# Finite Automata & Regular Expressions

https://xkcd.com/208/
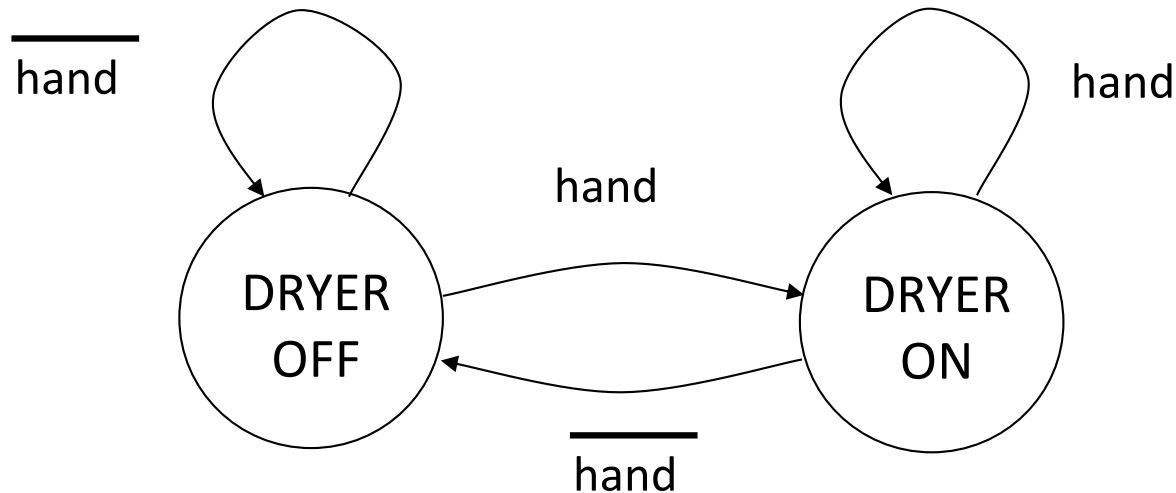
# What is a Computer?

- Theory of computation:
  - What are the fundamental capabilities and limitations of computers?

  - Complexity theory: what makes some problems computationally hard and others easy?
  - Computability theory: what makes some problems solvable and others unsolvable?
  - Automata theory: definitions and properties of mathematical models of computation

# Finite Automata

- Also known as finite state automata (FSA) or finite state machine (FSM)

- A simple mathematical model of a computer

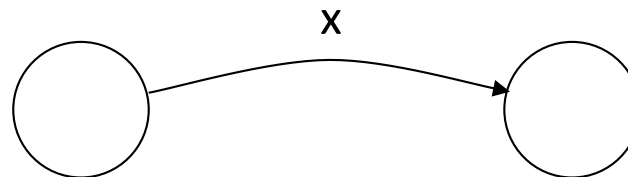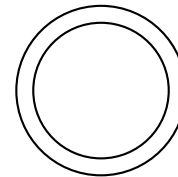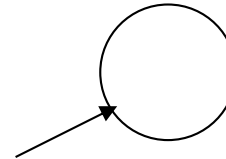- Applications: hardware design, compiler design, text processing

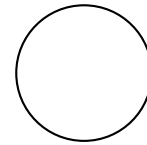# A First Example

- Touch-less hand-dryer

# Finite Automata State Graphs

- ## A state

- ## The start state

- ## An accepting state

- ## A transition

x

# Finite Automata

- Transition:  $s_1 \to^x s_2$
  - In state $s_1$ on input "x" go to state  $s_2$

- At end of input
  - If in accepting state => *accept*
  - Else => *reject*

- If no transition possible => reject

# Language of a FA

- Language of finite automaton M: set of all strings accepted by M

- Example:

( letter | digit )

letter

S

A

- Which of the following are in the language?
  - x, tmp2, 123, a?, 2apples

- A language is called a <u>regular language</u> if it is recognized by some finite automaton

# Example

- What is the language of this FA?

# Regular Expressions

- Regular expressions (regex) are used to describe regular languages

- Arithmetic expression example: (8+2)*3

- Regular expression example: (\+|-)?[0-9]+

# Example

- What is the language of this FA?
- Regular expression: (\+|-)?[0-9]+

# Three Equivalent Representations

Regular
expressions

Each
can
describe
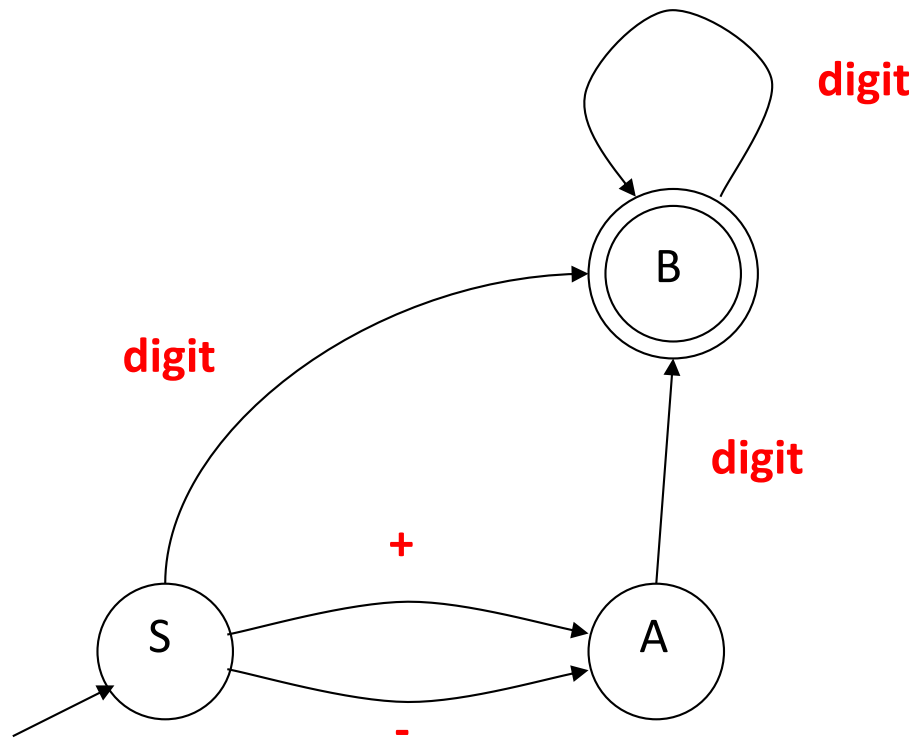the others

Finite
automata

Regular
languages

# Regex Rules

- Goal: match patterns
- String of characters matches the same string

| | |
|---|---|
| woodchuck | 'how much wood does a <u>woodchuck</u> chuck?' |
| e | 'you ar<u>e</u> a programm<u>e</u>r' |
| 206 | 'INFO<u>206</u> consists of <u>206</u>A and <u>206</u>B.' |
| ! | 'Keep it to yourself<u>!</u>' |

- . Wildcard matches any character at that position

| | |
|---|---|
| p.nt | '<u>pant</u>, <u>pint</u>, paint, print' |

# Regex Rules

- **? Zero or one occurrences of the preceding character/regex**
  woodchucks?      'how much wood does a <u>woodchuck</u> chuck?'
  behaviou?r      '<u>behaviour</u> is the British spelling of <u>behavior</u>'

- **\* Zero or more occurrences of the preceding character/regex**
  baa\*      ba, baa, baaa, baaaa …
  ba\*      b, ba, baa, baaa, baaaa …
  [ab]\*      $\varepsilon$, a, b, ab, ba, baaa, aaabbb, …
  [0-9][0-9]\*      any positive integer, or zero
  cat.\*cat      A string where 'cat' appears twice anywhere

- **+ One or more occurrences of the preceding character/regex**
  ba+      ba, baa, baaa, baaaa …

- **{n} Exactly n occurrences of the preceding character/regex**
  ba{3}      baaa

# Regex Rules

- ## * is greedy:

  <.*>      <u><a href="index.html">Home</a></u>

- ## Lazy (non-greedy) quantifier:

  <.*?>          <u><a href="index.html"</u>>Home<u></a></u>


  Similarly, +? is the non-greedy quantifier for +, and ?? is the non-greedy quantifier for ?

# Regex Rules

- ## [ ] Disjunction (Union)

  | | |
  |---|---|
  | [wW]ood | 'how much <u>wood</u> does a <u>Wood</u>chuck chuck?' |
  | [aeiou]* | '<u>yo</u>u <u>are</u> <u>a</u> pr<u>o</u>gr<u>a</u>mm<u>e</u>r' |
  | [A-Za-z0-9] | (any letter or digit) |
  | [A-Za-z]* | (any letter sequence) |

- ## | Disjunction (Union)

  | | |
  |---|---|
  | (cats?\|dogs?)+ | 'It is raining <u>cats</u> and a <u>dog</u>.' |

- ## ( ) Grouping

  | | |
  |---|---|
  | (gupp(y\|ies))* | 'His <u>guppy</u> is the king of <u>guppies</u>.' |

# Regex Rules

- **^  $  \b Anchors (start/end of input string; word boundary)**

  | | |
  |---|---|
  | ^The | 'The cat in the hat.' |
  | ^The end\.$ | 'The end.' |
  | ^The .* end\.$ | 'The bitter end.' |
  | (the)* | 'I saw him the other day.' |
  | (\bthe\b)* | 'I saw him the other day.' |

- **Special rule: when ^ is FIRST WITHIN BRACKETS it means NOT**

  | | |
  |---|---|
  | [^A-Z]* | (anything **not** an upper case letter) |

# Regex Rules

- **\ Escape characters**

  \.                'The + and \ characters are missing<u>.</u>'

  \+               'The <u>+</u> and \ characters are missing.'

  \\               'The + and <u>\</u> characters are missing.'

  and so on

# Operator Precedence

| Operator | Precedence |
|---|---|
| () | highest |
| * + ? {} | |
| sequences, anchors | |
| \| | lowest |

- What is the difference?
  - [a-z][a-z]|[0-9]*
  - [a-z]([a-z]|[0-9])*