# CMSC 330: Organization of Programming Languages

## Regular Expressions and Finite Automata

## How do regular expressions work?

- What we've learned
  - What regular expressions are
  - What they can express, and cannot
  - Programming with them
- · What's next: how they work
  - A great computer science result

## Languages and Machines



#### A Few Questions About REs

- How are REs implemented?
  - Given an arbitrary RE and a string, how to decide whether the RE matches the string?
- · What are the basic components of REs?
  - Can implement some features in terms of others
    - E.g., e+ is the same as ee\*
- · What does a regular expression represent?
  - Just a set of strings
    - This observation provides insight on how we go about our implementation
- ... next comes the math !

#### **Definition:** Alphabet

- An alphabet is a finite set of symbols
  - Usually denoted Σ
- Example alphabets:
  - Binary: Σ = {0,1}
  - Decimal: Σ = {0,1,2,3,4,5,6,7,8,9}
  - Alphanumeric:  $\Sigma = \{0-9, a-z, A-Z\}$

## **Definition: String**

- A string is a finite sequence of symbols from  $\boldsymbol{\Sigma}$ 
  - ε is the empty string ("" in Ruby)
  - |s| is the length of string s
    - |Hello| = 5,  $|\varepsilon| = 0$
  - Note
    - Ø is the empty set (with 0 elements)
    - $\square \emptyset \neq \{ \epsilon \} (and \emptyset \neq \epsilon)$
- Example strings over alphabet  $\Sigma = \{0,1\}$  (binary):
  - 0101
  - 0101110
  - 8

## **Definition: String concatenation**

• String concatenation is indicated by juxtaposition

 $s_1 = super$   $s_2 = hero$   $s_1s_2 = superhero$ •Sometimes also written  $s_1 \cdot s_2$ 



- For any string s, we have  $s\epsilon = s = \epsilon s$ 
  - You can concatenate strings from different alphabets; then the new alphabet is the union of the originals:

 $\label{eq:sigma_linear} If s_1 = super from \Sigma_1 = \{s, u, p, e, r\} and s_2 = hero from \Sigma_2 = \{h, e, r, o\}, then s_1s_2 = superhero from \Sigma_3 = \{e, h, o, p, r, s, u\}$ 

## **Definition: Language**

- A language L is a set of strings over an alphabet
- Example: All strings of length 1 or 2 over alphabet Σ = {a, b, c} that begin with a
  - L = { a, aa, ab, ac }
- Example: All strings over  $\Sigma = \{a, b\}$ 
  - L = {  $\epsilon$ , a, b, aa, bb, ab, ba, aaa, bba, aba, baa, ... }
  - Language of all strings written  $\boldsymbol{\Sigma}^{\star}$
- Example: All strings of length 0 over alphabet  $\Sigma$ 
  - L = { s | s  $\in \Sigma^*$  and |s| = 0 }

"the set of strings s such that s is from  $\Sigma^*$  and has length 0"

 $= \{ \epsilon \} \neq \emptyset$ 

## Definition: Language (cont.)

- Example: The set of phone numbers over the alphabet Σ = {0, 1, 2, 3, 4, 5, 6, 7, 9, (, ), -}
  - Give an example element of this language (123) 456-7890
  - Are all strings over the alphabet in the language? No
  - Is there a Ruby regular expression for this language?
     /\(\d{3,3}\)\d{3,3}-\d{4,4}/
- Example: The set of all valid (runnable) Ruby programs
  - Later we'll see how we can specify this language
  - (Regular expressions are useful, but not sufficient)

#### **Operations on Languages**

- Let  $\Sigma$  be an alphabet and let L, L<sub>1</sub>, L<sub>2</sub> be languages over  $\Sigma$
- Concatenation L<sub>1</sub>L<sub>2</sub> creates a language defined as
  - $L_1L_2 = \{ xy \mid x \in L_1 \text{ and } y \in L_2 \}$
- Union creates a language defined as
  - $L_1 \cup L_2 = \{ x \mid x \in L_1 \text{ or } x \in L_2 \}$
- Kleene closure creates a language is defined as
  - $L^* = \{ x \mid x = \varepsilon \text{ or } x \in L \text{ or } x \in LL \text{ or } x \in LLL \text{ or } ... \}$

#### **Operations Examples**

Let 
$$L_1 = \{a, b\}, L_2 = \{1, 2, 3\}$$
 (and  $\Sigma = \{a, b, 1, 2, 3\}$ )

- What is  $L_1L_2$ ?
  - { a1, a2, a3, b1, b2, b3 }
- What is L<sub>1</sub> U L<sub>2</sub>?
  { a, b, 1, 2, 3 }
- What is  $L_1^*$ ?
  - { ε, a, b, aa, bb, ab, ba, aaa, aab, bba, bbb, aba, abb, baa, bab, ... }

## Quiz 1: Which string is **not** in $L_3$

$$L_1 = \{a, ab, c, d, \epsilon\}$$
 where  $\Sigma = \{a, b, c, d\}$   
 $L_2 = \{d\}$   
 $L_3 = L_1 \cup L_2$ 

## Quiz 1: Which string is **not** in L<sub>3</sub>

$$L_1 = \{a, ab, c, d, \epsilon\}$$
 where  $\Sigma = \{a, b, c, d\}$   
 $L_2 = \{d\}$   
 $L_3 = L_1 \cup L_2$ 



### Quiz 2: Which string is **not** in L<sub>3</sub>



## Quiz 2: Which string is **not** in $L_3$



#### **Regular Expressions: Grammar**

- We can define a grammar for regular expressions R
  - $\mathbf{R} ::= \emptyset$  The empty language
    - ε The empty string
    - $\sigma$  A symbol from alphabet Σ
    - $R_1 R_2$  The concatenation of two regexps
    - $|R_1|R_2$  The union of two regexps
    - **R**<sup>\*</sup> The Kleene closure of a regexp

## **Regular Languages**

- Regular expressions denote languages. These are the regular languages
  - aka regular sets
- Not all languages are regular
  - Examples (without proof):
    - $\hfill\square$  The set of palindromes over  $\Sigma$
    - $[a^n b^n | n > 0] (a^n = sequence of n a's)$
- Almost all programming languages are not regular
  - But aspects of them sometimes are (e.g., identifiers)
  - Regular expressions are commonly used in parsing tools

## Semantics: Regular Expressions (1)

• Given an alphabet  $\Sigma$ , the regular expressions over  $\Sigma$  are defined inductively as follows



#### Constants

## Semantics: Regular Expressions (2)

 Let A and B be regular expressions denoting languages L<sub>A</sub> and L<sub>B</sub>, respectively. Then:

# regular expressiondenotes languageAB $L_A L_B$ A|B $L_A \cup L_B$ A\* $L_A^*$

#### **Operations**

• There are no other regular expressions over  $\Sigma$ 

#### Terminology etc.

- Regexps apply operations to symbols
  - Generates a set of strings (i.e., a language)
    - I (Formal definition shortly)
  - Examples
    - a generates language {a}
    - □ a|b generates language {a} U {b} = {a, b}
    - □ a<sup>\*</sup> generates language { $\epsilon$ } U {a} U {aa} U ... = { $\epsilon$ , a, aa, ... }
- If s ∈ language L generated by a RE r, we say that r accepts, describes, or recognizes string s

#### Precedence

- Order in which operators are applied is:
  - Kleene closure \* > concatenation > union
  - $ab|c = (ab)|c \rightarrow \{ab, c\}$
  - $ab^* = a (b^*) \longrightarrow \{a, ab, abb \dots\}$
  - $a|b^* = a | (b^*) \rightarrow \{a, \epsilon, b, bb, bbb \dots\}$
- · We use parentheses () to clarify
  - E.g., a(b|c), (ab)\*, (a|b)\*
  - Using escaped \( if parens are in the alphabet

## **Ruby Regular Expressions**

- Almost all of the features we've seen for Ruby REs can be reduced to this formal definition
  - /Ruby/ concatenation of single-symbol REs
  - /(Ruby|Regular)/ union
  - /(Ruby)\*/ Kleene closure
  - /(Ruby)+/ same as (Ruby)(Ruby)\*
  - /(Ruby)?/ same as (ε|(Ruby))
  - /[a-z]/ same as (a|b|c|...|z)
  - / [^0-9]/ same as (a|b|c|...) for a,b,c,...  $\in \Sigma$  {0..9}
  - ^, \$ correspond to extra symbols in alphabet

Think of every string containing a distinct, hidden symbol at its start and at CMSC 330 Sprinit and – these are written ^ and \$

## Implementing Regular Expressions

- We can implement a regular expression by turning it into a finite automaton
  - A "machine" for recognizing a regular language



## **Finite Automaton**



#### **Elements**

- States S (*start*, *final*)
- Alphabet Σ
- Transition edges δ



- Scan the next symbol  $\sigma \in \Sigma$  of the string s
- Take transition edge labeled with  $\sigma$
- String s is accepted if automaton is in final state when end of string s is reached

## **Finite Automaton: States**

- Start state
  - State with incoming transition from no other state
  - Can have only one start state



- Final states
  - States with double circle
  - Can have zero or more final states
  - Any state, including the start state, can be final











## Quiz 3: What Language is This?



A. All strings over {0, 1}
B. All strings over {1}
C. All strings over {0, 1} of length 1
D. All strings over {0, 1} that end in 1

## Quiz 3: What Language is This?



A. All strings over {0, 1}
B. All strings over {1}
C. All strings over {0, 1} of length 1
D. All strings over {0, 1} that end in 1 regular expression for this language is (0|1)\*1





















#### Quiz 4: Which string is **not** accepted?



#### Quiz 4: Which string is **not** accepted?





What language does this FA accept?

a\*b\*c\*

S3 is a dead state – a nonfinal state with no transition to another state - aka a trap state





#### Language? a\*b\*c\* again, so FAs are not unique

#### **Dead State: Shorthand Notation**

 If a transition is omitted, assume it goes to a dead state that is not shown



 Strings over {0,1,2,3} with alternating even and odd digits, beginning with odd digit



- Description for each state
  - S0 = "Haven't seen anything yet" OR "Last symbol seen was a b"
  - S1 = "Last symbol seen was an a"
  - S2 = "Last two symbols seen were ab"
  - S3 = "Last three symbols seen were abb"



Language as a regular expression?
(a|b)\*abb

## Quiz 5



- A. A string that contains a single b.
- B. Any string in {a,b}.
- c. A string that starts with b followed by a's.
- D. One or more b's, followed by zero or more a's.

## Quiz 5



- A. A string that contains a single b.
- B. Any string in {a,b}.
- c. A string that starts with b followed by a's.
- D. One or more b's, followed by zero or more a's.

- That accepts strings containing two consecutive 0s followed by two consecutive 1s
- That accepts strings with an odd number of 1s
- That accepts strings containing an even number of 0s and any number of 1s
- That accepts strings containing an odd number of 0s and odd number of 1s
- That accepts strings that DO NOT contain odd number of 0s and an odd number of 1s

• That accepts strings with an odd number of 1s

• That accepts strings with an odd number of 1s



 That accepts strings containing an even number of a's and any number of b's

 That accepts strings containing an even number of 0s and any number of 1s



 That accepts strings containing two consecutive 0s followed by two consecutive 1s

 That accepts strings containing two consecutive 0s very immediately (right after, no other things in between) followed by two consecutive 1s



 That accepts strings end with two consecutive 0s followed by two consecutive 1s

 That accepts strings end with two consecutive 0s followed by two consecutive 1s



 That accepts strings containing an odd number of 0s and odd number of 1s

 That accepts strings containing an odd number of 0s and odd number of 1s

![](_page_59_Figure_2.jpeg)

 That accepts strings that DO NOT contain odd number of 0s and an odd number of 1s

 That accepts strings that DO NOT contain odd number of 0s and an odd number of 1s

![](_page_61_Figure_2.jpeg)