# Tutorial 2

Computer Language Processing (COMP 4901U)

Monday September 20

- ▶ Recap on regular expression marching

- ▶ Exercise solving (Zoom breakout rooms available)

# Recap: Words

Let $A$ be an alphabet $\{a, b, c, \ldots\}$

We define words of length $n$, denoted $A^n$, as follows:

$A^0 = \{\varepsilon\}$ (only one word of length zero, always denoted $\varepsilon$)

For $n > 0$, $A^n = \{\, aw \mid w \in A^{n-1} \,\}$

Notation '$a\varepsilon$' abbreviated to '$a$'.

Concatenation: $u \cdot v$, also written $uv$, is associative.

We can decompose words arbitrarily, as in $w = ua$ (i.e., $w = u \cdot a\varepsilon$) if $|w| > 0$

Set of all words: $\qquad A^* = \bigcup_{n \geq 0} A^n$

which means: $w \in A^*$ if and only iff there exists $n$ such that $w \in A^n$.

# Recap: Languages

A *language* over alphabet $A$ is a set $L \subseteq A^*$.

Examples for $A = \{0, 1\}$:

- empty language $\varnothing$;
- finite languages like $L = \{1, 10, 1001\}$;
- language $L$ described by a characteristic function $f$, i.e., $L = \{w \in A^* \mid f(w)\}$

Operations on languages:

- Set operations: union ($\cup$), intersection ($\cap$), difference ($\backslash$), etc.
- Language concatenation: $L_1 \cdot L_2 = \{w_1 \cdot w_2 \mid w_1 \in L_1, w_2 \in L_2\}$, also written $L_1 L_2$
- Language exponentiation: $L^0 = \{\varepsilon\}$ and $L^{n+1} = L \cdot L^n$

# Recap: Regular Expressions

Syntax of regular expressions: $\quad e ::= \varnothing \mid \varepsilon \mid c \mid (e_1 \mid e_2) \mid e_1 e_2 \mid e^*$

The *semantics* of regular expression $e$ is the language $L(e)$ – also written $L_e$ – where:

$$
\begin{aligned}
L(\varnothing) &= \varnothing \\
L(\varepsilon) &= \{\varepsilon\} \\
L(c) &= \{c\} \qquad\qquad (c \in A) \\
L(e_1 \mid e_2) &= L(e_1) \cup L(e_2) \\
L(e_1 e_2) &= L(e_1) \cdot L(e_2) \\
L(e^*) &= L(e)^*
\end{aligned}
$$

Example: *letter*(*letter* | *digit*)* where *letter* = a | b | c | ... and *digit* = 0 | 1 | 2 ... | 9

# Recap: Closed Operations on Regular Expressions

Regular languages/expressions are closed under these operations:

- Shorthands for finite languages, such as $[a..z] = a \mid b \mid ... \mid z$
- Optionality $e^? = e \mid \varepsilon$
- Repeating at least once $e^+ = ee^*$
- Other repetitions $e^{k..*} = e^k e^*$ and $e^{p..q} = e^p (e^?)^{q-p}$
- Complementation: $!e$, denoting $A^* \setminus L(e)$ — Q: how to translate?
- Intersection $e_1 \& e_2 = !(!e \mid !e)$, denoting $L(e_1) \cap L(e_2)$

# Recap: Nullable

Formal definition of *nullable*: $\quad nullable(e) = \varepsilon \in L(e)$

Algorithm for computing *nullable*:

$$
\begin{aligned}
nullable(\emptyset) &= false \\
nullable(\varepsilon) &= true \\
nullable(a) &= false \\
nullable(e_1|e_2) &= nullable(e_1) \lor nullable(e_2) \\
nullable(e^*) &= true \\
nullable(e_1 e_2) &= nullable(e_1) \land nullable(e_2)
\end{aligned}
$$

# Recap: Lexers

Definition of a *lexer* (aka *tokenizer*):
an ordered set of *n* labelled token definitions

$$\langle\; Token_1 := e_1\,;\; Token_2 := e_2\,;\; \ldots\,;\; Token_n := e_n \;\rangle$$

Disambiguation rules:

- ▶ Longest-match
- ▶ First-match

# Recap: Regular Expression Derivative

### Definition (Brzozowski Derivative)

The *derivative* of regexp $e$ with respect to letter $c$, written $\delta^c(e)$, is defined as:

$$L(\delta^c(e)) = \{w \mid cw \in L(e)\}$$

# Recap: Regular Expression Derivative

### Definition (Brzozowski Derivative)

The *derivative* of regexp $e$ with respect to letter $c$, written $\delta^c(e)$, is defined as:

$$L(\delta^c(e)) = \{w \mid cw \in L(e)\}$$

A few derivative examples:

$\delta^a(\text{aaa}) = \text{aa}$

$\delta^a(\text{ab} \mid \text{ac} \mid \text{da}) = \text{b} \mid \text{c}$

$\delta^a((\text{ab})^*) = \text{b}(\text{ab})^*$

$\delta^a((\text{ab}|\text{c})^*\text{ad}) = \text{b}(\text{ab}|\text{c})^*\text{ad} \mid \text{d}$

## Recap: Regular Expression Derivative

Derivative of a regexp $e$ with respect to letter $c$, written $\delta^c(e)$, can be computed as:

$$\delta^c(\varnothing) = \varnothing$$

$$\delta^c(\varepsilon) = \varnothing$$

$$\delta^c(d) = \left\{ \begin{array}{ll} \varepsilon & \text{if } d = c \\ \varnothing & \text{if } d \neq c \end{array} \right.$$

$$\delta^c(e_1 \mid e_2) = \delta^c(e_1) \mid \delta^c(e_2)$$

$$\delta^c(e_1 \, e_2) = \left\{ \begin{array}{ll} \delta^c(e_1) \, e_2 \mid \delta^c(e_2) & \text{if } nullable(e_1) \\ \delta^c(e_1) \, e_2 & \text{otherwise} \end{array} \right.$$

$$\delta^c(e_1^*) = \delta^c(e_1) \, e_1^*$$

Exercises

# Exercise 1

Question:

Find a regular expression that generates all alternating sequences of 0 and 1 with arbitrary length (including lengths zero, one, two, etc.).

For example, the alternating sequences of length one are 0 and 1, length two are 01 and 10, length three are 010 and 101.

Note that no two adjacent character can be the same in an alternating sequence.

# Exercise 1

Question:

Find a regular expression that generates all alternating sequences of 0 and 1 with arbitrary length (including lengths zero, one, two, etc.).
For example, the alternating sequences of length one are 0 and 1, length two are 01 and 10, length three are 010 and 101.
Note that no two adjacent character can be the same in an alternating sequence.

Solution:

$$(10)^*1^? \,|\, (01)^*0^?$$

or, simpler:

$$0^?(10)^*1^?$$

# Exercise 2 – Integer Literals of Scala

Integer literals are in three forms in Scala: decimal, hexadecimal and octal. The compiler discriminates different classes from their beginning.

- ▶ Decimal integers are started with a non-zero digit.
- ▶ Hexadecimal numbers begin with 0x or 0X and may contain the digits from 0 through 9 as well as upper or lowercase digits from A to F.
- ▶ If the integer number starts with zero, it is in octal representation so it can contain only digits 0 through 7.
- ▶ l or L at the end of the literal shows the number is Long.

# Exercise 2 – Integer Literals of Scala

Integer literals are in three forms in Scala: decimal, hexadecimal and octal. The compiler discriminates different classes from their beginning.

▶ Decimal integers are started with a non-zero digit.

▶ Hexadecimal numbers begin with 0x or 0X and may contain the digits from 0 through 9 as well as upper or lowercase digits from A to F.

▶ If the integer number starts with zero, it is in octal representation so it can contain only digits 0 through 7.

▶ l or L at the end of the literal shows the number is Long.

Question: Write the corresponding regular expression.

# Exercise 2 – Integer Literals of Scala

Integer literals are in three forms in Scala: decimal, hexadecimal and octal. The compiler discriminates different classes from their beginning.

- ▶ Decimal integers are started with a non-zero digit.
- ▶ Hexadecimal numbers begin with 0x or 0X and may contain the digits from 0 through 9 as well as upper or lowercase digits from A to F.
- ▶ If the integer number starts with zero, it is in octal representation so it can contain only digits 0 through 7.
- ▶ l or L at the end of the literal shows the number is Long.

Question: Write the corresponding regular expression.

Solution: $((1..9)\text{digit}^* \mid 0(x \mid X)(\text{digit} \mid A..F)^* \mid 0(0..7)^*)(l \mid L)^?$
  where $\text{digit} = 0..9 = 0 \mid 1 \mid ... \mid 9$

# Exercise 2

Consider this lexer from last week where $A = \{a, b, c\}$:

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

# Exercise 2

Consider this lexer from last week where $A = \{a, b, c\}$:

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

Question: Construct the successive derivatives of each token for the given sequences:

- c
- ac
- cb
- bacacc

# Exercise 2 – Solutions

$$\left\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \right\rangle$$

$\delta^c(T_1) = \varnothing \quad \delta^c(T_2) = \varnothing \quad \delta^c(T_3) = ba \quad \delta^c(T_4) = c^*$

# Exercise 2 – Solutions

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

$$\delta^c(T_1) = \varnothing \quad \delta^c(T_2) = \varnothing \quad \delta^c(T_3) = \text{ba} \quad \delta^c(T_4) = c^*$$

$$\delta^a(T_1) = (ab)^* \quad \delta^a(T_2) = c(ac)^* \quad \delta^a(T_3) = \varnothing \quad \delta^a(T_4) = \varnothing$$

# Exercise 2 – Solutions

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

$\delta^c(T_1) = \varnothing \quad \delta^c(T_2) = \varnothing \quad \delta^c(T_3) = ba \quad \delta^c(T_4) = c^*$

$\delta^a(T_1) = (ab)^* \quad \delta^a(T_2) = c(ac)^* \quad \delta^a(T_3) = \varnothing \quad \delta^a(T_4) = \varnothing$

$\delta^{ac}(T_1) = \varnothing \quad \delta^{ac}(T_2) = (ac)^* \quad \delta^{ac}(T_3) = \varnothing \quad \delta^{ac}(T_4) = \varnothing$

# Exercise 2 – Solutions

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

$\delta^c(T_1) = \varnothing \quad \delta^c(T_2) = \varnothing \quad \delta^c(T_3) = ba \quad \delta^c(T_4) = c^*$

$\delta^a(T_1) = (ab)^* \quad \delta^a(T_2) = c(ac)^* \quad \delta^a(T_3) = \varnothing \quad \delta^a(T_4) = \varnothing$
$\delta^{ac}(T_1) = \varnothing \quad \delta^{ac}(T_2) = (ac)^* \quad \delta^{ac}(T_3) = \varnothing \quad \delta^{ac}(T_4) = \varnothing$

$\delta^{cb}(T_1) = \varnothing \quad \delta^{cb}(T_2) = \varnothing \quad \delta^{cb}(T_3) = a \quad \delta^{cb}(T_4) = \varnothing$

# Exercise 2 – Solutions

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

$$\delta^c(T_1) = \varnothing \quad \delta^c(T_2) = \varnothing \quad \delta^c(T_3) = ba \quad \delta^c(T_4) = c^*$$

$$\delta^a(T_1) = (ab)^* \quad \delta^a(T_2) = c(ac)^* \quad \delta^a(T_3) = \varnothing \quad \delta^a(T_4) = \varnothing$$
$$\delta^{ac}(T_1) = \varnothing \quad \delta^{ac}(T_2) = (ac)^* \quad \delta^{ac}(T_3) = \varnothing \quad \delta^{ac}(T_4) = \varnothing$$

$$\delta^{cb}(T_1) = \varnothing \quad \delta^{cb}(T_2) = \varnothing \quad \delta^{cb}(T_3) = a \quad \delta^{cb}(T_4) = \varnothing$$

…
$$\delta^{bacac}(T_1) = \varnothing \quad \delta^{bacac}(T_2) = (ac)^* \quad \delta^{bacac}(T_3) = \varnothing \quad \delta^{bacac}(T_4) = \varnothing$$

# Exercise 2 – Solutions

$$\langle T1 : a(ab)^*, T2 : b^*(ac)^*, T3 : cba, T4 : cc^* \rangle$$

$$\delta^c(T_1) = \varnothing \quad \delta^c(T_2) = \varnothing \quad \delta^c(T_3) = ba \quad \delta^c(T_4) = c^*$$

$$\delta^a(T_1) = (ab)^* \quad \delta^a(T_2) = c(ac)^* \quad \delta^a(T_3) = \varnothing \quad \delta^a(T_4) = \varnothing$$
$$\delta^{ac}(T_1) = \varnothing \quad \delta^{ac}(T_2) = (ac)^* \quad \delta^{ac}(T_3) = \varnothing \quad \delta^{ac}(T_4) = \varnothing$$

$$\delta^{cb}(T_1) = \varnothing \quad \delta^{cb}(T_2) = \varnothing \quad \delta^{cb}(T_3) = a \quad \delta^{cb}(T_4) = \varnothing$$

...
$$\delta^{bacac}(T_1) = \varnothing \quad \delta^{bacac}(T_2) = (ac)^* \quad \delta^{bacac}(T_3) = \varnothing \quad \delta^{bacac}(T_4) = \varnothing$$
$$\delta^{bacacc}(T_1) = \varnothing \quad \delta^{bacacc}(T_2) = \varnothing \quad \delta^{bacacc}(T_3) = \varnothing \quad \delta^{bacacc}(T_4) = \varnothing$$

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

▶ Sequence of open or closed parentheses of even length?
E.g. (), ((, )), )()))(, ...

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

▶ Sequence of open or closed parentheses of even length?
  E.g. ( ), (( , )), )()))( , ...
  ⇒ **YES**

▶ Sequence of balanced parentheses
  ( ( ( ) ) ( )) – balanced
  ( ) ) ( ( ) – not balanced

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

- ▶ Sequence of open or closed parentheses of even length?
  E.g. (), ((, )), )()))(, ...
  ⇒ **YES**
- ▶ Sequence of balanced parentheses
  ( ( ( ) ) ( )) – balanced
  ( ) ) ( ( ) – not balanced
  ⇒ **NO**
- ▶ As many digits before as after decimal point?

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

- Sequence of open or closed parentheses of even length?
    E.g. (), ((, )), )())) (, ...
  ⇒ **YES**
- Sequence of balanced parentheses
    ( ( ( ) ) ( )) – balanced
    ( ) ) ( ( ) – not balanced
  ⇒ **NO**
- As many digits before as after decimal point?    ⇒ **NO**
- Comments from // until LF

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

- ▶ Sequence of open or closed parentheses of even length?
  E.g. (), ((, )), )())(, ...
  ⇒ **YES**
- ▶ Sequence of balanced parentheses
  ( ( ( ) ) ( )) – balanced
  ( ) ) ( ( ) – not balanced
  ⇒ **NO**
- ▶ As many digits before as after decimal point? ⇒ **NO**
- ▶ Comments from // until LF ⇒ **YES**
- ▶ Nested comments like /* ... /* */ ... */

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

- Sequence of open or closed parentheses of even length?
  E.g. (), ((, )), )()))(,  ...
  ⇒ **YES**
- Sequence of balanced parentheses
  ( ( ( ) ) ( )) – balanced
  ( ) ) ( ( ) – not balanced
  ⇒ **NO**
- As many digits before as after decimal point?    ⇒ **NO**
- Comments from // until LF    ⇒ **YES**
- Nested comments like /* ... /* */ ... */    ⇒ **NO**

# Exercise 3 – On the expressiveness of regular expressions

For which of the following languages can you find an automaton or regular expression?

- ▶ Sequence of open or closed parentheses of even length?
  E.g. (), ((, )), )())(, ...
  ⇒ **YES**
- ▶ Sequence of balanced parentheses
  ( ( ( ) ) ( )) – balanced
  ( ) ) ( ( ) – not balanced
  ⇒ **NO**
- ▶ As many digits before as after decimal point?    ⇒ **NO**
- ▶ Comments from // until LF    ⇒ **YES**
- ▶ Nested comments like /* ... /* */ ... */    ⇒ **NO**

*Regular expressions cannot "count."*

# Exercise 4

Let *init* be a function that returns all the symbols of a string except the last one.

For example $init(mama) = mam$

*init* is undefined for an empty string.

If $L_1 \subseteq A^*$, then $INIT(L_1)$ applies the function to all non-empty words in $L_1$, ignoring $\varepsilon$ if it is in $L_1$:

$$INIT(L_1) = \{ v \in A^* \mid \exists a \in A.\ va \in L_1 \}$$

# Exercise 4

Let *init* be a function that returns all the symbols of a string except the last one.
For example $init(mama) = mam$
*init* is undefined for an empty string.
If $L_1 \subseteq A^*$, then $INIT(L_1)$ applies the function to all non-empty words in $L_1$, ignoring $\varepsilon$
if it is in $L_1$:

$$INIT(L_1) = \{ v \in A^* \mid \exists a \in A. \ va \in L_1 \}$$

Question: Give an algorithm that, given a regular expression $r$ for $L_1$, computes a
regular expression $rinit(r)$ for language $INIT(L_1)$

# Exercise 4

Let *init* be a function that returns all the symbols of a string except the last one.
For example $init(mama) = mam$
*init* is undefined for an empty string.
If $L_1 \subseteq A^*$, then $INIT(L_1)$ applies the function to all non-empty words in $L_1$, ignoring $\varepsilon$
if it is in $L_1$:

$$INIT(L_1) = \{ v \in A^* \mid \exists a \in A.\ va \in L_1 \}$$

Question: Give an algorithm that, given a regular expression $r$ for $L_1$, computes a
regular expression $rinit(r)$ for language $INIT(L_1)$

Solution: Just follow the same approach as for defining derivatives, using the other
direction for inductive decomposition of words.