



**Hasso
Plattner
Institut**

IT Systems Engineering | Universität Potsdam

Datenbanksysteme II
Indexstrukturen
(Kapitel 13)

5.5.2008

Felix Naumann

Klausur

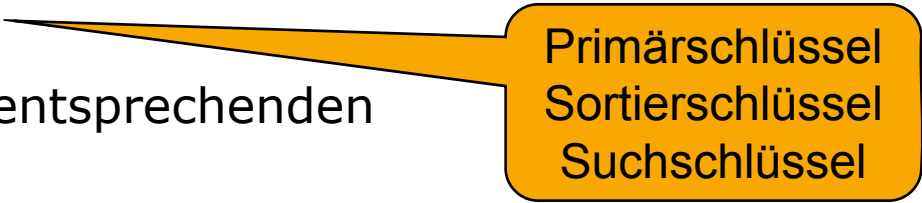
2

- Mittwoch, 23.7.
- 9 – 13 Uhr
 - 4 Stunden
 - Umfang auf 1,5 Stunden ausgelegt
- Keine Hilfsmittel

Motivation

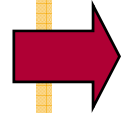
3

- Platzierung der Tupel in Blöcke
 - Naiv: Beliebig verteilen
 - ◇ Aber: `SELECT * FROM R`
 - ◇ Jeden Block untersuchen (Header-Datei)
 - Besser: Tupel einer Relation zusammenhängend speichern
 - ◇ Aber: `SELECT * FROM R WHERE a=10`
 - ◇ Alle Datensätze betrachten
 - Noch besser: Index
 - ◇ Input: Eigenschaften von Datensätzen (z.B: Feldwert)
 - „Suchschlüssel“
 - ◇ Schneller Output: Die entsprechenden Tupel
 - ◇ Nur wenige Datensätze werden betrachtet

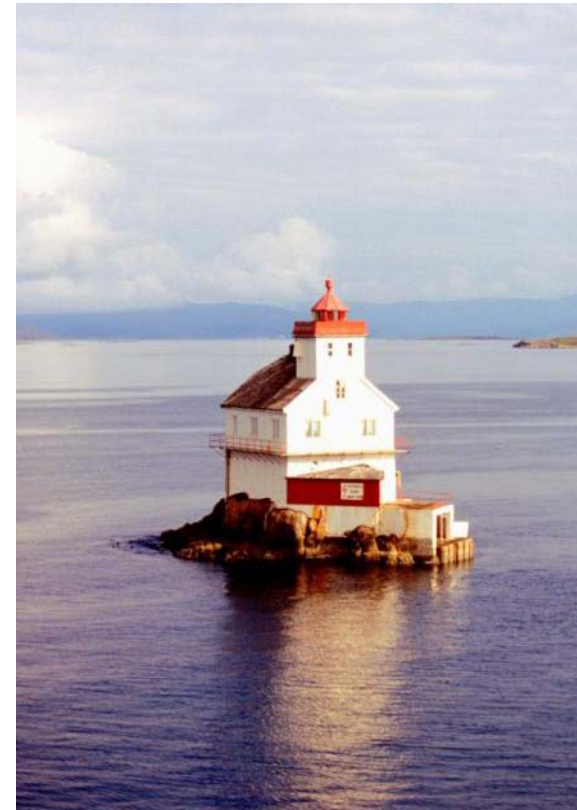


Primärschlüssel
Sortierschlüssel
Suchschlüssel

4



- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
- Hash-Tabellen



Einfachste Form eines Index

5

- Gegeben sortierte Datei (*data file*)
 - Sequenzielle Datei
- Indexdatei enthält Schlüssel-Pointer Paare
 - Schlüsselwert K ist mit einem Pointer verbunden
 - Pointer zeigt auf Datensatz, der den Schlüsselwert K hat.
- Dichtbesetzter Index
 - Ein Eintrag im Index für jeden Datensatz
- Dünnbesetzter Index
 - Nur einige Datensätze sind im Index repräsentiert.
 - Z.B. ein Eintrag pro Block

Sequenzielle Dateien

6

- Index kann sich auf Sortierung des Schlüsselattributs verlassen
 - Hier: Schlüssel ist Suchschlüssel
 - Oft: Suchschlüssel = Primärschlüssel

2 Tupel pro Block

Schlüsselfeld an erster Stelle

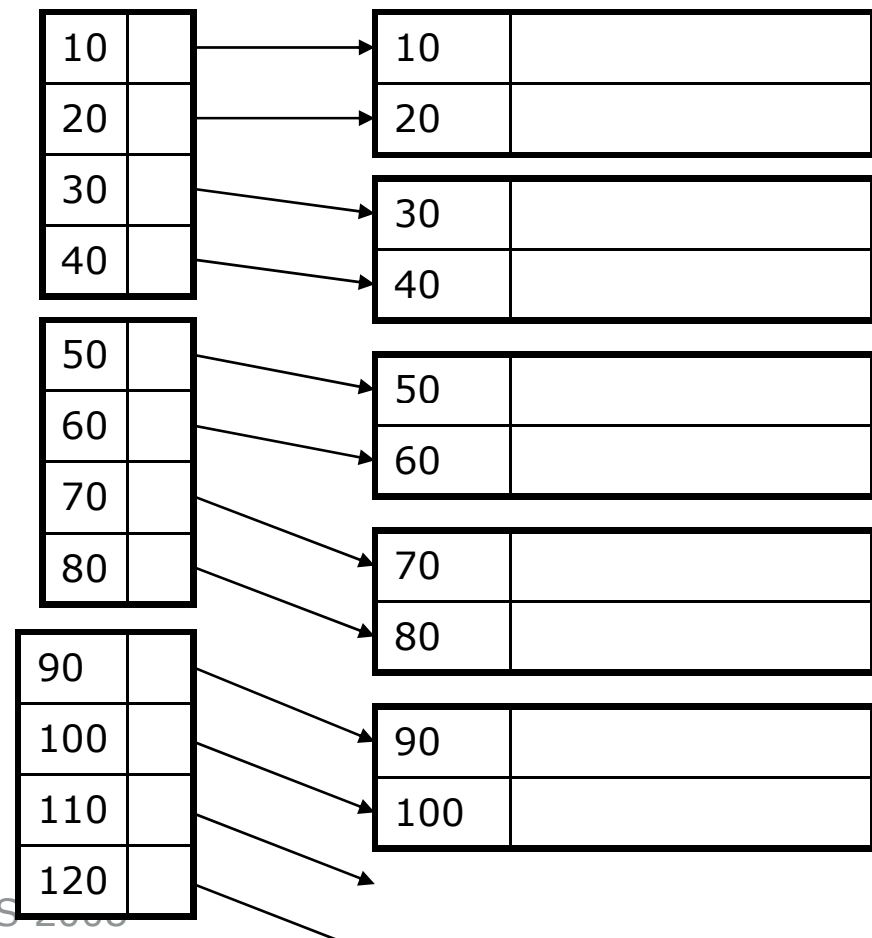
10	
20	
30	
40	
50	
60	
70	
80	
90	
100	

Dichtbesetzte Indizes

7

Indexdatei: Typischerweise Hunderte von Paaren pro Block

- Blocksequenz mit Schlüssel-Pointer Paaren
- Jeder Schlüssel der Daten ist durch ein Paar repräsentiert
 - Aber: Wesentlich kleinere Datenmenge
 - Passt womöglich in den Hauptspeicher
 - Nur ein I/O pro Zugriff
- Sortierung der Paare = Sortierung der Daten

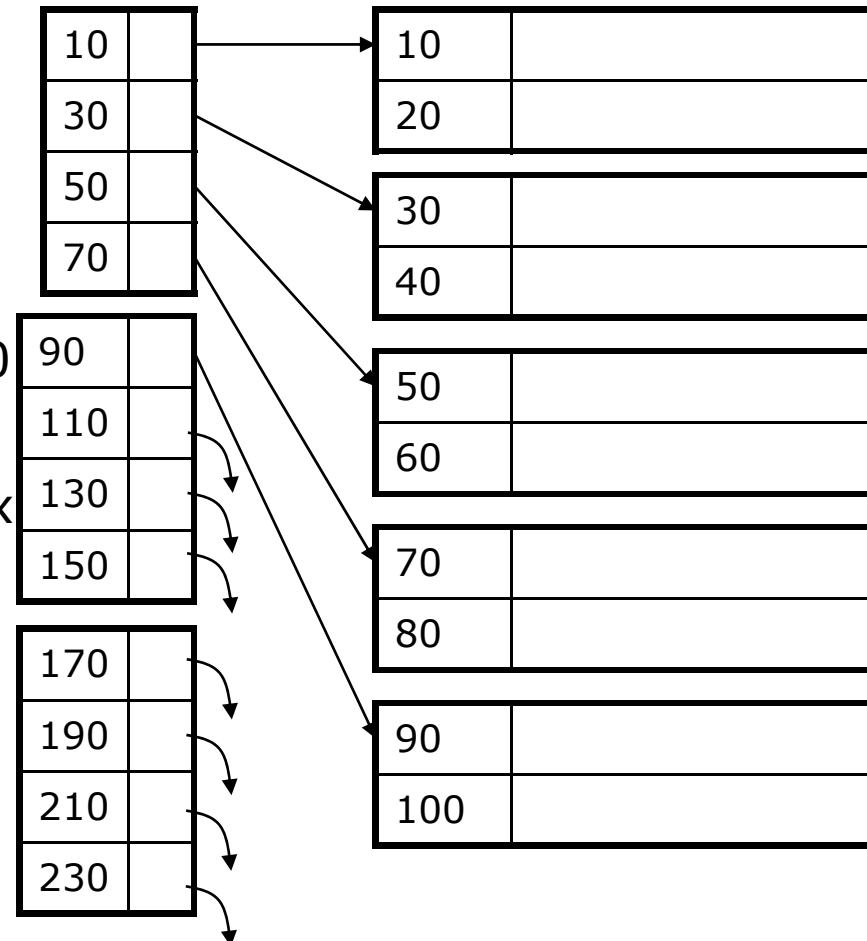


- Gegeben Suchschlüssel K
- Durchsuche Indexdatei nach K
- Folge Pointer
- Lade Block aus Datendatei
- Beschleunigung
 - Indexdatei hat nur wenige Blocks
 - ◇ Indexdatei im Hauptspeicher
 - Binäre Suche um K zu finden
- Beispiel: 1.000.000 Tupel
 - Block: 4096 Byte = 10 Tupel
 - Gesamtgröße 400 MB
 - Schlüsselfeld hat 30 Byte
 - Pointer hat 8 Byte
 - => 100 Paare pro Block
 - Dichtbesetzter Index: 10.000 Blöcke für Index
 - ◇ 40 MB => vielleicht OK im Hauptspeicher
 - Binäre Suche: $\log_2(10.000) \approx 13$
 - ◇ => 13-14 Blocks pro Suche
 - ◇ Wichtigsten Blöcke im Hauptspeicher reichen ($1/2, 1/4, 3/4, \dots$)

Dünnbesetzte Indizes

9

- Weniger Speicherbedarf
- Aber höherer Suchaufwand
- Nur ein Pointer pro Block
- Beispiel
 - 100.000 Datenblöcke, 100 Indexpaare pro Block
 - => 1.000 Blocks für Index = 4MB



Anfragebearbeitung mit dünnbesetzten Indizes

10

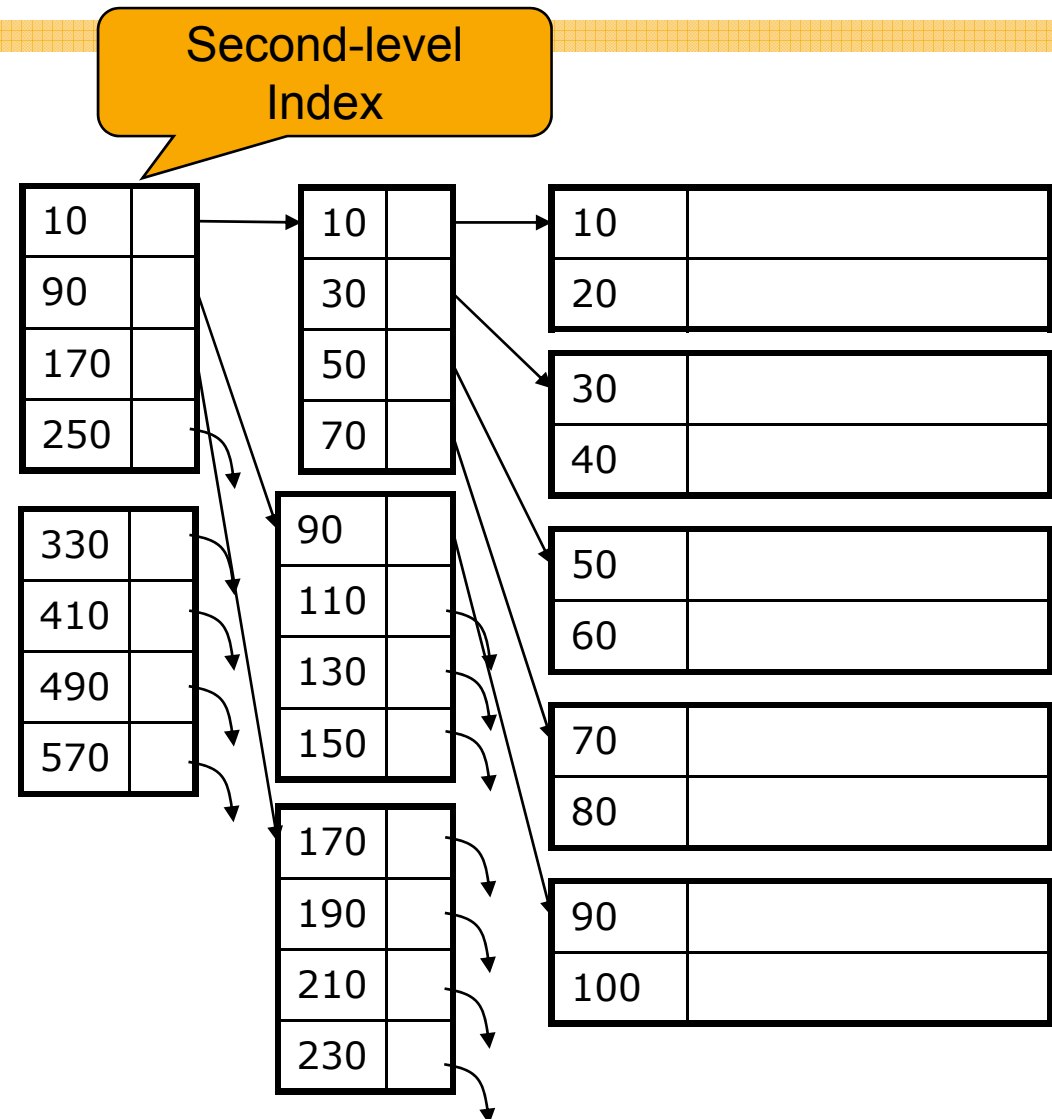
1. Suche im Index größten Schlüssel, der kleiner/gleich als Suchschlüssel ist.
 - Binäre Suche (leicht modifiziert)
 2. Hole assoziierten Datenblock
 3. Durchsuche Block nach Datensatz
-
- **Nachteil?**
 - `SELECT `TRUE` FROM R WHERE a=10`
 - Kann nicht ausschließlich mit Index beantwortet werden
 - ◇ Im dicht-besetzten Index schon...

Mehrstufiger Index

11

Auch ein Index kann unangenehm groß sein.

- Nimmt viel Speicher ein
- Kostet viel I/O
 - Auch bei binärer Suche
- Idee: Index auf den Index
 - Zweiter Index macht nur als dünn-besetzter Index Sinn.
- Theoretisch auch dritte, vierte, ... Ebene
 - B-Baum aber besser



Mehrstufiger Index – Beispiel

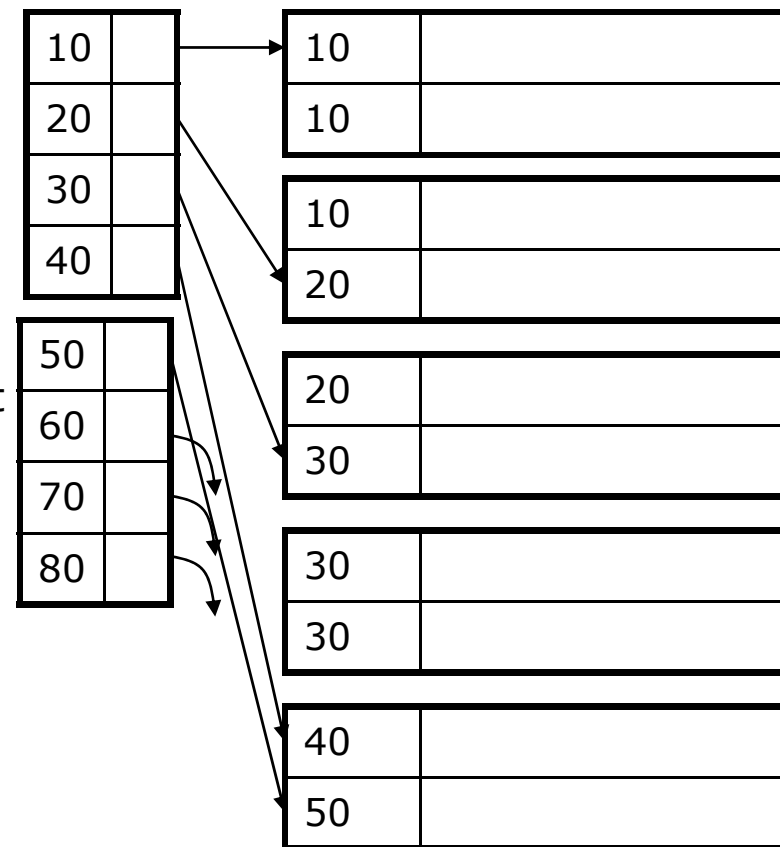
12

- 100.000 Datenblöcke, 100 Indexpaare pro Block
- => 1.000 Blocks für Index erster Stufe = 4MB
- => 10 Blocks für Index zweiter Stufe = 40KB
- Kann mit Sicherheit im Hauptspeicher verbleiben
- Vorgehen
 1. Suche im Index zweiter Stufe größten Schlüssel, der kleiner/gleich als Suchschlüssel ist.
 2. Hole entsprechenden Block im Index erster Stufe.
 - ◇ Eventuell schon im Hauptspeicher
 3. Suche in dem Block größten Schlüssel, der kleiner/gleich als Suchschlüssel ist.
 4. Hole entsprechenden Datenblock.
 5. Suche Datensatz (falls Index erster Stufe dünnbesetzt ist).
- Zusammen: 2 I/Os

Indizes für Nicht-eindeutige Suchschlüssel

13

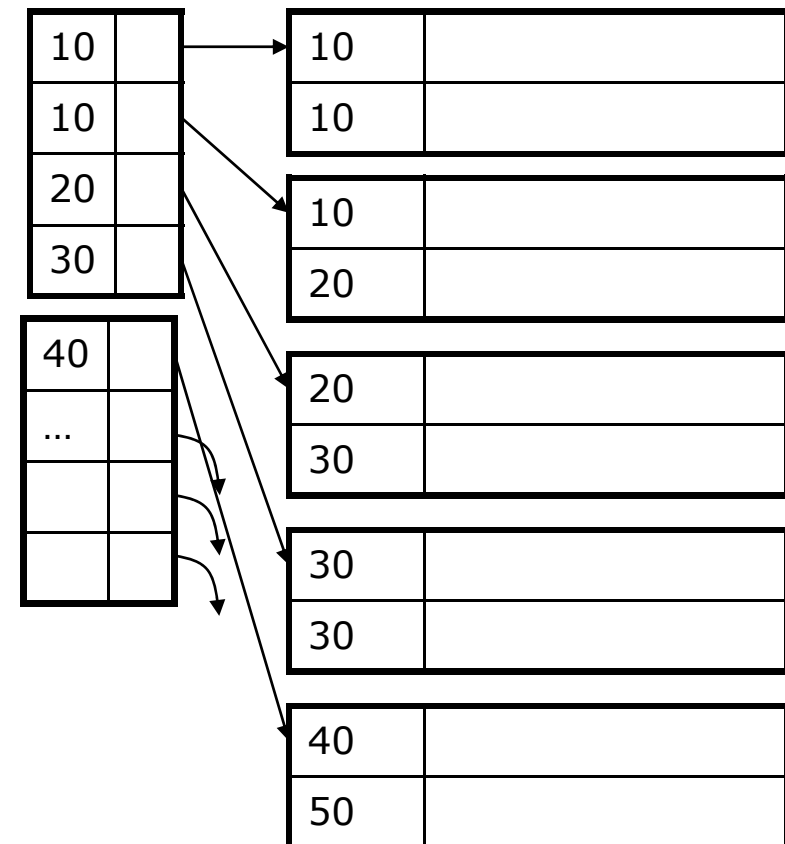
- Annahme bisher: Suchschlüssel ist auch Schlüssel bzw. eindeutig in der Relation
- Annahme weiter: Relation ist nach Suchschlüssel sortiert
- Idee 1: Dichtbesetzter Index: ein Paar im Index für jeden Datensatz
 - Anfragebearbeitung:
 - ◇ Suche erstes Paar mit K.
 - ◇ Wähle alle weiteren mit K (direkt dahinter)
 - ◇ Hole entsprechende Datensätze.
- Idee 2: Nur ein Indexpaar pro Schlüsselwert K. Der zeigt auf ersten Datensatz mit K.
 - Weitere Datensätze mit K folgen direkt.
 - Wichtig: Blöcke haben Pointer auf jeweils nächsten Block



Indizes für Nicht-eindeutige Suchschlüssel

14

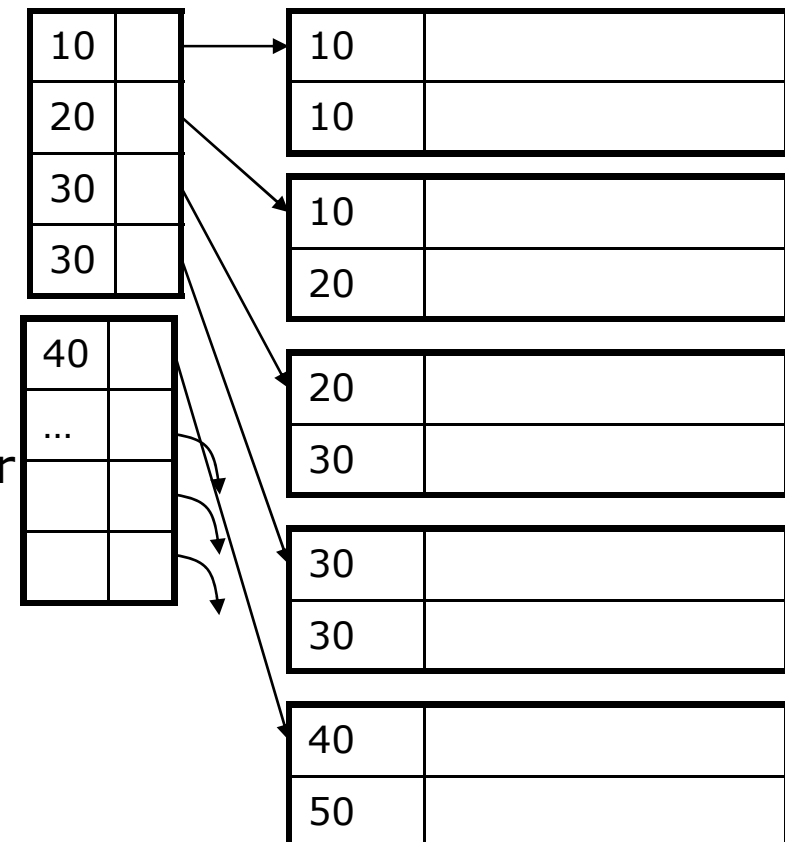
- Idee 3: Dünnbesetzter Index wie gehabt
 - Pointer jeweils auf Datensatz am Blockanfang
- Anfragebearbeitung
 - Suche letzten Eintrag E1 im Index, dessen Datenwert $\leq K$
 - Suche von dort im Index nach vorn bis zu einem Eintrag E2 mit Datenwert $< K$
 - Hole alle Datenblöcke zwischen und inklusive E1 und E2.
- Beispiel: $K = 20$



Indizes für Nicht-eindeutige Suchschlüssel

15

- Idee 4: Dünnbesetzter Index; aber
 - Datenwert im Index ist der kleinste neue Wert im entsprechenden Datenblock.
 - Falls kein neuer Wert im Block, dann den existierenden Wert.
- Anfragebearbeitung einfacher
 - Suche im Index nach Paar mit (Datenwert = K) oder ($< K$ aber nächste Wert ist $> K$).
 - Hole Datenblock und gegebenenfalls folgende Datenblöcke.



Änderungsoperationen

16

Daten ändern sich (Insert, Update, Delete).

- Annahme bisher: Daten füllen Blöcke perfekt und ändern sich nicht
- Änderungen im Datenblock: Siehe voriger Foliensatz
 - Overflow Blocks
 - ◇ In dünnbesetzten Indizes nicht repräsentiert
 - Neue Blöcke in der Sequenz
 - ◇ Benötigen neuen Indexeintrag
 - ◇ Indexänderungen bergen dieselben Probleme wie Datenänderungen.
 - Platzierung der Blocks
 - In Indizes höherer Stufe
 - Tupel verschieben
 - ◇ Index muss angepasst werden.
- Generelle Regel: Indizes können wie normale data files behandelt werden. Gleiche Strategien können angewendet werden.

Änderungsoperationen

17

Index für Datensätze

Index für Blöcke

Aktion	Dichtbesetzter Index	Dünnbesetzter Index
Erzeugung eines leeren Overflow Blocks		
Löschen eines leeren Overflow Blocks		
Erzeugen eines leeren sequenziellen Blocks		
Löschen eines leeren sequenziellen Blocks		
Datensatz einfügen		
Datensatz löschen		
Datensatz verschieben		

Änderungsoperationen

18

Index für Datensätze

Index für Blöcke

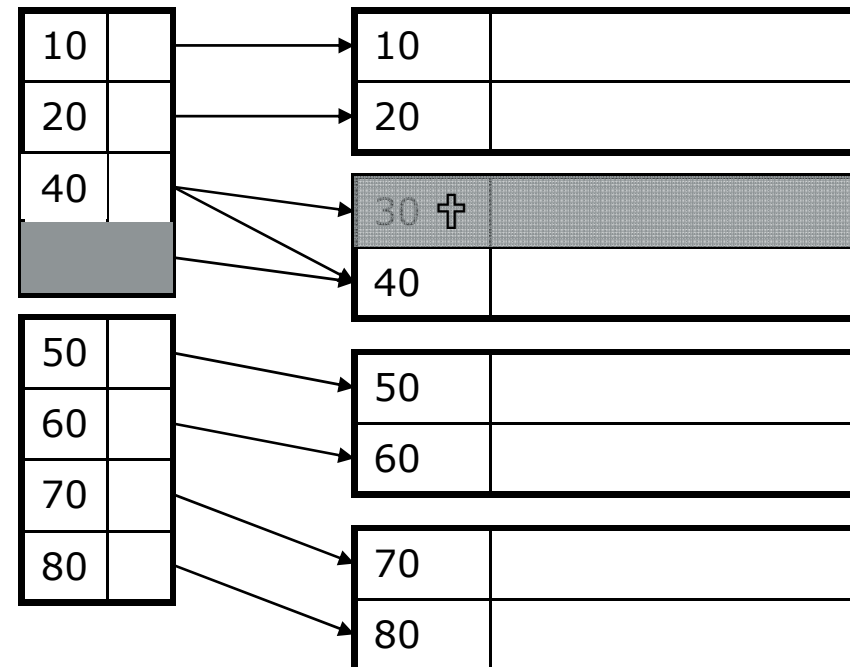
Aktion	Dichtbesetzter Index	Dünnbesetzter Index
Erzeugung eines leeren Overflow Blocks	%	%
Löschen eines leeren Overflow Blocks	%	%
Erzeugen eines leeren sequenziellen Blocks	%	Insert
Löschen eines leeren sequenziellen Blocks	%	Delete
Datensatz einfügen	Insert	Update?
Datensatz löschen	Delete	Update?
Datensatz verschieben	Update	Update?

Nur falls Datensatz erster im Block ist

Änderungsoperationen – Beispiele

19

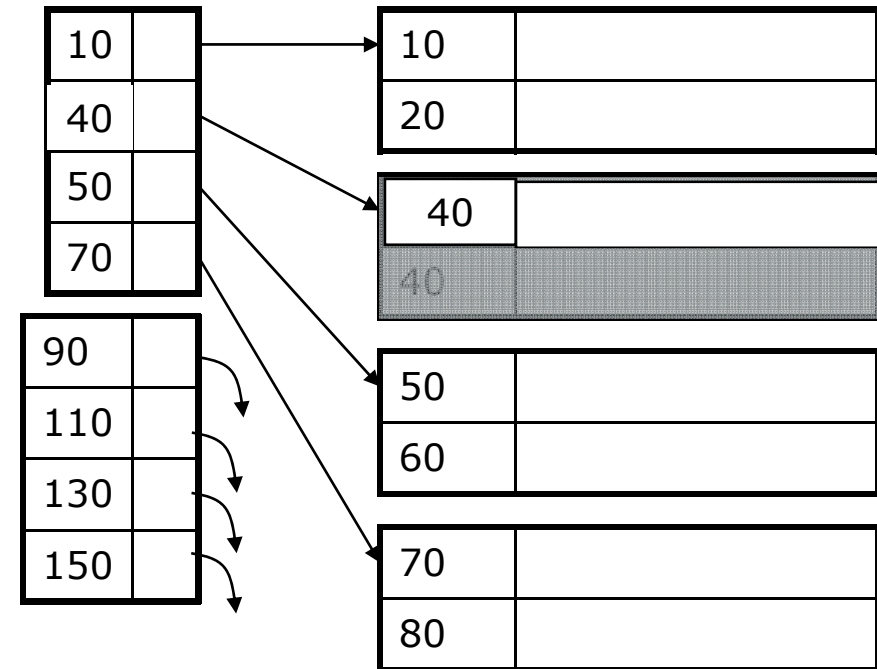
- Datensatz mit $K = 30$ wird gelöscht.
- Annahme: Block kann/soll nicht reorganisiert werden.
 - Ersatz durch *tombstone*
- Datensatz 40 wird nicht verschoben.
- Index kann reorganisiert werden.
 - Main memory



Änderungsoperationen – Beispiele

20

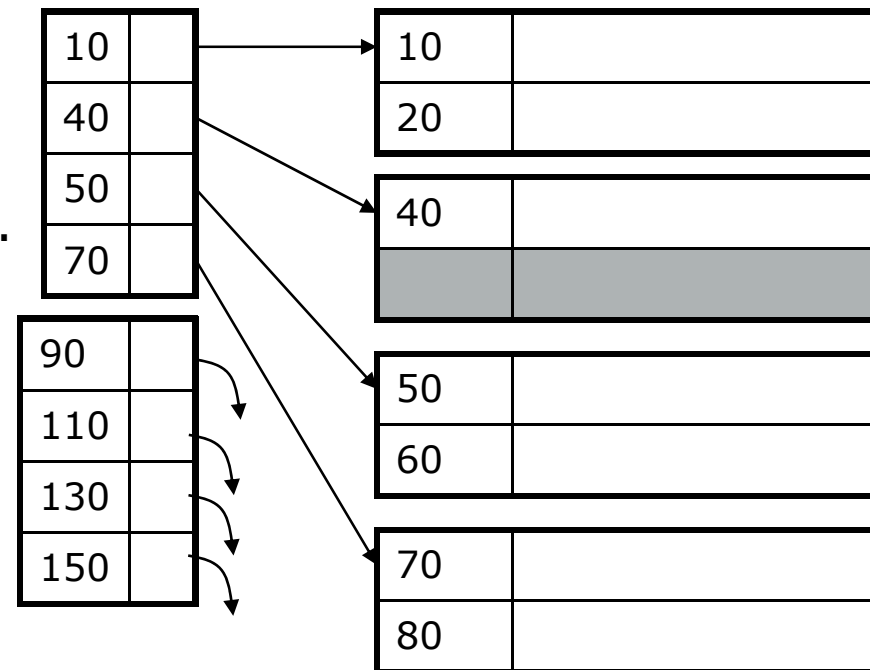
- Datensatz mit $K = 30$ wird gelöscht.
 - Annahme: Block kann reorganisiert werden.
 - Datensatz 40 wird verschoben
 - Index wird aktualisiert
- Nun auch Datensatz mit $K=40$ Löschen
 - Leerer Block entsteht
 - Index wird aktualisiert (löschen)
 - Index wird reorganisiert



Änderungsoperationen – Beispiele

21

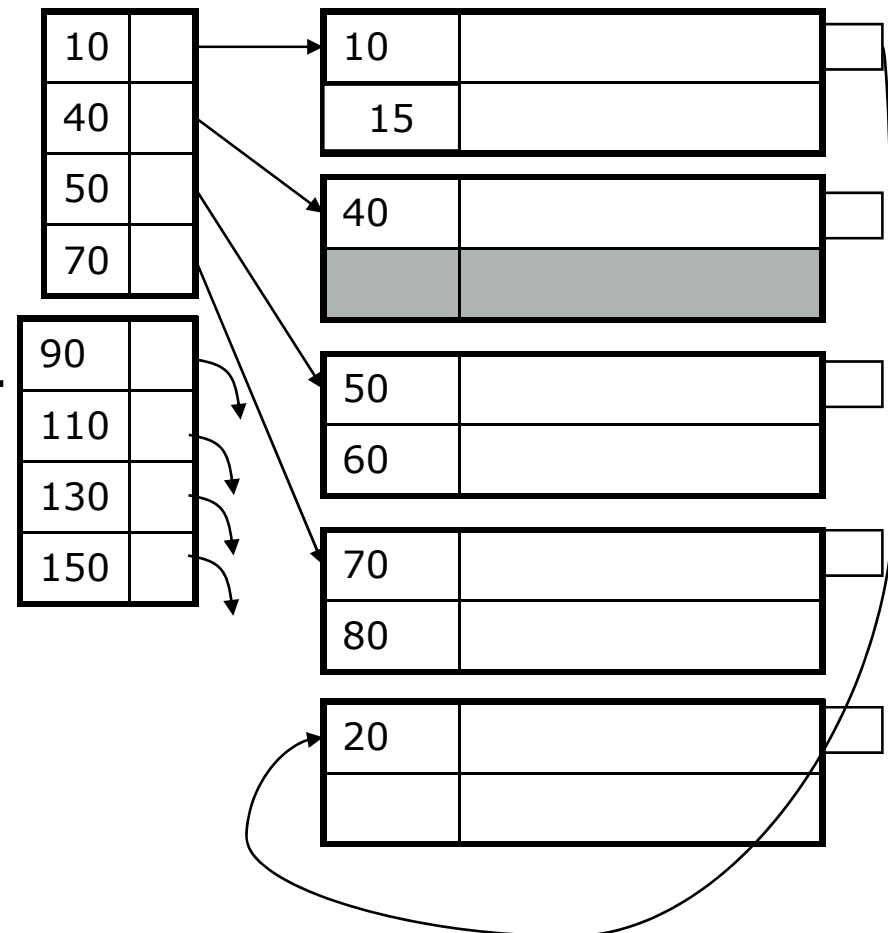
- Einfügen eines Datensatzes 15
 - Block 1 ist voll.
 - Datensatz 20 wird in nächsten Block verschoben.
 - ◇ Block wird reorganisiert.
 - Datensatz 15 wird eingefügt.
 - Index wird aktualisiert.
 - ◇ 20 statt 40



Änderungsoperationen – Beispiele

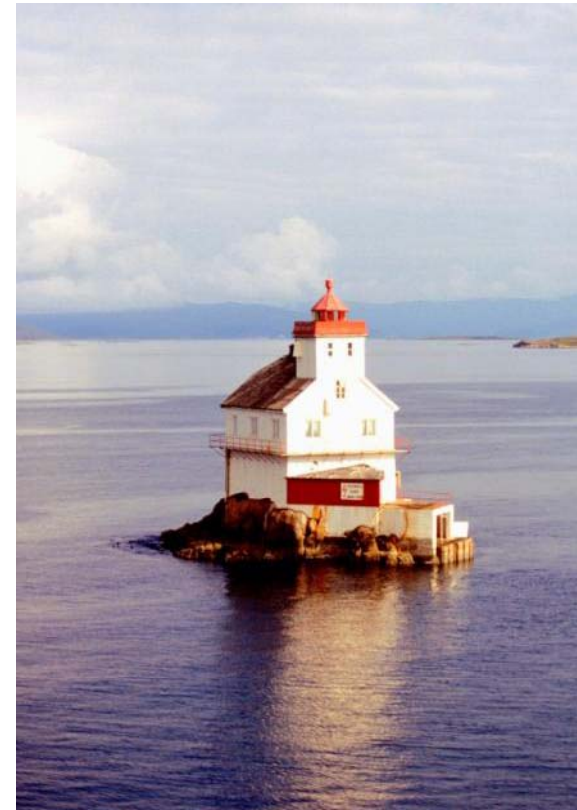
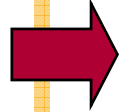
22

- Wieder Datensatz 15 einfügen
 - Diesmal mit Overflow Blocks
 - Block 1 ist voll.
 - Datensatz 20 wird in Overflow Block verschoben.
 - Datensatz 15 wird eingefügt.
 - Index bleibt gleich



23

- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
- Hash-Tabellen



Motivation

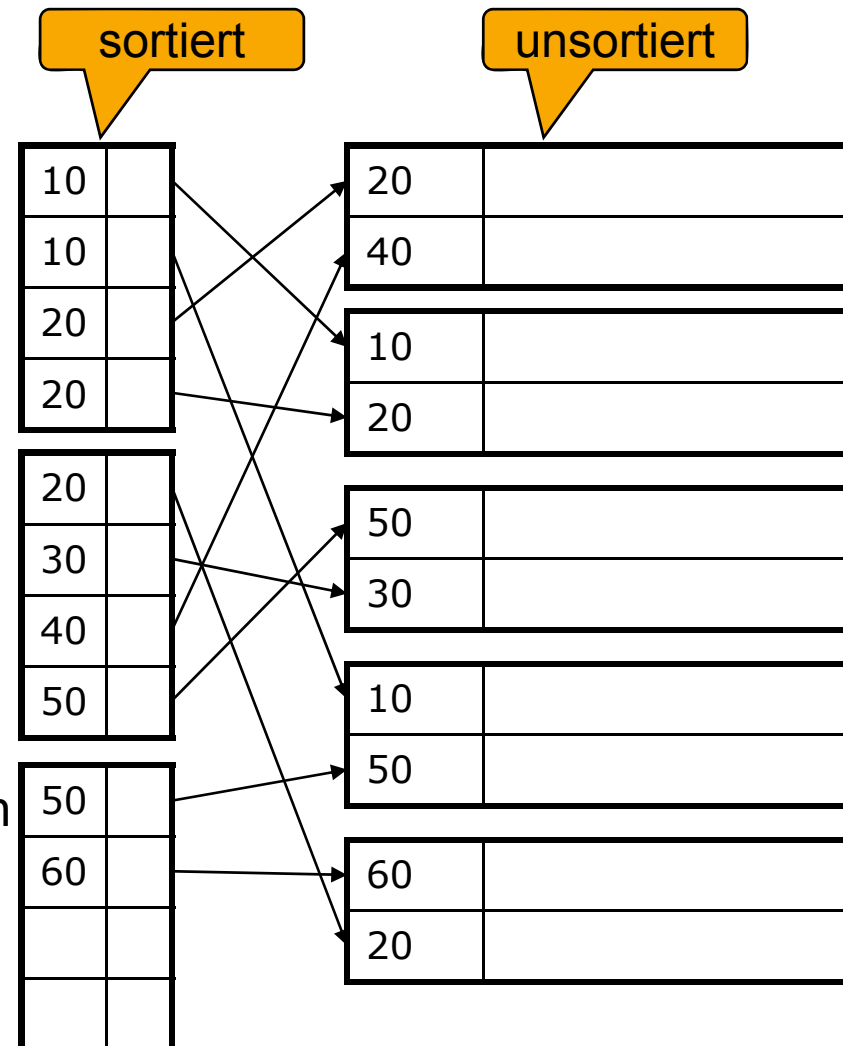
24

- Annahme bisher: Datensätze sind nach Schlüssel sortiert
 - „Primärindex“
- Oft sinnvoll: Mehrere Indizes pro Relation
- `Schauspieler(Name, Adresse, Geschlecht, Geburtstag)`
 - `Name` ist Primärschlüssel => Primärindex
 - `SELECT Name, Adresse`
`FROM Schauspieler`
`WHERE Geburtstag = DATE `1952-01-01``
- Sekundärindex auf `Geburtstag` beschleunigt Anfragebearbeitung.
 - `CREATE INDEX GEB_IDX ON Schauspieler(Geburtstag)`
- Sekundärindizes bestimmen nicht Platzierung der Datensätze, sondern geben Speicherort an.
 - Dünnbesetzte Sekundärindizes sind sinnlos.
 - => Sekundärindizes sind immer dichtbesetzt.

Aufbau von Sekundärindizes

25

- Dichtbesetzt; mit Duplikaten
- Schlüssel-Pointer Paare
- Schlüssel sind sortiert
- Index zweiter Stufe wäre wiederum dünnbesetzt
- Suche kostet idR mehr I/O
 - Beispiel: Suche nach „20“ muss 5 Blöcke lesen.
 - Ist nicht zu ändern: Daten sind halt nach einem anderen Schlüssel sortiert.



Anwendungen

26

- Unterstützung von Selektionsbedingungen auf Nicht-Primärschlüssel
- Datensätze liegen nicht sortiert vor.
 - Sekundärindex auf Primärschlüssel
- Datensätze aus zwei Relationen werden geclustered gespeichert.
 - N:1 Beziehung zwischen R und S
 - Speichere Datensätze aus R direkt beim entsprechenden Datensatz aus S.
 - *Clustered file*

Anwendung: Clustered file

27

- `Filme(Titel, Jahr, Länge, inFarbe, Studioname, Produzent)`
- `Studio(Name, Adresse, Präsident)`
- Häufige Anfrageform:
 - `SELECT Titel, Jahr`
`FROM Filme, Studio`
`WHERE Filme.Studioname = Studio.Name`
`AND Präsident = ?`

Studio1	Filme aus Studio1	Studio2	Filme aus Studio2	Studio3	Filme aus Studio3	...
---------	-------------------	---------	-------------------	---------	-------------------	-----

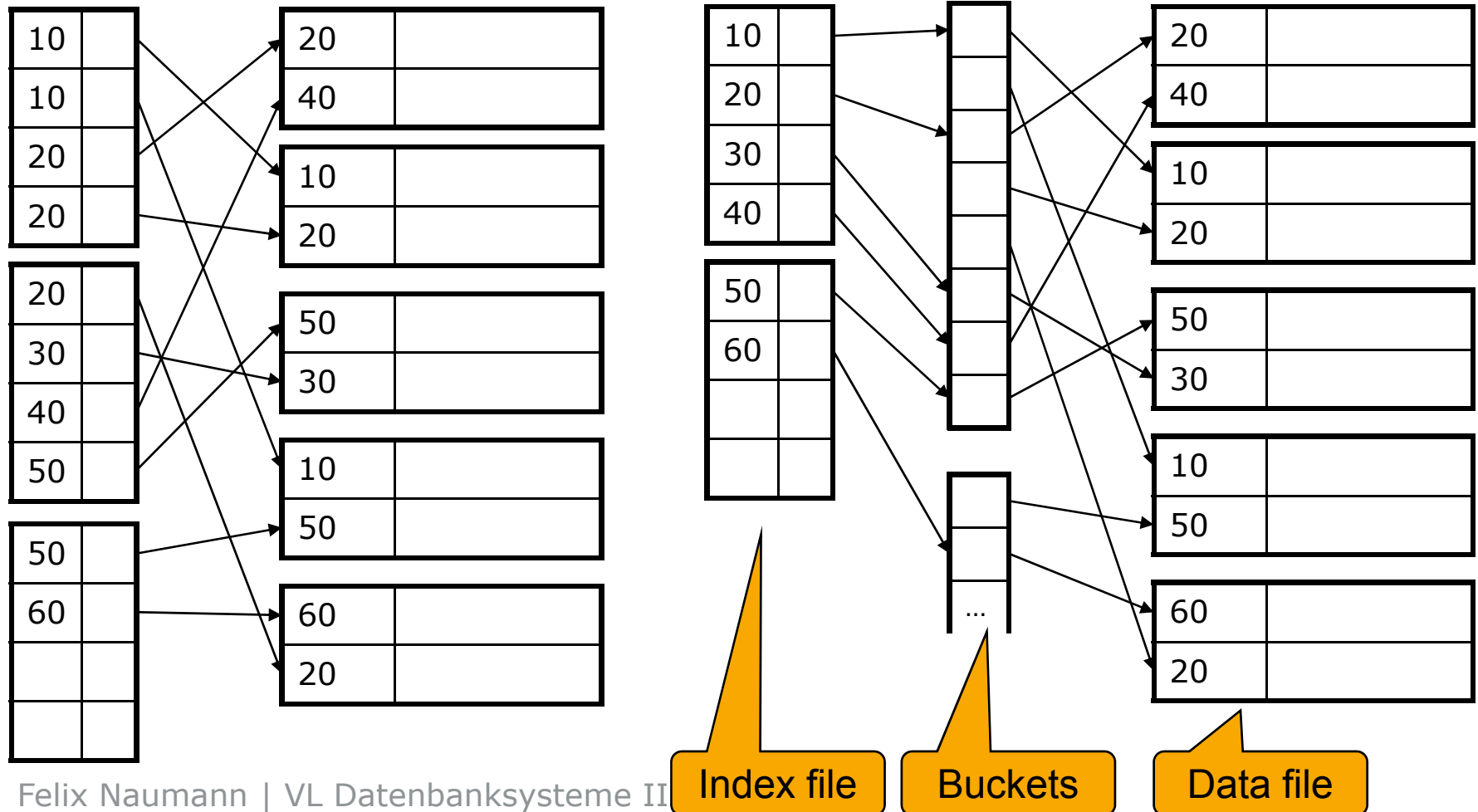
1. Index auf `Präsident` findet schnell Studio-Datensatz.
2. Entsprechende Filme folgen direkt.

Indirektion für Sekundärindizes

Können groß sein: Nicht Primärschlüssel

28

Bisherige Struktur verbraucht Platz: Datenwerte werden mehrfach gespeichert.



Indirektion für Sekundärindizes

29

- Spart Platz, falls
 - Suchschlüssel größer sind als Pointer.
 - Suchschlüssel im Durchschnitt mindestens zweimal auftauchen.
- Weiterer Vorteil: Bestimmte Anfragen können direkt anhand der Buckets beantwortet werden.
 - Mehrere Selektionsbedingungen, jeweils mit Sekundärindex: Schnittmenge der Pointer in Buckets
 - **Filme(Titel, Jahr, Länge, inFarbe, Studioname, Produzent)**
 - **SELECT Titel FROM Filme
WHERE StudioName = `Disney`
AND Jahr = 1995**

Invertierte Indizes

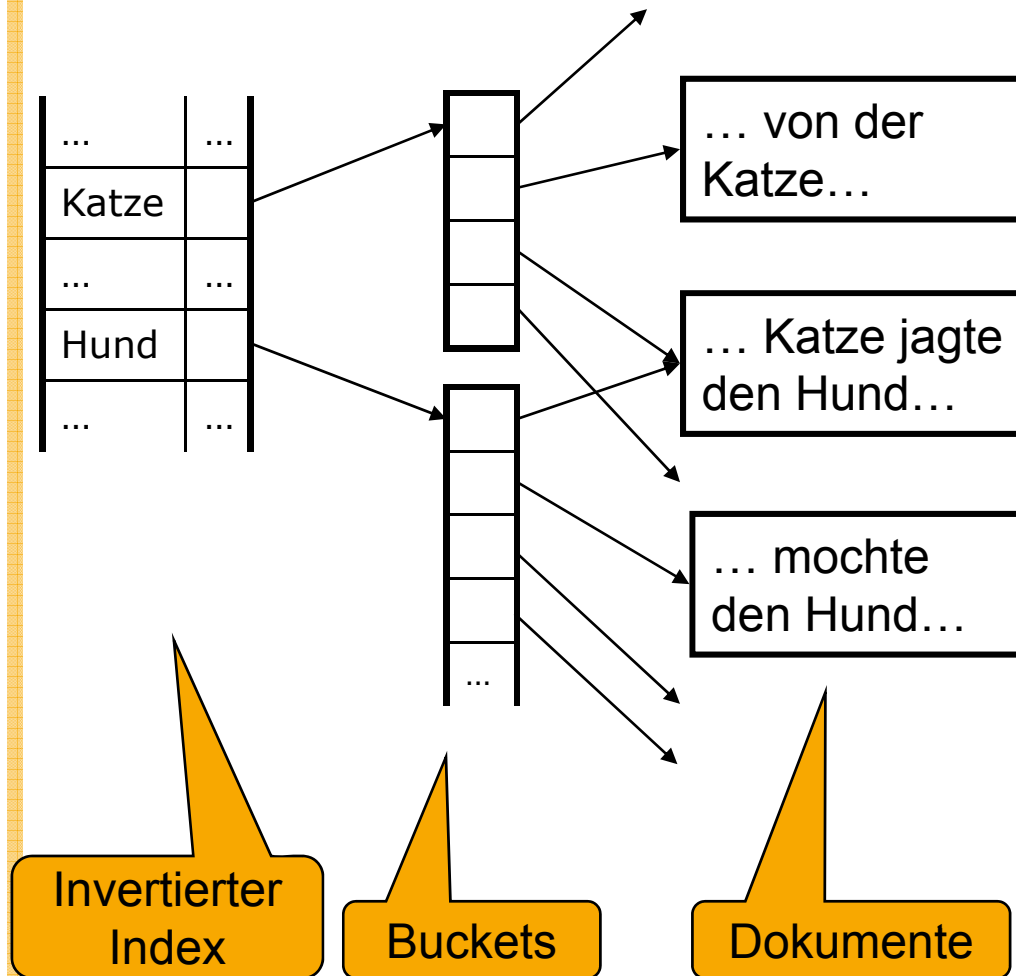
30

Motivation: Dokumenten-Retrieval

- Sichtweise Dokument als Relation
 - `Dok(hatKatze, hatHund, hatHaus, ...)`
 - Tausende Boolesche Attribute: `True` bedeutet das Dokument enthält das Wort.
 - Sekundärindex auf jedes Attribut
 - ◇ Aber: Nur die True-Werte werden indiziert
 - Alle Indizes in einen kombiniert: Invertierte Liste
 - ◇ Verwendet Indirektion

Invertierte Indizes

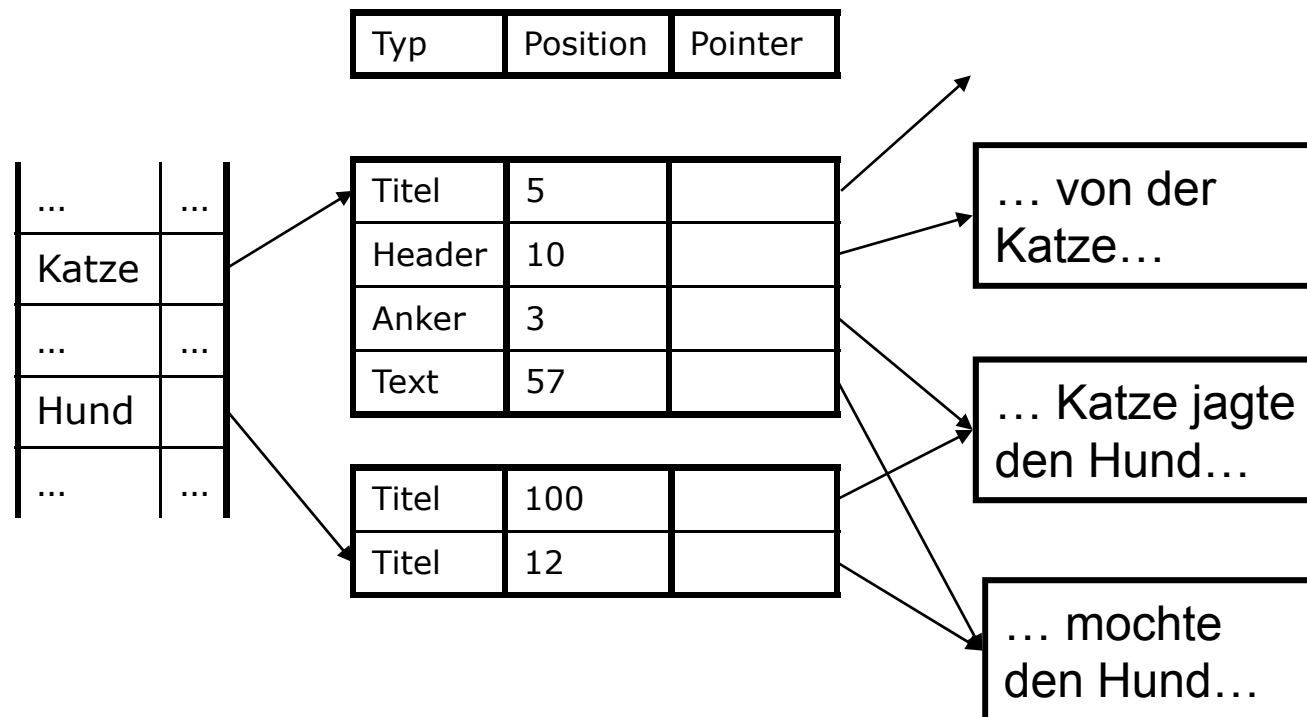
31



- Pointer in den Buckets
 - Auf ein Dokument
 - Auf eine Stelle im Dokument
- Erweiterung: Bucket speichert nicht nur Stelle sondern auch Metadaten
 - Art des Vorkommens (Titel, Abstract, Text, Tabelle, ...)
 - Satz (fett, kursiv, ...)
 - ...
- Anfragen: AND, OR, NOT
 - Durch Schneiden der Pointermengen

Invertierte Indizes

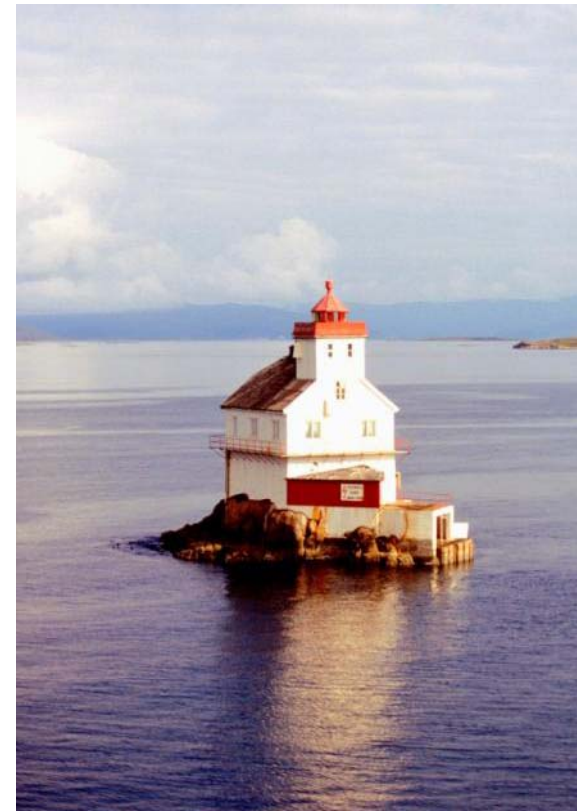
32



Suche nach Dokumenten, die Hunde und Katzen vergleichen

- Dokument erwähnt „Hund“ im Titel
- Dokument erwähnt „Katze“ in einem Anker (Link auf anderes Dokument)

- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- ➔ ■ B-Bäume
 - Aufbau
 - Suche
 - Updates
 - Effizienz
 - Varianten
- Hash-Tabellen



B-Bäume

34

- Bisher: Zweistufiger Index zur Beschleunigung des Zugriffs
- Allgemein: B Bäume (hier B+ Bäume)
 - So viele Stufen wie nötig
 - Blöcke sind mindestens zur Hälfte gefüllt
 - Overflow blocks nicht notwendig

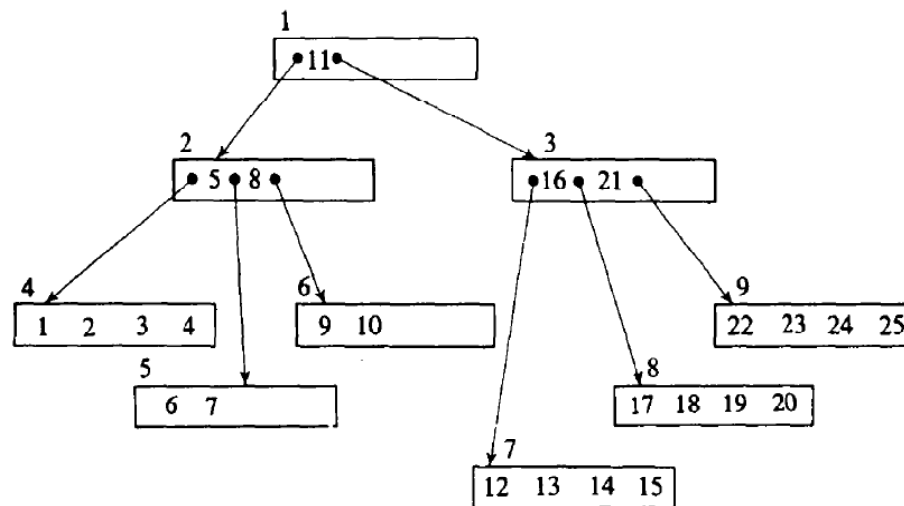
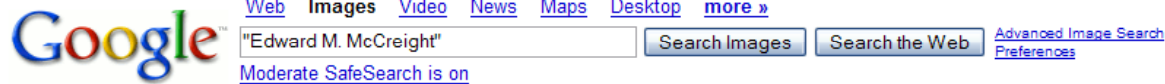


Fig. 2. A data structure in $\tau(2, 3)$ for an index

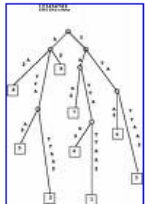


B-Bäume

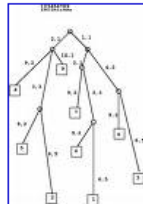
Rudolf Bayer und Edward McCreight



Images Showing: All image sizes Results 1 - 18 of about 55 for "Edward M. McCreight". (0.03 sec)



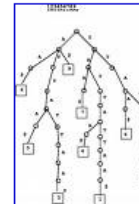
Trie versionis on ruumi raisatud.
459 x 689 - 12k - png
www.egeen.ee



That is, it spells out S[i..m].
461 x 689 - 11k - png
www.egeen.ee

M

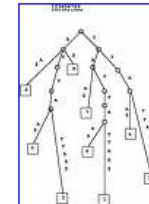
... the block size, ...
32 x 20 - 1k - gif
publications.csail.mit.edu



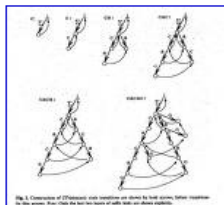
Vaatame näidet:
459 x 689 - 13k - png
www.egeen.ee



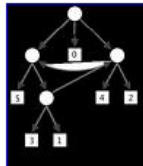
Sama näide sufiksipuul peal:
752 x 1050 - 24k - png
www.egeen.ee



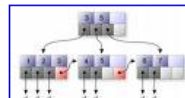
Alternatiivne esitusviis kus lehed
459 x 689 - 12k - png
www.egeen.ee



Kui vaja, jätkka seda teed tähega ...
1050 x 987 - 21k - png
www.egeen.ee



Suffix tree for the string BANANA ...
250 x 303 - 17k - png
www.answers.com



From Wikipedia, the free ...
400 x 220 - 28k - png
en.wikipedia.org

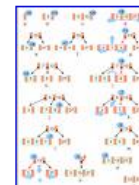
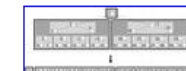


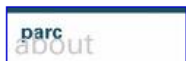
Abbildung 7: Evolution eines B- ...
320 x 448 - 57k - png
bayer-baum.meine-suchabfrage.de



... recht groben – Überblick über ...
400 x 736 - 81k
www.computerbase.de



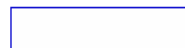
vEB layout
327 x 136 - 11k - jpg
publications.csail.mit.edu



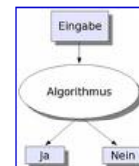
PARC Historical Publications
373 x 120 - 3k - gif
www.parc.xerox.com



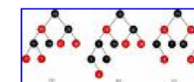
[Bearbeiten] Einordnung in die ...
480 x 208 - 57k - png
www.computerbase.de



PARC Historical Publications
377 x 89 - 1k - gif
www.parc.xerox.com



... Arbeitsweise von Automaten.
200 x 241 - 16k - png
www.computerbase.de



例: 赤黒木
480 x 194 - 22k - png
tiki.is.os-omicron.org



Logo
170 x 75 - 3k - gif
db4707.inf.tu-dresden.de

- Index-Blöcke in einem Baum organisiert
- Balanciert
 - Jeder Weg von Wurzel zu Blatt ist gleich lang.
- Parameter n
 - Jeder Block enthält bis zu n Suchschlüssel
 - Jeder Block enthält bis zu $n+1$ Pointer
 - Also wie Indexblock zuvor, aber ein zusätzlicher Pointer
- Wahl von n
 - n so groß wie möglich entsprechend der Blockgröße
 - 4096 Byte pro Block; 4 Byte pro Schlüssel; 8 Byte pro Pointer
 - $4n + 8(n+1) \leq 4096 \Rightarrow n = 340$

Regeln im B-Baum

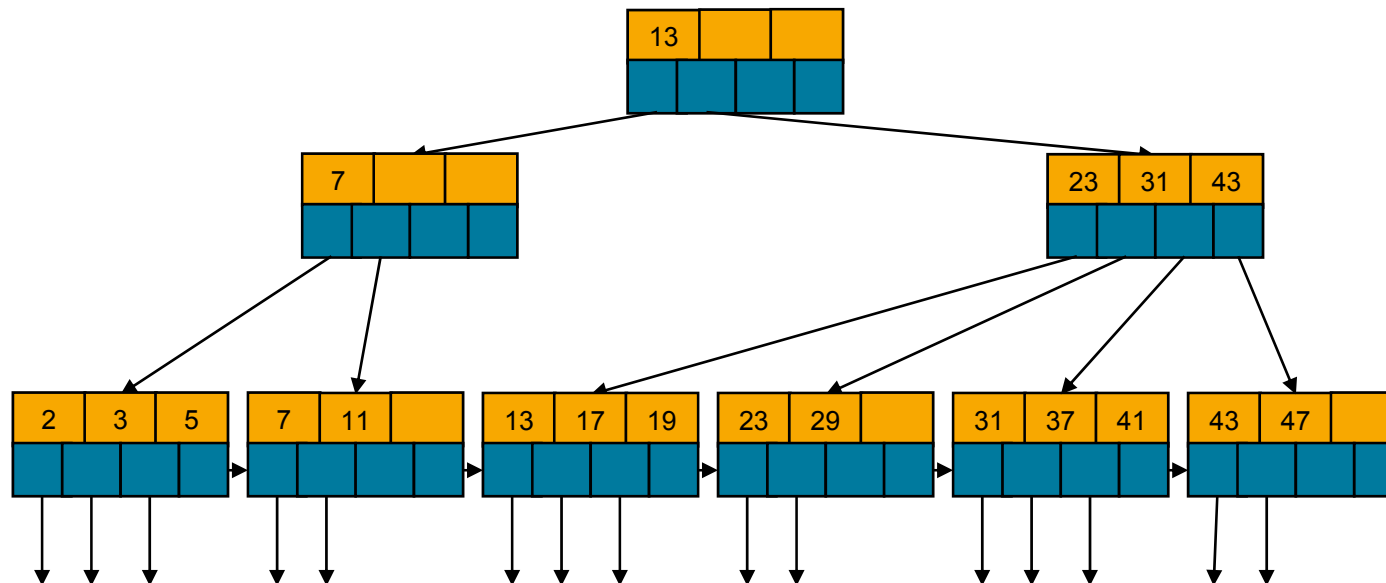
37

- Schlüssel in Blätter sind Schlüssel aus den Daten
 - Sortiert über alle Blätter verteilt (von links nach rechts)
- **Wurzel:** mindestens zwei verwendete Pointer.
 - Alle Pointer zeigen auf B-Baum Block in einer Ebene darunter
- **Blätter:** Der letzte Pointer zeigt auf das nächste Blatt (rechts)
 - Von den übrigen n Pointern werden mindestens $\lfloor (n+1)/2 \rfloor$ verwendet.
 - Zeigen auf Datenblöcke
- **Innere Knoten:** Pointer zeigen auf B-Baum Blöcke darunterliegender Ebenen
 - Mindestens $\lceil (n+1)/2 \rceil$ sind verwendet
 - Falls j Pointer verwendet werden, gibt es $j-1$ Schlüssel in dem Block
 - ◇ K_1, \dots, K_{j-1}
 - Erster Pointer zeigt auf Teilbaum mit Schlüsselwerten $< K_1$.
 - Zweiter Pointer auf Teilbaum mit Schlüsselwerten zwischen K_1 und K_2 .

Rechenbeispiele

38

- $n = 3$
 - Alle Knoten: Maximal 3 Suchschlüssel und 4 Pointer
 - Wurzel: Mindestens 1 Suchschlüssel und 2 Pointer
 - Innere Knoten: Mindestens 1 Suchschlüssel und 2 Pointer
 - Blätter: Mindestens 2 Suchschlüssel und 3 Pointer



Rechenbeispiele

39

- $n = 4$
 - Alle Knoten: Maximal 4 Suchschlüssel und 5 Pointer
 - Wurzel: Mindestens 1 Suchschlüssel und 2 Pointer
 - Innere Knoten: Mindestens $\lceil (n+1)/2 \rceil$ Pointer = 3 Pointer
 - ◇ \Rightarrow Mindestens 2 Suchschlüssel
 - Blätter:
 - ◇ 1 Pointer zum nächsten Blatt + mindestens $\lfloor (n+1)/2 \rfloor$ weitere Pointer = 3 Pointer
 - ◇ \Rightarrow Mindestens 2 Suchschlüssel
- $n = 5$
 - Alle Knoten: Maximal 5 Suchschlüssel und 6 Pointer
 - Wurzel: Mindestens 1 Suchschlüssel und 2 Pointer
 - Innere Knoten: Mindestens $\lceil (n+1)/2 \rceil$ Pointer = 3 Pointer
 - ◇ \Rightarrow Mindestens 2 Suchschlüssel
 - Blätter:
 - ◇ 1 Pointer zum nächsten Blatt + mindestens $\lfloor (n+1)/2 \rfloor$ weitere Pointer = 4 Pointer
 - ◇ \Rightarrow Mindestens 3 Suchschlüssel

Alternative Definition

40

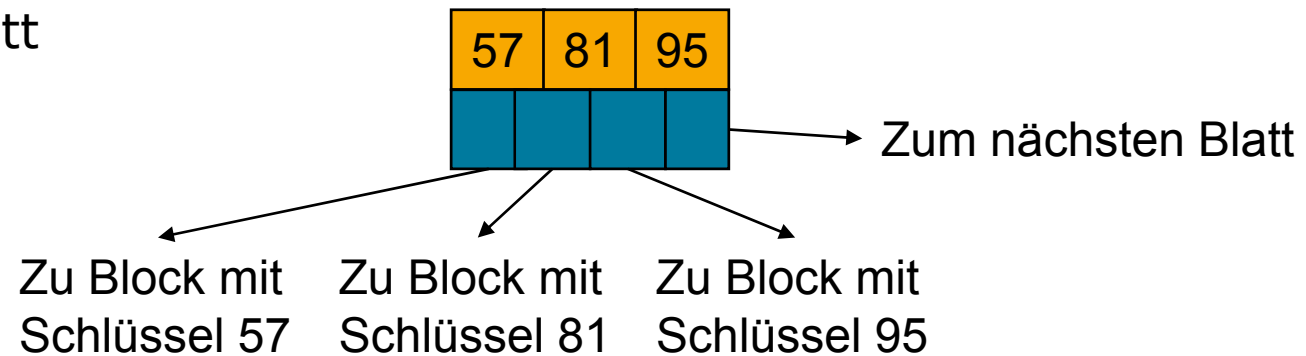
- Bisher: Parameter n
 - Block hat mindestens $\lfloor (n+1)/2 \rfloor$ Suchschlüssel
 - Block hat höchstens n Suchschlüssel
- Alternativ in Lehrbüchern: Parameter k
 - Block hat mindestens k Suchschlüssel
 - Block hat höchstens $2k$ Suchschlüssel
 - Block hat immer $x + 1$ Pointer (wie bisher)
- Immer
 - Ein innerer Block hat immer einen Pointer mehr als Anzahl Suchschlüssel
 - Ein Blatt hat immer ebenso viele Pointer wie Suchschlüssel
 - ◇ Plus verkettete Liste

Beispiel Blattknoten

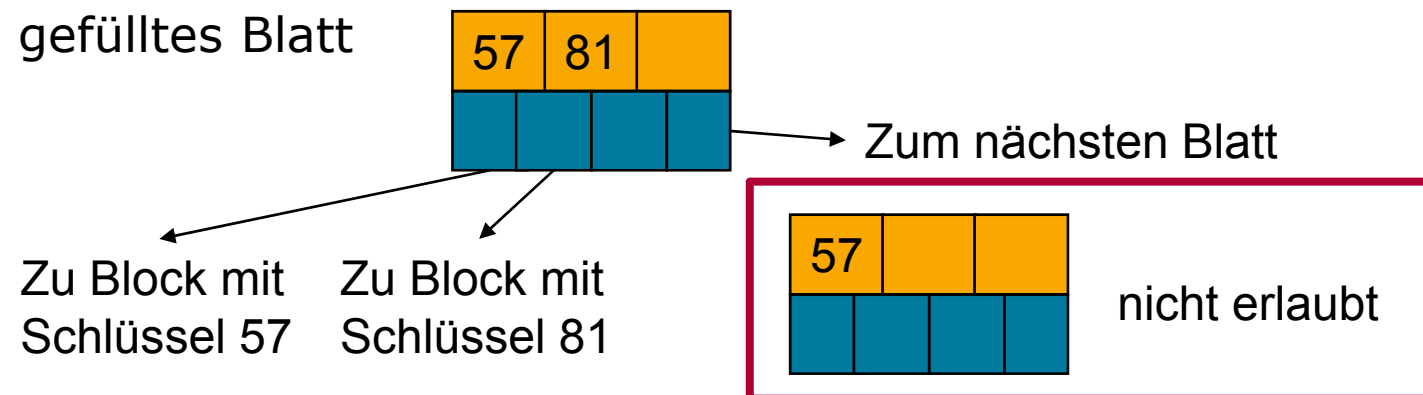
41

- $n = 3$
 - 3 Schlüssel und 4 Pointer

- Volles Blatt



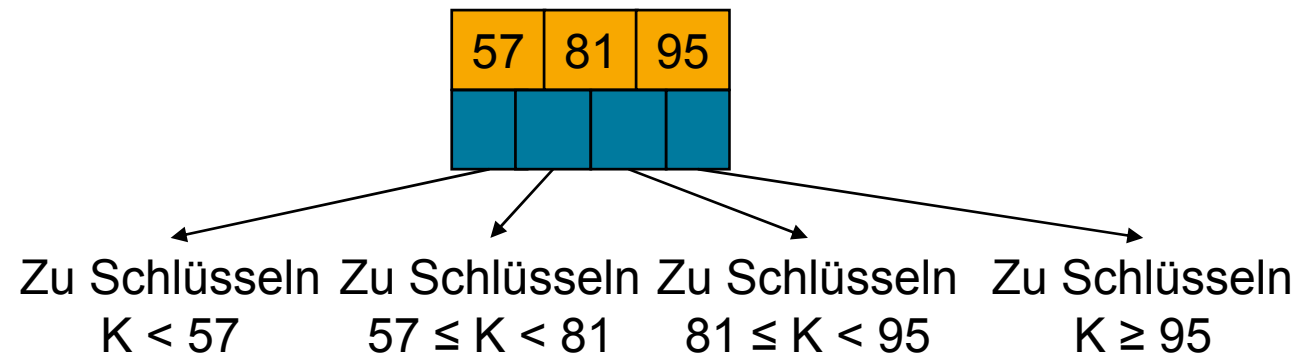
- Teilweise gefülltes Blatt



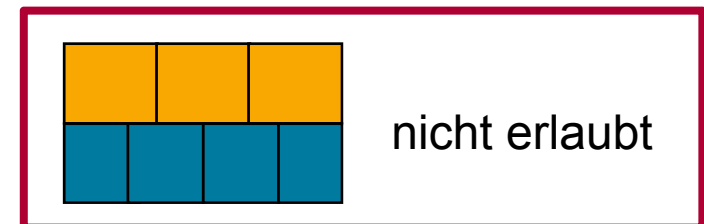
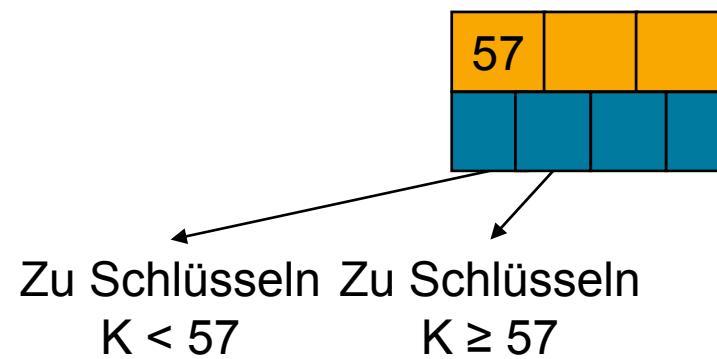
Beispiel innerer Knoten

42

- Voller innerer Knoten



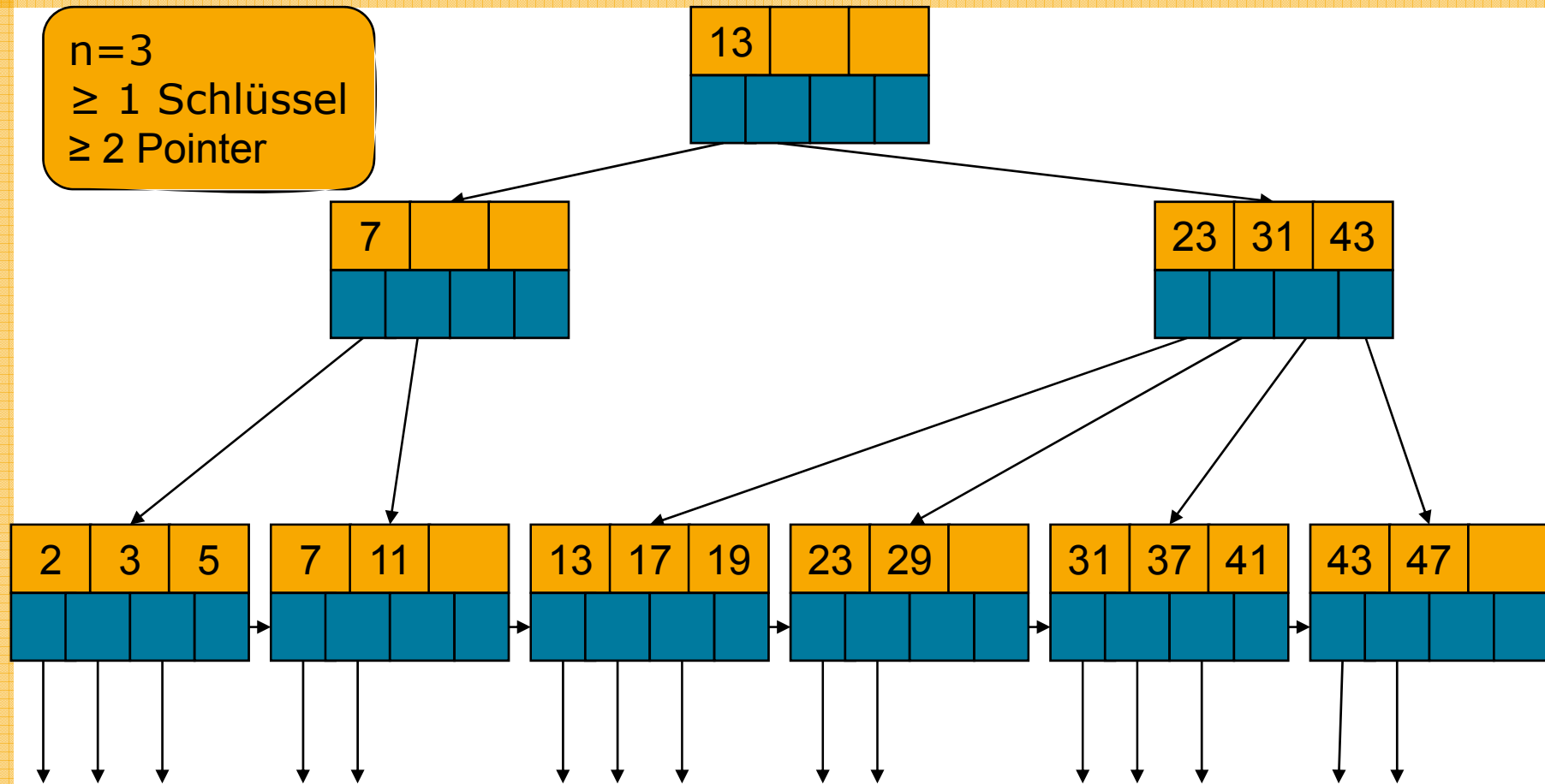
- Teilweise gefüllter innerer Knoten



Beispiel B-Baum

43

$n=3$
 ≥ 1 Schlüssel
 ≥ 2 Pointer



In den Blättern taucht sortiert jeder Schlüssel genau einmal auf.

Anwendungen von B-Bäumen

44

B-Bäume können verschiedene Index-Rollen übernehmen.

- Suchschlüssel ist Primärschlüssel; dicht-besetzter Index
 - Data file sortiert oder nicht
- Dünn-besetzter Index; data file ist sortiert
- Suchschlüssel ist nicht Primärschlüssel; Dicht-besetzter Index
 - Data file ist nach Suchschlüssel sortiert
 - Pointer zeigen auf jeweils ersten Wert
- ...

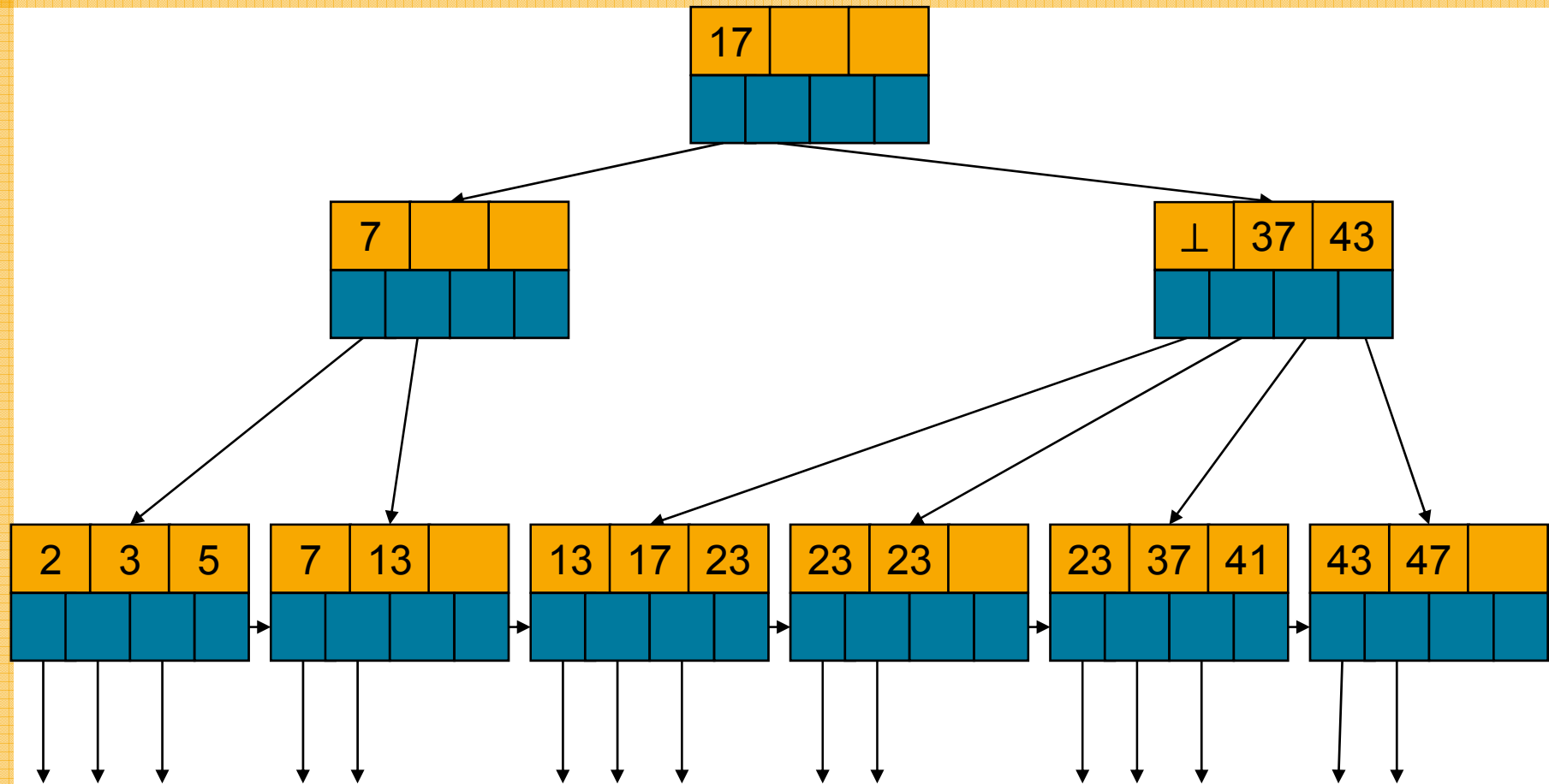
B-Bäume auf nicht-Primärschlüsseln

45

- Bedeutung der Pointer auf inneren Ebenen ändert sich.
 - Gegeben: Schlüssel K_1, \dots, K_j
 - $\Rightarrow K_i$ ist der kleinste neue Schlüsselwert, der vom $(i+1)$ -ten Pointer erreichbar ist.
 - ◇ D.h. es gibt keinen Schlüsselwert K_i im linken Teilbaum aber mindestens ein Vorkommen des Schlüsselwertes Teilbaum vom $(i+1)$ -ten Pointer an.
 - ◇ Problem: Es gibt nicht immer einen solchen Schlüssel

B-Bäume auf nicht-Primärschlüsseln

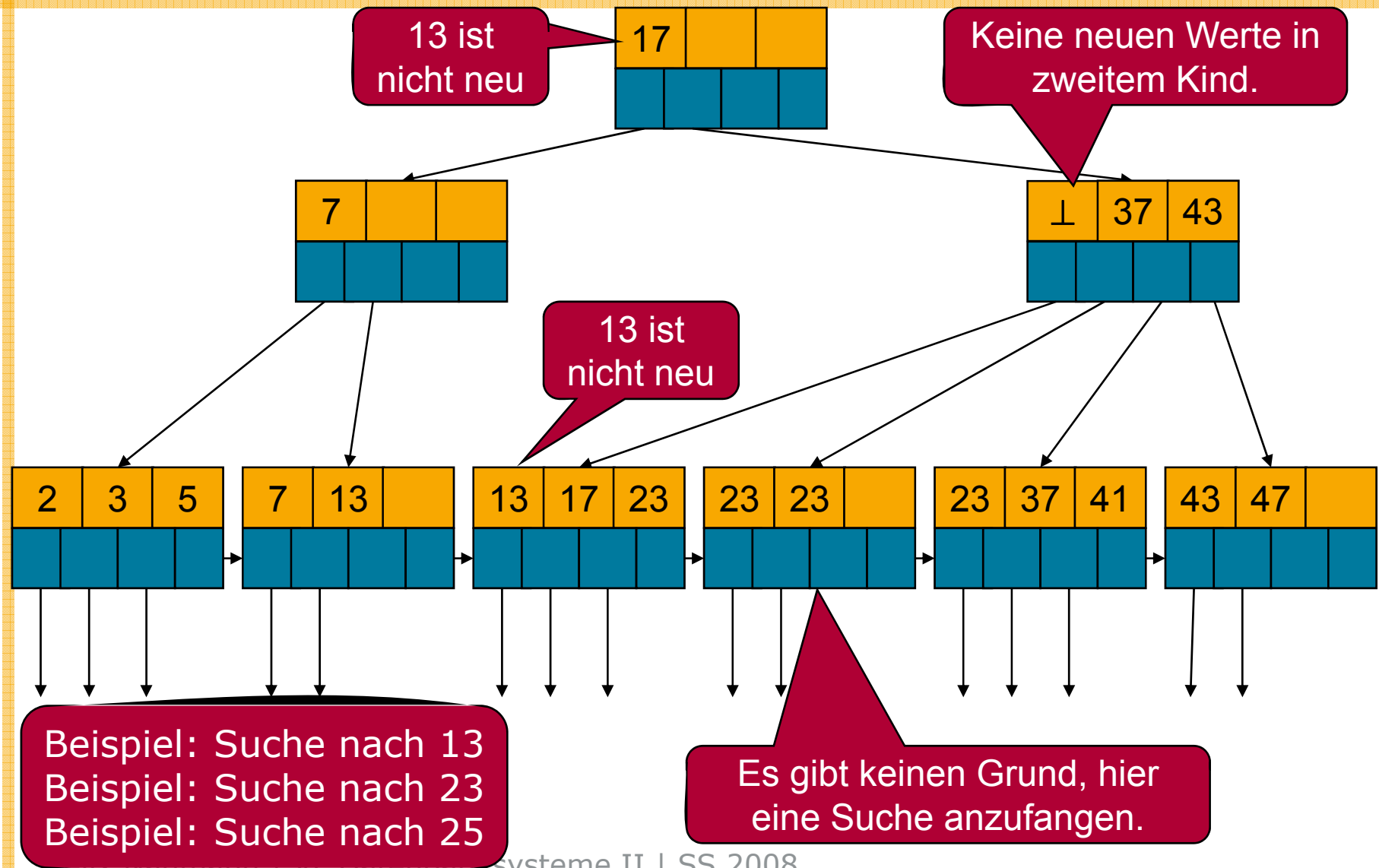
46



K_i ist der kleinste neue Schlüsselwert, der vom (i+1)-ten Pointer erreichbar ist.

B-Bäume auf nicht-Primärschlüsseln

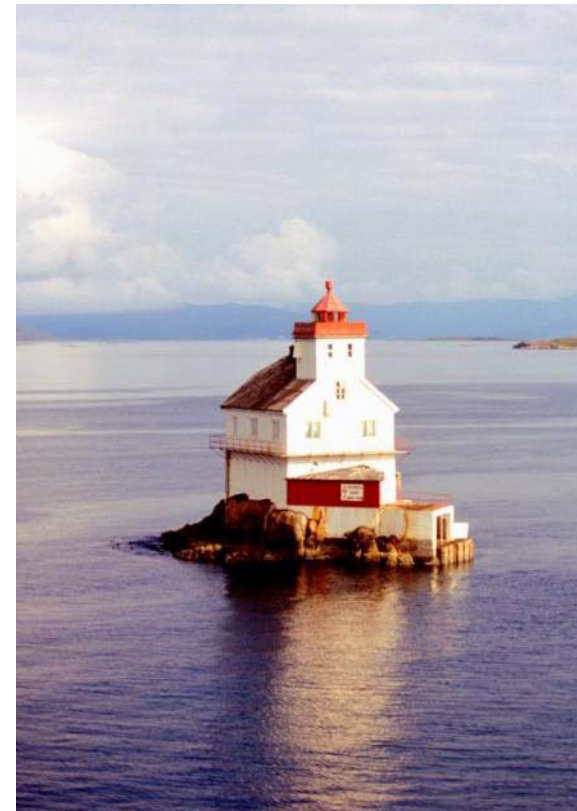
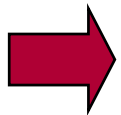
47



Überblick

48

- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
 - Aufbau
 - Suche
 - Updates
 - Effizienz
 - Varianten
- Hash-Tabellen



Suche in B-Bäumen

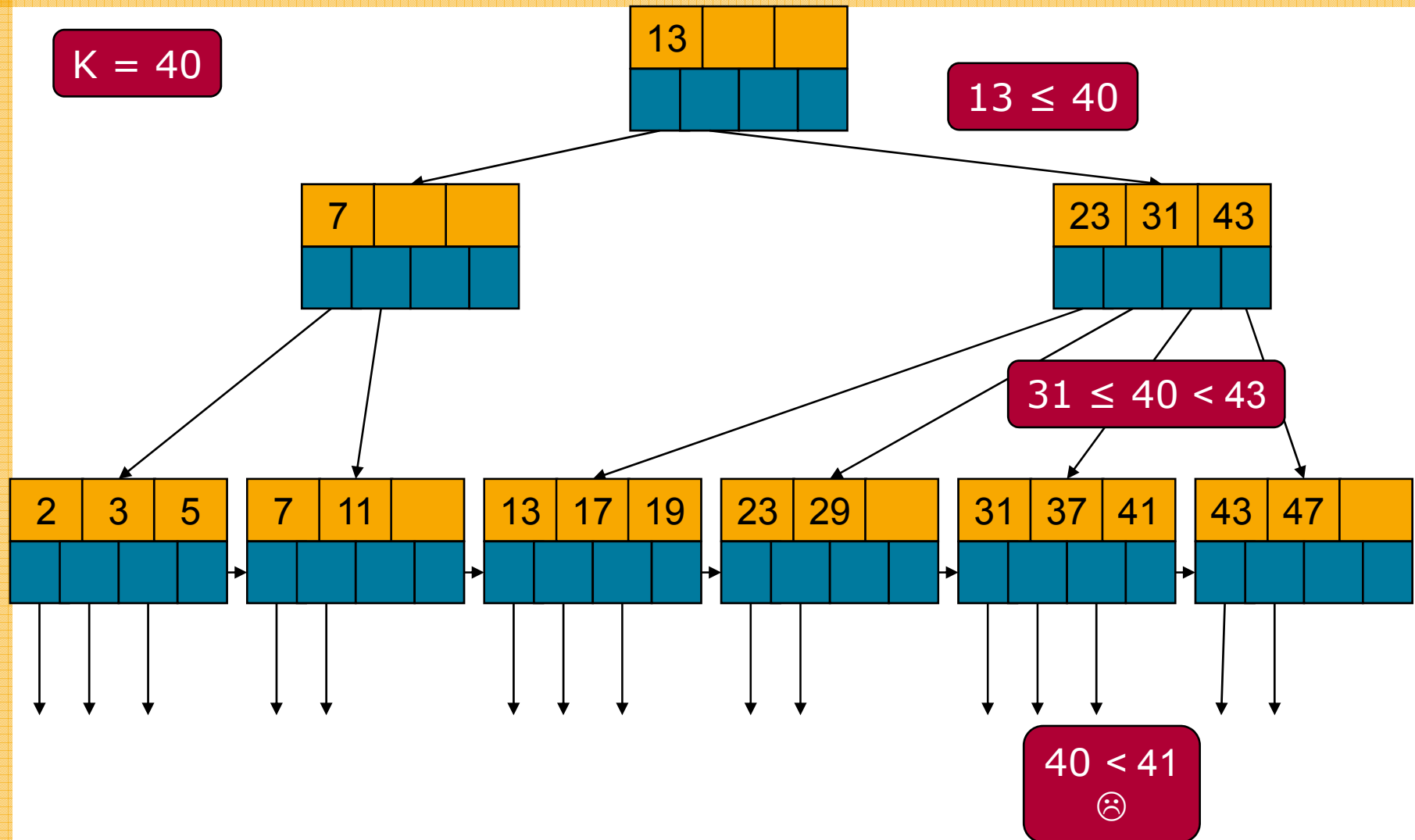
49

- Jetzt wieder: Suchschlüssel = Primärschlüssel
- Dicht-besetzter Index
 - Operationen für dünn-besetzte Indizes ähnlich
- Gesucht sei K .
 - Falls wir an einem Blattknoten sind:
 - ◇ Suche K auf dem Knoten.
 - Falls wir an einem inneren Knoten mit K_1, K_2, \dots, K_n sind:
 - ◇ Falls $K < K_1$ gehe zu erstem Kind
 - ◇ Falls $K_1 \leq K < K_2$ gehe zu zweitem Kind
 - ◇ ...
 - ◇ Falls $K_n \leq K$ gehe zu letztem Kind

Beispiel Suche im B-Baum

50

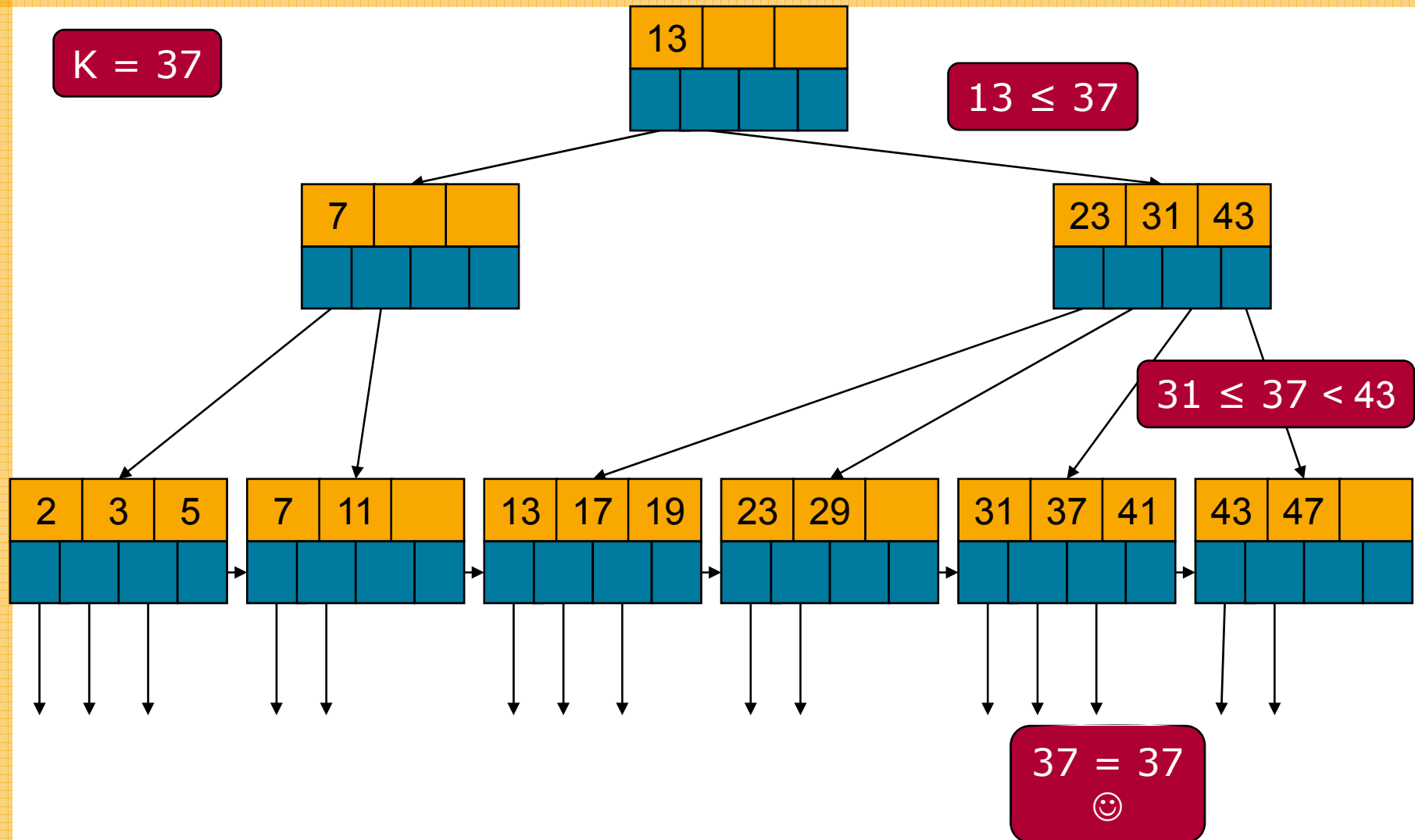
$K = 40$



Beispiel Suche im B-Baum

51

$K = 37$



Bereichsanfragen (*range queries*)

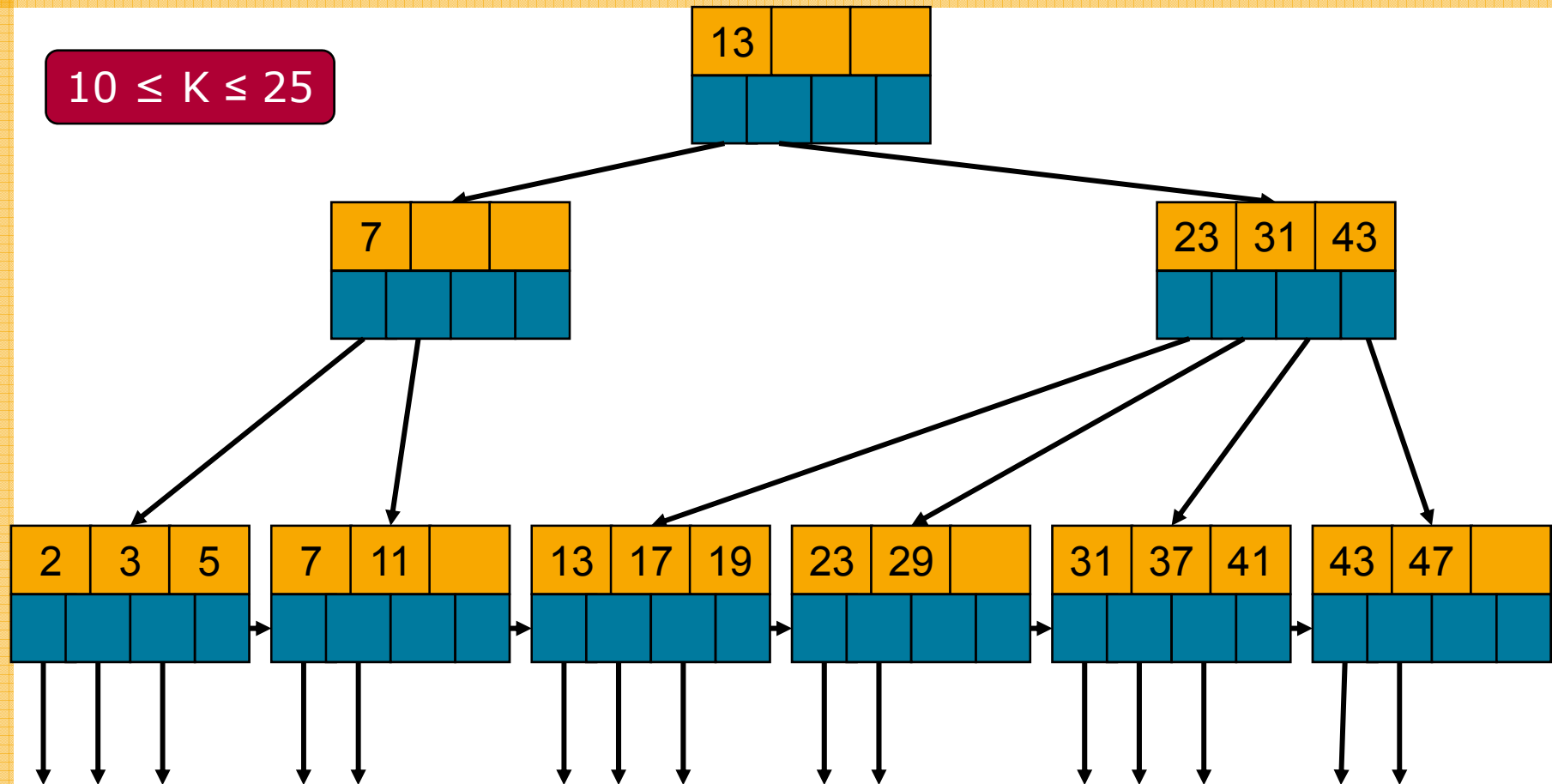
52

- Anfragen mit Ungleichheit in WHERE Klausel
 - `SELECT * FROM R`
`WHERE R.k > 40`
 - `SELECT * FROM R`
`WHERE R.k >= 10 AND R.k <= 25`
- Suche nach Bereich $[a, b]$
 1. Suche in B-Baum nach a .
 - ◇ Entsprechendes Blatt könnte a speichern.
 - ◇ Suche auf dem Blatt alle relevanten Schlüssel.
 2. Falls auf dem Blatt kein Wert $> b$.
 - ◇ Folge Pointer zu nächstem Knoten.
- Bei offenen Bereichen $[-\infty, b]$ bzw. $[a, \infty]$ ist Suche ähnlich.

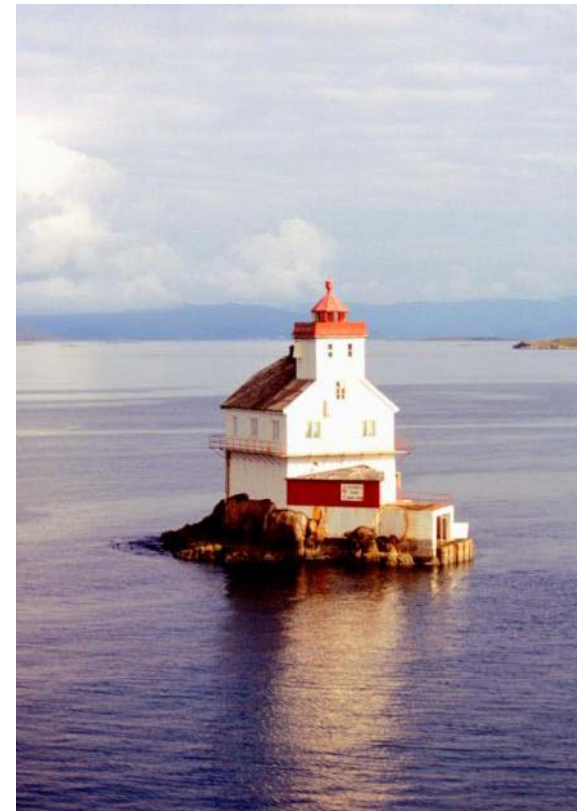
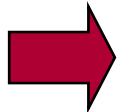
Bereichsanfragen (range queries)

53

$10 \leq K \leq 25$



- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
 - Aufbau
 - Suche
 - Updates
 - Effizienz
 - Varianten
- Hash-Tabellen



Einfügen in B-Bäume

55

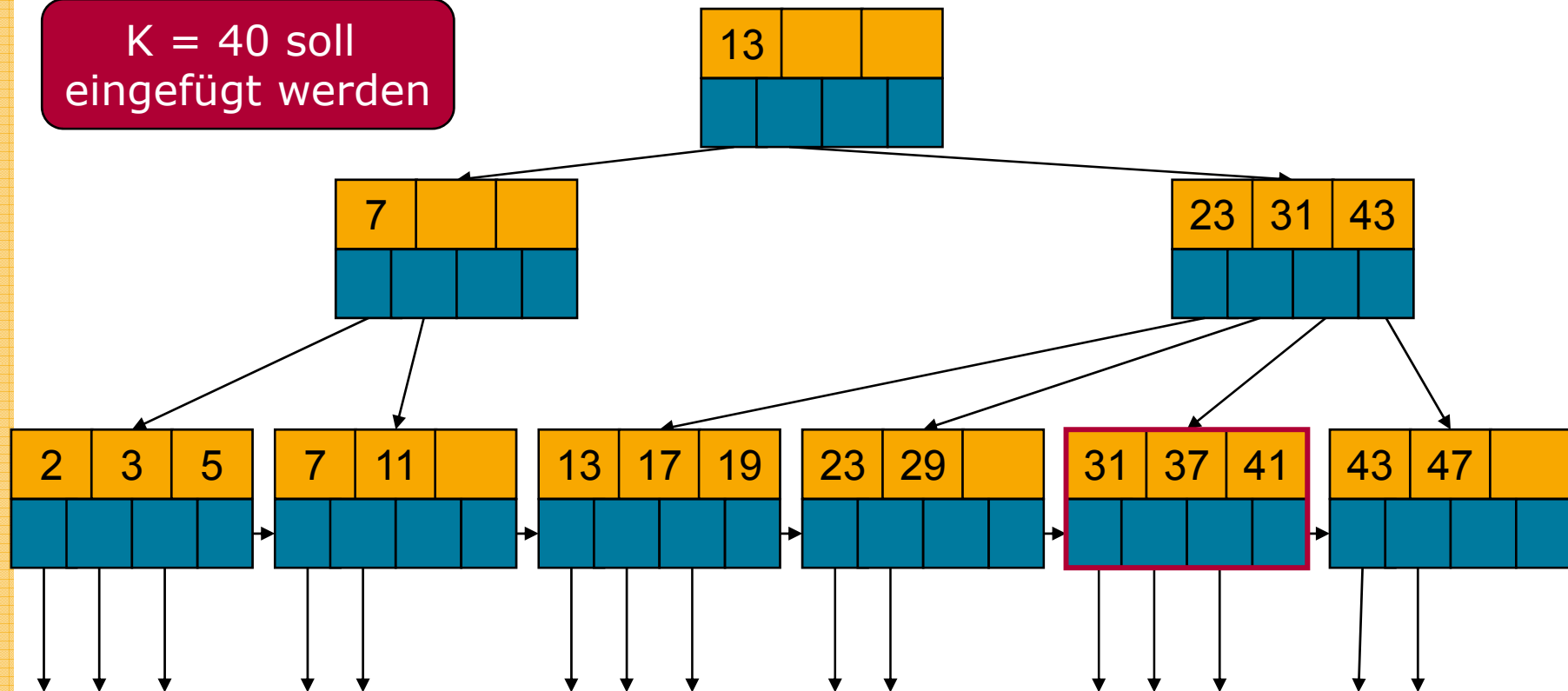
Rekursives Vorgehen:

- Suche entsprechendes Blatt.
 - Falls Platz herrscht, füge Schlüssel und Pointer ein.
- Falls kein Platz: Überlauf
 - Teile Blatt in zwei Teile und verteile Schlüssel gleichmäßig
 - „*split*“
- Teilung macht Einfügen eines neuen Schlüssel/Pointer-Paares im Elternknoten erforderlich
 - Gehe rekursiv aufwärts im Baum vor
- Ausnahme: Falls in Wurzel kein Platz
 - Teile Wurzel in zwei
 - Erzeuge neue Wurzel (mit nur einem Schlüssel)

Einfügen in B-Bäume – Beispiel

56

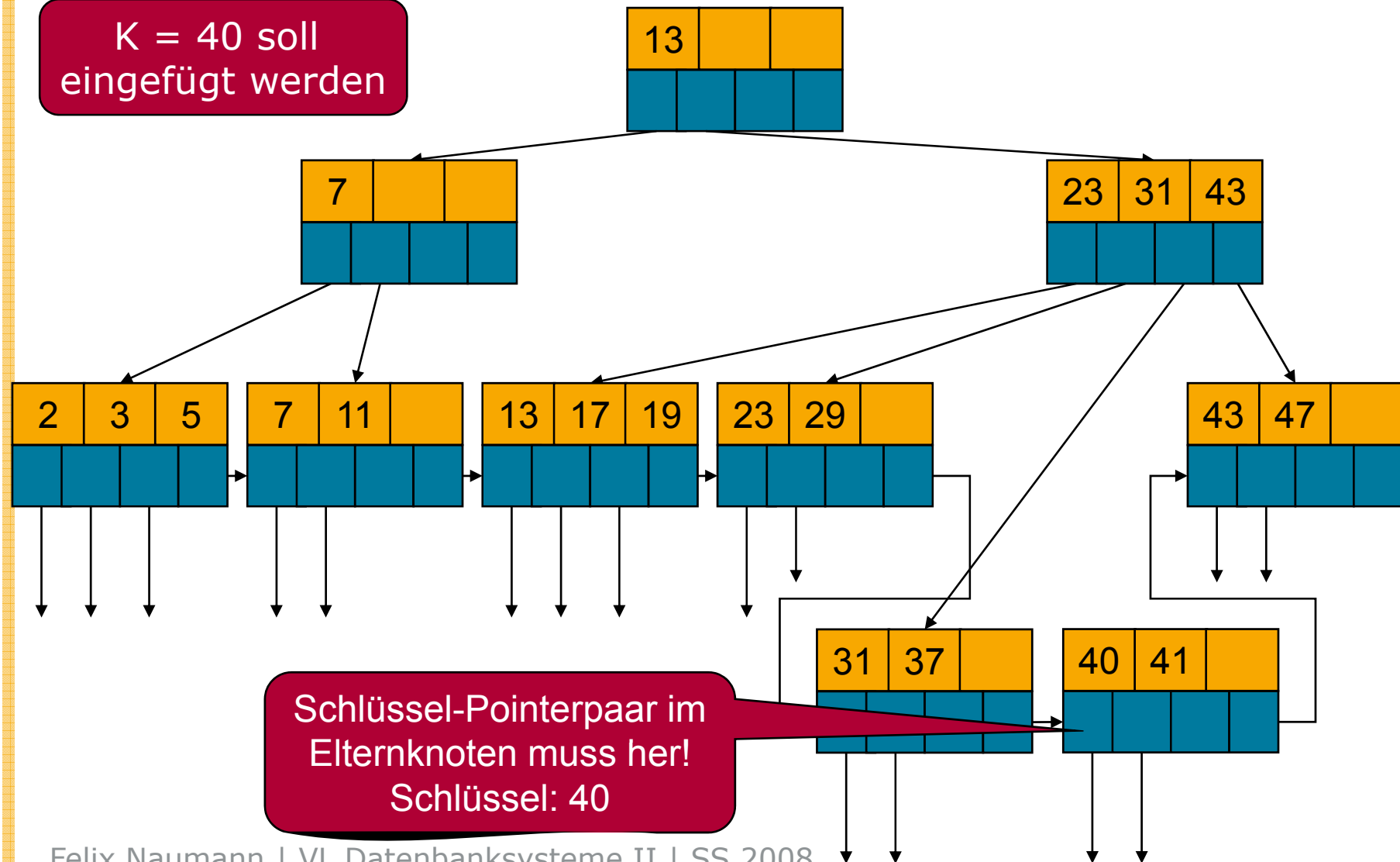
K = 40 soll eingefügt werden



Einfügen in B-Bäume – Beispiel

57

K = 40 soll eingefügt werden

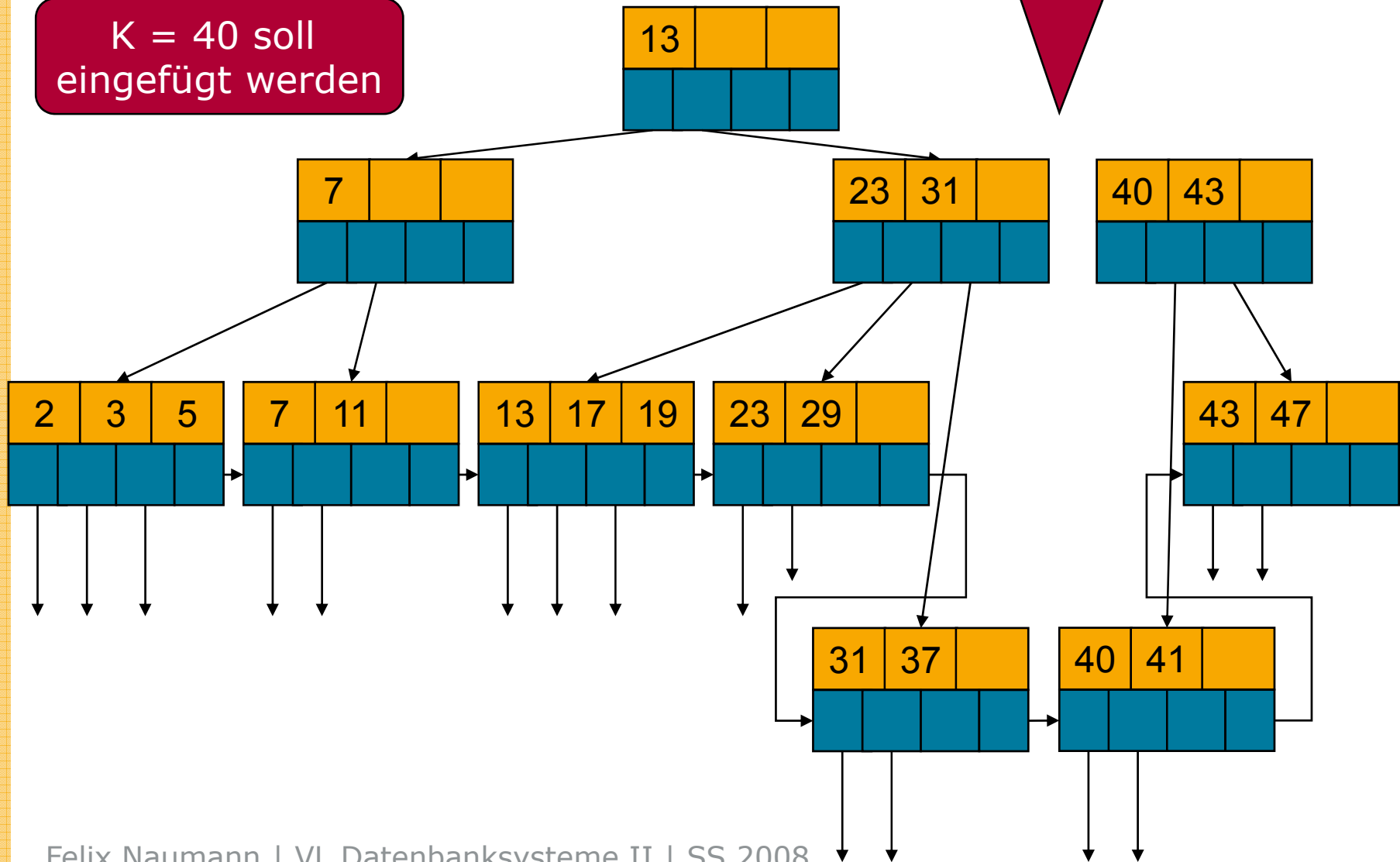


Einfügen in B-Bäume – Beispiel

Schlüssel-Pointerpaar
muss her
Schlüssel: 40

58

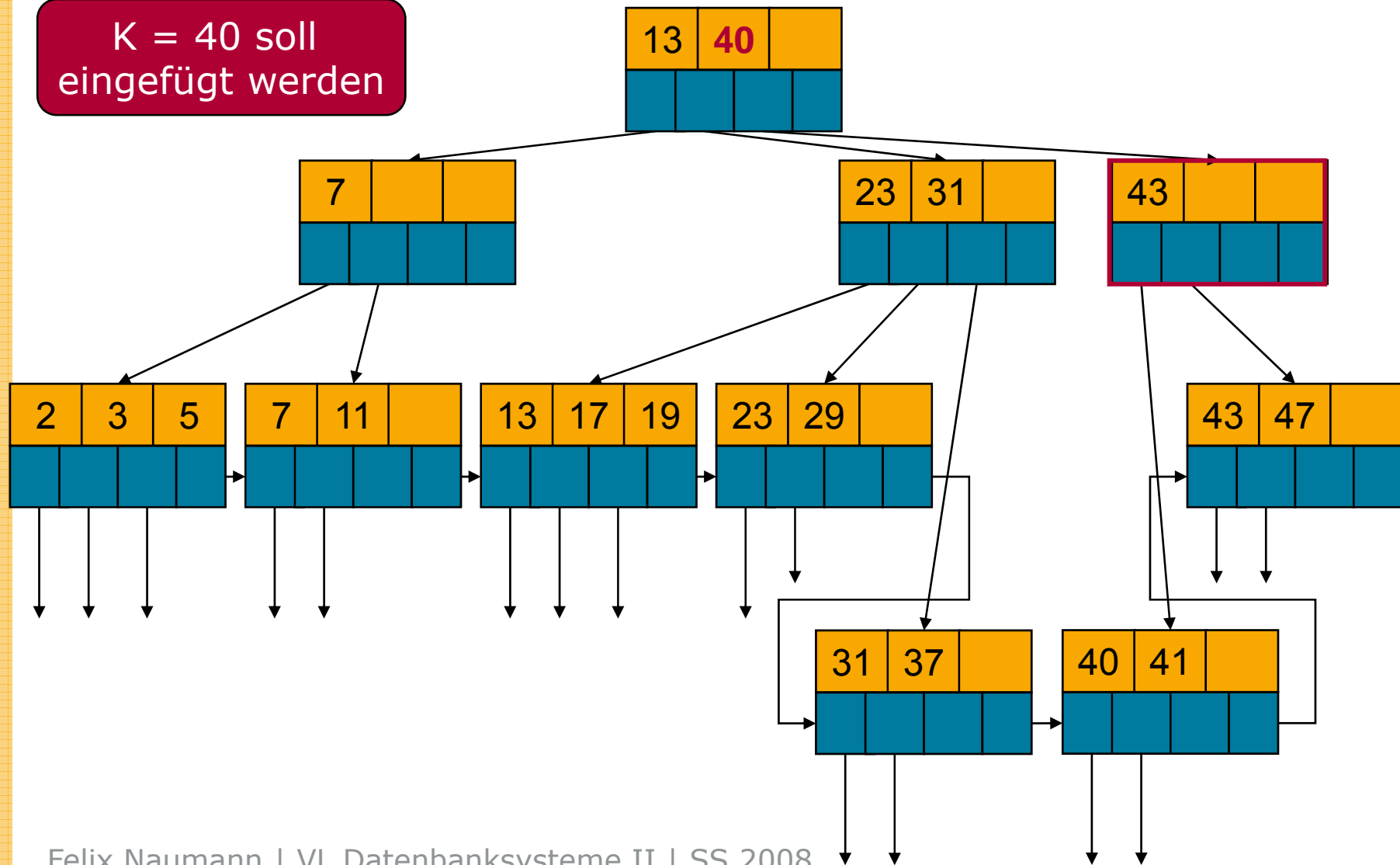
K = 40 soll
eingefügt werden



Einfügen in B-Bäume – Beispiel

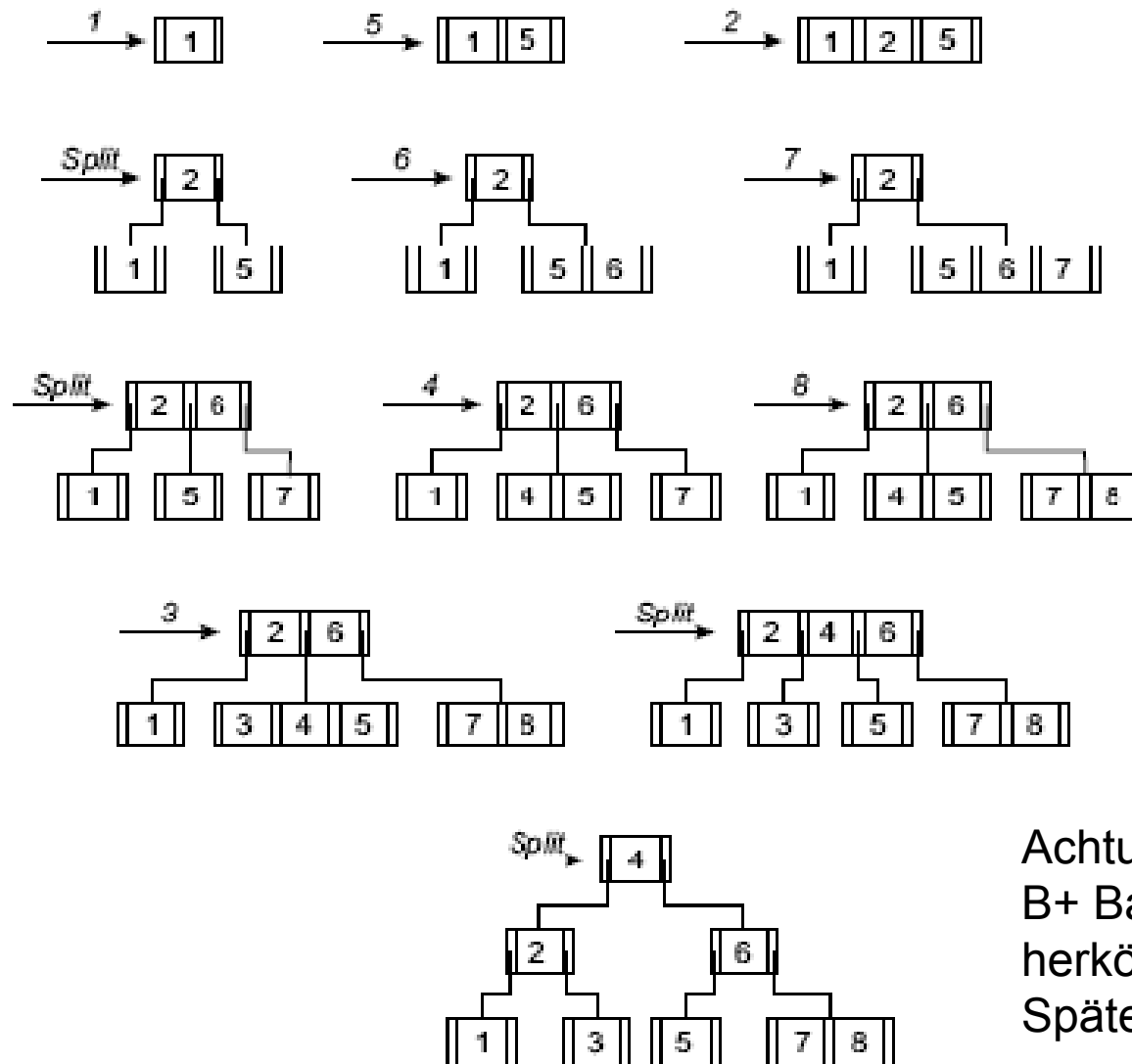
59

K = 40 soll eingefügt werden



Einfügen in B-Bäume – Beispiel

60



Achtung: Dies ist kein B+ Baum, sondern ein herkömmlicher B-Baum. Später mehr...

Kosten für Einfügen

61

- Sei h die Höhe des B-Baums
 - Meist $h = 3$
- Suche nach Blattknoten: h
- Falls keine Teilung nötig: Gesamtkosten $h + 1$
 - h Blöcke lesen, 1 Block schreiben
- Falls Teilung nötig
 - Worst-case: Bis zur Wurzel
 - Selbst Caching nützt nichts, da Knoten geschrieben werden müssen
 - Insgesamt: $3h + 1$
 - ◇ Auf jeder Ebene Suche und Überlaufbehandlung
 - ◇ + neue Wurzel schreiben

Löschen aus B-Bäumen

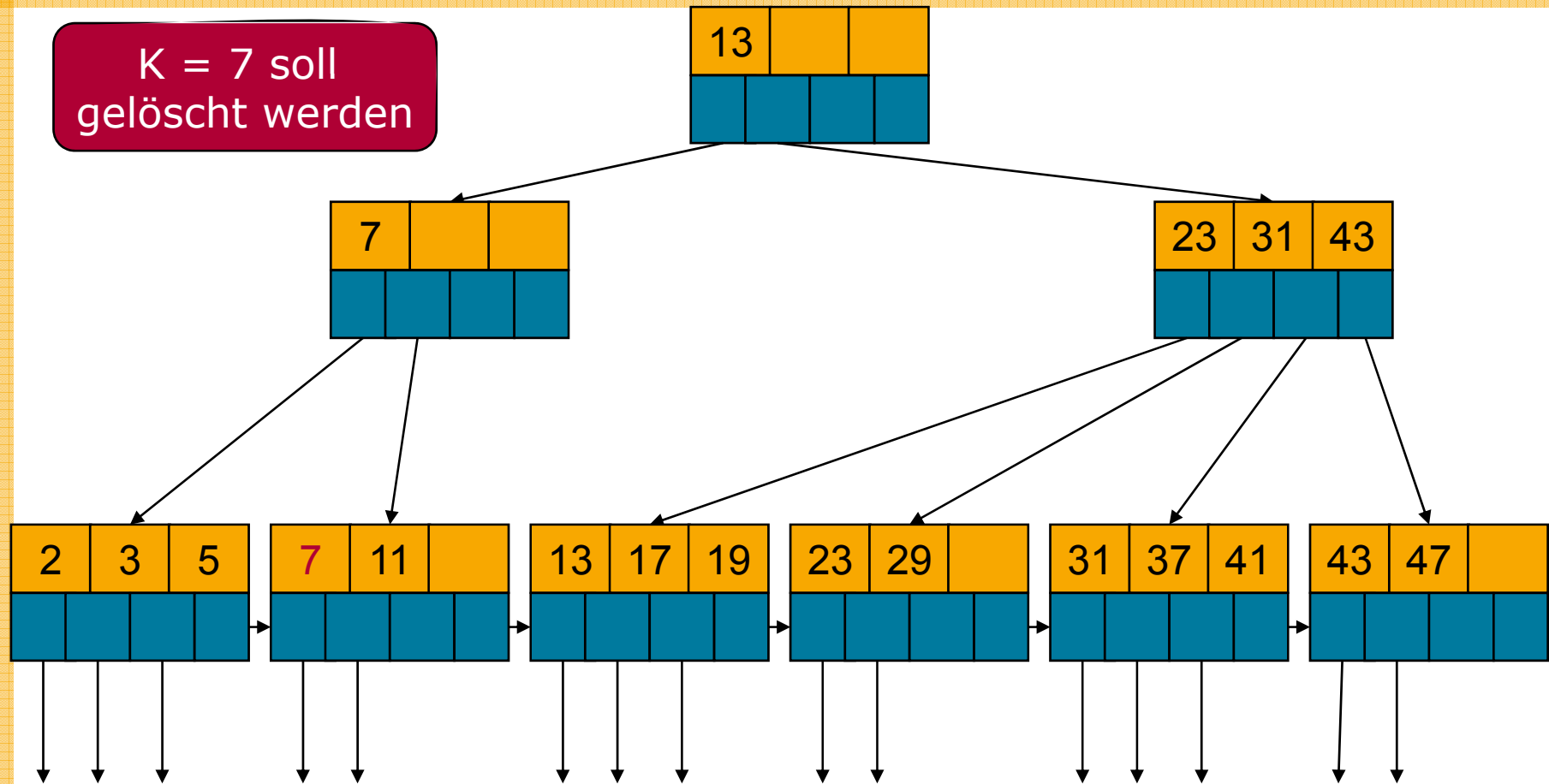
62

- Suche entsprechenden Knoten
- Lösche Schlüssel
- Falls immer noch minimale Menge an Schlüsseln im Knoten
 - Nichts tun
- Falls zu wenig Schlüssel im Knoten: *Merge* (Verschmelzen)
 - Falls ein Geschwisterknoten (links oder rechts) mehr als die minimale Schlüsselmenge hat, „klaue“ einen Schlüssel.
 - ◇ Gegebenenfalls Schlüsselwerte der Eltern anpassen
 - Falls nicht: Es existieren zwei Geschwister im Baum mit minimaler und sub-minimaler Schlüsselmenge
 - ◇ => Diese Knoten können vereinigt werden.
 - ◇ Schlüsselwerte der Eltern anpassen
 - gegebenenfalls rekursiv im Baum nach oben löschen

Löschen aus B-Bäumen – Beispiel

63

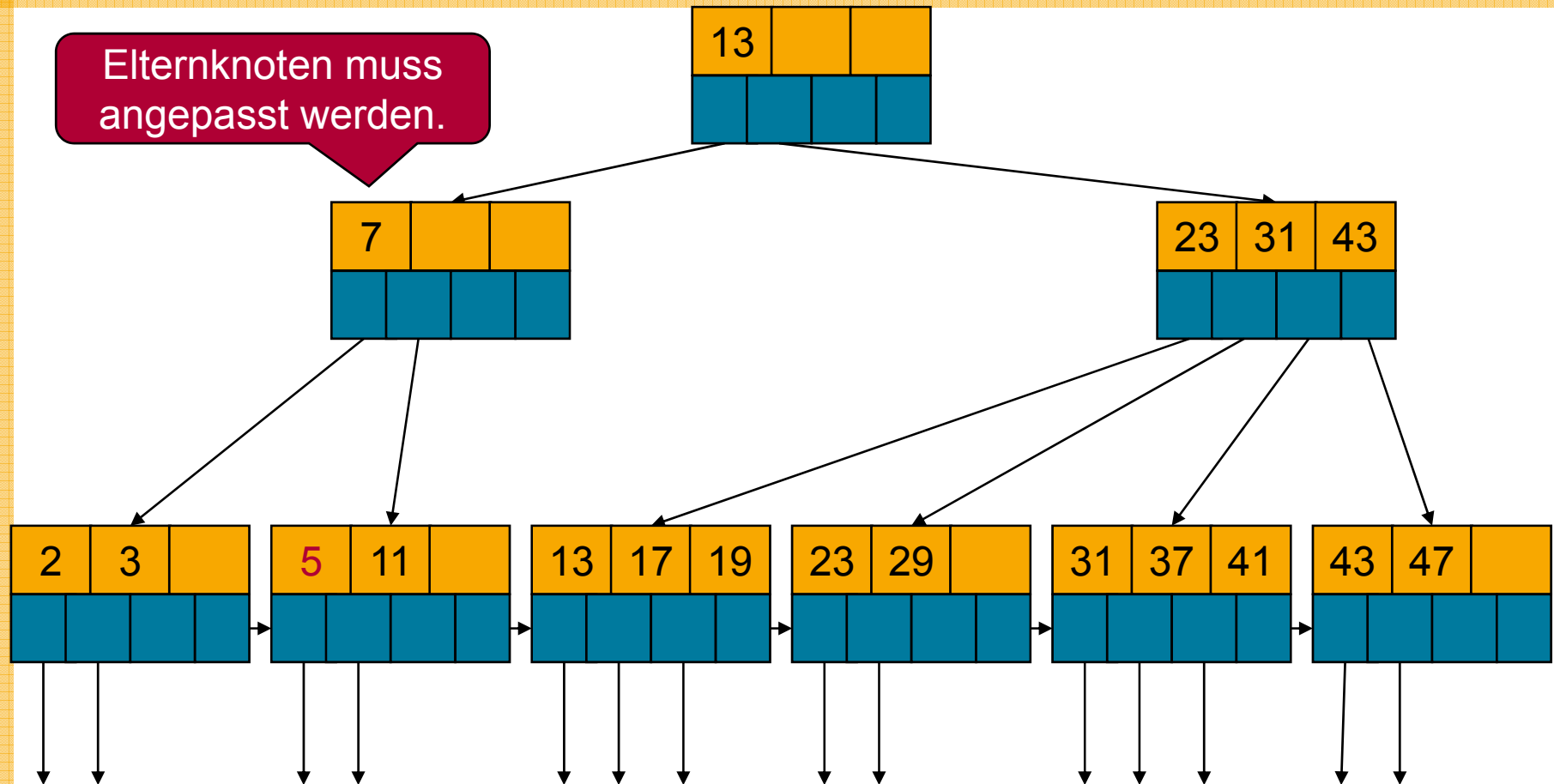
K = 7 soll gelöscht werden



Löschen aus B-Bäumen – Beispiel

64

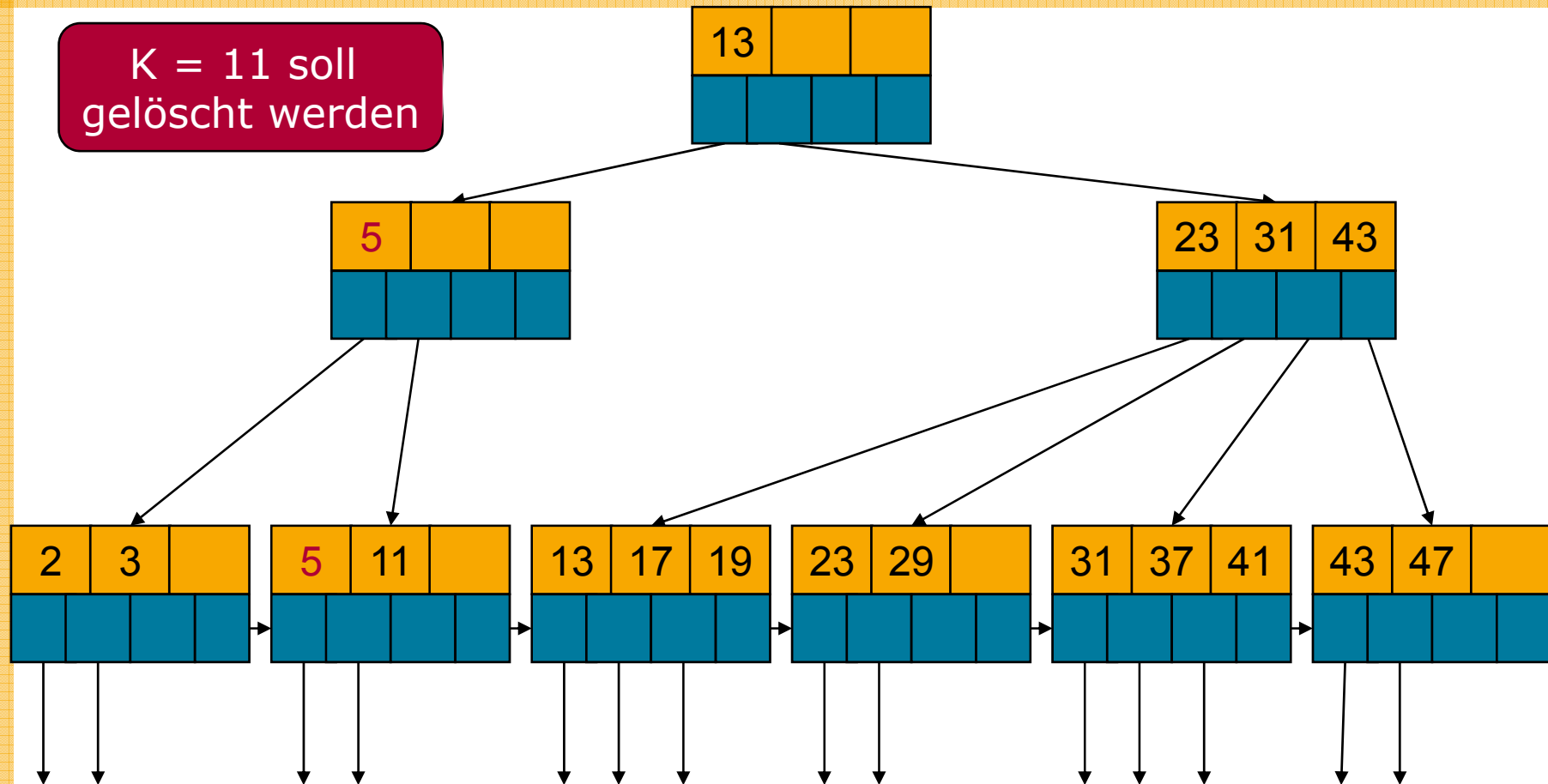
Elternknoten muss angepasst werden.



Löschen aus B-Bäumen – Beispiel

65

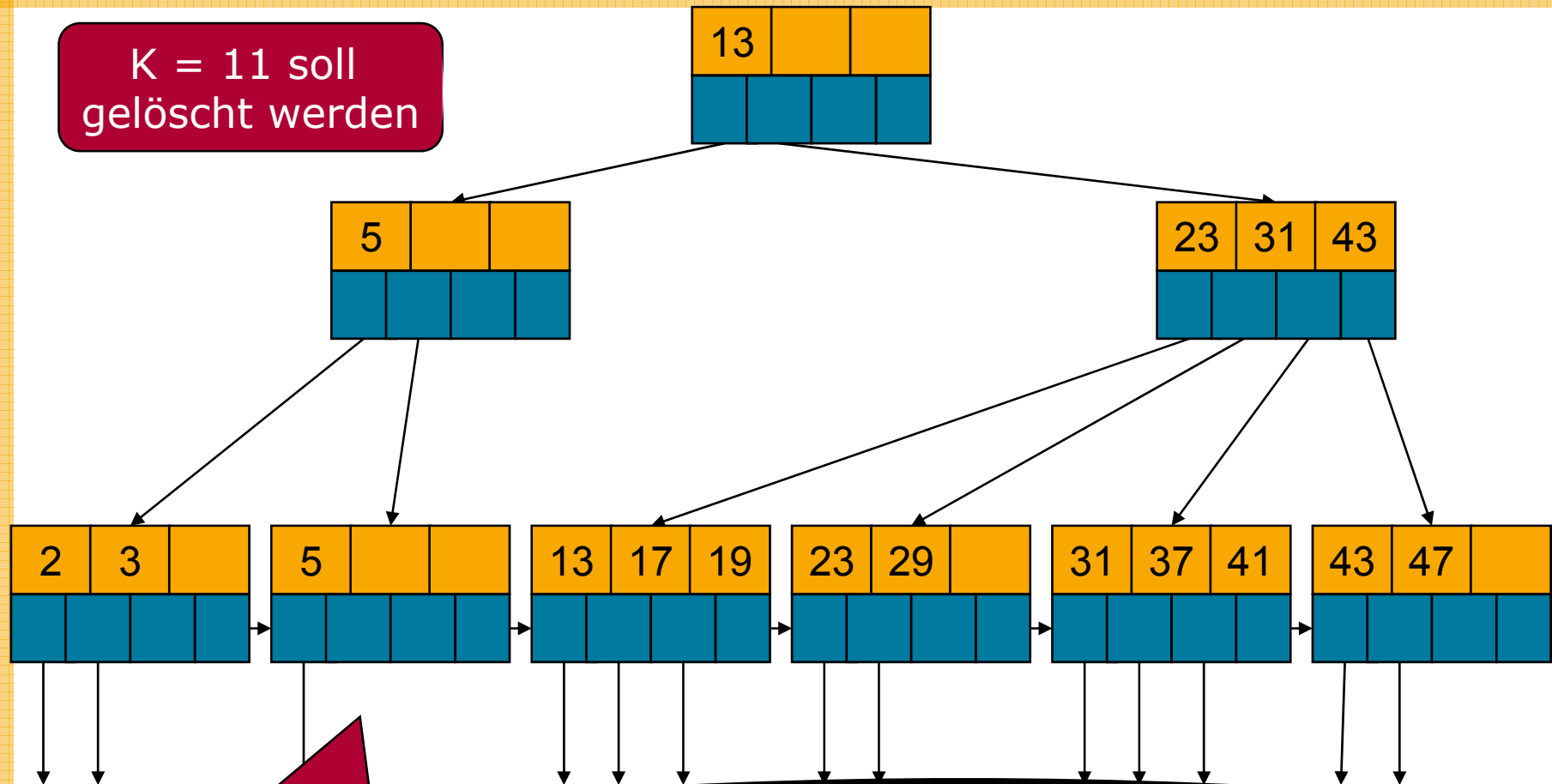
K = 11 soll gelöscht werden



Löschen aus B-Bäumen – Beispiel

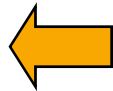
66

K = 11 soll gelöscht werden



Knoten zu klein.
Geschwisterknoten kann nichts mehr abgeben.

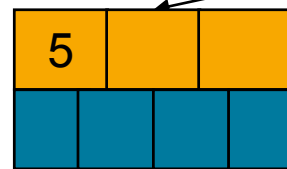
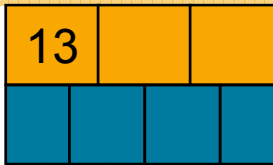
Variante 1: Von rechtem Nachbarn verschieben.
Aber: Kein Geschwister; sehr mühselig
Variante 2: Verschmelzen mit linkem Geschwister.



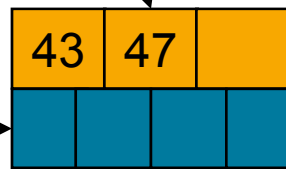
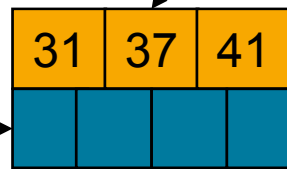
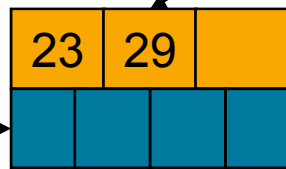
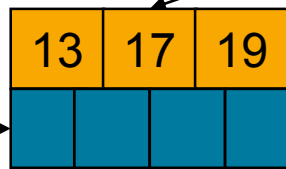
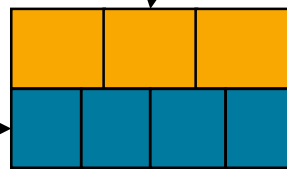
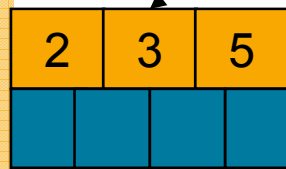
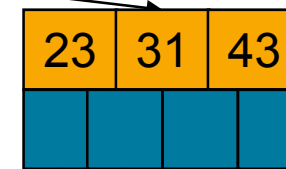
Löschen aus B-Bäumen – Beispiel

67

K = 11 soll gelöscht werden



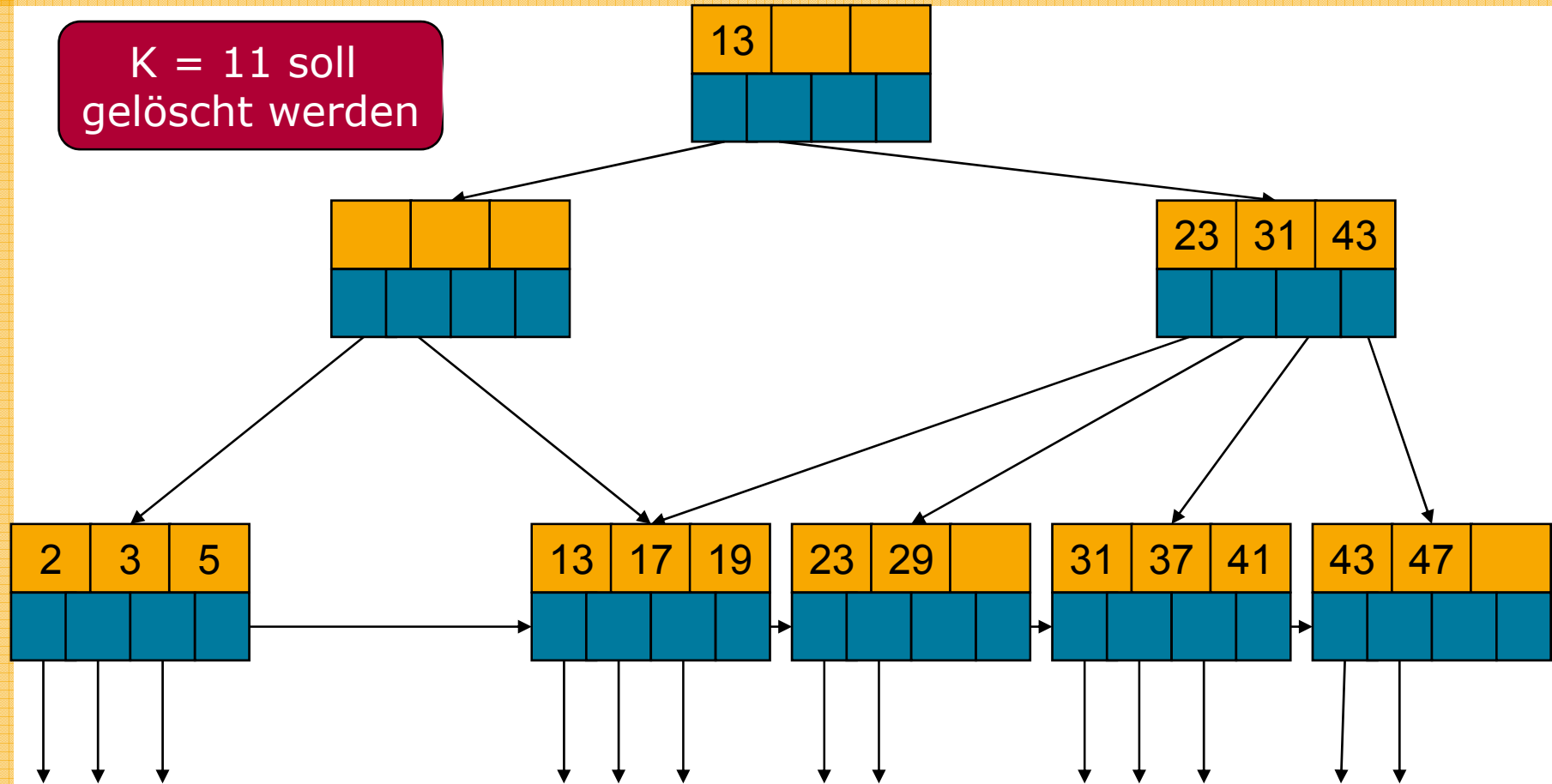
Knoten zu klein.



Löschen aus B-Bäumen – Beispiel

68

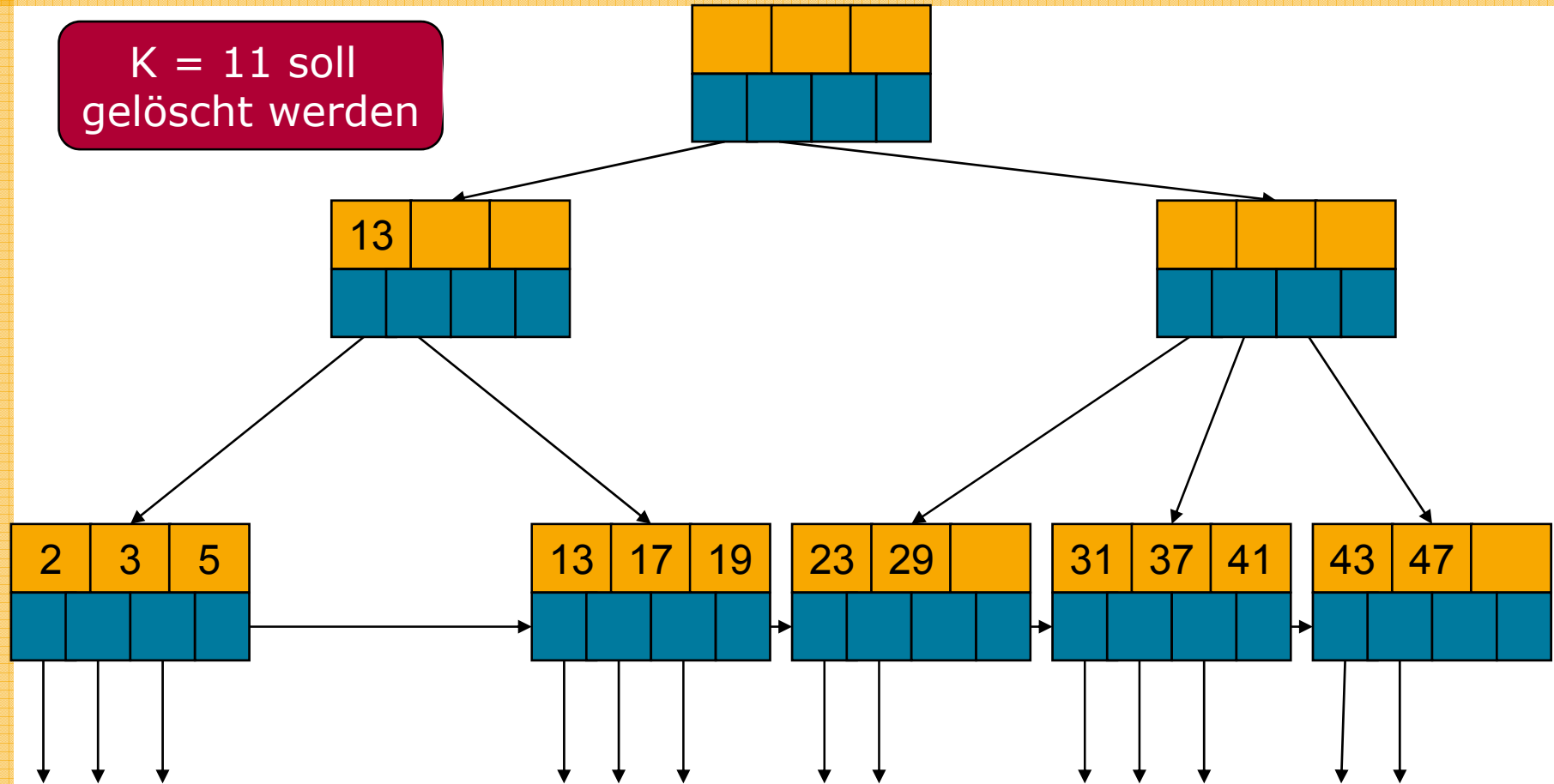
K = 11 soll gelöscht werden



Löschen aus B-Bäumen – Beispiel

69

K = 11 soll gelöscht werden



Kosten für Löschen

70

- Suche und lokales Löschen: $h + 2$
 - Schreibe Blattknoten
 - Schreibe Datenblock
- Bei merge mit Geschwisterknoten: $h + 6$
 - Prüfe rechts und links
 - Schreibe Block und veränderten Nachbarn
 - Schreibe Elternknoten
 - Schreibe Datenblock
- Bei merge bis zur Wurzel: $3h - 2$

Löschen aus B-Bäumen?

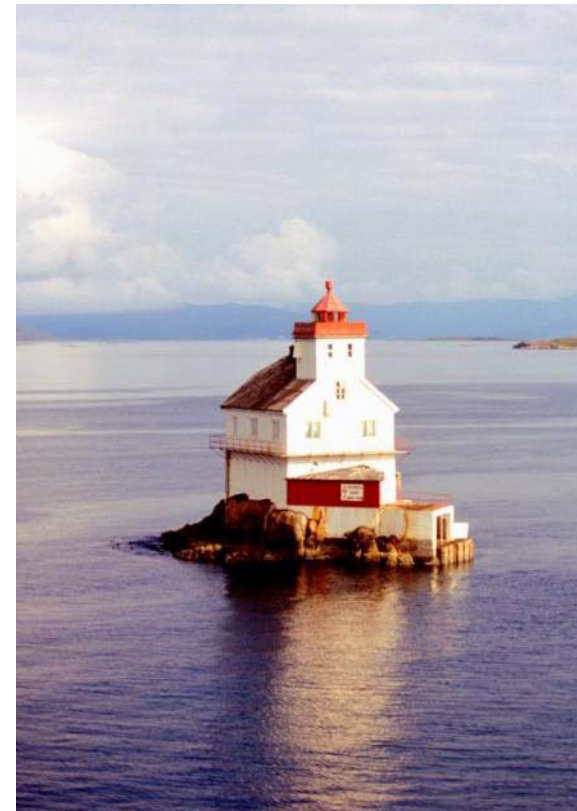
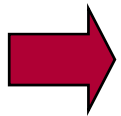
71

- Annahme: Tendenziell wachsen Datenmengen
- Folgerung: Nie Knoten des B-Baum löschen
 - Knoten, die durch Löschen zu klein werden, werden früher oder später wieder gefüllt.
 - Grabstein muss sowieso erhalten bleiben.
- Verbesserung: Grabstein im B-Baum statt im Block.

Überblick

72

- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
 - Aufbau
 - Suche
 - Updates
 - Effizienz
 - Varianten
- Hash-Tabellen



Effizienz von B-Bäumen

73

Suche, Einfügen und Löschen sollen möglichst wenig I/O-Operationen benötigen.

- Je größer n gewählt wird desto seltener müssen Blöcke verschmolzen oder getrennt werden.
 - Meist auf Blattknoten beschränkt
- Suche
 - Anzahl I/Os entspricht der Höhe des Baums
 - + 1 I/O auf den Daten für Lesen
 - + 3 I/O auf den Daten für Einfügen oder Löschen
- Typische Höhe eines B-Baums: 3
 - 340 Schlüssel-Pointer-Paare pro Block
 - Annahme: Füllstand durchschnittlich 255
 - \Rightarrow 255 innere Knoten $= > 255^2 = 65025$ Blätter $= 255^3$ Pointer
 - Insgesamt als 16.6 Mio Datensätze
- Zugriff mit 2 I/Os: Nur Wurzel im Hauptspeicher
- Zugriff mit 1 I/O: Wurzel und innere Knoten im Hauptspeicher

Bulk-loading

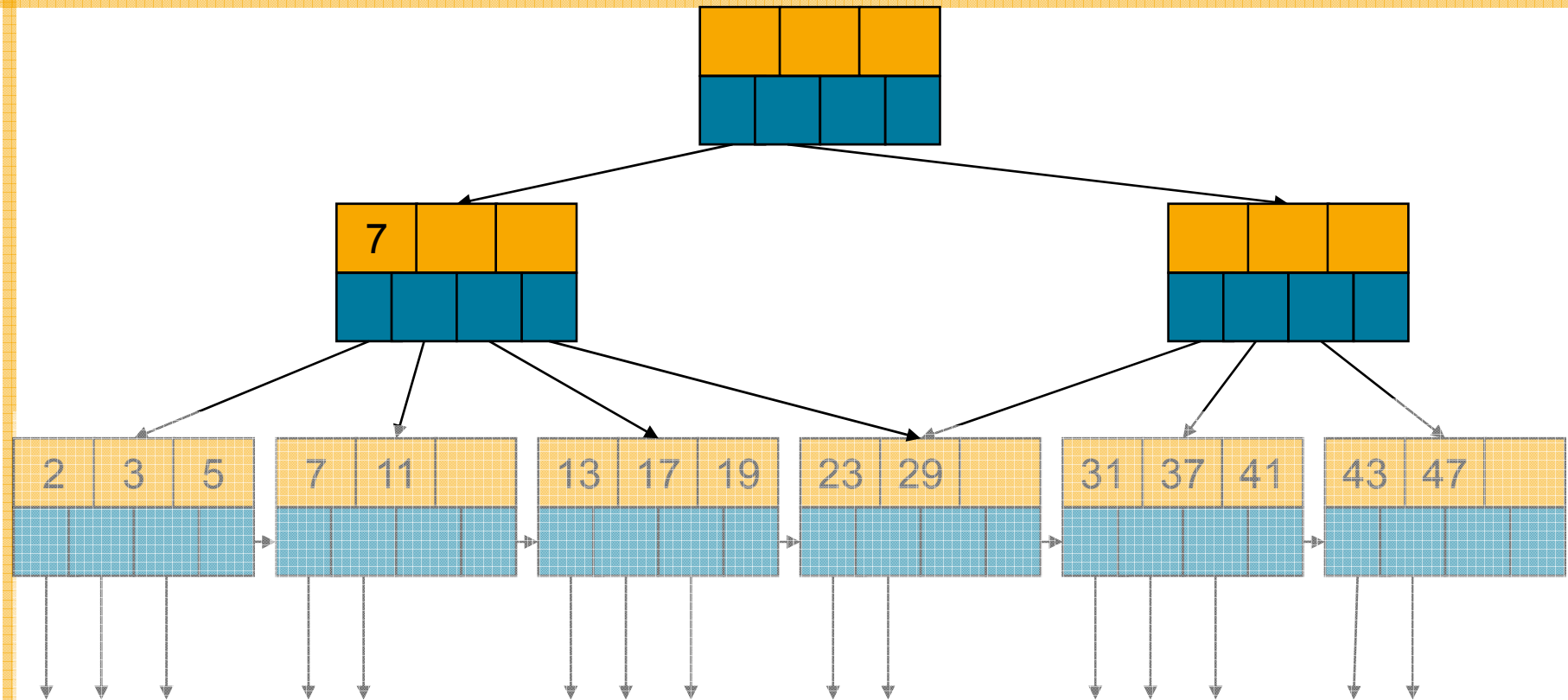
74

- Einfügevarianten
 - Einfügen von Daten, die bereits mit existierendem B-Baum indiziert sind.
 - Erstellung eines neuen B-Baums auf existierenden Daten
 - ◇ Sukzessives Einfügen jedes Datensatzes ist ineffizient.
 - ◇ Immer wieder über die Wurzel suchen

- Besser: Vorsortierung
 - Schritt 1: Erzeuge Schlüssel-Pointer Paare für alle Blöcke
 - Schritt 2: Sortierung der Paare nach Suchschlüssel
 - Schritt 3: Füge nun sukzessive die Paare ein

Bulk-loading

75



Wieso ist dies effizient?

- Nur Zugriff auf Indexblöcke im Cache
- Kosten: Einmal alle Knoten ausschreiben

Geht es noch besser? Besserer Füllstand?

Bulk-Loading mit hohem Füllstand

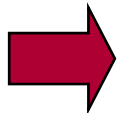
76

- Schritt 1: Erzeuge Schlüssel-Pointer Paare für alle Blöcke
- Schritt 2: Sortierung der Paare nach Suchschlüssel
- Schritt 3: Erzeuge volle B-Baum Blätter
- Schritt 4: Konstruiere innere Knoten aufgrund der Blätter

- Ergebnis: Perfekter Füllstand

- Anwendung: Insert-only Relationen

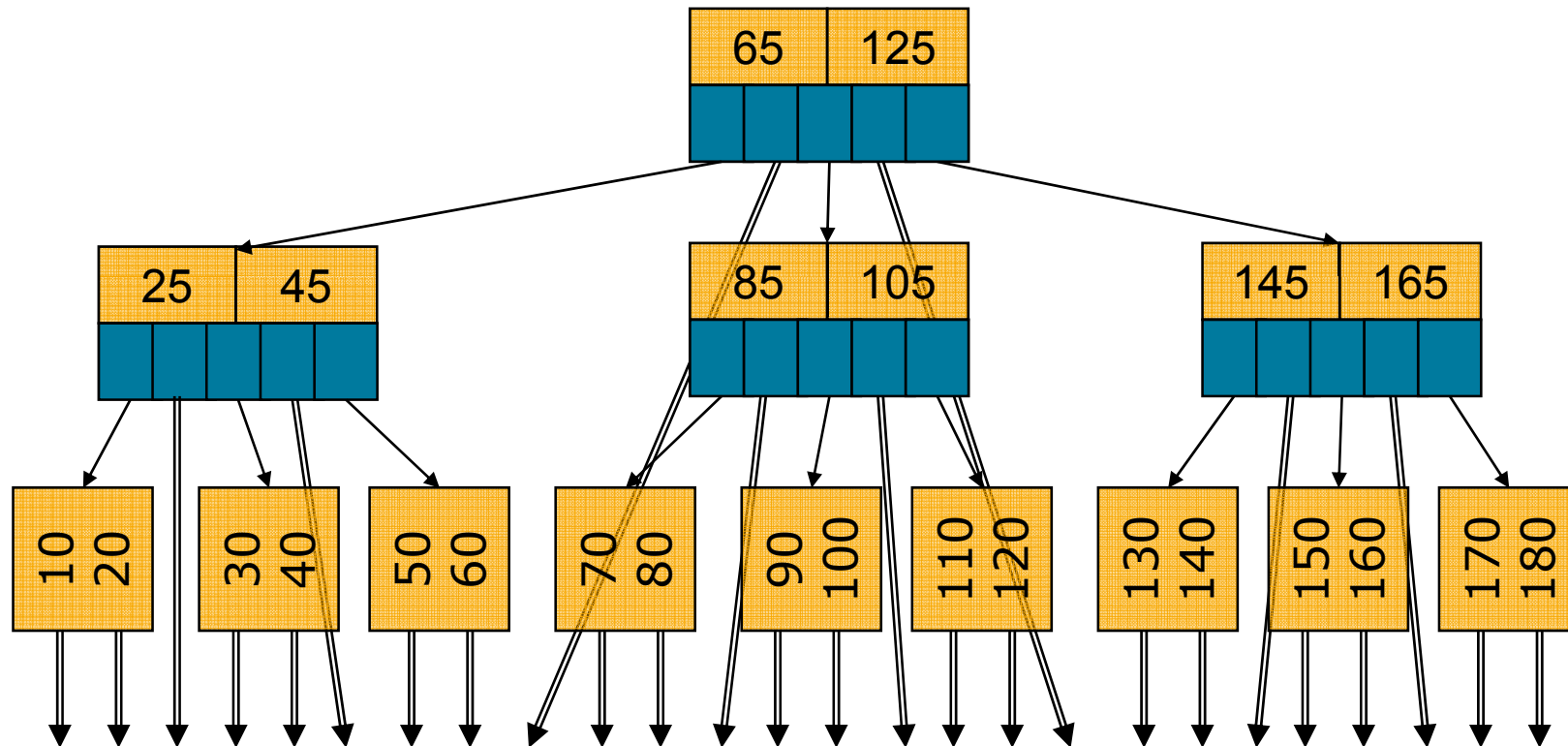
- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
 - Aufbau
 - Suche
 - Updates
 - Effizienz
 - Varianten
- Hash-Tabellen



B-Baum Varianten: B-Baum (ohne +)

78

- Bisher: B⁺-Baum
 - Pointer auf Datensätze nur in Blattknoten
- B-Baum
 - Pointer auf Datensätze in allen Knoten



B-Baum Varianten: B-Baum (ohne +)

79

- Vorteil: Durchschnittlich schnellere Suche als in B⁺-Bäumen
- Nachteil: Blätter und innere Knoten haben unterschiedliche Größe
 - Verwaltung schwieriger
- Nachteil: Weniger buschig; Platz wird für Pointer auf Datensätze „verschwendet“
- Nachteil: Löschen ist komplizierter
 - Löschung kann auch in inneren Knoten geschehen.
 - Schlüssel von einem Blatt muss nach oben wandern.

B-Baum Varianten: B^{*}-Baum

80

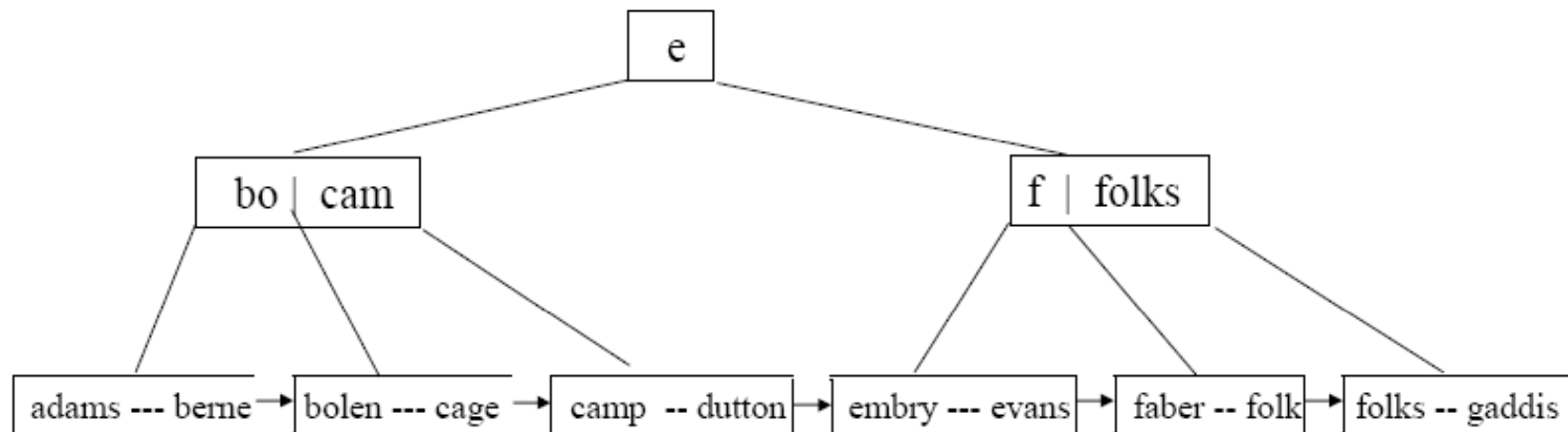
- B^{*}-und B[#]-Baum
 - Bei Überlauf werden nicht gleich zwei halb-leere Knoten erzeugt.
 - Stattdessen Neuverteilung über alle Blätter
 - Falls nicht möglich, erzeuge 3 neue Blätter aus 2 alten Blättern
 - Dadurch bessere Speicherausnutzung: Mindestens 66%

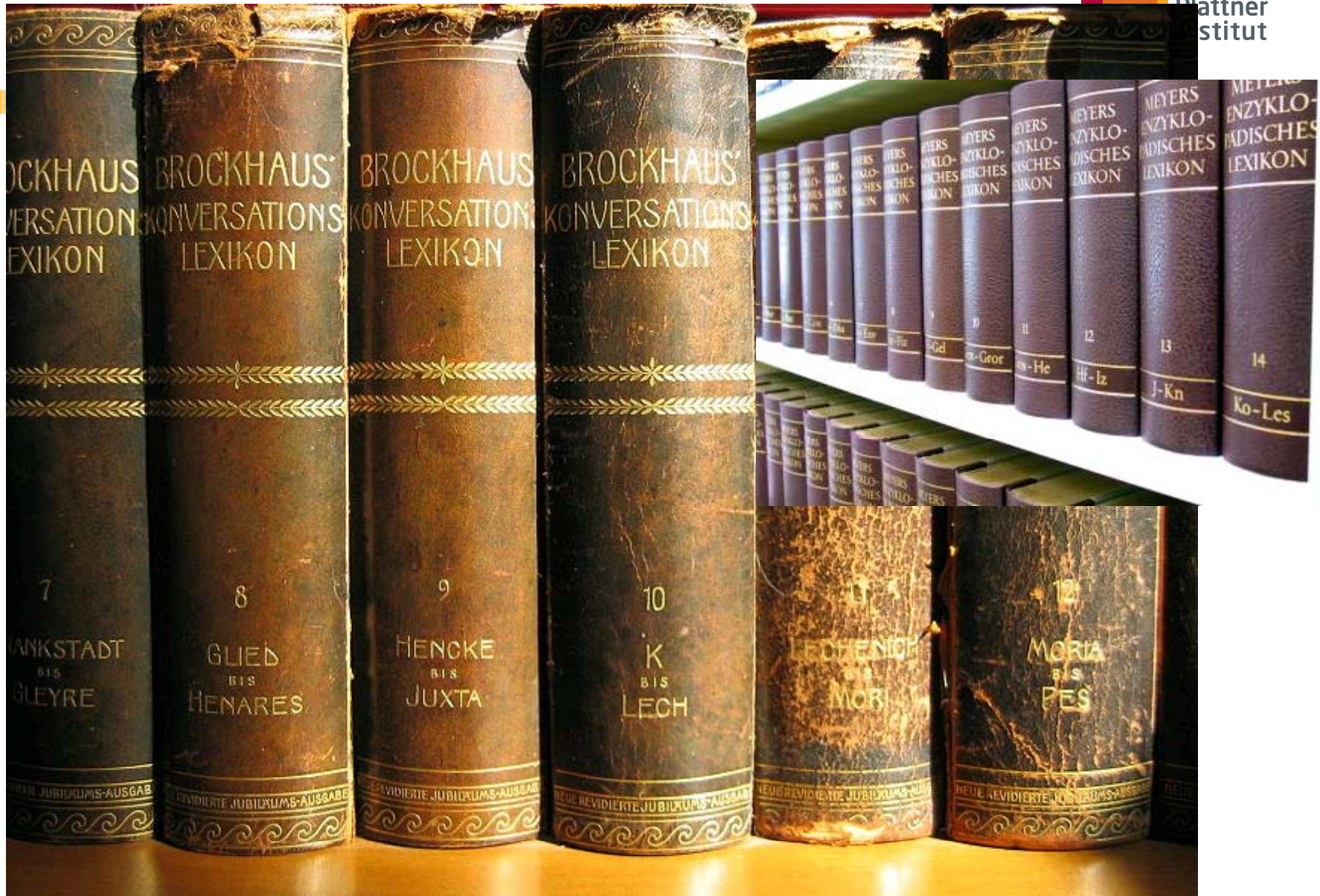
B-Baum Varianten: Präfix-B⁺-Baum

81

■ Präfix-B⁺-Baum

- Falls Suchschlüssel ein String ist => Hoher Speicherbedarf
- Besser: Speichere nur „Trennwert“
- Was trennt „Korn“ von „Loeser“?
- Am besten: Kleinster Trennwert
 - ◇ „L“
 - ◇ Präfix von „Loeser“

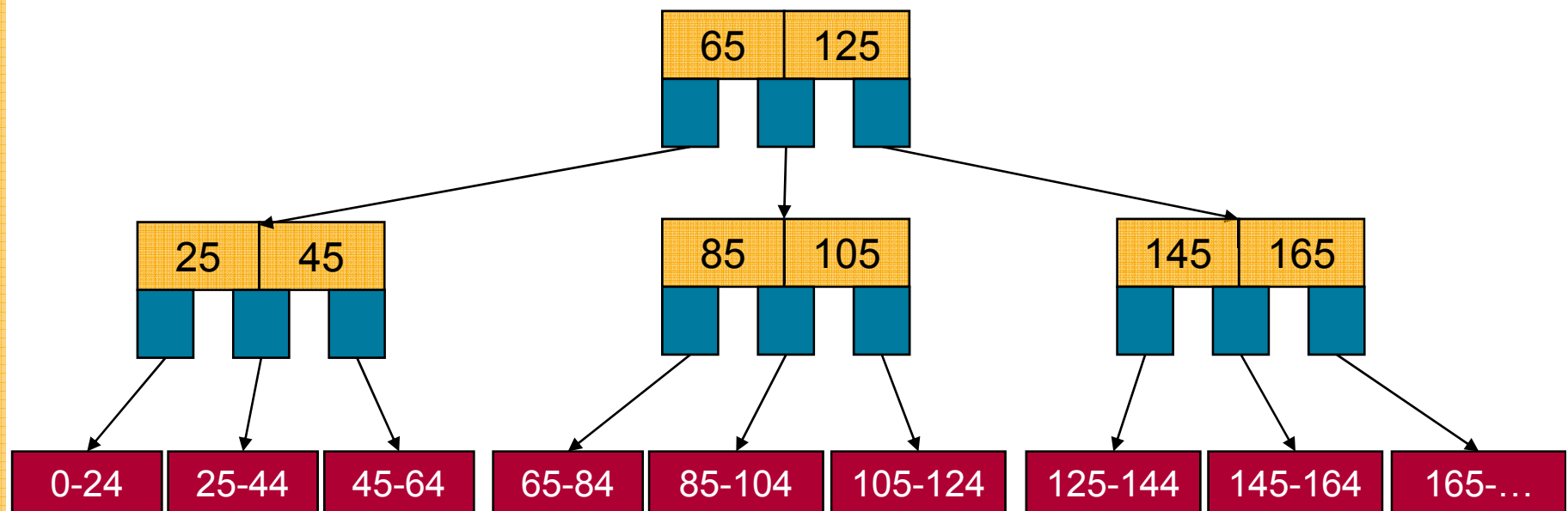




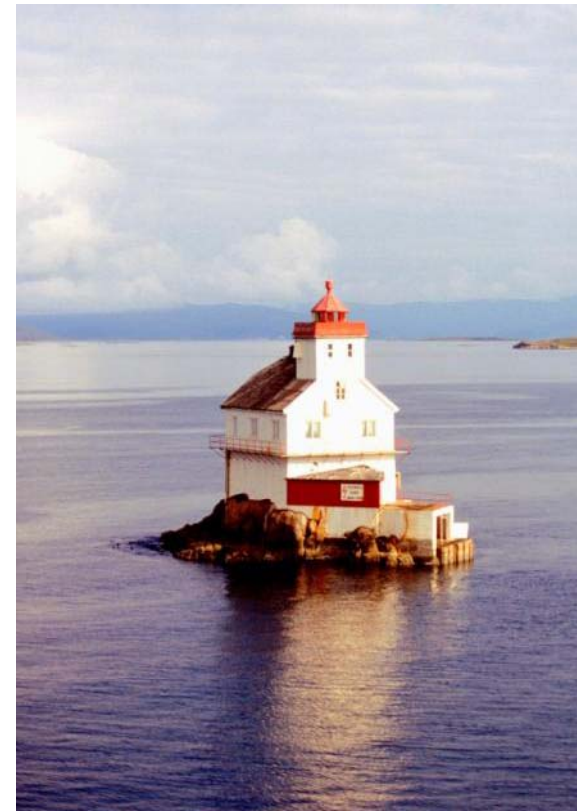
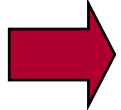
B⁺-Bäume für BLOBs

83

- Idee: Suchschlüssel repräsentiert Offsets im BLOB anstatt Suchschlüssel
 - Blätter zeigen auf Datenseiten des BLOBs
 - *Positional B-Tree*



- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
- Hash-Tabellen
 - Allgemeine Hashtabellen
 - Erweiterbare Hashtabellen
 - Lineare Hashtabellen



Hashtabellen – Grundprinzip

85

- Hashfunktion
 - Input: Suchschlüssel K (Hash-Schlüssel)
 - Output: Integer zwischen 0 und B-1
 - ◇ B = Anzahl Buckets
 - Oft z.B. $\text{MOD}(K/B)$
 - ◇ Bei Strings: Weise jedem Buchstaben Integer zu und summiere diese.
- Bucketarray
 - Array aus B Headern für B verkettete Listen

Hashtabellen auf Festplatten

86

- Bisher (im Studium): Hashtabellen im Hauptspeicher
- DBMS: Blöcke statt Pointer auf Listen
 - Datensätze, die auf einen gemeinsamen Wert gehashed werden, werden in dem entsprechenden Block gespeichert.
 - Overflowblocks können ergänzt werden.
- Annahme: Zuordnung eines Hashwerts zur Speicheradresse eines Blocks möglich
 - Z.B.: Offset

B = 4
2 Datensätze pro Block

0	D	

1	E	
	C	

2	B	

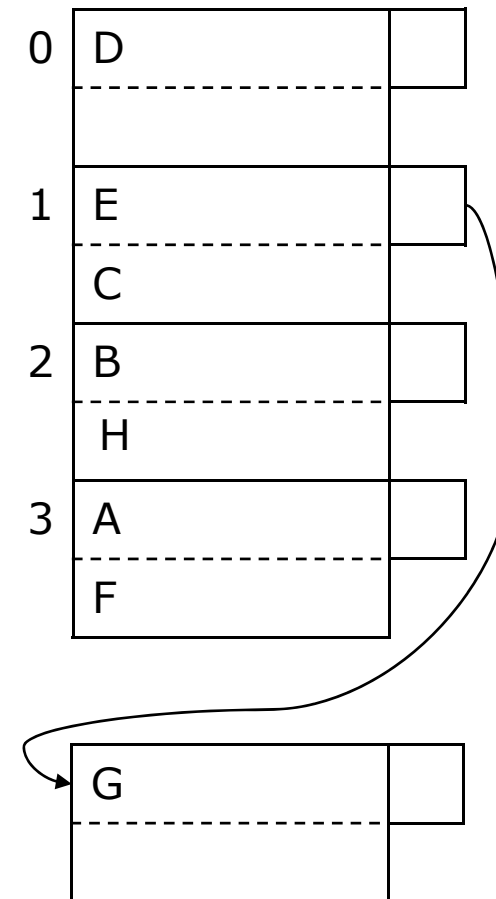
3	A	
	F	

Pointer für
Overflowblocks

Einfügen in Hashtabelle

87

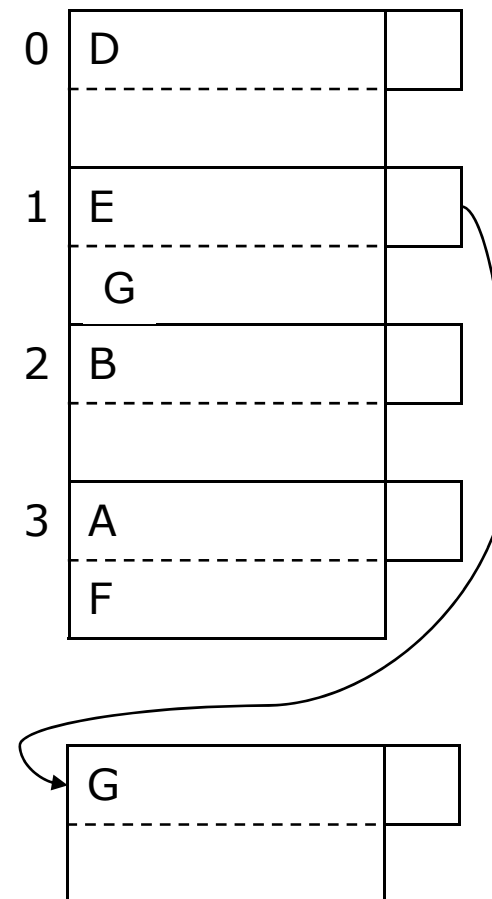
1. Berechne Hashwert
2. Falls Platz, füge Datensatz in entsprechenden Block ein
 - Sortierung in Bloc egal. Warum?
 - Oder in einen Overflowblock mit Platz
3. Falls kein Platz, erzeuge Overflowblock und füge dort ein
 - Beispiel: Füge H ein
 - ◇ $h(H) = 2$
 - Beispiel: Füge G ein
 - ◇ $h(G) = 1$



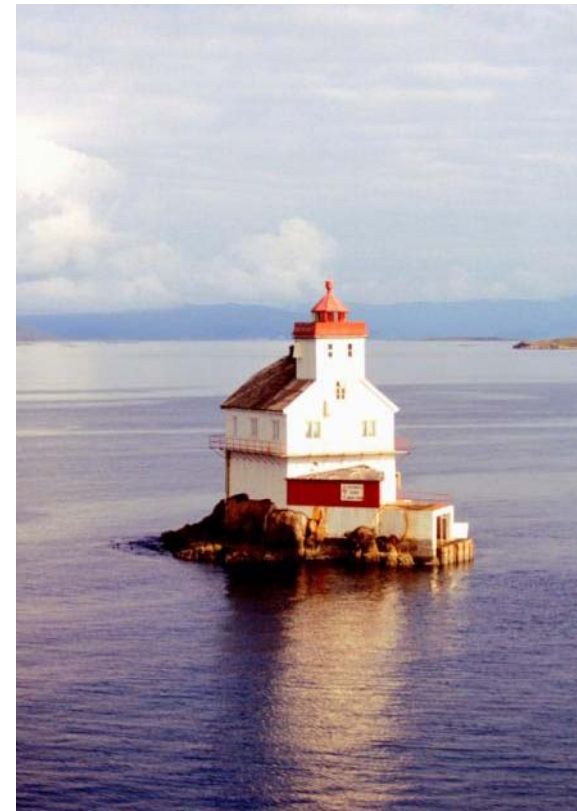
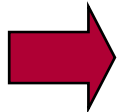
Löschen aus Hashtabellen

88

1. Suche Bucket (+ Overflowblocks)
2. Suche Datensatz / Datensätze
3. Lösche sie
4. Gegebenenfalls Reorganisation und Entfernung von Overflowblocks
 - Gefahr der Oszillation
 - Beispiel: Lösche C
 - G bewegen



- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
- Hash-Tabellen
 - Allgemeine Hashtabellen
 - Erweiterbare Hashtabellen
 - Lineare Hashtabellen



Effizienz statischer Hashtabellen

90

- Idealerweise: Pro Bucket nur ein Block
 - Zugriffszeit lesen: 1 I/O
 - Zugriffszeit Einfügen / Löschen: 2 I/O
 - Viel besser als Dicht- oder Dünnbesetzte Indizes
 - Besser als B-Bäume
 - ◇ Aber Nachteil: Bereichsanfragen
 - ◇ => Kein sequentielles Lesen von Disk
- Weiteres Problem: B wird einmalig festgelegt
 - Lange Listen von Overflowblöcken
 - „Statische Hashtabellen“
- Lösung: „Dynamische Hashtabellen“
 - Hashtabellen wachsen
 - $B \approx \text{Anzahl Datensätze} / \text{Datensätze pro Block}$
 - ◇ => Ca. 1 Block pro Bucket

Erweiterbare Hashtabellen

91

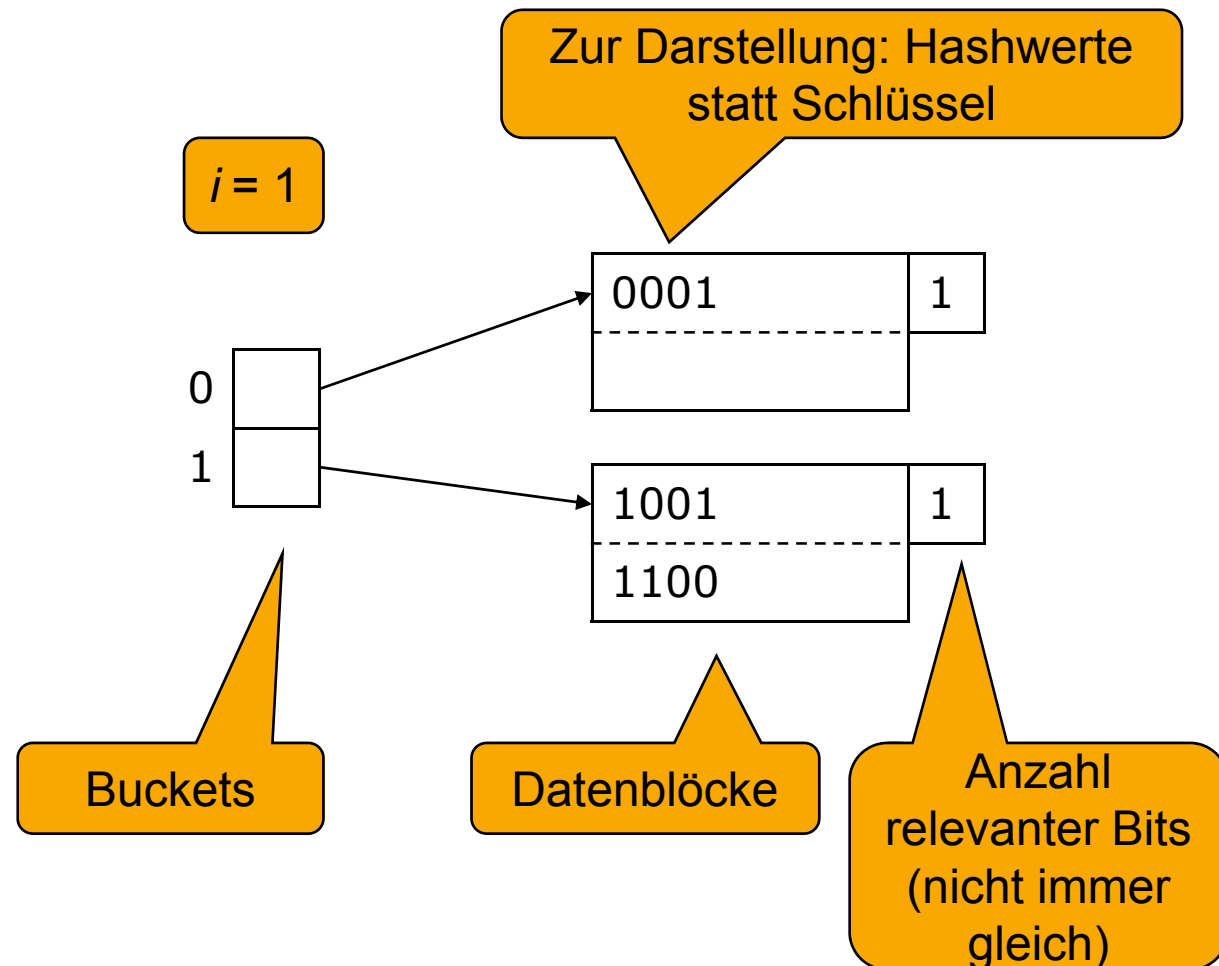
Neuerungen

- Indirektion
 - Bucket besteht aus Pointerarray statt Datenblock(s)
- Wachstum
 - Größe des Pointerarrays verdoppelt sich bei Bedarf
- Sharing
 - Buckets können sich Datenblöcke teilen...
 - ... wenn Platz ist.
- Hashfunktion
 - Berechnet (zu) großes Bitarray (k bit; z.B. $k = 32$)
 - Bucketarray verwendet nur die ersten i Bits ($i \leq k$)
 - $\Rightarrow 2^i$ buckets

Erweiterbare Hashtabellen

92

■ $k = 4$



Einfügen in erweiterbare Hashtabellen

93

- Ähnlich wie normale Hashtabelle
 - Berechne $h(K)$; wähle die ersten i Bits
 - Suche Eintrag im Bucketarray und lade entsprechenden Block
- Falls Platz: Füge neuen Datensatz ein.
- Falls kein Platz und $j < i$ (j ist aktuelle Bitanzahl des Blocks)
 - Spalte Block entzwei (*split*)
 - Verteile Datensätze gemäß des $(j+1)$ ten Bits
 - ◇ Falls 0: Verbleib
 - ◇ Falls 1: Verschieben in neuen Block
 - Setze $j+1$ als neue Bitanzahl
 - Pointer im Bucketarray aktualisieren
 - Was wäre Pech?
 - ◇ Alle Datensätze landen wieder im gleichen Block. Dann wieder $j++$.

Einfügen in erweiterbare Hashtabellen

94

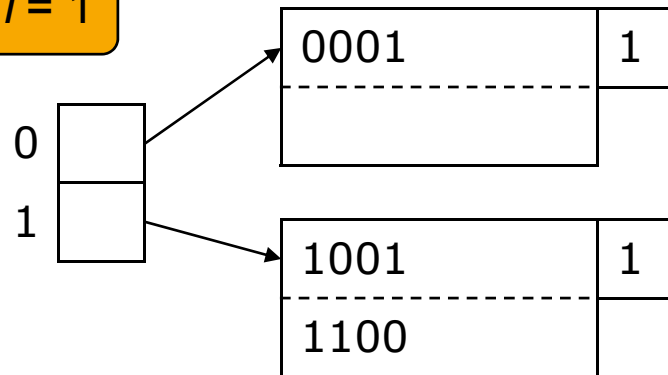
- Falls kein Platz und $j = i$ (j ist aktuelle Bitanzahl des Buckets)
 - $i++ \Rightarrow$ Länge des Bucketarrays verdoppelt sich
 - Datenblöcke bleiben unverändert
 - Zwei neue Pointer zeigen zunächst auf gleichen alten Block
 - Dann: Spalte relevanten Block entzwei (*split*)
 - ◇ Weiter wie zuvor (denn nun $j < i$)

Einfügen in erweiterbare Hashtabellen

95

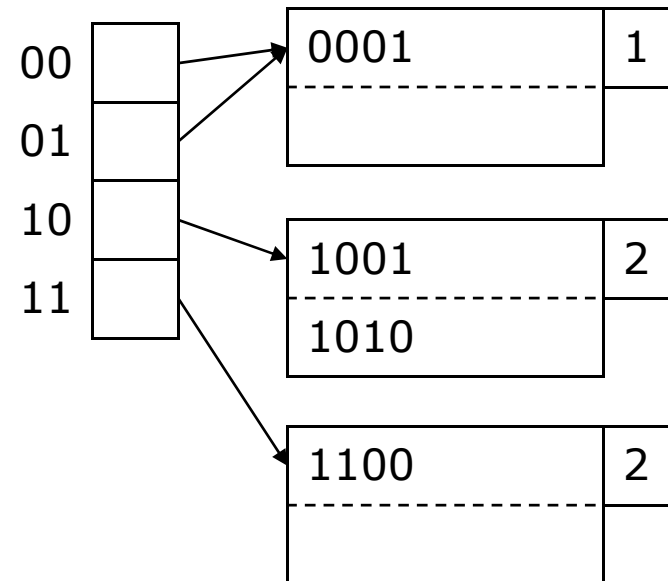
Füge ein: Datensatz mit Hashwert 1010

$i = 1$



Block #2
kein Platz => Split
 $j = i = 1 \Rightarrow i++$

$i = 2$

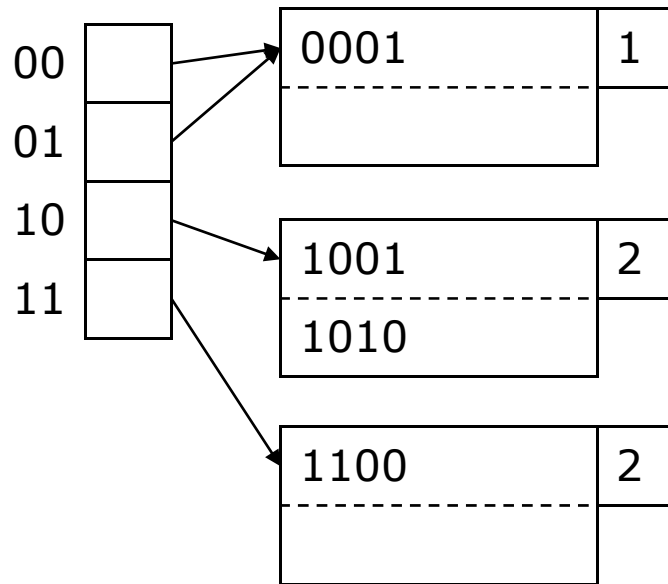


Einfügen in erweiterbare Hashtabellen

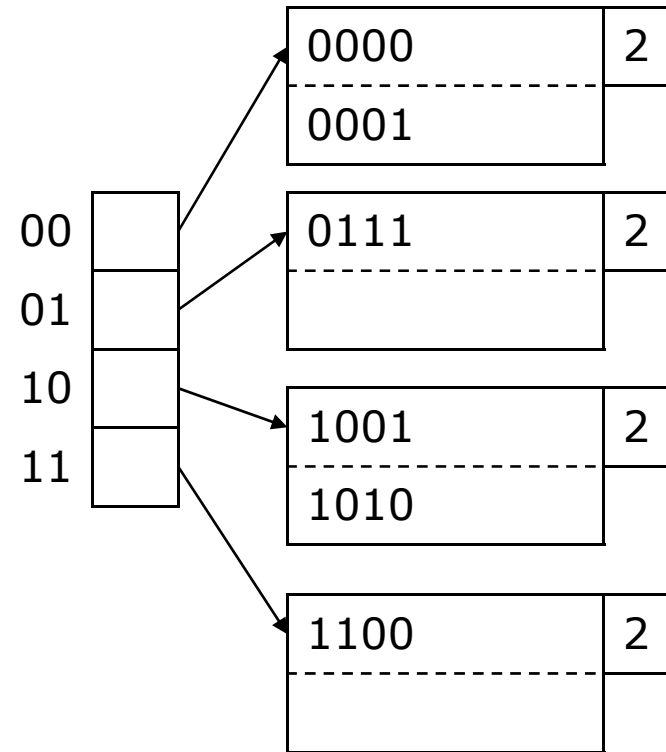
96

Füge ein: Datensätze mit Hashwerten 0000 und 0111

$i = 2$



Block #1: $j < i$

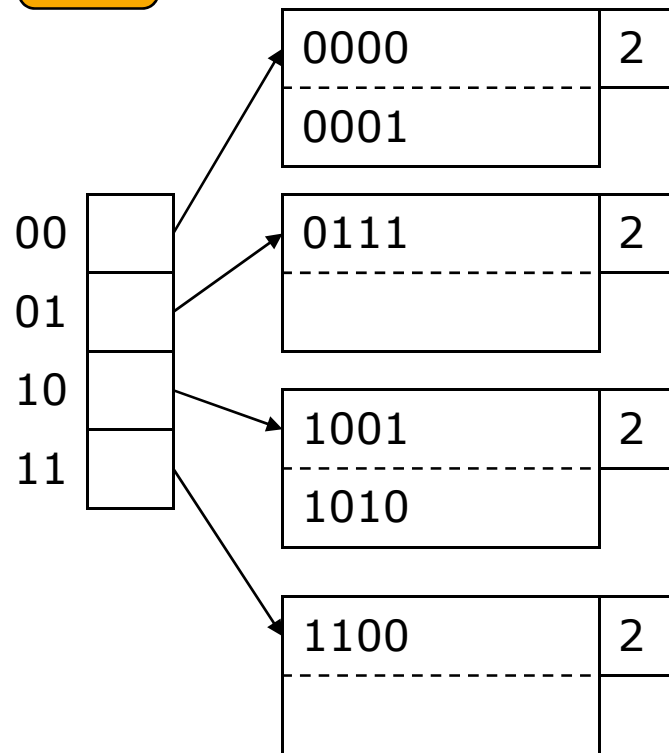


Einfügen in erweiterbare Hashtabellen

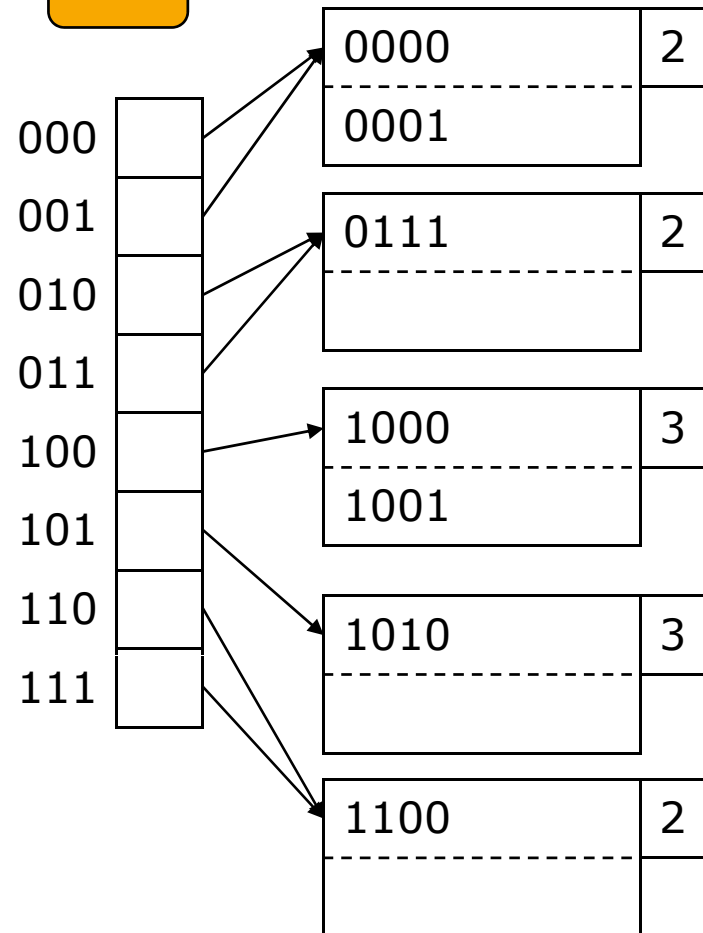
97

Füge ein: Datensatz mit Hashwert 1000

$i = 2$



$i = 3$

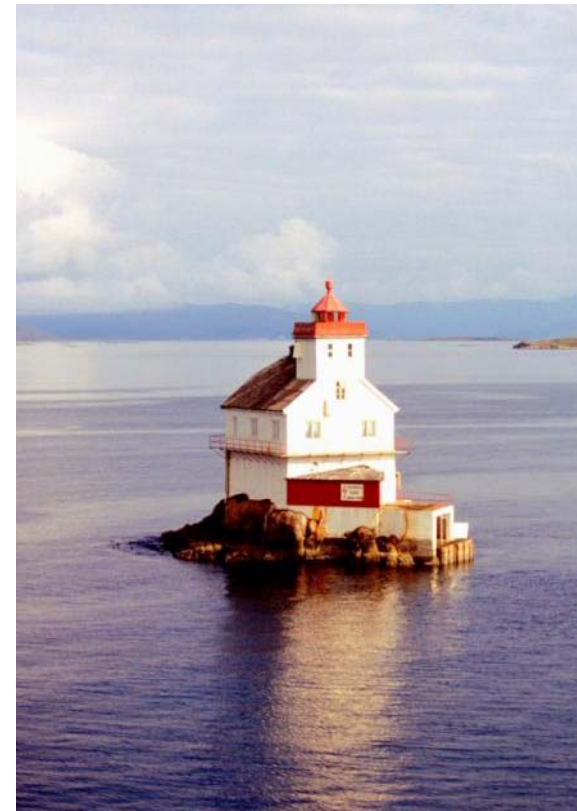
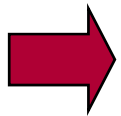


Analyse erweiterbarer Hashtabellen

98

- Vorteile
 - Bei Suche: Nie mehr als ein Datenblock betrachten
 - ◇ Keine Overflow Blocks
 - Bucketarray eventl. im Hauptspeicher
- Nachteile
 - Bei Verdopplung: Viel Arbeit
 - ◇ => Ab und zu dauert ein Einfügen sehr lang
 - Bucketarray wächst schnell
 - ◇ Passt eventuell nicht mehr in Hauptspeicher
 - Platzverschwendung bei wenigen Datensätzen pro Block
 - ◇ Beispiel: 2 Datensätze pro Block
 - ◇ Datensatz 1: 000000000000000000000001
 - ◇ Datensatz 2: 000000000000000000000010
 - ◇ Datensatz 3: 000000000000000000000011
 - ◇ => $i = 20$ (also $2^{20} = 1\text{Mio}$ buckets)

- Indizes auf sequenziellen Dateien
- Sekundärindizes auf nicht-sequenziellen Dateien
- B-Bäume
- Hash-Tabellen
 - Allgemeine Hashtabellen
 - Erweiterbare Hashtabellen
 - Lineare Hashtabellen



Lineare Hashtabellen

100

Idee: Anzahl Buckets wächst nur langsam.

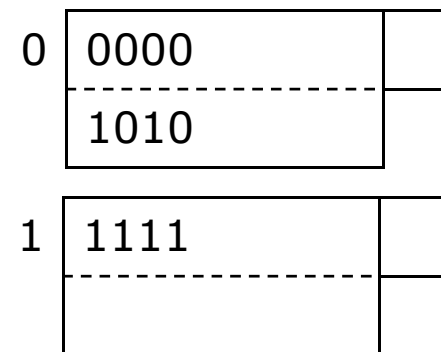
- Anzahl n der Buckets so gewählt, dass Datenblocks zu ca. 80% gefüllt sind.
- Overflow Blocks sind zugelassen, aber
 - durchschnittliche Anzahl Overflow Blocks pro Block: $\ll 1$
- $\log_2 n$ Bits zur Identifizierung der Buckets
 - Wähle die jeweils letzten Bits des Hashwerts

$i = 1$ (Anzahl relevanter Bits)

$n = 2$ (Anzahl Buckets)

$r = 3$ (Anzahl Datensätze)

Wähle n so, dass $r \leq 1,7 \cdot n$
(bei Blockgröße 2)



Lineare Hashtabellen – Einfügen

101

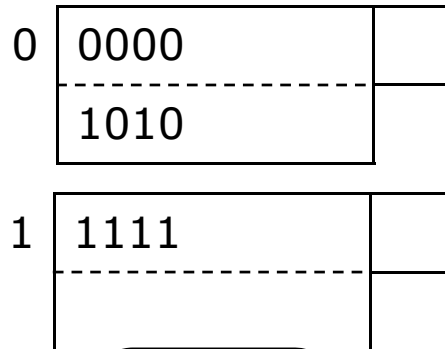
- Berechne $h(K)$
- Betrachte letzte i Bits; interpretiere als Integer m
- Falls $m < n$: Füge Datensatz in Bucket m ein.
- Falls $m \geq n$:
 - \Rightarrow Bucket m existiert noch nicht
- Füge Datensatz in Bucket $m - 2^{i-1}$ ein.
 - D.h. erstes Bit des Schlüssels wird zu 0 (vorher 1).
- Falls kein Platz: Erzeuge Overflow Block
- Berechne r/n
 - Falls zu hoch (z.B. $\geq 1,7$): Erzeuge neuen Bucket
 - Neuer Bucket hat nichts mit betroffenem Bucket zu tun.
 - Falls neuer Bucket $1xxx$ ist, *split* auf Bucket $0xxx$, je nach den xxx -Werten der Datensätze
- Falls nun $n > 2^i$: $i++$
 - D.h. alle Bitsequenzen erhalten eine 0 am Anfang
 - Physisch ändert sich nichts

Lineare Hashtabellen – Einfügen

102

$i = 1$
 $n = 2$
 $r = 3$

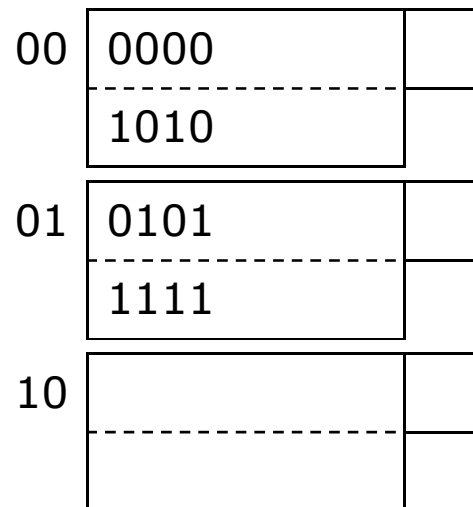
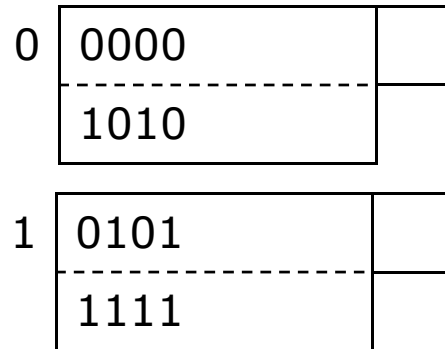
Füge ein: 0101



Lösung:
 $i = 2$
 $n = 3$
 $r = 4$

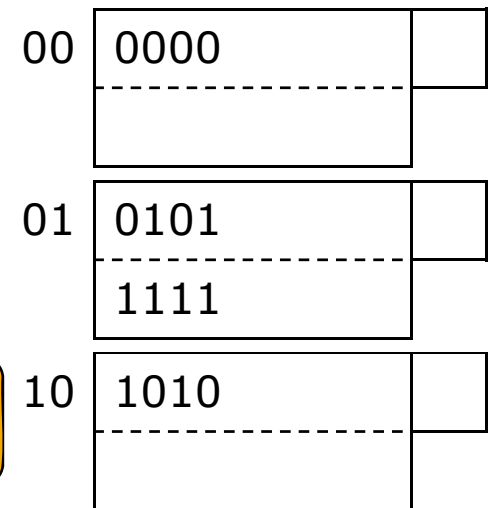
$i = 1$
 $n = 2$
 $r = 4$

Problem: $4/2 > 1,7$



Split 00

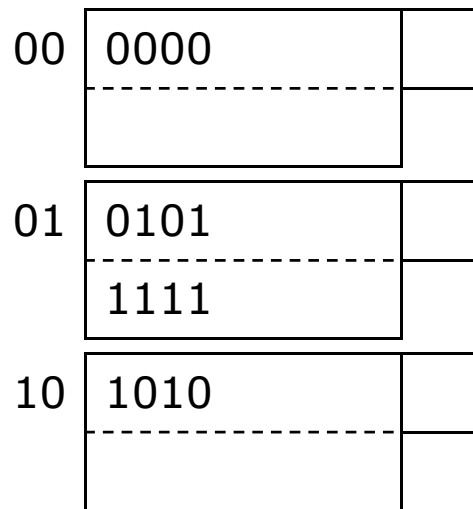
Unterscheidet sich von neuem Block nur um erstens Bit.



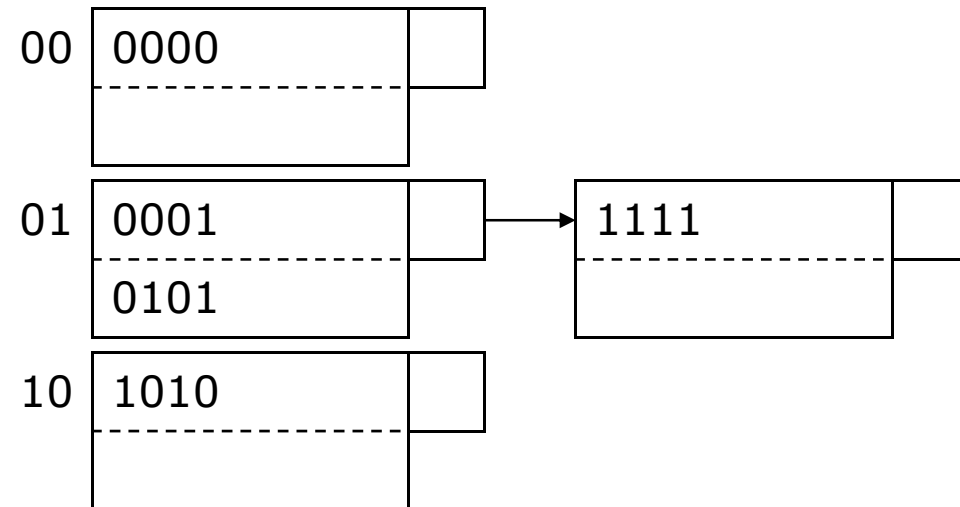
Lineare Hashtabellen – Einfügen

103

$i = 2$
 $n = 3$
 $r = 4$ Füge ein: 0001



$i = 2$
 $n = 3$
 $r = 5$



Lineare Hashtabellen – Einfügen

104

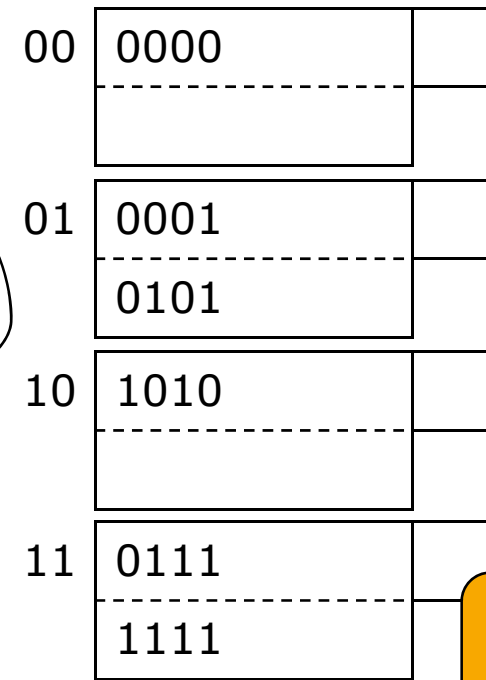
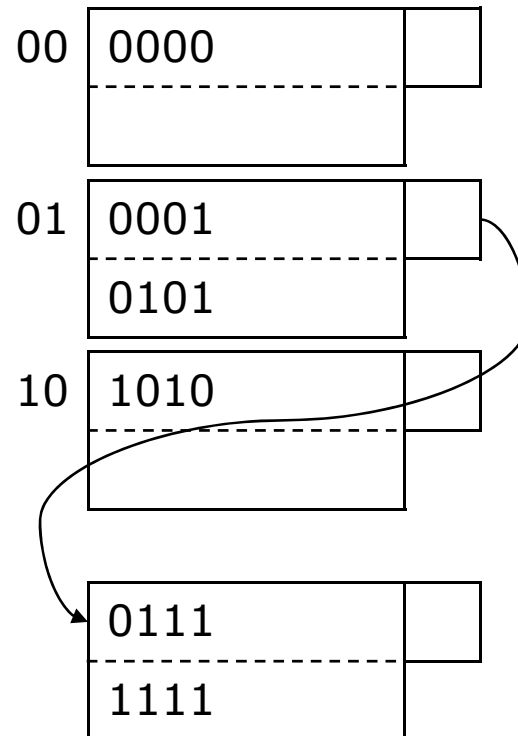
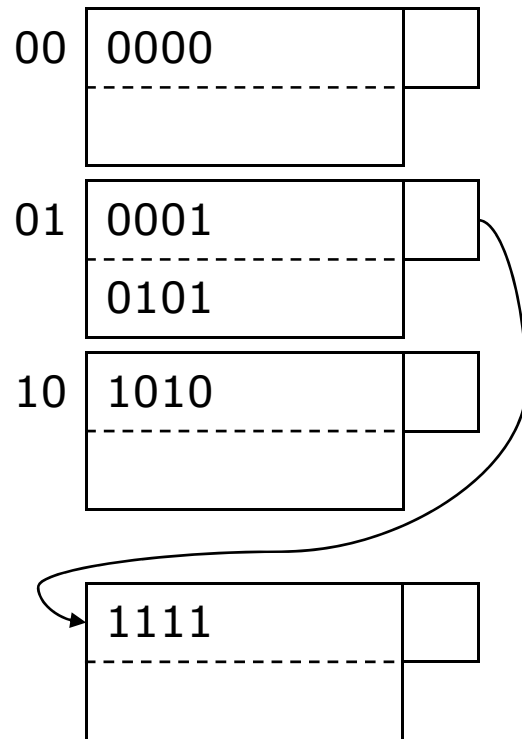
Füge ein: 0111

$i = 2$
 $n = 3$
 $r = 5$

Bucket 11 existiert noch nicht.
Deshalb wähle Bucket 01.

$i = 2$
 $n = 3$
 $r = 6$

Problem: $6/3 > 1,7$
 \Rightarrow Neuer Block 11
 \Rightarrow Split Block 01



$i = 2$
 $n = 4$
 $r = 6$

Hashing vs. B-Baum

105

- Hashing effizient zur Feststellung der Existenz eines Wertes
- B-Baum effizient für Bereichsanfragen

- Indexerzeugung in SQL:
 - `CREATE [UNIQUE] INDEX indexname
ON table (column [ASC | DESC]
[, column [ASC | DESC] ...]`
 - ◇ UNIQUE erlaubt NULL Werte
 - Keine Angabe über Art des Schlüssels
 - Keine Angabe über Parameter
 - Manchmal Hersteller-spezifische Syntax für Parameter

Zusammenfassung

106

- Sequentielle Dateien
 - Sortierung auf Festplatte
- Dicht-besetzte Indizes
 - Schlüssel-Pointer-Paar für jeden Datensatz
- Dünn-besetzte Indizes
 - Schlüssel-Pointer-Paar für jeden Block (erster Datensatz in Block)
- Mehrstufige Indizes
 - Indizes auf Indizes
- Einfügen von Daten
 - Overflow vs. Struktur erweitern
- Sekundärindizes
 - Auf nicht-sortierten Attributen
 - Immer dicht-besetzt
- Invertierte Indizes
 - Texte/Blobs
- B-Bäume
 - Mehrstufige mit günstigen Eigenschaften
- Hash-Tabellen
 - Hashwert eines Datensatzes zeigt auf Adresse des Blocks.
- Erweiterbares und lineares Hashing
 - Wachstum der Hashtabellen