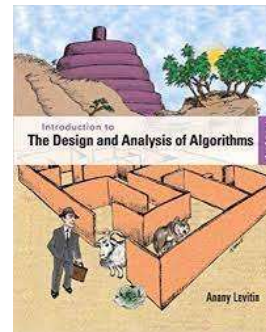


## 7-Space and Time Trade-Offs

A. Levitin "Introduction to the Design & Analysis of Algorithms," 3<sup>rd</sup> ed., Ch. 1 ©2012  
Pearson Education, Inc. Upper Saddle River,  
NJ. All Rights Reserved



1

## Space-for-time tradeoffs

Two varieties of space-for-time algorithms:

- input enhancement — preprocess the input (or its part) to store some info to be used later in solving the problem
  - counting sorts (lab)
  - string searching algorithms <
- prestructuring — preprocess the input to make accessing its elements easier
  - hashing (key, element)  $\rightarrow$   $\Theta(1)$
  - indexing schemes (e.g., B-trees)

30.11.2021

BLM19307E Algorithm Analysis and Design

2

2

## Review: String searching by brute force

*pattern*: a string of  $m$  characters to search for

*text*: a (long) string of  $n$  characters to search in

### Brute force algorithm

- Step 1 Align pattern at beginning of text
- Step 2 Moving from left to right, compare each character of pattern to the corresponding character in text until either all characters are found to match (successful search) or a mismatch is detected
- Step 3 While a mismatch is detected and the text is not yet exhausted, realign pattern one position to the right and repeat Step 2

30.11.2021

BLM19307E Algorithm Analysis and Design

3

3

## String searching by preprocessing

Several string searching algorithms are based on the input enhancement idea of preprocessing the pattern

- Knuth-Morris-Pratt (KMP) algorithm preprocesses pattern left to right to get useful information for later searching

④ Boyer-Moore algorithm preprocesses pattern right to left and store information into two tables

④ Horspool's algorithm simplifies the Boyer-Moore algorithm by using just one table

30.11.2021

BLM19307E Algorithm Analysis and Design

4

4

# Horspool's Algorithm

A simplified version of Boyer-Moore algorithm:

- preprocesses pattern to generate a shift table that determines how much to shift the pattern when a mismatch occurs
- always makes a shift based on the text's character *c* aligned with the last character in the pattern according to the shift table's entry for *c*

30.11.2021

BLM19307E Algorithm Analysis and Design

5

5

## How far to shift?

Look at first (rightmost) character in text that was compared:

- The character is not in the pattern

*not* ⇒  $\dots\dots\dots c \dots\dots\dots$  (*c* not in pattern)  
 $\begin{array}{c} \text{BAOBAB} \\ \text{BAOBAB} \end{array}$  (shift the pattern by its entire length)

- The character is in the pattern (but not the rightmost)

$\dots\dots\dots O \dots\dots\dots$  (*O* occurs once in pattern)  
 $\begin{array}{c} \text{BAOBAB} \\ \text{BAOBAB} \end{array}$  shift by 3.

$\dots\dots\dots A \dots\dots\dots$  (*A* occurs twice in pattern)  
 $\begin{array}{c} \text{BAOBAB} \\ \text{BAOBAB} \end{array}$  shift by 1

- The rightmost characters do match

$\begin{array}{c} \text{ABAB} \\ \text{BAOBAB} \end{array}$  shift by 2

30.11.2021

BLM19307E Algorithm Analysis and Design

6

6

## How far to shift?

Four possibilities may occur:

**Case 1** If there are no  $c$ 's in the pattern—e.g.,  $c$  is letter S in our example—we can safely shift the pattern by its entire length (if we shift less, some character of the pattern would be aligned against the text's character  $c$  that is known not to be in the pattern):

```

s0 ... S ... s_{n-1}
//
B A R B E R
B A R B E R
  
```

**Case 2** If there are occurrences of character  $c$  in the pattern but it is not the last one there—e.g.,  $c$  is letter B in our example—the shift should align the rightmost occurrence of  $c$  in the pattern with the  $c$  in the text:

```

s0 ... B ... s_{n-1}
//
B A R B E R
B A R B E R
  
```

**Case 3** If  $c$  happens to be the last character in the pattern but there are no  $c$ 's among its other  $m - 1$  characters—e.g.,  $c$  is letter R in our example—the situation is similar to that of Case 1 and the pattern should be shifted by the entire pattern's length  $m$ :

```

s0 ... M E R ... s_{n-1}
//
L E A D E R
L E A D E R
  
```

**Case 4** Finally, if  $c$  happens to be the last character in the pattern and there are other  $c$ 's among its first  $m - 1$  characters—e.g.,  $c$  is letter R in our example—the situation is similar to that of Case 2 and the rightmost occurrence of  $c$  among the first  $m - 1$  characters in the pattern should be aligned with the text's  $c$ :

```

s0 ... A R ... s_{n-1}
//
R E O R D E R
R E O R D E R
  
```

30.11.2021

BLM19307E Algorithm Analysis and Design

7

7

## Shift table

Shift sizes can be precomputed by the formula

$$t(c) = \begin{cases} \text{the pattern's length } m, & \text{if } c \text{ is not among the first } m - 1 \text{ characters of the pattern;} \\ \text{the distance from the rightmost } c \text{ among the first } m - 1 \text{ characters} & \\ \text{of the pattern to its last character, otherwise.} & \end{cases} \quad (7.1)$$

by scanning pattern before search begins and stored in a table called *shift table*

**ALGORITHM** *ShiftTable*( $P[0..m - 1]$ )

//Fills the shift table used by Horspool's and Boyer-Moore algorithms

//Input: Pattern  $P[0..m - 1]$  and an alphabet of possible characters

//Output:  $Table[0..size - 1]$  indexed by the alphabet's characters and

// filled with shift sizes computed by formula (7.1)

**for**  $i \leftarrow 0$  **to**  $size - 1$  **do**  $Table[i] \leftarrow m$

**for**  $j \leftarrow 0$  **to**  $m - 2$  **do**  $Table[P[j]] \leftarrow m - 1 - j$

**return**  $Table$

30.11.2021

BLM19307E Algorithm Analysis and Design

8

8

## Shift table

- Shift table is indexed by text and pattern alphabet  
Eg, for BAOBAB :

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	2	6	6	6	6	6	6	6	6	6	6	6	6	3	6	6	6	6	6	6	6	6	6	6	6

BAOBAB

BAOBAB

30.11.2021

BLM19307E Algorithm Analysis and Design

9

9

## Example of Horspool's alg. application

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	2	6	6	6	6	6	6	6	6	6	6	6	6	3	6	6	6	6	6	6	6	6	6	6	6

t: BARD LOVED BANANAS  
 p: BAOBAB shift by 6  
 BAOBAB shift by 2  
 BAOBAB shift by 6  
 BAOBAB (unsuccessful search)

30.11.2021

BLM19307E Algorithm Analysis and Design

10

10



## Boyer-Moore algorithm

Based on same two ideas:

- comparing pattern characters to text from right to left
- precomputing shift sizes in two tables
  - bad-symbol table indicates how much to shift based on text's character causing a mismatch
  - good-suffix table indicates how much to shift based on matched part (suffix) of the pattern

30.11.2021

BLM19307E Algorithm Analysis and Design

13

13

## Bad-symbol shift in Boyer-Moore algorithm

- If the rightmost character of the pattern doesn't match, BM algorithm acts as Horspool's
- If the rightmost character of the pattern does match, BM compares preceding characters right to left until either all pattern's characters match or a mismatch on text's character  $c$  is encountered after  $k > 0$  matches

$s_0$	...	$c$	$s_{i-k+1}$	...	$s_i$	...	$s_{n-1}$	text
		$\parallel$	$\parallel$		$\parallel$			
$p_0$	...	$p_{m-k-1}$	$p_{m-k}$	...	$p_{m-1}$			pattern

bad-symbol shift  $d_1 = \max\{t_1(c) - k, 1\}$

shift table.

(1)  $t_1(c) = 6$

(2)  $t_1(c) = 2$

$k=3$

$6-3=3$

$2-3=-1$

30.11.2021

BLM19307E Algorithm Analysis and Design

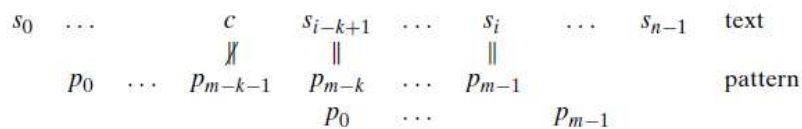
14

14

## Bad-symbol shift in Boyer-Moore algorithm

bad-symbol shift  $d_1 = \max\{t_1(c) - k, 1\}$

$t_1(c)$  is the entry in the precomputed table used by Horspool's algorithm and  $k$  is the number of matched characters



30.11.2021

BLM19307E Algorithm Analysis and Design

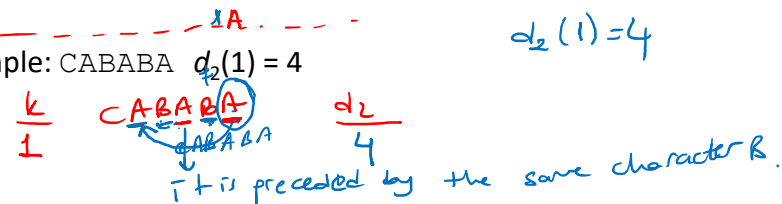
15

15

## Good-suffix shift in Boyer-Moore algorithm

- Good-suffix shift  $d_2$  is applied after  $0 < k < m$  last characters were matched
- $d_2(k)$  = the distance between matched suffix of size  $k$  and its rightmost occurrence in the pattern that is not preceded by the same character as the suffix

Example: CABABA  $d_2(1) = 4$



I'm looking at the occurrence of A where the preceding character is different than B.

30.11.2021

BLM19307E Algorithm Analysis and Design

16

16



<u>k</u>		<u>dL</u>	
1	CABABA	4	the preceding char is the same (char. A) → shift by the whole size of the pattern.
2	CABABA	6	
3	CABABA	2	
4	CA BABA	6	
5	CA BABA	6	

30.11.2021 BLM19307E Algorithm Analysis and Design 17

17

## Good-suffix shift in Boyer-Moore algorithm

- If there is no such occurrence, match the longest part of the  $k$ -character suffix with corresponding prefix; if there are no such suffix-prefix matches,  $d_2(k) = m$

Example: WOWWOW  $d_2(2) = 5$ ,  $d_2(3) = 3$ ,  $d_2(4) = 3$ ,  
 $d_2(5) = 3$

30.11.2021

BLM19307E Algorithm Analysis and Design

18

18

$k$		$d_1$	$k$		$d_2$
1	wowwow	2	1	BAOBAB	2
2	wowwow	5	2	BAOBAB	5
3	wowwow	3	3	BAOBAB	5
4	wowwow	3	4	BAOBAB	5
5	wowwow	3	5	BAOBAB	5

30.11.2021 BLM19307E Algorithm Analysis and Design 19

19

## Good-suffix shift in the Boyer-Moore alg. (cont.)

After matching successfully  $0 < k < m$  characters, the algorithm shifts the pattern right by

$$d = \max \{d_1, d_2\}$$

where  $d_1 = \max\{t_1(c) - k, 1\}$  is bad-symbol shift

$d_2(k)$  is good-suffix shift

30.11.2021

BLM19307E Algorithm Analysis and Design

20

20

## Boyer-Moore Algorithm (cont.)

- Step 1 Fill in the bad-symbol shift table  
 Step 2 Fill in the good-suffix shift table  
 Step 3 Align the pattern against the beginning of the text  
 Step 4 Repeat until a matching substring is found or text ends:  
 Compare the corresponding characters right to left.  
 If no characters match, retrieve entry  $t_1(c)$  from the bad-symbol table for the text's character  $c$  causing the mismatch and shift the pattern to the right by  $t_1(c)$ .  
 If  $0 < k < m$  characters are matched, retrieve entry  $t_1(c)$  from the bad-symbol table for the text's character  $c$  causing the mismatch and entry  $d_2(k)$  from the good-suffix table and shift the pattern to the right by
- $$d = \max \{d_1, d_2\}$$
- where  $d_1 = \max \{t_1(c) - k, 1\}$ .

30.11.2021

BLM19307E Algorithm Analysis and Design

21

21

## Example of Boyer-Moore alg. application

*t(c) shift table.*

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	-
1	2	6	6	6	6	6	6	6	6	6	6	6	6	3	6	6	6	6	6	6	6	6	6	6	6	6

*BAOBAB*

$d_1 = \max \{t_1(K) - k, 1\}$

$d_1 = t_1(K) = 6$

k	pattern	d <sub>2</sub>
1	BAOBAB	2
2	BAOBAB	5
3	BAOBAB	5
4	BAOBAB	5
5	BAOBAB	5

$d_1 = t_1(\_) - 2 = 4 = 6 - 2$

$d_2(2) = 5$

$d_1 = t_1(\_) - 1 = 5$

$d_2(1) = 2$

$\max \{4, 5\} = 5$

$\max \{5, 2\} = 5$

BAOBAB (success)

30.11.2021

BLM19307E Algorithm Analysis and Design

22

22

## Boyer-Moore example from their paper

Find pattern AT\_THAT in

*text* WHICH\_FINALLY\_HALTS. \_\_AT\_THAT

shift table

A	T	H	-	others
1	3	2	4	7

4 length

30.11.2021

BLM19307E Algorithm Analysis and Design

23

23

Horspool's Alg.

WHICH\_FINALY HALTS. \_\_AT\_THAT

AT-THAT shift by 7(F)  
 AT-THAT shift by 4(-)  
 AT-THAT shift by 3(T)  
 AT-THAT shift by 4(-)  
 AT-THAT shift by 4  
AT-THAT

30.11.2021

BLM19307E Algorithm Analysis and Design

24

24

Boyer-Moore  
 WHICH\_FINALY\_HALTS. \_\_AT\_THAT

shift table

A	T	H	.	others
1	3	2	4	7

4 length

good suffix.

k	AT-T#AT	ch
1	AT-T#AT	3
2	AT-T#AT	5
3	AT-T#AT	5
4	AT-T#AT	5
5	AT-T#AT	5
6	AT-T#AT	5

→ Bad symbol  
 $d_1 = \max(t_i(c) - k)$   
 13

30.11.2021 BLM19307E Algorithm Analysis and Design 25

25

Boyer-Moore  
 WHICH\_FINALY\_HALTS. \_\_AT\_THAT

AT-T#AT shift by 7 #  
 AT-T#AT shift by 4 #  
 AT-T#AT shift by 6 #  
 AT-T#AT shift by 5 #

k=1  
 $d_2(1)=3$   
 $d_1 = t_1(L) - k$   
 $d_1 = 7 - 1 = 6$   
 $\max\{d_1, d_2\} = 6$

k=2  
 $d_2(2)=5$   
 $d_1 = t_1(L) - k$   
 $d_1 = 4 - 2 = 2$   
 $\max\{d_1, d_2\} = 5$

AT-T#AT shift by 5 #  
 AT-T#AT shift by 5 #

Boyer-Moore = 7 + 7 = 14 comparisons  
 Horspool's = 6 + 7 = 13 comparisons.

30.11.2021 BLM19307E Algorithm Analysis and Design 26

26

# Hashing

- A very efficient method for implementing a *dictionary*, i.e., a set with the operations:
  - find
  - insert
  - delete
- Based on representation-change and space-for-time tradeoff ideas
- Important applications:
  - symbol tables
  - databases (*extendible hashing*)

30.11.2021

BLM19307E Algorithm Analysis and Design

27

27

## Hash tables and hash functions

The idea of *hashing* is to map keys of a given file of size  $n$  into a table of size  $m$ , called the *hash table*, by using a predefined function, called the hash function,

$$h: K \rightarrow \text{location (cell) in the hash table}$$

*key*

Example: student records, key = SSN. Hash function:

$$h(K) = K \bmod m \text{ where } m \text{ is some integer (typically, prime)}$$

If  $m = 1000$ , where is record with SSN= 314159265 stored?

$$\rightarrow h(K) = 265$$

Generally, a hash function should:

- be easy to compute
- distribute keys about evenly throughout the hash table

30.11.2021

BLM19307E Algorithm Analysis and Design

28

28

# Collisions



- If  $h(K1) = h(K2)$ , there is a collision
- Good hash functions result in fewer collisions, but some collisions should be expected (birthday paradox)
- Two principal hashing schemes handle collisions differently:
  - ⊗ Open hashing
    - each cell is a header of linked list of all keys hashed to it
  - ⊗ Closed hashing
    - one key per cell
    - in case of collision, finds another cell by
      - linear probing: use next free bucket
      - double hashing: use second hash function to compute increment

30.11.2021

BLM19307E Algorithm Analysis and Design

29

29

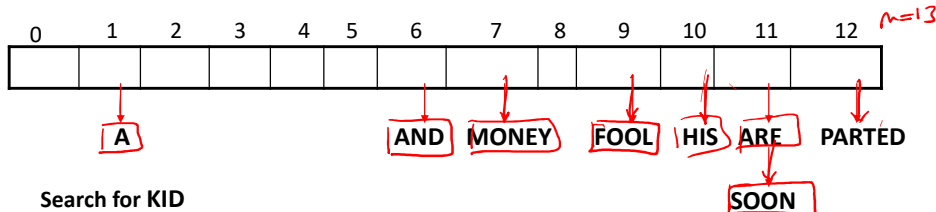
## Open hashing (Separate chaining)

Keys are stored in linked lists outside a hash table whose elements serve as the lists' headers.

Example: A, FOOL, AND, HIS, MONEY, ARE, SOON, PARTED

$h(K)$  = sum of K 's letters' positions in the alphabet MOD 13

Key	A	FOOL	AND	HIS	MONEY	ARE	SOON	PARTED
$h(K)$	1	9	6	10	7	11	11	12



30.11.2021

BLM19307E Algorithm Analysis and Design

30

30

## Open hashing (cont.)

- If the hash function distributes  $n$  keys among  $m$  cells of the hash table about evenly, each list will be about  $n/m$  keys long. This ratio ( $\alpha = n/m$ ) is called load factor.
- Average number of probes in successful,  $S$ , and unsuccessful searches,  $U$ :  

$$S \approx 1 + \alpha/2, \quad U = \alpha$$
- Load  $\alpha$  is typically kept small (ideally, about 1)
- Open hashing still works if  $n > m$

30.11.2021

BLM19307E Algorithm Analysis and Design

31

31

## Closed hashing (Open addressing)

- Keys are stored inside a hash table. *linear probing (next available slot)*

keys	A	FOOL	AND	HIS	MONEY	ARE	SOON	PARTED
hash addresses	1	9	6	10	7	11	11	12

	0	1	2	3	4	5	6	7	8	9	10	11	12
		A											
		A								FOOL			
		A				AND				FOOL			
		A				AND				FOOL	HIS		
		A				AND	MONEY			FOOL	HIS		
		A				AND	MONEY			FOOL	HIS	ARE	
		A				AND	MONEY			FOOL	HIS	ARE	SOON
		PARTED				AND	MONEY			FOOL	HIS	ARE	SOON

30.11.2021

BLM19307E Algorithm Analysis and Design

32

32



## Closed hashing (cont.)

- Does not work if  $n > m$
- Avoids pointers
- Deletions are not straightforward
- Number of probes to find/insert/delete a key depends on load factor  $\alpha = n/m$  (hash table density) and collision resolution strategy. For linear probing:
- $S = (\frac{1}{2}) (1 + 1/(1 - \alpha))$  and  $U = (\frac{1}{2}) (1 + 1/(1 - \alpha)^2)$
- As the table gets filled ( $\alpha$  approaches 1), number of probes in linear probing increases dramatically:

$\alpha$	$\frac{1}{2}(1 + \frac{1}{1-\alpha})$	$\frac{1}{2}(1 + \frac{1}{(1-\alpha)^2})$
50%	1.5	2.5
75%	2.5	8.5
90%	5.5	50.5

30.11.2021

BLM19307E Algorithm Analysis and Design

33

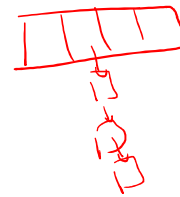
33

Linear probing.

$$h'(a_k, i) = (h(a_k) + i) \bmod m$$

$m$  : size

$i$  : # of collisions.



with Hashing : searching  
average insertion  
deletion

$$\left. \begin{array}{l} \text{searching} \\ \text{insertion} \\ \text{deletion} \end{array} \right\} \begin{array}{l} \Theta(?) \\ \Theta(1) \end{array}$$

① unlikely worst case

$$\left. \begin{array}{l} \text{searching} \\ \text{insertion} \\ \text{deletion} \end{array} \right\} \rightarrow \begin{array}{l} \Theta(?) \\ \underline{\underline{\Theta(n)}} \end{array}$$

30.11.2021

BLM19307E Algorithm Analysis and Design

34

34