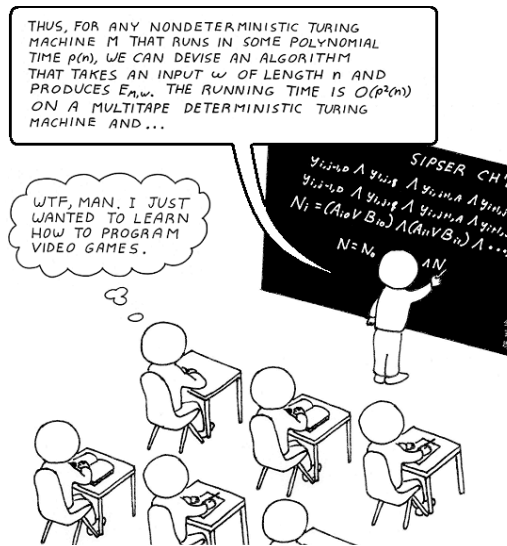


# Introduction to Theoretical CS



## Fundamental questions:

- Q. What can a computer do?
- Q. What can a computer do with limited resources?

## General approach.

- Don't talk about specific machines or problems.
- Consider minimal abstract machines.
- Consider general classes of problems.

1

## Why Learn Theory?

### In theory ...

- Deeper understanding of what is a computer and computing.
- Foundation of all modern computers.
- Pure science.
- Philosophical implications.

### In practice ...

- Web search: theory of pattern matching.
- Sequential circuits: theory of finite state automata.
- Compilers: theory of context free grammars.
- Cryptography: theory of computational complexity.
- Data compression: theory of information.

*"In theory there is no difference between theory and practice. In practice there is." – Yogi Berra*

## Regular Expressions

3

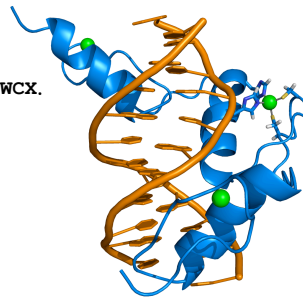
## Describing a Pattern

**PROSITE.** Huge database of protein families and domains.

**Q.** How to describe a protein motif?

**Ex.** [signature of the C<sub>2</sub>H<sub>2</sub>-type zinc finger domain]

1. C
2. Between 2 and 4 amino acids.
3. C
4. 3 more amino acids.
5. One of the following amino acids: LIVMFYWCX.
6. 8 more amino acids.
7. H
8. Between 3 and 5 more amino acids.
9. H



CAASC<sup>C</sup>GGPYACGGWAGY<sup>H</sup>HAGWH

5

## Pattern Matching Applications

**Test if a string matches some pattern.**

- Process natural language.
- Scan for virus signatures.
- Access information in digital libraries.
- Search-and-replace in a word processors.
- Filter text (spam, NetNanny, ads, Carnivore, malware).
- Validate data-entry fields (dates, email, URL, credit card).
- Search for markers in human genome using PROSITE patterns.

**Parse text files.**

- Compile a Java program.
- Crawl and index the Web.
- Read in data stored in TOY input file format.
- Automatically create Java documentation from Javadoc comments.

6

## Regular Expressions: Basic Operations

**Regular expression.** Notation to specify a set of strings.

operation	regular expression	matches	does not match
concatenation	<code>aabaab</code>	aabaab	every other string
wildcard	<code>.u.u.u.</code>	cumulus jugulum	succubus tumultuous
union	<code>aa   baab</code>	aa baab	every other string
closure	<code>ab*a</code>	aa abbba	ab ababa
parentheses	<code>a(a b)aab</code>	aaaab abaab	every other string
	<code>(ab)*a</code>	a ababababa	aa abbba

7

## Regular Expressions: Examples

**Regular expression.** Notation is surprisingly expressive.

regular expression	matches	does not match
<code>.*spb.*</code> <i>contains the trigraph spb</i>	raspberry crispbread	subspace subspecies
<code>a*   (a*ba*ba*ba*)*</code> <i>multiple of three b's</i>	bbb aaa bbbaababbaa	b bb baabbaa
<code>.*0....</code> <i>fifth to last digit is 0</i>	1000234 98701234	11111111 403982772
<code>gcg(cgg agg)*ctg</code> <i>fragile X syndrome indicator</i>	gcgctg gcgcgctg gcgcgaggctg	gcgcg cgcgcgctg gcgcgctg

8

## Generalized Regular Expressions

Regular expressions are a standard programmer's tool.

- Built in to Java, Perl, Unix, Python, ....
- Additional operations typically added for convenience.
  - Ex 1: `[a-e]+` is shorthand for `(a|b|c|d|e)(a|b|c|d|e)*`.
  - Ex 2: `\s` is shorthand for "any whitespace character" (space, tab, ...).

operation	regular expression	matches	does not match
one or more	<code>a(bc)+de</code>	abcde abcbcde	ade bcde
character class	<code>[A-Za-z][a-z]*</code>	lowercase Capitalized	camelCase 4illegal
exactly k	<code>[0-9]{5}-[0-9]{4}</code>	08540-1321 19072-5541	111111111 166-54-1111
negation	<code>[^aeiou]{6}</code>	rhythm	decade

9

## TEQ on REs 1

Q. Consider the RE

`a*bb(ab|ba)*`

Which of the following strings match (is in the set it describes)?

- abb
- abba
- aaba
- bbbaab
- cbb
- bbababbab

10

## TEQ on REs 2

Q. Give an RE that describes the following set of strings:

- characters are **A, C, T** or **G**
- starts with **ATG**
- length is a multiple of 3
- ends with **TAG, TAA, or TTG**

11

## Describing a Pattern

[PROSITE](#). Huge database of protein families and domains.

Q. How to describe a protein motif?

Ex. [signature of the  $C_2H_2$ -type zinc finger domain]

1. C
2. Between 2 and 4 amino acids.
3. C
4. 3 more amino acids.
5. One of the following amino acids: `LIVMFY`
6. 8 more amino acids.
7. H
8. Between 3 and 5 more amino acids.
9. H

A. `C.{2,4}C...[LIVMFYWC].{8}H.{3,5}H`



**CAASC**GGPY**ACGGWAGY**HAG**WH**

12

## REs in Java

public class String (Java's String library)

<code>boolean matches(String re)</code>	<i>does this string match the given regular expression?</i>
<code>String replaceAll(String re, String str)</code>	<i>replace all occurrences of regular expression with the replacement string</i>
<code>int indexOf(String r, int from)</code>	<i>return the index of the first occurrence of the string r after the index from</i>
<code>String[] split(String re)</code>	<i>split the string around matches of the given regular expression</i>

```
String re = C.{2,4}C...[LIVMFYWC].{8}H.{3,5}H;
String input = CAASC GG PYACGG AAGYHAGAH;
boolean test = input.matches(re);
```

is the input string in the set described by the RE?

13

## REs in Java

Validity checking. Is input in the set described by the re?

```
public class Validate
{
    public static void main(String[] args) {
        String re = args[0];
        String input = args[1];
        StdOut.println(input.matches(re));
    }
}
```

powerful string library method

```
% java Validate "C.{2,4}C...[LIVMFYWC].{8}H.{3,5}H" CAASC GG PYACGG AAGYHAGAH
true
% java Validate "[$_A-Za-z][$_A-Za-z0-9]*" ident123
true
% java Validate "[a-z]+@[a-z]+\.(edu|com)" wayne@cs.princeton.edu
true
```

*C<sub>2</sub>H<sub>2</sub> type zinc finger domain*  
*legal Java identifier*  
*valid email address (simplified)*  
*need quotes to "escape" the shell*

14

## REs in Java

public class String (Java's String library)

<code>boolean matches(String re)</code>	<i>does this string match the given regular expression?</i>
<code>String replaceAll(String re, String str)</code>	<i>replace all occurrences of regular expression with the replacement string</i>
<code>int indexOf(String r, int from)</code>	<i>return the index of the first occurrence of the string r after the index from</i>
<code>String[] split(String re)</code>	<i>split the string around matches of the given regular expression</i>

```
String s = StdIn.readAll();
s = s.replaceAll("\\s+", " ");
```

replace each sequence of at least one whitespace character with a single space

RE that matches any sequence of whitespace characters (at least 1).

Extra \ distinguishes from the string \s+

15

## REs in Java

public class String (Java's String library)

<code>boolean matches(String re)</code>	<i>does this string match the given regular expression?</i>
<code>String replaceAll(String re, String str)</code>	<i>replace all occurrences of regular expression with the replacement string</i>
<code>int indexOf(String r, int from)</code>	<i>return the index of the first occurrence of the string r after the index from</i>
<code>String[] split(String re)</code>	<i>split the string around matches of the given regular expression</i>

```
String s = StdIn.readAll();
String[] words = s.split("\\s+");
```

create an array of the words in StdIn

16

# DFAs

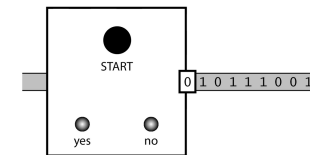
## Solving the Pattern Match Problem

Regular expressions are a concise way to describe patterns.

- How would you implement the method `matches()` ?
- Hardware: build a deterministic finite state automaton (DFA).
- Software: simulate a DFA.

DFA: simple machine that solves a pattern match problem.

- Different machine for each pattern.
- Accepts or rejects string specified on input tape.
- Focus on `true` or `false` questions for simplicity.

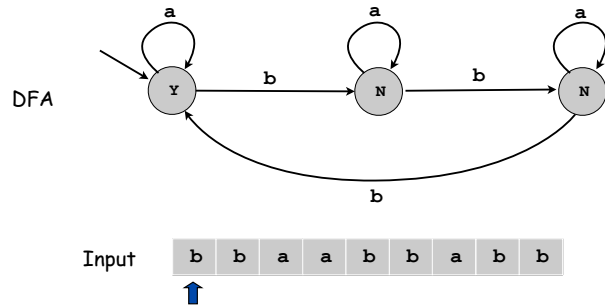


18

## Deterministic Finite State Automaton (DFA)

Simple machine with N states.

- Begin in start state.
- Read first input symbol.
- Move to new state, depending on current state and input symbol.
- Repeat until last input symbol read.
- Accept input string if last state is labeled Y.



19

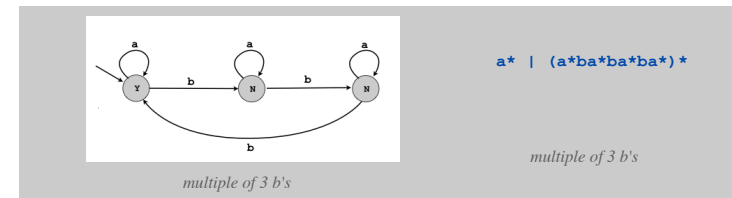
## DFA and RE Duality

RE. Concise way to describe a set of strings.

DFA. Machine to recognize whether a given string is in a given set.

Duality.

- For any DFA, there exists a RE that describes the same set of strings.
- For any RE, there exists a DFA that recognizes the same set.



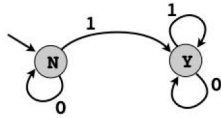
Practical consequence of duality proof: to match RE

- build DFA
- simulate DFA on input string.

20

### TEQ on DFAs 1

Q. Consider this DFA:



Which of the following sets of strings does it recognize?

- a. Bitstrings with at least one 1
- b. Bitstrings with an equal number of occurrences of 01 and 10
- c. Bitstrings with more 1s than 0s
- d. Bitstrings with an equal number of occurrences of 0 and 1
- e. Bitstrings that end in 1

21

### Implementing a Pattern Matcher

**Problem.** Given a RE, create program that tests whether given input is in set of strings described.

**Step 1.** Build the DFA.

- A compiler!
- See COS 226 or COS 320.

**Step 2.** Simulate it with given input.

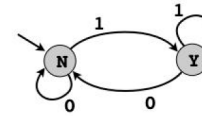
```

State state = start;
while (!StdIn.isEmpty())
{
    char c = StdIn.readChar();
    state = state.next(c);
}
StdOut.println(state.accept());
  
```

23

### TEQ on DFAs 2

Q. Consider this DFA:



Which of the following sets of strings does it recognize?

- a. Bitstrings with at least one 1
- b. Bitstrings with an equal number of occurrences of 01 and 10
- c. Bitstrings with more 1s than 0s
- d. Bitstrings with an equal number of occurrences of 0 and 1
- e. Bitstrings that end in 1

22

### Application: Harvester

Harvest information from input stream.

- Harvest patterns from DNA.

```

% java Harvester "gcg(cgg|agg)*ctg" chromosomeX.txt
gcgcgggcgggcgggcgggctg
gcgctg
gcgctg
gcgcgggcgggcgggcgggcgggctg
  
```

- Harvest email addresses from web for spam campaign.

```

% java Harvester "[a-z]+@[a-z]+\.(edu|com)" http://www.princeton.edu/~cos126
rs@cs.princeton.edu
maia@cs.princeton.edu
doug@cs.princeton.edu
wayne@cs.princeton.edu
  
```

24

## Application: Harvester

### Harvest information from input stream.

- Use `Pattern` data type to compile regular expression to NFA.
- Use `Matcher` data type to simulate NFA.

equivalent, but more efficient representation of a DFA

```
import java.util.regex.Pattern;
import java.util.regex.Matcher;

public class Harvester
{
    public static void main(String[] args)
    {
        String re      = args[0];
        In in          = new In(args[1]);
        String input    = in.readAll();
        Pattern pattern = Pattern.compile(re);
        Matcher matcher = pattern.matcher(input);

        while (matcher.find())
            StdOut.println(matcher.group());
    }
}
```

Annotations in the code:

- `Pattern.compile(re)`: create NFA from RE
- `pattern.matcher(input)`: create NFA simulator
- `matcher.find()`: look for next match
- `matcher.group()`: the match most recently found

25

## Application: Parsing a Data File

### Ex: parsing an NCBI genome data file.

```
LOCUS AC146846 128142 bp DNA linear HTG 13-NOV-2003
DEFINITION Ornithorhynchus anatinus clone CLM1-393H9,
ACCESSION AC146846
VERSION AC146846.2 GI:38304214
KEYWORDS HTG; HTGS_PHASE2; HTGS_DRAFT.
SOURCE Ornithorhynchus anatinus (platypus)
ORIGIN
 1 tgtatttoat ttgacogtgc tgttttttc oggtttttca gtaoggtgtt agggagocac
 61 gtgattctgt ttgttttatg ctgocgaata gctgctgat gaatctctgc atagacagct
121 gcgcagggga gaaatgacca gtttggatg acaaaatgta gaaaagctgt ttcttcataa
...
128101 gaaatggga cccocacgct aatgtacagc ttctttgatg tg
//
```

Annotations in the image:

- header info (points to the first line)
- line numbers (points to '1', '61', '121', '128101')
- spaces (points to spaces between words)
- comments (points to '// a comment')

Goal. Extract the data as a single `actg` string.

26

## Application: Parsing a Data File

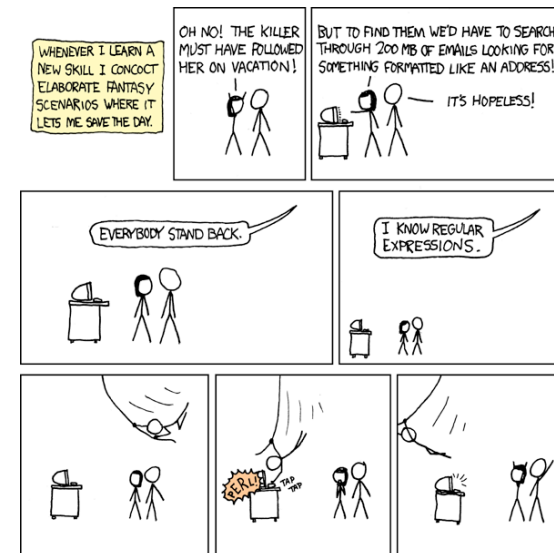
```
import java.util.regex.Pattern;
import java.util.regex.Matcher;

public class ParseNCBI
{
    public static void main(String[] args)
    {
        String re = "[ ]*[0-9]+([actg ]*).*";
        Pattern pattern = Pattern.compile(re);
        In in = new In(args[0]);
        String data = "";
        while (!in.isEmpty())
        {
            String line = in.readLine();
            Matcher matcher = pattern.matcher(line);
            if (matcher.find())
                data += matcher.group(1).replaceAll(" ", "");
        }
        System.out.println(data);
    }
}
```

```
LOCUS AC146846 128142 bp DNA linear HTG 13-NOV-2003
DEFINITION Ornithorhynchus anatinus clone CLM1-393H9,
ACCESSION AC146846
VERSION AC146846.2 GI:38304214
KEYWORDS HTG; HTGS_PHASE2; HTGS_DRAFT.
SOURCE Ornithorhynchus anatinus (platypus)
ORIGIN
 1 tgtatttoat ttgacogtgc tgttttttc oggtttttca gtaoggtgtt agggagocac
 61 gtgattctgt ttgttttatg ctgocgaata gctgctgat gaatctctgc atagacagct // a comment
121 gcgcagggga gaaatgacca gtttggatg acaaaatgta gaaaagctgt ttcttcataa
...
128101 gaaatggga cccocacgct aatgtacagc ttctttgatg tg
//
```

27

## Regular Expressions



<http://xkcd.com/208/>

28

## Summary

### Programmer.

- Regular expressions are a powerful pattern matching tool.
- Implement regular expressions with finite state machines.

### Theoretician.

- RE is a compact description of a set of strings.
- DFA is an abstract machine that solves RE pattern match problem.

**You.** Practical application of core CS principles.

29

## Basic Questions

**Q.** Are there patterns that **cannot** be described by any RE?

**A.** Yes.

- Bit strings with equal number of 0s and 1s.
- Strings that represent legal REs.
- Decimal strings that represent prime numbers.
- DNA strings that are Watson-Crick complemented palindromes.

30

## Basic Questions

**Q.** Are there languages that cannot be recognized by any **DFA**?

**A.** Yes.

- Bit strings with equal number of 0s and 1s.
- Strings that represent legal REs.
- Decimal strings that represent prime numbers.
- DNA strings that are Watson-Crick complemented palindromes.

31

## Basic Questions

**Q.** Are there languages that cannot be recognized by any **DFA**?

**A.** Yes: Bit strings with equal number of 0s and 1s.

### Proof sketch.

- Suppose that you have such a DFA, with N states.
- Give it N+1 0s followed by N+1 1s.
- Some state is revisited.
- Delete substring between visits.
- DFA recognizes that string, too.
- It does not have equal number of 0s and 1s.
- Contradiction.
- No such DFA exists.

0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1
0	1	3	5	6	8	7	3	5	.	.	.								
0	0	0	0	1	1	1	1	1	1	1	1	1	1	1					
0	1	3	5	.	.	.													

32



## Basic Questions

Q. Are there languages that cannot be recognized by any DFA?

A. Yes.

- Bit strings with equal number of 0s and 1s.
- Strings that represent legal REs.
- Decimal strings that represent prime numbers.
- DNA strings that are Watson-Crick complemented palindromes.

Fundamental problem: DFA lacks memory.

33

## Basic Questions

Q. Are there machines that are more powerful than a DFA?

A. Yes.

A 1-stack DFA can recognize

- Bit strings with equal number of 0s and 1s.
- Legal REs.
- Watson-Crick complemented palindromes.

34

## Basic Questions

Q. Are there machines that are more powerful than a 1-stack DFA?

A. Yes.

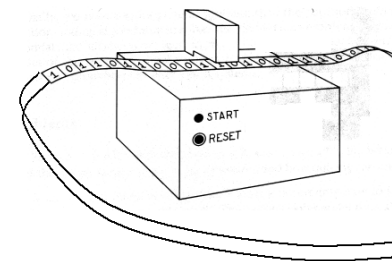
A 2-stack DFA can recognize

- Prime numbers.
- Legal Java Programs.

## Basic Questions

Q. Are there machines that are more powerful than a 2-stack DFA?

A. No! Not even a supercomputer!



2-stack DFAs are equivalent to Turing machines [stay tuned].

35

36