

Outline

- ✤ It's all about disks!
 - That's why we always draw databases as
 - And why the single most important metric in database processing is the number of disk I/O's performed
- ✤ Record layout
- ✤ Block layout



How far away is data?			
Location	Cycles	Location	Time
Registers	1		1 min.
On-chip cache	2		2 min.
On-board cache	10		10 min.
Memory	100		1.5 hr.
Disk	106		2 yr.
Tape	10 ⁹		2000 yr.
(Source: AlphaSort paper, 1995)			
I/O dominates—design your algorithms to reduce I/O!			











Disk access time

Sum of:

- Seek time: time for disk heads to move to the correct cylinder
- Rotational delay: time for the desired block to rotate under the disk head
- Transfer time: time to read/write data in the block
 (= time for disk to rotate over the block)

Random disk access

Seek time + rotational delay + transfer time

- ✤ Average seek time
 - Time to skip one half of the cylinders?
 - Not quite; should be time to skip a third of them (why?)
 - "Typical" value: 5 ms
- * Average rotational delay
 - Time for a half rotation (a function of RPM)
 - "Typical" value: 4.2 ms (7200 RPM)

Sequential disk access

Seek time + rotational delay + transfer time

✤ Seek time

- 0 (assuming data is on the same track)
- Rotational delay
 - 0 (assuming data is in the next block on the track)
- Easily an order of magnitude faster than random disk access!

Data layout strategy

Keep related things close together!

10

11

12

- ✤ Same sector/block
- ✤ Same track
- ✤ Same cylinder
- Adjacent cylinder

More performance tricks

- * Disk scheduling algorithm
- Example: "elevator" algorithm
- ✤ Track buffer
 - Read/write one entire track at a time
- * Double buffering
 - While processing the current block in memory, prefetch the next block from disk
- ♦ Parallel I/O
 - More disk heads working at the same time

Record layout

Record = row in a table

- ✤ Variable-format records
 - Rare in DBMS—table schema dictates the format
 - Relevant for semi-structured data such as XML
- * Focus on fixed-format records
 - With fixed-length fields only, or
 - With possible variable-length fields



• Add a bitmap at the beginning of the record





Block layout

How do you organize records in a block?

- * NSM (N-ary Storage Model)
 - Most commercial DBMS
- PAX (Partition Attributes Across)
 - Recent work (Ailamaki et al., VLDB 2001)

NSM

 $\boldsymbol{\bigstar}$ Store records from the beginning of each block

- * Use a directory at the end of each block
 - To locate records and manage free space
 - Necessary for variable-length records



Options

 Reorganize after every update/delete to avoid fragmentation (gaps between records) 18

- Need to rewrite half of the block on average
- What if records are fixed-length?
 - Reorganize after delete
 Only need to move one record
 - Need a pointer to the beginning of free space
 - Do not reorganize after update
 Need a bitmap indicating which slots are in use



- ♦ Query: SELECT SID FROM Student WHERE GPA > 2.0;
- Assumption: cache block size < record size</p>
- * Lots of cache misses
 - ID and GPA are not close enough by memory standards





20

21

19



Bart⊀

- * Most queries only access a few columns
- Cluster values of the same columns in each block
 - When a particular column of a row is brought into the cache, the same column of the next row is brought in together Reorganize after every update

4 (number of records) (for variable-length records only) 142 123 857 456 and delete to keep fields together lhouse tisa Ralph 10 10 8 8

1111

1111 (IS NOT NULL bitmap)

Summary

* Storage hierarchy

2.3 3.1 4.3 2.3

- Why I/O's dominate the cost of database operations
- Disk
 - Steps in completing a disk access
 - Sequential versus random accesses
- * Record layout
 - Handling variable-length fields
 - Handling NULL
 - Handling modifications
- * Block layout
 - NSM: the traditional layout
 - PAX: a layout that tries to improve cache performance