

Department of Computer Science

CS 44800: Introduction To Relational Database Systems

Storing the Data Prof. Chris Clifton 21 September 2021



Hardware: Key Takeaways

Department of Computer Science

- Database must reside on non-volatile storage
 Can cache in faster storage
- Non-volatile storage slow
 - But accessing a lot not much different than accessing a little
 - Therefore we read/write as large blocks (typically 4kb)
- Abstract performance as: α+βb
 - α is seek time (abstraction of read/write setup overhead)
 - $-\beta$ is transfer rate
 - b is block size
- Rotating media: seek can dominate (but caching, sequential reads reduce this)
- Solid state: transfer dominates
 - but erasure, protocol overheads make "seek" more than you'd expect
- · Writes typically worse than reads
 - Not "done" until safe in non-volatile storage, so reduces caching benefits

2

ndiana

Center for

Database

Systems



File Organization

- The database is stored as a collection of *files*. Each file is a sequence of *records*. A record is a sequence of fields.
- One approach
 - Assume record size is fixed
 - · Each file has records of one particular type only
 - · Different files are used for different relations

This case is easiest to implement; will consider variable length records later

1.4

• We assume that records are smaller than a disk block







Fixed-Length Records

Simple approach:

Database System Concepts - 7th Edition

- Store record *i* starting from byte n * (i 1), where *n* is the size of each record.
- · Record access is simple but records may cross blocks
 - Modification: do not allow records to cross block boundaries

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000
		1.6		

Fixed-Length Records Deletion of record *i*: alternatives: move records i + 1, ..., n to i, ..., n - 1 move record *n* to *i* do not move records, but link all free records on a free list • **Record 3 deleted** record 0 10101 Comp. Sci. 65000 Srinivasan record 1 12121 Wu Finance 90000 15151 Music 40000 record 2 Mozart record 4 32343 El Said History 60000 record 5 33456 Physics 87000 Gold record 6 45565 Katz Comp. Sci. 75000 58583 History 62000 record 7 Califieri 76543 80000 Finance record 8 Singh Biology 72000 record 9 76766 Crick record 10 83821 Brandt Comp. Sci. 92000 98345 record 11 Kim Elec. Eng. 80000 Database System Concepts - 7th Edition 1.7

©Silberschatz, Korth and Sudarshan



Fixed-Length Records

- Deletion of record i: alternatives:
 - move records i + 1, ..., n to i, ..., n 1
 - move record n to i
 - do not move records, but link all free records on a free list

Record 3 deleted and replaced by record 11

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 11	98345	Kim	Elec. Eng.	80000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000

1.8

Database System Concepts - 7th Edition

















Organization of Records in Files

- Heap record can be placed anywhere in the file where there is space
- Sequential store records in sequential order, based on the value of the search key of each record
- In a multitable clustering file organization records of several different relations can be stored in the same file
 - · Motivation: store related records on the same block to minimize I/O
- B⁺-tree file organization
 - Ordered storage even with inserts/deletes
 - More on this in Chapter 14
- Hashing a hash function computed on search key; the result specifies in which block of the file the record should be placed
 - More on this in Chapter 14

Database System Concepts - 7th Edition	1.16	©Silberschatz, Korth and Sudarshan

A	Heap File Organization					
	Records can be placed anywhere in the file where there is free space					
•	Records usually do not move once allocated					
	Important to be able to efficiently find free space within file					
	Free-space map					
	 Array with 1 entry per block. Each entry is a few bits to a byte, and records fraction of block that is free 					
	 In example below, 3 bits per block, value divided by 8 indicates fraction of block that is free 					
	Can have second-level free-space map					
	 In example below, each entry stores maximum from 4 entries of first-level free-space map 4 7 2 6 					
 Free space map written to disk periodically, OK to have wrong (old) values for some entries (will be detected and fixed) 						
Database System Co	uncents, 7 th Edition 117 @Silberschatz Korth and Sudarshan					
Database System CO						



Sequential File Organization

- Suitable for applications that require sequential processing of the entire file
- The records in the file are ordered by a search-key

10101	Srinivasan	Comp. Sci.	65000	-
12121	Wu	Finance	90000	-
15151	Mozart	Music	40000	-
22222	Einstein	Physics	95000	-
32343	El Said	History	60000	-
33456	Gold	Physics	87000	-
45565	Katz	Comp. Sci.	75000	-
58583	Califieri	History	62000	_
76543	Singh	Finance	80000	-
76766	Crick	Biology	72000	-
83821	Brandt	Comp. Sci.	92000	_
98345	Kim	Elec. Eng.	80000	_

```
Database System Concepts - 7th Edition
```

1.18









Multitable Clustering File Organization								
Store several relations in one file using a multitable clustering file organization								
department	dept_nai	me	buildin	g	budge	t		
	Comp.	Sci.	Taylor		100000)		
	ID	name	de	pt_nc	ame .	salary	,	
instructor	10101 33456 45565 83821	Sriniv Gold Katz Brand	rasan C Pi C It C	omp. nysic omp. omp.	. Sci. 6 s 8 . Sci. 7 . Sci. 9	55000 87000 75000 92000)))	
		<u>a i l</u>			100000			
multitable clustering	Comp. 10101	Sci.	Taylor	m	100000 Comp	Sci	65000	
of department and	45565		Katz		Comp.	Sci.	75000	
instructor	83821		Brandt		Comp.	Sci.	92000	
	Physics		Watson		70000			
33456 Gold Physics 87000								
Detabase System Concents - 7th Edition			1 23					@Silkerschatz Korth and Sudarshan
Database System Concepts - 7 th Edition			1.23					©Silberschatz, Korth and Sudarshan



Multitable Clustering File Organization (cont.)

- good for queries involving *department* ⋈ *instructor*, and for queries involving one single department and its instructors
- bad for queries involving only department
- results in variable size records
- Can add pointer chains to link records of a particular relation

Database System	Concepts	- 7 th	Edition

1.24





Department of Computer Science

CS 44800: Introduction To Relational Database Systems

Storing the Data Prof. Chris Clifton 23 September 2021



ndiana

Center for

Database

Svstems







Buffer Manager

- Programs call on the buffer manager when they need a block from disk.
 - If the block is already in the buffer, buffer manager returns the address of the block in main memory
 - · If the block is not in the buffer, the buffer manager
 - Allocates space in the buffer for the block
 - Replacing (throwing out) some other block, if required, to make space for the new block.
 - Replaced block written back to disk only if it was modified since the most recent time that it was written to/fetched from the disk.
 - Reads the block from the disk to the buffer, and returns the address of the block in main memory to requester.







Buffer Manager

- Buffer replacement strategy (details coming up!)
 - Pinned block: memory block that is not allowed to be written back to disk
 - Pin done before reading/writing data from a block
 - Unpin done when read /write is complete
 - Multiple concurrent pin/unpin operations possible
 - Keep a pin count, buffer block can be evicted only if pin count = 0
- Shared and exclusive locks on buffer
 - Needed to prevent concurrent operations from reading page contents as they are moved/reorganized, and to ensure only one move/reorganize at a time

1.32

- Readers get shared lock, updates to a block require exclusive lock
- Locking rules:
 - Only one process can get exclusive lock at a time
 - Shared lock cannot be concurrently with exclusive lock
 - · Multiple processes may be given shared lock concurrently

Database System Concepts - 7th Edition









Buffer-Replacement Policies

- Most operating systems replace the block least recently used (LRU strategy)
 - Idea behind LRU use past pattern of block references as a predictor of future references
 - LRU can be bad for some queries
- Queries have well-defined access patterns (such as sequential scans), and a database system can use the information in a user's query to predict future references
- Mixed strategy with hints on replacement strategy provided by the query optimizer is preferable
- Example of bad access pattern for LRU: when computing the join of 2 relations r and s by a nested loops

for each tuple *tr* of *r* do for each tuple *ts* of *s* do if the tuples *tr* and *ts* match ...

```
Database System Concepts - 7th Edition
```

1.37









Column-Oriented Storage

- Also known as columnar representation
- Store each attribute of a relation separately
- Example

10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

Database System Concepts - 7th Edition

1.41

©Silberschatz, Korth and Sudarshan



Columnar Representation

- Benefits:
 - · Reduced IO if only some attributes are accessed
 - Improved CPU cache performance
 - Improved compression
 - · Vector processing on modern CPU architectures
- Drawbacks
 - · Cost of tuple reconstruction from columnar representation
 - · Cost of tuple deletion and update
 - Cost of decompression
- Columnar representation found to be more efficient for decision support than row-oriented representation
- Traditional row-oriented representation preferable for transaction processing
- Some databases support both representations
 - Called hybrid row/column stores

Database System Concepts - 7th Edition

1.42



