

Lecture 13: Hashing

Last few days: Collision Resolution

- Open addressing
- Separate overflow

Today: Updating Sorted Files

- Differential Files
- Bloom Filters (a novel application of hashing)

See handout on hashing.

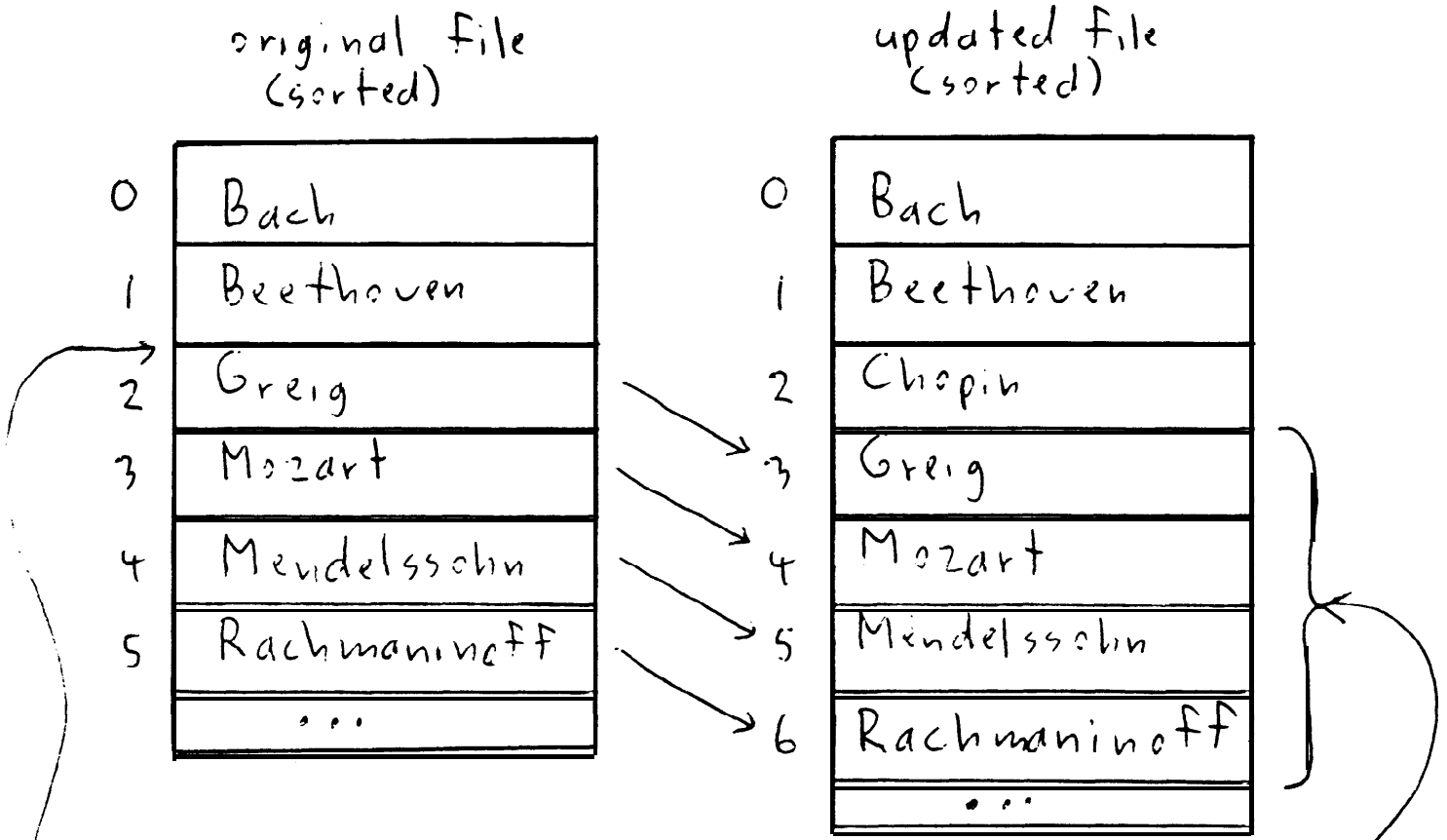
## Introduction

Goal: To efficiently update a file while keeping it sorted.

Problem: Inserts and deletes to sorted files are expensive, since much of the file may have to be shifted.

# Example

Insert a record into a sorted file.



Insert Chopin.

These records have all been shifted down by one position.

## Updating Sorted Files

Deletes: Physical deletion is expensive since much of the file must be shifted up by one position to fill in the gap created by the deleted record.

Instead, we can "mark" a record for deletion by putting a D in a special field, instead of physically deleting it.

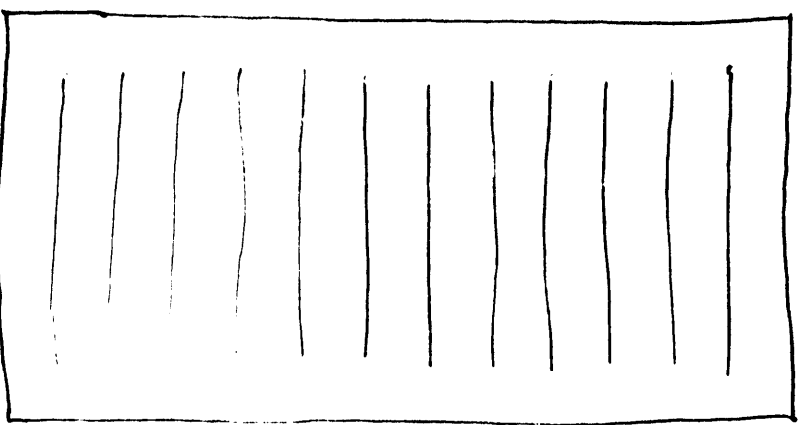
Modifications: These are easy.

eg. we can modify Greig's birth date or salary without shifting any records.

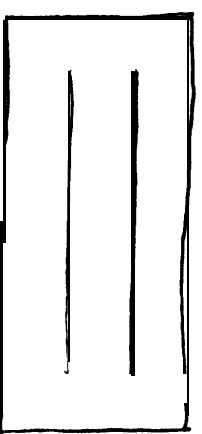
Inserts: These are a problem.

Solution: Differential Files

Master File  
(sorted)



Differential File  
(unsorted)



Records to be inserted into the master file are appended to the differential file instead.

## Retrieval

- First, search the differential file (sequentially).  
If the record is found, then we are done.  
If the record is not found, then ...
- Search the Master file (possibly using binary search).  
If the record is not found, then it was never inserted into the file.

Note: Searching the diff'l file requires disk accesses

- Also, most records are not in the diff'l file.
- So, we do not want to search the diff'l file unless we have good reason to think the record is there

## Solution: Bloom Filters

- Bloom Filtering is a hashing-based technique to avoid unnecessary searches of a differential file.
- It is an (imperfect) test for the presence of a key in a diff'l file.
- The test outcome is either:
  - No, the <sup>record</sup>  $\checkmark$  is not in the file,
  - OR
  - Maybe the record is in the file.

Bloom Filters: Inserting records

- Keep an array,  $B(n)$ , of  $N$  bits in main memory
- Use a number of hash functions,  
 $h_1, h_2, \dots, h_t$   
where  $0 \leq h_i(\text{key}) \leq N-1$  for each  $i$ .
- When inserting a record with key  $K$ ,  
first set bits  $h_1(K), \dots, h_t(K)$  of the  
array to 1.  
Then insert the record into the diff'l file.



## Bloom Filters: Retrieving records

To retrieve a record with key  $K$ ,

- First test bits  $b_1(K), \dots, b_t(K)$  of the array

There are two cases:

(a) If any of these bits is 0, then the record is not in the differential file. So, do not search the differential file, but search the master file instead.

(b) If all these bits are 1, then the record may be in the diff'l file, so we must search it.

If the record is found in the diff'l file, then we are done; otherwise we must also search the master file.

# Bloom Filters : Example

Master File  
(disk)

Bach
Beethoven
Greig
Mendelssohn
Rachmaninoff
Tchaikovsky
Vivaldi

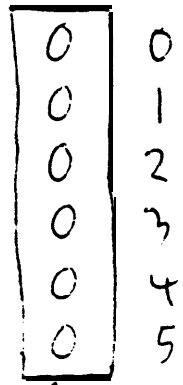
sorted

Differential File  
(disk)



initially empty

Array, B  
(mm)



All bits are initially zero

<u>key</u>	<u><math>h_1(\text{key})</math></u>	<u><math>h_2(\text{key})</math></u>
Bach	0	3
Chopin	3	5
Mendelssohn	1	5
Mozart	1	3
Rachmaninoff	2	4

Program Filters: Example (cont.)Insert

Mozart

Chopin

RetrieveIn Diff. file?Search Diff. File?

Mozart

maybe

yes

Bach

no

no

Mendelssohn

maybe

yes

DeleteIn Diff. File?Search Diff. File?

Rachmaninoff

no

no

Mendelssohn

maybe

yes

Chopin

maybe

yes

Bloom Filters: Example (cont.)

Master File  
(disk)

	Bach _____
	Beethoven _____
D	Greig _____
D	Mendelssohn _____
D	Rachmaninoff _____
	Tchaikovsky _____
	Vivaldi _____

sorted

Differential File  
(disk)

	Mozart _____
D	Chopin _____

unsorted

Array, B  
(M)

0	0
1	1
2	0
3	1
4	0
5	1

New File and array contents

## Merging

- When most of the bits in the array,  $B$ , are 1, then the Bloom filter is useless since you will have to search the differential file most of the time.
- At this point, the differential file is merged with the master file, to produce a new, sorted master file, in which records marked D are physically deleted.
- Also, a new, empty differential file is created, and all array bits are set to 0.

Bloom Filters: Example (cont.)

New Master File  
(disk)

Bach
Beethoven
Greig
Mozart
Tchaikovsky
Vivaldi

sorted

New Diff'l File  
(disk)



empty

Array B  
(M/M)

0	0
0	1
0	2
0	3
0	4
0	5

Final File and array contents, after merging