

# The Five-Minute Rule

for trading memory for disc accesses

Advanced Topics in Databases  
Marcel Taeumel



# Agenda

2

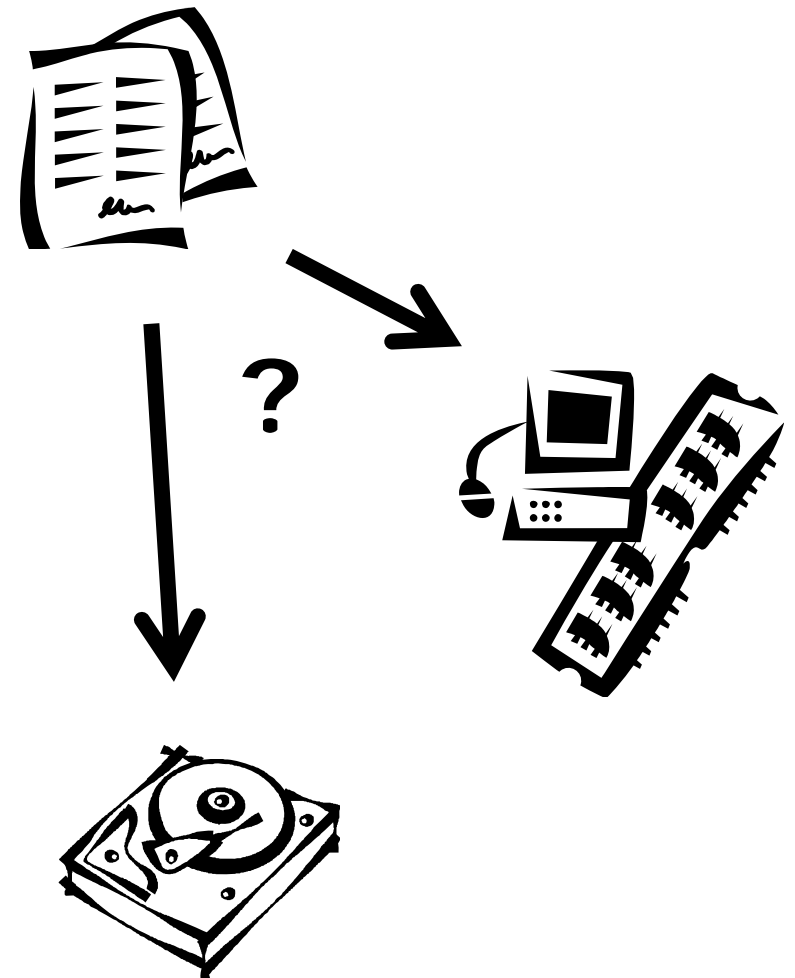
## 1. Autoren

## 2. "Five-Minute Rule"

- Herleitung der Regel
- Fallstudie
- 10 Jahre später – 1997
- 20 Jahre später – 2007

## 3. Weitere Faustregeln

- Größe von Indexseiten
- "One-Minute Sequential Rule"
- Weitere Themenbereiche

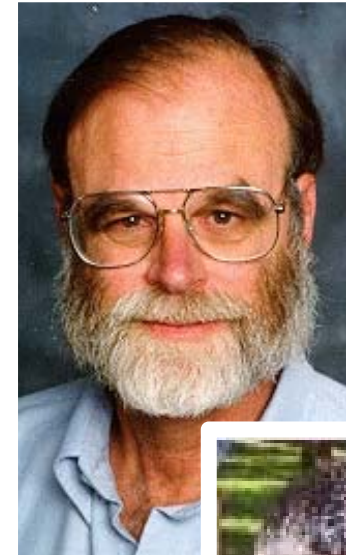


# Autoren

3

## Jim Gray

- 1966 B. Sc. in Mathematik und Statistik
- 1969 Promotion in Informatik
- Entwickler bei IBM, Tandem, DEC



## Franco Putzolu

- Entwickler bei IBM (u.a. DB2)
- Chef-Designer für NonStop SQL (Tandem)
- Senior Database Architect (Oracle)



## Goetz Graefe

- regelbasierte, Top-Down Anfrageoptimierung
- 12 Jahre bei Microsoft tätig (SQL Server)
- seit 2007 bei HP Labs



Erste Gedanken...

1987

# Wichtige Größen (1987)

5

- Kosten pro Festplatte
  - inkl. Betriebskosten (CPU, Controller)
- Zugriffsgeschwindigkeit
  - Anzahl Zugriffe pro Sekunde (wahllos)
- Kosten pro Megabyte Hauptspeicher
- Größe des Datensatzes
- Größe der Datenblöcke
  - maßgebend, falls Datensatz zu groß
  - Fragmentierung

$$\$_{disc} = 30.000 \frac{\$}{disc}$$

$$t_{acc} \approx 66 ms$$

$$\Rightarrow 15 \frac{access}{second}$$

$$\$_{RAM} = 5.000 \frac{\$}{MB}$$

- **Zeit zwischen zwei Zugriffen auf gleiches Datum**

# Vorüberlegung

6

- Kosten für I/O (Festplatte)

$$\$_{I/O} = \frac{\textit{PricePerAccessPerSecond}}{\textit{ReferenceInterval}}$$

- Hauptspeicherkosten für Datensatz/-block

$$\$_{RAM,Record} = \$_{RAM} \cdot \textit{RecordSize}$$

- Kosteneinsparung durch Halten der Daten im Hauptspeicher

$$\textit{Savings} = \$_{I/O} - \$_{RAM,Record}$$

# Herleitung

7

- Break-Even-Point

$$0 = \$_{I/O} - \$_{RAM,Record}$$

- Beispielrechnung (1KB-Seiten)

$$ReferenceInterval = \frac{PagesPerMBofRAM}{AccessPerSecondPerDisc} \cdot \frac{PricePerDiskDrive}{PricePerMBofRAM}$$

$$ReferenceInterval = \frac{1000}{15} \cdot \frac{30000}{5000} \cdot \frac{sec}{acc} = 400 \frac{sec}{acc} \approx \underline{\underline{5 \frac{min}{acc}}}$$

# "Five-Minute Rule"

8

**Daten, welche innerhalb von 5 Minuten\*  
wiederholt referenziert werden, sollten im  
Hauptspeicher verbleiben.**

\* abhängig von der Seitengröße (hier: 1KB)



# Fallstudie

9

- Ein Kunde möchte seine Datenbank vollständig im Hauptspeicher halten. Lohnt sich das?
  
- Fakten
  - **500.000** Datensätze à **1000** Bytes
  - fast alle Zugriffe auf einzelne Datensätze
  - Spitzenlast: **600** Transaktionen pro Sekunde
  - geforderte Antwortzeit: 1 Sekunde

# Fallstudie: "All-In-Main-Memory"

10

- 36 VLX-Prozessoren à 16 MB Hauptspeicher
  
- 2 Festplatten (gespiegelt) für:
  - komplette Datenbank
  - Indizes
  - Programme
  
- Festplatten im Betrieb nicht benötigt
  
- Kosten: **ca. 11,7 Mio \$**

# Fallstudie: "Disc-Memory-Tradeoff"

11

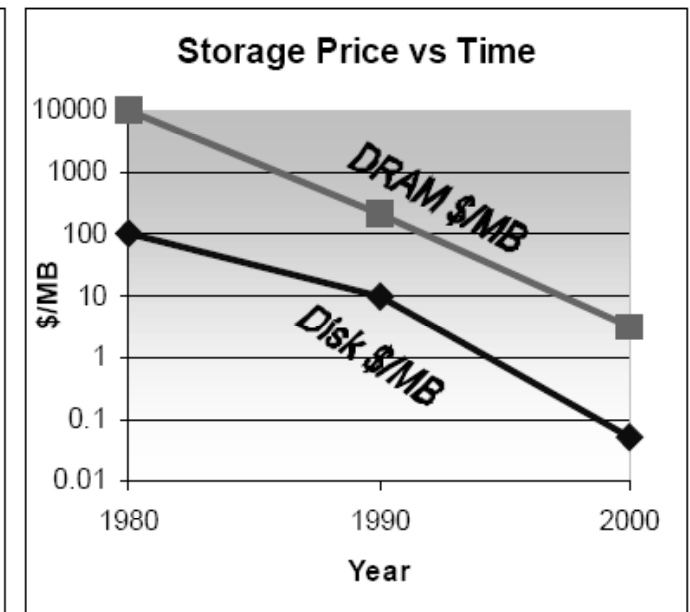
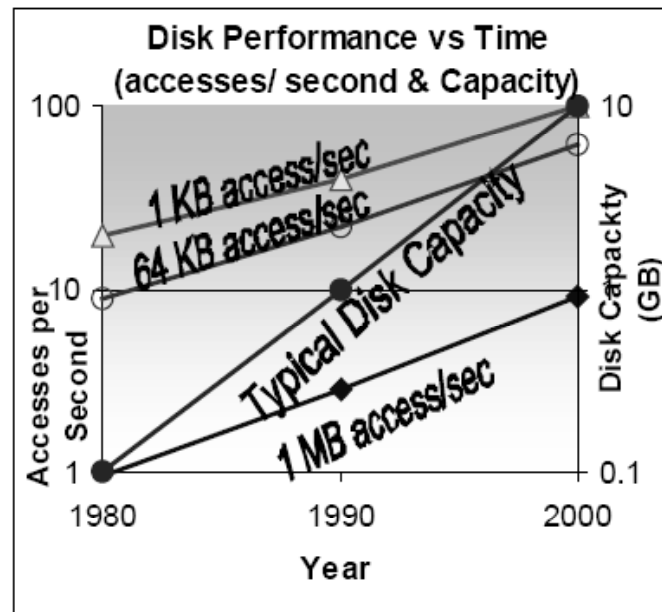
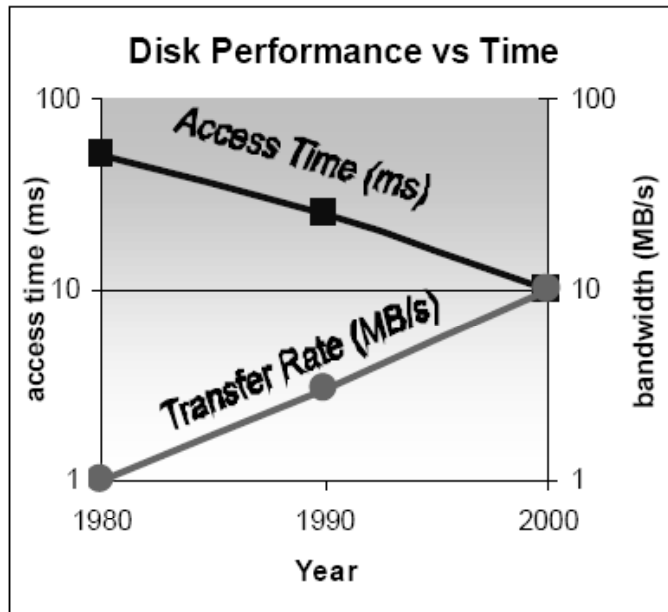
- nur 30.000 der 500.000 Datensätze innerhalb von **5 Minuten** wiederholt referenziert
  - Spitzenlast berücksichtigt
- 96% der Zugriffe für diese 6% der Daten
  - statistische Abschätzungen
- 2 gespiegelte Festplatten für 24 Zugriffe pro Sekunde genügen
  - 4% von 600TPS
  
- 30 VLX-Prozessoren
- 106 MB Hauptspeicher
- **Kosten:** 8,2 Mio \$

## Einsparungen

6 VLX-Prozessoren  
470 MB Hauptspeicher  
3,5 Mio \$ (~30%)

Zehn Jahre später...

1997



- höhere Bandbreiten und günstigerer Hauptspeicher ermöglichen **größere Seiten**

# Technologie und Wirtschaftlichkeit

14

$$ReferenceInterval = \frac{PagesPerMBofRAM}{AccessPerSecondPerDisc} \cdot \frac{PricePerDiskDrive}{PricePerMBofRAM}$$



$$ReferenceInterval = TechnologyRatio \cdot EconomicRatio$$

10x 

10x 

- aktualisierte Beispielrechnung (8KB-Seiten)

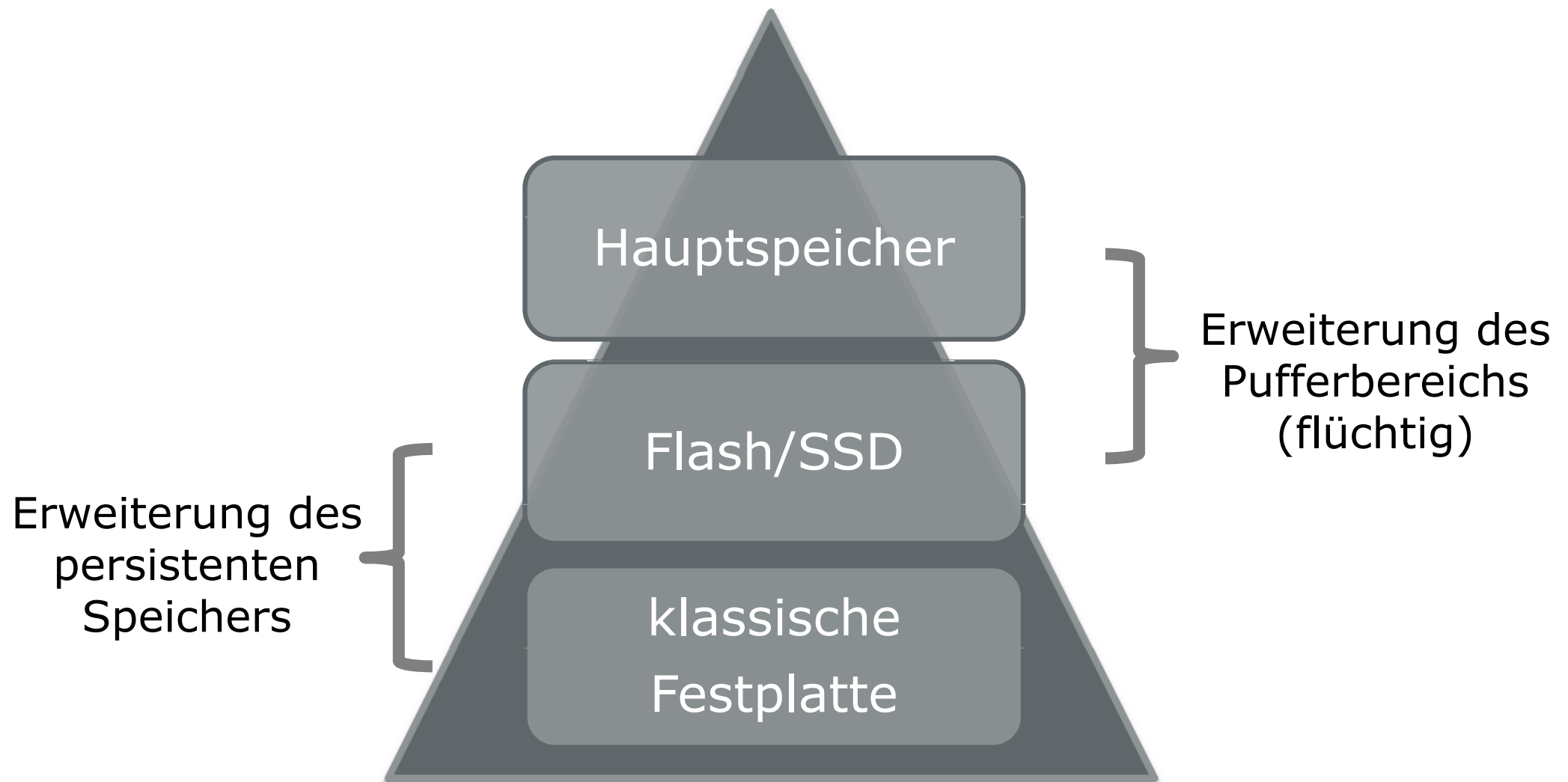
$$ReferenceInterval = \frac{128}{64} \cdot \frac{2000}{15} \cdot \frac{sec}{acc} = 266 \frac{sec}{acc} \approx \underline{\underline{5 \frac{min}{acc}}}$$

- Vergleich (1987)
  - 1KB-Seite, 400 sec/acc
  - 4KB-Seite, 100 sec/acc

Nochmal zehn Jahre später...  
2007

# Neue Ebene der Speicherhierarchie

16





# Neue Werte – Alte Formel

17

- Festplatte (4KB-Seiten), moderner Hauptspeicher:

$$\frac{256}{83} \cdot \frac{80}{0,047} \cdot \text{sec/acc} = 5248 \text{ sec/acc} \approx \underline{\underline{90 \text{ min/acc}}}$$

- alte Regel nur mit 64KB-Seiten noch gültig (**334 sec/acc**)

- SSD (4KB-Seiten), moderner Hauptspeicher:

$$\frac{256}{6300} \cdot \frac{999}{0,047} \cdot \text{sec/acc} = 876 \text{ sec/acc} \approx \underline{\underline{15 \text{ min/acc}}}$$

# Zwei neue Regeln

18

- billigere SSDs bereits verfügbar (ca. 400\$)
- 4KB-Seiten zwischen Hauptspeicher und SSD

$$\frac{256}{6300} \cdot \frac{400}{0,047} \cdot \text{sec/acc} = 351 \text{sec/acc} \approx \underline{\underline{6 \text{min/acc}}}$$

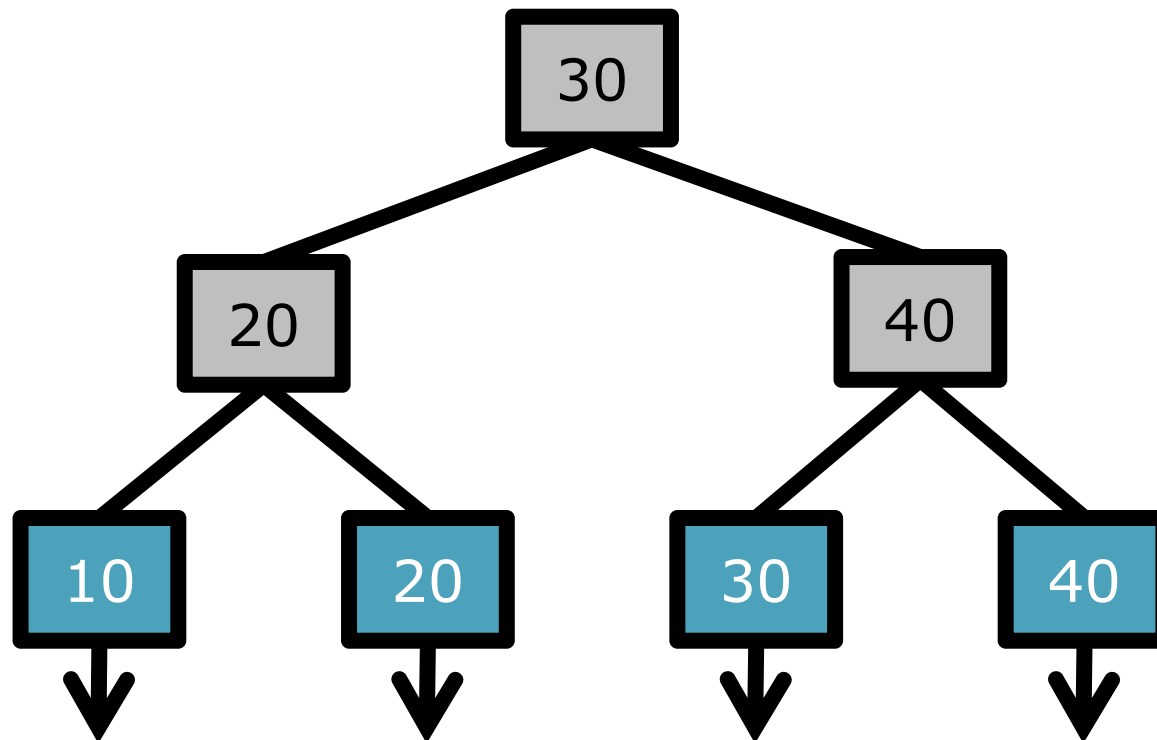
- 256KB-Seiten zwischen SSD und Festplatte

$$\frac{4}{83} \cdot \frac{80}{0,012} \cdot \text{sec/acc} = 321 \text{sec/acc} \approx \underline{\underline{5 \text{min/acc}}}$$

# Weitere Faustregeln der Speicherverwaltung

# Binäre Indexbäume

20



# Größe einer Indexseite

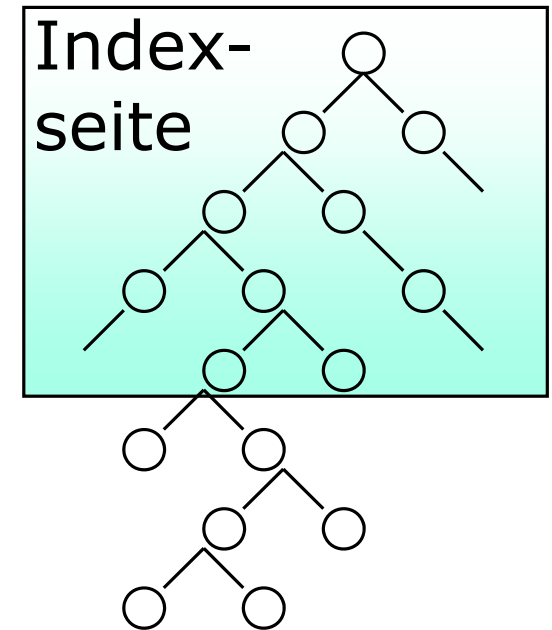
21

- **Bereitstellungskosten**
  - Lesen von Festplatte
- **"Fan-Out"**
  - Einträge pro Seite

- **Höhe** des Baumes in Indexseiten

$$IndexHeight \approx \frac{\log_2(N)}{\log_2(EntriesPerPage)}$$

- **Kosten? Nutzen? Optimale Größe?**



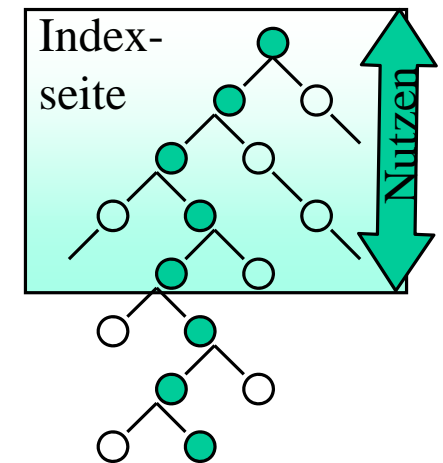
# "Index Page Utility"

22

- Wie viele **Schritte** bringt mich eine **Seite** näher zum **Ziel**?

$$IndexPageUtility = \log_2(EntriesPerPage)$$

- Anzahl der Ebenen im Baum pro Seite



- Beispiel

- Größe pro Eintrag: **20 Bytes**
- Größe einer Seite: **2000 Bytes**
- Nutzen einer Seite: **6.2**



ca. 70 Einträge  
bei 70% Füllstand

- **Zugriffskosten?**

# Kosten-Nutzen-Verhältnis

23

- Zugriffskosten (Lesen der Seite von Festplatte)

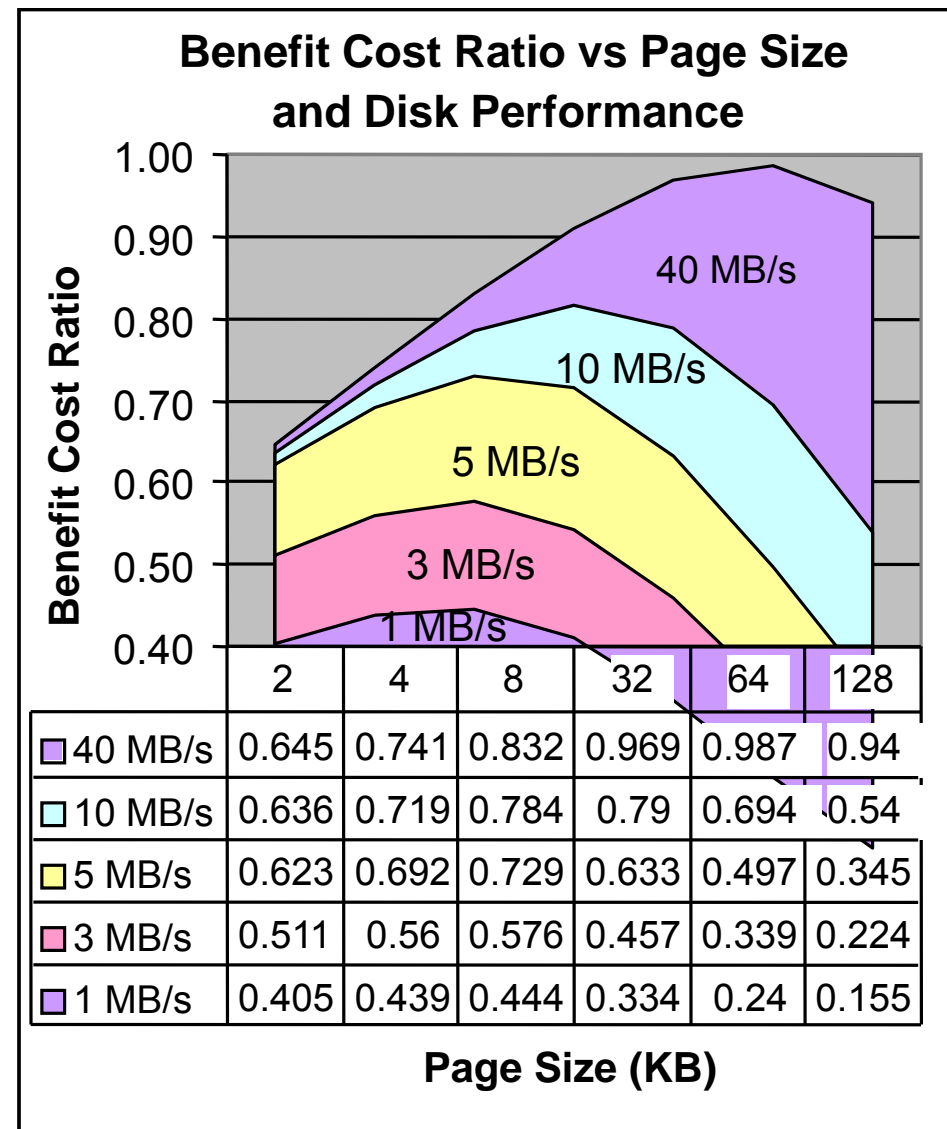
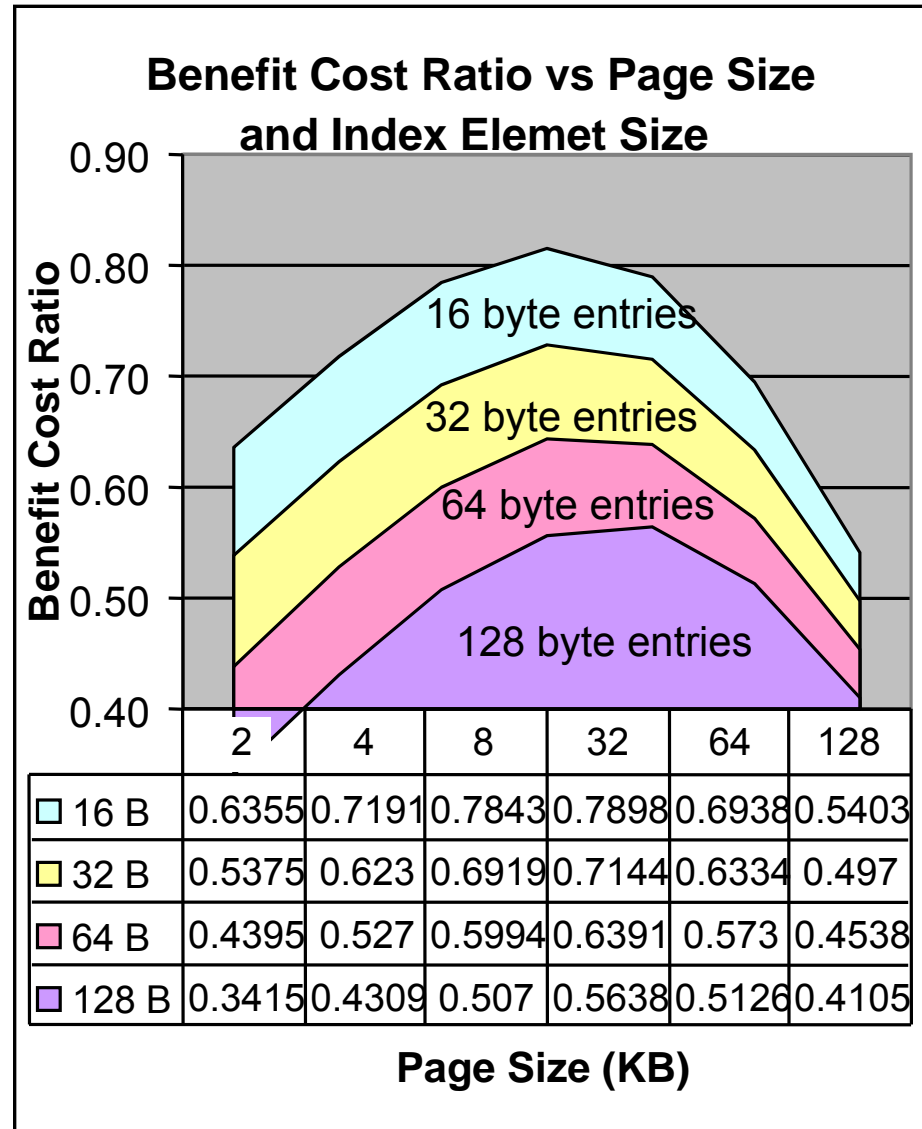
$$\mathit{IndexPageAccessCost} = \mathit{DiskLatency} + \frac{\mathit{PageSize}}{\mathit{DiskTransferRate}}$$

- Indikator für den wirtschaftlichen Wert der Seitengröße

$$\mathit{BenefitCostRatio} = \frac{\mathit{IndexPageUtility}}{\mathit{IndexPageAccessCost}}$$

# Optimale Größe einer Indexseite



24





# Aktuelle Werte (2007)

25

- Eintrag: 20 Bytes, Seite: 2000 Bytes
  
- Indexseiten auf **Festplatte**
  - Optimale Größe: **256 KB**  Neue 5MR für 256-KB-Seiten
  - Nutzen: **13**
  - Kosten-Nutzen-Verhältnis: **1.01**
  
- Indexseiten auf **Flashspeicher**
  - Optimale Größe: **2 KB**  Neue 5MR für 4-KB-Seiten
  - Nutzen: **6**
  - Kosten-Nutzen-Verhältnis: **46.1**
  
- **SB-Bäume** berücksichtigen zwei Seitengrößen, siehe [oneil92]

# Sequentieller Zugriff

26

- **höhere Bandbreite** bei sequentiellem Lesen der Daten
  - 10x langsamer bei wahlfreiem Zugriff
  
- andere Betrachtung für **sequentielle I/O-Operationen**
  - z. B. **Sort**, Cube, Hash-Join, Rollup
  
- **TPMMS** vs. One-Pass-Sort
  - doppelt so viel I/O
  - viel weniger Hauptspeicher
  
- **Wann ist es ratsam, einen One-Pass-Sort zu nutzen?**

# Neue Regel - Herleitung

27

- neuer Break-Even-Point (64-KB-Seiten)
  - höhere Transferrate (ca. 5 MB/s)
  - doppelte Intervallgröße für **Lesen und Schreiben**

$$ReferenceInterval = 2 \cdot \frac{PagesPerMBofRAM}{AccessPerSecondPerDisc} \cdot \frac{PricePerDiskDrive}{PricePerMBofRAM}$$

$$ReferenceInterval = 2 \cdot \frac{16}{80} \cdot \frac{2000}{15} \cdot \frac{sec}{acc} = 2 \cdot 26 \frac{sec}{acc} \approx \underline{\underline{1 \frac{min}{acc}}}$$

- bei hohem Datentransfer sind Zugriffskosten vernachlässigbar
  - nur Transferrate ausschlaggebend

# "One-Minute Sequential Rule" (1997)

28

**Sequentielle Operationen**, welche Daten innerhalb von **einer Minute wiederholt referenzieren**, sollten diese Daten im **Hauptspeicher** halten.

## ■ Beispiel

- One-Pass-Sort verarbeitet **5 GB/Minute**
  - ◇ variiert je nach Algorithmus
- **ab 5 GB** Relationsgröße **Two-Pass-Sort** nutzen

# Weitere Themenbereiche

29

- **Energiekosten** berücksichtigen
  - nicht nur Anschaffungskosten
- Hauptspeicher vs. CPU-Zyklen
  - "**10-Byte Rule**" (1987)
- Richtlinien für Datenbewegung in **Speicherhierarchie**
  - Manuell (Admin) vs. Automatisch (LRU)
- Nebeneffekte des Flashspeichers
  - Auswirkungen auf **paralleles Programmieren**
- Skalierbarkeit von **In-Memory-Datenbanken**
- generative **Garbage-Collection**

- [gray87] J. Gray, F. Putzolu, "The 5 Minute Rule for Trading Memory for Disc Access and The 10 Byte Rule for Trading Memory for CPU Time", ACM, 1987
- [gray97] J. Gray, G. Graefe, "The Five-Minute Rule Then Years Later and Other Computer Storage Rules of Thumb", SIGMOD Record, 1997
- [graefe08] G. Graefe, "The Five-minute rule – 20 Years Later and How Flash Memory Changes the Rules", ACM Queue, 2008
- [oneil92] P. E. O'Neil, "The SB-tree: An Index-Sequential Structure for High-Performance Sequential Access", Acta Informatica, 1992