



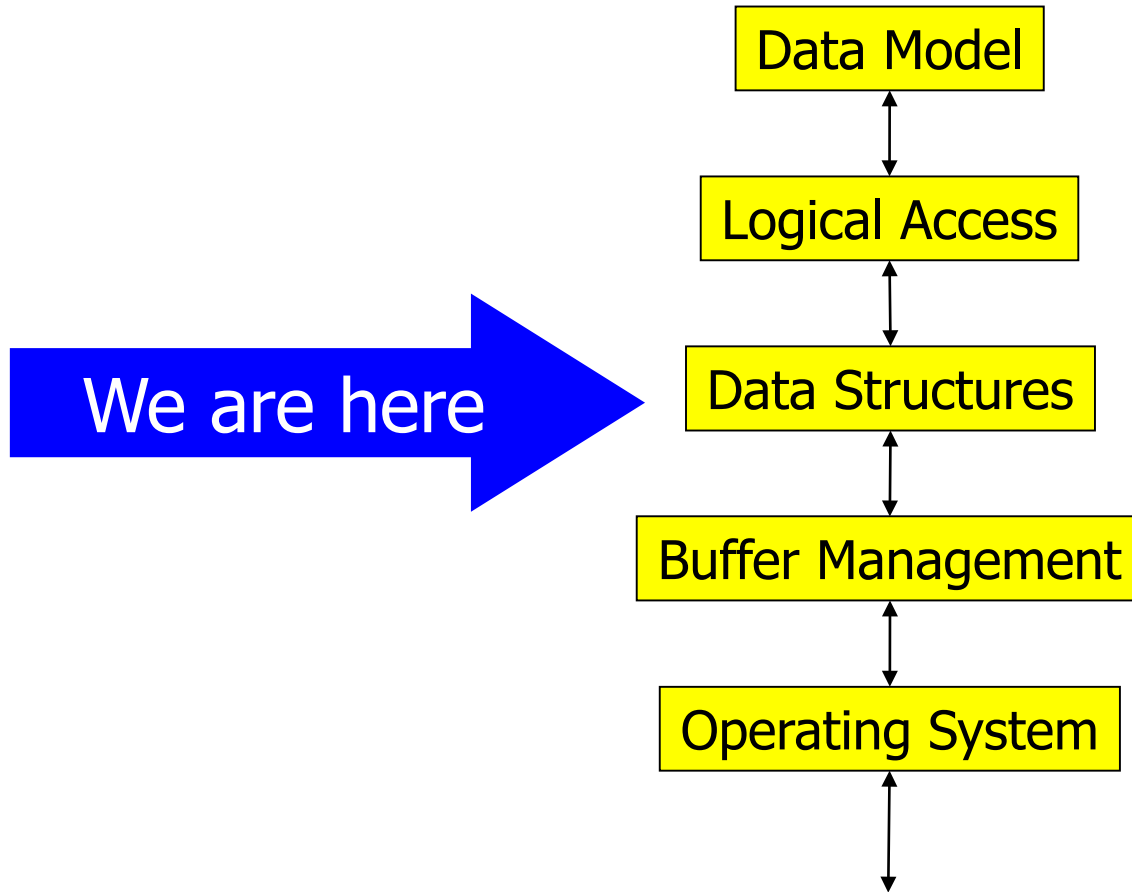
Datenbanksysteme II: File Structures

Ulf Leser

Content of this Lecture

- File structure
 - Heap files
 - Sorted files
 - Index Files
 - Hierarchical Index Files
 - B*-Trees

5 Layer Architecture

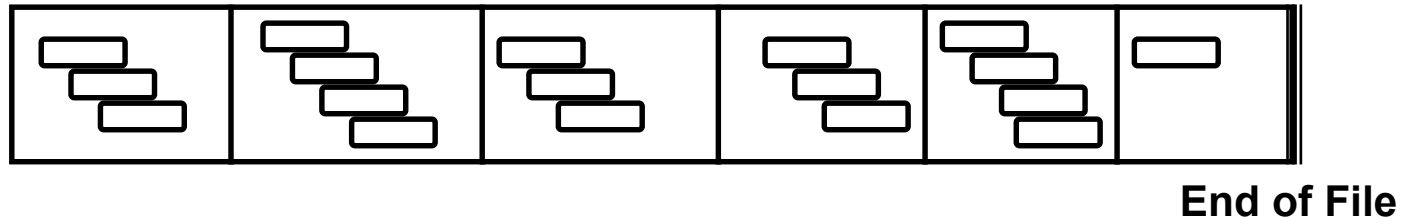


Files and Storage Structures

- We have
 - Records are stored in blocks **without particular order**
 - Makes INSERTs and DELETEs faster
 - Blocks are managed/cached by the buffer manager
 - Access **records by TID** through cache manager with adr-translation
- DBs mainly search **records with certain properties**
 - `SELECT * FROM COSTUMER`
`WHERE Name = "Bond"`
 - `SELECT * FROM ACCOUNT`
`WHERE Account# < 1000`
- This is **not "access by TID"**
- Do we **always need to scan** all records in all blocks?

Sequential (Heap) File

- Records are **stored sequentially** in the order of inserts



- Insert always add to end of file
- "Holes" occur if records are deleted
 - And can be reused by clever free-space management
- Minimal number of blocks : $b = \lceil n / R \rceil$
 - n = number of records, R = number of records per block
- Better to **keep some space free** for growing records
 - Fraction depends on expected read/write ratio

Operations on Heap Files

- In the following: We assume **highly selective** searches
 - Only a few records qualify
 - If most records are selected, scanning is hard to beat – see later
- Search by value of any attribute
 - **$b/2$ block IO** in case of successful searching a PK (on average)
 - b block IO in case of failure or searching **non-unique** values
- Insert record without duplicate checking
 - Remember: relational model is per-se **duplicate-free**
 - Simple case: read last block, add, write last block: 2 IO
 - **Free list management** makes things more complicated
- Insert record with duplicate checking / delete record
 - $b/2$: for successful search and no insert (on average)
 - $b+1$: in case of search without success and insert

Content of this Lecture

- File structure
 - Heap files
 - Sorted files
 - Index Files
 - Hierarchical Index Files

Sorted Files

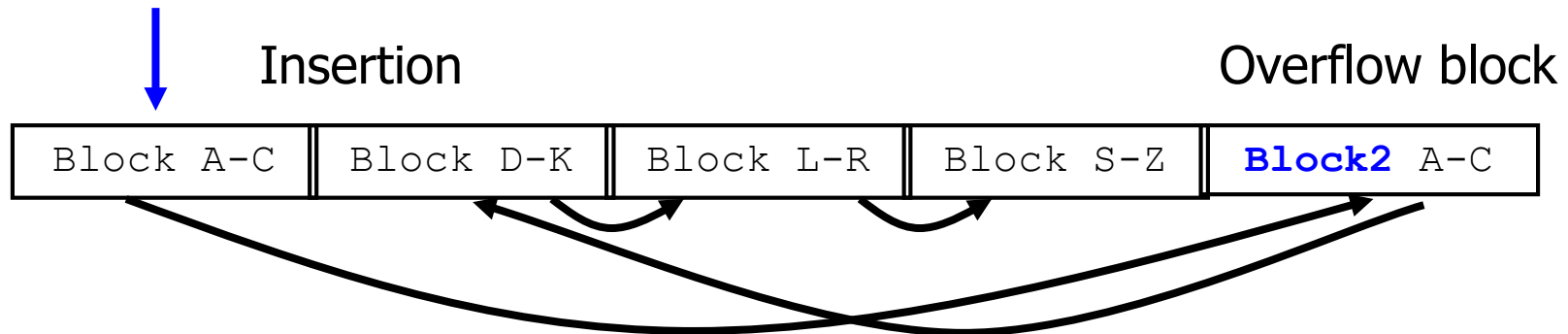
- Sort records in file according to some attribute
 - Fast searching when this attribute is search key
 - More complex management – order must be preserved
- Operations and associated costs
 - Search (using binsearch on blocks)
 - $\log(b)$ IO; searching in block is free (as always)
 - But: That's mostly random-access IO
 - Change / delete record based on value
 - First search in $\log(b)$
 - Write changes / mark space as free
 - Insert record
 - First search correct position in $\log(b)$
 - Then do what?

Inserting in a Sorted File

- General: **Reserve free space** in every new blocks
 - Don't fill blocks to 100% when allocated first time
 - Chances increase that later insertions can be handled in the block
- Option 1: Use space available in block
 - 1 additional IO for writing
- Option 2: Check neighbors
 - See X blocks down and X blocks up in the file (usually X=1)
 - When space is found, in-between records need to be moved
 - Add change block translation table
 - Cost: depends on **how far we need/want to look**
 - Typical: +4 IO if X=1
- Option 3: ...

Overflow Blocks

- Option 3: Generate **overflow blocks**
 - Create a new, “orthogonal” overflow block and insert record
 - When blocks are connected by pointers
 - Sorted table scan still possible as blocks are chained in correct order
 - New block will **not be in sequential physical order**
 - When block is added at end of file
 - Sequential table scan still possible, but **not in order** of attribute



Disadvantages Sorted Files

- Additional cost for keeping order
 - INSERT requires $\log(b)$ search first, management of overflow blocks, more random-access IO ...
- We can sort by only **one search key**
 - Searching on other attributes requires linear scans
 - With more random-access
 - Many ideas: See multi-dimensional indexes
- Search time grows only logarithmically with b
 - For 10.000.000 blocks, we need ~ 23 IO
- Can we do better?

Idea 1: Interpolated search; Build Histograms

- Partition key value range into buckets
- Count number of **keys in each bucket**
- Searching: Start at **estimated position** of search key
 - Example: Search “Immel”, [A-C]=7500, [D-F]=6200, [G-I]=3300
 - Estimated position: $7500 + 6200 + (3300/3) * 2 + \dots$
 - Continue with local search around estimated position
- Advantages
 - Very little IO if data is **uniformly distributed** – exact estimates
 - Small space consumption when few buckets are used
 - But: the more buckets (higher granularity), the better the estimates
- Disadvantages (see later for ideas)
 - Histograms (statistics) need to be **maintained**
 - Potential bottleneck for concurrent update operations in same bucket
 - Choosing **optimal bucket number** and range is difficult

Content of this Lecture

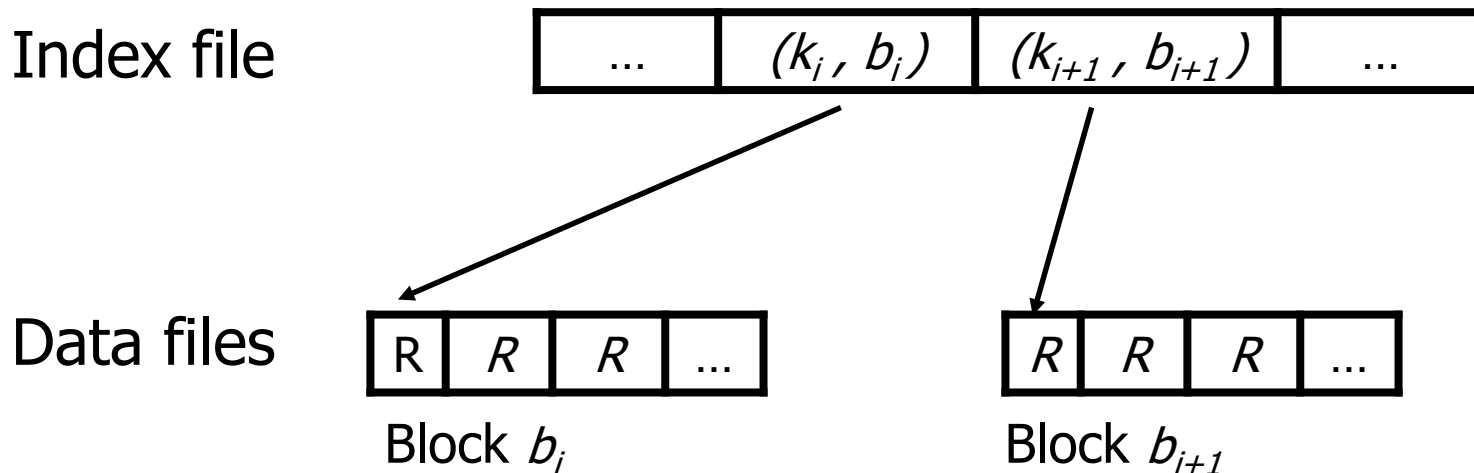
- File structure
 - Heap files
 - Sorted files
 - Index Files
 - Hierarchical Index Files

Idea 2: Decrease b

- Keep only Essential Info in less Blocks
- Use **additional file (index)** storing **only keys and TIDs**
- Searching: (Bin-)search index, then access data by TID
- Advantages
 - Data file **need not be sorted** any more:
 - Faster inserts in data file, but additional cost for **updating index**
 - Integer keys: Fixed-length index entries; strings: Use fixed-length prefix
 - Faster search due to smaller records and **less blocks**: $b_{\text{index}} < b_{\text{records}}$
 - **Several indexes** can be build for different attributes
 - More flexibility, more update cost
- Disadvantages
 - More files to manage, lock, recover, ...
 - But no more fast sorted scans of entire table

Further Decrease b: Index Sequential Files

- Data file has records **sorted on key**
- Index stores pairs (first key, pointer) for **each data block**
 - **Sparse index**: Only put **first key per block** in index
- Constraint (k_i, ptr) : For all k in ptr^\uparrow : $k_i \leq k \leq k_{i+1}$



Searching in Index-Sequential Files

- Search key in index using binsearch, then access by TID
- Advantages
 - Index has only few keys: $b_{\text{index}} \ll b_{\text{records}}$
 - Assume 10.000.000 records of size 200, $|\text{blockID}|=10$, $|\text{search key}|=20$, block size=4096
 - Number of blocks $b = 10.000.000 * 200 / 4096 = 500.000$
 - Access if kept sorted: $\log(500.000) \sim 19$ IO
 - Index-seq file: $\log(500.000 * (10 + 20) / 4096) \sim 12$ IO +1 for data
 - Chances that **index fits into main memory**
- Disadvantages
 - Only possible for one attribute (data file must be sorted)
 - More administration (compared to heap file)

Index-Sequential Files: Other Operations

- Insert record r with key k
 - Search for block b_i with $k_i \leq k \leq k_{i+1}$
 - Free space in block? Insert r ; done
 - Else, either check neighbors
 - Index needs to be updated, as block's first keys change
 - ... or create overflow blocks
 - Option 1: New block not represented in index; index not updated
 - More IO when searching data, as overflow blocks need to be followed
 - Option 2: Index is updated (more IO at time of insertion)
 - We need to insert into the index – leave free space in index blocks!
- Ideas for improving search further?

Multi-Level Index Files

Sparse
2nd level

10	—
90	
170	
250	

330	
410	
490	
570	

Sparse
1st level

10	—
30	
50	
70	

90	
110	
130	
150	

170	
190	
210	
230	

Sorted File

10	
20	

30	
40	

50	
60	

70	
80	

90	
100	

Hierarchical Index-Sequential files

- Build a sparse, **second-level index** on the first-level index
 - First level may be sparse or dense
 - All but the first level are sparse
- Advantages
 - Access time is reduced further
 - Assume 10.000.000 records of size 200, $|\text{blockID}|=10$, $|\text{search key}|=20$, block size=4096, $b = 500.000$
 - Index-seq file: $\log(500.000 \cdot (10+20)/4096) = 12+1$ block IO
 - With second level: $\log(3662 \cdot (10+20)/4096) = \mathbf{5+2}$ blocks IO
 - With three levels: $\log(28 \cdot (10+20)/4096) = 1+3$
 - Higher levels are very small – **cache permanently**
- With more than one level, inserting becomes tricky
 - Either degradation (overflows) or costly reorganizations
 - Alternative: **B-trees** (later)

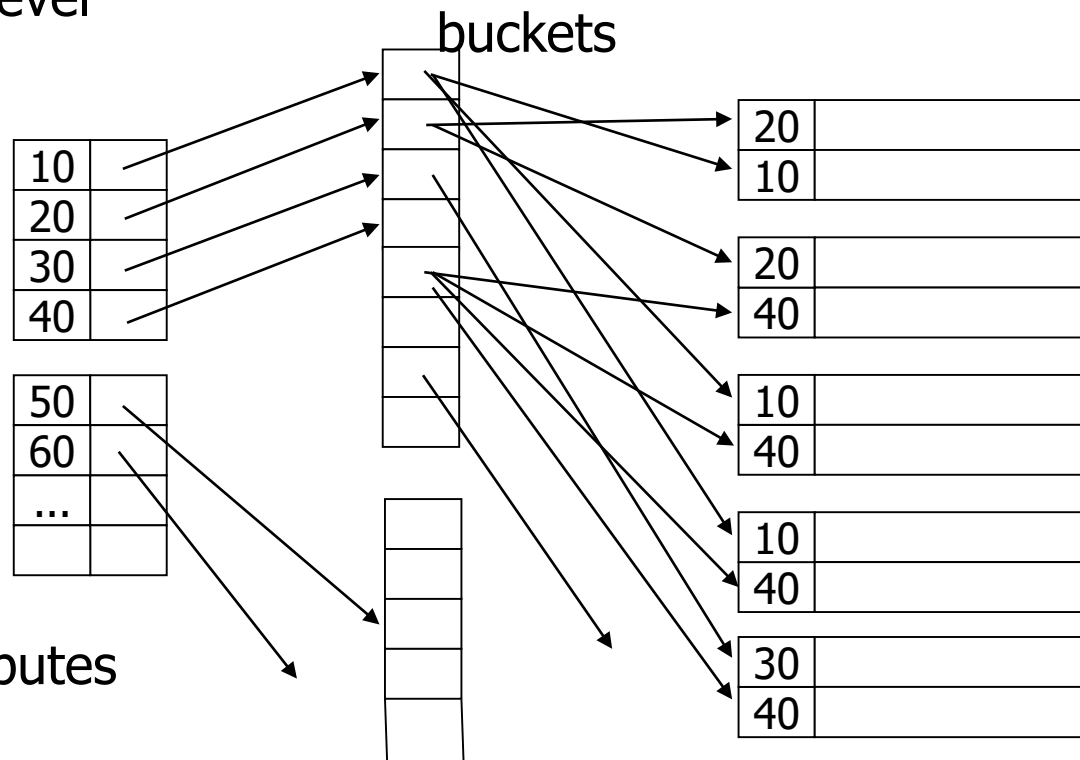
Index Files and Duplicates

- What happens if search key is not unique?
- Index file may
 - Store duplicates: one pointer for each record
 - Ignore duplicates: one pointer for **each distinct value**
 - Smaller index file
 - Requires sorted data file
 - **“Semi-sparse” index**
- **Index degradation**
 - If only **few distinct values** exist, every search selects many TID
 - E.g. index on Boolean attributes – index has only two different entries
 - Semi-sparse index leads to less IO
 - But selects blocks in random IO – scan might be cheaper

Secondary Index Files

- Primary ind.: Index on attribute on which **data file is sorted**
- Secondary index: Index on any **other attribute**
 - Cannot exploit order in data file
 - **Must be dense** at first level

- Improvement:
Use **intermediate buckets only for TIDs**
 - Buckets hold TIDs sorted by index key
 - Buckets don't store key values
 - Advantageous for **low cardinality** attributes



Indexes in Oracle

- Data files are heap files
 - Exception: **Index-organized tables (IOT)**
 - Recommended only for “read-only” tables
 - Every primary key is indexed
 - Every UNIQUE attribute is indexed
 - Default: B* tree
 - Alternatives: Multidim index, hash index, bitmap index, ...
- **Join index**: Index on attribute of foreign table with FK/PK
- **Cluster index** – cluster two tables and index common key
 - Example: Cluster department and employee on common depNum
 - Tuples with same depNum will go into same data block
 - Cluster index: Create index on depNum (~ persistent join)
 - Oracle has no **clustered indexes** – use index-organized tables

Content of this Lecture

- File structure
 - Heap files
 - Sorted files
 - Index Files
 - Hierarchical Index Files
 - Excursion: Indexing texts

Excursion: Indexing Text

- Information retrieval
 - Searching documents with keywords
 - Typically, each document is represented as “bag of words”
 - Queries search for documents containing a set of words
- Naïve relational database way fails
 - Indexed varchar2(64KB) attribute containing text
 - Not efficient for keyword queries (INSTR())
 - We cannot store each word in an extra column
- Alternatives?

Inverted Lists

- Build a **secondary, bucketed index on the words**
- Find documents by intersecting buckets
 - Enables AND, NOT or OR

