# **Physical Data Organisation**
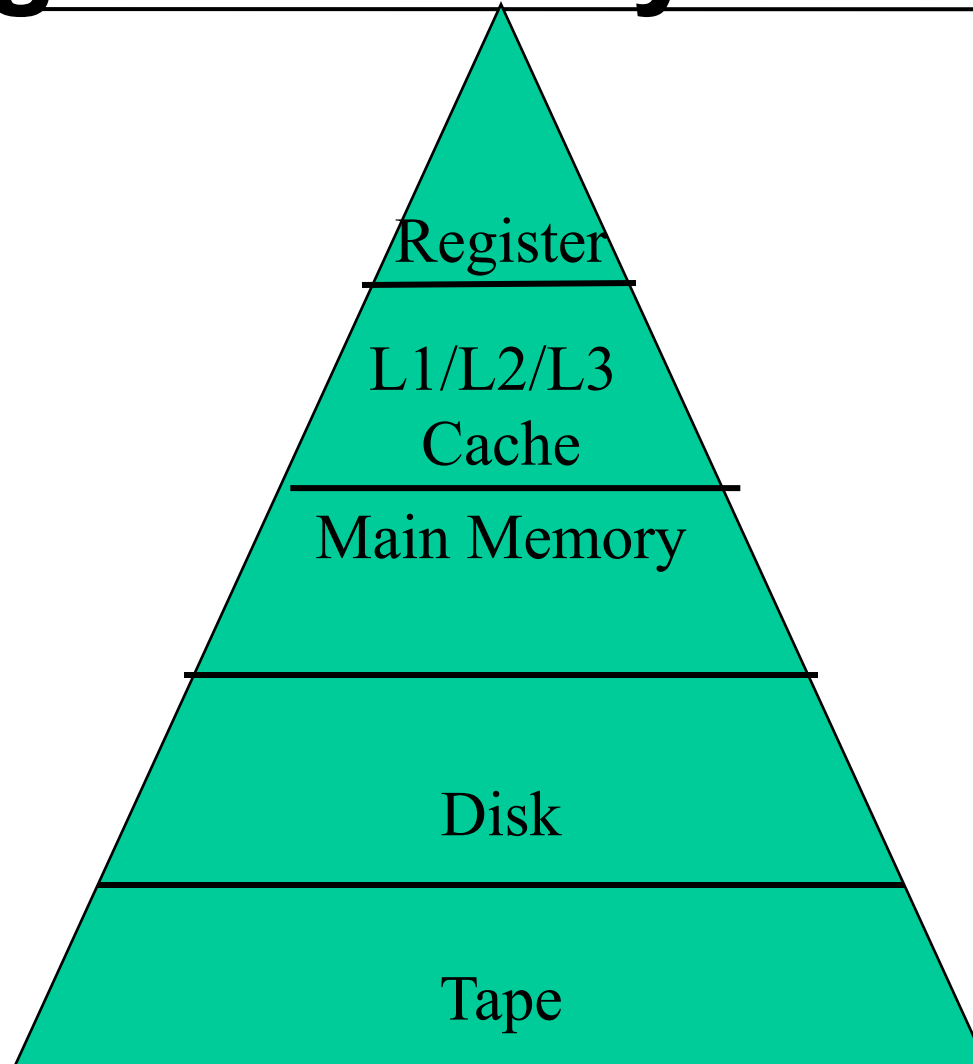
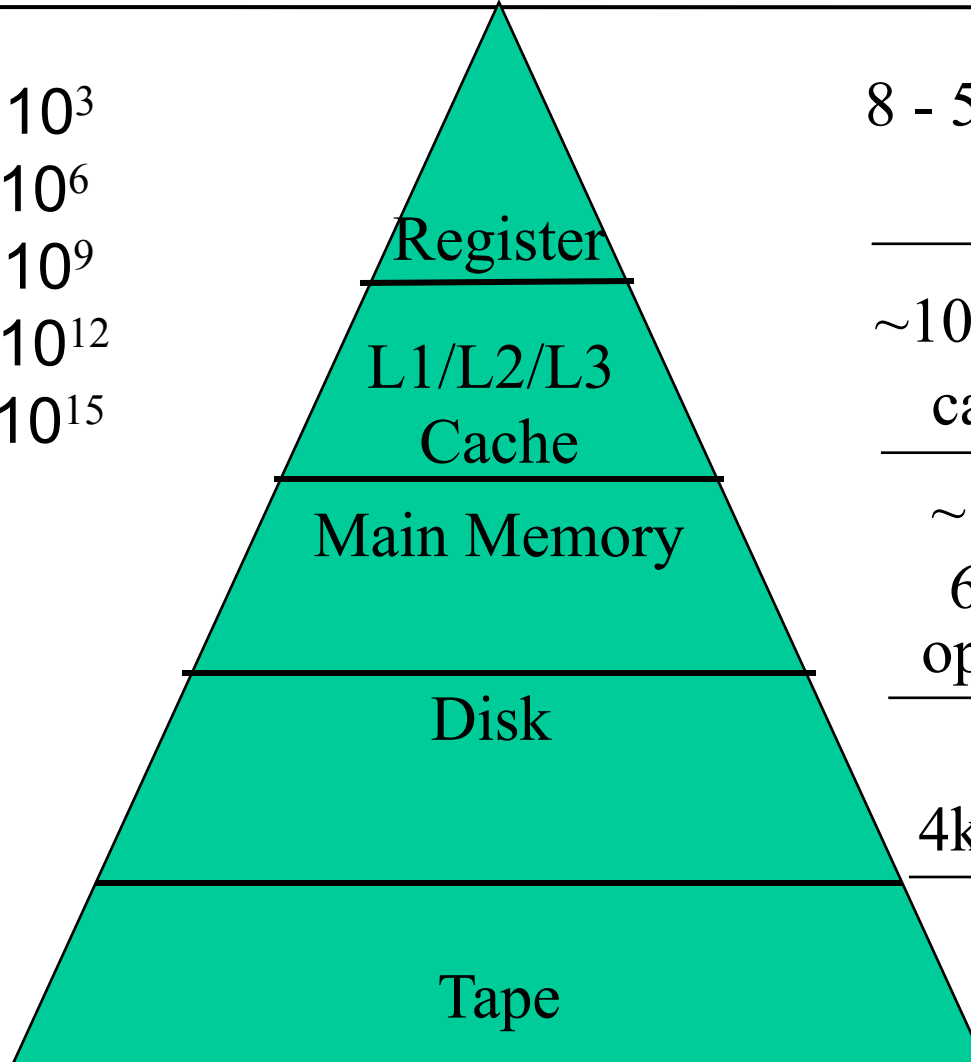Topics:

- Storage hierarchy
- Storage structures
- ISAM
- B-Trees
- Hashing
- Clustering

# Storage Hierarchy

Register

L1/L2/L3 Cache

Main Memory

Disk

Tape

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Storage Hierarchy

1 K (Kilo)   = $10^3$
1 M (Mega) = $10^6$
1 G (Giga)  = $10^9$
1 T (Tera)   = $10^{12}$
1 P (Peta)  = $10^{15}$

Register

L1/L2/L3 Cache

Main Memory

Disk

Tape

8 - 512 Byte/Register
Compiler

~10 MB Byte/Cache
cache-controller

~100 GB-range,
64B block size
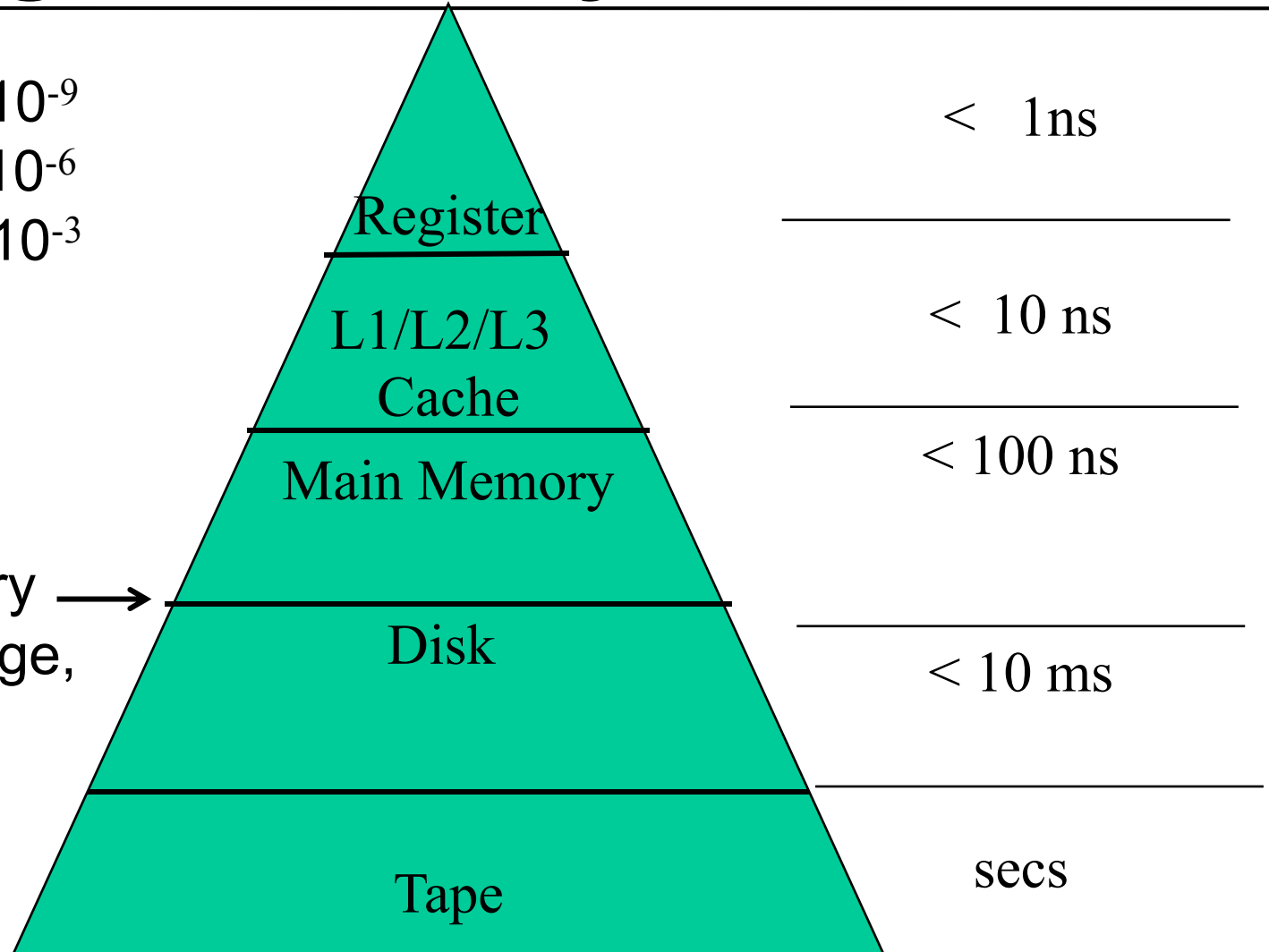operating system

~1 TB-range
4kB blocks | user

PB-range
user

Database System Concepts for Non-
Computer Scientists WS 2018/2019

# Storage Hierarchy

1 n (nano)  = $10^{-9}$
1 μ (micro) = $10^{-6}$
1 m (milli)  = $10^{-3}$

Register

L1/L2/L3 Cache

Main Memory

(Flash-Memory ⟶
Lower TB-range,
< 100 μs)

Disk

Tape

<   1ns

<  10 ns

< 100 ns

< 10 ms

secs

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Storage Hierarchy

Room (1min)

Building (10min)

City (1.5h)

(Mars (2 month))

Pluto (9 years)

Andromeda

(2000 years)

< 1ns

register

< 10ns

L1/L2/L3 Cache

< 100ns

Main Memory

< 10 ms

Disk

secs

Tape

Factor 100.000

# Storage Hierarchy

Room (1min)

Building (10min)

City (1.5h)

(Mars (2 month))

Pluto (9 years)

Andromeda

(2000 years)

< 1ns

register

< 10ns

L1/L2/L3 Cache

< 100ns

Main Memory

< 10 ms

Disk

secs

Tape

# Buffer Management



Main Memory

Disk

# Buffer Management



Main Memory

replace     fill

disk ~ persistent DB

Disk

Database System Concepts for Non-Computer Scientists WS 2018/2019
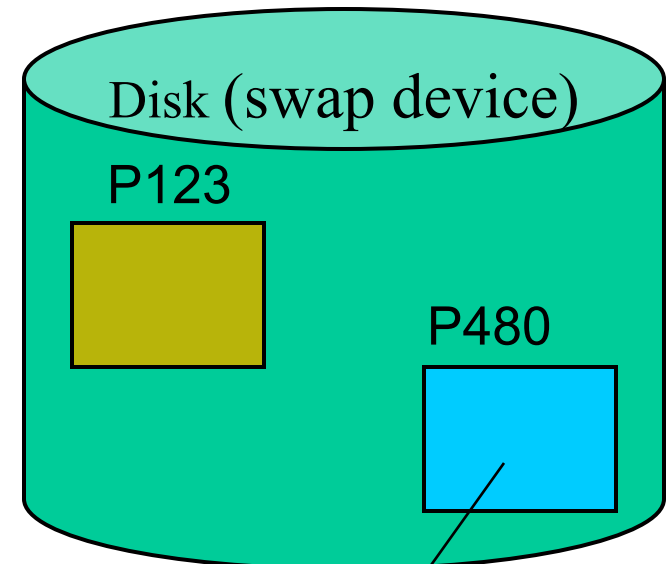
# Fill and replace pages

- System buffer is divided in frames of equal size
- A frame can be filled with one page
- Overflow pages are swapped on disk

Main Memory

| 0 | 4K | 8K | 12K |
|-----|-----|-----|-----|
| 16K | 20K | 24K | 28K |
| 32K | 36K | 40K | 44K |
| 48K | 52K | 56K | 60K |

Disk (swap device)

P123

P480

(Buffer) Frames

Page

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Addressing tuples on disk

Database System Concepts for Non-
Computer Scientists WS 2018/2019

# Moving within a page

# Moving from one page to another



Forward

# Moving from one page to another

TID

4711 | 2

1 2 3

5001 ∘ Grundzüge ∘ . . .

5041 ∘ Ethik ∘ . . .

4812 | 3

Page 4711

1 2 3

4052 ∘ Mathematische Logik für Informatiker ∘ . . .

Page 4812

With the next move the „Forward" on page 4711 is altered (no more Forward to page 4812)

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Storage Summary



### Storage Hierarchy

- Huge/slow storage vs small/fast memory
- Very important for DBMSs design
- Algorithms need to be aware of performance difference and place data optimally

### Buffer Management

- A component of the database
- Migrates pages between disc and main memory
- Keeps hot pages in DRAM and cold ones on disc

### Tuple IDs

- Used to locate a tuple
- Composed of page identifier and a page-local record identifier

# Storage Operations

Full table scan: Retrieve all tuples from a table
**select** * **from** students;


Point query: Find one specific tuple
**select** * **from** students
       **where** studNr =26120;


Range query: Find one specific tuple
**select** * **from** students
       **where** 26000 <= studNr
        **and** studNr < 27000;

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Index Motivation

RAM

--------------------------------------------

Disk

| 2755 Schopenhauer | 26830 Aristoxenos | 24002 Xenokrate |
| 29120 Theophrastos | 29555 Feuerbach | 26120 Fichte |
| | 28106 Carnap | 25403 Jonas |

Page 1        Page 2        Page 3

# Index Motivation

Full table scan: Load all pages from disk, one by one.

Point query: When we sort the data, we can find a key more easily. Compare to a dictionary: You can look at the page in the middle to determine in which half the word you are looking for is and then continue this process with the new found half. This way you need to look at much less pages then reading it front to back.

Range query: When the data is stored in a sorted fashion, a range query can be processed as a combination of a point query (to find the starting point) and then a scan.

# Index Motivation

RAM

--------------------------------

Disk

| | | |
|---|---|---|
| 24002 Xenokrate<br>25403 Jonas<br>26120 Fichte | 26830 Aristoxenos<br>2755 Schopenhauer | 28106 Carnap<br>29120 Theophrastos<br>29555 Feuerbach |
| Page 1 | Page 2 | Page 3 |

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Data Transfer

Simple query execution:

**select** * **from** students **where** studNr=26120;

Get one tuple/page after the other to the main memory and evaluate predicates.

→ Most expensive way ☹

→ Mostly only a small fraction of the tuples fulfills the query

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Index Structures

- Index structures are used to keep the data volume to be transferred from disk to main memory small

- Only that part of the data which is really needed to answer the query is transferred

- Two main indexing methods:
    - o Hierarchical (trees)
    - o Partitioning (Hashing)

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Hierarchical Indexes

We consider two hierarchical index structures:
- ISAM (Index-Sequential Access Method)
- B-Trees

- ISAM is the predecessor of B-Trees
- Main idea: sort tuples on the indexed attribute and create an index file on it
- Similar to a thumb index in a book

# ISAM Example

Index pages

⑨
① 25403 ③ 26830 ⑤ ────────→ ⑥ 27550 ④ 29120 ②

⑤

① 24002    25403
  Xenokrates  Jonas

③ 26120    26830
  Fichte   Aristoxenos

⑥ 27550
  Schopenhauer

④ 28106    29120
  Carnap  Theophrastos

② 29555
  Feuerbach

Sorted →

Datapages

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Example cont.

- Student with student number 28106 is searched
- Go through the **index pages** (starting at ⑨) and look where 28106 fits => the first larger (or equal) key is 29120 (on page ⑤)
- From there you get the referenced **data page** => ④
- **Advantage**: Number of index pages is **much** less than number of data pages, i.e. you save I/O
- You can also answer **range queries**, e.g. all StudNr between 26830 and 29120 : find as a start the first fitting data page for 26830 and from there on you can go **sequentially** through the data pages until StudNr 29120

# **Problems with ISAM**

Simple and fast search but **maintenance of index** is expensive:

- Inserting a tuple in a full data page: need to make room in **dividing data page into two** (because we need to keep the tuples sorted)

- This creates a **new entry** on an **index page**

- Inserting an entry in a full index page leads to **shifting the entries** to make room

- Although the number of index pages is smaller than the number of data pages **going through the index pages** can nevertheless **take a long time**

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Exercise 1

# Exercise 2

# Advancement

Idea:

Why not have **index pages for the index pages**?

→ This is in principle the idea of a **B-Tree**

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Trees

All trees in computer science

… have nodes

… have edges

… have a root (at the top!)

… have leaves (at the bottom!)

… are often balanced

(otherwise in extreme cases rather a chain)

Root

Schematic depiction of a
balanced tree:

Leaves

# B-Tree (degree = 2)

6
57
Rudi
5 9 - - -

5
33 44
John Jana
1 7 3 - -

9
65 85
Silke Flo
8 2 4 - -

1
24 25 31
Andre Luke Varun

7
40 42
Max Chris

3
46 55 56
Michi Alf Steffi

8
58 60 62 63
Andy Alex Nina Chris

2
79 81 83
Jessi Tom Harry

4
89 93
Max Marta

Database System Concepts for Non-Computer Scientists WS 2018/2019

# B-Tree (degree = 2)

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Node Structure



Inner nodes                                                    Leaf nodes

Tree properties:
- One node is one page
- Tree is balanced
- Node utilization at least 50%

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Properties of a B-Tree

B-Tree of degree *i* has following properties:

- Every path from the root to a leaf has the same length
- Every node – except the root – has at least *i* and at most *2i entries* (in the example above *i=2*)
- Entries in every node are sorted
- Every node – except the leaves – with *n* entries has *n+1* children

# Properties of a B-Tree

- *Let
  $p_0, k_1, p_1, k_2, \ldots k_n, p_n$
  be entries in a node ($p_j$ are page identifier, $k_j$ keys)*

Then the following holds:
1. Sub-tree in $p_0$ contains only keys smaller than $k_1$
2. $p_j$ has a sub-tree with keys between $k_j$ and $k_j+1$
3. Sub-tree being referenced by $p_n$ contains only keys greater than $k_n$

| $p_0$ | $\dfrac{k_1}{d_1}$ | ... | $\dfrac{k_j}{d_j}$ | $p_j$ | $\dfrac{k_{j+1}}{d_{j+1}}$ | ... | $\dfrac{k_n}{d_n}$ | $p_n$ |
|---|---|---|---|---|---|---|---|---|

# B-Tree (degree = 2)

Database System Concepts for Non-Computer Scientists WS 2018/2019

# B-Tree (degree = 2)

Database System Concepts for Non-Computer Scientists WS 2018/2019

# B-Tree Size

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Insert Algorithm

1. Find the proper leaf node to insert new key
2. Insert key there
3. If node full
   i.  Divide node into two and extract median
   ii. Insert all keys smaller than median into left node, all keys greater than median into right node
   iii. Insert median in parent node and adapt pointers
4. If parent node full
   i.  If root node then create new root node, insert median, and adapt pointers
   ii. Otherwise repeat 3. with parent node

Database System Concepts for Non-Computer Scientists WS 2018/2019

# B-Tree (degree = 2)



|    |    |   |   |
|------|------|---|---|
| 57 |    |   |   |
| Rudi |    |   |   |

| 33 | 44 |   |   |
|------|------|---|---|
| John | Jana |   |   |

| 65 | 85 |   |   |
|------|------|---|---|
| Silke | Flo |   |   |

| 24 | 25 | 31 |   |
|------|------|------|---|
| Andre | Luke | Varun |   |

| 40 | 42 |   |   |
|------|------|---|---|
| Max | Chris |   |   |

| 46 | 55 | 56 |   |
|------|------|------|---|
| Michi | Alf | Steffi |   |

| 58 | 60 | 62 | 63 |
|------|------|------|------|
| Andy | Alex | Nina | Chris |

| 79 | 81 | 83 |   |
|------|------|------|---|
| Jessi | Tom | Harry |   |

| 89 | 93 |   |   |
|------|------|---|---|
| Max | Marta |   |   |

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Gradual Assembly of a B-Tree of Degree i=2

See:
https://db.in.tum.de/teaching/ws1819/DBSandere/BTreeExample.pdf

In the internet there are a number of animation programs for B-Trees – **no warranty!**

https://www.cs.usfca.edu/~galles/visualization/BTree.html
looks quite good, but uses a different notation for the maximal node size and does not handle node underflows.

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Delete Algorithm

Read the literature or example on lecture website

Database System Concepts for Non-Computer Scientists WS 2018/2019

# B+-Trees

- Performance of a B-Tree heavily depends on height: on average $\log_K(n)$ page accesses to read one data element
  (k=degree of branching, n=number of indexed data elements)
  $\rightarrow$ preferably high degree of branching of the inner nodes
- Storing data in the inner nodes reduces branching degree
- B+-Trees only store reference keys in inner nodes – data itself is stored in leaf nodes
- Usually leaf nodes are bidirectionally linked in order to enable fast sequential search

Database System Concepts for Non-Computer Scientists WS 2018/2019

# B+-Tree



Index pages

Leaf pages

| 57 | | | |
|----|--|--|--|

| 33 | 44 | | |
|----|----|--|--|

| 63 | 83 | | |
|----|----|--|--|

| 24 | 25 | 31 | 33 |
|----|----|----|----|
| Andre | Luke | Varun | John |

| 40 | 42 | 44 | |
|----|----|----|--|
| Max | Chris | Jana | |

| 46 | 55 | 56 | 57 |
|----|----|----|----|
| Michi | Alf | Steffi | Rudi |

| 58 | 60 | 62 | 63 |
|----|----|----|----|
| Andy | Alex | Nina | Chris |

| 65 | 79 | 81 | 83 |
|----|----|----|----|
| Silke | Jessi | Tom | Harry |

| 85 | 89 | 93 | |
|----|----|----|--|
| Flo | Max | Marta | |

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Prefix B+-Trees

- Further Improvement by use of prefixes of reference keys, e.g. with long strings as keys
- You only have to find a reference key which separates the left and the right sub-tree:
  - ➢ Disestablishment   <= E < Incomprehensibility
  - ➢ Systemprogram     <= ? < Systemprogrammer

# Several Indexes on the same Data

Primary index – Secondary index

| Students | | |
|---|---|---|
| StudNr | Name | Semester |
| 25403 | Jonas | 12 |
| 29120 | Theophrastos | 2 |
| 29555 | Feuerbach | 2 |
| 27550 | Schopenhauer | 6 |
| ⋮ | ⋮ | ⋮ |

When

- Index on StudNr?
- Index on Name?
- Index on Semester?

# Secondary indexes



Primary index

45

Index pages

9  29  35

. . .

. . .

| 1 | 2 |
|---|---|
| D1 | D2 |

Page 1

| 15 | 17 | 20 |
|---|---|---|
| D15 | D17 | D20 |

Page 2

Page n

Data pages, sorted, bidirectionally linked

. . .

Index pages

Secondary index

# DDL: Create Index

CREATE [UNIQUE] INDEX index_name
ON table_name (column_name1 [, column_name2, …])


Example:

CREATE INDEX full_name
ON Person (Last_Name, First_Name)

# Partitioning
# What is Hashing?

- (to hash = zerhacken)

- Storing tuples in a defined memory area

- Hash function: mapping tuples (key values)
  to a fixed set of function values (memory area)

- Optimal hash function:
  - injective (no identical function values for different arguments)
  - surjective (no waste of memory)

- Typical hash function h: h (x) = x mod N
  set of function values thereby {0,..., N-1}

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Hashing Example

*hash(x) = x mod 3*

25403, 'Jonas', 12
27550, 'Schopenhauer', 3
26120, 'Fichte', 10

| pid=0 | |
|-------|--|
|       |  |
| pid=1 | |
|       |  |
| pid=2 | |
|       |  |

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Hashing Example

*hash(x) = x mod 3*

25403, 'Jonas', 12
27550, 'Schopenhauer', 3
26120, 'Fichte', 10

| pid=0 | |
|---|---|
| | |
| pid=1 | (27550, 'Schopenhauer', 3) |
| | |
| pid=2 | (25403, 'Jonas', 12) |
| | (26120, 'Fichte', 10) |

Database System Concepts for Non-Computer Scientists WS 2018/2019

# Collisions

Collision handling

| pid=0 | |
|---|---|
| pid=1 | (27550, 'Schopenhauer', 3) |
| | |
| pid=2 | (25403, 'Jonas', 12) |
| | (26120, 'Fichte', 10) |

| (26123, 'Fichte', 10) | |
|---|---|
| (28106, 'Carnap', 3) | ● |

…

Could be inefficiently: hash table is too small, a bad hash function is used, unlucky, adversary input
Way out: extensible (dynamic) Hashing

# Advantages / Disadvantages Hashing

+ Few accesses to external storage
  constant cost: O(1), generally 1-2
+ Simple implementation


– Collision handling necessary
– Pre-allocation of memory area
– Not dynamic resp. only with adjustment
– **No range queries, only point queries**