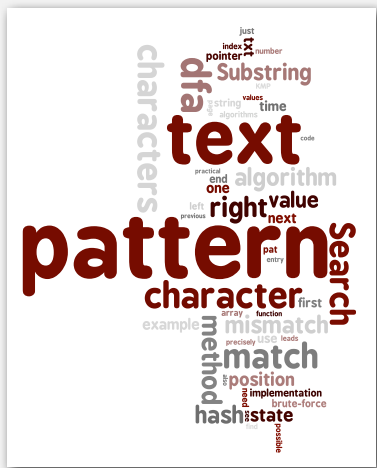


5.3 Substring Search



- ▶ brute force
- ▶ Knuth-Morris-Pratt
- ▶ Boyer-Moore
- ▶ Rabin-Karp

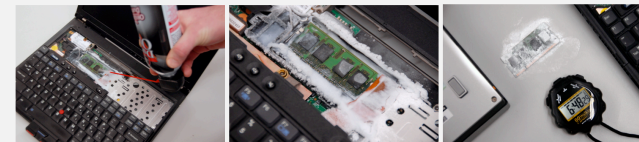
Substring search

Goal. Find pattern of length M in a text of length N .

typically $N \gg M$



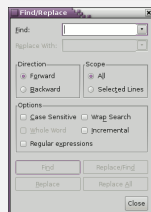
Computer forensics. Search memory or disk for signatures, e.g., all URLs or RSA keys that the user has entered.



<http://citp.princeton.edu/memory>

Applications

- Parsers.
- Spam filters.
- Digital libraries.
- Screen scrapers.
- Word processors.
- Web search engines.
- Electronic surveillance.
- Natural language processing.
- Computational molecular biology.
- FBI's Digital Collection System 3000.
- Feature detection in digitized images.
- ...



Application: spam filtering

Identify patterns indicative of spam.

- PROFITS
- LOSE WEIGHT
- herbal Viagra
- There is no catch.
- LOW MORTGAGE RATES
- This is a one-time mailing.
- This message is sent in compliance with spam regulations.



Application: electronic surveillance



Need to monitor all internet traffic. (security)

No way! (privacy)



Well, we're mainly interested in "ATTACK AT DAWN"



OK. Build a machine that just looks for that.



"ATTACK AT DAWN" substring search machine found

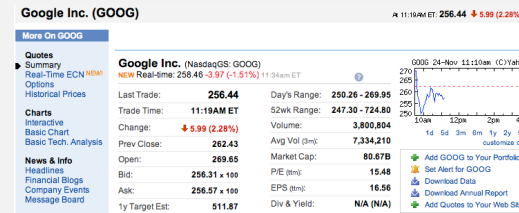


5

Application: screen scraping

Goal. Extract relevant data from web page.

Ex. Find string delimited by `` and `` after first occurrence of pattern `Last Trade:`.



<http://finance.yahoo.com/q?s=goog>

```
...  
<tr>  
<td class= "yfnctablehead1"  
width= "48%">  
Last Trade:  
</td>  
<td class= "yfnctabledatal1">  
<big><b>452.92</b></big>  
</td></tr>  
<td class= "yfnctablehead1"  
width= "48%">  
Trade Time:  
</td>  
<td class= "yfnctabledatal1">  
...  

```

6

Screen scraping: Java implementation

Java library. The `indexOf()` method in Java's string library returns the index of the first occurrence of a given string, starting at a given offset.

```
public class StockQuote  
{  
    public static void main(String[] args)  
    {  
        String name = "http://finance.yahoo.com/q?s=";  
        In in = new In(name + args[0]);  
        String text = in.readAll();  
        int start = text.indexOf("Last Trade:", 0);  
        int from = text.indexOf("<b>", start);  
        int to = text.indexOf("</b>", from);  
        String price = text.substring(from + 3, to);  
        StdOut.println(price);  
    }  
}
```

```
% java StockQuote goog  
564.35
```

```
% java StockQuote msft  
26.04
```

7

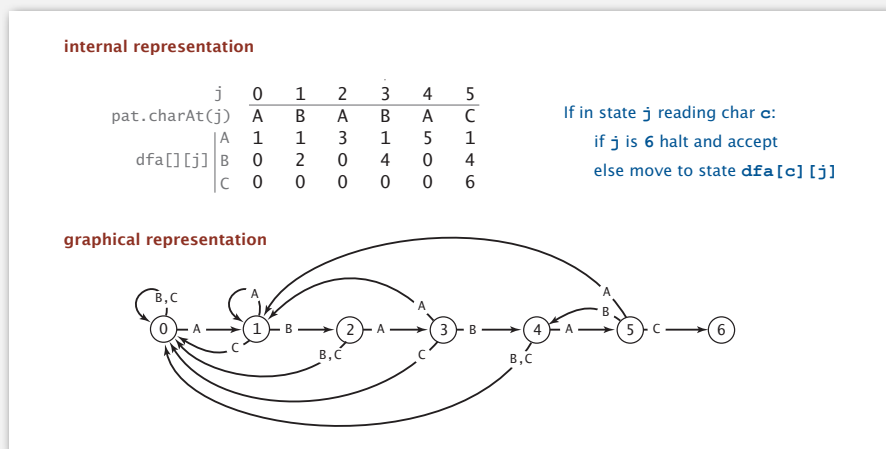
- ▶ brute force
- ▶ Knuth-Morris-Pratt
- ▶ Boyer-Moore
- ▶ Rabin-Karp

8

Deterministic finite state automaton (DFA)

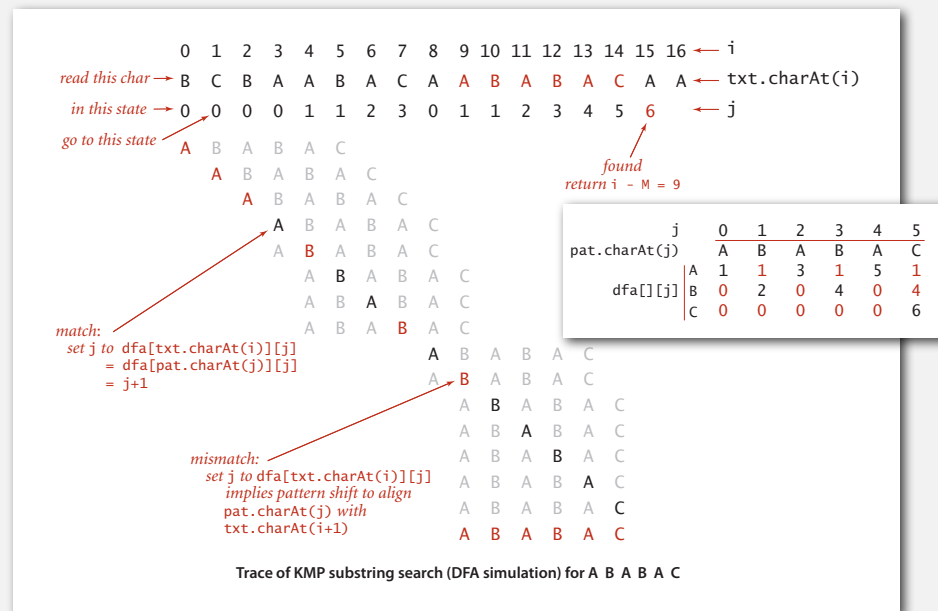
DFA is abstract string-searching machine.

- Finite number of states (including start and halt).
- Exactly one transition for each char in alphabet.
- Accept if sequence of transitions leads to halt state.



17

KMP substrig search: trace



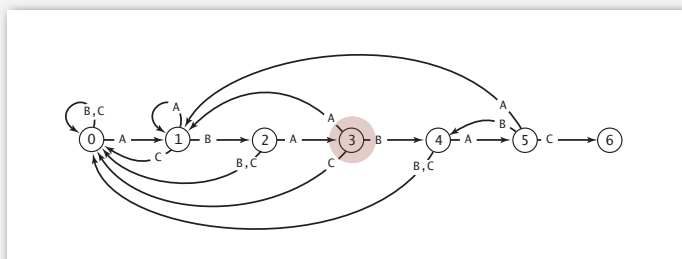
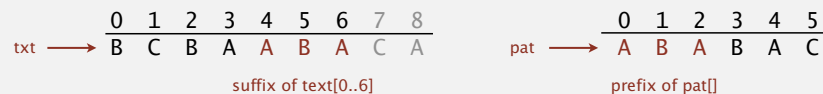
18

Interpretation of Knuth-Morris-Pratt DFA

Q. What is interpretation of DFA state after reading in $txt[i]$?

A. State = number of characters in pattern that have been matched.
 (length of longest prefix of $pat[]$ that is a suffix of $txt[0..i]$)

Ex. DFA is in state 3 after reading in character $txt[6]$.



19

KMP search: Java implementation

Key differences from brute-force implementation.

- Text pointer i never decrements.
- Need to precompute $dfa[][]$ from pattern.

```
public int search(String txt)
{
    int i, j, N = txt.length();
    for (i = 0, j = 0; i < N && j < M; i++)
        j = dfa[txt.charAt(i)][j];
    if (j == M) return i - M;
    else return N;
}
```

← no backup

Running time.

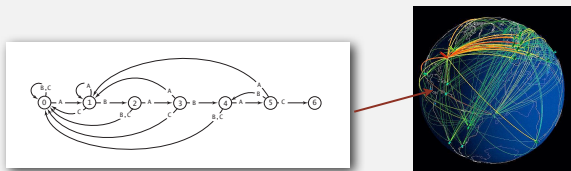
- Simulate DFA on text: at most N character accesses.
- Build DFA: how to do efficiently? [warning: tricky algorithm ahead]

20

Key differences from brute-force implementation.

- Text pointer i never decrements.
- Need to precompute $dfa[][]$ from pattern.
- Could use **input stream**.

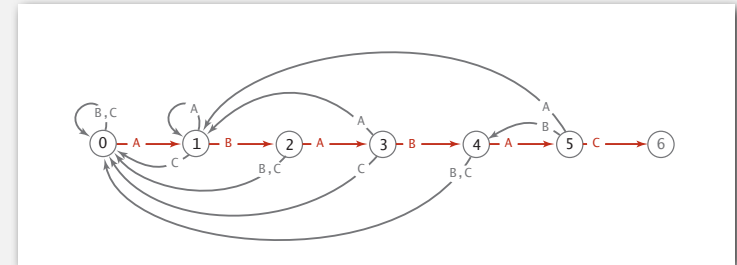
```
public int search(In in)
{
    int i, j;
    for (i = 0, j = 0; !in.isEmpty() && j < M; i++)
        j = dfa[in.readChar()][j];
    if (j == M) return i - M;
    else return NOT_FOUND;
}
```



Match transition. If in state j and next char $c == pat.charAt(j)$, then go to state $j+1$.

now first $j+1$ characters of pattern have been matched
 first j characters of pattern have already been matched
 next char matches

j	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C



Mismatch transition. If in state j and next char $c != pat.charAt(j)$, then the last j characters of input are $pat[1..j-1]$, followed by c .

To compute $dfa[c][j]$: Simulate $pat[1..j-1]$ on DFA and take transition c .

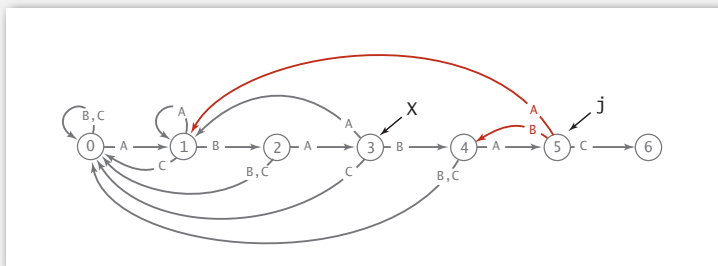
Running time. Seems to require j steps.

↑
still under construction (!)

Ex. $dfa['A'][5] = 1$; $dfa['B'][5] = 4$

simulate BABA (state X); take transition 'A'
 simulate BABA (state X); take transition 'B'

j	0	1	2	3	4	5
pat.charAt(j)	A	B	A	B	A	C



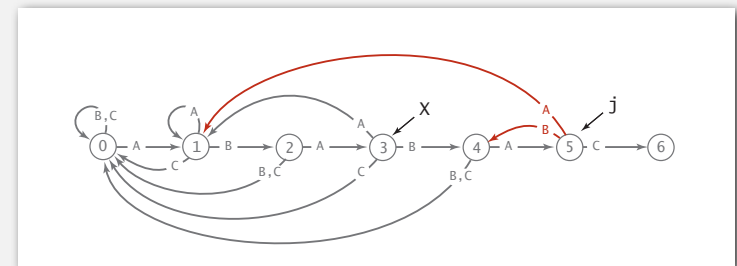
Mismatch transition. If in state j and next char $c != pat.charAt(j)$, then the last j characters of input are $pat[1..j-1]$, followed by c .

To compute $dfa[c][j]$: Simulate $pat[1..j-1]$ on DFA and take transition c .

Running time. Takes only constant time if we know state X . (!)

Ex. $dfa['A'][5] = 1$; $dfa['B'][5] = 4$; $X' = 0$

from state X, take transition 'A' = $dfa['A'][X]$
 from state X, take transition 'B' = $dfa['B'][X]$
 from state X, take transition 'C' = $dfa['C'][X]$



Constructing the DFA for KMP substring search: example

pat.charAt(j)	0	1	2	3
A	1			
B	0			
C	0			

pat.charAt(j)	0	1
A	B	
B	0	
C	0	

copy dfa[] [X] to dfa[] [j]
 dfa[pat.charAt(j)][j] = j+1;
 X = dfa[pat.charAt(j)][X];

pat.charAt(j)	0	1	2
A	B	A	
B	0	2	
C	0	0	

pat.charAt(j)	0	1	2	3
A	B	A	B	
B	0	2	0	
C	0	0	0	

Constructing the DFA for KMP substring search for A B A B A C

25

Constructing the DFA for KMP substring search: example

pat.charAt(j)	0	1	2	3
A	B	A	B	
B	0	2	0	
C	0	0	0	

pat.charAt(j)	0	1	2	3	4
A	B	A	B	A	
B	0	2	0	4	
C	0	0	0	0	

pat.charAt(j)	0	1	2	3	4	5
A	B	A	B	A	C	
B	0	2	0	4	0	
C	0	0	0	0	6	

Constructing the DFA for KMP substring search for A B A B A C

26

Constructing the DFA for KMP substring search: Java implementation

For each state j :

- Copy $dfa[] [X]$ to $dfa[] [j]$ for mismatch case.
- Set $dfa[pat.charAt(j)][j]$ to $j+1$ for match case.
- Update x .

```
public KMP(String pat)
{
    this.pat = pat;
    M = pat.length();
    dfa = new int[R][M];
    dfa[pat.charAt(0)][0] = 1;
    for (int X = 0, j = 1; j < M; j++)
    {
        for (int c = 0; c < R; c++)
            dfa[c][j] = dfa[c][X];
        dfa[pat.charAt(j)][j] = j+1;
        X = dfa[pat.charAt(j)][X];
    }
}
```

← copy mismatch cases
 ← set match case
 ← update restart state

Running time. M character accesses (but space proportional to RM).

27

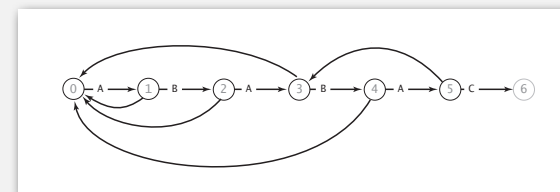
KMP substring search analysis

Proposition. KMP substring search accesses no more than $M + N$ chars to search for a pattern of length M in a text of length N .

Pf. Each pattern char accessed once when constructing the DFA; each text char accessed once (in the worst case) when simulating the DFA.

Proposition. KMP constructs $dfa[] []$ in time and space proportional to RM .

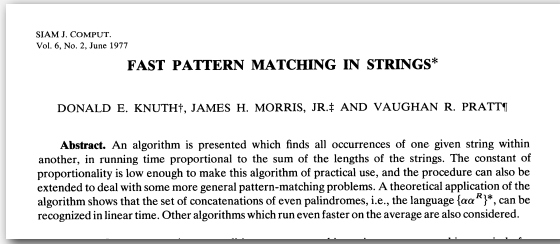
Larger alphabets. Improved version of KMP constructs $nfa[] []$ in time and space proportional to M .



28

Knuth-Morris-Pratt: brief history

- Independently discovered by two theoreticians and a hacker.
 - Knuth: inspired by esoteric theorem, discovered linear-time algorithm
 - Pratt: made running time independent of alphabet size
 - Morris: built a text editor for the CDC 6400 computer
- Theory meets practice.



Don Knuth



Jim Morris



Vaughan Pratt

Boyer-Moore: mismatched character heuristic

Intuition.

- Scan characters in pattern from right to left.
- Can skip M text chars when finding one not in the pattern.

i	j	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
$text \rightarrow$		F	I	N	D	I	N	A	H	A	Y	S	T	A	C	K	N	E	E	D	L	E	I	N	A
		0	5																						
		5	5																						
		11	4																						
		15	0																						

$return\ i = 15$

- › brute force
- › Knuth-Morris-Pratt
- › Boyer-Moore
- › Rabin-Karp



Robert Boyer



J. Strother Moore

Boyer-Moore: mismatched character heuristic

Q. How much to skip?

A. Compute $right[c]$ = rightmost occurrence of character c in pat .

```
right = new int[R];
for (int c = 0; c < R; c++)
    right[c] = -1;
for (int j = 0; j < M; j++)
    right[pat.charAt(j)] = j;
```

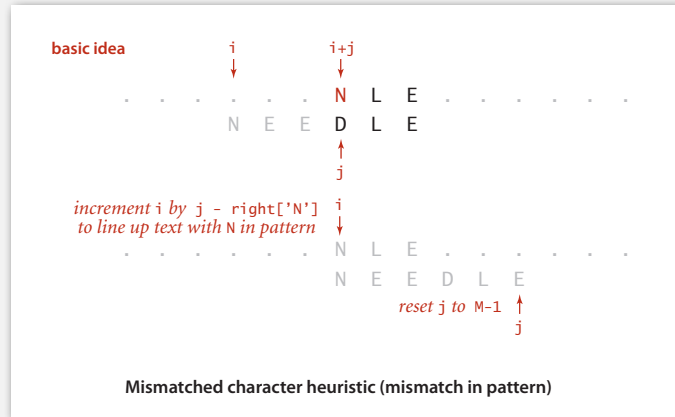
c	N	E	E	D	L	E	$right[c]$
A	-1	-1	-1	-1	-1	-1	-1
B	-1	-1	-1	-1	-1	-1	-1
C	-1	-1	-1	-1	-1	-1	-1
D	-1	-1	-1	-1	3	3	3
E	-1	-1	1	2	2	5	5
...							-1
L	-1	-1	-1	-1	-1	4	4
M	-1	-1	-1	-1	-1	-1	-1
N	-1	0	0	0	0	0	0
...							-1

Boyer-Moore skip table computation

Boyer-Moore: mismatched character heuristic

Q. How much to skip?

A. Compute `right[c]` = rightmost occurrence of character `c` in `pat`.

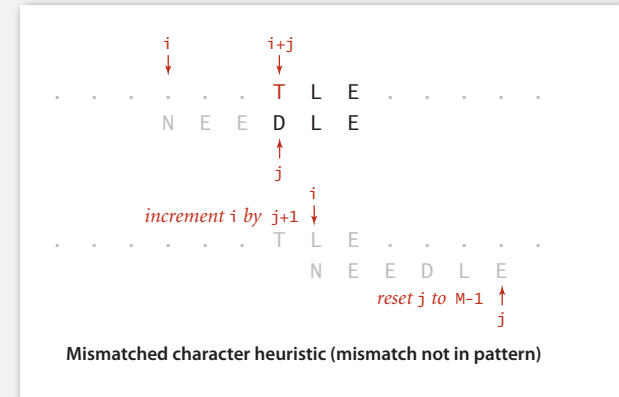


34

Boyer-Moore: mismatched character heuristic

Q. How much to skip?

A. Compute `right[c]` = rightmost occurrence of character `c` in `pat`.



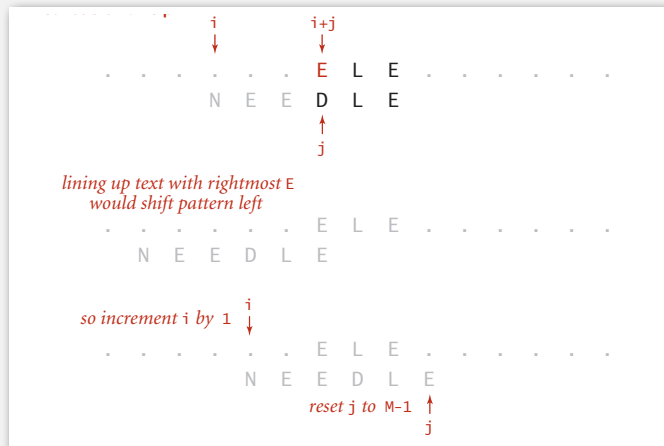
Character not in pattern? Set `right[c]` to `-1`.

35

Boyer-Moore: mismatched character heuristic

Q. How much to skip?

A. Compute `right[c]` = rightmost occurrence of character `c` in `pat`.



Heuristic no help? Increment i and reset j to $M-1$

36

Boyer-Moore: Java implementation

```
public int search(String txt)
{
    int N = txt.length();
    int M = pat.length();
    int skip;
    for (int i = 0; i <= N-M; i += skip)
    {
        skip = 0;
        for (int j = M-1; j >= 0; j--)
        {
            if (pat.charAt(j) != txt.charAt(i+j))
            {
                skip = Math.max(1, j - right[txt.charAt(i+j)]);
                break;
            }
        }
        if (skip == 0) return i;
    }
    return N;
}
```

compute skip value

match

37

Boyer-Moore: analysis

Property. Substring search with the Boyer-Moore mismatched character heuristic takes about $\sim N/M$ character compares to search for a pattern of length M in a text of length N . *sublinear*

Worst-case. Can be as bad as $\sim MN$.

i	skip	0	1	2	3	4	5	6	7	8	9
		B	B	B	B	B	B	B	B	B	B
0	0	A	B	B	B	B					
1	1		A	B	B	B					
2	1			A	B	B	B				
3	1				A	B	B	B			
4	1					A	B	B	B		
5	1						A	B	B	B	

Boyer-Moore variant. Can improve worst case to $\sim 3N$ by adding a KMP-like rule to guard against repetitive patterns.

38

- brute force
- Knuth-Morris-Pratt
- Boyer-Moore
- Rabin-Karp



Michael Rabin, Turing Award '76
and Dick Karp, Turing Award '85

39

Rabin-Karp fingerprint search

Basic idea = modular hashing.

- Compute a hash of pattern characters 0 to $M - 1$.
- For each i , compute a hash of text characters i to $M + i - 1$.
- If pattern hash = text substring hash, check for a match.

i	pat.charAt(i)	
0	2	$\% 997 = 613$
1	6	
2	5	
3	3	
4	5	

i	txt.charAt(i)	
0	3	$\% 997 = 508$
1	4	
2	1	
3	5	
4	9	
5	2	$\% 997 = 201$
6	6	
7	5	
8	3	
9	5	
10	8	$\% 997 = 715$
11	9	
12	2	
13	6	
14	5	
15	9	$\% 997 = 971$
16	7	
17	9	
18	3	
19	3	
20	2	$\% 997 = 442$
21	6	
22	5	
23	3	
24	5	
25	9	$\% 997 = 929$
26	2	
27	6	
28	5	
29	3	
30	5	$\% 997 = 613$
31	2	
32	6	
33	5	
34	3	

6 ← return i = 6

40

Efficiently computing the hash function

Modular hash function. Using the notation t_i for `txt.charAt(i)`, we wish to compute

$$x_i = t_i R^{M-1} + t_{i+1} R^{M-2} + \dots + t_{i+M-1} R^0 \pmod{Q}$$

Intuition. M -digit, base- R integer, modulo Q .

Horner's method. Linear-time method to evaluate degree- M polynomial.

i	pat.charAt(i)	
0	2	$\% 997 = 2$
1	6	
2	5	
3	3	
4	5	
1	2	$\% 997 = (2*10 + 6) \% 997 = 26$
2	6	
3	5	
4	3	
5	5	
2	2	$\% 997 = (26*10 + 5) \% 997 = 265$
3	6	
4	5	
5	3	
6	5	
3	2	$\% 997 = (265*10 + 3) \% 997 = 659$
4	6	
5	5	
6	3	
7	5	
4	2	$\% 997 = (659*10 + 5) \% 997 = 613$
5	6	
6	5	
7	3	
8	5	

Computing the hash value for the pattern with Horner's method

```
// Compute hash for M-digit key
private int hash(String key, int M)
{
    int h = 0;
    for (int j = 0; j < M; j++)
        h = (R * h + key.charAt(j)) % Q;
    return h;
}
```

41

Rabin-Karp analysis

Theory. If Q is a sufficiently large random prime (about MN^2), then the probability of a false collision is about $1/N$.

Practice. Choose Q to be a large prime (but not so large as to cause overflow). Under reasonable assumptions, probability of a collision is about $1/Q$.

Monte Carlo version.

- Always runs in linear time.
- Extremely likely to return correct answer (but not always!).

Las Vegas version.

- Always returns correct answer.
- Extremely likely to run in linear time (but worst case is MN).

46

Rabin-Karp fingerprint search

Advantages.

- Extends to 2d patterns.
- Extends to finding multiple patterns.

Disadvantages.

- Arithmetic ops slower than char compares.
- Poor worst-case guarantee.
- Requires backup.

Q. How would you extend Rabin-Karp to efficiently search for any one of P possible patterns in a text of length N ?



47

Substring search cost summary

Cost of searching for an M -character pattern in an N -character text.

algorithm	version	operation count		backup in input?	correct?	extra space
		guarantee	typical			
brute force	—	MN	$1.1N$	yes	yes	1
Knuth-Morris-Pratt	full DFA (Algorithm 5.6)	$2N$	$1.1N$	no	yes	MR
	mismatch transitions only	$3N$	$1.1N$	no	yes	M
	full algorithm	$3N$	N/M	yes	yes	R
Boyer-Moore	mismatched char heuristic only (Algorithm 5.7)	MN	N/M	yes	yes	R
Rabin-Karp [†]	Monte Carlo (Algorithm 5.8)	$7N$	$7N$	no	yes [†]	1
	Las Vegas	$7N$ [†]	$7N$	yes	yes	1

[†] probabilistic guarantee, with uniform hash function

48