

Physical Data Organization

CS348 Spring 2023

Instructor: Sujaya Maiyya

Sections: **002 & 004 only**

Announcements

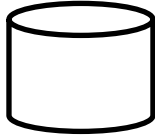
- **Assignment 2**
 - Due today!
- **Milestone 1**
 - Due Thursday, June 22nd
- **Midterm Exam**
 - Monday, June 26th, 7-8:50 PM
 - **Covers lectures 1-10 (except lecture 6)**

Where are we?

- Relational model (lecture 2)
 - SQL (lectures 3-6)
 - Database design (lectures 7-10)
- Conceptual/Logical level
- **Storage management** & indexing (lectures 11-12)
 - Query processing & optimizations (lectures 13-14)
 - Transaction management (lectures 15-17)

Physical Data Organization

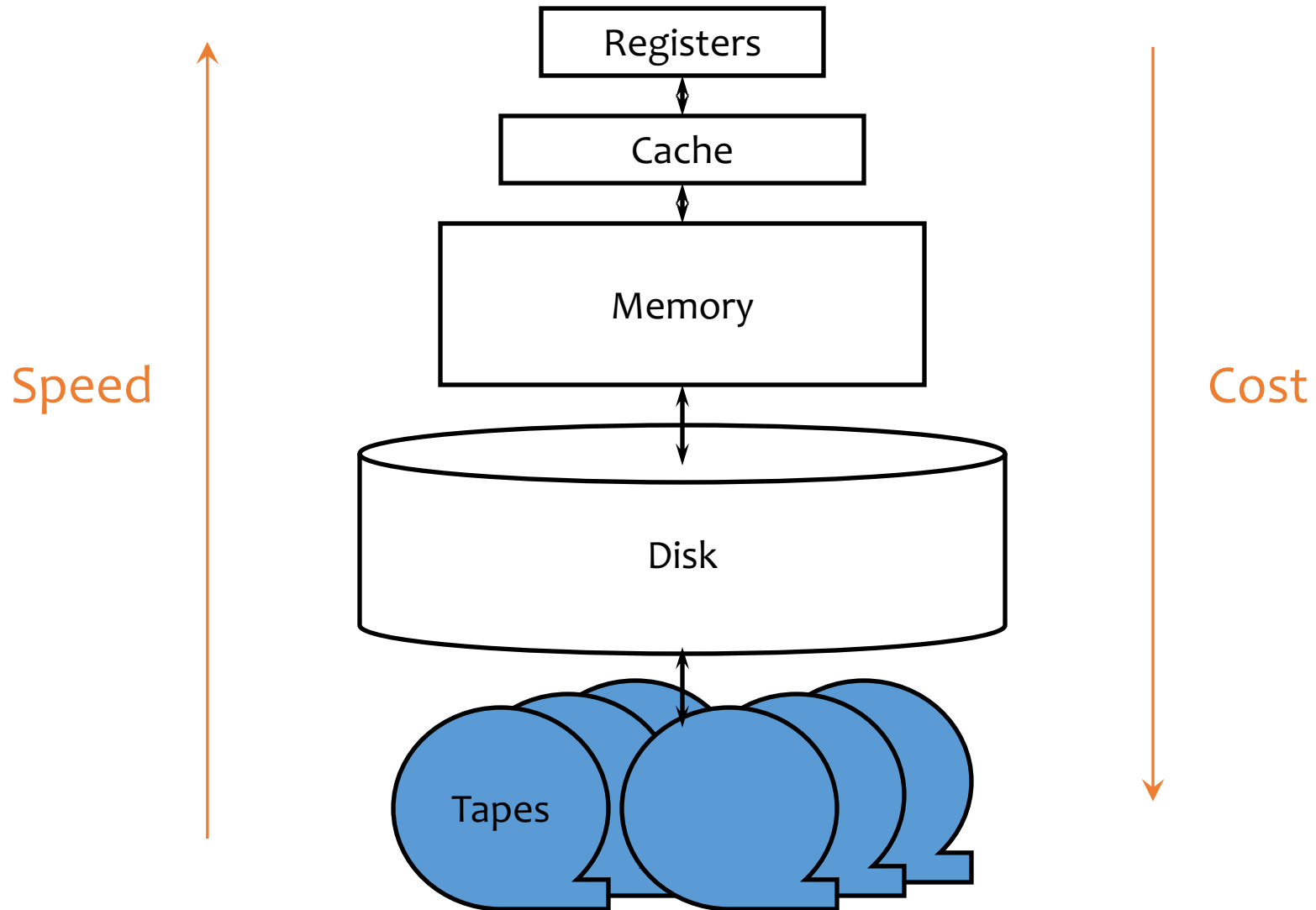
- It's all about disks!

- That's why we always draw databases as 
- And why one of the most important metric in database processing is the number of disk I/O's performed

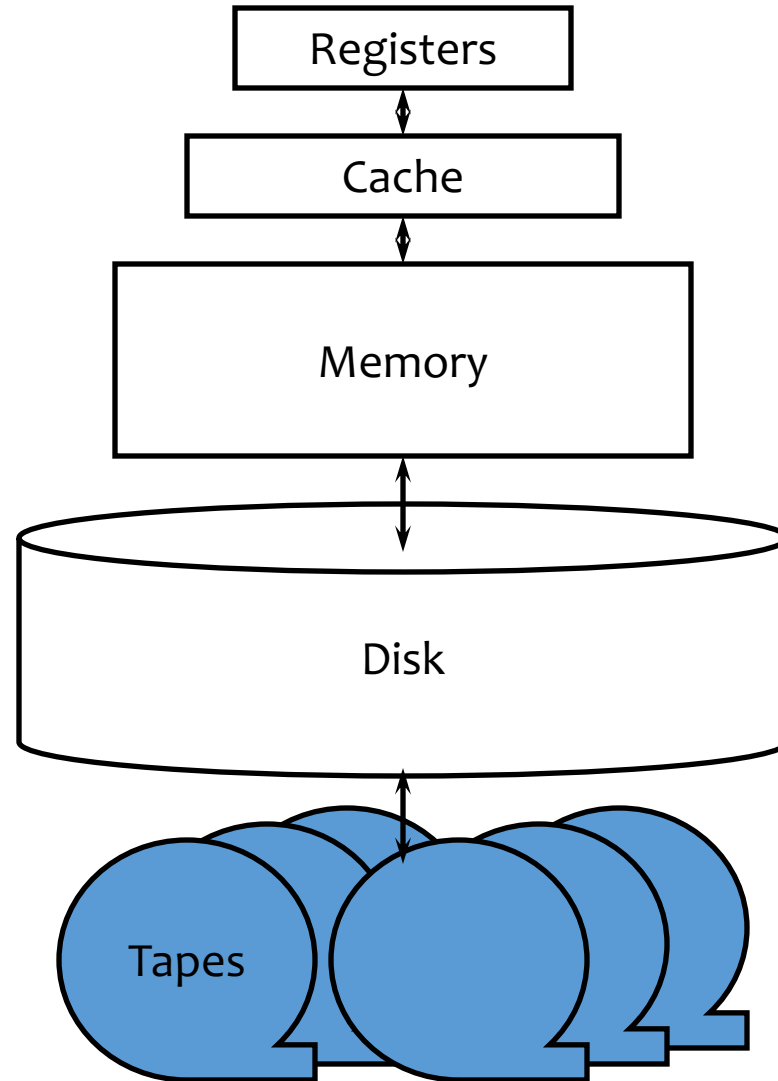
- Storing data on a disk

- Record layout
- Block layout
- Column stores

Storage hierarchy



Storage hierarchy



How far away is data?

<u>Location</u>	<u>Cycles</u>	<u>Location</u>	<u>Time</u>
Registers	1	This room	1-2 min.
On-chip cache	2	Waterloo campus	10 min.
On-board cache	10	Toronto	1.5 hr.
Memory	100	Pluto	2 yr.
Disk	10^6	Andromeda	2000 yr.
Tape	10^9		

(Source: AlphaSort paper, 1995)
The gap has been widening!

👉 I/O dominates—design your algorithms to reduce I/O!

Latency Numbers Every Programmer Should Know

Latency Comparison Numbers

L1 cache reference	0.5 ns				
Branch mispredict	5 ns				
L2 cache reference	7 ns				14x L1 cache
Mutex lock/unlock	25 ns				
Main memory reference	100 ns				20x L2 cache, 200x L1 cache
Compress 1K bytes with Zippy	3,000 ns		3 us		
Send 1K bytes over 1 Gbps network	10,000 ns		10 us		
Read 4K randomly from SSD*	150,000 ns		150 us		~1GB/sec SSD
Read 1 MB sequentially from memory	250,000 ns		250 us		
Round trip within same datacenter	500,000 ns		500 us		
Read 1 MB sequentially from SSD*	1,000,000 ns		1,000 us	1 ms	~1GB/sec SSD, 4X memory
Disk seek	10,000,000 ns		10,000 us	10 ms	20x datacenter roundtrip
Read 1 MB sequentially from disk	20,000,000 ns		20,000 us	20 ms	80x memory, 20X SSD
Send packet CA->Netherlands->CA	150,000,000 ns		150,000 us	150 ms	

Notes

 1 ns = 10⁻⁹ seconds
 1 us = 10⁻⁶ seconds = 1,000 ns
 1 ms = 10⁻³ seconds = 1,000 us = 1,000,000 ns

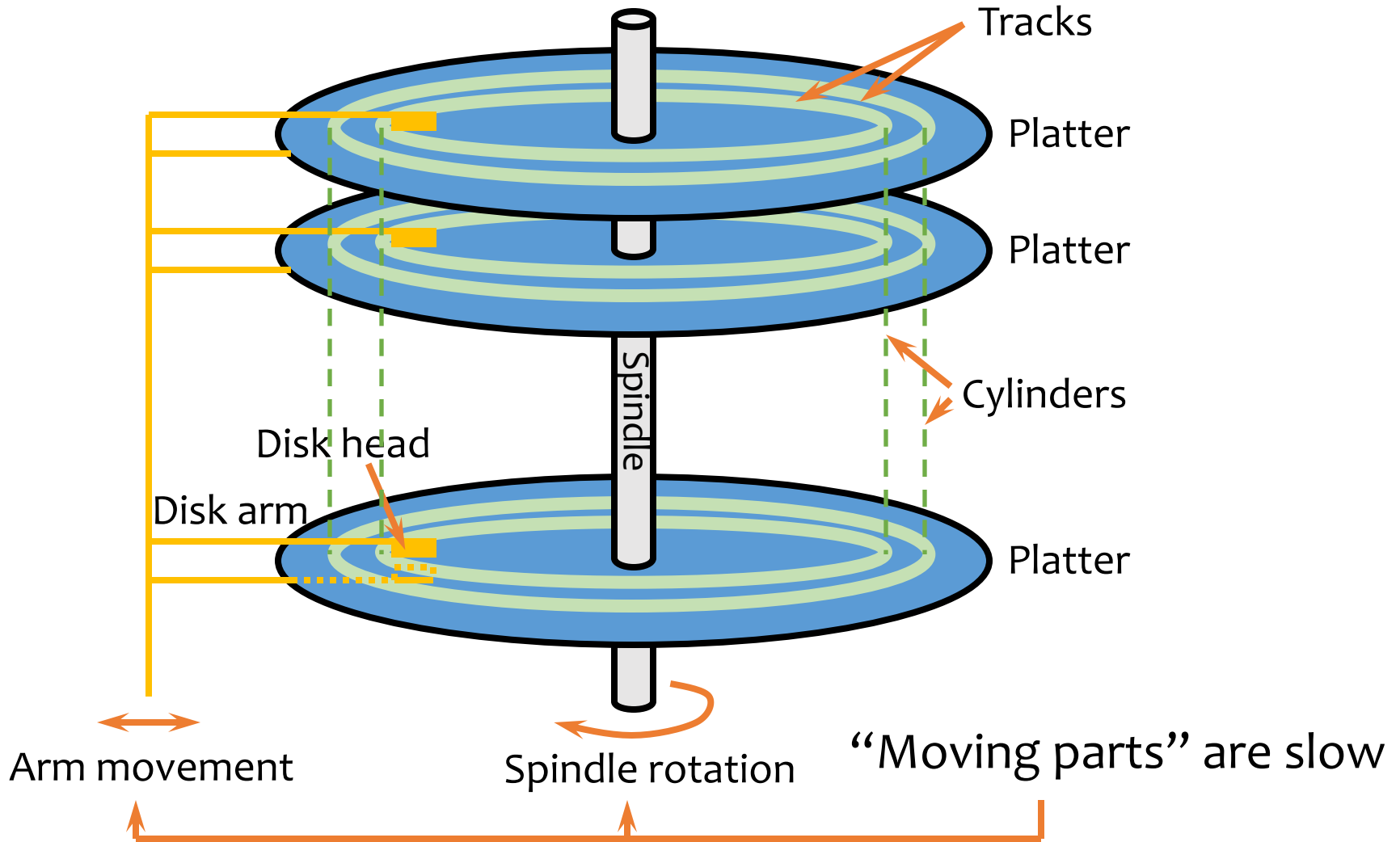
Credit

 By Jeff Dean: <http://research.google.com/people/jeff/>
 Originally by Peter Norvig: <http://norvig.com/21-days.html#answers>

A typical hard drive

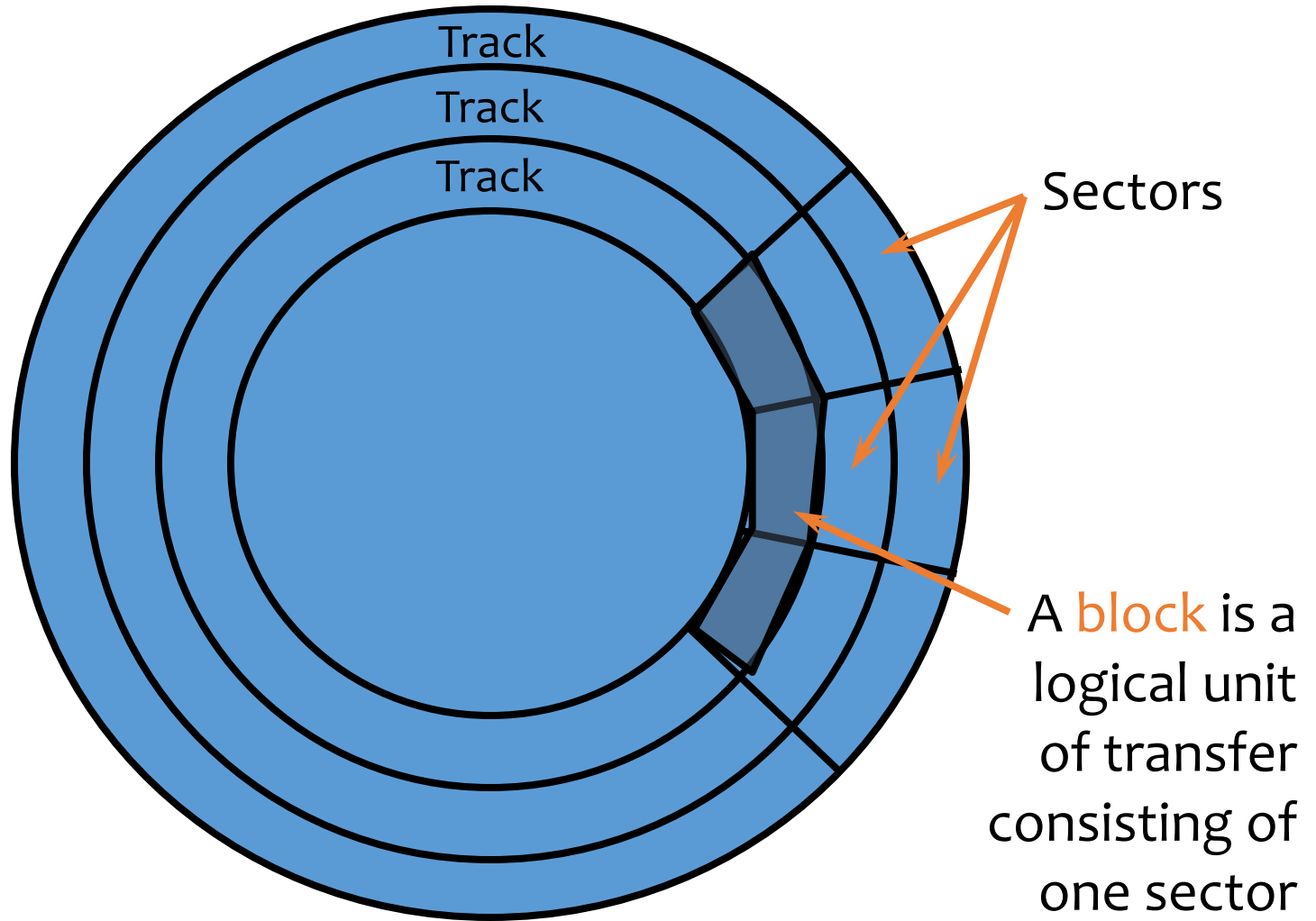


A typical hard drive



Top view

“Zoning”: more sectors/data on outer tracks



Disk access time

Disk access time: time from when a read or write request is issued to when data transfer begins

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)
- Total data access time = seek time + rotational delay + transfer time

Random disk access

→ Successive requests are for blocks that are randomly located on disk

Delay = Seek time + rotational delay + transfer time

- Average seek time
 - Seek the right cylinder for each access
 - “Typical” value: 5 ms
- Average rotational delay
 - Rotate for the right block for each access
 - “Typical” value: 4.2 ms (7200 RPM)

Sequential disk access

→ Successive requests are for successive block numbers, which are on the same track, or on adjacent tracks

Delay = Seek time + rotational delay + transfer time

- Seek time
 - 1 time delay: seek the right cylinder once
- Rotational delay
 - 1 time delay: rotate to the right block once
- Easily an order of magnitude faster than random disk access!

What about SSD (solid-state drives)?



No mechanical parts (flash-based)



What about SSD (solid-state drives)?

- 1-2 orders of magnitude **faster random access** than hard drives (under 0.1ms vs. several ms)
- Little difference between random vs. sequential read performance
- Random writes still hurt
 - In-place update would require erasing the whole “erasure block” and rewriting it!

Important consequences

- It's all about reducing I/O's!
- Cache blocks from stable storage in memory
 - DBMS maintains a memory **buffer pool** of blocks
 - Reads/writes operate on these memory blocks
 - Dirty (updated) memory blocks are “flushed” back to stable storage
- Sequential I/O is much faster than random I/O

Performance tricks

- Disk layout strategy: keep related things close
- Prefetching
- Parallel I/O: multiple disk heads
- Track buffer: read/write one entire track at a time

Where are we?

- Storage hierarchy: I/O cost
- Disk: Sequential versus random accesses
- Record layout

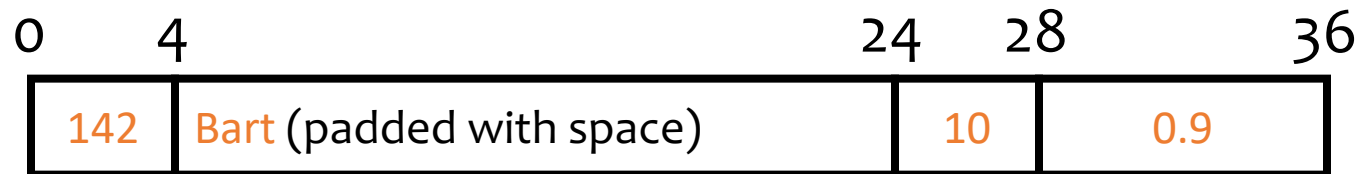
Record layout

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

Fixed-length fields

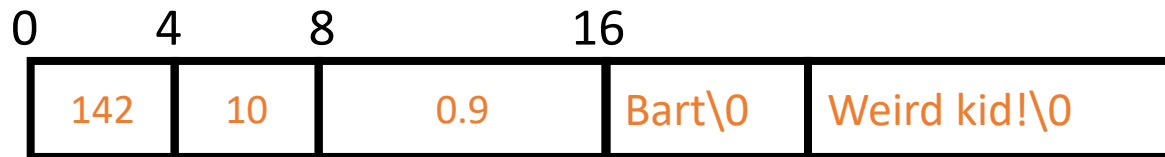
- All field lengths and offsets are constant
 - Computed from schema, stored in the system catalog
- Example: CREATE TABLE User(uid INT, name CHAR(20), age INT, pop FLOAT);



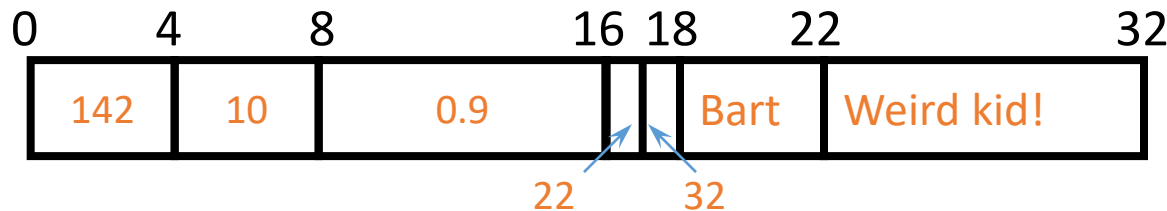
- If block size \neq 36, one row maybe split across multiple blocks or move to next block & leave the remaining space empty
- What about NULL?
 - Add a bitmap at the beginning of the record

Variable-length records

- Example: `CREATE TABLE User(uid INT, name VARCHAR(20), age INT, pop FLOAT, comment VARCHAR(100));`
- Put all variable-length fields at the end
- Approach 1: use field delimiters ('`\0`' okay?)



- Approach 2: use an offset array



- Scheme update is messy if it changes the length of a field

BLOB fields

- Example: CREATE TABLE User(uid INT, name CHAR(20), age INT, pop FLOAT, picture BLOB(32000));
- User records get “de-clustered”
 - Bad because most queries do not involve picture
- Decomposition (automatically and internally done by DBMS without affecting the user)
 - (uid, name, age, pop)
 - (uid, picture)

Where are we?

- Storage hierarchy: I/O cost
- Disk: Sequential versus random accesses
- Record layout: fixed length v.s. variable length
- Block layout

Block layout

How do you organize records in a block?

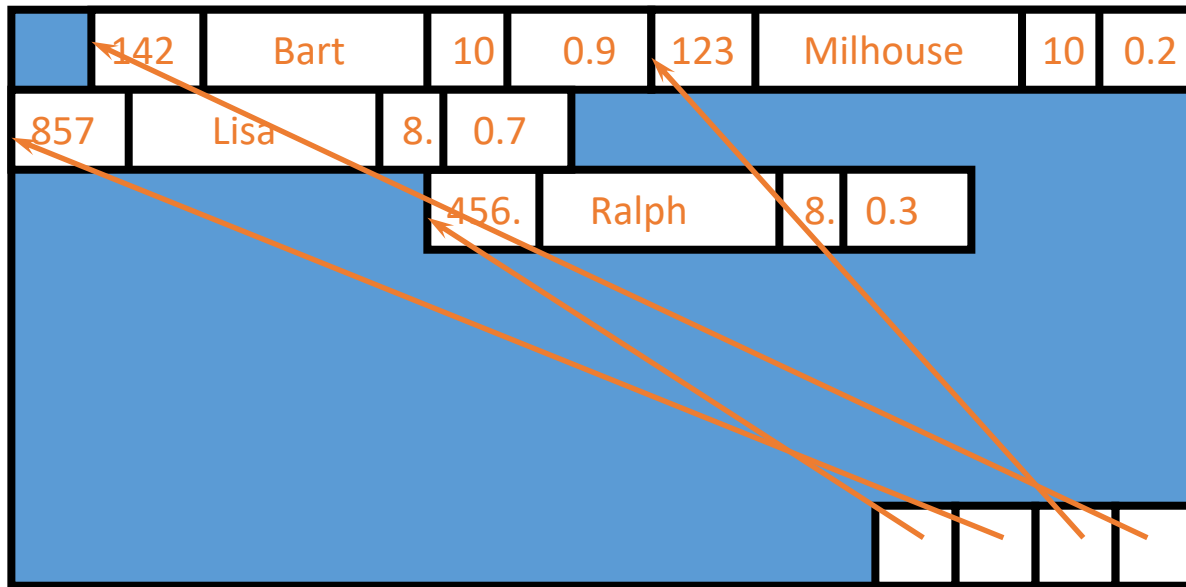
- **NSM** (N-ary Storage Model)
 - Most commercial DBMS
- **PAX** (Partition Attributes Across)
 - Ailamaki et al., *VLDB* 2001

Options

- Reorganize after every update/delete to avoid fragmentation (gaps between records)
 - Need to rewrite half of the block on average
- A special case: What if records are fixed-length?
 - Option 1: reorganize after delete
 - Only need to move one record
 - Need a pointer to the beginning of free space
 - Option 2: do not reorganize after update
 - Need a bitmap indicating which slots are in use

Cache behavior of NSM

- Query: `SELECT uid FROM User WHERE pop > 0.8;`
- Assumptions: no index, and cache line size < record size
- Lots of cache misses



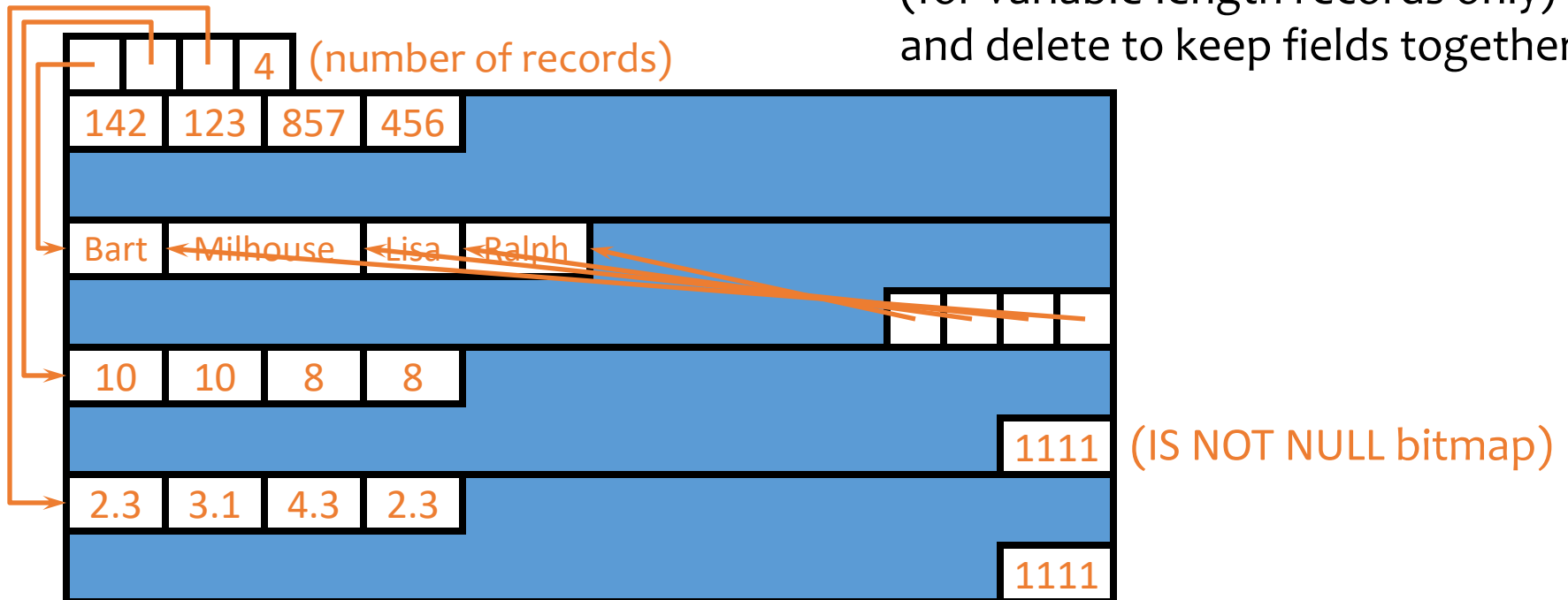
142 Bart 10
0.9 123 Milhouse
10 0.2 857 Lisa
8 0.7
456 Ralph 8
0.3

Cache

PAX

- Most queries only access a few columns
- Cluster values of the same columns in each block
- Better sequential reads for queries that read a single column

Reorganize after every update
(for variable-length records only)
and delete to keep fields together



Beyond block layout: column stores

- Store tables by columns instead of rows
 - Better cache performance
 - Fewer I/O's for queries involving many rows but few columns
 - Aggressive compression to further reduce I/O's
- More disruptive changes to the DBMS architecture are required than PAX
 - Not only storage, but also query execution and optimization

Column vs. row oriented db

User:

uid	name	pop	age
1	Bart	.6	12
2	Lisa	.9	10
3	Abe	.3	65

Row oriented

1	Bart	.6	12
2	Lisa	.9	10
3	Abe	.3	65

Column oriented

1	2	3
Bart	Lisa	Abe
.6	.9	.3
12	10	65

Summary

- Storage hierarchy: I/O cost
- Disk: Sequential versus random accesses
- Record layout: fixed length v.s. variable length
- Block layout: NSM v.s. PAX
- Column stores: NSM transposed, beyond blocks