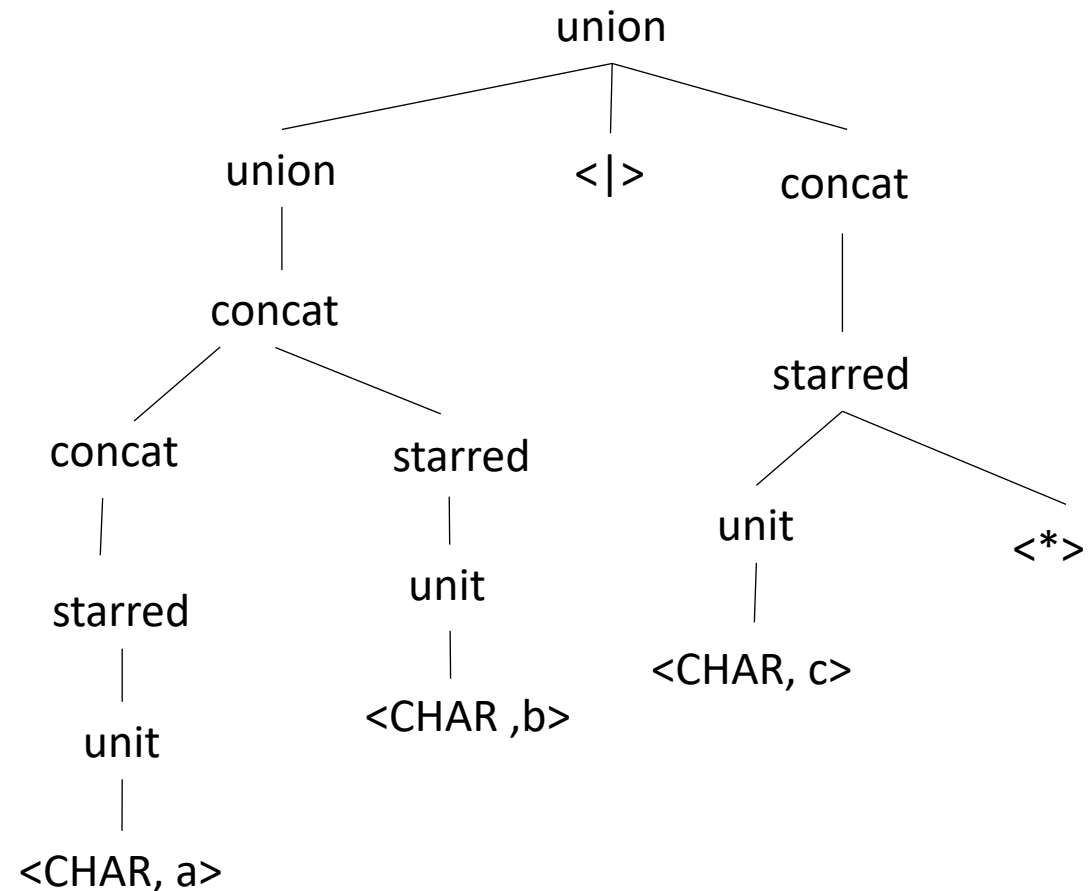
	$\delta_{c_1}(re)$	$\delta_{c_2}(\delta_{c_1}(re)) = \delta_{c_1, c_2}(re)$	$\delta_s(re)$	<div style="border: 1px solid black; padding: 5px; display: inline-block;">If this is true, Then <math>re</math> matches <math>s</math></div>
$L(re) = \{.. s ..\}$	$L(\delta_{c_1}(re)) = \{.. s[1:] ..\}$	$L(\delta_{c_1, c_2}(re)) = \{.. s[2:] ..\}$	$L(\delta_s(re)) = \{.. \epsilon ..\}$	$NULL(\delta_s(re)) == \{\epsilon\}$

# Implementations

Create a parse tree for regular expressions

input: a.b | c\*

Operator	Name	Productions
	union	: union \  concat   concat
.	concat	: concat . starred   starred
*	starred	: starred *   unit
()	unit	: (union)   CHAR

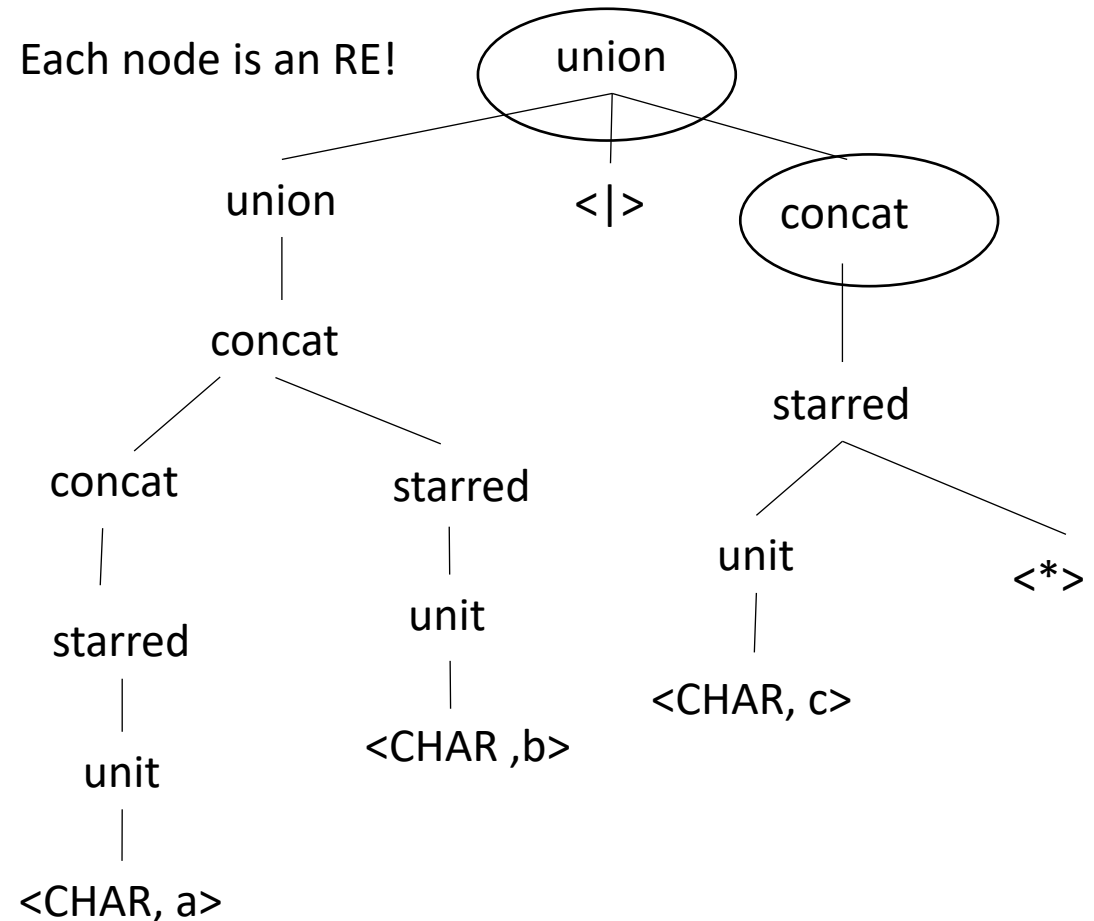


# Implementations

Create a parse tree for regular expressions

input: a.b | c\*

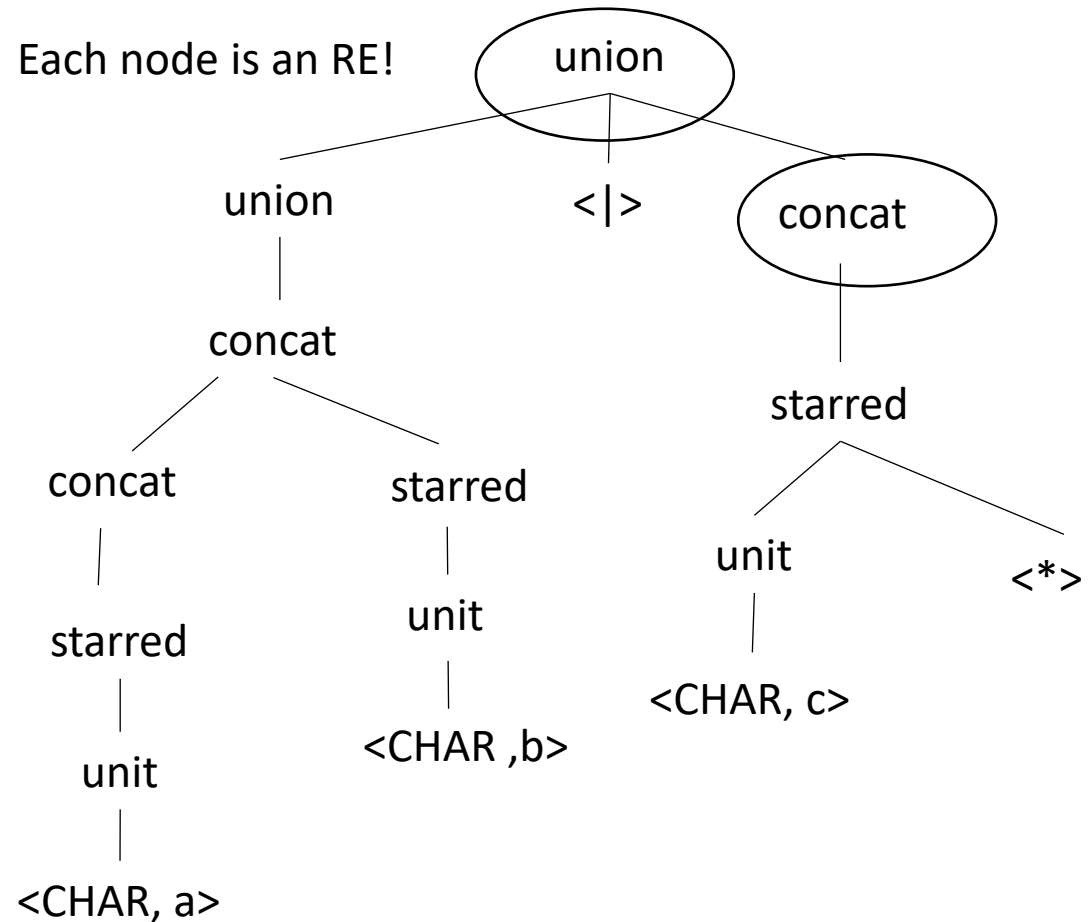
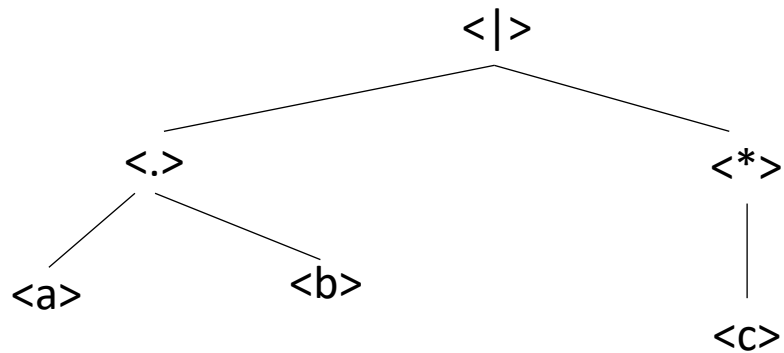
Operator	Name	Productions
	union	: union \  concat   concat
.	concat	: concat . starred   starred
*	starred	: starred *   unit
()	unit	: (union)   CHAR



# Implementations

Create a parse tree for regular expressions

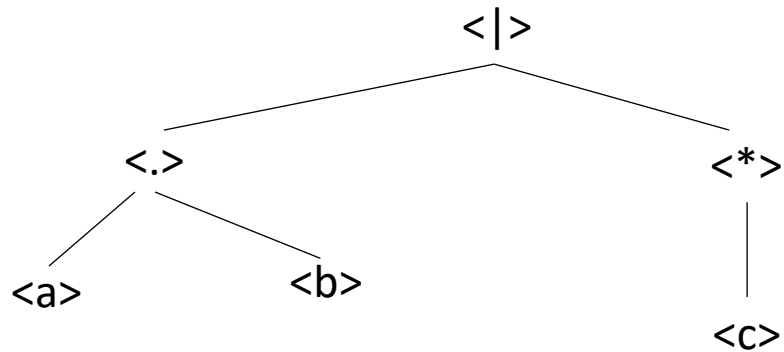
input: a.b | c\*



# Implementations

Create a parse tree for regular expressions

input:  $a.b \mid c^*$

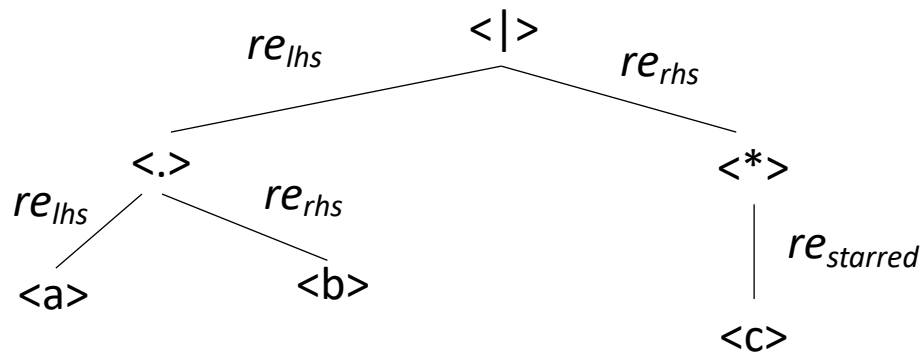


- *regular expression =*
  - |  $\{ \}$
  - |  $\epsilon$
  - |  $a$  (single character)
  - |  $re_{lhs} \mid re_{rhs}$
  - |  $re_{lhs} \cdot re_{rhs}$
  - |  $re_{starred}^*$

# Implementations

Create a parse tree for regular expressions

input:  $a.b \mid c^*$



- *regular expression =*
  - |  $\{ \}$
  - |  $\epsilon$
  - |  $a$  (single character)
  - |  $re_{lhs} \mid re_{rhs}$
  - |  $re_{lhs} \cdot re_{rhs}$
  - |  $re_{starred}^*$

# Implementations

Homework overview



# Next Lecture

- Starting on module 2: Optimizations and Flow analysis
- Start on homework!
  - Visit office hours tomorrow!
  - Post questions/comments on Canvas!