





<http://algs4.cs.princeton.edu>

## 5.4 REGULAR EXPRESSIONS

---

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



<http://algs4.cs.princeton.edu>

## 5.4 REGULAR EXPRESSIONS

---

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

allinurl:spreadsheets.google.com\*formkey



**Web**

Images

Maps

Shopping

More ▾

Search tools

About 414,000 results (0.34 seconds)

[Response summary - \[ 4chan Users' Survey \] - Google Docs](#)

<https://spreadsheets.google.com/viewanalytics?formkey...>

Summary. 8106 responses. Note 1: Spreadsheet of Responses up to 6500 (7100+ are gamed) and 1-1700 have been processed for obvious trolls ...

[Welcome to Google Docs](#)

<https://spreadsheets.google.com/viewform?formkey...>

Upload your files from your desktop: It's easy to get started and it's free!  
Access anywhere: Edit and view your docs from any computer or smart phone.

[Optimal PID Settings For Your MWC v2.1 - Google Docs](#)

<https://spreadsheets.google.com/viewanalytics?formkey...>

Summary. 214 responses. Your RC Groups Name.  
Shikrapilotguy95hugyjhoexpPoint65abc123itbsjbBill ...

[My Little Pony: Friendship is Magic Survey - Google Docs](#)

<https://spreadsheets.google.com/viewform?formkey...>

This poll is intended for the audiences of My Little Pony: Friendship is magic.  
Please ONLY do this survey IF you have watched My Little Pony: Friendship  
Is ...

# My Little Pony: Friendship is Magic Survey

This poll is intended for the audiences of My Little Pony: Friendship is magic.

Please ONLY do this survey IF you have watched My Little Pony: Friendship Is Magic  
And Please do not submit multiple surveys. Thank you

\*\*Survey result can be found here

[https://spreadsheets.google.com/viewanalytics?  
hl=en&formkey=dDB5dy0yTjhlMnlKdEsxRGFzWTNLUIE6MQ](https://spreadsheets.google.com/viewanalytics?hl=en&formkey=dDB5dy0yTjhlMnlKdEsxRGFzWTNLUIE6MQ)

\* Required

**What is your sex? \***

- Male
- Female

**What is your age (range)? \***

- 5 - 11
- 12 - 14
- 15 - 17
- 18 - 24
- 25 - 30
- 31 - 40
- 41- 50
- 51 - 60
- 60 or above

**What is your exact age? (optional)**

---

**Who is your Favorite Character \***

- Twilight Sparkle
- Pinkie Pie
- Fluttershy
- Applejack
- Rainbow Dash
- Rarity
- Princess Celestia
- Derpy Hooves
- Spike
- None of the above/ others
- Luna

**Who is your Second Favorite Character \***

- Twilight Sparkle
- Pinkie Pie
- Fluttershy

Google docs Senate reconciliation whip count

File Edit View Insert Format Form Tools Help

10pt B Abc [Color Picker] [Background Color] [Text Color] [List Bullets] [List Numbers] [Sum]

	A	B	C	D	E	F	G
1	State	Senator	D.C. Phone #	Open to using reconciliation to finish health reform?	Sign Bennet letter on public option?	Call Status	Link to statement (if there is one)
2				Totals	Totals		
3			YES	34	20		
4			MAYBE	5	9		
5			NO	1	5		
6			??	19	25		
7							
8	State	Senator	D.C. Phone #	Open to using reconciliation to finish health reform?	Sign Bennet letter on public option?	Call Status	Link to statement (if there is one)
9	Alaska	Mark Begich	(202) 224-3004	??	??	Call made, awaiting response	
10	Arkansas	Mark Pryor	(202) 224 2353	MAYBE	??	Russ A.	<a href="http://blog.healthcareforamericamoving-towards-reconciliation-to-finish-health-reform/">http://blog.healthcareforamericamoving-towards-reconciliation-to-finish-health-reform/</a>
11	Arkansas	Blanche Lincoln		NO	NO	Done	
12	California	Barbara Boxer		YES	YES	done	<a href="http://whipcongress.com/">http://whipcongress.com/</a>
13	California	Diane Feinstein		YES	YES	Done	<a href="http://whipcongress.com/">http://whipcongress.com/</a>
14	Colorado	Michael Bennet		YES	YES	Done	<a href="http://whipcongress.com/">http://whipcongress.com/</a>
15	Colorado	Mark Udall	(202) 224 5941	??	??	Zapp and Jeff J.	
16	Connecticut	Chris Dodd	(202) 224 2823	??	??	Call made, awaiting response	
17	Connecticut	Joe Lieberman	(202) 224 4041	??	NO	Call made, awaiting response	
18	Delaware	Tom Carper	(202) 224 2441	YES	??	Dan S.	<a href="http://blog.healthcareforamericamoving-towards-reconciliation-to-finish-health-reform/">http://blog.healthcareforamericamoving-towards-reconciliation-to-finish-health-reform/</a>
19	Delaware	Ted Kaufman	(202) 224-5042	??	??	call made	
20	Florida	Bill Nelson	(202) 224-5274	??	??	Call placed, will hear tomorrow	
21	Hawaii	Daniel Akaka	(202) 224-6361	??	??	Left message, will follow up tomorrow	
22							

# 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);  
    }  
}
```

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	



# 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

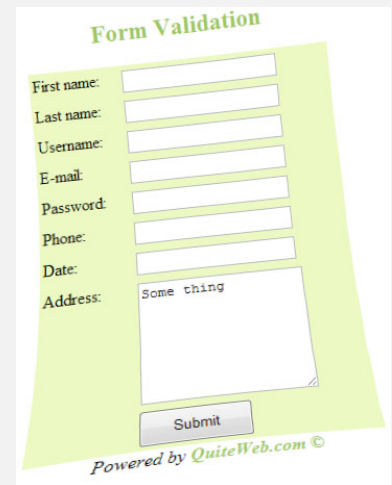
<http://code.google.com/p/chromium/source/search>

# Pattern matching: applications

---

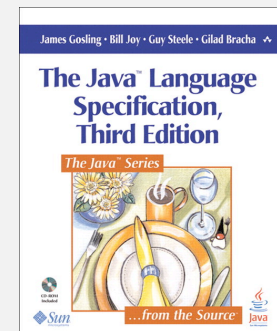
## Test if a string matches some pattern.

- 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.
- ...



# Regular expressions

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

↑  
possibly infinite

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

# Regular expression shortcuts

---

Additional operations are often added for convenience.

operation	example RE	matches	does not match
wildcard	.U.U.U.	CUMULUS JUGULUM	SUCCUBUS TUMULTUOUS
character class	[A-Za-z][a-z]*	word Capitalized	camelCase 4illegal
at least 1	A(BC)+DE	ABCDE ABCBCDE	ADE BCDE
exactly k	[0-9]{5}-[0-9]{4}	08540-1321 19072-5541	111111111 166-54-111

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

# 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.

# Regular expression examples

---

regular expression

(CAT|SUP\*)

CATU

CATUP

CATSUPPP

SU

[pollEv.com/jhug](http://pollEv.com/jhug)

text to **37607**

Q: How many of the strings match the regular expression?

- |      |          |      |          |
|------|----------|------|----------|
| A. 0 | [406437] | D. 3 | [406446] |
| B. 1 | [406438] | E. 4 | [406795] |
| C. 2 | [406439] |      |          |

# Regular expression examples

---

regular expression

(CAT|SUP\*)

CATU

CATUP

CATSUPPP

SU

×

×

×

[pollEv.com/jhug](http://pollEv.com/jhug)

text to **37607**

Q: How many of the strings match the regular expression?

B. 1 [406438]

# Regular expression examples

---

regular expression

`(A|(B|C*))*`

A

ABC

BBBCBCCB

[pollEv.com/jhug](http://pollEv.com/jhug)

text to **37607**

Q: How many of the strings match the regular expression?

A. 0 [424621]

B. 1 [424623]

C. 2 [424625]

D. 3 [424765]

E. 4 [424770]



# Regular expression examples

---

regular expression

`(A|(B|C*))*`



Q: How many of the strings match the regular expression?

D. 3



## Regular expression caveat

---

Writing a RE is like writing a program.

- Need to understand programming model.
- Can be easier to write than read.
- Can be difficult to debug.



*“ Some people, when confronted with a problem, think  
‘I know I’ll use regular expressions.’ Now they have  
two problems. ”*

*— Jamie Zawinski (flame war on alt.religion.emacs)*

**Bottom line.** REs are amazingly powerful and expressive,  
but using them in applications can be amazingly complex and error-prone.



<http://algs4.cs.princeton.edu>

## 5.4 REGULAR EXPRESSIONS

---

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

# 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.

See COS487 next Fall (requires COS340)

**Kleene's theorem.**

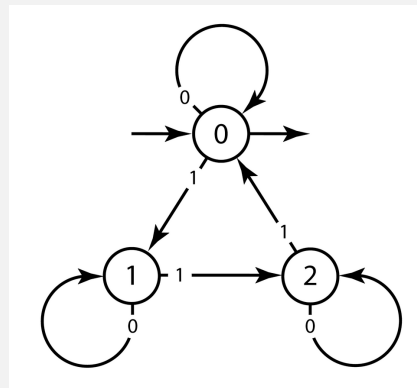
- For any DFA, there exists a RE that describes the same set of strings.
- For any RE, there exists a DFA that recognizes the same set of strings.

RE

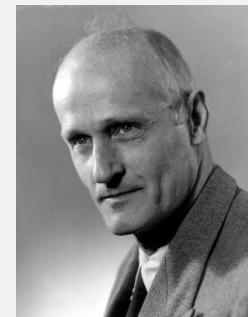
$0^* \mid (0^*10^*10^*10^*)^*$

number of 1's is a multiple of 3

DFA



number of 1's is a multiple of 3

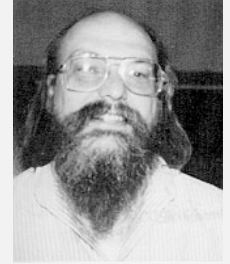


Stephen Kleene  
Princeton Ph.D. 1934

# Pattern matching implementation: basic plan (first attempt)

Overview is the same as for KMP.

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

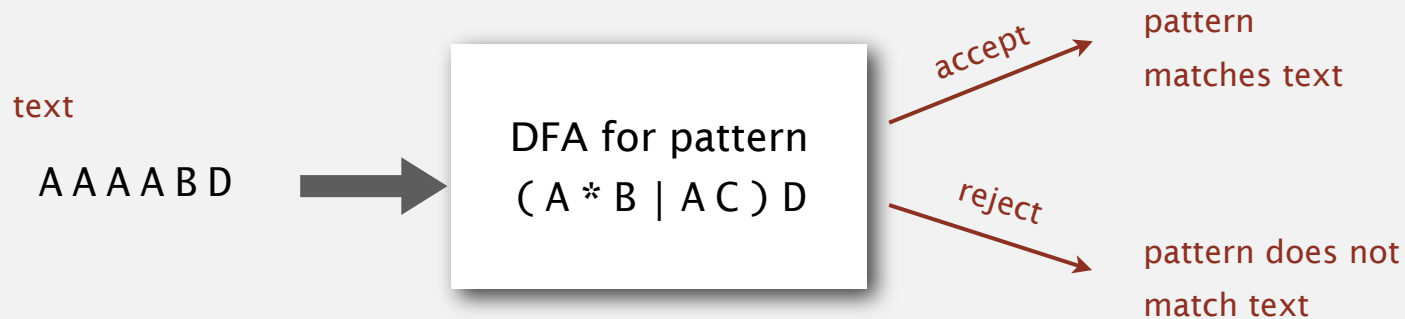


Ken Thompson  
Turing Award '83

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

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

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

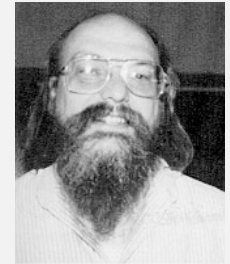


**Bad news.** Basic plan is infeasible (DFA may have exponential # of states).

# Pattern matching implementation: basic plan (revised)

Overview is similar to KMP.

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

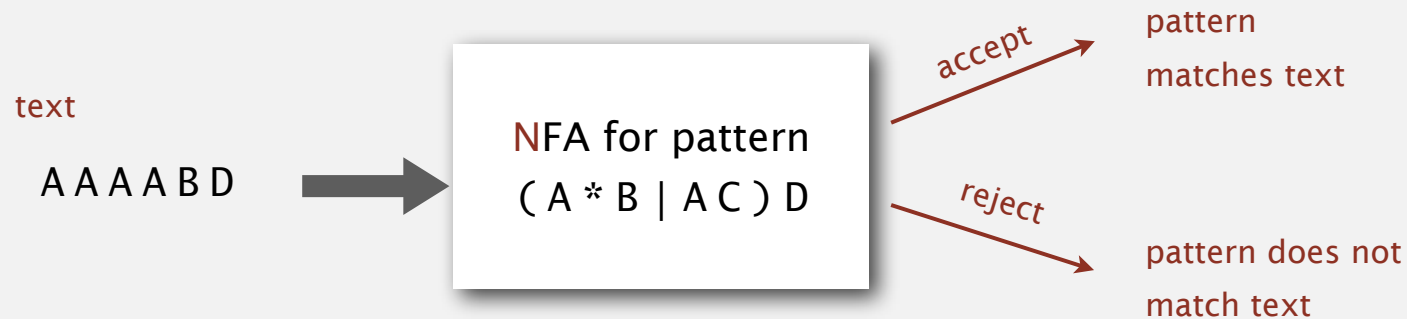


Ken Thompson  
Turing Award '83

Underlying abstraction. **N**ondeterministic 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?

# Nondeterministic finite-state automata

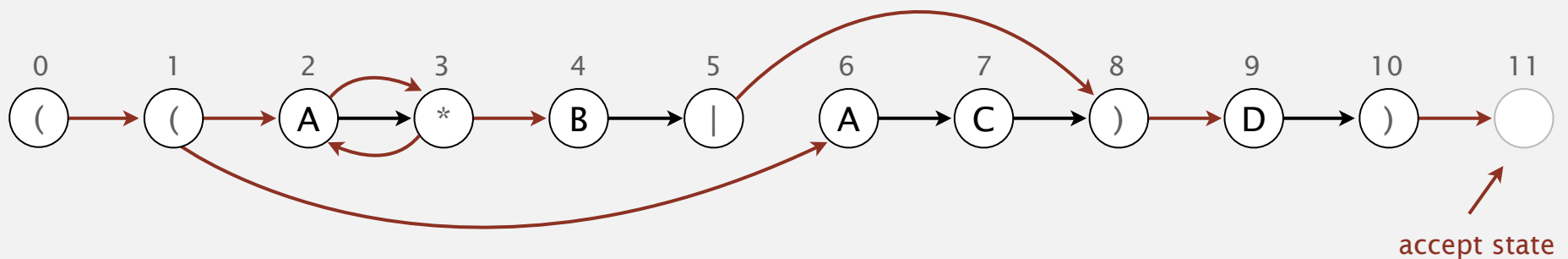
## Regular-expression-matching NFA.

- 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.

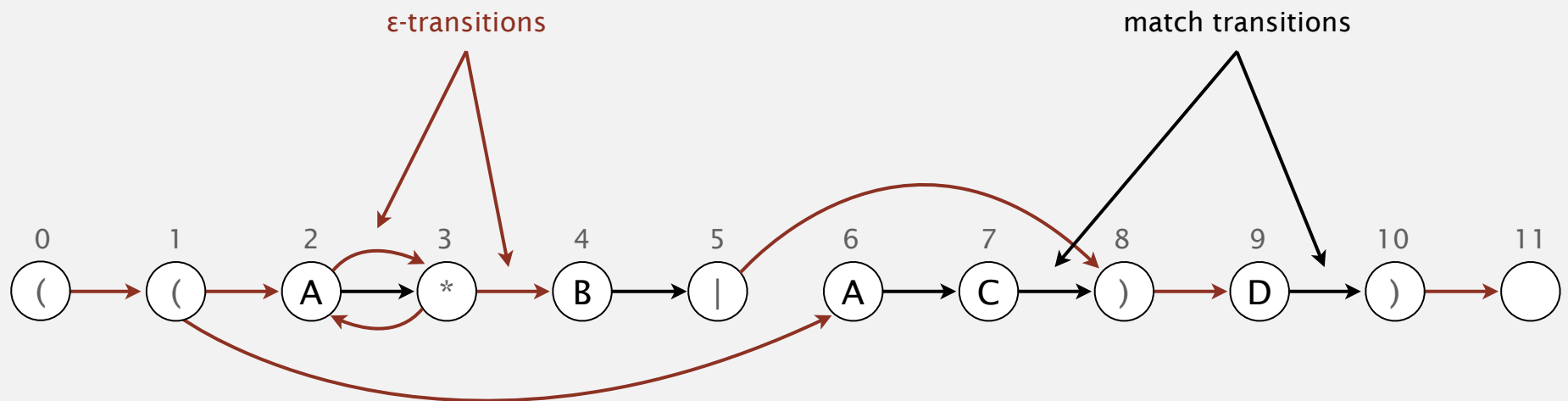
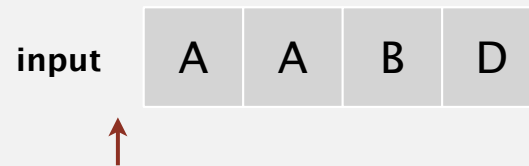


NFA corresponding to the pattern  $((A^*B|AC)D)$



# NFA simulation demo

Goal. Check whether input matches pattern.

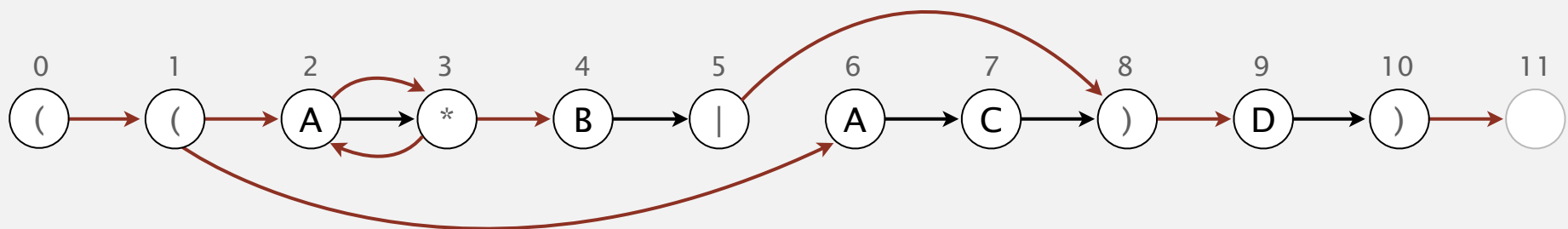
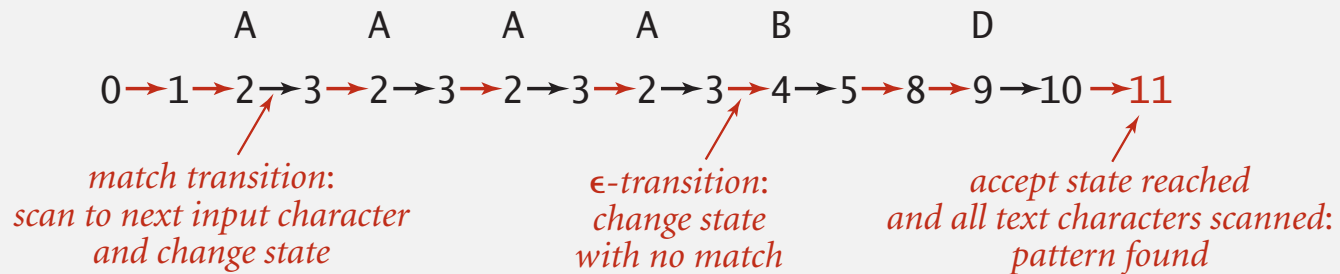


NFA corresponding to the pattern  $((A * B | A C) D)$

# Nondeterministic finite-state automata

Q. Is A A A B D matched by NFA?

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



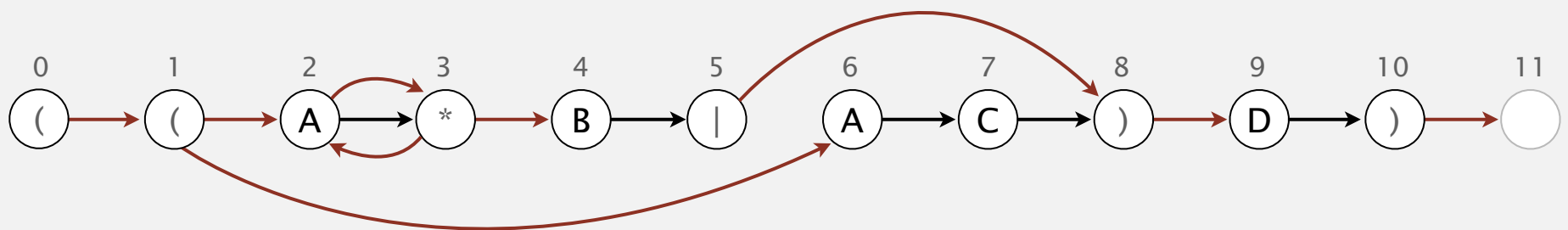
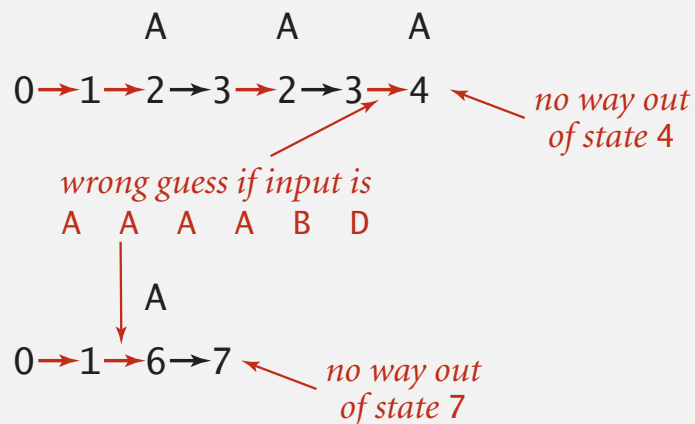
NFA corresponding to the pattern ( ( A \* B | A C ) D )

# Nondeterministic finite-state automata

Q. Is A A A B D 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 ]



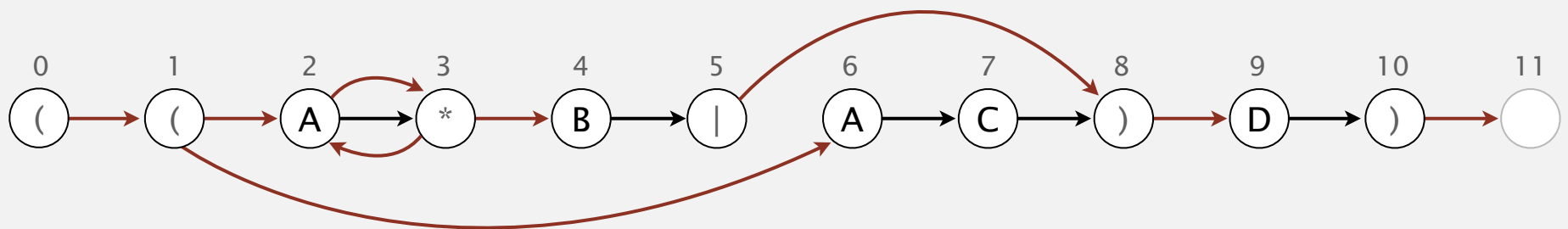
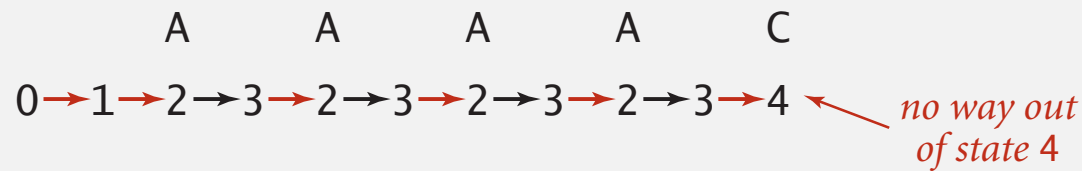
NFA corresponding to the pattern ( ( A \* B | A C ) D )

# Nondeterministic finite-state automata

Q. Is A A A C 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 )

# Nondeterminism

---

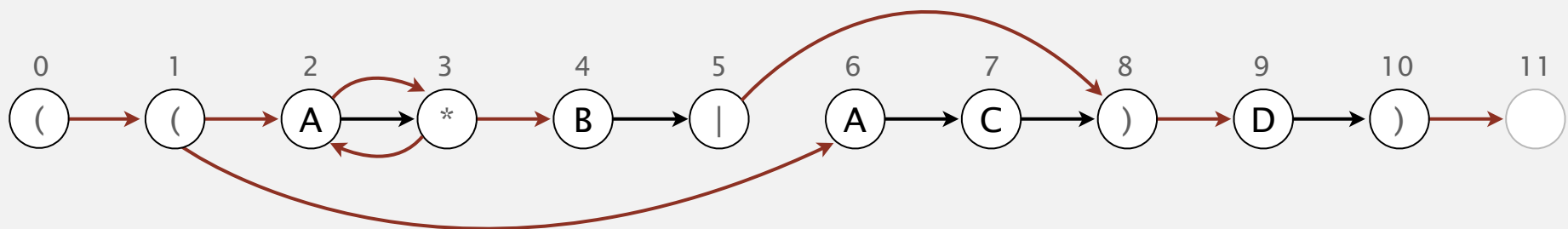
Q. How to determine whether a string is matched by an automaton?

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.



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

# NFA representation

State names. Integers from 0 to  $M$ .

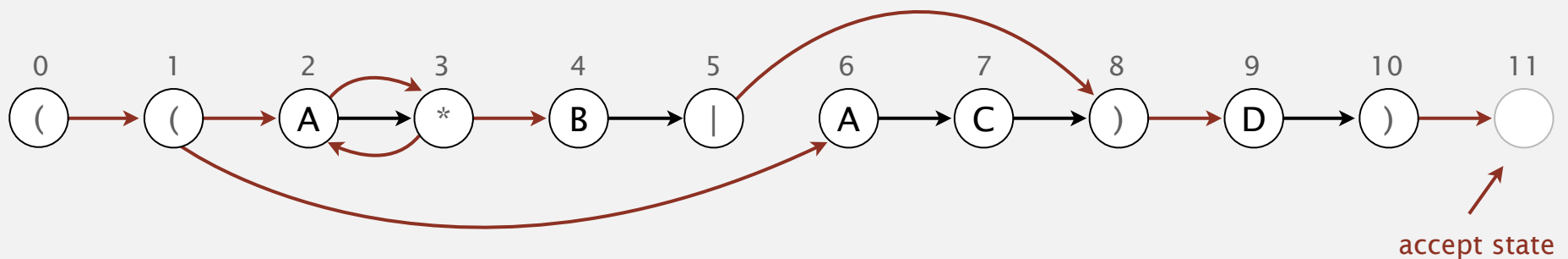
number of symbols in RE

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

e.g. `re[4] = 'B'`

$\epsilon$ -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 )`

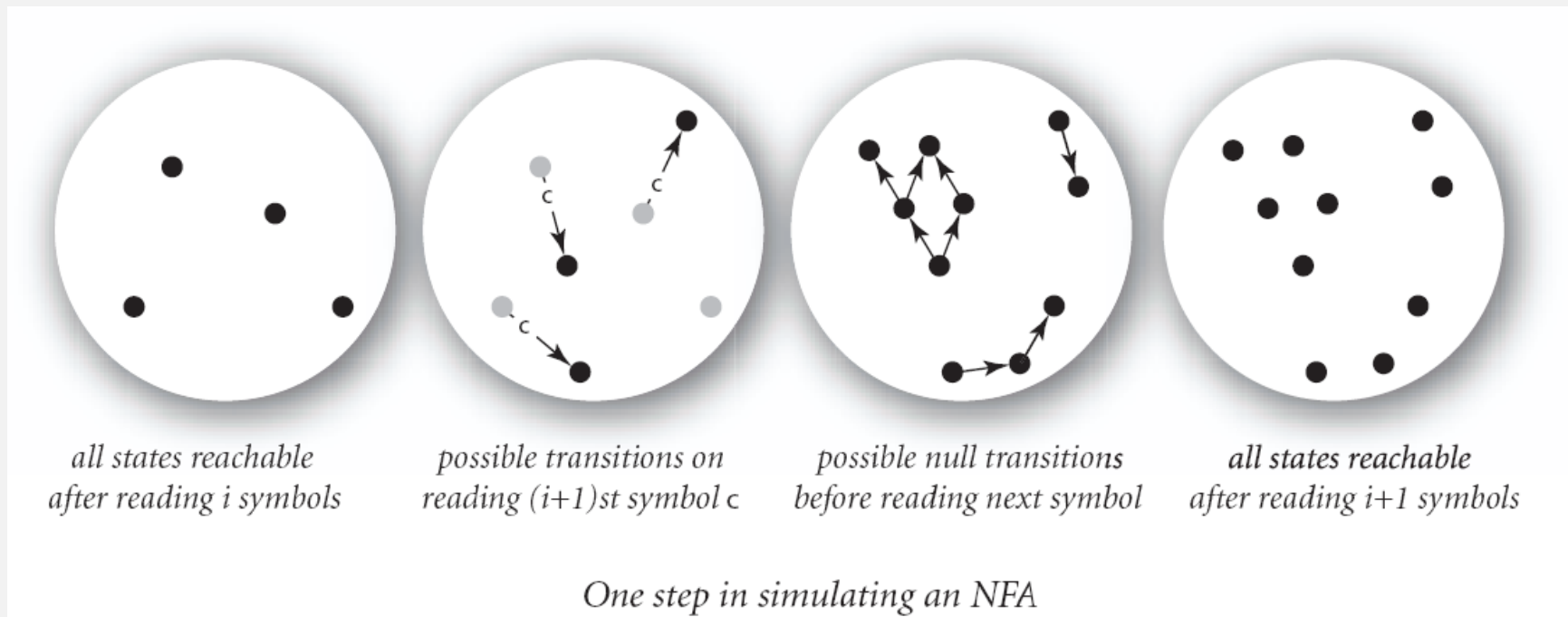
accept state

## 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?

# 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  $E + V$ .



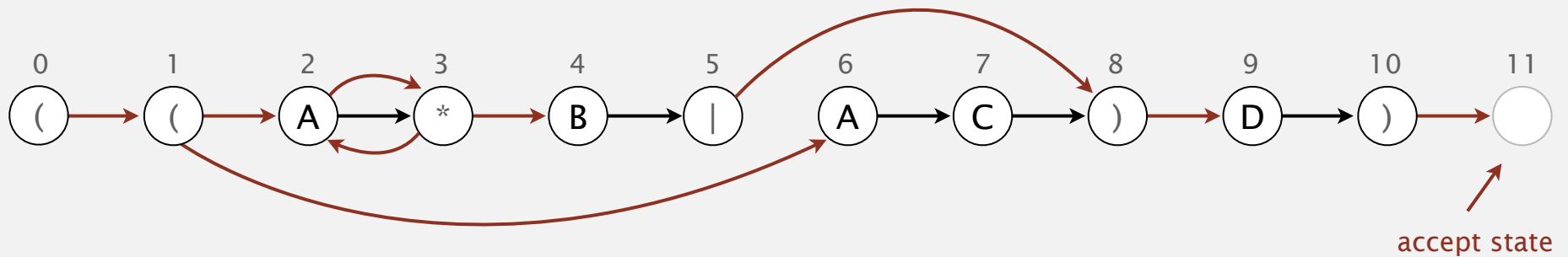
## Runtime for NFA simulation

### NFA simulation steps for each text character:

- Find states reachable by match transitions.
- Find states reachable by  $\epsilon$ -transitions

### Let:

- $N$  = length of text
- $M$  = length of pattern (also number of DFA states)
- $R$  = alphabet size



pollEv.com/jhug

text to 37607

Q: What is the runtime to complete NFA simulation in the worst case?

- |         |          |                |          |
|---------|----------|----------------|----------|
| A. $N$  | [428594] | D. $N + M$     | [428616] |
| B. $M$  | [428595] | E. $N + M + R$ | [428617] |
| C. $NM$ | [428615] | F. $NMR$       | [428634] |

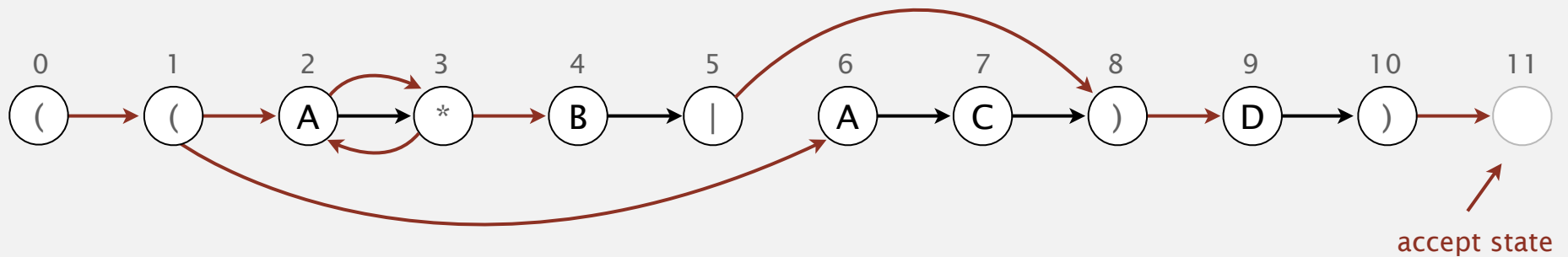
## Runtime for NFA simulation

### NFA simulation steps for each text character:

- Find states reachable by match transitions.
- Find states reachable by  $\varepsilon$ -transitions

### Let:

- $N$  = length of text
- $M$  = length of pattern (also number of DFA states)
- $R$  = alphabet size



Q: What is the runtime to complete NFA simulation in the worst case?

C.  $NM$

Finding states reachable by match transitions:  $O(M)$

Finding states reachable by match transitions:  $V + E$

Number of vertices:  $M$ . Total number of epsilon edges:  $O(M)$  (stay tuned)

# 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 = buildEpsilonTransitionsDigraph();
    }

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

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

← stay tuned (next segment)

# 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);
```

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

```
    for (int i = 0; i < txt.length(); i++)
    {
```

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

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

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

← follow  $\epsilon$ -transitions

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

← accept if can end in state M

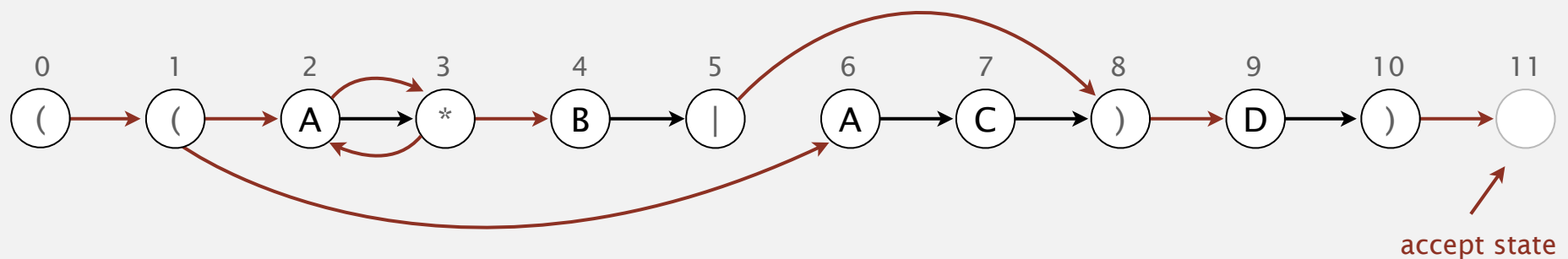
## NFA simulation: analysis

**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  $\leq 3M$ .]

**Alternate Pf.** See PollEverywhere question.



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



<http://algs4.cs.princeton.edu>

## 5.4 REGULAR EXPRESSIONS

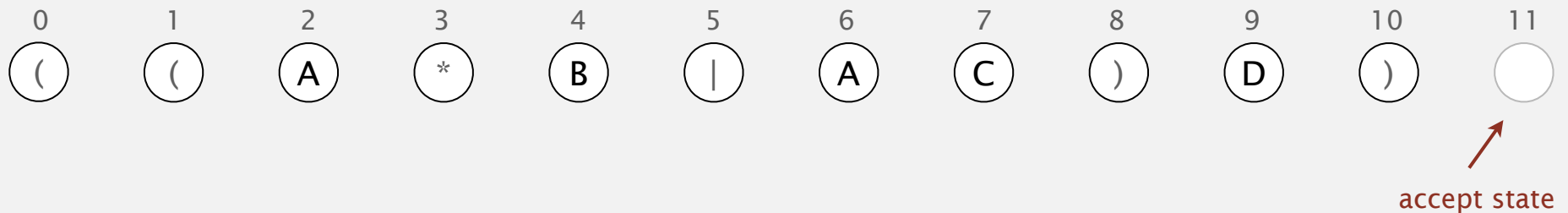
---

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

# Building an NFA corresponding to an RE

---

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



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

# 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.** ( ) . \* |



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



# Building an NFA corresponding to an RE

---

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

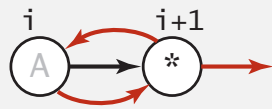


NFA corresponding to the pattern  $( ( A * B | A C ) D )$

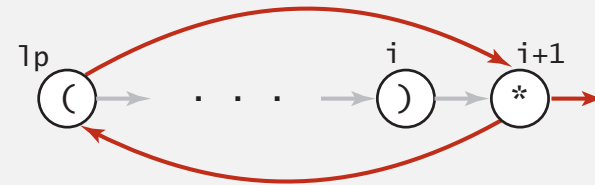
# Building an NFA corresponding to an RE

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

single-character closure



closure expression

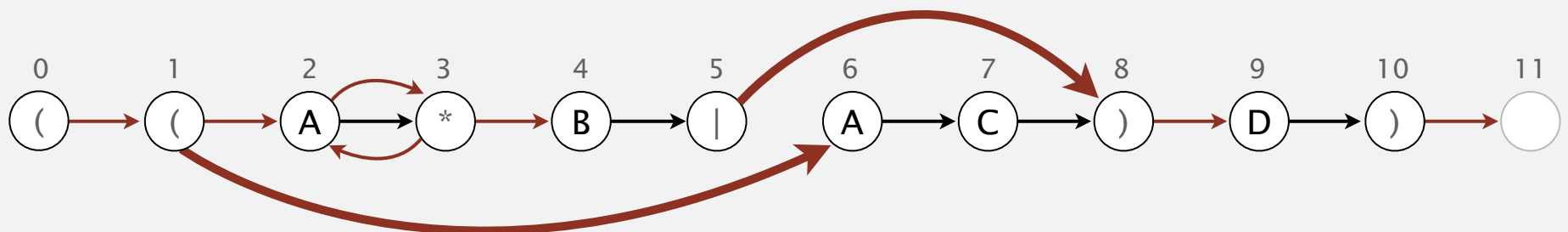
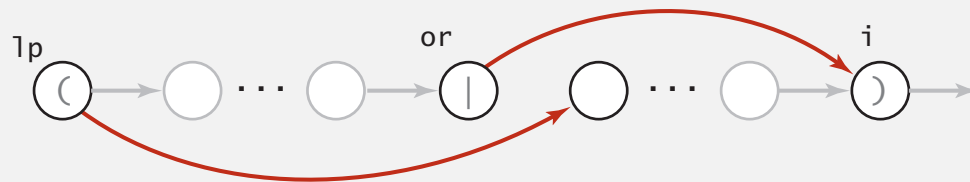


NFA corresponding to the pattern  $((A * B | A C) D)$

# Building an NFA corresponding to an RE

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

or expression



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

# NFA construction: implementation

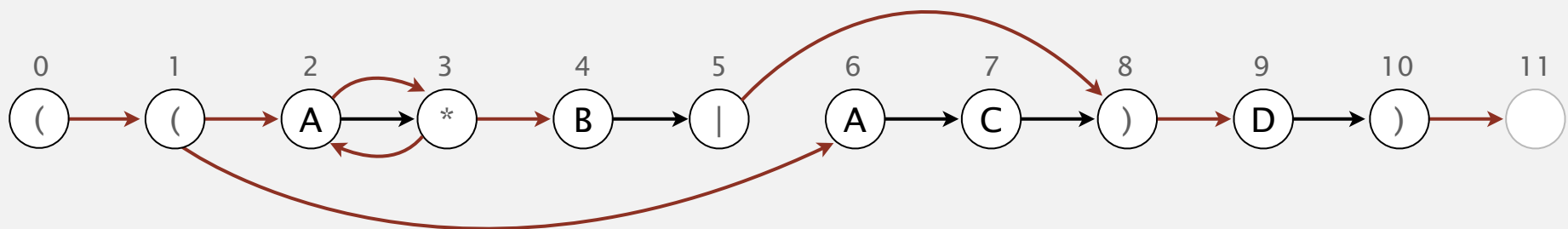
---

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

# NFA construction demo

---

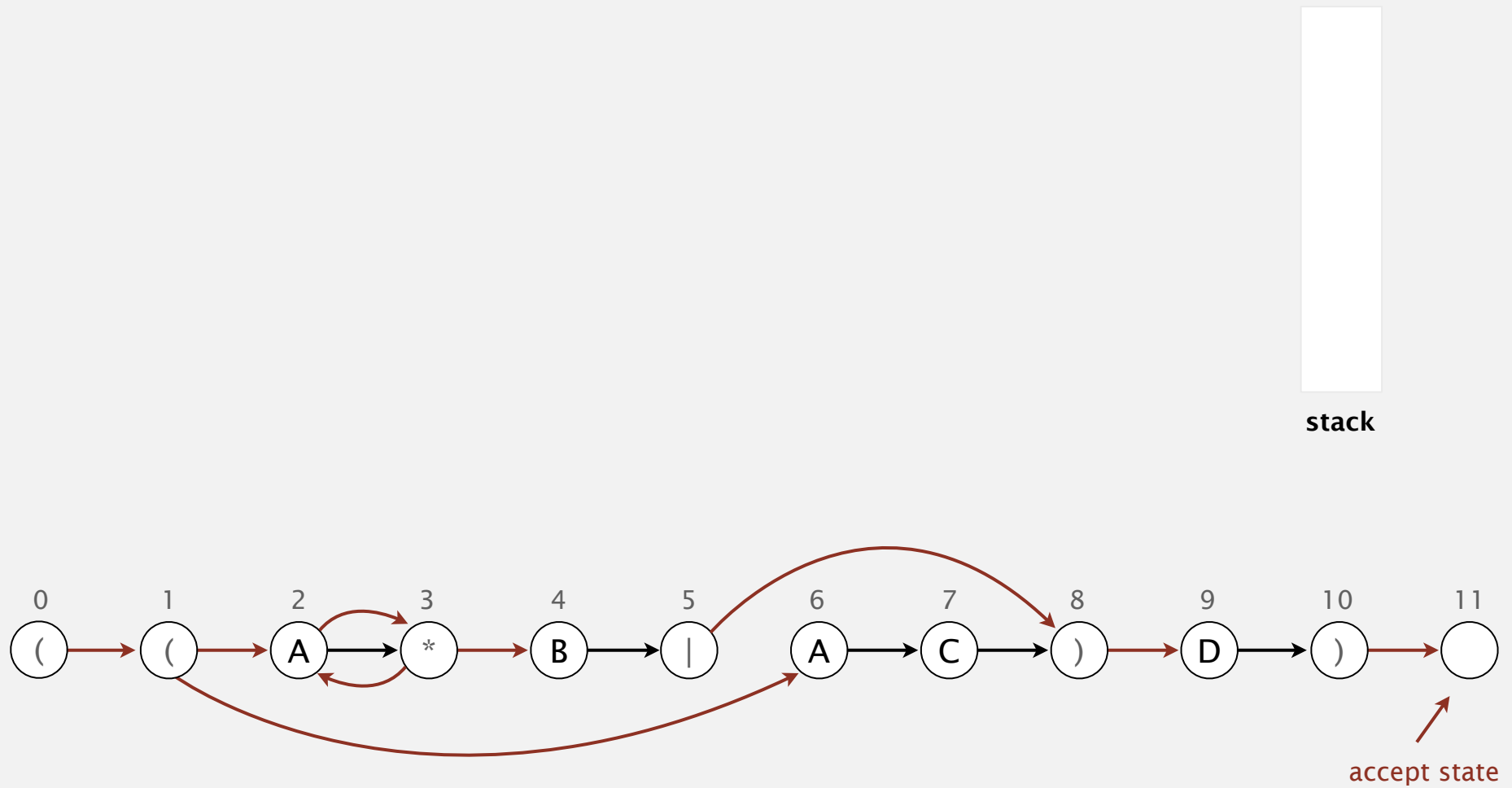


stack

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

# NFA construction demo

---



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

# 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;
}
```

← left parentheses and |

← or

← closure  
(needs 1-character lookahead)

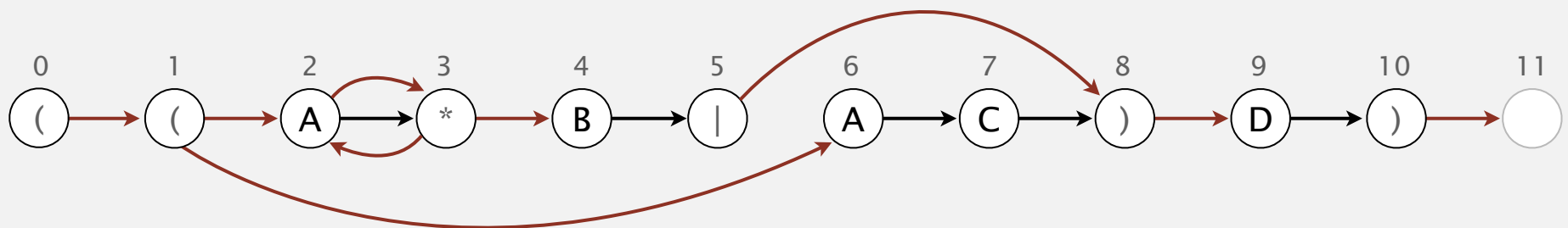
← metasympols

## 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  $( ( A * B | A C ) D )$





<http://algs4.cs.princeton.edu>

## 5.4 REGULAR EXPRESSIONS

---

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

## Generalized regular expression print

---

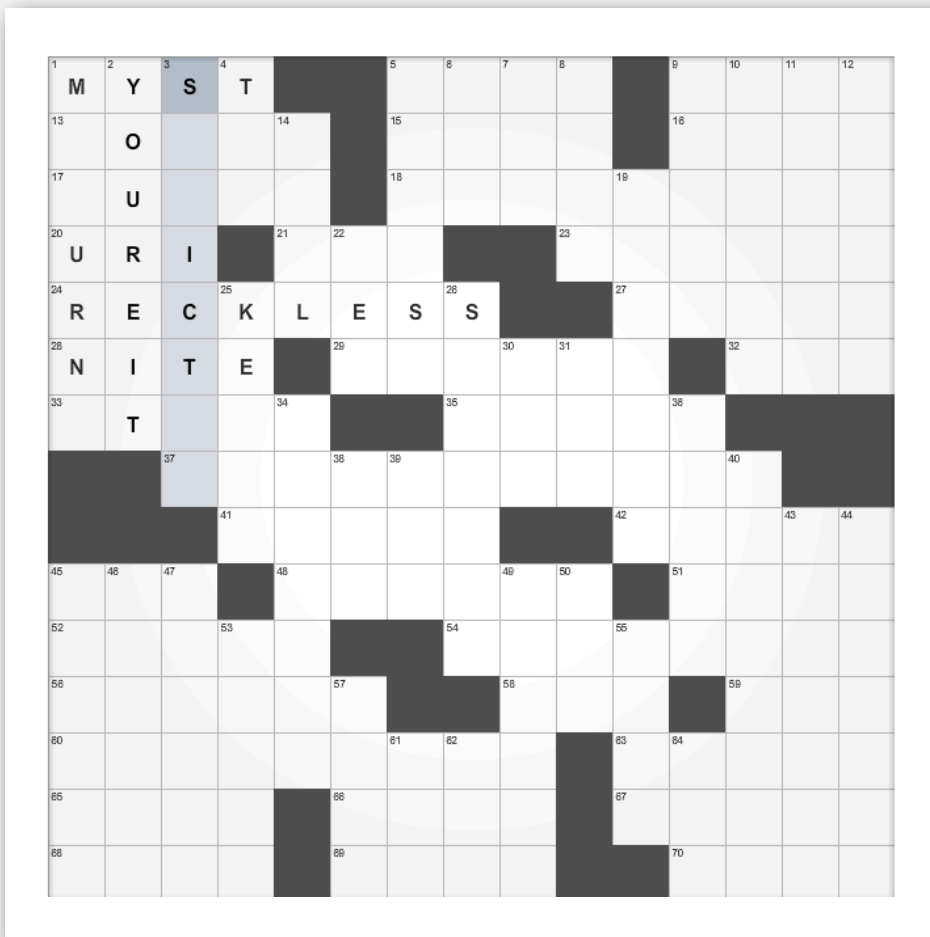
**Grep.** Take a RE as a command-line argument and print the lines from standard input having a **substring that is matched** by the RE.

```
public class GREP
{
    public static void main(String[] args)
    {
        String re = "(.*" + args[0] + ".*)";
        NFA nfa = new NFA(re);
        while (StdIn.hasNextLine())
        {
            String line = StdIn.readLine();
            if (nfa.recognizes(line))
                StdOut.println(line);
        }
    }
}
```

← contains RE  
as a substring

**Bottom line.** Worst-case for grep (proportional to  $MN$ ) is the same as for brute-force substring search.

# Typical grep application: crossword puzzles



```
% more words.txt
```

```
a  
aback  
abacus  
abalone  
abandon
```

```
...
```

```
% grep "s..ict.." words.txt
```

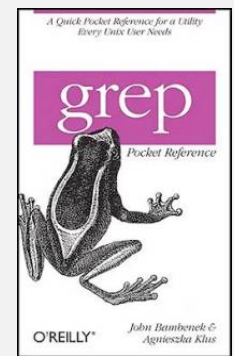
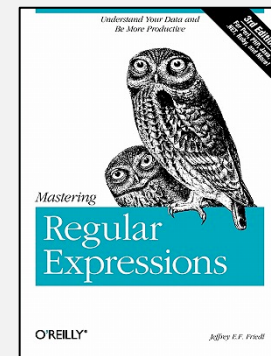
```
constrictor  
stricter  
stricture
```

dictionary  
(standard in Unix)  
also on booksite

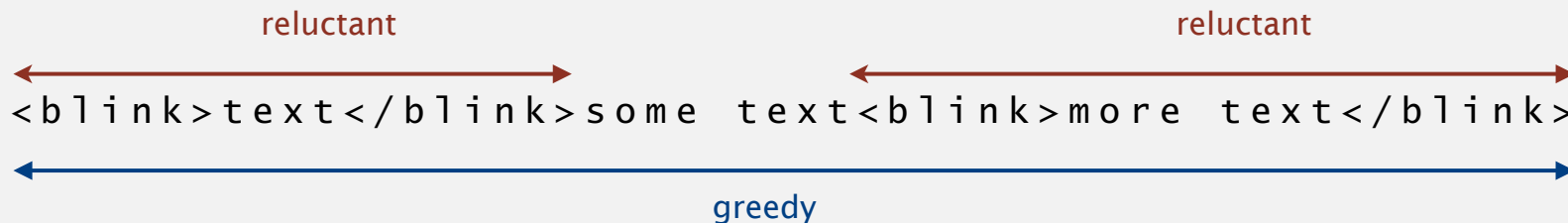
# Industrial-strength grep implementation

## To complete the implementation:

- Add character classes. `[0-9][0-9]`
- Handle metacharacters. `a\.b`
- Add capturing capabilities.
- Extend the closure operator. `a+b`
- Error checking and recovery. `*|*|*`
- Greedy vs. reluctant matching. `to.*?be`



Ex. Which substring(s) should be matched by the RE `<blink>.*</blink>` ?



# 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
true
```

← legal Java identifier

```
% java Validate "[a-z]+@([a-z]+\.)+(edu|com)" rs@cs.princeton.edu
true
```

← valid email address  
(simplified)

```
% java Validate "[0-9]{3}-[0-9]{2}-[0-9]{4}" 166-11-4433
true
```

← Social Security number

# 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



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

# 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()`

# Algorithmic complexity attacks

---

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



Unix grep, Java, Perl

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



# Not-so-regular expressions

---

## Back-references.

- `\1` notation matches subexpression that was matched earlier.
- Supported by typical RE implementations, but not our NFA approach.

```
(.+)\1          // beriberi couscous  
1?$|^(11+?)\1+ // 1111 111111 111111111
```

## Some non-regular languages.

- Strings of the form  $w w$  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: atttcggaat.

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

# 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.

# 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

# Regular expressions in other languages

---

## Broadly applicable programmer's tool.

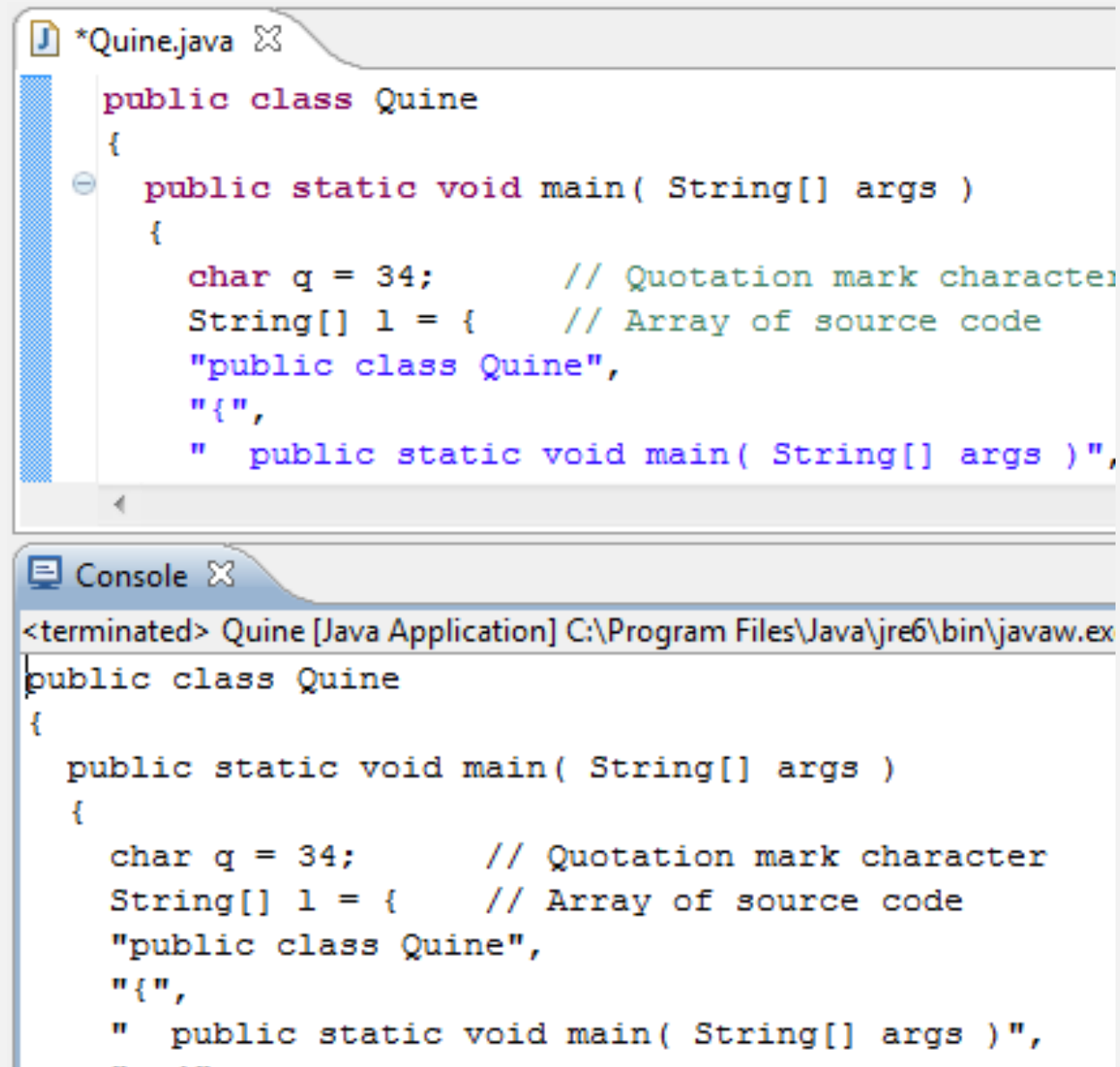
- Originated in Unix in the 1970s.
- Popularized through early tools developed by Ken Thompson.
- Many languages support extended regular expressions.
- Built into grep, awk, emacs, Perl, PHP, Python, JavaScript, ...

```
% grep 'NEWLINE' */*.java
```

← print all lines containing NEWLINE which occurs in any file with a .java extension

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

## Some programs produce themselves as output



```
*Quine.java X
public class Quine
{
    public static void main( String[] args )
    {
        char q = 34;        // Quotation mark character
        String[] l = {      // Array of source code
            "public class Quine",
            "{",
            "    public static void main( String[] args )",
            "    {"
        };
    }
}

Console X
<terminated> Quine [Java Application] C:\Program Files\Java\jre6\bin\javaw.exe
public class Quine
{
    public static void main( String[] args )
    {
        char q = 34;        // Quotation mark character
        String[] l = {      // Array of source code
            "public class Quine",
            "{",
            "    public static void main( String[] args )",
            "    {"
```

**What about programs that take themselves as input?**

**“Reflections on Trusting Trust” - Ken Thompson (1983 Turing lecture)**