

# CMSC 330: Organization of Programming Languages

## Regular Expressions and Finite Automata

CMSC 330

1

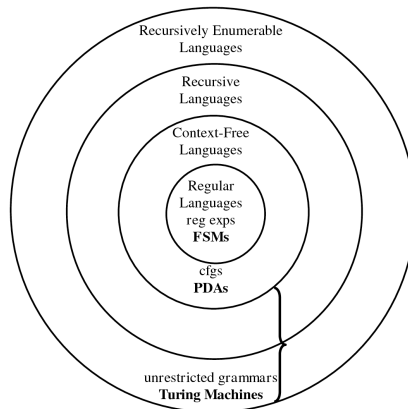
## How do regular expressions work?

- ▶ What we've learned
  - What regular expressions are
  - What they can express, and cannot
  - Programming with them
- ▶ What's next: how they work
  - Mixture of a very practical tool (string matching with REs) and some nice theory
  - A great computer science result

CMSC 330

2

## Languages and Machines



Chapter 1 - 3

## A Few Questions About REs

- ▶ How are REs implemented?
  - Implementing a one-off RE is not so hard
    - ▶ How to do it in general?
  - We'll see how to build a structure to parse REs
- ▶ What are the basic components of REs?
  - Can implement some features in terms of others
    - ▶ E.g., we saw that  $e^+$  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

CMSC 330

4

## Definition: Alphabet

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

CMSC 330

5

## Definition: String

- ▶ A **string** is a finite sequence of symbols from  $\Sigma$ 
  - $\epsilon$  is the empty string (" in Ruby)
  - $|s|$  is the length of string  $s$ 
    - ▶  $|\text{Hello}| = 5, |\epsilon| = 0$
  - Note
    - ▶  $\emptyset$  is the empty set (with 0 elements)
    - ▶  $\emptyset \neq \{\epsilon\} \neq \epsilon$
- ▶ Example strings:
  - $0101 \in \Sigma = \{0,1\}$  (binary)
  - $0101 \in \Sigma = \text{decimal}$
  - $0101 \in \Sigma = \text{alphanumeric}$

CMSC 330

6

## Definition: String concatenation

- ▶ String **concatenation** is indicated by juxtaposition
  - If  $s_1 = \text{super}$  and  $s_2 = \text{hero}$ , then  $s_1s_2 = \text{superhero}$
  - Sometimes also written  $s_1 \cdot s_2$
  - For any string  $s$ , we have  $s\epsilon = \epsilon s = s$
  - You **can** concatenate strings from different alphabets; then the new alphabet is the union of the originals:
    - ▶ If  $s_1 = \text{super} \in \Sigma_1 = \{s,u,p,e,r\}$  and  $s_2 = \text{hero} \in \Sigma_2 = \{h,e,r,o\}$ , then  $s_1s_2 = \text{superhero} \in \Sigma_3 = \{e,h,o,p,r,s,u\}$



CMSC 330

7

## Definition: Language

- ▶ A **language**  $L$  is a set of strings over an alphabet
- ▶ Example: The set of phone numbers over the alphabet  $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 9, (, ), -\}$ 
  - 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?  
$$\wedge (\d{3,3}) \d{3,3} - \d{4,4} /$$
- ▶ Example: The set of all strings over  $\Sigma$ 
  - Often written  $\Sigma^*$

CMSC 330

8

## Definition: Language (cont.)

- ▶ Example: The set of strings of length 0 over the alphabet  $\Sigma = \{a, b, c\}$ 
  - $L = \{s \mid s \in \Sigma^* \text{ and } |s| = 0\} = \{\epsilon\} \neq \emptyset$
- ▶ Example: The set of all valid Ruby programs
  - Is there a Ruby regular expression for this language?

No. Matching (an arbitrary number of) brackets so that they are balanced is impossible using REs  $\{\{\dots\}\}$
- ▶ REs represent languages, but not all of them
  - The languages represented by regular expressions are called, appropriately, the **regular languages**

CMSC 330

9

## Definition: Regular Expressions

- ▶ Given an alphabet  $\Sigma$ , the **regular expressions** over  $\Sigma$  are defined inductively as

regular expression	denotes language
$\emptyset$	$\emptyset$
$\epsilon$	$\{\epsilon\}$
each element $\sigma \in \Sigma$	$\{\sigma\}$

**Constants**

CMSC 330

10

## Definition: Regular Expressions (cont.)

- ▶ Let **A** and **B** be regular expressions denoting languages  $L_A$  and  $L_B$ , respectively

regular expression	denotes language
<b>AB</b>	$L_A L_B$
<b>(A B)</b>	$L_A \cup L_B$
<b>A*</b>	$L_A^*$

**Operations**

- ▶ There are no other regular expressions over  $\Sigma$

CMSC 330

11

## Operations on Languages

- ▶ Let  $\Sigma$  be an alphabet and let  $L, L_1, L_2$  be languages over  $\Sigma$
- ▶ Concatenation  $L_1 L_2$  is defined as
  - $L_1 L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}$
- ▶ Union is defined as
  - $L_1 \cup L_2 = \{x \mid x \in L_1 \text{ or } x \in L_2\}$
- ▶ **Kleene closure** is defined as
  - $L^* = \{x \mid x = \epsilon \text{ or } x \in L \text{ or } x \in LL \text{ or } x \in LLL \text{ or } \dots\}$

CMSC 330

12

## Regular Expressions Denote Languages

- ▶ By applying operations on constants
  - Generates a set of strings (i.e., a language)
  - Examples
    - >  $a \rightarrow \{“a”\}$
    - >  $a|b \rightarrow \{“a”\} \cup \{“b”\} = \{“a”, “b”\}$
    - >  $a^* \rightarrow \{\epsilon\} \cup \{“a”\} \cup \{“aa”\} \cup \dots = \{\epsilon, “a”, “aa”, \dots\}$
- ▶ If  $s \in$  language generated by a RE  $r$ , we say that  $r$  **accepts**, **describes**, or **recognizes** string  $s$

CMSC 330

13

## Precedence

- ▶ Order in which operators are applied
  - In arithmetic
    - > Multiplication  $\times$  > addition  $+$
    - >  $2 \times 3 + 4 = (2 \times 3) + 4 = 10$
  - In regular expressions
    - > Kleene closure  $*$  > concatenation > union  $|$
    - >  $ab|c = (ab) | c = \{“ab”, “c”\}$
    - >  $ab^* = a(b^*) = \{“a”, “ab”, “abb” \dots\}$
    - >  $a|b^* = a | (b^*) = \{“a”, “”, “b”, “bb”, “bbb” \dots\}$
  - Can change order using parentheses  $( )$ 
    - > E.g.,  $a(b|c)$ ,  $(ab)^*$ ,  $(a|b)^*$

CMSC 330

14

## Regular Languages

- ▶ The languages that can be described using regular expressions are the **regular languages** or **regular sets**
- ▶ Not all languages are regular
  - Examples (without proof):
    - > The set of palindromes over  $\Sigma$ 
      - reads the same backward or forward
    - >  $\{a^n b^n \mid n > 0\}$  ( $a^n$  = 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

CMSC 330

15

## Ruby Regular Expressions

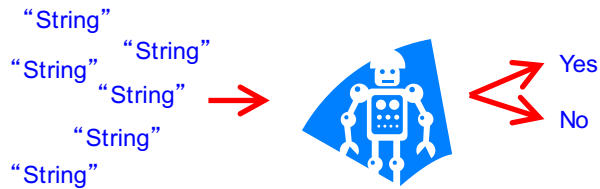
- ▶ Almost all of the features we've seen for Ruby REs can be reduced to this formal definition
  - $/\text{Ruby}/$  – concatenation of single-character REs
  - $/(\text{Ruby}|\text{Regular})/$  – union
  - $/(\text{Ruby})^*/$  – Kleene closure
  - $/(\text{Ruby})+/$  – same as  $(\text{Ruby})(\text{Ruby})^*$
  - $/(\text{Ruby})?/$  – same as  $(\epsilon | (\text{Ruby}))$  ( $/$  is  $\epsilon$ )
  - $/[a-z]/$  – same as  $(a|b|c|\dots|z)$
  - $/[^\wedge0-9]/$  – same as  $(a|b|c|\dots)$  for  $a, b, c, \dots \in \Sigma - \{0..9\}$
  - $^\wedge$ ,  $^\$$  – correspond to extra characters in alphabet

CMSC 330

16

## Implementing Regular Expressions

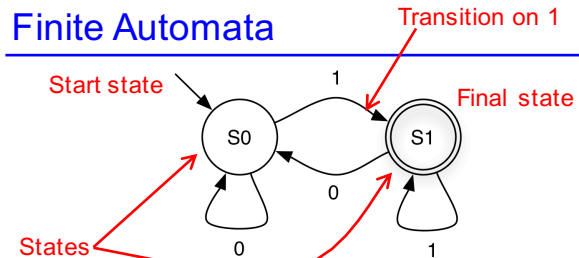
- ▶ We can implement a regular expression by turning it into a **finite automaton**
  - A “machine” for recognizing a regular language



CMSC 330

17

## Finite Automata



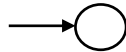
- ▶ Machine starts in **start** or **initial** state
- ▶ Repeat until the end of the string is reached
  - Scan the next **symbol s** of the string
  - Take transition edge labeled with **s**
- ▶ String is **accepted** if automaton is in **final** state when end of string reached

CMSC 330

18

## Finite Automata: States

- ▶ **Start state**
  - State with incoming transition from no other state
  - Can have only 1 start state



- ▶ **Final states**

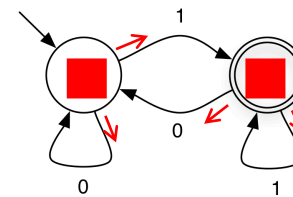
- States with double circle
- Can have 0 or more final states
- Any state, including the start state, can be final



CMSC 330

19

## Finite Automaton: Example 1



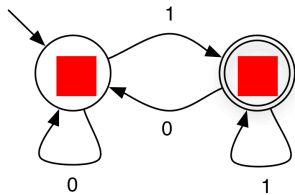
0 0 1 0 1 1

accepted

CMSC 330

20

## Finite Automaton: Example 2



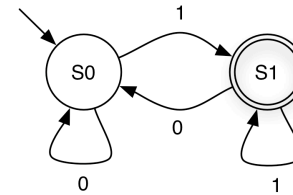
0 0 1 0 1 0

not accepted

CMSC 330

21

## What Language is This?

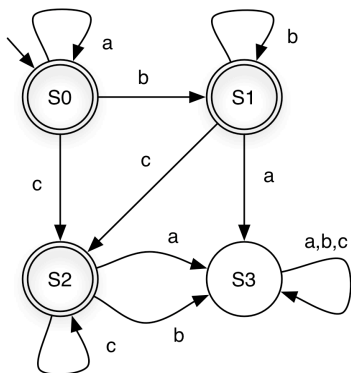


- ▶ All strings over  $\{0, 1\}$  that end in 1
- ▶ What is a regular expression for this language?  
 $(0|1)^*1$

CMSC 330

22

## Finite Automaton: Example 3



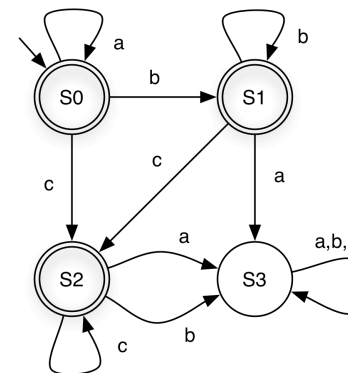
string	state at end	accepts ?
aabcc	S2	Y
acc	S2	Y
bbc	S2	Y
aabbb	S1	Y
aa	S0	Y
$\epsilon$	S0	Y
acba	S3	N

(a,b,c notation shorthand for three self loops)

CMSC 330

23

## Finite Automaton: Example 3 (cont.)



What language does this FA accept?

$a^*b^*c^*$

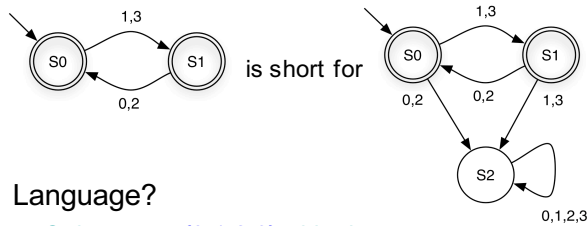
S3 is a **dead state** – a nonfinal state with **no** transition to another state

CMSC 330

24

## Dead State: Shorthand Notation

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



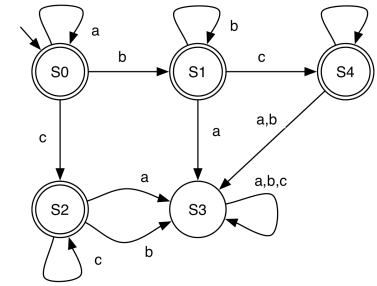
- ▶ Language?

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

CMSC 330

25

## Finite Automaton: Example 4

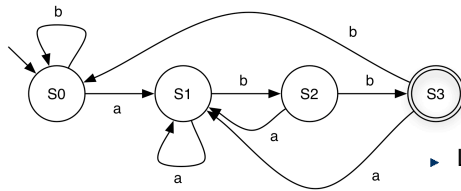


$a^*b^*c^*$  again, so DFAs are not unique

CMSC 330

26

## Finite Automaton: Example 5



- ▶ Language?
  - $(a|b)^*abb$

- ▶ Description for each state

- $S_0$  = "Haven't seen anything yet" OR "seen zero or more b's" OR "Last symbol seen was a b"
- $S_1$  = "Last symbol seen was an a"
- $S_2$  = "Last two symbols seen were ab"
- $S_3$  = "Last three symbols seen were abb"

## The Questions (for next time)

- ▶ Are FAs **equivalent** to regular expressions?
  - Every FA can be translated into an RE
  - Every RE can be translated into an FA
  - Yes!
- ▶ How can we **generate** an FA for a given RE?
  - If we can do this, we can implement RE matching
- ▶ How can we **optimize** an FA?
  - Many FAs can implement the same language
  - Some might be more efficient than others

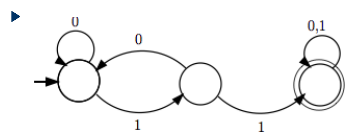
CMSC 330

28

## Practice

Give the English descriptions and the DFA or regular expression of the following languages:

- ▶  $((0|1)(0|1)(0|1)(0|1)(0|1))^*$ 
  - All strings with length a multiple of 5
- ▶  $(01)^*|(10)^*|(01)^*0|(10)^*1$ 
  - All alternating binary strings



All binary strings containing the substring "11"

CMSC 330

29

## Practice

▶ Give the regular expressions and finite automata for the following languages

- You and your neighbors' names
- All protein-coding DNA strings (including only ATCG and appearing in multiples of 3)
- All binary strings containing an even length substring of all 1's
- All binary strings containing exactly two 1's
- All binary strings that start and end with the same number

CMSC 330

30

## Review

- ▶ Languages
  - Sets of strings
  - Operations on languages
- ▶ Regular expressions
  - Constants
  - Operators
  - Precedence
- ▶ Finite automata
  - States
  - Transitions
  - Accept strings

CMSC 330

31