

17. Introduction to Theoretical CS

<http://introcscs.princeton.edu>

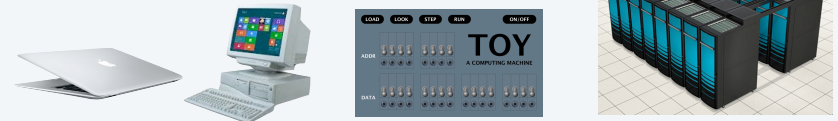
Introduction to theoretical computer science

Fundamental questions

- What can a computer do?
- 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.



Surprising outcome. Sweeping and relevant statements about *all* computers.

2

Why study theory?

In theory...

- Deeper understanding of computation.
- 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

3

17. Introduction to Theoretical CS

- Regular expressions
- DFAs
- Applications
- Limitations

Pattern matching

Pattern matching problem. Is a given string a member of a given set of strings?

Example 1 (from genomics)

A **nucleic acid** is represented by one of the letters a, c, t, or g.

A **genome** is a string of nucleic acids.

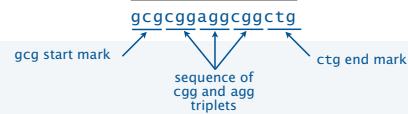
A **Fragile X Syndrome pattern** is a genome having an occurrence of gcg, followed by any number of cgg or agg triplets, followed by ctg.

Note. The number of triplets correlates with Fragile X Syndrome, a common cause of mental retardation.

Q. Does this genome contain a such a pattern?

g c g g c g t g t g t g c g a g a g a g t g g g t t t a a a g c t g g c g c g g a g g c g g c t g g c g c g g a g g c t g

A. Yes.



5

Pattern matching

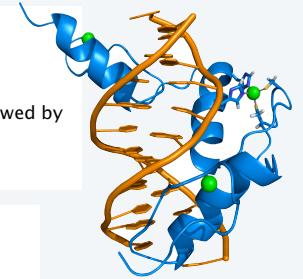
Example 2 (from computational biochemistry)

An **amino acid** is represented by one of the characters CAVLIMCRKH DENQSTYFWP.

A **protein** is a string of amino acids.

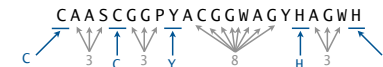
A **C₂H₂-type zinc finger domain signature** is

- C followed by 2, 3, or 4 amino acids, followed by
- C followed by 3 amino acids, followed by
- L, I, V, M, F, Y, W, C, or X followed by 8 amino acids, followed by
- H followed by 3, 4, or 5 amino acids, followed by
- H.



Q. Is this protein in the C₂H₂-type zinc finger domain?

A. Yes.



6

Pattern matching

Example 3 (from commercial computing)

An **e-mail address** is

- A sequence of letters, followed by
- the character "@", followed by
- the character ".", followed by a sequence of letters, followed by
- [any number of occurrences of the previous pattern]
- "edu" or "com" (others omitted for brevity).

Q. Which of the following are e-mail addresses?

	A.
rs@cs.princeton.edu	✓
not an e-mail address	✗
wayne@cs.princeton.edu	✓
eve@airport	✗
rs123@princeton.edu	✗

Ooops, need to fix description →

Challenge. Develop a precise description of the set of strings that are legal e-mail addresses.

7

Regular expressions

A **regular expression** (RE) is a notation for specifying sets of strings.

An RE is

- A sequence of letters or "."
- The *union* of two REs
- The *closure* of an RE (any number of occurrences)
- May be delimited by ().

operation	example RE	matches (IN the set)	does not match (NOT in the set)
concatenation	aabaab	aabaab	every other string
wildcard	.u.u.u.	cumulus jugulum	succubus tumultuous
union	aa baab	aa baab	every other string
closure	ab*a	aa abbba	ab ababa
parentheses	a(a b)aab	aaaaab abaab	every other string
	(ab)*a	a ababababa	aa abbba

8

More examples of regular expressions

The notation is surprisingly expressive.

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

9

Generalized regular expressions

Additional operations further extend the utility of REs.

operation	example RE	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	11111111 166-54-1111
negation	<code>[^aeiou]{6}</code>	rhythm	decade
white space	<code>\s</code>	any whitespace char (space, tab, newline...)	every other character

Note. These operations are all *shorthand*.
They are very useful but not essential.

RE: `(a|b|c|d|e)(a|b|c|d|e)*`
shorthand: `(a-e)+`

10

Example of describing a pattern with a generalized RE

A C₂H₂-type zinc finger domain signature is

- C followed by 2, 3, or 4 amino acids, followed by
- C followed by 3 amino acids, followed by
- L, I, V, M, F, Y, W, C, or X followed by 8 amino acids, followed by
- H followed by 3, 4, or 5 amino acids, followed by
- H.



Q. Give a generalized RE for all such signatures.

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

"Wildcard" matches any of the letters
CAVLIMCRKHDEHQSTYFW



11

Example of a real-world RE application: PROSITE

12

Another example of describing a pattern with a generalized RE

An e-mail address is

- A sequence of letters, followed by
- the character "@", followed by
- the character ".", followed by a sequence of letters, followed by
- [any number of occurrences of the previous pattern]
- "edu" or "com" (others omitted for brevity).

Q. Give a generalized RE for e-mail addresses.

A. $[a-z]^+@([a-z]+\.)^+(edu|com)$

Exercise. Extend to handle rs123@princeton.edu, more suffixes such as .org, and any other extensions you can think of.

Next. Determining whether a given string matches a given RE.

13

Self-assessment 1 on REs

Q. Which of the following strings match the RE $a^*bb(ab|ba)^*$?

↑
is in the set
it describes

1. abb
2. aaba
3. abba
4. bbbaab
5. cbb
6. bbababbab

14

Self-assessment 2 on REs

Q. Give an RE for *genes*

- Characters are a, c, t or g.
- Starts with atg (a *start codon*).
- Length is a multiple of 3.
- Ends with tag, taa, or ttg (a *stop codon*).



15

17. Introduction to Theoretical CS

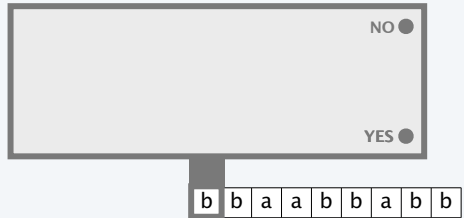
- Regular expressions
- DFAs
- Applications
- Limitations

Deterministic finite state automata (DFA)

A **DFA** is an abstract machine that solves a pattern matching problem.

- A string is specified on an input tape (no limit on its length).
- The DFA reads each character on input tape once, moving left to right.
- The DFA lights "YES" if it *recognizes* the string, "NO" otherwise.

Each DFA defines a *set* of strings (all the strings that it recognizes).

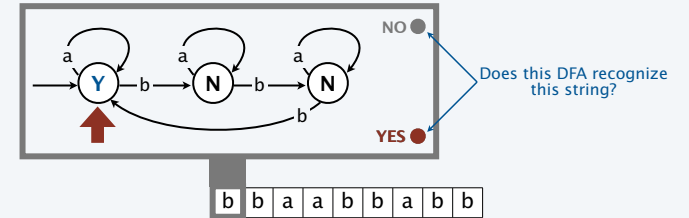


17

Deterministic finite state automata details and example

A **DFA** is an abstract machine with a finite number *states*, each labeled Y or N, and *transitions* between states, each labelled with a symbol. One state is the *start* state.

- Begin in the *start* state (denoted by an arrow from nowhere).
- Read an input symbol and move to the indicated state.
- Repeat until the last input symbol has been read.
- Turn on the "YES" or "NO" light according to the label on the current state.

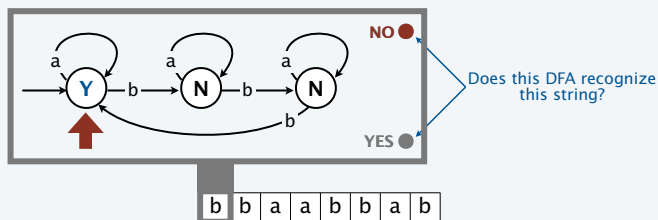


18

Deterministic finite state automata details and example

A **DFA** is an abstract machine with a finite number *states*, each labeled Y or N and *transitions* between states, each labelled with a symbol. One state is the *start* state.

- Begin in the *start* state.
- Read an input symbol and move to the indicated state.
- Repeat until the last input symbol has been read.
- Turn on the "YES" or "NO" light according to the label on the current state.



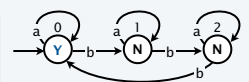
19

Simulating the operation of a DFA

```
public class DFA
{
    private int state;
    private int start;
    private String[] action;
    private ST<Character, Integer>[] next;
    public DFA(In in)
    { /* Fill in data structures */ }
    public String simulate(String input)
    {
        state = start;
        for (int i = 0; i < input.length(); i++)
            state = next[state].get(input.charAt(i));
        return action[state];
    }
    public static void main(String[] args)
    {
        DFA dfa = new DFA(new In(args[0]));
        while (!StdIn.isEmpty())
        {
            input = StdIn.readString();
            StdOut.println(dfa.simulate(input));
        }
    }
}
```

symbol table to map
chars a, b, ... to next
state 0, 1, ...

action[]	next[]	a	b
0 Yes	0 0 1		
1 No	1 1 2		
2 No	2 2 0		

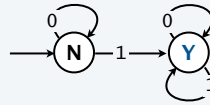


```
# states → % more b3.txt
alphabet → 3
start state → ab
Yes 0 1
No 1 2
No 2 0
% java DFA b3.txt
bababa
Yes
bb
No
abbababababaaa
Yes
abbabababba
No
```

20

Self-assessment 1 on DFAs

Q. Which of the following strings does this DFA accept?

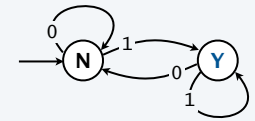


1. Bitstrings that end in 1
2. Bitstrings with an equal number of occurrences of 01 and 10
3. Bitstrings with more 1s than 0s
4. Bitstrings with an equal number of occurrences of 0 and 1
5. Bitstrings with at least one 1

21

Self-assessment 2 on DFAs

Q. Which of the following strings does this DFA accept?



1. Bitstrings with at least one 1
2. Bitstrings with an equal number of occurrences of 01 and 10
3. Bitstrings with more 1s than 0s
4. Bitstrings with an equal number of occurrences of 0 and 1
5. Bitstrings that end in 1

22

Kleene's theorem

Two ways to define a set of strings

- Regular expressions (REs).
- Deterministic finite automata (DFAs).

Remarkable fact. DFAs and REs are *equivalent*.

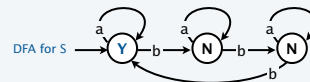
Equivalence theorem (Kleene)

Given any RE, there exists a DFA that accepts the same set of strings.
 Given any DFA, there exists an RE that matches the same set of strings.

Consequence: A way to solve the RE pattern matching problem

- Build the DFA corresponding to the given RE.
- Simulate the operation of the DFA.

$S =$ the set of ab strings where the number of occurrences of b is a multiple of 3



RE for S $a^* | (a^*ba^*ba^*ba^*)^*$



Steven Kleene
1909–1994

23

17. Introduction to Theoretical CS

- Regular expressions
- DFAs
- **Applications**
- Limitations

GREP: a solution to the RE pattern matching problem

An algorithm for the RE pattern matching problem?

- Build the DFA corresponding to the given RE.
- Simulate the operation of the DFA.

Practical difficulty: The DFA might have *exponentially* many states.

A more efficient algorithm: use Nondeterministic Finite Automata (NFA)

- Build the NFA corresponding to the given RE.
- Simulate the operation of the NFA.



Interested in details? Take a course in algorithms.

"GREP" (Generalized Regular Expression Pattern matcher).

- Developed by Ken Thompson, who designed and implemented Unix.
- Indispensable programming tool for decades.
- Found in most development environments, including Java.



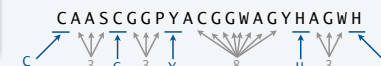
Ken Thompson
1983 Turing Award

REs in Java

Java's String class implements GREP.

public class String	
...	
boolean matches(String re)	<i>does this string match the given RE?</i>
...	

```
String re = "C.{2,4}C...[LIVMFYWC].{8}H.{3,5}H";
String zincFinger = "CAASCGGPYACGGWAGYHAGAH";
boolean test = zincFinger.matches(re);
true!
```



Java RE client example: Validation

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

Does a given string match a given RE?

- Take RE from command line.
- Take strings from StdIn.

```
% java Validate "C.{2,4}C...[LIVMFYWC].{8}H.{3,5}H"
CAASCGGPYACGGWAGYHAGAH
true
CAASCGGPYACGGWAGYHGAH
false
% java Validate "[$_A-Za-z][$_A-Za-z0-9]*"
ident123
true
123ident
false
% java Validate "[a-z]+@[a-z]+\.(edu|com)"
wayne@cs.princeton.edu
true
eve@airport
false
```

Annotations: "need quotes to 'escape' the shell" (pointing to the first RE), "C₂H₂ type zinc finger domain" (pointing to the first RE), "legal Java identifier" (pointing to the second RE), "valid email address (simplified)" (pointing to the third RE).

Applications

- Scientific research.
- Compilers and interpreters.
- Internet commerce.
- ...

Beyond matching

Java's String class contains other useful RE-related methods.

- RE search and replace
- RE delimited parsing

public class String	
...	
String replaceAll(String re, String to)	<i>replace all occurrences of substrings matching RE with to</i>
String[] split(String re)	<i>split the string around matches of the given RE</i>
...	

Examples using the RE "\\s+" (matches one or more whitespace characters).
 Tricky notation (typical in string processing): \ signals "special character" so "\\\" means "\" and "\\s" means "\s"

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

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

Create an array of the words in StdIn (basis for StdIn.readAllStrings() method)

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


Java pattern matcher real-world example: 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]);
        while (in.hasNext Line())
        {
            String line = in.readLine();
            Matcher matcher = pattern.matcher(line);
            if (matcher.find())
                StdOut.print(matcher.group(1).replaceAll(" ", ""));
        }
        StdOut.println();
    }
}
```

remove the spaces

```
% java ParseNCBI platypus.txt
tgtatttcatttgacctgctgtttttccgg
tttttcagtaacggtgtaaggagcaacgtgatt
ctgtttgtttatgctgccaatagctgctgga
tgaatcctgcatagacagctgccgaggagaa
aatgaccagtttgtgatgacaaaatgtaggaaa
gctgtttttcataa...
```

33

Applications of REs

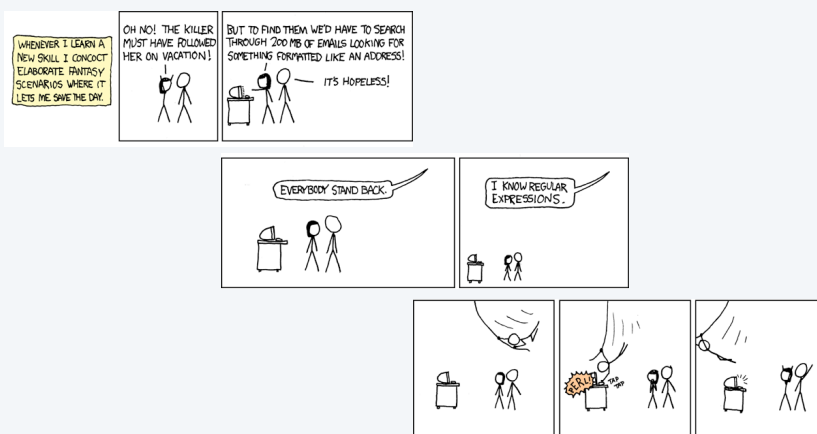
Pattern matching and beyond.

- Compile a Java program.
- Scan for virus signatures.
- Crawl and index the Web.
- Process natural language.
- Access information in digital libraries.
- Search-and-replace in a word processors.
- Process NCBI and other scientific data files.
- 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.
- Automatically create Java documentation from Javadoc comments.

GREP and related facilities are built in to Java, Unix shell, PERL, Python ...

virtually every computing environment

34



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

35

17. Introduction to Theoretical CS

- Regular expressions
- DFAs
- Applications
- Limitations

Summary

Programmers

- Regular expressions are a powerful pattern matching tool.
- Equivalent DFA/NFA paradigm facilitates implementation.
- Combination greatly facilitates real-world string data processing.



Theoreticians

- REs provide compact descriptions of sets of strings.
- DFAs are abstract machines with equivalent descriptive power.
- Are there languages and machines with more descriptive power?



You

- CS core principles provide useful tools that you can exploit now.
- REs and DFAs provide an introduction to theoretical CS.



37

Basic questions

Q. Are there sets of strings that cannot be described by *any* RE?

A. Yes.

- Bitstrings 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.
- ...

Q. Are there sets of strings that cannot be described 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.
- ...

The *same* question, by Kleene's theorem

38

A limit on the power of REs and DFAs

Proposition. There exists a set of strings that cannot be described by any RE or DFA.

Proof sketch. No DFA can recognize the set of bitstrings with equal number of 0s and 1s.

- Assume that you have such a DFA, with N states.
- It recognizes the string with $N + 1$ 0s followed by $N + 1$ 1s.
- Some state is revisited when recognizing that string.
- Delete the substring between visits.
- DFA recognizes that string, too.
- It does not have equal number of 0s and 1s.
- *Proof by contradiction:* the assumption that such a DFA exists must be false.

Ex. $N = 10$

0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1
0	3	5	9	8	7	5
0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
0	3	5

39

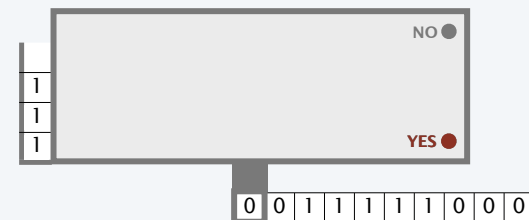
Another basic question

Q. Are there abstract machines that are more powerful than DFAs?

A. Yes. A 1-stack DFA can recognize

- Bitstrings with equal number of 0s and 1s.
- Strings that represent legal REs.

Proof. [details omitted]



40

Yet another basic question

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

- A. Yes. A 2-stack DFA can recognize
- Decimal strings that represent prime numbers.
 - Strings that represent legal Java programs.
 - ...

[stay tuned for next lecture]



41

One last basic question

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

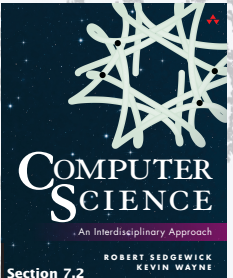
- A. No! Not even a roomful of supercomputers (!!!)

[stay tuned for next lecture]



42

COMPUTER SCIENCE
SEDGEWICK / WAYNE



17. Introduction to Theoretical CS

<http://introc.cs.princeton.edu>