

# Data-Warehouse-Technologien

## Exercise sheet 9

**Assignment 1:** Classify the index structures discussed in the lecture!

**Assignment 2:** Define the properties of an ideal index for data warehouses.

**Assignment 3:** Compare the resource consumption of B+- and standard bitmap indices for the following data:

1. number of tuples: 5000000; number of key values: 100; resource consumption for 1 TID in Byte: 1;
2. number of tuples: 5000000; number of key values: 3; resource consumption for 1 TID in Byte: 1;
3. number of tuples: 5000000; number of key values: 3; resource consumption for 1 TID in Byte: 4;

**Assignment 4:** Transform the postal code 39106 into a bitmap index using the Standard-Bitmap-Index and the Multi-Component-Bitmap-Index. How can Bitmap-Indexes be used to support range queries for postal codes?

**Assignment 5:** Construct a R-Tree for the following data. Insert the data in the order of the given list. A node has at most 3 entries, thus a MBR is described. Outline the insertion process. Assume an approach to minimize the number of boxes.

(10 Z), (20 K), (30 A), (40 L), (50 M), (60 N), (70 D), (80 X)

**Assignment 6:** What is a Gridfile? How is it created? what are its advantages compared to tree based approaches? Use the information of exercise 5 to create a grid file.

# Assignment 1

## Classify index structures discussed in the lecture

**Index structure:** - organisator: Mechanismen, die genutzt werden, um die Leistung von Abfragen auf großen Datensätzen zu verbessern  
 → Konzepte: wie Indizes organisiert werden!

### Kriterien

**Clustering:** - data that are likely to often processed together will be stored physically close to each other  
 (Gruppierung von Daten)

introduces ordering based on similarity → **Tuple clustering:** storing tuples on same physical page (z.B. alle Daten eines Kunden liegen auf 1 Seite)

→ **Page clustering:** storing related pages close together  
 if I want to add other insertion overhead  
 in secondary storage (z.B. alle Daten eines Monats nacheinander gespeichert)

**Dimensionality:** - specify how many attributes (dimensions) of the underlying relation for calculation of the index key can be used → how many dimensions can be indexed  
 (Anteil genutzter Spalten in Index)  
 → one-dimensional index  
 → multi-dimensional index

→ does order of indexed values make a query performance difference  
**Symmetry:** - performance is independent of the order of the index attributes (order does not matter)  
 (Reihenfolge der Spalten in Index)  
 - symmetrical index structures or asymmetrical → yes  
 order makes difference on performance

**Tuple references:** - type of tuple references within the index structure  
 (wie Index auf Daten zeigt)  
 ↳ reference via TID or Store values directly  
 (direkt: Index zeigt auf Daten | indirekt: Index zeigt auf zw. Speicher, wo Daten drin sind)

**dynamic behavior:** - effort to update index structure for insert, update, delete  
 (Was passiert bei Änderungen)  
 ↳ Anpassung Index, wenn Daten geändert, gelöscht, eingefügt werden  
 ↳ how flexible Insert, Update, Delete is supported  
 ↳ scales from static to dynamic → static  
 ↳ needs frequent rebuilds  
 ↳ by shrinking and increasing

# Classification

- B-Trees
- Bitmap Indexes
- multi-dimensional Index structure

## 1) B - Trees

⇒ one dimensional tree structure

clustering:

clustered (tuple clustering)

dimensionality:

one dimensional

symmetry:

asymmetric → order important

tuple references:

references are stored

dynamic behaviour:

dynamic (efficient update algorithms)

+ exam question

what is UB tree

• uses space-filling curve  
↳ 4-curve   
↳ to generate 1-dim. key  
↳ key is indexed via B-Tree

## 2.) UB - Tree

clustering:

clustered (curve puts data into order)

dimensionality:

multidimensional

symmetry:

symmetric (no preference of single dimension)

tuple references:

references are stored

dynamic behaviour:

dynamic (e.g. grows dynamically)

## 2.) Bit-Map Index

inhl. Standard bitmap, Multicomponent BM, Range Encoded BM, Interval Encoded BM

clustering:

not clustered

[scan whole bitmap  
→ similar values are not stored/clustered  
closed together)

dimensionality:

one-dimensional

symmetry:

symmetric

tuple references:

tuple references

dynamic behaviour:

high update effort  
(syn. between static and dynamic)

## Assignment 2

### Ideal Index for data warehouse

- clustering:
- (page) clustering
  - enables faster reading of large data blocks through prefetching

- dimensionality:
- multidimensional
  - e.g. Bitmap excels in handling multiple dimensions efficiently and are highly compressible

- symmetry:
- symmetric index structures
  - ensures that order of columns in query does not impact performance

- tuple references:
- direct references
  - reduce number of intermediate steps, when accessing data

- dynamic behaviour:
- resilient to updates
  - Bitmap or B-Tree
    - ↳ for frequent updates
  - Should dynamically grow

## Assignment 3 (not important for exam) !!!

### B<sup>+</sup> Tree Index Resource Consumption

- requires storing a key value and corresponding TIDs (tuple IDs)
- Cost formula:

$$B^+ \text{-Tree - cost} = (\text{Number of key values} \times \text{Average TID list size per key}) \times \text{Resource consumption per TID}$$

Average TID list size:  $\frac{\text{Number of tuples}}{\text{Number of key values}}$

### Bitmap Index structure

- each key value in a bitmap index has a bit vector representing the presence of each tuple in the dataset
- Cost formula:

$$\text{Bitmap Cost} = (\text{Number of key values} \times \text{Number of tuples}) \times \text{Resource consumption per bit}$$

$\underbrace{\qquad\qquad\qquad}_{1/8 \text{ byte (as 1 byte - 8 bits)}}$

A) 5,000,000 tuples, 100 key values, resource consumption for TID: 1 Byte

### B<sup>+</sup> Tree Cost

$$\text{Average TID list size} = \frac{5000.000}{100} = 50.000$$

$$B^+ \text{Tree Cost} = 100 \times 50.000 \times 1 = 5,000,000 \text{ bytes}$$

### Bitmap Cost

$$\text{Bitmap Cost} = 100 \times 5.000.000 \times \frac{1}{8} = 62,500,000 \text{ bytes}$$

B) 5,000,000 tuples, 3 key values, resource consumption per TID: 1 byte

### B<sup>+</sup> Tree Cost

$$\text{Cost} = 3 \times \frac{5000.000}{3} \times 1 = 5,000,001 \text{ bytes}$$

### Bitmap Cost

$$\text{Cost} = 3 \times 5.000.000 \times \frac{1}{8} = 18,750,000 \text{ bytes}$$

C) 5000,000 tuples, 3 key values, Resource consumption per TID: 4 bytes

### B<sup>+</sup> Tree Cost

$$\text{Cost} = 3 \times \frac{5000.000}{3} \times 4 = 20,000,004 \text{ bytes}$$

### Bitmap Cost

$$\text{Cost} = 3 \times 5.000.000 \times \frac{1}{8} = 18,750,000 \text{ bytes}$$

## Assignment 4

Transform postcode 39106 to Bitmap index structure

Standard Bitmap Index → 100.000 Bitmap Indexes (000000-999999)  
→ one for each value

ID	city	PLZ	Bitmap-39106	Bitmap-39106
0	C1	39106	0	1
1	C2	39106	0	1
2	C3	39104	1	0
3	C4	39104	1	0
4	C5	39106	0	1

→ for exact match: one Bitmap

→ Range query: [n,m]: m-n Bitmaps

Multi-Component-Bitmap-Index       $x = y + z$

each digit is one component

x	y					z					
M	B-5-1	B-4-1	B-3-1	B-2-1	B-1-1	B-0-1	B-4-0	B-3-0	B-2-0	B-1-0	B-0-0
3	0	0	1	0	0	0	0	0	0	0	1
9	1	0	0	0	0	0	1	0	0	0	0
1	0	0	0	0	0	1	0	0	0	1	0
0	0	0	0	0	0	1	0	0	0	0	0
6	0	0	1	0	0	0	0	1	0	0	0

→ 80 Bitmaps    50 - 10 (0-9 values)

→ exact match: 5 Bitmaps

→ Range: 24 Bitmaps

↙  $39\overline{106} - 39\overline{100}$  → two ranges on the digits  
 $1+1+2+10+10$

3 and 9 are always in range,  
391-392 → two digits

better storage consumption but  
complexer than  
standard AMI

How can Bitmap Indexes be used to support range queries for postal codes?

- selection of relevant range for each position
- less storage consumption
- efficient support of range queries

## Range Encoded Bitmap Index

$\leq 1$	$\leq 2$	$\leq 3$
1	1	1
0	1	1
1	1	1
0	0	1
1	1	1

Bitmap of 2 and not  
 $BM - 1$

$BM - 2 \cap BM - 1$

gg.99S Bitmap Indexes

exact match: 2

Range: 2

## Interval Encoded BMI

00000 - 00100

(one Bitmap)

00001 - 00101

00002 - 00102

BMaps: gg.99S

Exact Match: 2

Range: 1  $\rightarrow$  if query window fits range window

Know: How many Bitmap you need for specific query

# Assignment 5

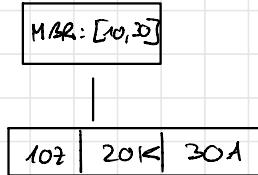
## Construct a R-Tree

- Insert Data in order of given list
- node has at most 3 entries
- outline insertion process
- Ziel: MBR (Minimum Bounding Rectangles) sollte minimal sein, in Größe u.  
Wuschneidung

*exam relevant*

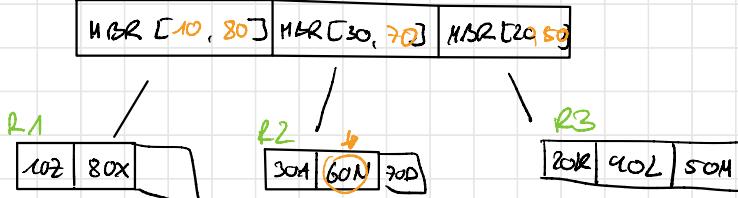
## R-Tree

- 1.) Insert 10z, 20k u. 30t



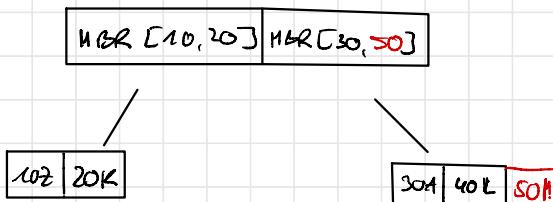
- 2.) Insert 40L exceed space limit → split

R4



→ don't want to have another box that's why we add point to irgendinem node

- 3.) Insert 50M



## R-Tree

### • MBR

↳ encloses indexed points taking minimal space

→ similar to B-Tree

• balanced tree

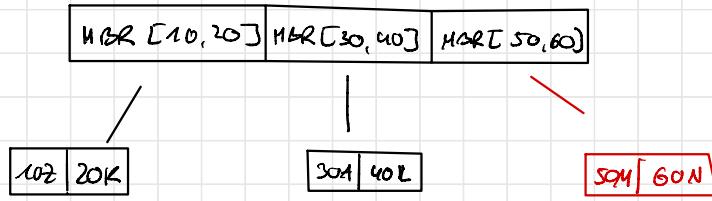
• balancing algorithm

→ rectangles can overlap

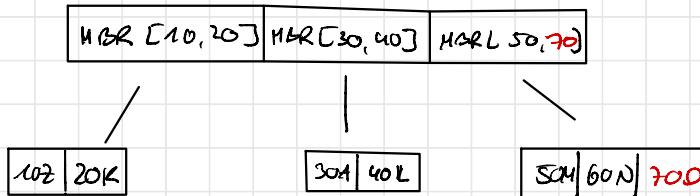
↗ defined by min & max of the MBR diagonal

) sollte minimal sein, in Größe u.

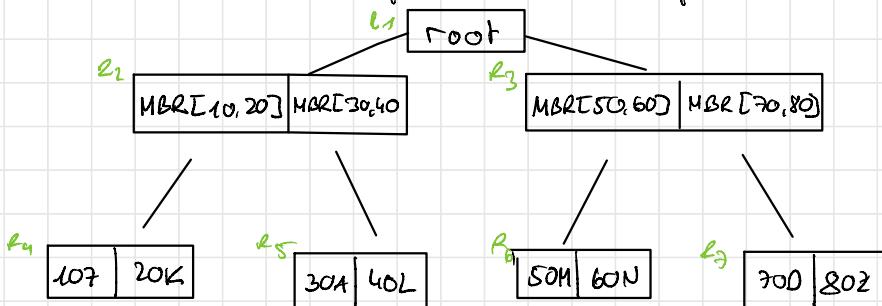
4) Insert 60N exceeds space limit  $\rightarrow$  split

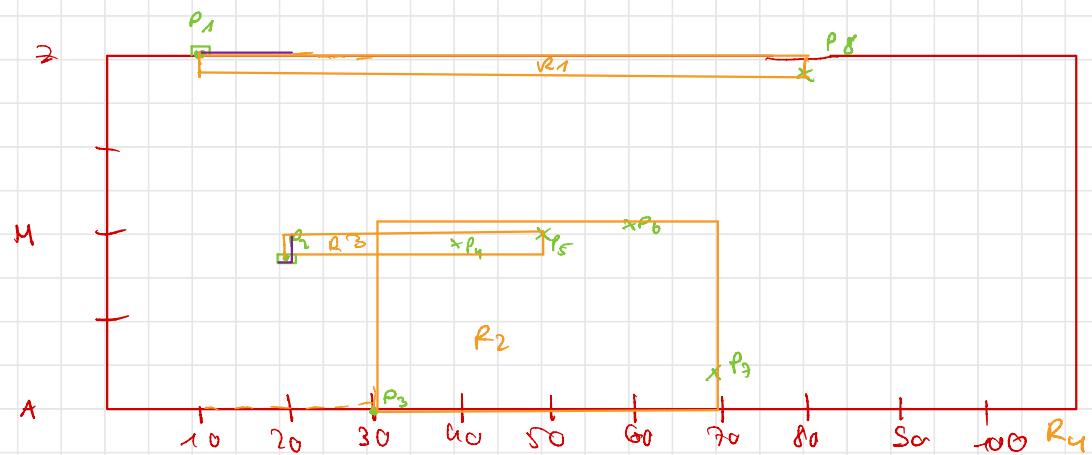


5.) Insert 70D



6) Add 80X exceeds space limit  $\rightarrow$  split





## Assignment 6

### Grid file

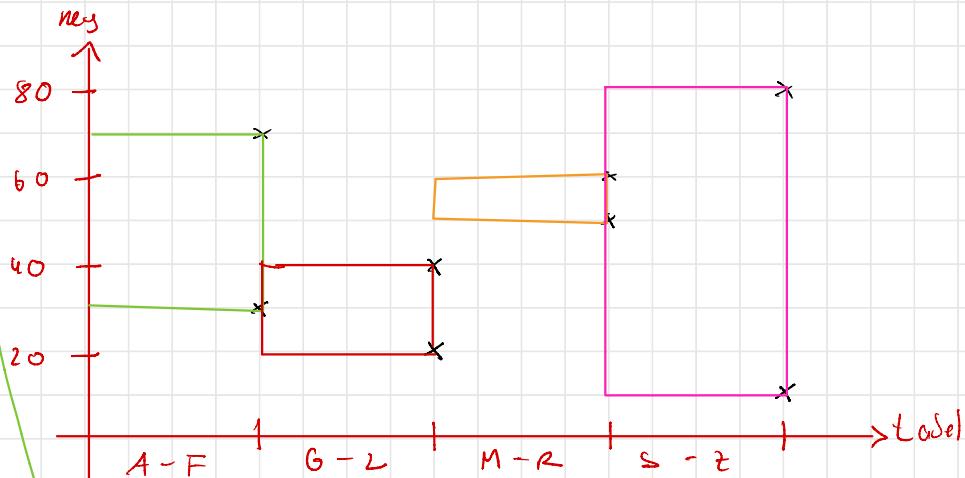
- scales: 1 per dimension
- grid directory: formed by partition of scales
- cells: 1 or more cells form regions
- 1 region = 1 record page

- multidimensional form of data organization  
(combination of key transformation elements and Index File)  
    ↳ Hash                          ↳ Tree

- data space is divided into a grid of cells
- Each cell corresponds to a range of values
- effective for multidimensional data (2D, 3D)

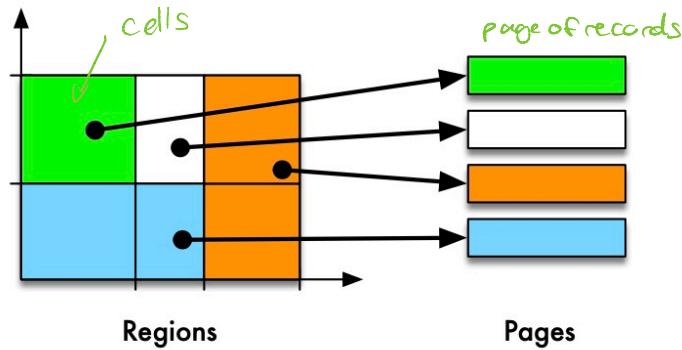
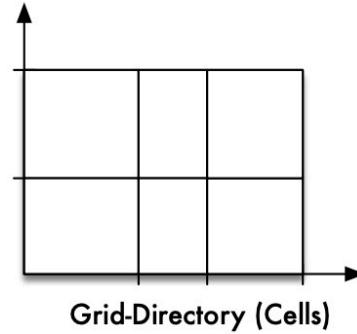
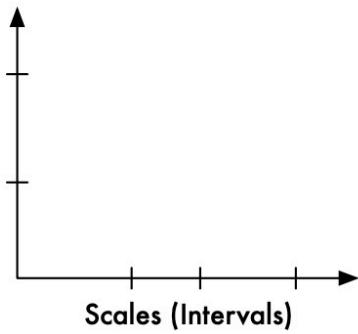
### Creation

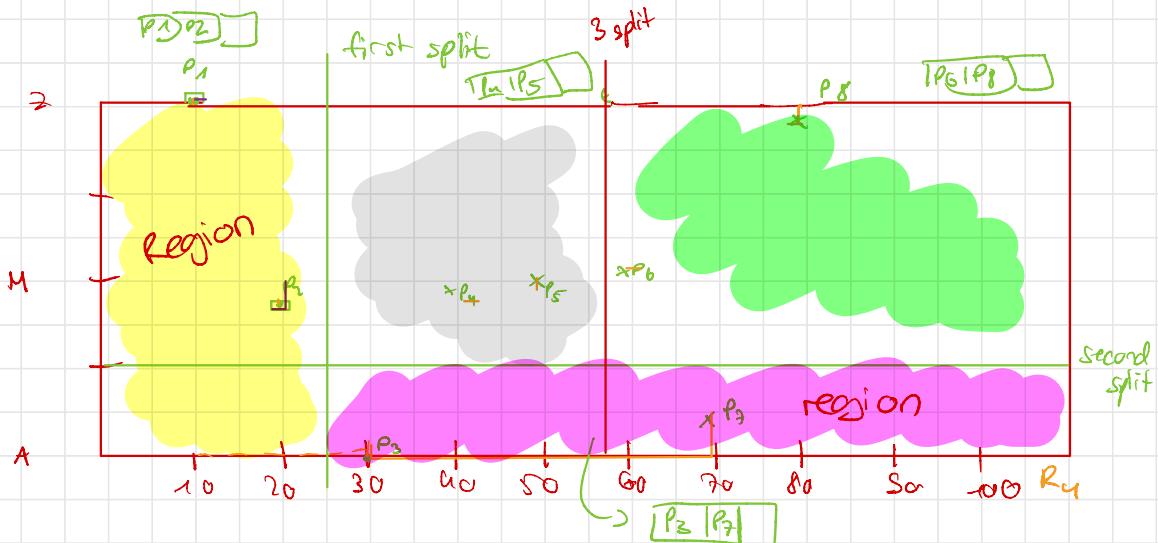
- 1.) Define Dimensions
- 2.) Construct Grid Directory
- 3.) Allocate Buckets



## Advantages

- direct access to data in buckets
- efficient for range queries
- dynamic and flexible
  - ↳ easily add more grid lines or split buckets
- simple and compact (in comparison to a tree structure)





$\Rightarrow$  3 entries per page

$\hookrightarrow$  split: start split on y-axis,  $\rightarrow$  es. equal number of all pages  
 next split on x-axis  
 $\rightarrow$  then again

\* 2. split at y-axis: (globally split between  $P_3, P_4, P_5, P_6$ )

3. split x-axis (again)

$\Rightarrow$  6 cells with  $\alpha$  regions

$\rightarrow$  can dynamically grow a tree

$\rightarrow$  split depends on heuristic (choose myself)  
 (e.g. Median value between overloaded node)