

9. Mehrdimensionale Zugriffspfade

- **Ziele**

- Entwurfsprinzipien für Zugriffspfade auf die Sätze einer Tabelle, bei denen mehrere Suchkriterium symmetrisch unterstützt werden
- Indexierung von Punktobjekten und räumlich ausgedehnten Objekten

- **Klassifikation der Anfragen**

- **Grundprobleme**

- Organisation räumlicher Daten
- Erhaltung der Topologie (Clusterbildung)
- Objektdarstellung

- **Organisation der Datensätze**

- Quad-Tree
- Multi-Key-Hashing

- **Organisation des umgebenden Datenraums**

- k-d-Baum
- Grid-File

- **Zugriffspfade für ausgedehnte räumliche Objekte**

- R-Baum
- R^+ -Baum

- **Eindimensionale Intervalle**

Klassifikation der Anfragetypen

• Definitionen:

- Eine Datei ist eine Sammlung von N Sätzen des Typs $R = (A_1, \dots, A_n)$, wobei jeder Satz ein Punktobjekt durch ein geordnetes n -Tupel $t = (a_1, a_2, \dots, a_n)$ von Werten darstellt. Die Attribute A_1, \dots, A_k ($k \leq n$) seien Schlüssel.
- Eine Anfrage Q spezifiziert einige Bedingungen, die von den Schlüsselwerten der Sätze in der Treffermenge erfüllt sein müssen.
- **Schnittbildende Anfragen** (*intersection queries*): sich qualifizierende Objekte überlappen mit dem Anfragebereich
- **Enthaltenseins- oder Umschließungsanfragen** (*containment, enclosure queries*): sich qualifizierende Objekte sind ganz im Anfragebereich enthalten oder enthalten den Anfragebereich vollständig

• Klassifikation der schnittbildenden Anfragen

1. **Exakte Anfrage** (*exact match query*):

spezifiziert für jeden Schlüssel einen Wert

$$Q = (A_1 = a_1) \wedge (A_2 = a_2) \wedge \dots \wedge (A_k = a_k)$$

2. **Partielle Anfrage** (*partial match query*):

spezifiziert $s < k$ Schlüsselwerte

$$Q = (A_{i_1} = a_{i_1}) \wedge (A_{i_2} = a_{i_2}) \wedge \dots \wedge (A_{i_s} = a_{i_s})$$

$$\text{mit } 1 \leq s < k \text{ und } 1 \leq i_1 < i_2 < \dots < i_s \leq k$$

3. **Bereichsanfrage** (*range query*):

spezifiziert einen Bereich $r_i = [l_i \leq a_i \leq u_i]$ für jeden Schlüssel A_i

$$Q = (A_1 = r_1) \wedge \dots \wedge (A_k = r_k)$$

$$\equiv (A_1 \geq l_1) \wedge (A_1 \leq u_1) \wedge \dots \wedge (A_k \geq l_k) \wedge (A_k \leq u_k)$$

4. **Partielle Bereichsanfrage** (*partial range query*):

spezifiziert für $s < k$ Schlüssel einen Bereich

$$Q = (A_{i_1} = r_{i_1}) \wedge \dots \wedge (A_{i_s} = r_{i_s})$$

$$\text{mit } 1 \leq s < k \text{ und } 1 \leq i_1 < \dots < i_s \leq k \text{ und } r_{ij} = [l_{ij} \leq a_{ij} \leq u_{ij}], 1 \leq j \leq s$$

Klassifikation der Anfragetypen (2)

- **Allgemeine Bereichsanfrage**

- genauer Bereich $[l_i = a_i = u_i]$
- unendlicher Bereich $[-\infty \leq a_i \leq \infty]$

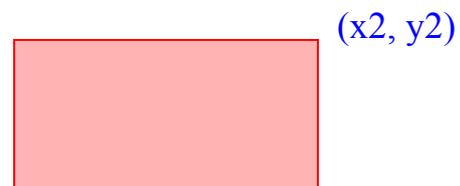
➔ Alle 4 Fragetypen der schnittbildenden Anfragen lassen sich als allg. Bereichsanfrage ausdrücken

- **Klassifikation der Enthaltenseinsanfragen**

- **Punktanfrage** (*point query*): Gegeben ist ein Punkt im Datenraum D. Finde alle Objekte, die ihn enthalten.
- **Gebietsanfrage** (*region query*): Gegeben ist ein Anfragegebiet. Finde alle Objekte, die es schneiden (es umschließen, in ihm enthalten sind).

- **Beispiele: Suche nach Rechtecken**

RECTANGLES (x_1, y_1, x_2, y_2)



(x_1, y_1)

a) Bestimmung aller Rechtecke, die den Punkt (2,5) enthalten

```
SELECT    x1, y1, x2, y2
FROM      RECTANGLES
WHERE     x1 <= 2 AND x2 >= 2 AND y1 <= 5 AND y2 >= 5
```

b) Bestimmung der Rechtecke mit Punkt (1,3) als Eckpunkt links unten

```
SELECT    x1, y1, x2, y2
FROM      RECTANGLES
WHERE     x1 = 1 AND y1 = 3
```

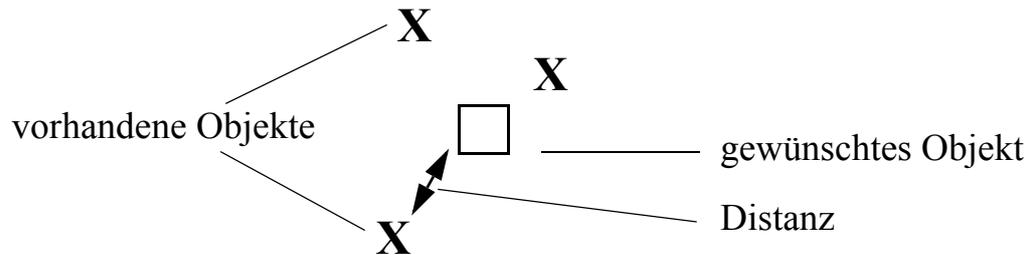
➔ Das sind alles **harmlose Fragen**, die sogar im Relationenmodell mühelos beantwortet werden können

Klassifikation der Anfragetypen (3)

- **Nächster-Nachbar-Anfragen** (*best match query, nearest neighbor query*)

- gewünschtes Objekt nicht vorhanden

➔ Frage nach **möglichst ähnlichen** Objekten



- "best" wird bestimmt über verschiedene Arten von Distanzfunktionen

- **Beispiele:**

- Objekt erfüllt nur 8 von 10 geforderten Eigenschaften
- Objekt ist durch Synonyme beschrieben

- **Bestimmung des nächsten Nachbarns:**

D = Distanzfunktion

B = Sammlung von Punkten im k-dimensionalen Raum

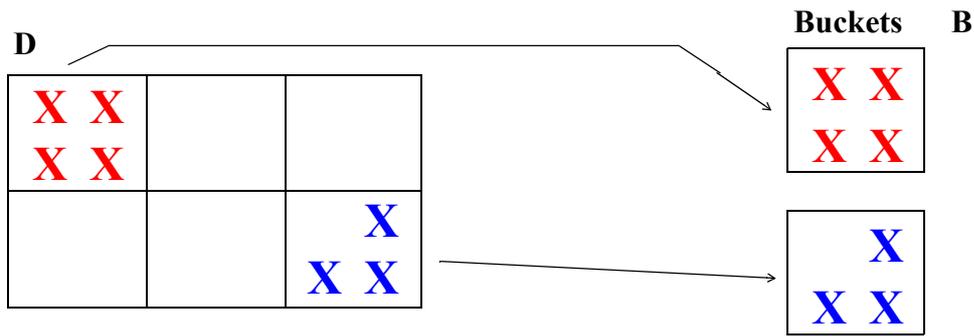
Gesucht: nächster Nachbar von p (in B)

Der nächste Nachbar ist q, wenn

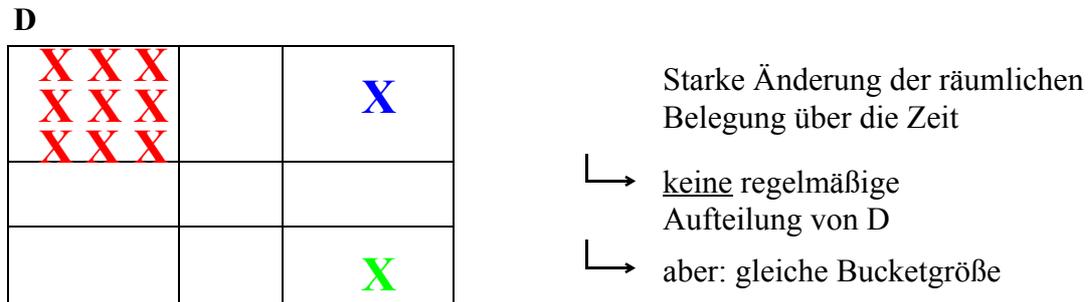
$$(\forall r \in B) \{r \neq q \Rightarrow [D(r, p) \geq D(q, p)]\}$$

Grundprobleme

1. Erhaltung der topologischen Struktur



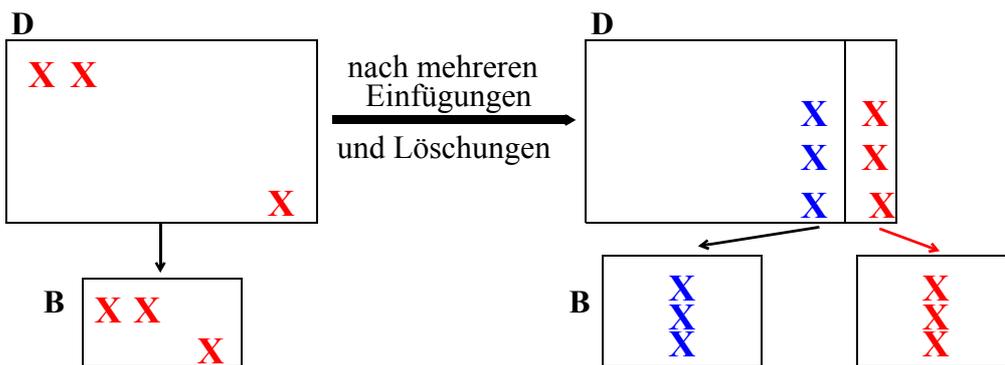
2. Stark variierende Dichte der Objekte



3. Objektdarstellung

- Punktoobjekte
- Objekte mit Ausdehnung

4. Dynamische Reorganisation



5. Balancierte Zugriffsstruktur

- beliebige Belegungen und Einfüge-/Löschreihenfolgen
- Garantie eines gleichförmigen Zugriffs

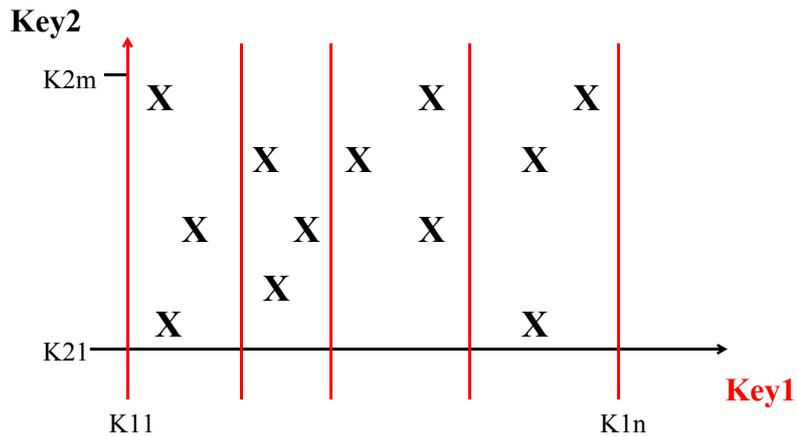
➔ **2 oder 3 Externspeicherzugriffe**

Nutzung eindimensionaler Zugriffspfade

- **Bisher:**

Indexierung (Invertierung) einer Dimension, z. B. B*-Baum

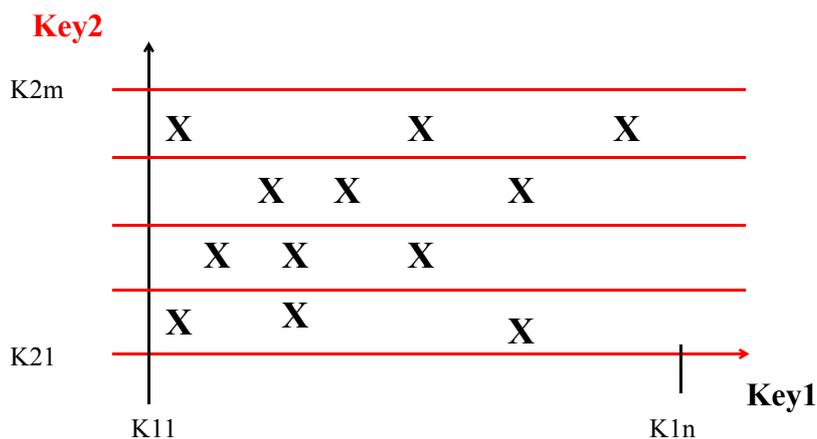
- **Zerlegungsprinzip des Schlüsselraumes beim B*-Baum (2-dim.)**



B*-Baum (Key1)

➔ Partitionierung des Raumes nach Werten von Key1

- **Zusätzlicher B*-Baum (Key2) möglich:**



B*-Baum (Key2)

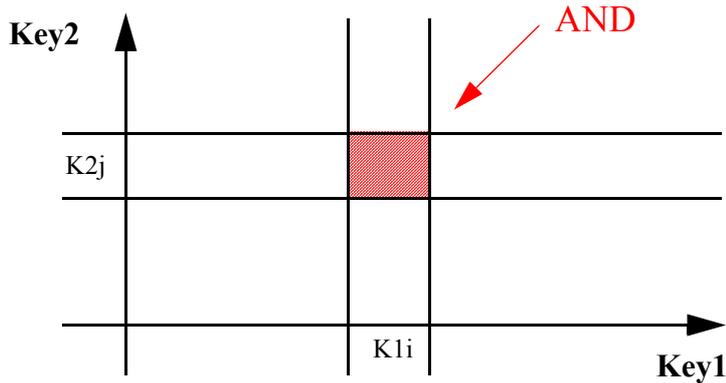
➔ Partitionierung des Raumes nach Werten von Key2

Mehrattributsuche

- Zugriff nach $(\text{Key1} = \text{K1i}) \left\{ \begin{array}{c} \text{OR} \\ - \\ \text{AND} \end{array} \right\} (\text{Key2} = \text{K2j})$

- Zeigerliste für K1i : aus B*-Baum (Key1)
- Zeigerliste für K2j : aus B*-Baum (Key2)

➔ Mischen + Zugriff auf Ergebnistupel



➔ große Zeigerlisten und Zwischenergebnisse!

- Simulation des mehrdimensionalen Zugriffs mit **einem B*-Baum?**

Idee:	Konkatenierte Schlüssel: Konkatenierte Werte:	Key1 K11 K11 K11 K12 K12 K12 K13 K1n	Key2 K21 K22 K2m K21 K22 K2m K21 K2m
--------------	--	--	---

Suchoperationen:

- $(\text{Key1} = \text{K1i}) \text{ AND } (\text{Key2} = \text{K2j})$?
- $\text{Key2} = \text{K2j}$?
- $\text{Key1} = \text{K1i}$?
- OR-Verknüpfung ?

Vergleich	Quad-Tree	Multi-Key Hashing	k-d-Baum	Grid-File	R-Baum
Dynamische Reorganisation					
Erhaltung der topologischen Struktur (Clusterbildung)					
Bestimmung der Nachbarschaft: Bereichsanfragen, best match					
Trennung von Zugriffsstruktur und Daten,					
Exakte Anfrage (balanc. Struktur)					
Objektdarstellung					

Quad-Tree (Quadranten-Baum)

- **Speicherungsstruktur**

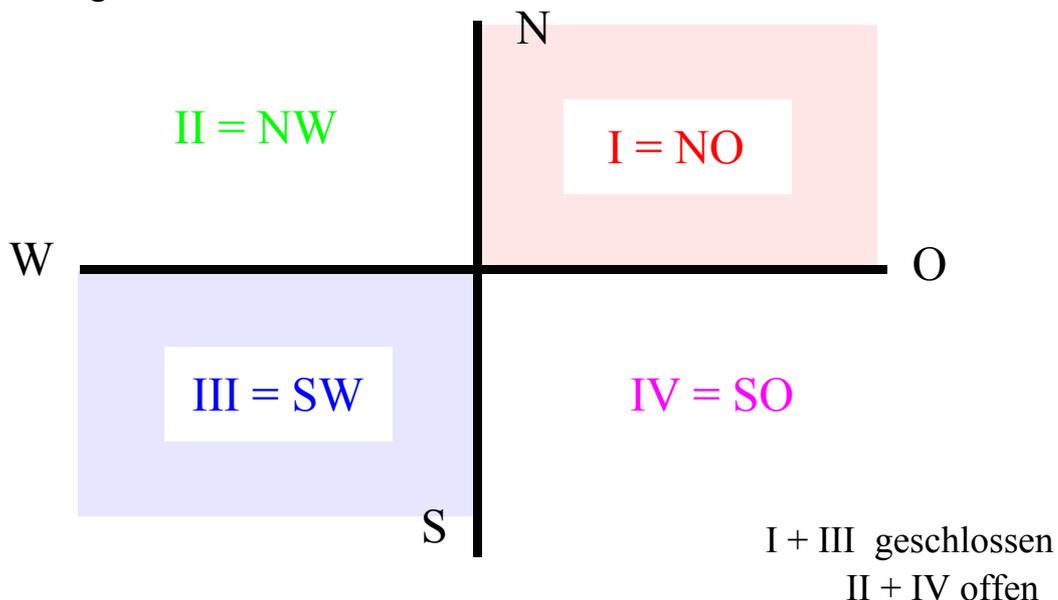
- für 2-dimensionalen Mehrattributzugriff
- Zerlegungsprinzip des Datenraumes D:
rekursive Partitionierung durch Quadranten

- **Realisierung als Generalisierung des Binärbaumes**

- jeder Knoten enthält einen Satz
- Außengrad eines Knotens: max. 4
- Wurzel teilt 2-dimensionalen Raum in 4 Quadranten auf
- rekursive Aufteilung jedes Quadranten durch Wurzel eines Unterbaumes
- i-ter Unterbaum eines Knotens enthält Punkte im i-ten Quadranten

- **Ziel: Berücksichtigung von Nachbarschaftsbeziehungen**

- Beispiel: geographische Daten mit Koordinaten x und y
- Aufteilung:

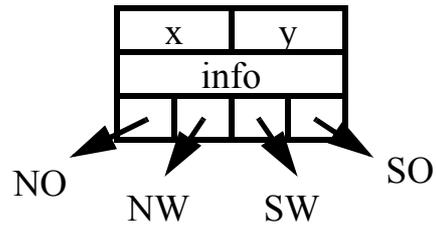


- **Verallgemeinerung:**

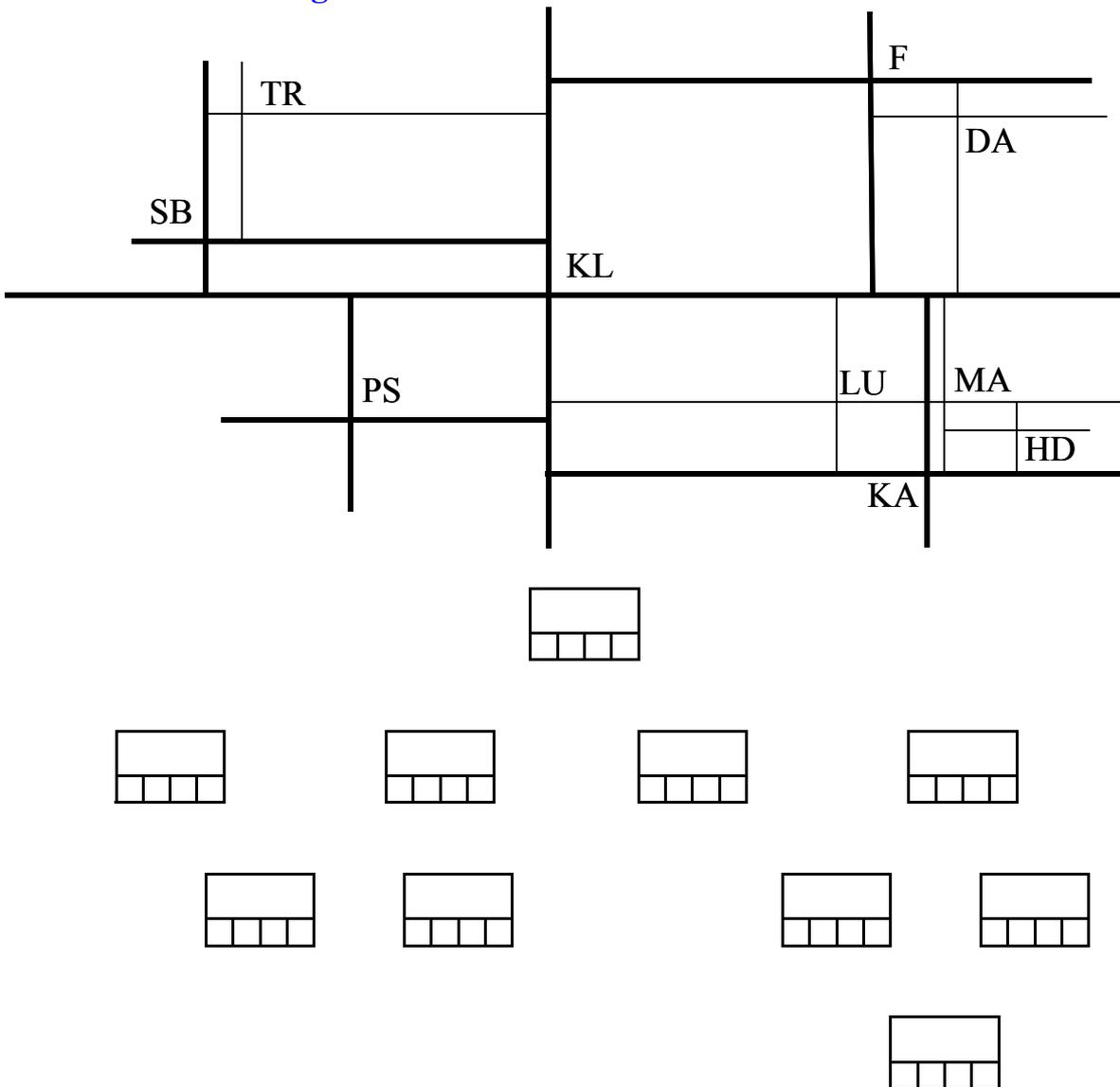
- k-dimensionaler Schlüssel \Rightarrow Außengrad jedes Knotens: 2^k
- k=3: Octtree, k=4: Hextree

Quad-Tree (2)

- Knotenformat**



- Räumliche Aufteilung**



- Eigenschaften**

- Baumstruktur abhängig von Einfügereihenfolge (unbalanciert)
- aufwendiges Löschen (Neueinfügen der Unterbäume)
- keine Abbildung auf Seiten

Multi-Key-Hashing (partitioned hashing)

- **Zerlegungsprinzip von D:**

Partitionierung durch Hashfunktionen in jeder Dimension:

wird realisiert durch Aufteilung der Bucketadresse in k Teile
(k = Anzahl der Schlüssel)

- **Für jeden Schlüssel i :**

eigene Hash-Funktion h_i , die Teil der Bucketadresse bestimmt

- Anzahl der Buckets sei 2^B (Adresse = Folge von B Bits)

- jede Hash-Funktion h_i liefert b_i Bits mit $B = \sum_{i=1}^k b_i$

- Satz $t = (a_1, a_2, \dots, a_k, \dots)$ gespeichert in Bucket mit

$$\text{Adr.} = h_1(a_1) | h_2(a_2) | \dots | h_k(a_k)$$

- **Vorteile:**

- kein Index, geringer Speicherplatzbedarf und Änderungsaufwand

- Exact-Match-Anfragen: Gesamtadresse bekannt

↳ Zugriff auf 1 Bucket

- Partial-Match-Anfrage: Eingrenzung des Suchraumes

($A_i = a_i$): Anzahl zu durchsuchender Buckets reduziert sich um 2^{b_i}

↳ $N_B = 2^B / 2^{b_i} = 2^{B-b_i}$

- **Nachteile / Probleme:**

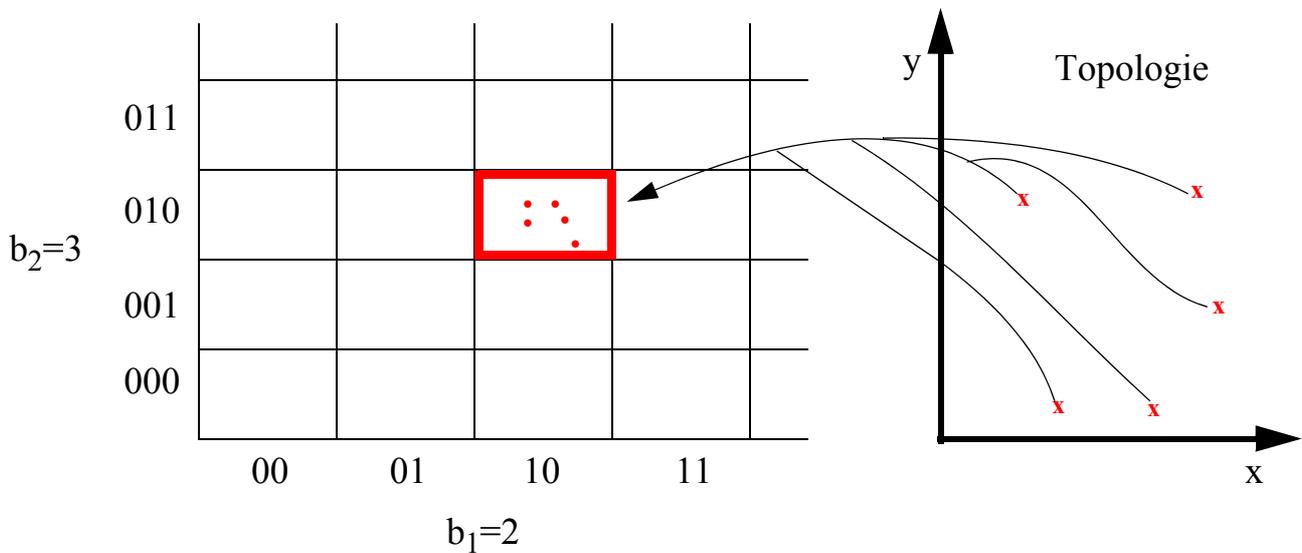
- topologische Struktur der Daten bleibt *nicht* erhalten

- keine Unterstützung von Bereichs- und Best-Match-Anfragen

- Optimale Zuordnung der b_i zu A_i abhängig von Fragehäufigkeiten

Multi-Key-Hashing: Beispiel (1)

Multi-Key-Hashing: Beispiel (2)



Bucket mit Adresse '10010' enthält alle Sätze mit

$$h_1(a_1) = '10' \quad \text{und} \quad h_2(a_2) = '010'$$

• Anwendungsbeispiel:

Pnr: INT (5) $b_1 = 4$

Svnr: INT (9) $b_2 = 3$

Aname: CHAR (10) $b_3 = 2$

↪ $B=9$ (512 Buckets)

$$h_1(\text{Pnr}) = \text{Pnr} \bmod 16$$

$$h_1(58651) = 11 \quad \rightarrow 1011$$

$$h_2(\text{Svnr}) = \text{Svnr} \bmod 8$$

$$h_2(130326734) = 6 \quad \rightarrow 110$$

$$h_3(\text{Aname}) = L(\text{Aname}) \bmod 4$$

$$h_3(\text{XYZ55}) = 1 \quad \rightarrow 01$$

↪ B-Adr. = 101111001

• Anzahl der Zugriffe

Exact-Match-Anfragen: Zugriff auf 1 Bucket

Partial-Match-Anfragen:

$$\text{Pnr} = 58651$$

↪ $N_B = 2^9 / 2^4 = 32$

$$(\text{Pnr} = 73443) \text{ AND } (\text{Svnr} = 2332)$$

↪ $N_B = 2^9 / 2^{4+3} = 4$

Mehrdimensionale binäre Suchbäume

- **Erweiterung des Binärbaumes**

- Berücksichtigung von k Schlüsseln: k -d-Baum
- alle Datensätze werden mit Hilfe der Baumstruktur organisiert:
knotenorientiert (homogen)
- Wartungsoperationen wie beim binären Suchbaum, aber
- **auf jeder Ebene erfolgt Schlüsselvergleich für einen der k Schlüssel**

- **Diskriminator legt Schlüssel auf jeder Ebene fest**

- zyklische Variation des Diskriminators d :
für alle Knoten der Baumebene i gilt: $d = (i \bmod k) + 1$
- der linke (rechte) Nachfolger zu einem Knoten enthält alle Sätze mit
kleineren (größerem) Werten für das Diskriminatorattribut
 $\forall Q \in \text{LOSON}(P): K_d(Q) \leq K_d(P)$
 $\forall R \in \text{HISON}(P) : K_d(R) > K_d(P)$

- **Baumstruktur ist nicht balanciert – verschiedene Erweiterungsvorschläge**

- zu einem Satz werden $m-1$ nachfolgende Ebenen des Baumes komplett
in eine Seite abgebildet:
 - Seite muß bis zu $2^m - 1$ Sätze aufnehmen können
 - jeder Satz auf Ebene $1, m+1, 2m+1, \dots$ bildet Wurzel eines Teil-
baumes, der in einer Seite liegt

- **Weitere Variante des k -d-Baumes**

- blattorientiert (heterogen, Speicherung der Sätze in Buckets)

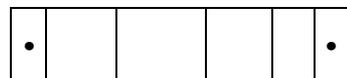
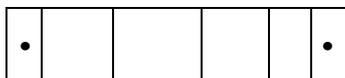
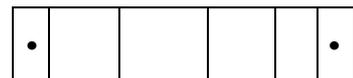
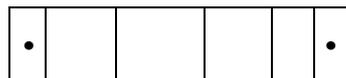
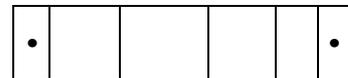
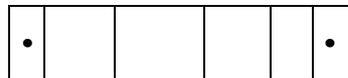
Beispiel: 3-d-Baum (knotenorientiert)

Alter	Gehalt	Ort
35	17K	KL
28	40K	F
29	15K	DA
25	45K	KL
40	12K	SB
29	16K	DA
30	17K	F
42	100K	F
29	14K	DA

Diskr.



Alter Gehalt Ort



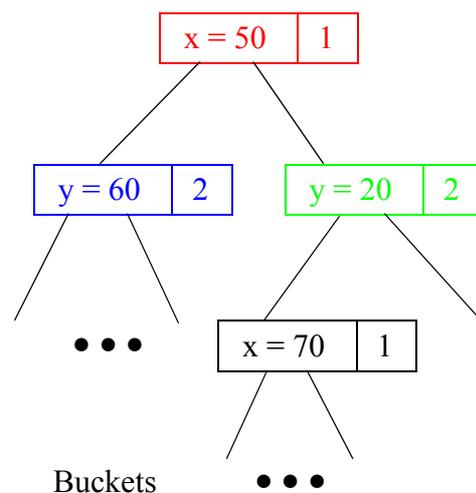
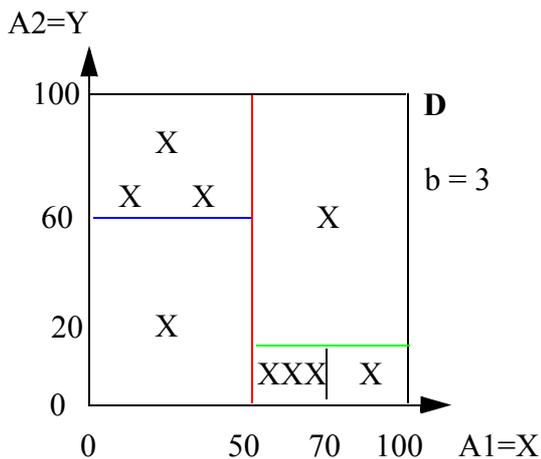
- Baumstruktur abhängig von Einfügereihenfolge (unbalancierter Baum)
- Eingrenzung des Suchraumes für Partial-Match- und Bereichsanfragen
- Löschen ist sehr schwierig

Organisation des umgebenden Datenraums – Divide and Conquer

- Zerlegungsprinzip von D**

- D wird dynamisch in Zellen aufgeteilt
- Die Objekte einer Zelle werden als Sätze in Buckets abgelegt
- bei Bucket-Überlauf: lokale Zellverfeinerung
- ➔ **Divide and Conquer**
- abschnittsweise Erhaltung der Topologie (Clusterbildung)
- Baum als Zugriffsstruktur für die Buckets hat nur Wegweiserfunktion

- Beispiel: Heterogener k-d-Baum – k=2:**



- Eigenschaften des k-d-Baumes**

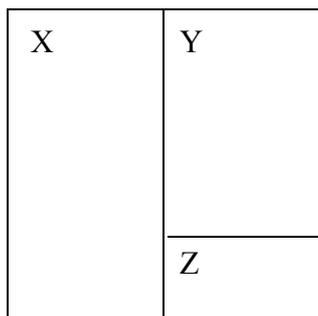
- Clusterbildung durch Buckets ist Voraussetzung für praktischen Einsatz
- Eingebauter Balancierungsmechanismus ist nicht vorhanden
- Wie werden Aktualisierungsoperationen (Löschen!) durchgeführt?
- Wie werden die verschiedenen Anfragetypen unterstützt?

k-d-B-Baum¹

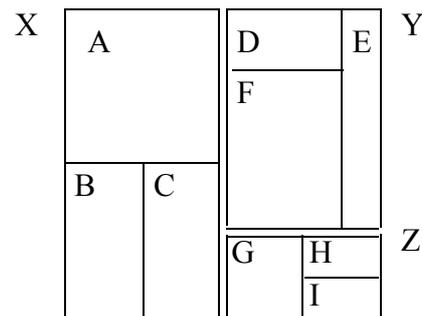
- **Lösungsidee: Kombination von k-d- und B*-Baum**

- k-d-B-Baum paginiert k-d-Bäume und ordnet ihren Teilbäumen Buckets (Seiten) zu, wie das bei B*-Bäumen der Fall ist
- auf jeder Baumebene wird der k-dimensionale Datenraum in schachtelförmige Zellen oder Regionen (bei k=2 in Rechtecke) partitioniert, wobei eine Region jeweils die Zellen/Regionen eines Knotens der darunterliegenden Ebene zusammenfaßt und repräsentiert
- alle Datensätze sind in den Blättern (Buckets) gespeichert
- die inneren Knoten (Index- oder Directory-Seiten) besitzen nur Wegweiser
- operationales Verhalten wie beim B*-Baum

- **Beispiel:**



Ebene 0 (Wurzel),
eine Indexseite



Ebene 1,
drei Indexseiten X, Y, Z

Ebenenweise Organisation der Regionen von D (k=2)

- **Sehr komplexe Regeln für die dynamische Reorganisation**

- Bucket-Überlauf erzwingt Split der Zelle – optimierte Aufteilung der Zelle möglich (z. B. Median-Split)
- Fortpflanzung des Split zur Wurzel hin möglich
- bei Split von Indexseiten ist manchmal eine Neuaufteilung von Regionen erforderlich (sonst können sich sehr ungünstige Strukturen ergeben), was einen sich zu den Blattknoten hin fortpflanzenden Split (forced split) auslöst!

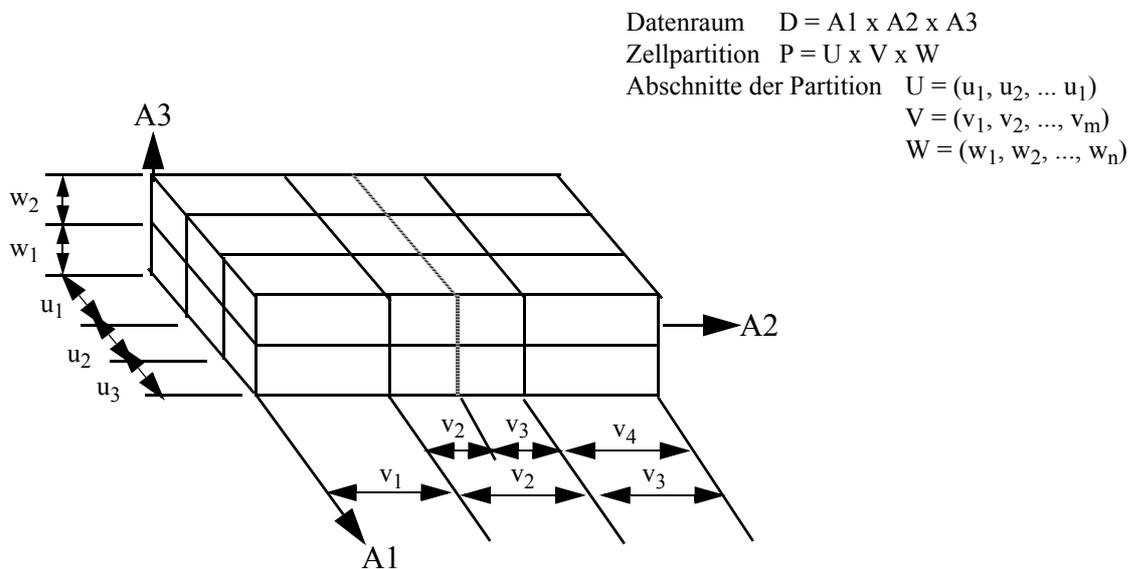
1. Robinson, J.T.: The k-d-B-Tree: a search structure for large multidimensional dynamic indexes, in: Proc. ACM SIGMOD Conf., 1981, pp. 10-18

Organisation des umgebenden Datenraums – Dimensionsverfeinerung

• Prinzip

- Datenraum D wird dynamisch durch **ein orthogonales Raster** (grid) partitioniert, so daß k -dimensionale Zellen (Grid-Blöcke) entstehen
- Die in den Zellen enthaltenen Objekte werden in Buckets gespeichert
- Eine Zelle ist deshalb eindeutig einem Bucket zuzuordnen
- Die klassenbildende Eigenschaft dieser Verfahren ist das Prinzip der Dimensionsverfeinerung, bei dem ein Abschnitt in der ausgewählten Dimension durch einen vollständigen Schnitt durch D verfeinert wird

• Beispiel



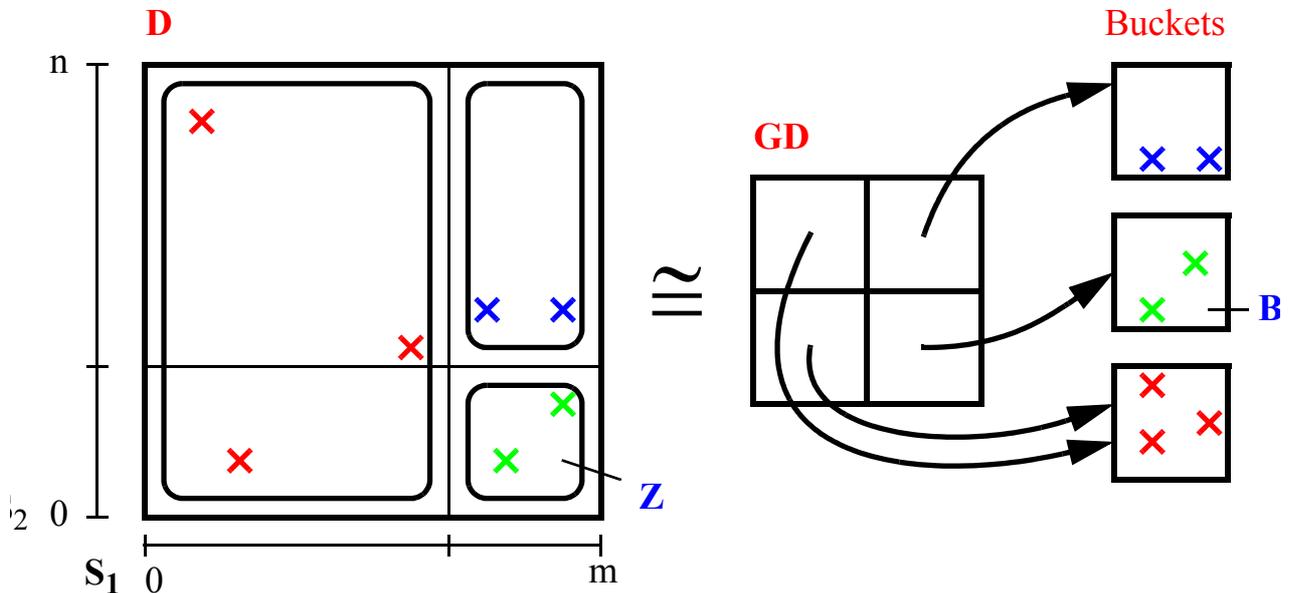
Dreidimensionaler Datenraum D mit Zellpartition P ; Veranschaulichung eines Split-Vorganges im Intervall v_2

• Probleme der Dimensionsverfeinerung

- Wieviele neue Zellen entstehen jedesmal?
- Was folgt für die Bucketzuordnung?
- Welche Abbildungsverfahren können gewählt werden?
- Gibt es Einschränkungen bei der Festlegung der Dimensionsverfeinerung?

Grid-File²: Idee

- Zerlegungsprinzip von D: Dimensionsverfeinerung



- **Komponenten**

- k Skalierungsvektoren (Scales) definieren die Zellen (Grid) auf k-dim. Datenraum D
- Zell- oder Grid-Directory GD: dynamische k-dim. Matrix zur Abbildung von D auf die Menge der Buckets
- Bucket: Speicherung der Objekte einer oder mehrerer Zellen (Bucketbereich BB)

- **Eigenschaften**

- 1:1-Beziehung zwischen Zelle Z_i und Element von GD
- Element von GD = Ptr. zu Bucket B
- n:1-Beziehung zwischen Z_i und B

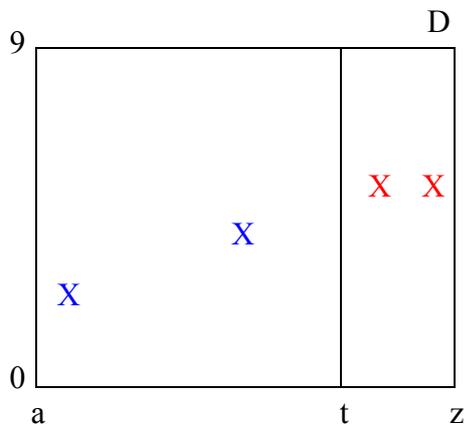
- **Ziele**

- Erhaltung der Topologie
- effiziente Unterstützung aller Fragetypen
- vernünftige Speicherplatzbelegung

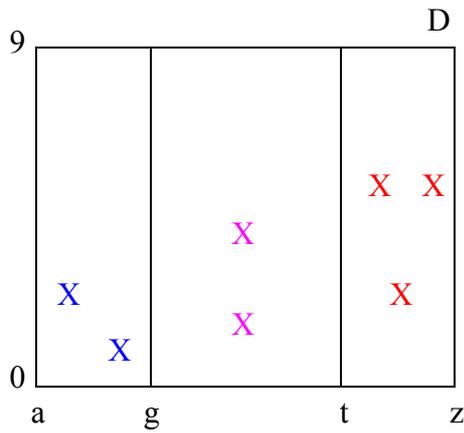
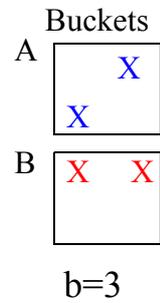
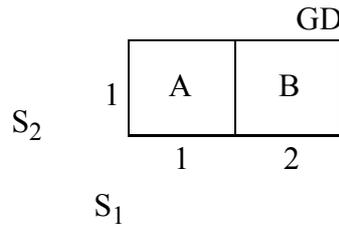
2. Nievergelt, J. et al.: The Grid File: An Adaptable, Symmetric Multikey File Structure, ACM Trans. Database System, 1984, pp. 38-71

Grid-File – Beispiel

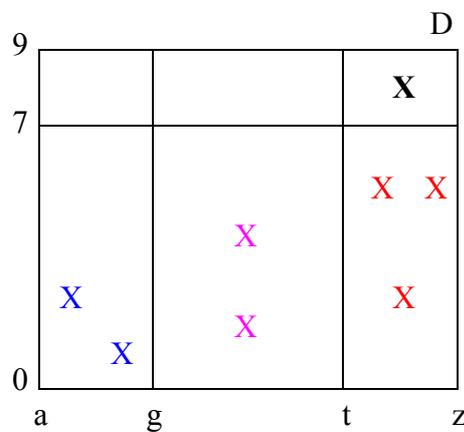
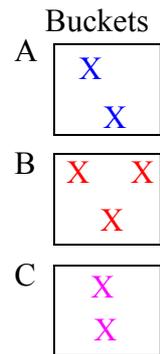
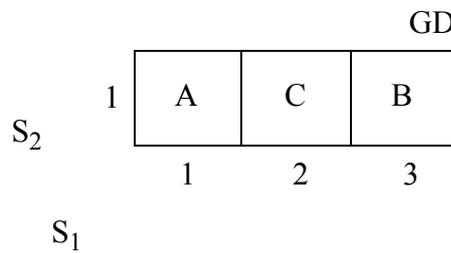
- Schrittweise Entwicklung eines GF



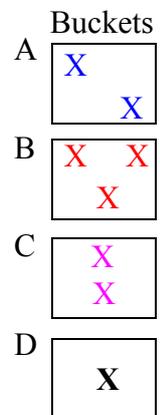
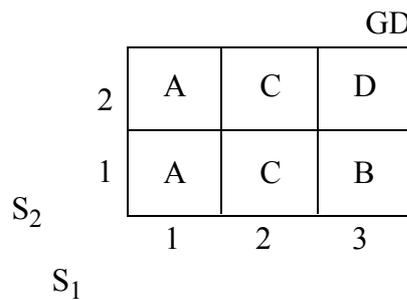
Situation a



Situation b



Situation c



Zentrale Datenstruktur: Grid-Directory

- **Anforderungen**

- Prinzip der zwei Plattenzugriffe
unabhängig von Werteverteilungen, Operationshäufigkeiten
und Anzahl der gespeicherten Sätze
- Split- und Mischoperationen jeweils nur auf zwei Buckets
- Speicherplatzbelegung
 - durchschnittliche Belegung der Buckets nicht beliebig klein
 - schiefe Verteilungen vergrößern nur GD

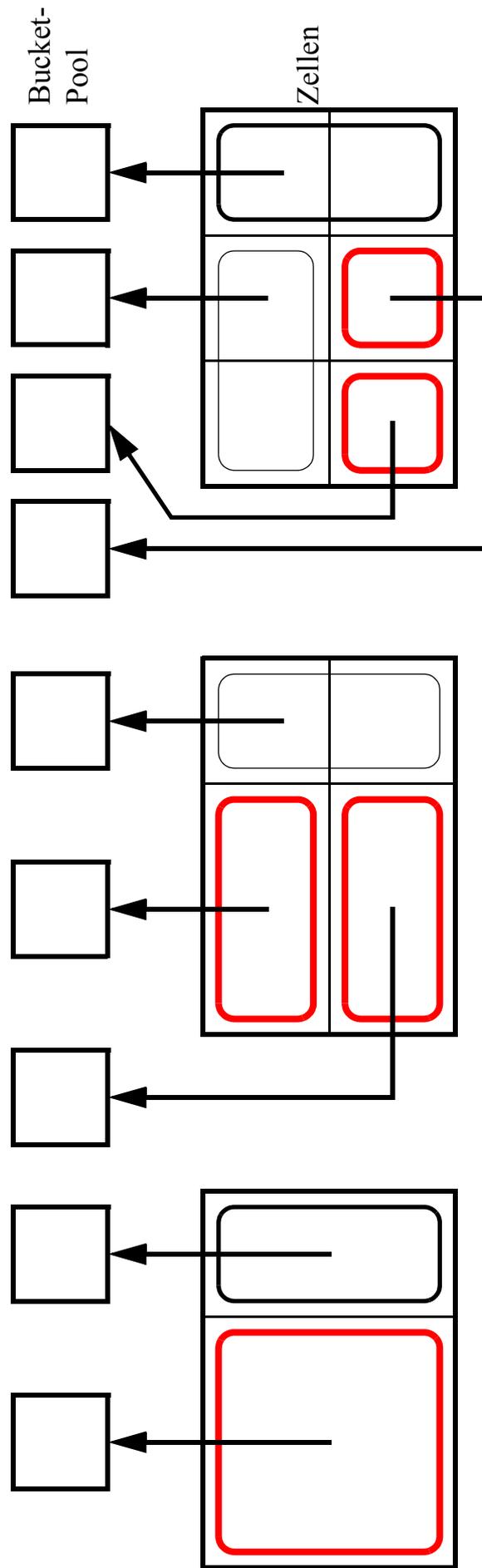
- **Entwurf einer Directory-Struktur**

- dynamische k -dim. Matrix GD (auf Externspeicher)
- k eindim. Vektoren S_i (im Hauptspeicher)

- **Operationen auf GD**

- direkter Zugriff auf einen GD-Eintrag
- relativer Zugriff (NEXTABOVE, NEXTBELOW)
- Mischen zweier benachbarter Einträge einer Dimension
(mit Umbenennung der betroffenen Einträge)
- Splitting eines Eintrages einer Dimension (mit Umbenennung)

Schachelförmige Zuweisung von Zellen zu Buckets

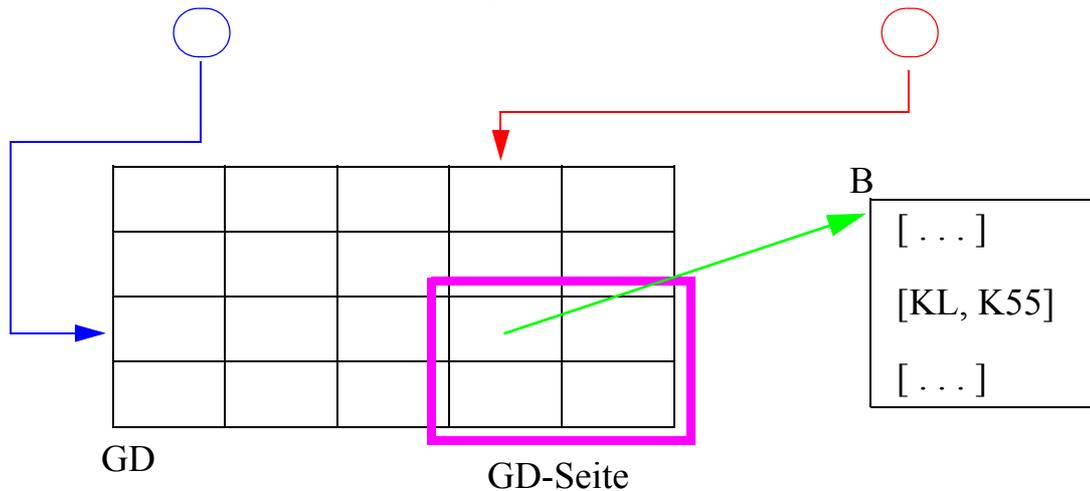


Grid-File – Suchfragen

- **Exakte Anfrage (exact match)**

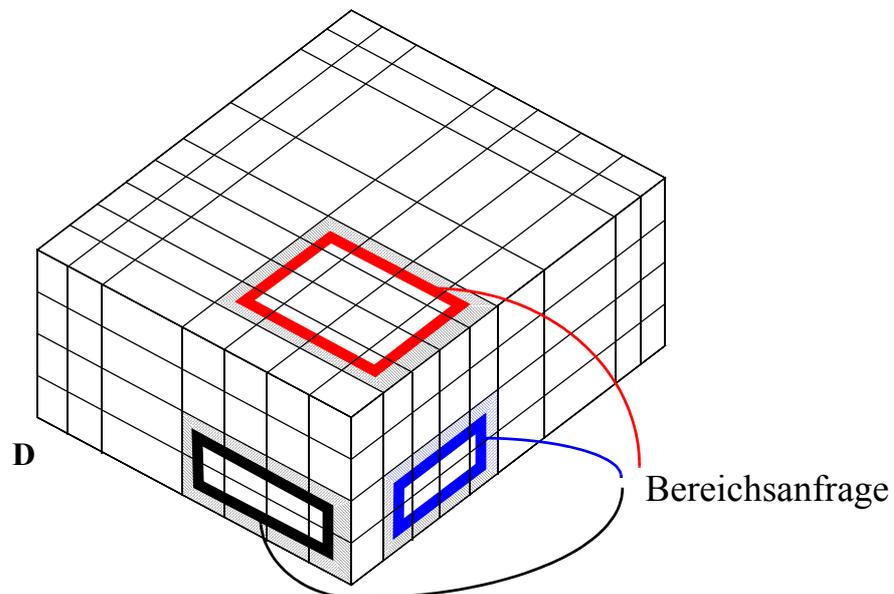
```
SELECT *  
FROM PERS  
WHERE ORT = 'KL' AND ANR = 'K55'
```

S₁: AA, DA, FR, MK, ZZ S₂: K00, K17, K39, K52, K89, K99



- **Bereichsanfrage**

- Bestimmung der Skalierungswerte in jeder Dimension
- Berechnung der qualifizierten GD-Einträge
- Zugriff auf die GD-Seite(n) und Holen der referenzierten Buckets



Grid-File – Leistungsbetrachtungen

- **Wachstum von GD ist überlinear**

- bei gleichförmigen Objektverteilungen: $O(N^{1+(k-1)/(k*b)})$
- bei schiefen Objektverteilungen bis zu:

- **Verfeinerung in Dimension k**

- Anzahl der hinzukommenden GD-Einträge: $n = (m_1 * m_2 * \dots * m_{k-1})$
- erforderliche Plattenzugriffe ($E = \#GD\text{-Einträge}$):

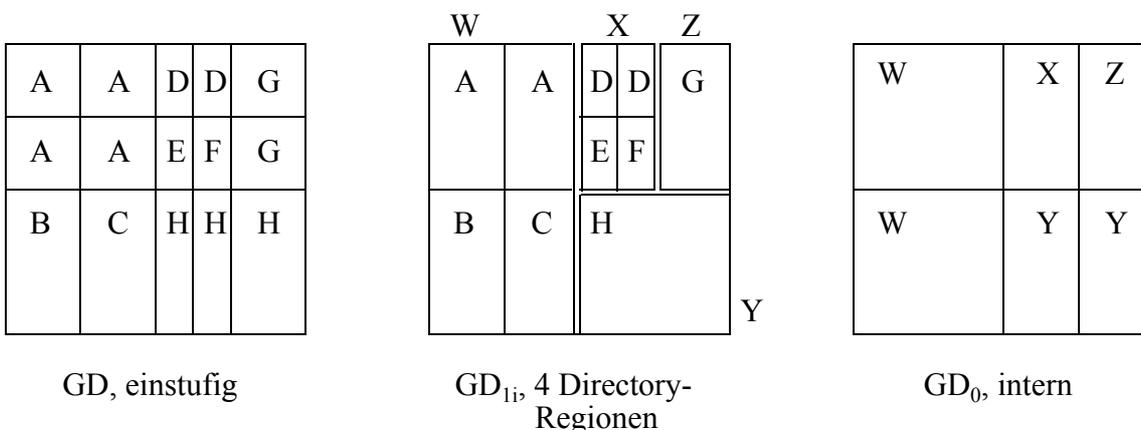
- **Speicherplatzausnutzung für Buckets:**

- **Problemfälle bei der Abbildung:**

- geringe Selektivität und ungünstige Verteilung
 $S_1 : \geq a, \geq d, \geq d, \geq d, \geq f, \geq s, \dots$
 $S_2 : \geq 0, \geq 5, \geq 5, \geq 5, \geq 7, \geq 9, \dots$
- Wieviele GD-Einträge kann eine exakte Anfrage liefern?

- **Organisation von GD als Grid-File**

- Regelfall: zweistufige Realisierung



Gegenüberstellung von einstufiger und zweistufiger Grid-Directory-Realisierung

- Einsparungen bei Bereichsanfragen und Zellabbildung

Zugriffspfade für ausgedehnte räumliche Objekte

- **Ausgedehnte räumliche Objekte besitzen**

- allgemeine Merkmale wie Name, Beschaffenheit, . . .
- Ort und Geometrie (Kurve, Polygon, . . .)

- **Indexierung des räumlichen Objektes**

- genaue Darstellung?
- Objektapproximation durch schachtelförmige Umhüllung – effektiv!
 - ➔ dadurch werden Fehltreffer möglich

- **Probleme**

- neben Objektdichte muß Objektausdehnung bei der Abbildung und Verfeinerung berücksichtigt werden
- Objekte können andere enthalten oder sich gegenseitig überlappen

- **Klassifikation der Lösungsansätze**

- sich gegenseitig überlappende Regionen (R-Baum)
- Clipping (R^+ -Baum)
- Transformationsansatz

➔ bildet ausgedehnte räumliche Objekte funktional auf höherdimensionale Punkte ab – begrenzte Anwendbarkeit und Tauglichkeit!

R-Baum³

- **Ziel: Effiziente Verwaltung räumlicher Objekte**
(Punkte, Polygone, Quader, ...)
- **Anwendungen:**
 - Kartographie:
Speicherung von Landkarten,
effiziente Beantwortung “geometrischer” Fragen
 - CAD:
Handhabung von Flächen, Volumina und Körpern
(z.B. Rechtecke beim VLSI-Entwurf)
 - Computer-Vision und Robotics
- **Hauptoperationen:**
 - Punktanfragen (point queries):
Finde alle Objekte, die einen gegebenen Punkt enthalten
 - Gebietsanfragen (region queries):
Finde alle Objekte, die mit einem gegebenen Suchfenster
überlappen (es umschließen, in ... vollständig enthalten sind)
- **Ansatz: Speicherung und Suche von achsenparallelen Rechtecken**
 - Objekte werden durch Datenrechtecke repräsentiert und müssen
durch kartesische Koordinaten beschrieben werden
 - Repräsentation im R-Baum erfolgt durch minimale begrenzende
(k-dimensionale) Rechtecke/Regionen
 - Suchanfragen beziehen sich ebenfalls auf Rechtecke/Regionen

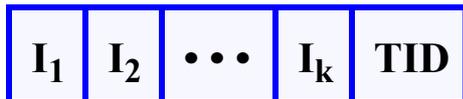
3. A. Guttman: R-Trees: A Dynamic Index Structure for Spatial Searching,
in: Proc. ACM SIGMOD Conf., 1984, pp. 47-57

R-Baum (2)

- **R-Baum ist höhenbalancierter Mehrwegbaum**

- jeder Knoten entspricht einer Seite
- pro Knoten maximal M , minimal m ($\geq M/2$) Einträge

Blattknoteneintrag:



kleinstes umschreibendes Rechteck (Datenrechteck) für TID

Zwischenknoteneintrag:



Intervalle beschreiben kleinste umschreibende Datenregion für alle in PID enthaltenen Objekte

I_j = geschlossenes Intervall bzgl. Dimension j

TID: Verweis auf Objekt

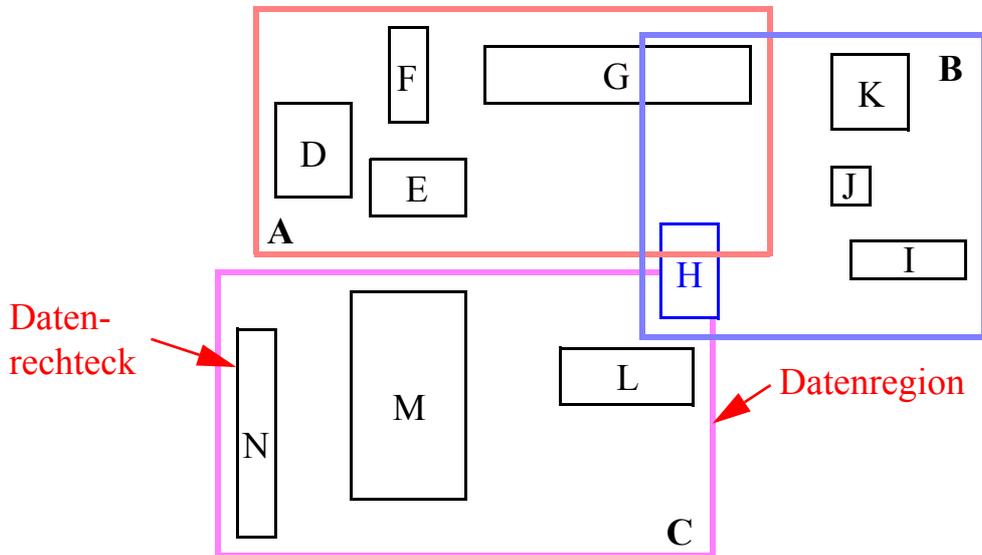
PID: Verweis auf Sohn

- **Eigenschaften**

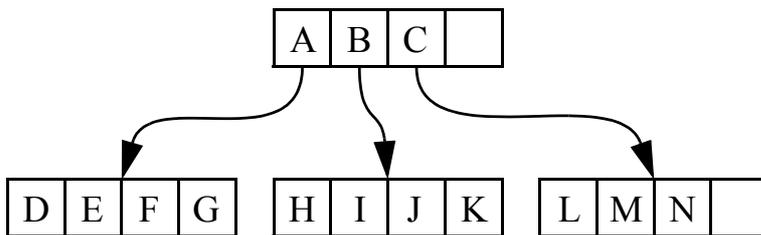
- starke Überlappung der umschreibenden Rechtecke/Regionen auf allen Baumebenen möglich
- bei Suche nach Rechtecken/Regionen sind ggf. mehrere Teilbäume zu durchlaufen
- Änderungsoperationen ähnlich wie bei B-Bäumen

Abbildung beim R-Baum

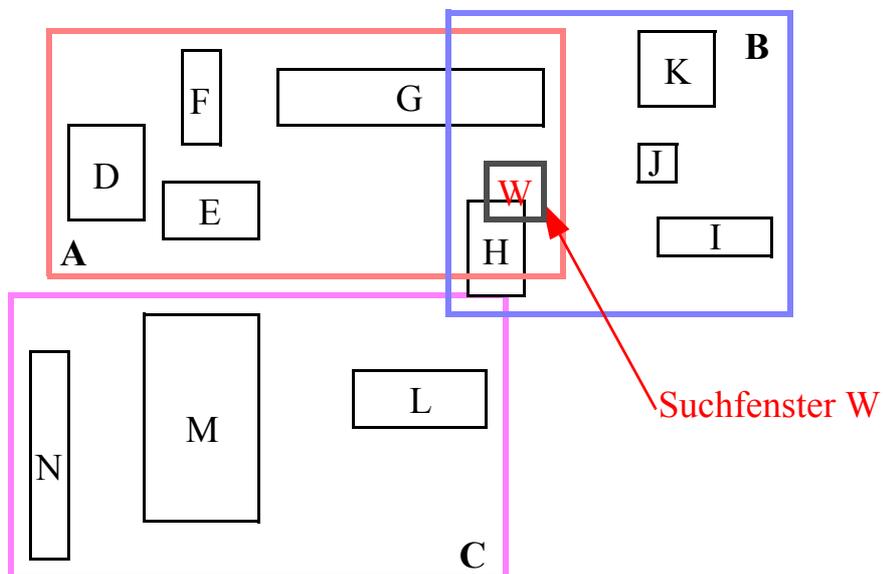
- Aufteilung des Datenraumes



- R-Baum

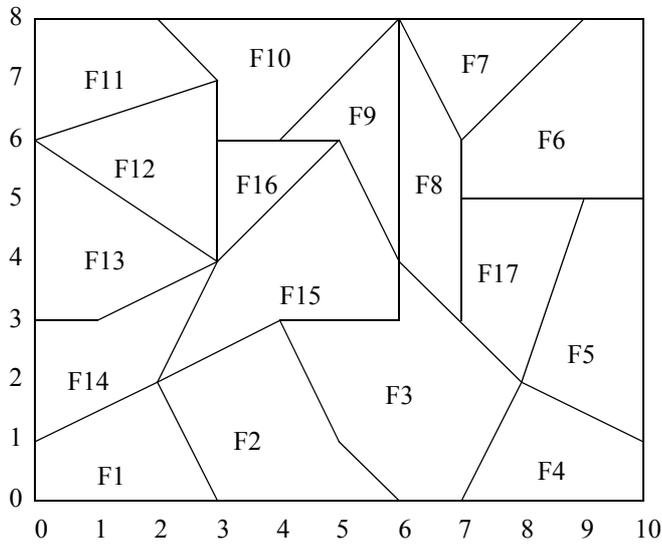


- Beispiel für ein „schlechtes“ Suchfenster

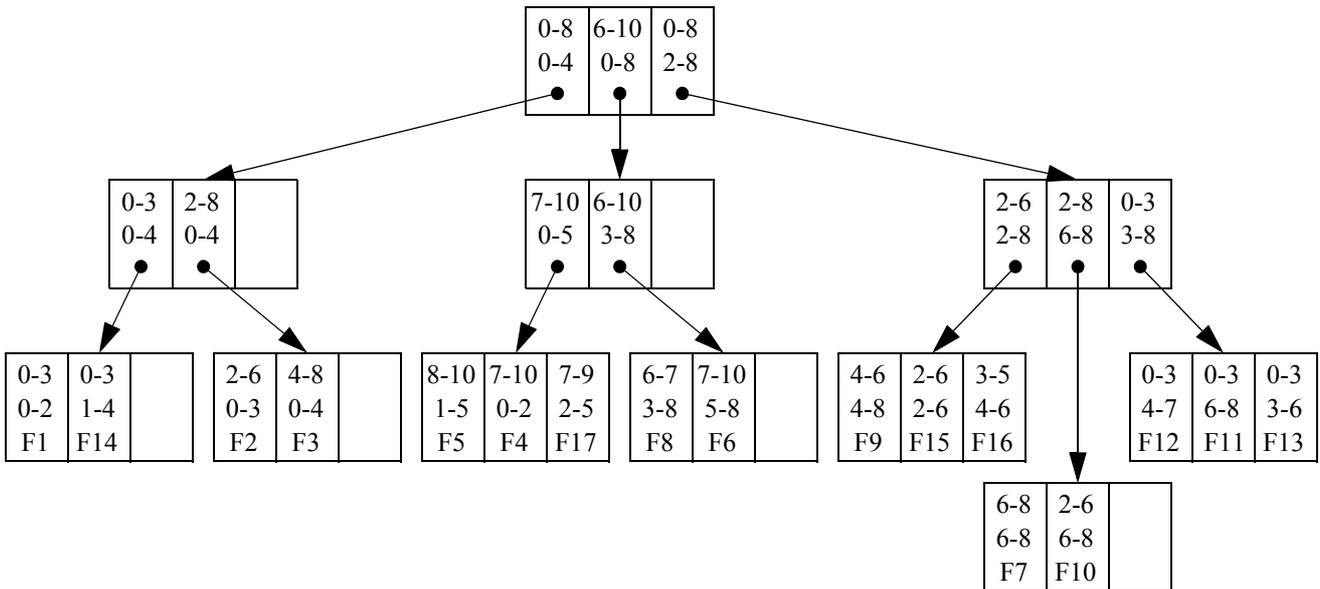


R-Baum: Beispiel

- Abzuspeichernde Flächenobjekte



- Zugehöriger R-Baum



Suchoptimierung durch den R^+ -Baum

- **Überdeckung und Überlappung** bei einer Ebene des R-Baumes

- **Überdeckung** (coverage) ist der gesamte Bereich, um alle zugehörigen Rechtecke zu überdecken
- **Überlappung** (overlap) ist der gesamter Bereich, der in zwei oder mehr Knoten enthalten ist
- Minimale Überdeckung reduziert die Menge des „toten Raumes“ (leere Bereiche), der von den Knoten des R-Baumes überdeckt wird.
- Minimale Überlappung reduziert die Menge der Suchpfade zu den Blättern (noch kritischer für die Zugriffszeit als minimale Überdeckung)

➔ Effiziente Suche erfordert minimale Überdeckung und Überlappung

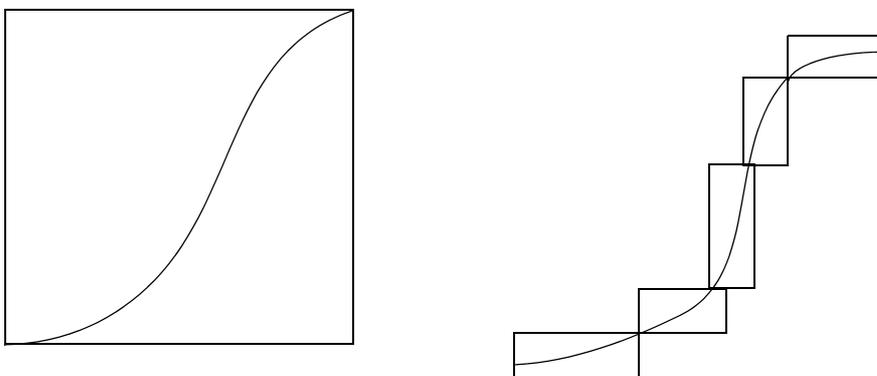
- **Idee**

Es sind Partitionierungen erlaubt, die Datenrechtecke „zerschneiden“ (Clipping)

➔ Vermeidung von Überlappungen bei Zwischenknoten

- **Konsequenz**

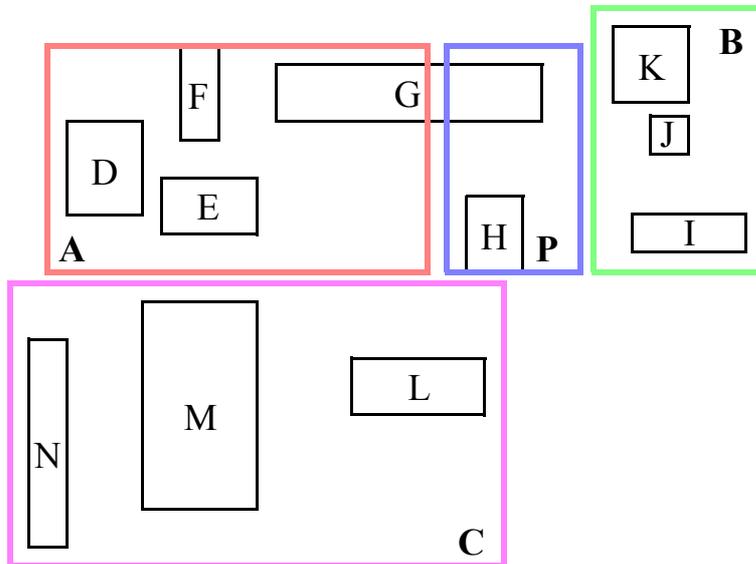
Daten-Rechteck wird ggf. in eine Sammlung von disjunkten Teilrechtecken zerlegt und auf der Blattebene in verschiedenen Knoten gespeichert.



Aufteilungsmöglichkeiten eines langen Linienobjektes im R^+ -Baum

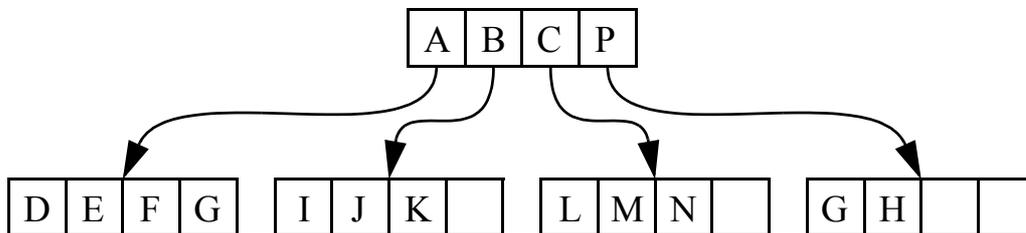
R⁺-Baum⁴

- Aufteilung des Datenraumes



➔ Höhere Flexibilität durch Partitionierung von Daten-Rechtecken

- R⁺-Baum



- Eigenschaften

- Überlappung von Datenregionen wird vermieden
- Überdeckungsproblematik wird wesentlich entschärft
- Clusterbildung der in einer Region zusammengefaßten Objekte kann durch Clipping verhindert werden
- komplexere Algorithmen für bestimmte Anfragen (Enthaltensein)
- schwierigere Wartung, keine Leistungsvorteile gegenüber R-Baum

4. Sellis, T. et al.: The R⁺-Tree: A Dynamic Index for Multi-Dimensional Objects, in: Proc.14th Int. Conf. on Very Large Data Bases, Brighton,1987, pp. 507-518

Eindimensionale Intervalle

- **Anwendungsbeispiele für Intervallanfragen**

- Temporale Datenbanken
- Fuzzy-Logik/Matching
- Intervallmodellierung ist oft hilfreich bei Profilabgleich (Personalisierung von Web-Seiten) oder Qualitätsklassen (Qualitätskriterien in Intervalldarstellung)

- **Mehrdimensionale Indexstrukturen sind nicht spezialisiert** auf die

- effiziente Verwaltung von 1D-Objekten bzw. Intervallen
- optimierte Auswertung entsprechender Anfragen

- **Wie lassen sich Intervalle darstellen?**

- in B*-Bäumen, Quad-Tree, Grid-File, ...
- **Idee:** Endpunkt-Transformation resultiert in einem Parameterraum mit 2 Dimensionen von gleichen Wertebereichen

$$[13, 27] \xrightarrow[\text{Transformation}]{\text{Parameterraum}} (13, 27)$$

1D-Intervall

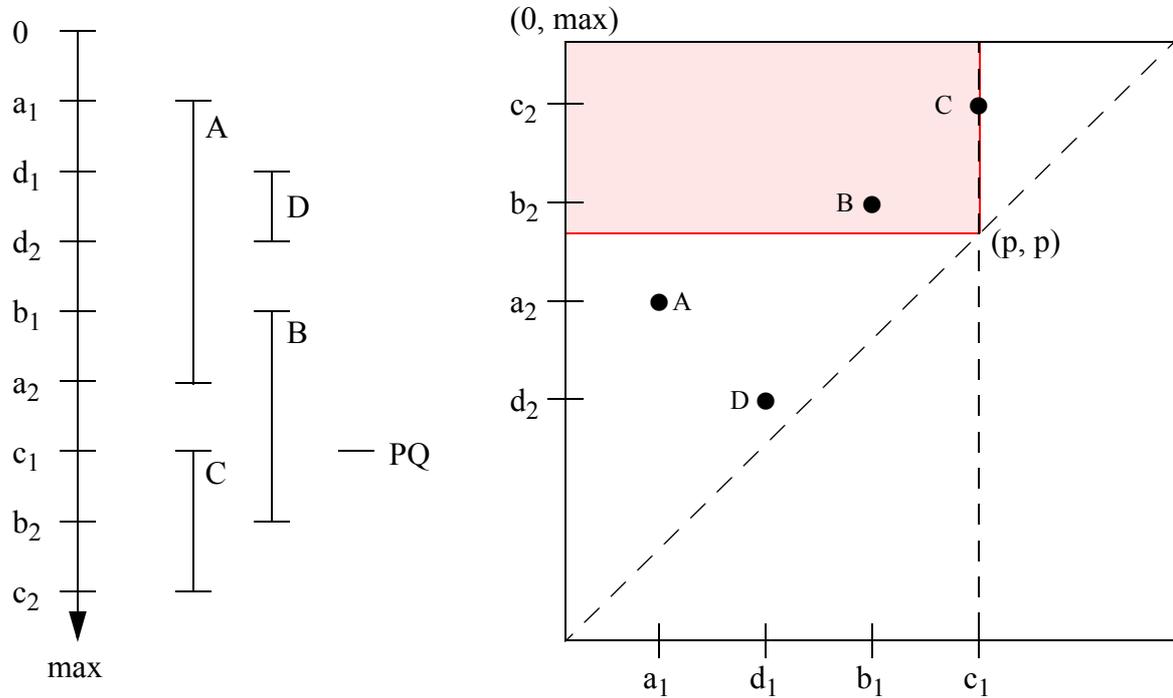
2D-Punkt

- **Welche Anfragetypen sind zu unterstützen?**

- Enthaltensein eines Punktes (PQ)
- Überlappung von Intervallen
 - IQ: Intersection Query (Durchschnitt)
 - EQ: Enclosure Query (Umschließung)
 - CQ: Containment Query (Enthaltensein)

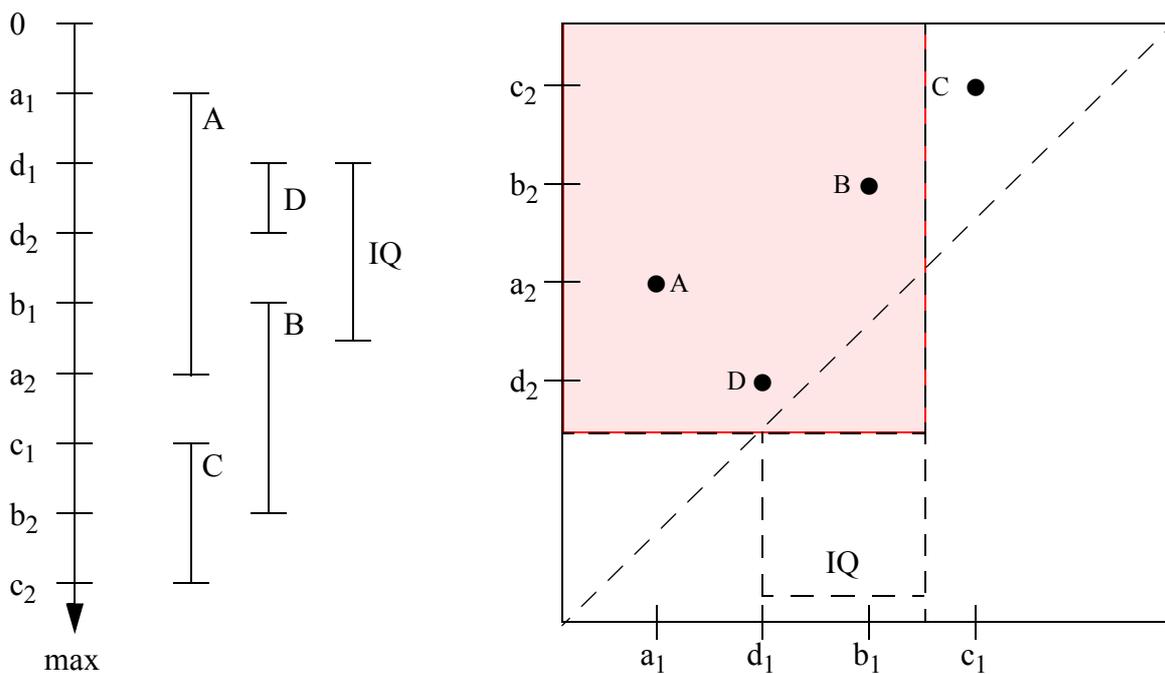
Eindimensionale Intervalle (2)

- Wertebereich für Intervalle $U = [0, \max]$
- Point Query (PQ): $[p, p]$ mit $(p \in U)$



- Anfragebereich: $[(0, \max), (p, p)]$

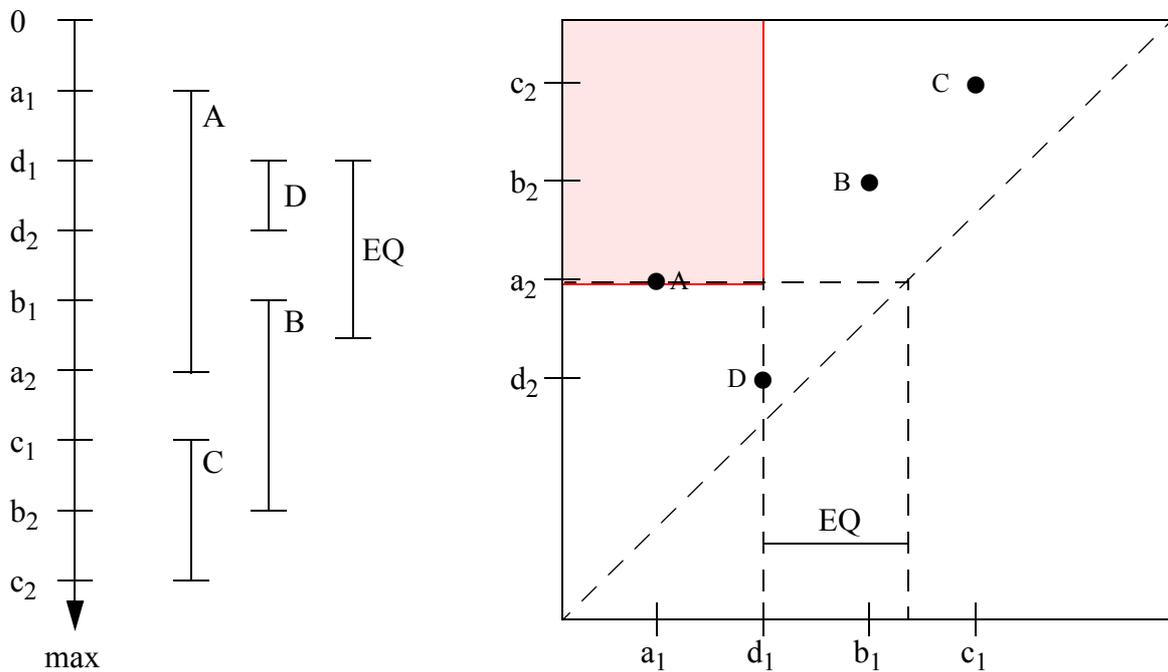
- Intersection Query (IQ): $[s, e]$ mit $(s, e \in U)$



- Anfragebereich: $[(0, \max), (e, s)]$

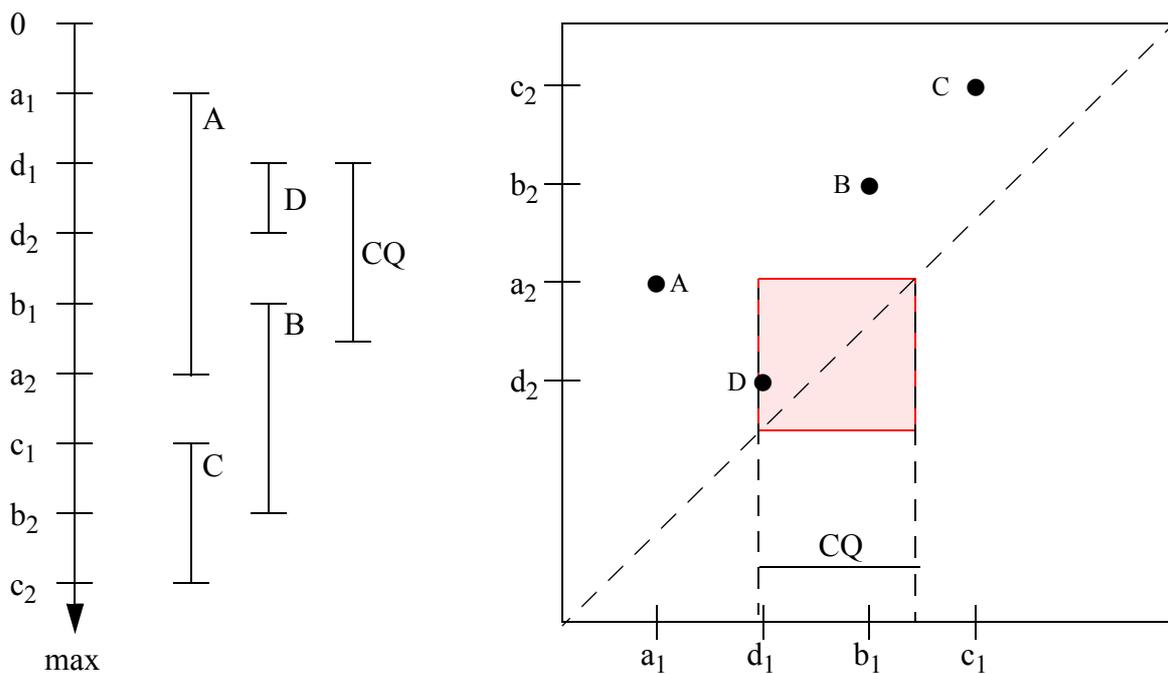
Eindimensionale Intervalle (3)

- **Enclosure Query (EQ):** $[s, e]$ mit $(s, e \in U)$



- Anfragebereich: $[(0, \max), (s, e)]$

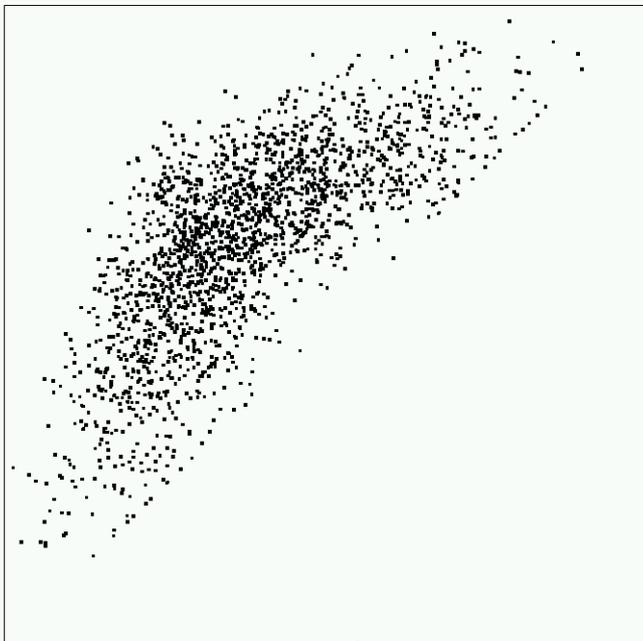
- **Containment Query (CQ):** $[s, e]$ mit $(s, e \in U)$



- Anfragebereich: $[(s, e), (e, s)]$

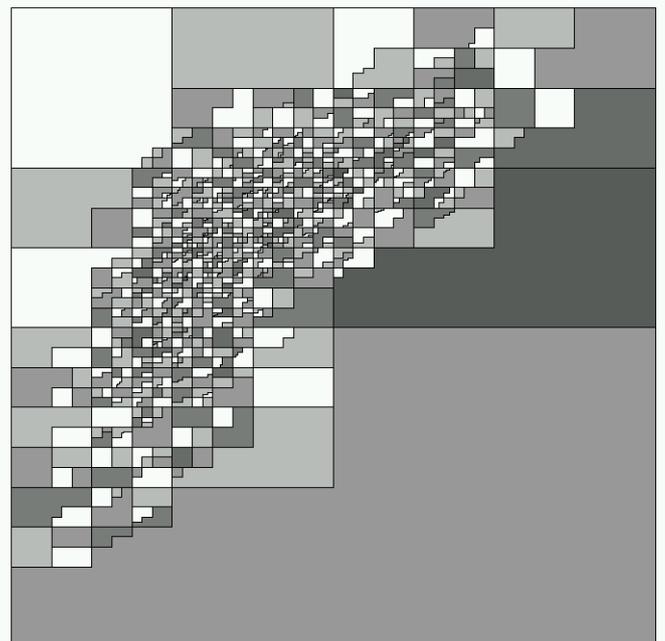
Eindimensionale Intervalle (4)

- **Abbildung auf UB-Baum** ➔ **Vorteile des B*-Baumes**
 - hervorragende Externspeicherverwaltung, Skalierbarkeit
 - logarithmische ‘worst case’-Garantien für alle wichtigen Operationen
 - Antwortzeit ist linear zur Größe der Ergebnismenge
 - Unabhängigkeit von Eingabedaten und Einfügereihenfolge
 - geringe Speicherplatzkosten für den Index
- **Beispiel mit 2000 Gauß-verteiltern Intervallen**



Parameterraum

UB-Baum



Bayer, R.:
The Universal B-Tree for Multidimensional
Indexing: General Concepts, in: Proc. Int.
Conf. on Worldwide Computing and its
Applications, LNCS 1274, Springer, 1997,
pp. 198-209

Anwendung für Geo-DBS

- **Was sind Geo-Objekte?**
- **Darstellung eines bestimmten Ausschnitts der Erdoberfläche**
 - geometrische Komponente (Ort, Ausdehnung)
 - thematische Komponente (qualitative oder quantitative Beschreibung)
- **Modellierung**
 - topologische Grundelemente: Punkt, Linie, Fläche
 - linienartige Gebilde: Polygonzüge
 - flächige Objekte: Rastermodell, Vektormodell
 - Vektormodell wird für Geo-DBS favorisiert
 - ↳ gute Skalierbarkeit, geringerer Speicherplatzbedarf
 - Objektklasse zur Flächenmodellierung: einfache Polygone mit Löchern (EPL)
- **Thematische Komponente**
 - einfache Datentypen zusammengefaßt in einem Attributvektor (Tupel)
 - Beispiel: Flächennutzung, Niederschlagsmenge
- **Charakteristika der Objekte**
 - Anzahl der Objekte: bis zu 10^9 Datensätze
 - Datenvolumen: bis zu 1 TeraByte
 - Variabilität der Objekte und Objektmengen
 - Objektausdehnung: $1 : 10^6$
 - Objektform
 - Speicherplatzbedarf von Polygonen: 0.5 KB – 1 GB
 - Verteilung der Objekte in der Ebene (Objektdichte): $1 : 10^4$
 - ↳ es gibt keine Obergrenzen für diese Eigenschaften

Anfragen an räumliche Objekte – Aspekte der Erweiterung der DB-Sprache

- **Kein Raumbezug und keine räumlichen Operatoren** im Relationenmodell

- Hohe Komplexität bei einfachen Beispielen
- Beispiel: **Darstellung beliebig gelagerter Rechtecke** in der Ebene

R-ECK (RE-NR, X1, Y1, X2, Y2, X3, Y3, X4, Y4)

Anfrage: Finde alle Rechtecke, die ein Rechteck mit den Punkten (PA, PB, PC, PD) echt umschließen.

```
SELECT RE-NR FROM R-ECK
WHERE  XA > X1 AND YA < G(P1, P2, XA)
        AND YA > G(P4, P1, XA) AND
        XB < Y2 AND YB < G(P1, P2, XB)
        AND XB < G(P2, P3, XB) AND
        XC < X3 AND YC < G(P2, P3, XC)
        AND YC > G(P3, P4, XC) AND
        YD > Y4 AND YD > G(P3, P4, XD)
        AND YD > G(P1, P4, XD);
```

mit Abkürzungen: z.B. G(P1, P2, XA) für

$$((Y1 - Y2) * XA) / (X1 - X2) + (Y2 * X1 - Y1 * X2) / (X1 - X2)$$

- **Objekt-relationale Datenmodelle – Schlüsseleigenschaft:**

- Erweiterbarkeit mit benutzerdefinierten Funktionen und Operatoren
- Anfragen mit Operatoren wie Überlappung, Schnitt, Entfernung, Enthaltensein («) ...

➔ **SELECT RE-NR FROM R-ECK X
WHERE [PA, PB, PC, PD] « X;**

Zusammenfassung

- **Wichtige Eigenschaften mehrdimensionaler Zugriffspfade**
 - Erhaltung der topologischen Struktur des Datenraumes
 - Adaption an die Objektdichte
 - Dynamische Reorganisation
 - Balancierte Zugriffsstruktur
 - Unterstützung von verschiedenen Anfragetypen
(*intersection queries, best match queries*)
- **Darstellung von räumlichen Objekten**
 - Abstraktion zu punktförmiger Repräsentation ist Regelfall
 - „Ausgedehnte“ Darstellung nur in grober Annäherung:
 - ↳ **Verarbeitungsprobleme**
- **Vorteile neuerer Konzepte**
 - Organisation des umgebenden Datenraumes – Grid File
 - Darstellung räumlich ausgedehnter Objekte – R- und R⁺-Baum
 - ↳ **Jedoch:** Es existieren sehr viele Vorschläge und Konzepte, die sehr speziell und praktisch nicht erprobt sind
- **Notwendigkeit der DBS-Integration** von „konventionellen“ (eindimensionalen) und mehrdimensionalen Zugriffspfaden
 - Zugriffsmöglichkeiten über räumliche und zeitliche Dimensionen
 - Unterstützung von „Data Warehouse“-Anwendungen, Geographischen Informationssystemen, . . .
- **Erweiterung objekt-relationaler DBS** durch mehrdimensionale Zugriffspfadstrukturen (vor allem Grid File und R-Baum)