course: Database Systems (A7B36DBS)

# Implementation of database structures

Doc. RNDr. Irena Holubova, Ph.D.

Acknowledgement: The slides were kindly lent by Doc. RNDr. Tomas Skopal, Ph.D., Department of Software Engineering, Charles University in Prague

# Today's lecture outline

- disk management, paging, buffer managerdatabase files organization
- indexing
  - B<sup>+</sup>-tree
  - bitmaps
  - hashing

### Three layers of database modeling

#### abstraction conceptual layer

- models a part of the "structured" real world relevant for applications built on top of your database
  - real world part
     real-world entities and relationships between them
- different conceptual models (e.g. ER, UML)
- logical layer
  - specifies how conceptual components are represented in logical machine interpretable data structures
  - different logical models (e.g. object, relational, object-relational, XML, graph, etc.)
- physical model
  - specifies how logical database structures are implemented in a specific technical environment

implementation

data files, index structures (e.g. B+ trees), etc.

# Introduction

- relations/tables are stored in files on the disk
  we need to organize table records within a file
  - efficient storage, update and access

Example:

Employees (name char(20), age integer, salary integer)

# Paging

- records are stored in disk pages of fixed size (a few kB)
  - reason: hardware
    - assuming a magnetic disk based on rotational plates and reading heads
  - the data organization must be adjusted w.r.t. this mechanism
- the HW firmware can only access entire pages
  - I/O operations reads, writes
- real time for I/O operations =
  - = seek time + rotational delay + data transfer time
- sequential access to pages is much faster than random access
  - the seek time and rotational delay not needed

Example: reading 4 KB could take 8 + 4 + 0,5 ms = 12,5 ms; i.e., the reading itself takes only 0,5 ms = 4% of the real time!!!

# Paging

- I/O is a unit of time cost
- the page is divided into slots, that are used to store records
  - a record is identified by rid (record id) = page id + slot id
- a record can be stored
  - in multiple pages  $\Rightarrow$ 
    - better space utilization
    - need for more I/Os for record manipulation
  - in a single page (assuming it fits)  $\Rightarrow$ 
    - a part of page may not be used
    - less I/Os
  - ideally: records fit the entire page

# Paging

- only fixed-size data types are used in the record ⇒ fixed record size
- also variable-size data types are used in the record
   ⇒ variable size of the records,
  - e.g., types varchar (X), BLOB, ...
- fixed-size records = fixed-size slots
- variable-size records = need for slot directory in the page header

### Fixed-size page organization, example



### Variable-size page organization, example



# **Buffer management**

- buffer = a piece of main memory for temporary storage of disk pages
  - disk pages are mapped into memory frames 1:1
- every frame has 2 flags:
  - pin\_count = number of references to the page in frame
  - **dirty** = indication of containing a modified record
- buffer manager
  - implements read and write operations for higher DBMS logic
- read: retrieves the page from buffer + increasing pin\_count
  - if it is not there, it is first fetched from the disk
- write: puts the page into the buffer + setting dirty
- if the buffer is full (during read or write), some page must be replaced ⇒ various policies, e.g., LRU (least-recently-used),
  - if the replaced page is **dirty**, it must be stored

## **Buffer management**



# Database storage

- data files contain table data
- index files speed up processing of queries
- system catalogue contains metadata
  - table schemas
  - index names
  - integrity constraints, keys, etc.

## Data files

- 1. heap
- 2. sorted file
- 3. hashed file

Observing average I/O cost of simple operations:

- 1) sequential access to records
- 2) searching records based on equality (w.r.t search key)
- 3) searching records based on range (w.r.t search key)
- 4) record insertion
- 5) record deletion

Cost model:

N = number of pages, R = records per page

## Simple operations, SQL examples

- sequential reading of pages
   SELECT \* FROM Employees
- searching on equality
   SELECT \* FROM Employees WHERE age = 40
- searching on range
   SELECT \* FROM Employees
   WHERE salary > 10000 AND salary < 20000</li>
- record insertion
   INSERT INTO Employees VALUES (...)
- record deletion based on rid
   DELETE FROM Employees WHERE rid = 1234
- record deletion
   DELETE FROM Employees WHERE salary < 5000</li>

# Heap file

- records stored in pages are <u>not ordered</u> (e.g., according to key)
  - they are stored in the order of insertion
- page search can only be achieved by sequential scan (GetNext operation)
- quick record insertion (at the end of file)
- deletion problems: "holes" (pieces of not utilized space)

## Maintenance of empty heap pages

- double linked list
  - header + lists of full and non full pages
- page directory
  - linked list of directory pages
  - every item in the directory refers to a data page
    - flag = item utilization

## Maintenance of empty heap pages



# Heap, cost of simple operations

- sequential access = N
- search on equality = N
- search on range = N
- record insertion = 1
- record deletion
  - 2,

assuming **rid** based search costs 1 I/O

N or 2\*N,

if deleted based on equality or range

# Sorted file

- records stored in pages based on an <u>ordering according to</u> <u>a search key</u>
  - single or multiple attributes
- file pages maintained contiguous, i.e., no "holes"
- fast: search on equality and/or range
- slow: insertion and deletion
  - "moving" the rest of pages
- in practice:
  - sorted file at the beginning
  - each page has an overhead space where to insert
  - if the overhead space is full, update pages are used (linked list)
  - a reorganization needed from time to time
    - i.e., sorting

## Sorted file, cost of simple operations

- sequential access = N
- search on equality = log<sub>2</sub>N
- search on range =  $\log_2 N + M$ 
  - where M is the number of relevant pages
- record insertion = N
- record deletion = log<sub>2</sub>N + N (based on key)

# **Hashed file**

- organized in <u>K buckets</u>
  - a bucket is extensible to multiple disk pages
- a record is inserted into/read from a bucket determined by <u>hashing function f</u> applied on search key
  - bucket id = f(key)
- if the bucket is full, new pages are allocated and linked to the bucket (linked list)
- fast search / deletion on equality
- higher space overhead, problems with chained pages (solved by dynamic hashed techniques)

## Hashed file



### Hashed file, cost of simple operations

- sequential access = N
- search on equality = N/K (best case)
  - K = number of buckets
- search on range = N
- record insertion = N/K (best case)
- deletion on equality = N/K + 1 (best case)

# Indexing

- index is a <u>helper structure</u> that provides fast search based on search key(s)
- organized into disk pages (like data files)
  - usually different file than data files
- usually contains only search keys and links to the respective records
  - i.e., rid
- need much less space than data files
  - e.g., 100x less

# Indexing, principles

- index item can contain
  - the whole record (then index and data file are the same)
  - pair <key, rid>
  - pair <key, rid-list>, where rid-list is a list of links to records with the same search key value
- **1. clustered:** ordering of index items is (almost) the same as ordering in the data file
  - tree-based index, index containing the entire records, hashed index,
     ...
  - primary key = search key used in clustered index
- 2. **unclustered:** the order of search keys is not preserved

# Indexing, principles

#### **CLUSTERED INDEX**



#### UNCLUSTERED INDEX



#### **Clustered index:**

Pros: huge speedup when searching on range – result record pages are read sequentially from data file

Cons: large overhead for keeping the data file sorted

#### B<sup>+</sup>-tree



- extends B-tree
  - balanced tree-based index
- provides logarithmic complexity for insertion, search on equality (no duplicates), deletion on equality (no duplicates)
- guarantees 50% node (page) utilization
- B<sup>+</sup>-tree extends B-tree by
  - all keys are in the leaves inner nodes contain indexed intervals
  - linking leaf pages for efficient range queries

### B<sup>+</sup>-tree, schema



# **Hashed index**

- similar to hashed data file
  - i.e., buckets + hashing function
- buckets contain only key values together with the rids
- same pros/cons

# **Bitmaps**

- suitable for indexing attributes of low-cardinality data types
  - e.g., attribute FAMILY\_STATUS = {single, married, divorced, widow}
- for each value *h* of an indexed attribute *a* a bitmap (binary vector) is constructed, where 1 on *i*<sup>th</sup> position means the value *h* appears in the *i*<sup>th</sup> record (in the attribute *a*), while it holds
  - bitwise OR = 1 (every attribute has a value)
  - bitwise AND = o (attribute values are deterministic)

Name	Address	Family status	
John Smith	London	single	
Rostislav Drobil	Prague	married	
Franz Neumann	Munich	married	
Fero Lakatoš	Malacky	single	
Sergey Prokofjev	Moscow	divorced	

single	married	divorced	widow
1	ο	ο	0
0	1	ο	0
0	1	ο	0
1	0	Ο	0
ο	0	1	0

# **Bitmaps**

- query evaluation
  - bitwise operations with attribute bitmaps
  - resulting bitmap marks the queried records
- example
  - Which single or divorced people did not complete the military service? (bitmap(single) OR bitmap(divorced)) AND not bitmap(YES)



# **Bitmaps**

#### pros

- efficient storage, could be also compressed
- fast query processing, bitwise operations are fast
- easy parallelization
- CONS
  - suitable only for attributes with small cardinality domain
  - range queries get slow linearly with the number of values in the range (bitmaps for all the values must be processed)