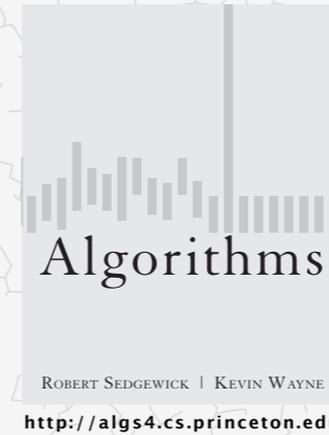


## 5.4 REGULAR EXPRESSIONS

- ▶ *regular expressions*
- ▶ *REs and NFAs*
- ▶ *NFA simulation*
- ▶ *NFA construction*
- ▶ *applications*



## 5.4 REGULAR EXPRESSIONS

- ▶ *regular expressions*
- ▶ *REs and NFAs*
- ▶ *NFA simulation*
- ▶ *NFA construction*
- ▶ *applications*

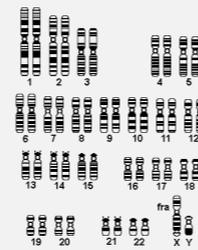
### Pattern matching

**Substring search.** Find a single string in text.

**Pattern matching.** Find one of a **specified set** of strings in text.

**Ex.** [genomics]

- Fragile X syndrome is a common cause of mental retardation.
- A human's genome is a string.
- It contains triplet repeats of CGG or AGG, bracketed by GCG at the beginning and CTG at the end.
- Number of repeats is variable and is correlated to syndrome.



**pattern** GCG(CGG|AGG)\*CTG

**text** GCGCGTGTGTGCGAGAGAGTGGGTTTAAAGCTGCGCGGAGGGCGGCTGGCGGGAGGCTG

### Syntax highlighting

```

/*****
 * Compilation: javac NFA.java
 * Execution: java NFA regexp text
 * Dependencies: Stack.java Bag.java Digraph.java DirectedDFS.java
 *
 * % java NFA "(A*B|AC)D" AAAABD
 * true
 *
 * % java NFA "(A*B|AC)D" AAAAC
 * false
 *****/

public class NFA
{
    private Digraph G;           // digraph of epsilon transitions
    private String regexp;      // regular expression
    private int M;              // number of characters in regular expression

    // Create the NFA for the given RE
    public NFA(String regexp)
    {
        this.regexp = regexp;
        M = regexp.length();
        Stack<Integer> ops = new Stack<Integer>();
        G = new Digraph(M+1);
        ...
    }
}
    
```

GNU source-highlight 3.1.4

input	output
Ada	HTML
Asm	XHTML
Applescript	LATEX
Awk	MediaWiki
Bat	ODF
Bib	TEXINFO
Bison	ANSI
C/C++	DocBook
C#	
Cobol	
Caml	
Changelog	
Css	
D	
Erlang	
Flex	
Fortran	
GLSL	
Haskell	
Html	
Java	
Javalog	
Javascript	
Latex	
Lisp	
Lua	
:	

## Google code search

**Search public source code**

Search via regular expression, e.g. `^java/.*\.java$`

Search Options	In Search Box
Package <input type="text"/>	package:linux-2.6
Language <input type="text" value="Any language"/>	lang:c++
File Path <input type="text"/>	file:(code [^or]g)search
Class <input type="text"/>	class:HashMap
Function <input type="text"/>	function:toString
License <input type="text" value="Any license"/>	license:mozilla
Case Sensitive <input type="text" value="No"/>	case:yes

<http://code.google.com/p/chromium/source/search>

5

## Pattern matching: applications

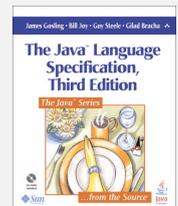
Test if a string matches some pattern.

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



Parse text files.

- Compile a Java program.
- Crawl and index the Web.
- Read in data stored in ad hoc input file format.
- Create Java documentation from Javadoc comments.



6

## Regular expressions

A **regular expression** is a notation to specify a set of strings.

↑  
possibly infinite

operation	order	example RE	matches	does not match
<b>concatenation</b>	3	AABAAB	AABAAB	<i>every other string</i>
<b>or</b>	4	AA   BAAB	AA BAAB	<i>every other string</i>
<b>closure</b>	2	AB*A	AA ABBBBBBBBA	AB ABABA
<b>parentheses</b>	1	A(A B)AAB	AAAAB ABAAB	<i>every other string</i>
		(AB)*A	A ABABABABABA	AA ABBA

7

## Regular expression shortcuts

Additional operations are often added for convenience.

operation	example RE	matches	does not match
<b>wildcard</b>	.U.U.U.	CUMULUS JUGULUM	SUCCUBUS TUMULTUOUS
<b>character class</b>	[A-Za-z][a-z]*	word Capitalized	camelCase 4illegal
<b>at least 1</b>	A(BC)+DE	ABCDE ABCBCDE	ADE BCDE
<b>exactly k</b>	[0-9]{5}-[0-9]{4}	08540-1321 19072-5541	11111111 166-54-111

Ex. `[A-E]+` is shorthand for `(A|B|C|D|E)(A|B|C|D|E)*`

8

## Regular expression examples

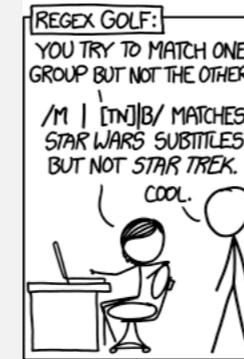
RE notation is surprisingly expressive.

regular expression	matches	does not match
<code>.*SPB.*</code> <i>(substring search)</i>	RASPBERRY CRISPBREAD	SUBSPACE SUBSPECIES
<code>[0-9]{3}-[0-9]{2}-[0-9]{4}</code> <i>(U. S. Social Security numbers)</i>	166-11-4433 166-45-1111	11-55555555 8675309
<code>[a-z]+@[a-z]+\.(edu com)</code> <i>(simplified email addresses)</i>	wayne@princeton.edu rs@princeton.edu	spam@nowhere
<code>[\$_A-Za-z][\$_A-Za-z0-9]*</code> <i>(Java identifiers)</i>	ident3 PatternMatcher	3a ident#3

REs play a well-understood role in the theory of computation.

9

## Regular expression golf



<http://xkcd.com/1313>

yes	no
obama	romney
bush	mccain
clinton	kerry
reagan	gore
...	...
washington	clinton

Ex. Match elected presidents but not opponents (unless they later won).

RE. `bu|[rn]t|[coy]e|[mtg]a|j|i|so|[h]l|[ae]d|lev|sh|[l]nd|i|[po]o|ls`



10

## Illegally screening a job candidate

“ [First name]! and pre/2 [last name] w/7  
bush or gore or republican! or democrat! or charg!  
or accus! or criticiz! or blam! or defend! or iran contra  
or clinton or spotted owl or florida recount or sex!  
or controvers! or fraud! or investigat! or bankrupt!  
or layoff! or downsiz! or PNTR or NAFTA or outsourc!  
or indict! or enron or kerry or iraq or wmd! or arrest!  
or intox! or fired or racis! or intox! or slur!  
or controvers! or abortion! or gay! or homosexual!  
or gun! or firearm! ”

— LexisNexis search string used by Monica Goodling  
to illegally screen candidates for DOJ positions

<http://www.justice.gov/oig/special/s0807/final.pdf>



11

## Can the average web surfer learn to use REs?

Google. Supports \* for full word wildcard and | for union.



12



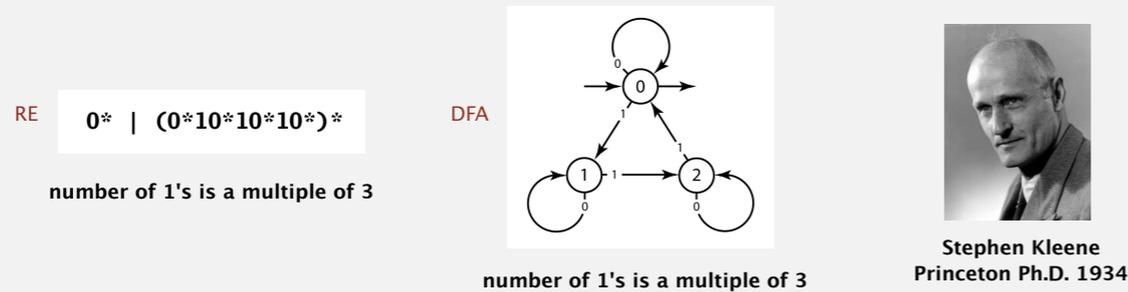
## Duality between REs and DFAs

**RE.** Concise way to describe a set of strings.

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

**Kleene's theorem.**

- 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 of strings.



17

## Pattern matching implementation: basic plan (first attempt)

Overview is the same as for KMP.

- No backup in text input stream.
- Linear-time guarantee.

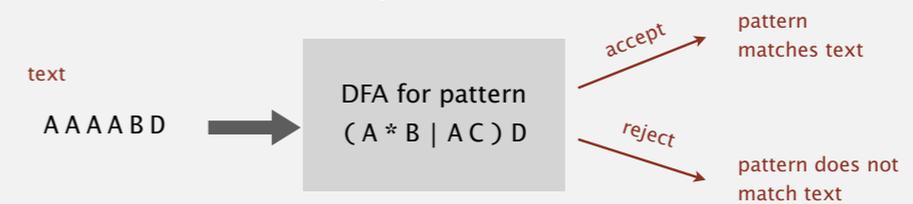


Ken Thompson  
Turing Award '83

**Underlying abstraction.** Deterministic finite state automata (DFA).

**Basic plan.** [apply Kleene's theorem]

- Build DFA from RE.
- Simulate DFA with text as input.



**Bad news.** Basic plan is infeasible (DFA may have exponential # of states).

18

## Pattern matching implementation: basic plan (revised)

Overview is similar to KMP.

- No backup in text input stream.
- **Quadratic-time guarantee** (linear-time typical).

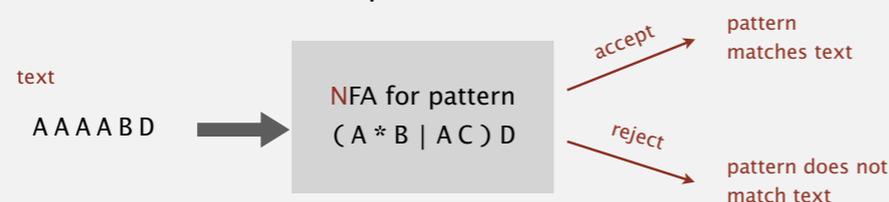


Ken Thompson  
Turing Award '83

**Underlying abstraction.** Nondeterministic finite state automata (NFA).

**Basic plan.** [apply Kleene's theorem]

- Build **NFA** from RE.
- Simulate **NFA** with text as input.



Q. What is an NFA?

19

## Nondeterministic finite-state automata

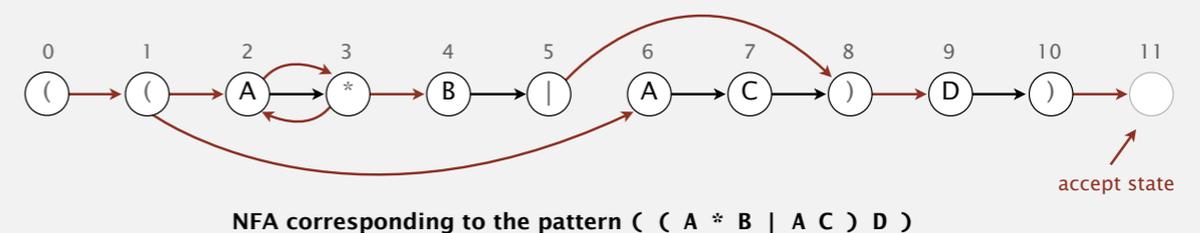
**Regular-expression-matching NFA.**

- We assume RE enclosed in parentheses.
- One state per RE character (start = 0, accept =  $M$ ).
- Red  $\epsilon$ -transition (change state, but don't scan text).
- Black match transition (change state and scan to next text char).
- Accept if **any** sequence of transitions ends in accept state.

← after scanning all text characters

**Nondeterminism.**

- One view: machine can guess the proper sequence of state transitions.
- Another view: sequence is a proof that the machine accepts the text.

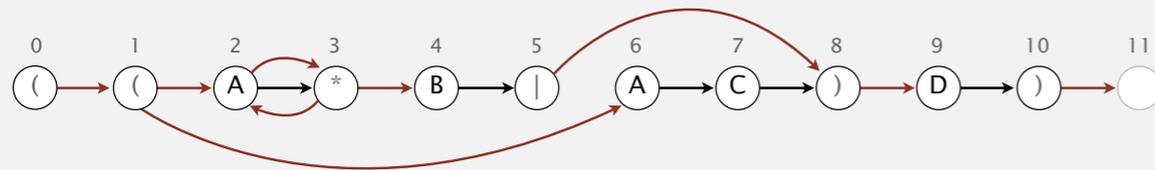
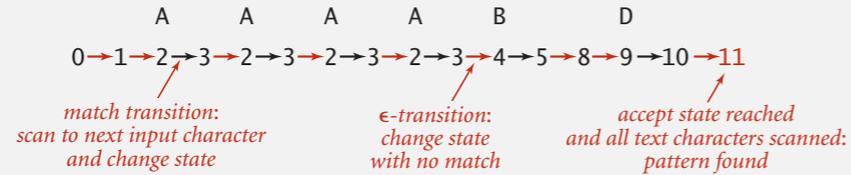


20

## Nondeterministic finite-state automata

Q. Is AAAABD matched by NFA?

A. Yes, because **some** sequence of legal transitions ends in state 11.



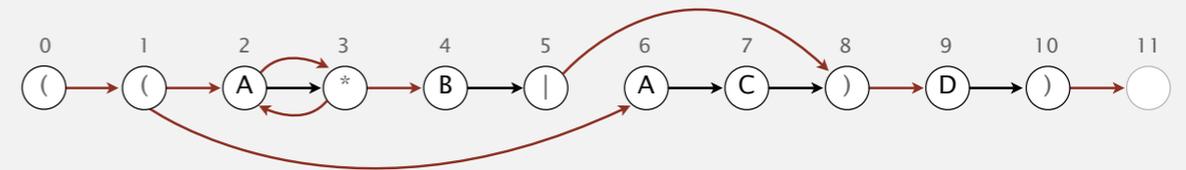
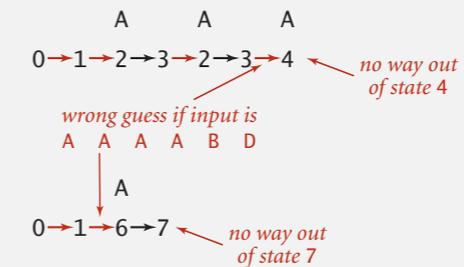
NFA corresponding to the pattern  $((A * B | A C) D)$

21

## Nondeterministic finite-state automata

Q. Is AAAABD matched by NFA?

A. Yes, because **some** sequence of legal transitions ends in state 11.  
[ even though some sequences end in wrong state or stall ]



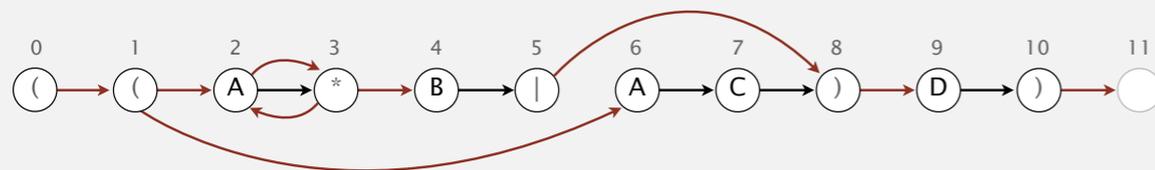
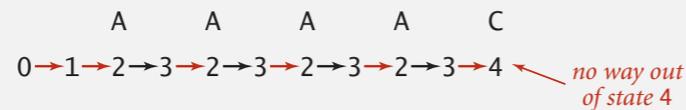
NFA corresponding to the pattern  $((A * B | A C) D)$

22

## Nondeterministic finite-state automata

Q. Is AAAC matched by NFA?

A. No, because **no** sequence of legal transitions ends in state 11.  
[ but need to argue about all possible sequences ]



NFA corresponding to the pattern  $((A * B | A C) D)$

23

## Nondeterminism

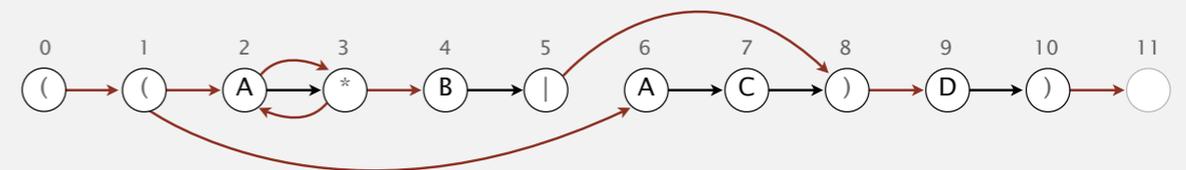
Q. How to determine whether a string is matched by an automaton?

DFA. Deterministic  $\Rightarrow$  easy because exactly one applicable transition.

NFA. Nondeterministic  $\Rightarrow$  can be several applicable transitions; need to select the right one!

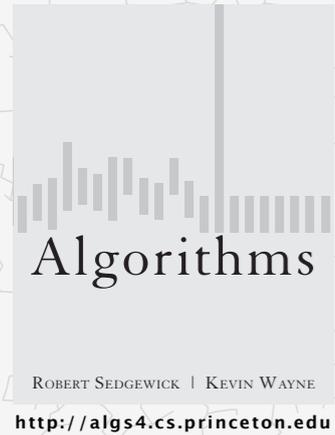
Q. How to simulate NFA?

A. Systematically consider **all** possible transition sequences. [stay tuned]



NFA corresponding to the pattern  $((A * B | A C) D)$

24



## 5.4 REGULAR EXPRESSIONS

- ▶ regular expressions
- ▶ REs and NFAs
- ▶ NFA simulation
- ▶ NFA construction
- ▶ applications

## NFA representation

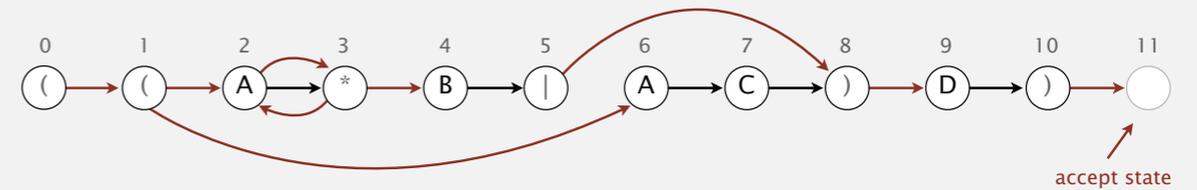
State names. Integers from 0 to  $M$ .

number of symbols in RE

Match-transitions. Keep regular expression in array `re[]`.

$\epsilon$ -transitions. Store in a digraph  $G$ .

$0 \rightarrow 1, 1 \rightarrow 2, 1 \rightarrow 6, 2 \rightarrow 3, 3 \rightarrow 2, 3 \rightarrow 4, 5 \rightarrow 8, 8 \rightarrow 9, 10 \rightarrow 11$

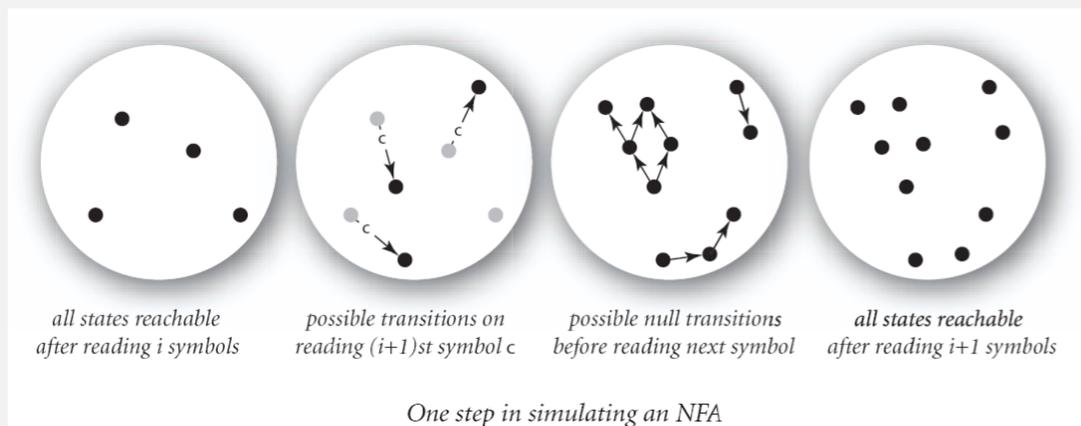


NFA corresponding to the pattern `(( A * B | A C ) D )`

## NFA simulation

Q. How to efficiently simulate an NFA?

A. Maintain set of **all** possible states that NFA could be in after reading in the first  $i$  text characters.



Q. How to perform reachability?

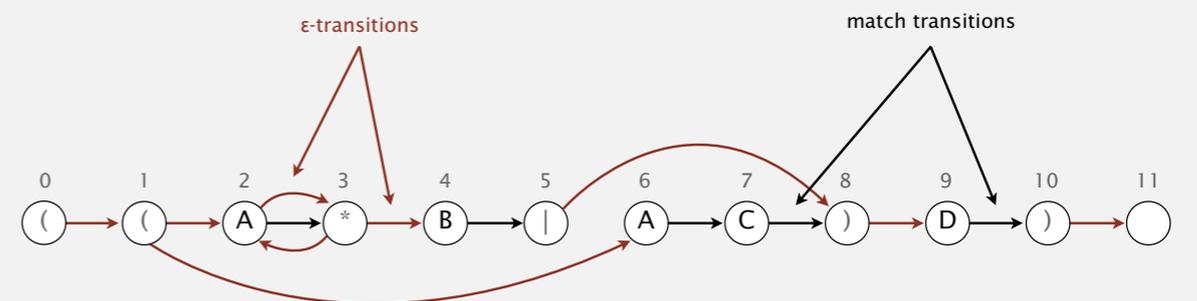
## NFA simulation demo

Goal. Check whether input matches pattern.



input 

A	A	B	D
---	---	---	---



NFA corresponding to the pattern `(( A * B | A C ) D )`

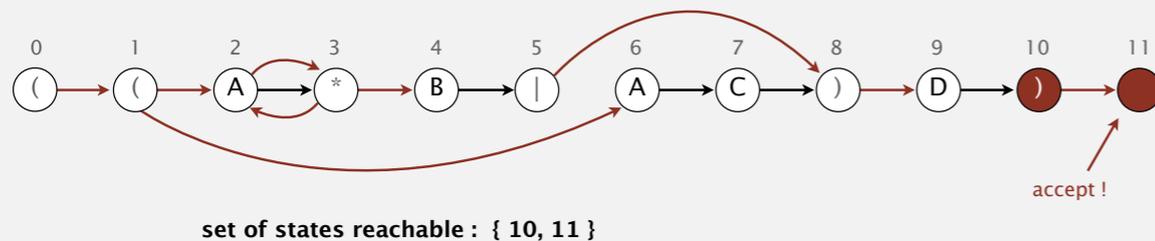
## NFA simulation demo

When no more input characters:

- Accept if any state reachable is an accept state.
- Reject otherwise.

input 

A	A	B	D
---	---	---	---



29

## Digraph reachability

**Digraph reachability.** Find all vertices reachable from a given source or set of vertices. recall Section 4.2

```
public class DirectedDFS
{
    DirectedDFS(Digraph G, int s)           find vertices reachable from s
    DirectedDFS(Digraph G, Iterable<Integer> s) find vertices reachable from sources
    boolean marked(int v)                  is v reachable from source(s)?
}
```

**Solution.** Run DFS from each source, without unmarking vertices.

**Performance.** Runs in time proportional to  $E + V$ .

30

## NFA simulation: Java implementation

```
public class NFA
{
    private char[] re;           // match transitions
    private Digraph G;          // epsilon transition digraph
    private int M;              // number of states

    public NFA(String regexp)
    {
        M = regexp.length();
        re = regexp.toCharArray();
        G = buildEpsilonTransitionDigraph();
    }

    public boolean recognizes(String txt)
    { /* see next slide */ }

    public Digraph buildEpsilonTransitionDigraph()
    { /* stay tuned */ }
}
```

← stay tuned (next segment)

31

## NFA simulation: Java implementation

```
public boolean recognizes(String txt)
{
    Bag<Integer> pc = new Bag<Integer>();
    DirectedDFS dfs = new DirectedDFS(G, 0);
    for (int v = 0; v < G.V(); v++)
        if (dfs.marked(v)) pc.add(v);

    for (int i = 0; i < txt.length(); i++)
    {
        Bag<Integer> states = new Bag<Integer>();
        for (int v : pc)
        {
            if (v == M) continue;
            if ((re[v] == txt.charAt(i)) || re[v] == '.')
                states.add(v+1);
        }

        dfs = new DirectedDFS(G, states);
        pc = new Bag<Integer>();
        for (int v = 0; v < G.V(); v++)
            if (dfs.marked(v)) pc.add(v);
    }

    for (int v : pc)
        if (v == M) return true;
    return false;
}
```

← states reachable from start by  $\epsilon$ -transitions

← set of states reachable after scanning past `txt.charAt(i)`

← not necessarily a match (RE needs to match full text)

← follow  $\epsilon$ -transitions

← accept if can end in state M

32

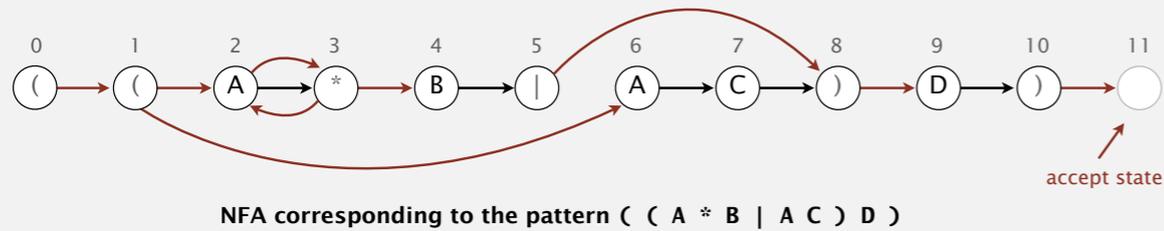


## NFA simulation: analysis

**Proposition.** Determining whether an  $N$ -character text is recognized by the NFA corresponding to an  $M$ -character pattern takes time proportional to  $MN$  in the worst case.

**Pf.** For each of the  $N$  text characters, we iterate through a set of states of size no more than  $M$  and run DFS on the graph of  $\epsilon$ -transitions.

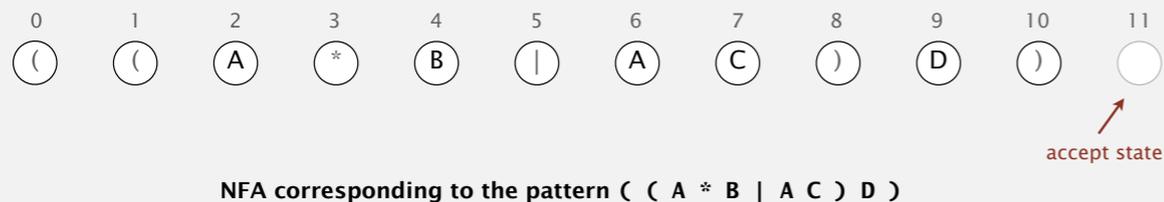
[The NFA construction we will consider ensures the number of edges  $\leq 3M$ .]



33

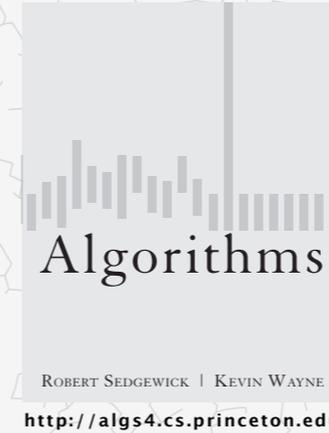
## Building an NFA corresponding to an RE

**States.** Include a state for each symbol in the RE, plus an accept state.



35

## 5.4 REGULAR EXPRESSIONS



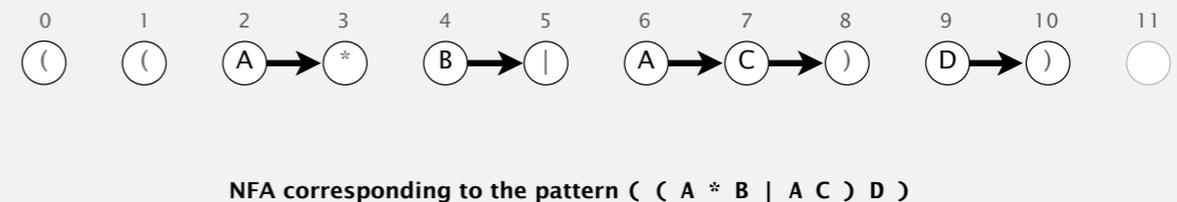
- ▶ regular expressions
- ▶ REs and NFAs
- ▶ NFA simulation
- ▶ NFA construction
- ▶ applications

## Building an NFA corresponding to an RE

**Concatenation.** Add match-transition edge from state corresponding to characters in the alphabet to next state.

**Alphabet.** A B C D

**Metacharacters.** ( ) . \* |



36

## Building an NFA corresponding to an RE

**Parentheses.** Add  $\epsilon$ -transition edge from parentheses to next state.



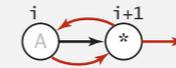
NFA corresponding to the pattern  $( ( A * B | A C ) D )$

37

## Building an NFA corresponding to an RE

**Closure.** Add three  $\epsilon$ -transition edges for each \* operator.

single-character closure



closure expression



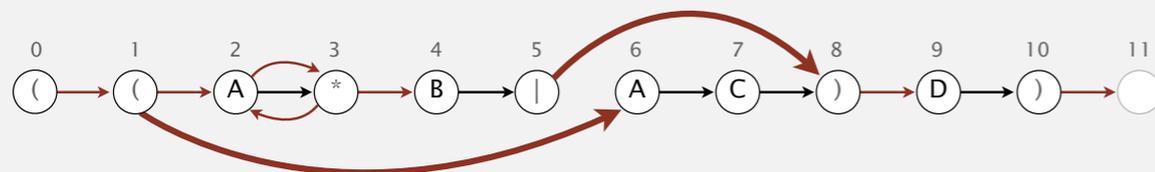
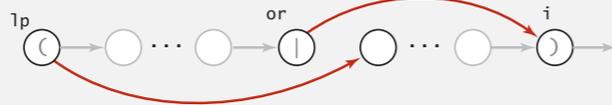
NFA corresponding to the pattern  $( ( A * B | A C ) D )$

38

## Building an NFA corresponding to an RE

**Or.** Add two  $\epsilon$ -transition edges for each | operator.

or expression



NFA corresponding to the pattern  $( ( A * B | A C ) D )$

39

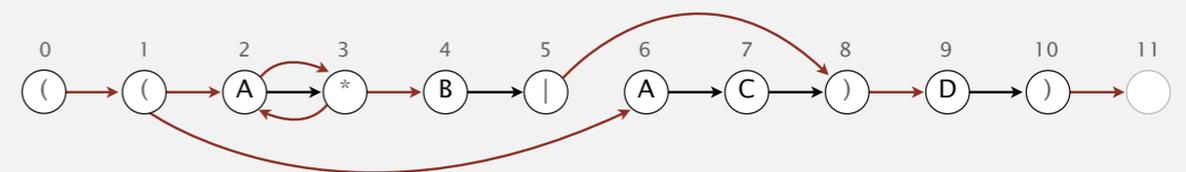
## NFA construction: implementation

**Goal.** Write a program to build the  $\epsilon$ -transition digraph.

**Challenges.** Remember left parentheses to implement closure and or; remember | to implement or.

**Solution.** Maintain a stack.

- ( symbol: push ( onto stack.
- | symbol: push | onto stack.
- ) symbol: pop corresponding ( and any intervening |; add  $\epsilon$ -transition edges for closure/or.



NFA corresponding to the pattern  $( ( A * B | A C ) D )$

40

## NFA construction demo



stack

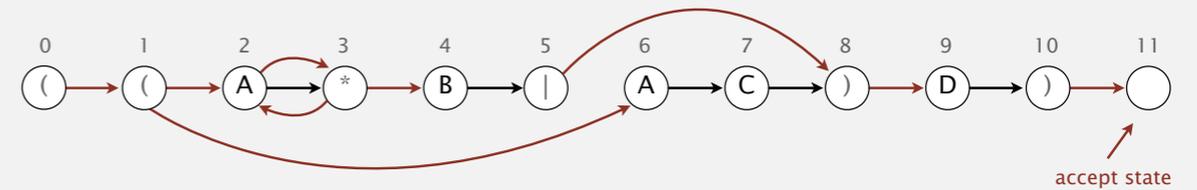
( ( A \* B | A C ) D )

41

## NFA construction demo



stack



NFA corresponding to the pattern ( ( A \* B | A C ) D )

42

## NFA construction: Java implementation

```
private Digraph buildEpsilonTransitionDigraph() {
    Digraph G = new Digraph(M+1);
    Stack<Integer> ops = new Stack<Integer>();
    for (int i = 0; i < M; i++) {
        int lp = i;

        if (re[i] == '(' || re[i] == '|') ops.push(i);

        else if (re[i] == ')') {
            int or = ops.pop();
            if (re[or] == '|') {
                lp = ops.pop();
                G.addEdge(lp, or+1);
                G.addEdge(or, i);
            }
            else lp = or;
        }

        if (i < M-1 && re[i+1] == '*') {
            G.addEdge(lp, i+1);
            G.addEdge(i+1, lp);
        }

        if (re[i] == '(' || re[i] == '*' || re[i] == ')')
            G.addEdge(i, i+1);
    }
    return G;
}
```

Annotations:

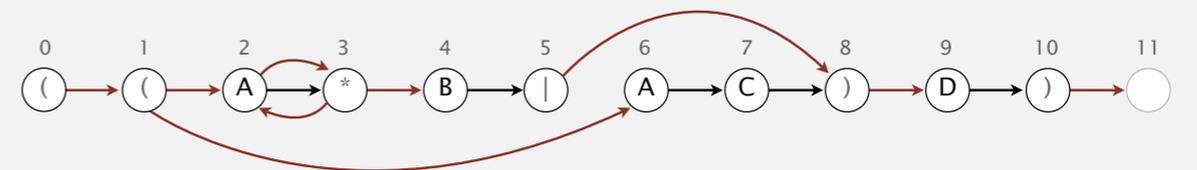
- left parentheses and |
- 2-way or
- closure (needs 1-character lookahead)
- metasymbols

43

## NFA construction: analysis

**Proposition.** Building the NFA corresponding to an  $M$ -character RE takes time and space proportional to  $M$ .

**Pf.** For each of the  $M$  characters in the RE, we add at most three  $\epsilon$ -transitions and execute at most two stack operations.



NFA corresponding to the pattern ( ( A \* B | A C ) D )

44



## Regular expressions in other languages

### Broadly applicable programmer's tool.

- Originated in Unix in the 1970s.
- Many languages support extended regular expressions.
- Built into grep, awk, emacs, Perl, PHP, Python, JavaScript, ...

```
% grep 'NEWLINE' *.*.java ← print all lines containing NEWLINE which occurs in any file with a .java extension
```

```
% egrep '^[qwertyuiop]*[zxcvbnm]*$' words.txt | egrep '.....'  
typewritten
```

### PERL. Practical Extraction and Report Language.

```
% perl -p -i -e 's|from|to|g' input.txt ← replace all occurrences of from with to in the file input.txt
```

```
% perl -n -e 'print if /^[A-Z][A-Za-z]*$/' words.txt ← print all words that start with uppercase letter  
↑  
do for each line
```

49

## Regular expressions in Java

### Validity checking. Does the input match the re?

**Java string library.** Use `input.matches(re)` for basic RE matching.

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

```
% java Validate "[$_A-Za-z][$_A-Za-z0-9]*" ident123 ← legal Java identifier  
true  
  
% java Validate "[a-z]+@[a-z]+\.(edu|com)" rs@cs.princeton.edu ← valid email address (simplified)  
true  
  
% java Validate "[0-9]{3}-[0-9]{2}-[0-9]{4}" 166-11-4433 ← Social Security number  
true
```

50

## Harvesting information

**Goal.** Print all substrings of input that match a RE.

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

↑  
harvest patterns from DNA  
  
↓  
harvest links from website

```
% java Harvester "http://(\w+\.)*(\w+)" http://www.cs.princeton.edu  
http://www.princeton.edu  
http://www.google.com  
http://www.cs.princeton.edu/news
```

51

## Harvesting information

RE pattern matching is implemented in Java's `java.util.regex.Pattern` and `java.util.regex.Matcher` classes.

```
import java.util.regex.Pattern;  
import java.util.regex.Matcher;  
  
public class Harvester  
{  
    public static void main(String[] args)  
    {  
        String regexp = args[0];  
        In in = new In(args[1]);  
        String input = in.readAll();  
        Pattern pattern = Pattern.compile(regexp);  
        Matcher matcher = pattern.matcher(input);  
        while (matcher.find())  
        {  
            StdOut.println(matcher.group());  
        }  
    }  
}
```

compile() creates a Pattern (NFA) from RE  
matcher() creates a Matcher (NFA simulator) from NFA and text  
find() looks for the next match  
group() returns the substring most recently found by find()

52

## Algorithmic complexity attacks

**Warning.** Typical implementations do **not** guarantee performance!

↖  
Unix grep, Java, Perl, Python

```
% java Validate "(a|aa)*b" aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaac 1.6 seconds
% java Validate "(a|aa)*b" aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaac 3.7 seconds
% java Validate "(a|aa)*b" aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaac 9.7 seconds
% java Validate "(a|aa)*b" aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaac 23.2 seconds
% java Validate "(a|aa)*b" aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaac 62.2 seconds
% java Validate "(a|aa)*b" aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaac 161.6 seconds
```

**SpamAssassin regular expression.**

```
% java RE "[a-z]+@[a-z]+([a-z\.\.]+\.)+[a-z]+" spammer@x.....
```

- Takes exponential time on pathological email addresses.
- Troublemaker can use such addresses to DOS a mail server.

53

## Not-so-regular expressions

**Back-references.**

- \1 notation matches subexpression that was matched earlier.
- Supported by typical RE implementations.

```
(.+)\1 // beriberi couscous
1?$|^(11+?)\1+ // 1111 111111 1111111111
```

**Some non-regular languages.**

- Strings of the form  $ww$  for some string  $w$ : beriberi.
- Unary strings with a composite number of 1s: 111111.
- Bitstrings with an equal number of 0s and 1s: 01110100.
- Watson-Crick complemented palindromes: atttcggaaat.

**Remark.** Pattern matching with back-references is intractable.

54

## Context

**Abstract machines, languages, and nondeterminism.**

- Basis of the theory of computation.
- Intensively studied since the 1930s.
- Basis of programming languages.

**Compiler.** A program that translates a program to machine code.

- KMP string  $\Rightarrow$  DFA.
- grep RE  $\Rightarrow$  NFA.
- javac Java language  $\Rightarrow$  Java byte code.

	KMP	grep	Java
pattern	string	RE	program
parser	unnecessary	check if legal	check if legal
compiler output	DFA	NFA	byte code
simulator	DFA simulator	NFA simulator	JVM

55

## Summary of pattern-matching algorithms

**Programmer.**

- Implement substring search via DFA simulation.
- Implement RE pattern matching via NFA simulation.



**Theoretician.**

- RE is a compact description of a set of strings.
- NFA is an abstract machine equivalent in power to RE.
- DFAs, NFAs, and REs have limitations.



**You.** Practical application of core computer science principles.

**Example of essential paradigm in computer science.**

- Build intermediate abstractions.
- Pick the right ones!
- Solve important practical problems.

56