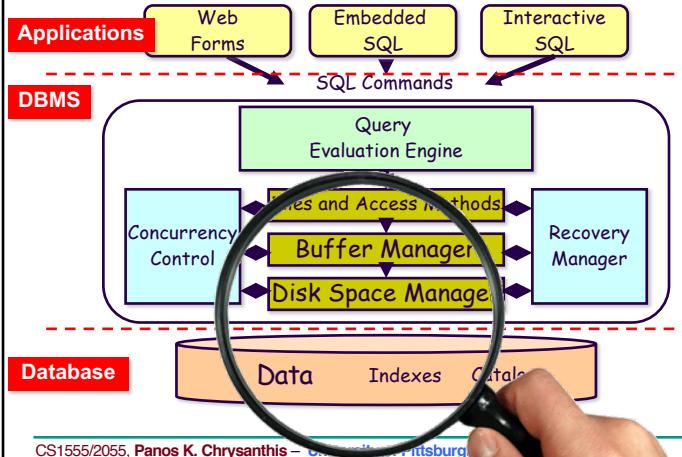


Data Storage

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

1

Database Management System (DBMS)



2

Data Storage

- ❑ Two DBMS fundamental questions?
 1. How do we store and manage very large volumes of data?
 2. What representation and data structures best support efficient manipulation of data?
 - RAM model vs. I/O model of Computation

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

3

Storage Hierarchy

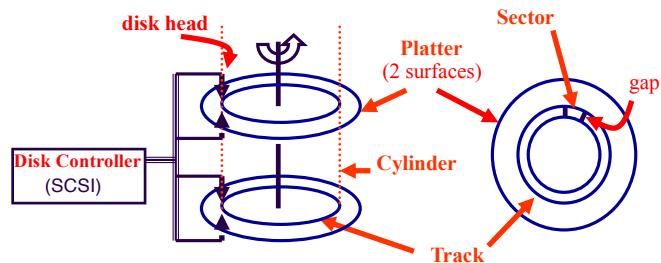
Speed; \$\$

- ❑ Primary Storage: random access; volatile
 - Cache - on board or level-2 cache
 - Main memory
- ❑ Secondary storage: random access; non-volatile
 - Flash-based or Solid State Disk
 - Magnetic disk
 - Virtual Memory (Main-memory DBS), File System, DBMS
- ❑ Tertiary storage: non-volatile
 - Optical disk / juke boxes – random access
 - Magnetic cartridge / tape silos – seq. access

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

4

Magnetic Disks



- ❑ *Disk block* is the unit of transfer: a disk *block* is a continuous sequence of sectors from a single track of one platter.
- ❑ Physical address of $x: <\text{volume}\#, \text{platter}\#, \text{track}\#, \text{sector}\#>$
- ❑ Block address of $x: <\text{volume}\#, \text{sector}/\text{block}\#>$

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

8

Accessing a Disk Block

- ❑ Time to access (read/write) a disk block:
 1. **seek time:**
 - moving arms to position disk head on track
 2. **rotational delay:**
 - waiting for sector to rotate under head
 3. **transfer time:**
 - actually moving data to/from disk surface

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

9

Performance Measures of Disks

For any block of data,

Average block access time =

seek time + rotational delay + transfer time

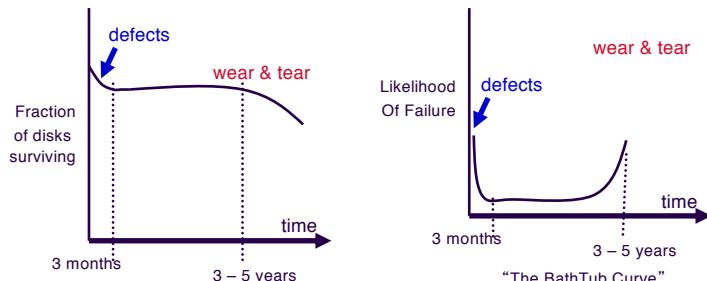
- ❑ Seek time and rotational delay dominate!
- ❑ **Disk Access** = *Seek time* + *Rotation latency*
 - <2000 8 - 20 ms = 2-10 ms + 4-10 ms
 - now 2 -10 ms = 1-6 ms + 1-4 ms
- ❑ **Data Transfer Rate** = 4-8MB/sec (max 25-40MB/sec)
 - Block transfer is less than 0.5 msec
 - Multiple disks with share interface can support high rates: 33MB/s(Ultra-DMA), 40MB/s(SCSI-3), 400MB/s(FibreChannel)

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

10

Disk Failures

- ❑ Mean Time Between Failure (MTTF)
 - Vendor's claim: 30K hrs (3.6 yrs) – 1.2M hrs (136 yrs)
 - Practice: 1000 disks with 1.2M hrs MTTF then on average 1 disk will fail in 120 hrs.



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

11

But first, how is data stored on disk?



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

12

Data Elements

- ❑ **Field**: a database attribute (sequence of bytes)
- ❑ **Record**: sequence of fields

```
CREATE TABLE Student (
    Sid INTEGER,
    Name CHAR (24),
    Age INTEGER
);
```

0	3	4	27	28	31
Sid	Name				Age

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

13

Data Elements

- ❑ **Field**: a database attribute (sequence of bytes)
- ❑ **Record**: sequence of fields

```
CREATE TABLE Student (
    Sid INTEGER,
    Name CHAR (24),
    Age INTEGER
);
```

0	3	4	27	28	31
Sid	Name				Age

- ❑ **Block**: sequence of records
- ❑ **File**: sequence of blocks

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

14

Blocks & Files

```
CREATE TABLE Student (
    Sid INTEGER,
    Name CHAR (24),
    Age INTEGER
);
```

SID	Name	Age
546007	Peter	18
546100	Bob	19
546107	Ann	21
546207	Jane	20
546240	John	24
546350	Ben	18
546420	Suzy	27
546500	Peter	20

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

15

Blocks & Files

```
CREATE TABLE Student (
    Sid INTEGER,
    Name CHAR(24),
    Age INTEGER
);
```

	SID	Name	Age
Record 0	546007	Peter	18
	546100	Bob	19
	546107	Ann	21
	546207	Jane	20
Record 4	546240	John	24
Record 5	546350	Ben	18
Record 6	546420	Suzy	27
Record 7	546500	Peter	20

Block Header	Record 0	Record 1	Record 2	Record 3
--------------	----------	----------	----------	----------

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

16

Improving Access Time ?

- ❑ What are the costs in I/O requests?
- ❑ I/O cost = *disk access* + *queuing delays*
- ❑ Disk access is reduced with *data organization*
- ❑ Queuing delays are reduced with *scheduling*



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

19

Improving Access Time

- ❑ Data access time could be improved by:
 1. Scheduling
 2. Block transfer
 3. Buffering and Prefetching
 4. Cylinder-based Organization
 5. Multiple disks
 6. Record Placement



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

20

Queuing Delay



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

21

Queuing Delay

- ❑ **Queuing Delay** occurs when:
 - multiple queries are executed at the same time, and
 - those queries access the same disk at the same time (i.e., send I/O requests)
- ❑ The disk head can only serve one request at a time!
- ❑ **Solution:**
 - Smart scheduling of serving disk requests
 - Parallel processing (multiple disks organizations)

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

22

Disk scheduling: Elevator Algorithm

- ❑ **Idea** : Schedule reading of requested blocks in the order in which they appear under the head
- ❑ **Elevator Algorithm (Scan)** - Works like an elevator
 - The arm moves towards **one direction** serving all requests on the visited tracks
 - Reverses direction when no more pending requests
- ❑ **Advantage:**
 - Reduces average access time for unpredictable requests
- ❑ **Problem:**
 - It is more effective in many disk-access requests
 - The benefit is not uniform among requests

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

24

Variations of Elevator Algorithm

- ❑ **Scan or Elevator Algorithm**
 - services both directions and reverses direction when reaches the edge of the disk
- ❑ **C-Scan or Circular Elevator Algorithm**
 - all requests are serviced in only one direction
 - more equal performance for all head positions
- ❑ **Look**
 - Scan but changes direction when no more pending requests
- ❑ **C-Look**
 - Combination of Look and C-Scan
- ❑ **FCScan**
 - prevents "arm stickiness" by using two queues (Q1 & Q2)
 - When Scan uses Q1 with old requests, new requests go to Q2

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

25

Improving Access Time

- ❑ Data access time could be improved by:
 1. Scheduling
 2. Block transfer
 3. Buffering and Prefetching
 4. Cylinder-based Organization
 5. Multiple disks
 6. Record Placement



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

26

Block Transfer

- ❑ To minimize access time: data is always transferred from/to disks by **blocks** of bytes
- ❑ A disk block is an adjacent sequence of sectors from a single track of one platter
- ❑ Block size typically ranges from 512 bytes to 4KBytes.



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

27

Improving Access Time

- ❑ Data access time could be improved by:
 1. Scheduling
 2. Block transfer
 3. Buffering and Prefetching
 4. Cylinder-based Organization
 5. Multiple disks
 6. Record Placement



Improving Access Time

- ❑ Data access time could be improved by:
 1. Scheduling
 2. Block transfer
 3. Buffering and Prefetching
 4. Cylinder-based Organization
 5. Multiple disks
 6. Record Placement

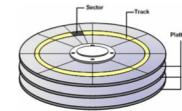


CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

31

Cylinder-based Organization

- ❑ **Observation:** Data such as in relations is likely to be accessed together
- ❑ **Idea:** Store such data on the same cylinder
 - Blocks are "**Next**" to each other
 - Blocks on a track or on a cylinder effectively involve only the first seek time and first rotational latency
- ❑ **Advantage:**
 - Approaches the theoretical transfer rate
 - Excellent if access can be predicted in advance
 - Only one process/transaction is using the disk



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

32

Multiple Disks

- ❑ **Observation:**
 - Disk drives continue to become smaller and cheaper
- ❑ **Idea:**
 - Use multiple disks to support **parallel** access
 - Disks with independent heads connected to a single disk controller
 - 2X20 GB drives is faster than a single 40GB drive
 - But more expensive...
 - More efficient if they are kept busy!
 - It also enhances system reliability thru redundancy!!



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

33

Improving Access Time

- ❑ Data access time could be improved by:
 1. Scheduling
 2. Block transfer
 3. Buffering and Prefetching
 4. Cylinder-based Organization
 5. Multiple disks
 6. Record Placement



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

34

Multiple Disks: Organizations

- ❑ Data **partitioning** over several disks
 - **Pros:** increased access rate
 - **Problem:** if popular data is stored on same disk → request collisions (hot spots)
 - **Cons:** the cost of several small disks is greater than a single one with the same capacity
- ❑ **Mirror** disks: Disks hold identical copies
 - **Pros:** Doubles read rate (disk 1 **OR** 2)
 - **Does not** have the problem of request collisions
 - **Does not** slow down write requests (disk 1 **AND** 2)
 - **Cons:** Pay cost for two disks to get the storage of one



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

35

RAID Technology

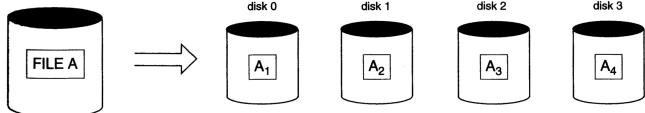


CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

36

RAID Technology

- ❑ A major advance in secondary storage technology is represented by the development of **RAID**
 - **Redundant Arrays of Inexpensive (independent) Disks**
- ❑ Acts as a single high-performance large disk
- ❑ **Data Striping:** distributes data transparently over multiple disks to make them appear as a single large, fast disk



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

37

Data Striping

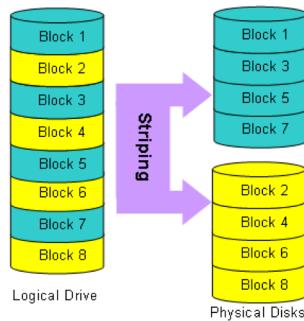
- ❑ **Bit-level striping:** Split groups of bits over different disks
 - Example: split each bit of a bytes over **8** disks
 - Bit number x is written on disk number x
 - Every disk participates in every access
 - every access reads **8 times** as many data
- ❑ **Block-level** striping for blocks of a file
- ❑ **Sector-level** striping for sectors of a block
- ❑ **RAID Goals:**
 - Parallelize accesses so the **access time is reduced**
 - **Combining** Striping, Mirroring and Reliability → levels of RAID



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

38

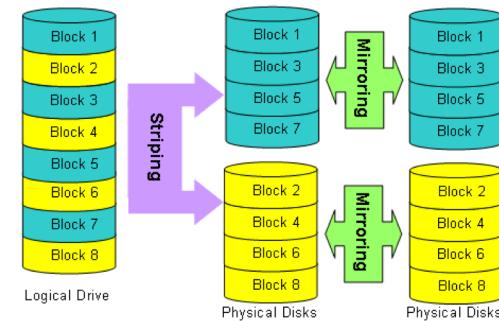
Raid 0



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

39

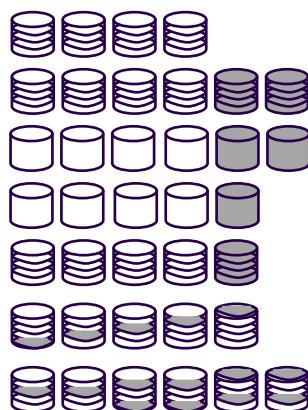
Raid 0-1



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

40

Levels of RAID



0: Non-Redundant

1: Mirrored

2: Memory Style ECC

3: Bit-Interleaved Parity

4: Block Interleaved Parity

5: Block-Interleaved Distribution-Parity

6: P+Q Redundancy

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

41

Solid State Disk

- ❑ **Idea:** A SSD is essentially an array of NAND-flash memory devices which include a controller that emulates a hard disk drive



- ❑ **Organization (flash):**
 - One or more dies (chips), each 8192 blocks
 - Dies are organized in planes (e.g., 4 planes of 2018 blocks)
 - Blocks is a set of pages 128–256KB (i.e., 32–64 pages)
 - Page is the unit of read and write (program), 512B–4KB
 - 128 bytes are used for metadata
 - Page or Block is the unit of storage system write

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

42

Constraints of Solid State Disks

Constraints:

- **Erase-Constraint** (or erase-before-write): A page cannot be over-written unless the block that contains the page has been erased
- **Write-Constraint**: Pages must be written out sequentially within a block, from low to high addresses.
 - if p_i is written, p_{i-1} cannot be written until the block is erased
- **Wear-Constraint**: The number of times a block can be erased is limited – typically 10,000 to 100,000 times.
 - Cleaning and wear-leveling schemes are important
- **Asymmetry between operations**: Write is more expensive than read (~3 times).



43

Properties of Solid State Disks

Properties:

- Random and Sequential Read Access with equal cost
- Faster than HDD
- parallel requests, interleaving, striping.. 100K – 300K I/Os
- Energy Efficient.. $\sim 25/750 \mu\text{J}$ Read/Write, $\sim 220 \mu\text{J}/\text{s}$ Idle
- Resilient & reliable
 - Shock resistant and extreme temperature fluctuations
 - immune to strong magnetic fields
 - With good wear-leveling can last 10 years
- Small factor and silent (unless fan is used)



44

Problem:

- Higher cost per megabyte of storage

- Capacity doubles every year

Improving Access Time

- ### Data access time could be improved by:
1. Scheduling
 2. Block transfer
 3. Buffering and Prefetching
 4. Cylinder-based Organization
 5. Multiple disks
 6. Record Placement



45

Two Store Approaches

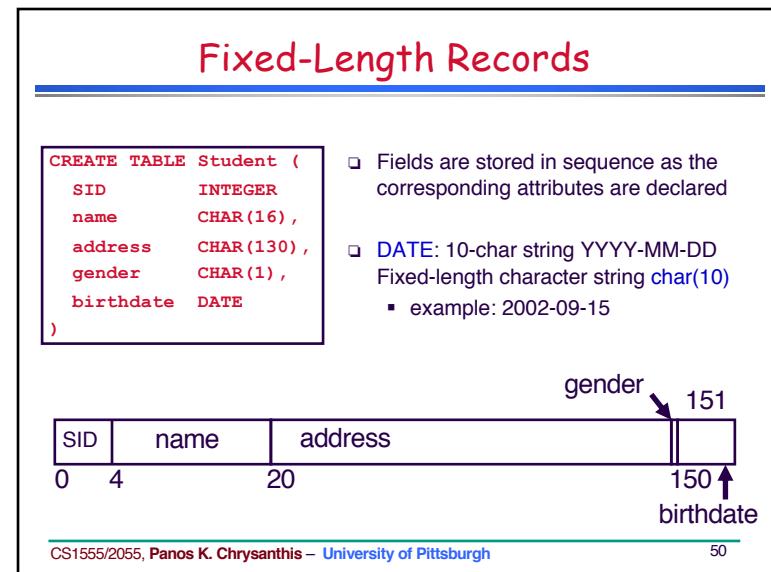
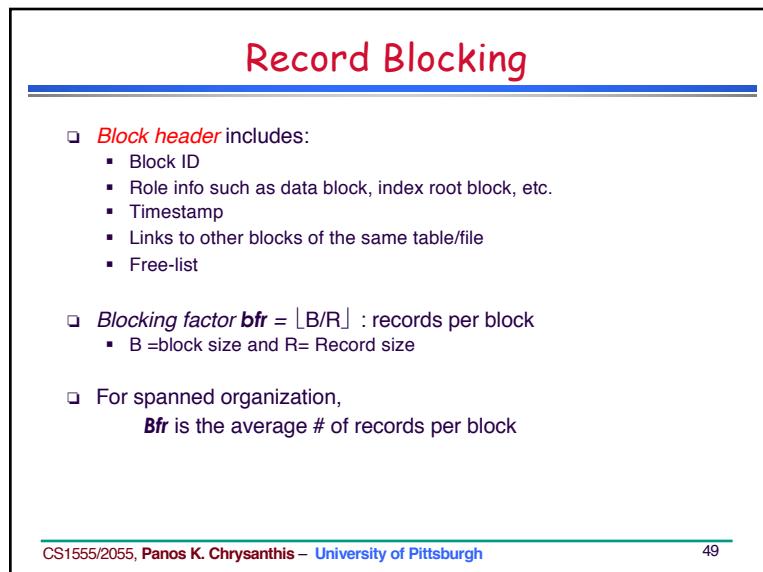
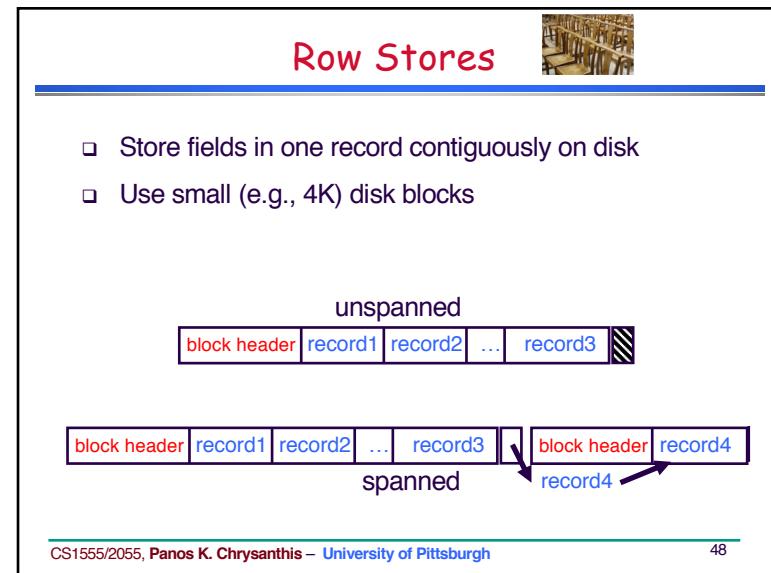
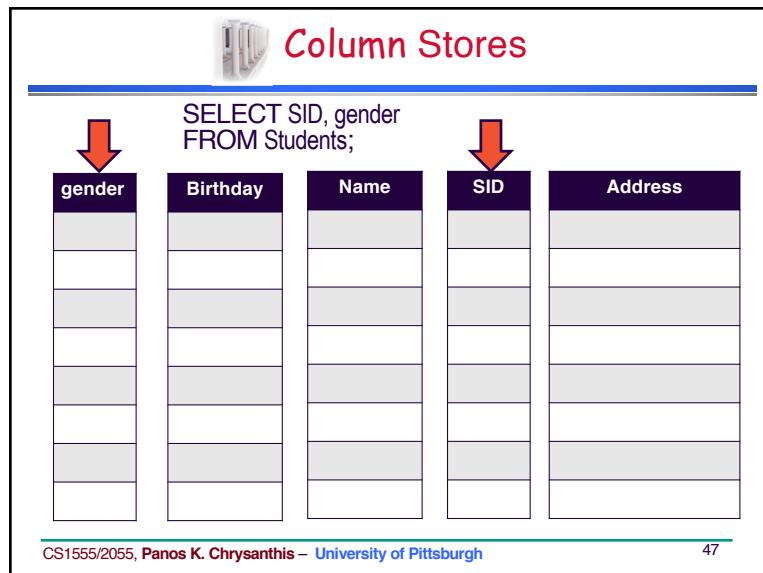


Row Stores



Column Stores

46



Fixed-Length Records - Alignment

```
CREATE TABLE Student (
    SID      INTEGER,
    name     CHAR(16),
    address  CHAR(120),
    gender   CHAR(1),
    birthdate DATE
)
```

- ❑ Each record within a block starts at a byte that is multiple of 4
- ❑ Each field within a record starts at a byte off-set from the beginning of the record that is multiple of 4



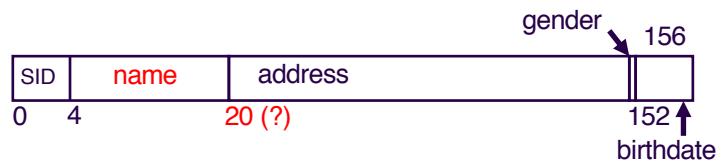
CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

51

Variable-Length Records

```
CREATE TABLE Student (
    SID      INTEGER,
    name     VARCHAR(16),
    address  CHAR(120),
    gender   CHAR(1),
    birthdate DATE
)
```

- ❑ Fields are stored in sequence as the corresponding attributes are declared
- ❑ DATE: 10-char string YYYY-MM-DD
- ❑ Fixed-length character string char(10)
 - example: 2002-09-15



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

52

Variable-Length Attributes

- ❑ example: type **VARCHAR(16)**
 - length + data:



- special end-of-file symbol terminated:
 - Any special character that does not appear in any field value, E.g., ¶, \$, ■



- maximum length:



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

53

File Types

- ❑ Unordered files
- ❑ Ordered files
- ❑ Clustered files
- ❑ Hash files

File header or descriptor includes:

- field names and their data types
- address of the file block on disk



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

55

Unordered Files

- ❑ The simplest file structure: records are stored in no particular order
- ❑ Also called: **Heap**, Pile, or Random File
- ❑ New records are inserted at the **end** of file
 1. The last disk block is copied into buffer (i.e., memory)
 2. New record is added
 3. Block is rewritten back to disk
- ❑ Record **insertion** is quite efficient



Example of Heap File

```
CREATE TABLE deposit (
    account_number CHAR(10),
    branch_name   CHAR(22),
    balance        REAL
)
```

Block 1	Record 0	A-102	Oakland	400
	Record 1	A-305	Shadyside	350
	Record 2	A-101	Downtown	700
	Record 3	A-222	Squirrel Hill	500
Block 2	Record 4	A-217	Shadyside	900
	Record 5	A-110	Waterfront	340
	Record 6	A-257	Oakland	600
	Record 7	A-403	Downtown	250

Properties of Unordered Files

- ❑ File records are inserted at the end of the file or in any file block with free space.
- ❑ Thus, insertion is efficient.
- ❑ To search for a record, a **linear search** through the file records is necessary which is quite expensive.
- ❑ Reading the records in order of any field requires sorting the file records.



Ordered Files

- ❑ Also called **sequential** files.
- ❑ File records are kept sorted by the value of an **ordering** key which has unique value (e.g., primary key)

```
CREATE TABLE deposit (
    account_number CHAR(10),
    branch_name   CHAR(22),
    balance        REAL
)
```

Record 0	A-101	Downtown	700
Record 1	A-102	Oakland	400
Record 2	A-205	Shadyside	350
Record 3	A-217	Shadyside	900
Record 4	A-222	Squirrel Hill	500
Record 5	A-310	Waterfront	340
Record 6	A-357	Oakland	600
Record 7	A-403	Downtown	250

Properties of Ordered Files

- ❑ Insertion is expensive: records must be inserted in the correct order.
- ❑ Deletion?
- ❑ Search for a record on its ordering field value is quite efficient (**binary search algorithm**).
- ❑ Search for a record on a non-ordering field?
- ❑ Reading the records in order of the ordering field is also quite efficient.
- ❑ Reading the records in any order?

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

67

Clustered Files

- ❑ Order files with ordering field which is **not a key**
- ❑ Ordering field has not a unique value

DEPTNUM	1 1 1 2	2 3 3 3	3 3 4 4
NAME			
SSN			
JOB			
BIRTHDATE			
SALARY			
5	5 5 5 5	6 6 6 6	6 8 8 8
6			
7			
8			
9			
10			
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			
21			
22			
23			
24			
25			
26			
27			
28			
29			
30			
31			
32			
33			
34			
35			
36			
37			
38			
39			
40			
41			
42			
43			
44			
45			
46			
47			
48			
49			
50			
51			
52			
53			
54			
55			
56			
57			
58			
59			
60			
61			
62			
63			
64			
65			
66			
67			
68			
69			
70			
71			
72			
73			
74			
75			
76			
77			
78			
79			
80			
81			
82			
83			
84			
85			
86			
87			
88			
89			
90			
91			
92			
93			
94			
95			
96			
97			
98			
99			
100			
101			
102			
103			
104			
105			
106			
107			
108			
109			
110			
111			
112			
113			
114			
115			
116			
117			
118			
119			
120			
121			
122			
123			
124			
125			
126			
127			
128			
129			
130			
131			
132			
133			
134			
135			
136			
137			
138			
139			
140			
141			
142			
143			
144			
145			
146			
147			
148			
149			
150			
151			
152			
153			
154			
155			
156			
157			
158			
159			
160			
161			
162			
163			
164			
165			
166			
167			
168			
169			
170			
171			
172			
173			
174			
175			
176			
177			
178			
179			
180			
181			
182			
183			
184			
185			
186			
187			
188			
189			
190			
191			
192			
193			
194			
195			
196			
197			
198			
199			
200			
201			
202			
203			
204			
205			
206			
207			
208			
209			
210			
211			
212			
213			
214			
215			
216			
217			
218			
219			
220			
221			
222			
223			
224			
225			
226			
227			
228			
229			
230			
231			
232			
233			
234			
235			
236			
237			
238			
239			
240			
241			
242			
243			
244			
245			
246			
247			
248			
249			
250			
251			
252			
253			
254			
255			
256			
257			
258			
259			
260			
261			
262			
263			
264			
265			
266			
267			
268			
269			
270			
271			
272			
273			
274			
275			
276			
277			
278			
279			
280			
281			
282			
283			
284			
285			
286			
287			
288			
289			
290			
291			
292			
293			
294			
295			
296			
297			
298			
299			
300			
301			
302			
303			
304			
305			
306			
307			
308			
309			
310			
311			
312			
313			
314			
315			
316			
317			
318			
319			
320			
321			
322			
323			
324			
325			
326			
327			
328			
329			
330			
331			
332			
333			
334			
335			
336			
337			
338			
339			
340			
341			
342			
343			
344			
345			
346			
347			
348			
349			
350			
351			
352			
353			
354			
355			
356			
357			
358			
359			
360			
361			
362			
363			
364			
365			
366			
367			
368			
369			
370			
371			
372			
373			
374			
375			
376			
377			
378			
379			
380			
381			
382			
383			
384			
385			
386			
387			
388			
389			
390			
391			
392			
393			
394			
395			
396			
397			
398			
399			
400			
401			
402			
403			
404			
405			
406			
407			
408			
409			
410			
411			
412			
413			
414			
415			
416			
417			
418			
419			
420			
421			
422			
423			
424			
425			
426			
427			
428			
429			
430			
431			
432			
433			
434			
435			
436			
437			
438			
439			
440			
441			
442			
443			
444			
445			
446			
447			
448			</td

Static Hashing

- ❑ Hashing converts the key of a record into an address in which the record is stored.
- ❑ An external **hash** function maps the key to the relative address of a **bucket** in which the record is stored.
 - if a file is allocated s buckets, the hash function must convert a key k into the relative address of the block:
$$h(k) \in \{0, \dots, s-1\}$$
- ❑ A bucket is either one disk block or a cluster of contiguous blocks
- ❑ A table stored in the header of the file maps relative bucket numbers to disk block addresses.

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

74

Collision

- ❑ Insertion of a new record may lead to **collision**.
 - No space in $B = h(k)$
- ❑ *Probing or conflict resolution:*
 - *Open Addressing (Rehashing):* If bucket $h(k)$ is full, use another hash function until a bucket with a free space is found.
E.g., *linear probing*
$$\alpha = h(\text{key}) = \text{key mod } s$$

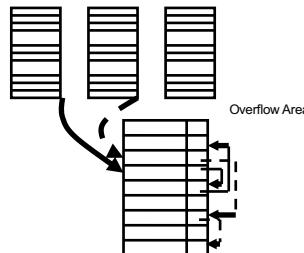
$$rh(\alpha) = h(\alpha + 1)$$
 - Not a very good technique for databases (Why?)
- *Chaining:* Use overflow buckets.

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

75

Handling of the Overflow Area

- ❑ Three approaches:
- 1. **Common** overflow area for all blocks in the file.
 - Each block has a pointer to its first record in the overflow area.
 - Records belonging to the same block are linked by pointers.

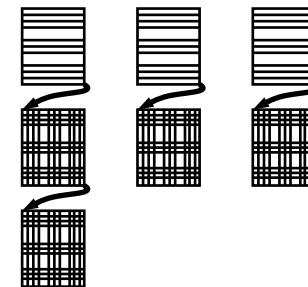


CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

76

Handling of the Overflow Area

- 2. **Share** overflow area by some blocks
 - e.g., residing on the same cylinder.
- 2. **Individual** overflow areas
 - Each block has its own overflow
 - Records are not linked



CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

77

Hashing Functions

- ❑ A good hash functions must
 - 1) be computed efficiently
 - 2) minimize the number of collisions by spreading keys around the file as evenly and uniform as possible.
- ❑ Example of good functions
 - truncation
 - division: $h(\text{key}) = \text{key} \bmod s$
 - Mid-square
 - Folding or partitioning
 - Other ad hoc methods
- ❑ Order preserving hash functions:
 - Maintain records in order of hash field values

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

78

Pros and Cons of Hashing

- ❑ Excellent performance for searching on **equality on the key** used for hashing (assuming low density).
- ❑ Records are not ordered (heap files).
 - ⇒ Any search other than on equality is very expensive (linear search or involves sorting).
- ❑ Prediction of total number of buckets is difficult.
 - allocate a large space.
 - estimate a ``reasonable'' size and periodically reorganize.

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

79

Dynamic Hashing Methods

- ❑ Allow the file size to change as records are added or deleted.
 - Linear Hashing
 - No additional structure
 - Extendible Hashing
 - Binary Hashing

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

80

Linear Hashing

- The file size grows linearly, bucket by bucket.
 - The file starts with s main buckets, where s is a power of 2, and k overflow buckets.
 - Buckets are numbered from 0 to $s-1$
 - ⇒ initial hashing function $h_0(\text{key}) = \text{key} \bmod s$.
 - Collisions are handled thru **chaining**.
- We keep the following info:
 - B_{last} : A pointer to the **current** last bucket.
Initially, $B_{last} = s-1$.
 - B_{split} : A pointer to the bucket that should be split next.
Initially, $B_{split} = 0$

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

81

Insertion of a Record (Method 1)

- ❑ if there is a collision, push tuple to an overflow bucket.
- ❑ if there are no overflow buckets available, proceed as follows until one becomes available:
 - A new bucket is appended at the end of the hash table and $B_{last} = B_{last} + 1$
 - Records in B_{split} bucket are hashed again using $h_1(key) = key \bmod (2s)$,
 - these will either remain in B_{split} or stored in B_{last}
 - ⇒ this may free an overflow bucket.
 - B_{split} becomes $B_{split} + 1$.
- ❑ if $B_{last} = 2s - 1$
 - set $s = 2s$, $B_{last} = s - 1$, $B_{split} = 0$, $h_0(key) = h_1(key)$
- ❑ proceed as above.

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

82

Insertion of a Record (Method 2)

- ❑ if there is a collision, push tuple to an overflow bucket.
- ❑ if there are no overflow buckets available, proceed as follows until one becomes available:
 - A new bucket is appended at the end of the hash table (but B_{last} does not change as in Method 1)
 - Records in B_{split} bucket are hashed again using $h_1(key) = key \bmod (2s)$,
 - these will either remain in B_{split} or stored in B_{last}
 - ⇒ this may free an overflow bucket.
 - B_{split} becomes $B_{split} + 1$.
- ❑ if $B_{last} < B_{split}$
 - set $s = 2s$, $B_{last} = s - 1$, $B_{split} = 0$, $h_0(key) = h_1(key)$
- ❑ proceed as above.

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

83

Search for a record

- Search for a Record
 - $I = h_0(key)$
 - if $I < B_{split}$ then $I = h_1(key)$
 - search the bucket whose hash value is I (and its overflow, if any).

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

84

Variations

- ❑ Split a bucket either when there is a need for an overflow bucket or based on some threshold u , e.g., $u > 0.9$
- ❑ u is an upper bound of the file load factor or number of records that current buckets may store,
$$I = r / (bfr * N)$$
 $r = \text{current number of records}, N = \text{current number of buckets}/$
- ❑ Combination can be triggered when, e.g., $u < 0.7$
- ❑ Use a single hashing function:
$$A = k \bmod 2s$$
 $\text{if } A > B_{last}$ $\quad \text{then } A = A - s$ $\quad \text{else } A$

CS1555/2055, Panos K. Chrysanthis – University of Pittsburgh

85

File Types

- ❑ Unordered files
- ❑ Ordered files
- ❑ Clustered files
- ❑ Hash files

File header or descriptor includes:

- field names and their data types
- address of the file block on disk

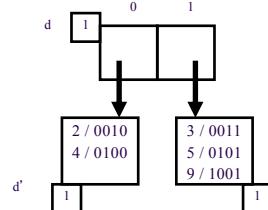


Dynamic Hashing Methods

- ❑ Allow the file size to change as records are added or deleted.
 - Linear Hashing
 - No additional structure
 - Extendible Hashing
 - Binary Hashing

Extendible Hashing

- ❑ The file is structured into two levels: **directory** and buckets.



Extendible Hashing

- ❑ Directory contains only pointers to buckets (no keys).
- ❑ Use some hash function to generate a **pseudocodey** of b-bits (typically b=32)
- ❑ Use the first (or last) **d** bits to find the offset of the bucket pointer in the directory.
- ❑ **d** is called the (global) **directory depth**. It may be stored in the directory.
- ❑ In each bucket a local depth **d'** is stored indicating the most (or least) significant bits common to all keys in that bucket (**d' ≤ d**).

Extendible Hashing - File Growth

- ❑ No overflow buckets. New buckets are added one by one.
- ❑ When a bucket B with local depth d' overflows,
B is split into two buckets and all keys are rehash using $d'+1$ bits.
- ❑ If after a split, $d' > d$, double the size of the directory ($d=d+1$) to accommodate the growth.