# CSE211: Compiler Design

Oct. 8, 2021

- **Topic**: Parsing regular expressions with derivatives

- **Questions**:
  - *How do you parse a regular expression?*
    *How do you parse a context free grammar?*

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

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

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

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

# Announcements

- Homework 1 is out
  - Due on the 18$^{th}$
  - Get started early!
  - Today we will do parsing with derivatives

- Reading for today:
  - first 7 pages
  - https://www.ccs.neu.edu/home/turon/re-deriv.pdf
  - **Not optional! It will make part 2 of the homework much easier**

- End of module 1, starting module 2 next week

# CSE211: Compiler Design

Oct. 8, 2021

- **Topic**: Parsing regular expressions with derivatives

- **Questions**:
  - *How do you parse a regular expression?*
    *How do you parse a context free grammar?*

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

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

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

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

# Parsing RE's with Derivatives

- A simple regular expression parser implementation
  - Given an RE AST, you can parse with very few lines of code

- Think recursively!

# Language Derivatives

- A language is a (potentially infinite) set of strings $\{s_1, s_2, s_3, s_4, ...\}$

- A language is regular if it can be captured using a regular expression

- Examples of regular languages:
  - *{"a"}, {"+"}, {"+", "-", "*", "\"}*
  - *{"1", "1+1", "1+1+1"}*
  - *{""}, 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 $\delta_c(L)$

- We'll go over some examples in the next slides

# Language Derivatives Examples

- $L = \{\text{"a"}\}$

- $\delta_a(L) = \{\text{""}\}$

- $\delta_b(L) = \{\}$

# Language Derivatives Examples

- $L = \{``+\text{''}, ``-\text{''}, ``*\text{''}, ``/\text{''}\}$

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

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

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

# Language Derivatives Examples

- $L = \{\text{"1"}, \text{"1+1"}, \text{"1+1+1"}, \text{"1+1+1+1"}, \ldots\}$

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

- $\delta_1 (L) = \{\text{""}, \text{"+1"}, \text{"+1+1"}, \text{"+1+1+1"}, \ldots\}$

- $\delta_{1+} (L) = \{\text{"1"}, \text{"1+1"}, \text{"1+1+1"}, ..\} = L$

# Language Derivatives Examples

- $L = \{\text{"aaa"}, \text{"ab"}, \text{"ba"}, \text{"bba"}\}$

- $\delta_a (L) = \{\text{"aa"}, \text{"b"}\}$

- $\delta_{aa} (L) = \{\text{"a"}\}$

- $\delta_b (L) = \{\text{"a"}, \text{"ba"}\}$

- $\delta_{ba} (L) = \{\text{""}\}$

# Regular Expressions

Recall we defined regular expressions recursively:

The three base cases: a character literal

- The RE for a character "a" is given by "a". It matches only the character "a"

- The RE for the empty string is is given by "" or $\varepsilon$

- The RE for the empty set is given by {}

# Regular Expressions

three recursive definitions

- The concatenation of two REs x and y is given by x.y and matches the strings of RE x **concatenated** with the strings of RE y

- The union of two REs x and y is given by x|y and matches the strings of RE x **or** the strings of RE y

- The Kleene star of an RE x is given by x* and matches the strings of RE x **repeated** 0 or more times

# Regular expressions recursive definition

re =
   |{}
   | ""
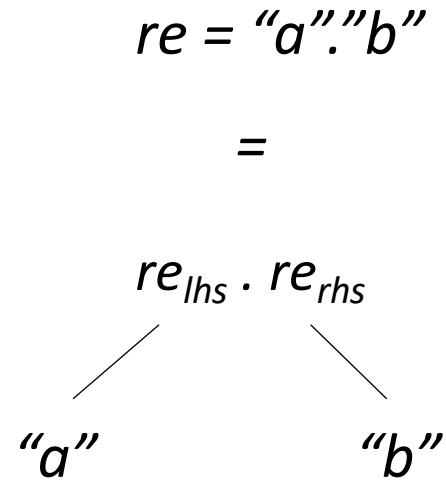   | $c$   (single character)
   | $re_{lhs}$ | $re_{rhs}$
   | $re_{lhs}$ . $re_{rhs}$
   | $re_{starred}$ *

# Regular expressions recursive definition

re =
   |{}
   | ""
   | $c$   (single character)
   | $re_{lhs}$ | $re_{rhs}$
   | $re_{lhs}$ . $re_{rhs}$
   | $re_{starred}$ *

$re = \text{"a"."b"}$

$=$

$re_{lhs} . re_{rhs}$

"a"         "b"

# parse tree for a regular expression

input: "a"."b" | "c"*

| Operator | Name | Productions |
|---|---|---|
| \| | union | : union PIPE concat<br>\| concat |
| . | concat | : concat CONCAT starred<br>\| starred |
| * | starred | : starred STAR<br>\| unit |
|  | unit | : CHAR<br>\| "" |

Excluding special cases for {}

# parse tree for a regular expression

input: "a"."b" | "c"*

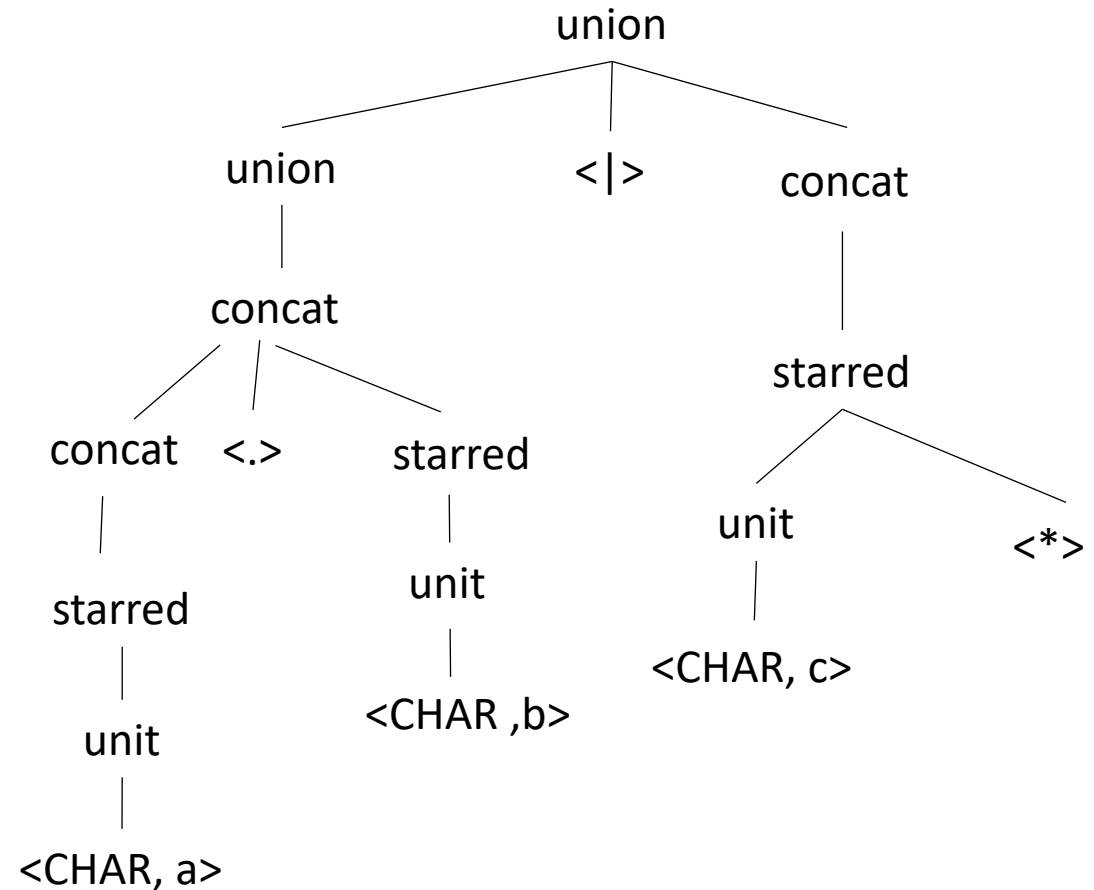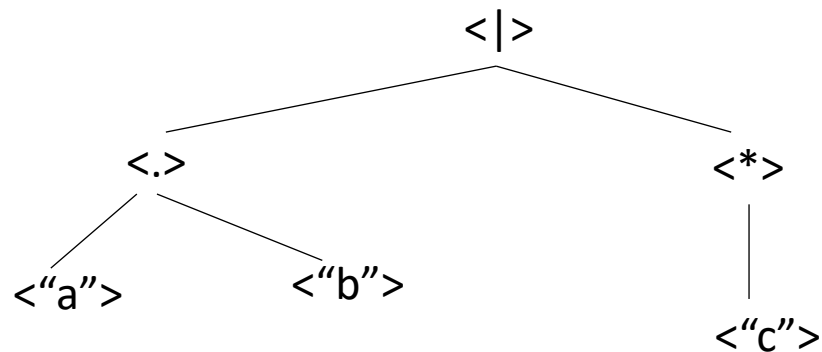| Operator | Name | Productions |
|----------|------|-------------|
| \| | union | : union PIPE concat<br>\| concat |
| . | concat | : concat CONCAT starred<br>\| starred |
| * | starred | : starred STAR<br>\| unit |
| | unit | : CHAR<br>\| "" |

Excluding special cases for {}

# parse tree for a regular expression

input: "a"."b" | "c"*
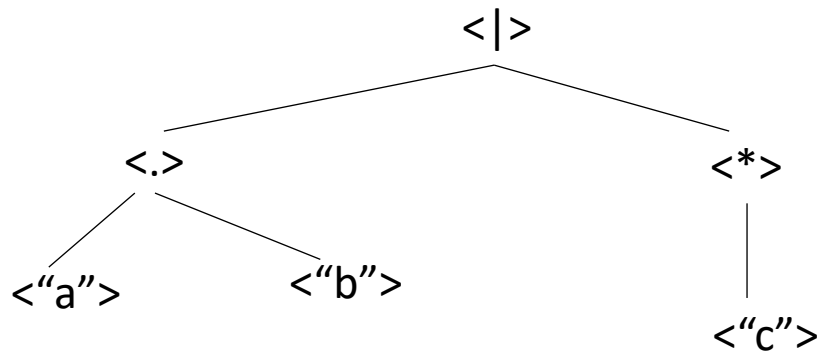
abstract syntax tree

# parse tree for a regular expression
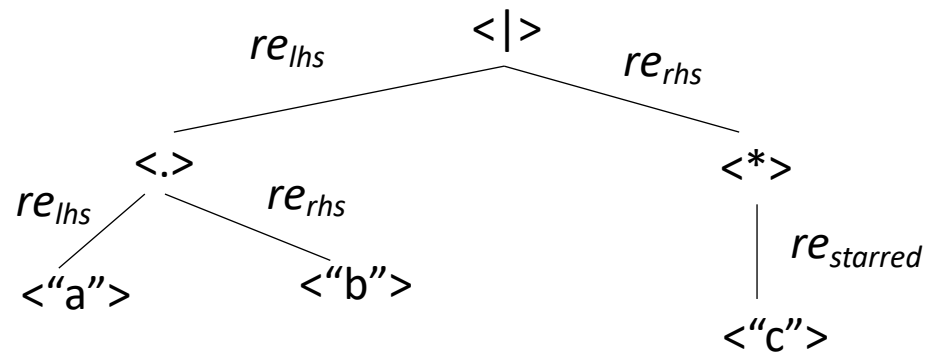
input: "a"."b" | "c"*

abstract syntax tree

```
            <|>
          /      \
       <.>         <*>
      /   \          |
  <"a">   <"b">    <"c">
```

- re =
  $|\{\}$
  $|$ ""
  $|$ a (single character)
  $|$ $re_{lhs}$ $|$ $re_{rhs}$
  $|$ $re_{lhs}$ . $re_{rhs}$
  $|$ $re_{starred}$ *

# parse tree for a regular expression

`input: "a"."b" | "c"*`

abstract syntax tree



- re =
  - |{}
  - | ""
  - | a (single character)
  - | $re_{lhs}$ | $re_{rhs}$
  - | $re_{lhs}$ . $re_{rhs}$
  - | $re_{starred}$ *

# parse tree for a regular expression

input: "a"."b" | "c"*

abstract syntax tree
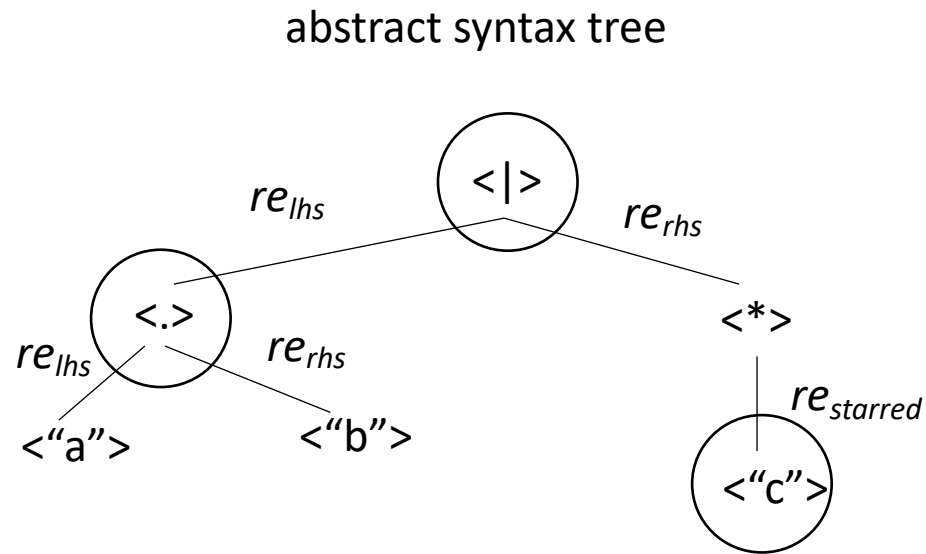
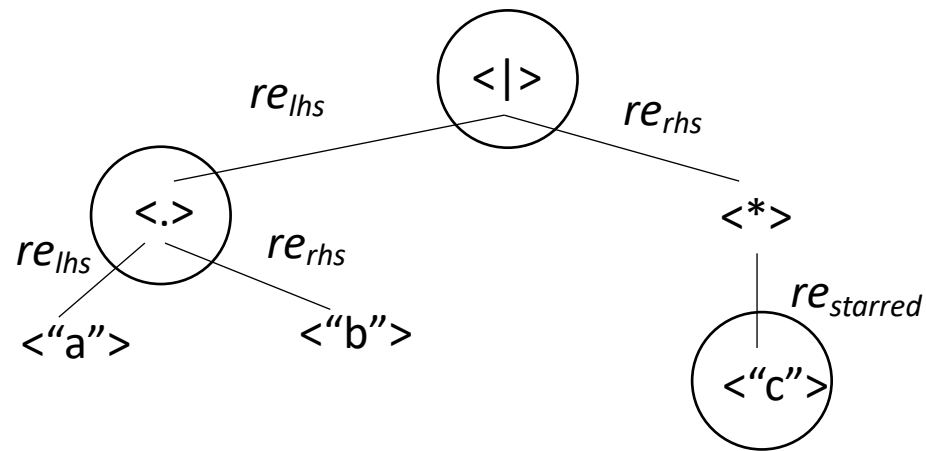each node is
also a regular expression!

- re =
  |{}
  | ""
  | a (single character)
  | $re_{lhs}$ | $re_{rhs}$
  | $re_{lhs}$ . $re_{rhs}$
  | $re_{starred}$ *

# parse tree for a regular expression

input: "a"."b" | "c"*

abstract syntax tree



each node is
also a regular expression!

- *Check homework code to see AST construction*

- *Question: given a regular expression AST, how check if a string is in the language?*

- *parsing with derivatives!*

# Regular expressions are closed under derivatives

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

- *Let's try some!*

# Regular expressions are closed under derivatives

- *re = "a"*

- {"a"}

- $\delta_a(re) = \{""\} = re("")$

- $\delta_b(re) = \{\}$

# Regular expressions are closed under derivatives

- *re = "a"*

- *L(re) = {"a"}*

- *$\delta_a$(re) = ""*

- *$\delta_b$(re) = {}*

# Regular expressions are closed under derivatives

- *re = "a" | "b"*

- *{"a","b"}*

- $\delta_a(re) = \{""\}$

- $\delta_b(re) = \{""\}$

# Regular expressions are closed under derivatives

- *re = "a" | "b"*

- *L(re) = {"a", "b"}*

- $\delta_a$*(re) = ""*

- $\delta_b$*(re) = ""*

# Regular expressions are closed under derivatives

- *re = "a"."a" | "a"."b"*

*{"aa", "ab"}*

- $\delta_a(re) = \{"a","b"\} = "a" \mid "b"$

- $\delta_b(re) = \{\}$

# Regular expressions are closed under derivatives

- *re = "a"."a" | "a"."b"*

- *L = {"aa", "ab"}*

- $\delta_a(re) = ??$

- $\delta_b(re) = ??$

# Regular expressions are closed under derivatives

- *re = "a"."a" | "a"."b"*

- *L = {"aa", "ab"}*

- $\delta_a(re) = \{$*"a", "b"*$\} = ??$

- $\delta_b(re) = \{\}$

# Regular expressions are closed under derivatives

- *re = "a"."a" | "a"."b"*

- *L = {"aa", "ab"}*

- $\delta_a(re) = \{$*"a", "b"$\}$ = *"a" | "b"*

- $\delta_b(re) = \{\}$

# Regular expressions are closed under derivatives

- *re = ("a"."b"."c")\**

- *{"", "abc", "abcabc", "abcabcabc", …}*

- $\delta_a$*(re) = {"bc", "bcacb", "bcabcabc" …} = "b"."c".("a"."b"."c")\**

# Regular expressions are closed under derivatives

- *re = ("a"."b"."c")\**

- *L = {"", "abc", "abcabc", "abcabcabc" ...}*

- *$\delta_a$(re) = ??*

# Regular expressions are closed under derivatives

- *re = ("a"."b"."c")\**

- *L = {"", "abc", "abcabc", "abcabcabc" …}*

- $\delta_a$*(re) = {"bc", "bcabc", "bcabcabc", …} = ??*

# Regular expressions are closed under derivatives

- *re = ("a"."b"."c")\**

- *L = {"", "abc", "abcabc", "abcabcabc" ...}*

- *$\delta_a$(re) = {"bc", "bcabc", "bcabcabc", ...} = "b"."c".("a"."b"."c")\**

# What is a method for computing the derivative?

Consider the base cases

- $\delta_c$ *(re)* = match re with:

  - {}
      return {}

  - ""
      return {}

  - *a* (single character)
      if a == c then return $\{\varepsilon\}$
      else return {}

- re =
  | {}
  | $\varepsilon$
  | a (single character)
  | $re_{lhs}$ | $re_{rhs}$
  | $re_{lhs}$ . $re_{rhs}$
  | $re_{starred}$ *

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ *(re)* = match re with:

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

    return ? ?

  - $re_{starred}$*

    return ? ?

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

    return    ? ?

- re =
    |{}
    | ε
    | a (single character)
    | $re_{lhs} \mid re_{rhs}$
    | $re_{lhs} \cdot re_{rhs}$
    | $re_{starred}$ *

# Regular expressions are closed under derivatives

- *re = "a"."a" | "a"."b"*

- *L = {"aa", "ab"}*

- *$\delta_a$(re) = {"a", "b"} = "a" | "b"*

- *$\delta_b$(re) = {}*

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ *(re) =* match re with:

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

                    return ??

    - $re_{starred}*$

                    return ??

    - $re_{lhs} . re_{rhs}$

                    return   ??

- re =

    |{}

    | ε

    | a (single character)

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

    | $re_{lhs} . re_{rhs}$

    | $re_{starred}*$

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c\ (re)$ = match re with:

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

    $\qquad$ return $\delta_c(re_{lhs})\ |\ \delta_c\ (re_{rhs})$

  - $re_{starred}{}^*$

    $\qquad$ return ? ?

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

    $\qquad$ return   ? ?

- re =
  - $|\{\}$
  - $|\ \varepsilon$
  - $|$ a (single character)
  - $|\ re_{lhs}\ |\ re_{rhs}$
  - $|\ re_{lhs}\ .\ re_{rhs}$
  - $|\ re_{starred}{}^*$

# Regular expressions are closed under derivatives

- *re = ("a"."b"."c")\**

- *L = {"", "abc", "abcabc", "abcabcabc" ...}*

- $\delta_a$*(re) = {"bc", "bcabc", "bcabcabc", ...} = "b"."c".("a"."b"."c")\**

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ *(re)* = match re with:

  - $re_{lhs}$ | $re_{rhs}$

    $\qquad$ return $\delta_c(re_{lhs})$ | $\delta_c (re_{rhs})$

  - $re_{starred}*$

    $\qquad$ return ? ?

  - $re_{lhs} . re_{rhs}$

    $\qquad$ return   ? ?

- re =
  - |{}
  - | ε
  - | a (single character)
  - | $re_{lhs}$ | $re_{rhs}$
  - | $re_{lhs}$ . $re_{rhs}$
  - | $re_{starred}*$

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ *(re)* = match re with:

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

    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 ??

- re =
    $\mid \{\}$
    $\mid \varepsilon$
    $\mid$ a (single character)
    $\mid re_{lhs} \mid re_{rhs}$
    $\mid re_{lhs} \cdot re_{rhs}$
    $\mid re_{starred}{}^*$

# Derivative Recursive Cases

Let's look at concatenation:

- $\delta_c$ *(re)* = match re with:

  - $re_{lhs} . re_{rhs}$

    return ??

*Example:*

*re = "a"."b"*

$\delta_a(re) = $ *"b"*

# Derivative Recursive Cases

Let's look at concatenation:

- $\delta_c\ (re)$ = match re with:

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

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

---

*Example:*

*re = "a"."b"*

$\delta_a(re) = $ *"b"*

# Derivative Recursive Cases

Let's look at concatenation:

- $\delta_c$ *(re)* = match re with:

  - $re_{lhs}$ . $re_{rhs}$

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

What about?

Example:

re = "c"*."a"."b"

$\delta_a(re)$ = "b"

# Derivative Recursive Cases

Let's look at concatenation:

- $\delta_c \, (re) =$ match re with:

  - $re_{lhs} \, . \, re_{rhs}$

    return $\quad \delta_c(re_{lhs}) \, . \, re_{rhs}$ |

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

Example:

re = "c"*."a"."b"

$\delta_a(re) =$ "b"

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ *(re)* = match re with:

  - $re_{lhs}$ | $re_{rhs}$

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

  - $re_{starred}*$

    return $\delta_c(re_{starred})$ . $re_{starred}*$

  - $re_{lhs}$ . $re_{rhs}$

    return  $\delta_c(re_{lhs})$ . $re_{rhs}$ |

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

- re =
  |{}
  | ε
  | a (single character)
  | $re_{lhs}$ | $re_{rhs}$
  | $re_{lhs}$ . $re_{rhs}$
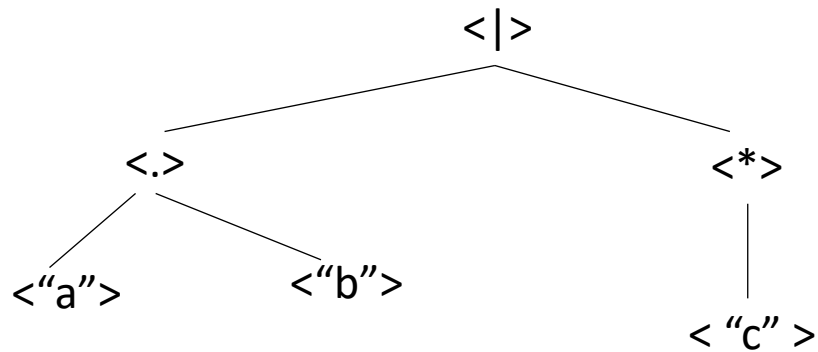  | $re_{starred}$ *

# Nullable operator

- NULL(re) =

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

# Nullable operator

- NULL(re) =

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

implement over a RE abstract syntax tree



- re =
  - |{}
  - | ""
  - | a (single character)
  - | $re_{lhs}$ | $re_{rhs}$
  - | $re_{lhs}$ . $re_{rhs}$
  - | $re_{starred}$ *

# What is a method for computing NULL?

Consider the base cases

- NULL(*re)* = match re with:

  - {}

    return {}

  - ""

    return ""

  - *a* (single character)

    return {}

- re =
  |{}
  | ""

  | a (single character)
  | $re_{lhs}$ | $re_{rhs}$
  | $re_{lhs}$ . $re_{rhs}$
  | $re_{starred}$ *

# What is a method for computing NULL?

Consider the recursive cases:

- NULL(*re) =* match re with:

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

    return NULL(re_lhs) | NULL(re_rhs)

  - $re_{starred}*$

    return ""

  - $re_{lhs} . re_{rhs}$

    return NULL(re_lhs) . NULL(re_rhs)

- re =
  |{}
  | ε
  | a (single character)
  | $re_{lhs} \mid re_{rhs}$
  | $re_{lhs} . re_{rhs}$
  | $re_{starred}*$

# What is a method for computing NULL?

Consider the recursive cases:

- NULL(*re)* = match re with:

  - $re_{lhs}$ | $re_{rhs}$

    return NULL($re_{lhs}$) | NULL($re_{rhs}$)

  - $re_{starred}$*

    return ""

  - $re_{lhs}$ . $re_{rhs}$

    return NULL($re_{lhs}$) . NULL($re_{rhs}$)

- re =
  - |{}
  - | ε
  - | a (single character)
  - | $re_{lhs}$ | $re_{rhs}$
  - | $re_{lhs}$ . $re_{rhs}$
  - | $re_{starred}$*

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ *(re)* = match re with:

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

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

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

# Derivative Recursive Cases

Consider the recursive cases:

- $\delta_c$ (re) = match re with:

  - $re_{lhs}$ | $re_{rhs}$

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

  - $re_{starred}*$

    return $\delta_c(re_{starred})$ . $re_{starred}*$

  - $re_{lhs}$ . $re_{rhs}$

    return    $\delta_c(re_{lhs})$ . $re_{rhs}$ |

    $NULL(re_{lhs})$ . $\delta_c(re_{rhs})$

- re =
  - {}
  - $\varepsilon$
  - a (single character)
  - $re_{lhs}$ | $re_{rhs}$
  - $re_{lhs}$ . $re_{rhs}$
  - $re_{starred}*$

# 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 $s = c_1 . c_2 . c_3 ...$ (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 $s = c_1 . c_2 . c_3 \ldots$ (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 $s = c_1 . c_2 . c_3 ...$ (concat of characters)

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

$\delta_{c1} (re)$

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

$L(\delta_{c1} (re)) = \{.. s[1:] ..\}$

# 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 $s = c_1 . c_2 . c_3 ...$ (concat of characters)

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

| | $\delta_{c1} (re)$ | $\delta_{c2} (\delta_{c1} (re)) = \delta_{c1,c2} (re)$ |
|---|---|---|
| $L(re) = \{.. s ..\}$ | | |
| | $L(\delta_{c1} (re)) = \{.. s[1:] ..\}$ | $L(\delta_{c1,c2} (re)) = \{.. s[2:] ..\}$ |

# 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 $s = c_1 . c_2 . c_3 ...$ (concat of characters)

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

| | $\delta_{c1}$ *(re)* | $\delta_{c2}$ ($\delta_{c1}$ (re) ) $= \delta_{c1,c2}$ *(re)* | $\delta_s(re)$ |
|---|---|---|---|
| $L(re) = \{.. s ..\}$ | | | |
| | $L(\delta_{c1}$ (re)) $= \{.. s[1:] ..\}$ | $L(\delta_{c1,c2}$ (re)) $= \{.. s[2:] ..\}$ | $L(\delta_s(re)) = \{.. \varepsilon ..\}$ |

# 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 $s = c_1 . c_2 . c_3 ...$ (concat of characters)

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

| $L(re) = \{.. \, s ..\}$ | $\delta_{c1} \, (re)$ | $\delta_{c2} \, (\delta_{c1} \, (re) \,) = \delta_{c1,c2} \, (re)$ | $\delta_s(re)$ | $NULL(\delta_s(re)) ==$ "" |
|---|---|---|---|---|
| | $L(\delta_{c1} \, (re)) = \{.. \, s[1:] ..\}$ | $L(\delta_{c1,c2} \, (re)) = \{.. \, s[2:] ..\}$ | $L(\delta_s(re)) = \{.. \, "" \, ..\}$ | |

# Have a good weekend!

*Take a look at part 2 of the homework, everything we discussed today is implemented there, with a few missing pieces for you to implement!*

*Next week we start module 2!*