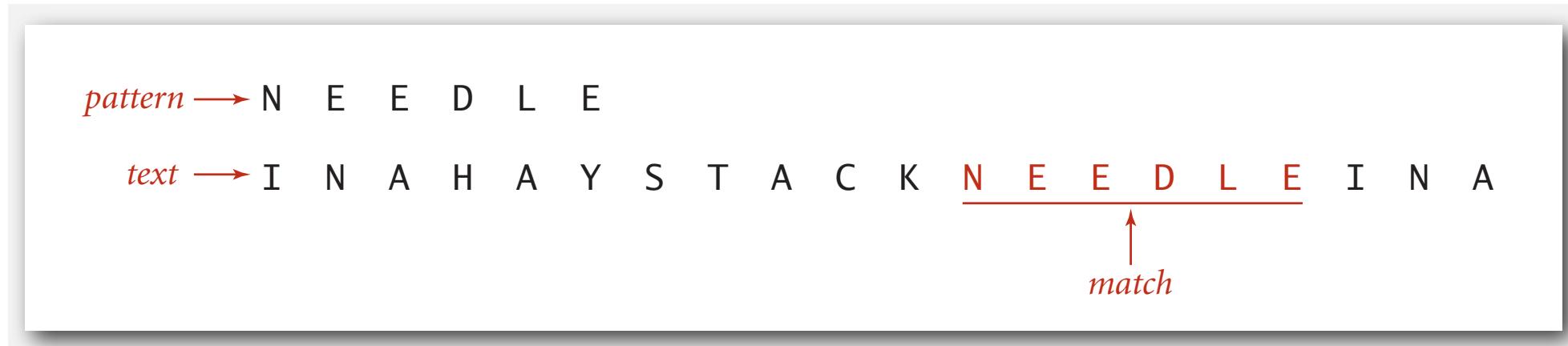


Substring Search

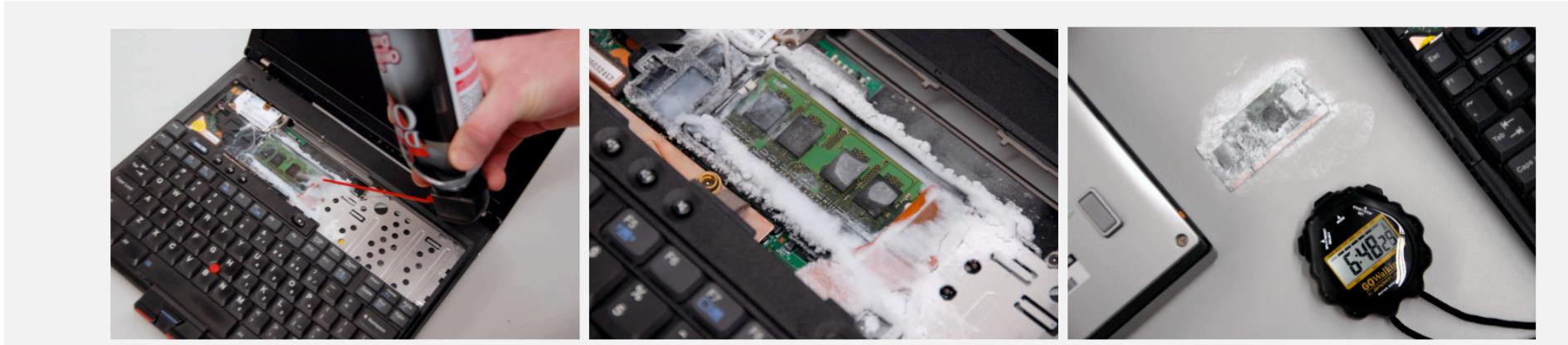
Substring Search

- Goal : Find pattern of length M in a text of length N



Application

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



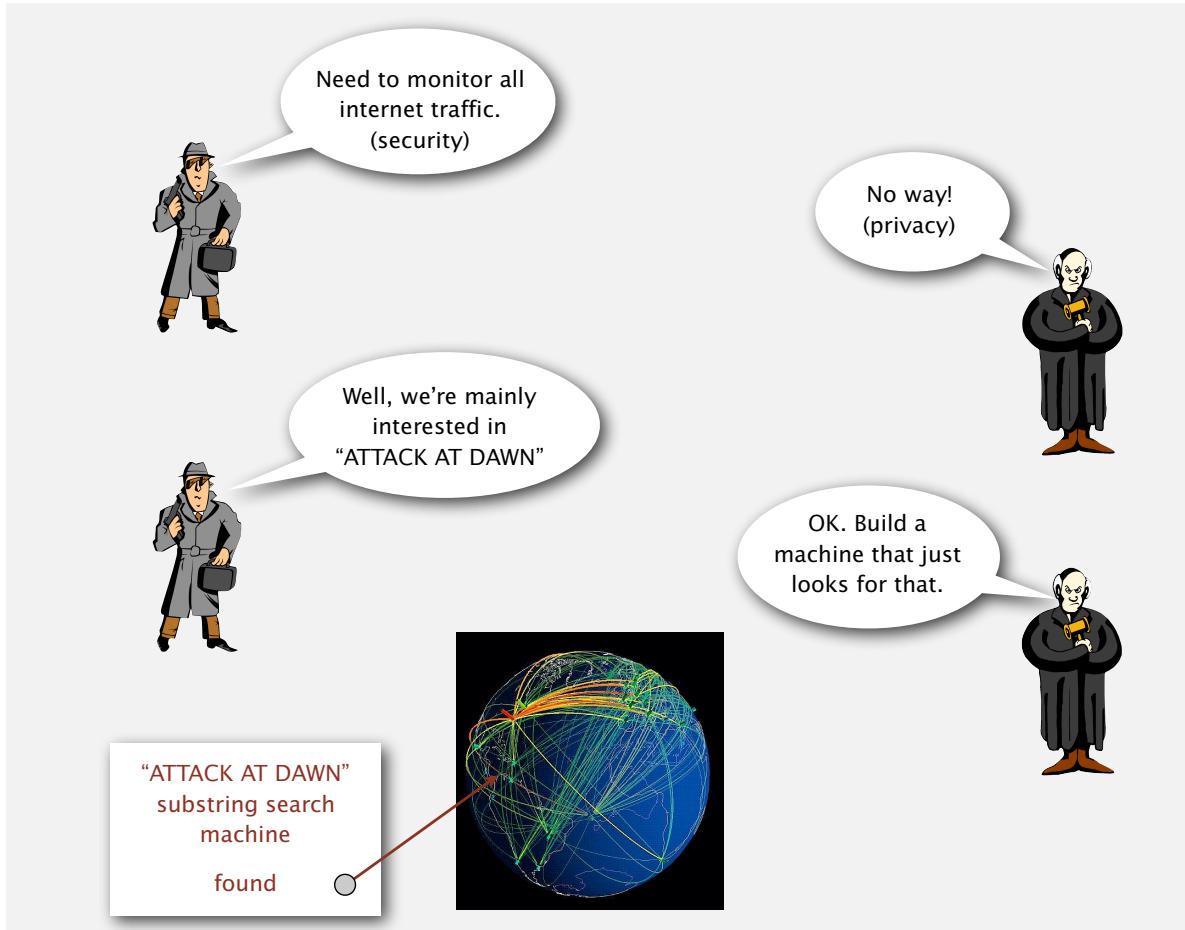
Application

- *Spam Filtering:* Identify patterns indicative of spam.

- 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*



Brute-force substring search

- Check for pattern starting at each text position.

i	j	i+j	0	1	2	3	4	5	6	7	8	9	10
			txt → A	B	A	C	A	D	A	B	R	A	C
0	2	2	A	B	R	A							
1	0	1		A	B	R	A						
2	1	3			A	B	R	A					
3	0	3				A	B	R	A				
4	1	5					A	B	R	A			
5	0	5						A	B	R	A		
6	4	10							A	B	R	A	

entries in red are mismatches

entries in gray are for reference only

entries in black match the text

return i when j is M

match

Brute-force substring search : worst case

- Brute-force algorithm can be slow if text and pattern are repetitive.
- Worst case happens when each time all $m-1$ characters match and the last one does not
- *character compare :*

- *typical:* N
- *worst case :* $(N - M + 1) * M$
- *worst case :* $N * M$

<i>i</i>	<i>j</i>	<i>i+j</i>	0	1	2	3	4	5	6	7	8	9
			txt → A	A	A	A	A	A	A	A	A	B
0	4	4	A	A	A	A	B	← pat				
1	4	5		A	A	A	A	B				
2	4	6			A	A	A	A	B			
3	4	7				A	A	A	A	B		
4	4	8					A	A	A	A	B	
5	5	10						A	A	A	A	B

↑
match

Brute-force substring search

```
int search( char PAT[], char TEXT[])
{
    int i, j, N, M;
    N = strlen(TEXT);
    M = strlen(PAT);
    for( i = 0; i <= N-M ; i++)
    {
        j = 0;
        while(( j < M ) && (PAT[ j ] == TEXT[ i+j ] ))
            j++;
        if(j==M)
            return i; // index in text where pattern starts
    }
    return -1; // not found
}
```



Sublinear search algorithms

- move right to left in pattern ←
- shift left to right in text →
- Search time proportional to N/M for typical problems

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
		F	I	N	D	I	N	A	H	A	Y	S	T	A	C	K	N	E	E	D	L	E	I	N	A
		<i>text</i> →																							
0	5	N	E	E	D	L	E	← pattern																	
5	5						N	E	E	D	L	E													
11	4												N	E	E	D	L	E							
15	0												N	E	E	D	L	E							

return i = 15

Boyer-Moore Algorithm

.....

i= 0;

J = M - 1 ;

.....

while ((j >= 0) && (p[j] == t[i+j])) //scan text from right to left

{

 j--;

}

.....

1	J	0	1	2	3	4	5	6	7	8	9	10

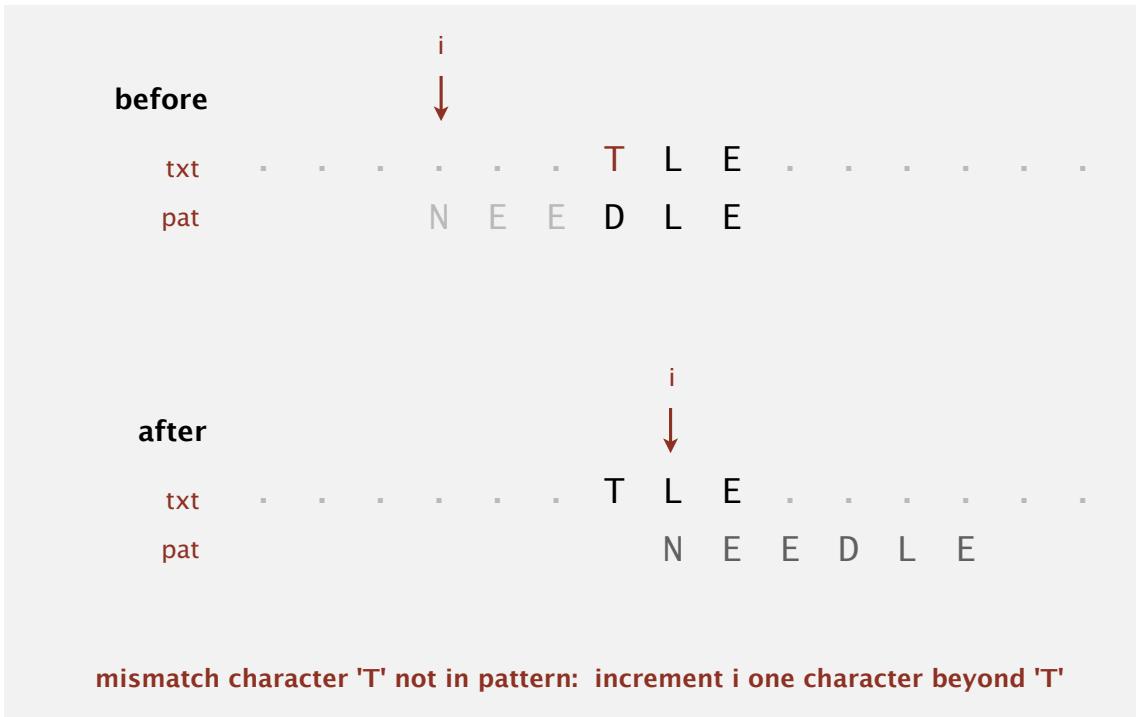
text → F I N D I **N** A H A Y S

0 5 N E E D L **E** ← *pattern*

Mismatch character - How much to skip ?

- *Case 1* : Mismatch character not in pattern

$$i = i + (j + 1)$$



5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21
N	A	H	A	Y	S	T	A	C	K	N	E	E	D	L	E	I
N	E	E	D	L	E	N	E	E	D	L	E					

before : $i = 8, j = 3$
after : $i = 8 + 3 + 1 = 12$

Mismatch character - How much to skip ?

- Case 2-a : Mismatch character in the pattern

$$i = i + (j - \text{rightMostOccurrenceInPATTERN})$$

	i	
before	↓	
txt	.	N L E
pat	N E E D L E	
after	↓	
txt	.	N L E
pat	N E E D L E	

mismatch character 'N' in pattern: align text 'N' with rightmost pattern 'N'

0	1	2	3	4	5	6	7	8	9	10
F	I	N	D	I	N	A	H	A	Y	S
N	E	E	D	L	E					

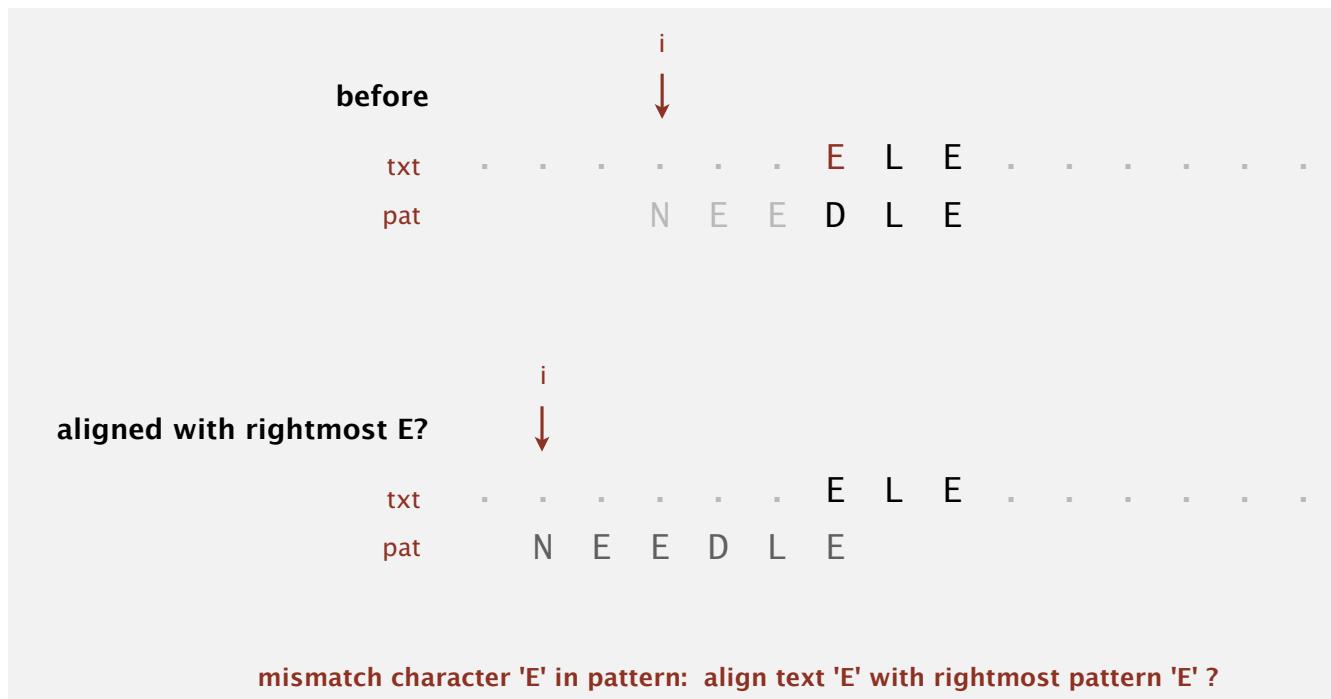
E ← pattern

before : $i=6, j=3$
after : $i = 6 + 3 - 0 = 9$

Mismatch character - How much to skip ?

- Case 2-b : Mismatch character in pattern, but skip is negative ??

$$i = i + (j - \text{rightMostOccurrenceInPATTERN})$$



before : $i = 16, j = 3$

after : $i = 16 + (3 - 5) ?$

Boyer-Moore Algorithm - How much to skip ?

- *Case 2-b* : Mismatch character in pattern, but shift would **not increase i**

before	i	↓												
txt	E	L	E
pat	N	E	E	D	L	E								

just increment i by 1

$$i = i + 1$$

aligned with rightmost E?	i	↓												
txt	E	L	E
pat	N	E	E	D	L	E								

mismatch character 'E' in pattern: align text 'E' with rightmost pattern 'E' ?

after	i	↓												
txt	E	L	E
pat							N	E	E	D	L	E		

mismatch character 'E' in pattern: increment i by 1

Boyer-Moore Algorithm- How much to skip ?

- Do not search every time for an occurrence of a character in the pattern
 - Precompute index of rightmost occurrence of character **c** in pattern

NEEDLE →

Skip Table Computation Algorithm

```
void shiftTable(int PAT[], int right[], int offset, int size)
{
    int i, M = strlen(PAT);
    for(i=0; i < size ; i++)
        right[i] = -1; // -1 if character not in pattern
    for(i=0; i < M; i++)
        right[P[i]-offset] = i ; // offset is 'A' for this example
```

}

NEEDLE →

Skip Table Computation Algorithm

c	N	E	E	D	L	E	right[c]
	0	1	2	3	4	5	
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	2	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 Algorithm

```
int search( char PAT[], char TEXT[])
{
    int i, j, skip=0, N = strlen(TEXT), M = strlen(PAT);
    for( i = 0; i <= N-M ; i += skip)
    {
        j = M - 1;
        while(( j >= 0 ) && (PAT[ j ] == TEXT[ i+j ] ))
            j--;
        if(j >= 0)
        {
            skip = j - right [ TEXT[ i + j ]];
            if( skip < 0)
                skip = 1;
        }
        else
            return i; // index in text where pattern starts
    }
    return -1; // not found
}
```

Example

```

if(j >= 0)
{
    skip = j - right [ TEXT[ i + j ]];
    if( skip < 0)
        skip = 1;
}

```

c	N	E	E	D	L	E	right[c]
	0	1	2	3	4	5	
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	2	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

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
		F	I	N	D	I	N	A	H	A	Y	S	T	A	C	K	N	E	E	D	L	E	I	N	A
text →																									
0	5	N	E	E	D	L	E	← pattern																	
5	5																								
11	4																								
15	0																								

return i = 15

Boyer - Moore Analysis

- 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 .
- worst case : $N * M$

i	skip	0	1	2	3	4	5	6	7	8	9
	txt →	B	B	B	B	B	B	B	B	B	B
0	0	A	B	B	B	B	← pat				
1	1		A	B	B	B	B				
2	1			A	B	B	B	B			
3	1				A	B	B	B	B		
4	1					A	B	B	B	B	
5	1						A	B	B	B	B

Rabin – Karp Fingerprint Algorithm

- We want to convert an M-digit base-R number to an int value between 0 to Q-1.
- Modular Hashing :

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

t : text M : length of text R : constant integer

Rabin – Karp Fingerprint Algorithm

- 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

pat.charAt(i)																
i	0	1	2	3	4											
	2	6	5	3	5											
$\% 997 = 613$																
txt.charAt(i)																
i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
	3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3
0	3	1	4	1	5	$\% 997 = 508$										
1		1	4	1	5	9	$\% 997 = 201$									
2			4	1	5	9	2	$\% 997 = 715$								
3				1	5	9	2	6	$\% 997 = 971$							
4					5	9	2	6	5	$\% 997 = 442$						
5						9	2	6	5	3	$\% 997 = 929$					
6	$\leftarrow \text{return } i = 6$					2	6	5	3	5	$\% 997 = 613$		match			

Compute hash for M digit key

```
int hash(char PAT[], int R, int Q, int offset)
{
    int i, h = 0, M = strlen(PAT);
    for(i=0; i<M; i++)
        h = ( R * h + (PAT[i] – offset)) ;
    h = h % Q ;
    return h;
}
```

Efficiently computing the hash function

- How to efficiently compute $x[i+1]$ given that we know $x[i]$.

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

$$x_{i+1} = t_{i+1} R^{M-1} + t_{i+2} R^{M-2} + \dots + t_{i+M} R^0$$

Efficiently computing the hash function

- How to efficiently compute $x[i+1]$ given that we know $x[i]$.

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

$$x_{i+1} = t_{i+1} R^{M-1} + t_{i+2} R^{M-2} + \dots + t_{i+M} R^0$$

- Key idea : Can update hash function in constant time!

Efficiently computing the hash function

- Key idea : Can update hash function in constant time!

$$x_{i+1} = (x_i - t_i R^{M-1}) R + t_{i+M}$$

can precompute

i	...	2	3	4	5	6	7	...
current value	1	4	1	5	9	2	6	5
new value		4	1	5	9	2	6	5

↔ text

$$\begin{array}{r} 4 & 1 & 5 & 9 & 2 & \text{current value} \\ - 4 & 0 & 0 & 0 & 0 \\ \hline 1 & 5 & 9 & 2 & & \text{subtract leading digit} \\ & * & 1 & 0 & & \text{multiply by radix} \\ 1 & 5 & 9 & 2 & 0 & \\ & & & + & 6 & \text{add new trailing digit} \\ 1 & 5 & 9 & 2 & \textcolor{red}{6} & \text{new value} \end{array}$$

Rabin – Karp Substring Search Example

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
	3	1	4	1	5	9	2	6	5	3	5	8	9	7	9	3	
0	3	1	%	997	=	3											
1	3	1	%	997	=	(3*10 + 1) % 997	=	31									
2	3	1	4	%	997	=	(31*10 + 4) % 997	=	314								
3	3	1	4	1	%	997	=	(314*10 + 1) % 997	=	150							
4	3	1	4	1	5	%	997	=	(150*10 + 5) % 997	=	508						
5	1	4	1	5	9	%	997	=	((508 + 3*(997 - 30)) * 10 + 9) % 997	=	201						
6	4	1	5	9	2	%	997	=	((201 + 1*(997 - 30)) * 10 + 2) % 997	=	715						
7	1	5	9	2	6	%	997	=	((715 + 4*(997 - 30)) * 10 + 6) % 997	=	971						
8	5	9	2	6	5	%	997	=	((971 + 1*(997 - 30)) * 10 + 5) % 997	=	442						
9	9	2	6	5	3	%	997	=	((442 + 5*(997 - 30)) * 10 + 3) % 997	=	929						
10	← return i-M+1 = 6	2	6	5	3	5	%	997	=	((929 + 9*(997 - 30)) * 10 + 5) % 997	=	613					

Q

RM

R

match

Rabin – Karp Fingerprint Algorithm Analysis

- Theory : If Q is a sufficiently large random prime (about $M * N^2$) 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 $1 / Q$.

Rabin – Karp Fingerprint Algorithm Analysis

- Monte Carlo version : Return match if hash match.
 - Always runs in linear time N
 - Extremely likely to return correct answer (but not always!).
- Las Vegas version : Check for substring match if hash match; continue search if false collision
 - Always returns correct answer.
 - Extremely likely to run in linear time N (but worst case is $M*N$).

Summary

algorithm	version	operation count		backup in input?	correct?	extra space
		guarantee	typical			
brute force	—	MN	$1.1N$	yes	yes	1
	<i>full algorithm</i>	$3N$	N/M	yes	yes	R
Boyer-Moore	<i>mismatched char heuristic only</i> <i>(Algorithm 5.7)</i>	MN	N/M	yes	yes	R
	<i>Monte Carlo</i> <i>(Algorithm 5.8)</i>	$7N$	$7N$	no	yes^\dagger	1
Rabin-Karp [†]	<i>Las Vegas</i>	$7N^\dagger$	$7N$	yes	yes	1

[†] probabilistic guarantee, with uniform hash function