

Verteilte Datenbanken

- **Motivation:**
geographisch verteilte Organisationsform einer Bank mit ihren Filialen
- **Filialen sollen Daten lokaler Kunden bearbeiten können**
- **Zentrale soll Zugriff auf alle Daten haben (z.B. für Kontogutschreibungen)**



Terminologie

Sammlung von Informationseinheiten, verteilt auf mehreren Rechnern, verbunden mittels Kommunikationsnetz → nach *Ceri & Pelagatti* (1984)

Kooperation zwischen autonom arbeitenden Stationen, zur Durchführung einer globalen Aufgabe

2

Kommunikationsmedien

LAN: local area network, z.B. Ethernet, Token-Ring oder FDDI-Netz

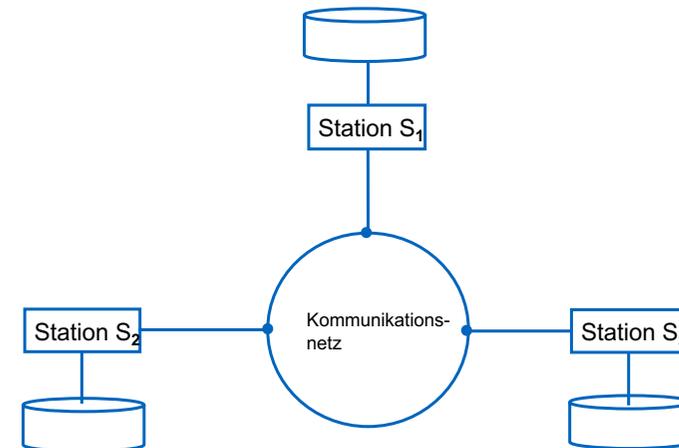
WAN: wide area network, z.B. das Internet

Telefonverbindungen, z.B. ISDN oder einfache Modem-Verbindungen

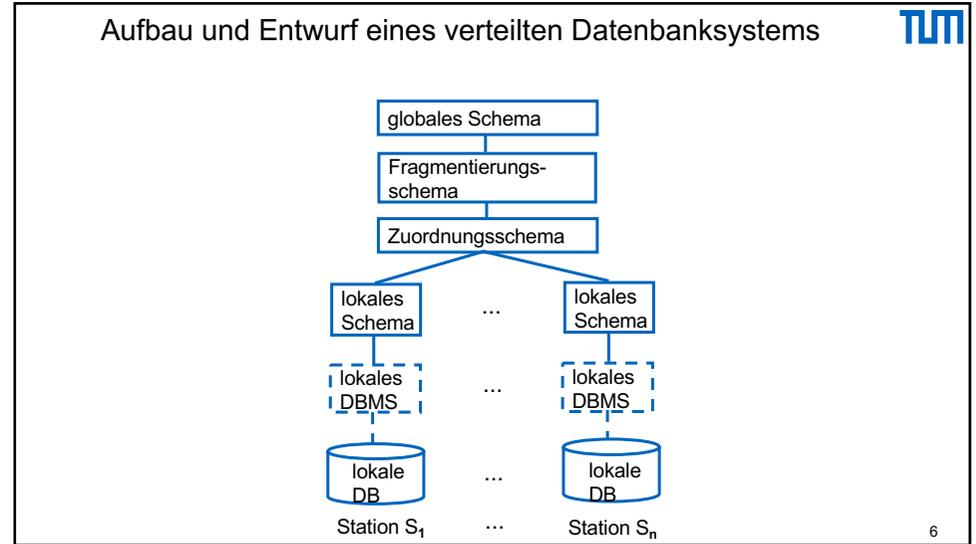
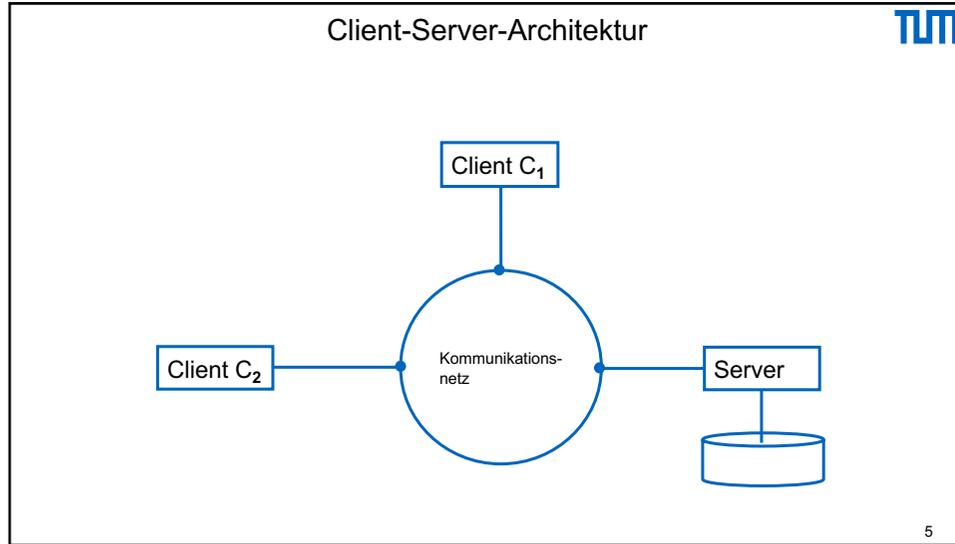


3

Verteiltes Datenbanksystem



4



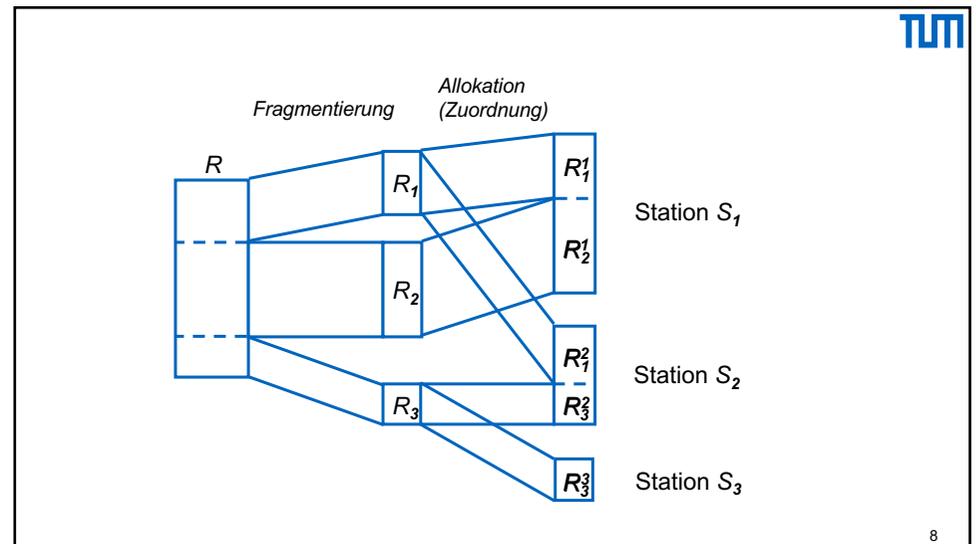
Fragmentierung und Allokation einer Relation

Fragmentierung: Fragmente enthalten Daten mit gleichem Zugriffsverhalten

Allokation: Fragmente werden den Stationen zugeordnet

- mit Replikation (redundanzfrei)
- ohne Replikation

7



Fragmentierung



horizontale Fragmentierung: Zerlegung der Relation in disjunkte Tupelmengen

vertikale Fragmentierung: Zusammenfassung von Attributen mit gleichem Zugriffsmuster

kombinierte Fragmentierung: Anwendung horizontaler und vertikaler Fragmentierung auf dieselbe Relation

Korrektheits-Anforderungen



Rekonstruierbarkeit

Vollständigkeit

Disjunktheit

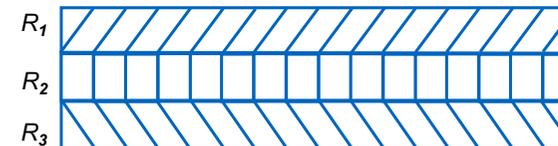
Beispielrelation Professoren



Professoren						
PersNr	Name	Rang	Raum	Fakultät	Gehalt	Steuerklasse
2125	Sokrates	C4	226	Philosophie	85000	1
2126	Russel	C4	232	Philosophie	80000	3
2127	Kopernikus	C3	310	Physik	65000	5
2133	Popper	C3	52	Philosophie	68000	1
2134	Augustinus	C3	309	Theologie	55000	5
2136	Curie	C4	36	Physik	95000	3
2137	Kant	C4	7	Philosophie	98000	1

abstrakte Darstellung:

R



Für 2 Prädikate p_1 und p_2 ergeben sich 4 Zerlegungen:

$$R_1 := \sigma_{p_1 \wedge p_2}(R)$$

$$R_2 := \sigma_{p_1 \wedge \neg p_2}(R)$$

$$R_3 := \sigma_{\neg p_1 \wedge p_2}(R)$$

$$R_4 := \sigma_{\neg p_1 \wedge \neg p_2}(R)$$



n Zerlegungsprädikate p_1, \dots, p_n ergeben 2^n Fragmente

sinnvolle Gruppierung der Professoren nach Fakultätszugehörigkeit: 

➔ 3 Zerlegungsprädikate:

$p_1 \equiv \text{Fakultät} = \text{'Theologie'}$
 $p_2 \equiv \text{Fakultät} = \text{'Physik'}$
 $p_3 \equiv \text{Fakultät} = \text{'Philosophie'}$

TheolProfs' := $\sigma_{p_1 \wedge \neg p_2 \wedge \neg p_3}$ (Professoren) = σ_{p_1} (Professoren)
 PhysikProfs' := $\sigma_{\neg p_1 \wedge p_2 \wedge \neg p_3}$ (Professoren) = σ_{p_2} (Professoren)
 PhiloProfs' := $\sigma_{\neg p_1 \wedge \neg p_2 \wedge p_3}$ (Professoren) = σ_{p_3} (Professoren)
 AndereProfs' := $\sigma_{\neg p_1 \wedge \neg p_2 \wedge \neg p_3}$ (Professoren)

13

Abgeleitete horizontale Fragmentierung 

Beispiel *Vorlesungen* aus dem Universitätsschema:
 Zerlegung in Gruppen mit gleicher SWS-Zahl

2SWSVorls := $\sigma_{\text{SWS}=2}$ (Vorlesungen)
 3SWSVorls := $\sigma_{\text{SWS}=3}$ (Vorlesungen)
 4SWSVorls := $\sigma_{\text{SWS}=4}$ (Vorlesungen)

➔ für Anfragebearbeitung schlechte Zerlegung

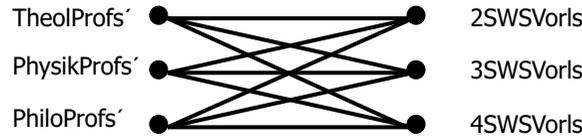
14

select Titel, Name
from Vorlesungen, Professoren
where gelesenVon = PersNr;

resultiert in:

$\Pi_{\text{Titel, Name}}((\text{TheolProfs}' \bowtie \text{2SWSVorls}) \cup$
 $(\text{TheolProfs}' \bowtie \text{3SWSVorls}) \cup$
 $\dots \cup (\text{PhiloProfs}' \bowtie \text{4SWSVorls}))$

Join-Graph zu diesem Problem:



15

Andere Join-Arten 

• natürlicher Join

L		
A	B	C
a ₁	b ₁	c ₁
a ₂	b ₂	c ₂

⋈

R		
C	D	E
c ₁	d ₁	e ₁
c ₃	d ₂	e ₂

=

Resultat				
A	B	C	D	E
a ₁	b ₁	c ₁	d ₁	e ₁

16

Semi-Joins

- Semi-Join von R mit L

L		
A	B	C
a ₁	b ₁	c ₁
a ₂	b ₂	c ₂

 \times

R		
C	D	E
c ₁	d ₁	e ₁
c ₃	d ₂	e ₂

 $=$

Resultat		
C	D	E
c ₁	d ₁	e ₁

- Semi-Join von L mit R

L		
A	B	C
a ₁	b ₁	c ₁
a ₂	b ₂	c ₂

 \times

R		
C	D	E
c ₁	d ₁	e ₁
c ₃	d ₂	e ₂

 $=$

Resultat		
A	B	C
a ₁	b ₁	c ₁

17

→ Lösung: abgeleitete Fragmentierung

TheolProf's' ●————● TheolVorles

PhysikProf's' ●————● PhysikVorles

PhiloProf's' ●————● PhiloVorles

TheolVorles := Vorlesungen $\bowtie_{\text{gelesenVon=PersNr}}$ TheolProf's'

PhysikVorles := Vorlesungen $\bowtie_{\text{gelesenVon=PersNr}}$ PhysikProf's'

PhiloVorles := Vorlesungen $\bowtie_{\text{gelesenVon=PersNr}}$ PhiloProf's'

$\Pi_{\text{Titel, Name}}((\text{TheolProf's}' \bowtie_p \text{TheolVorles}) \cup (\text{PhysikProf's}' \bowtie_p \text{PhysikVorles}) \cup (\text{PhiloProf's}' \bowtie_p \text{PhiloVorles}))$

mit $p \equiv (\text{PersNr} = \text{gelesenVon})$

18

Vertikale Fragmentierung

abstrakte Darstellung:

R

R_1 κ R_2

19

Vertikale Fragmentierung

Beliebige vertikale Fragmentierung gewährleistet **keine Rekonstruierbarkeit**

2 mögliche Ansätze, um Rekonstruierbarkeit zu garantieren:

- jedes Fragment enthält den Primärschlüssel der Originalrelation. Aber: Verletzung der *Disjunktheit*
- jedem Tupel der Originalrelation wird ein eindeutiges **Surrogat** (= künstlich erzeugter Objektindikator) zugeordnet, welches in jedes vertikale Fragment des Tupels mit aufgenommen wird

20

Vertikale Fragmentierung (Beispiel)



für die Universitätsverwaltung sind PersNr, Name, Gehalt und Steuerklasse interessant:

ProfVerw := Π PersNr, Name, Gehalt, Steuerklasse (Professoren)

für Lehre und Forschung sind dagegen PersNr, Name, Rang, Raum und Fakultät von Bedeutung:

Profs := Π PersNr, Name, Rang, Raum, Fakultät (Professoren)

Rekonstruktion der Originalrelation *Professoren*:

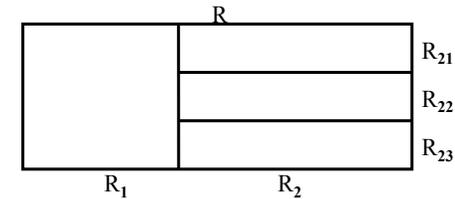
Professoren = ProfVerw $\bowtie_{\text{ProfVerw.PersNr} = \text{Profs.PersNr}}$ Profs

21

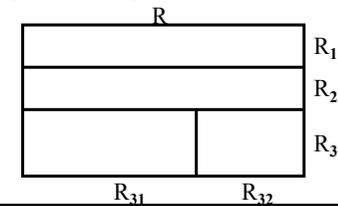
Kombinierte Fragmentierung



a) horizontale Fragmentierung nach vertikaler Fragmentierung



b) vertikale Fragmentierung nach horizontaler Fragmentierung



22

Rekonstruktion nach kombinierter Fragmentierung



Fall a)

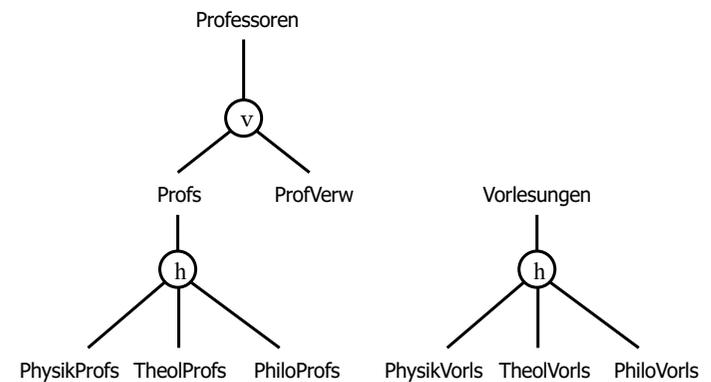
$$R = R_1 \bowtie_p (R_{21} \cup R_{22} \cup R_{23})$$

Fall b)

$$R = R_1 \cup R_2 \cup (R_{31} \bowtie_{R_{31,k} = R_{32,k}} R_{32})$$

23

Baumdarstellung der Fragmentierungen (Beispiel)



24

Allokation



Dasselbe Fragment kann mehreren Stationen zugeordnet werden

Allokation für unser Beispiel ohne Replikationen ⇒
redundanzfreie Zuordnung

Station	Bemerkung	zugeordnete Fragemente
S _{Verw}	Verwaltungsrechner	{ <i>ProfVerw</i> }
S _{Physik}	Dekanat Physik	{ <i>PhysikVorls, PhysikProfs</i> }
S _{Philo}	Dekanat Philosophie	{ <i>PhiloVorls, PhiloProfs</i> }
S _{Theol}	Dekanat Theologie	{ <i>TheolVorls, TheolProfs</i> }

25

Transparenz in verteilten Datenbanken



Grad der Unabhängigkeit den ein VDBMS dem Benutzer beim Zugriff auf verteilte Daten vermittelt

verschiedene Stufen der Transparenz:

- ◆ Fragmentierungstransparenz
- ◆ Allokationstransparenz
- ◆ Lokale Schema-Transparenz

26

Fragmentierungstransparenz



Beispielanfrage, die Fragmentierungstransparenz voraussetzt:

```
select Titel, Name
from Vorlesungen, Professoren
where gelesenVon = PersNr;
```

Beispiel für eine Änderungsoperation, die Fragmentierungstransparenz voraussetzt:

```
update Professoren
  set Fakultät = ‚Theologie‘
  where Name = ‚Sokrates‘;
```

27

Fortsetzung Beispiel



Ändern des Attributwertes von *Fakultät*

- Transferieren des Sokrates-Tupels aus Fragment *PhiloProfs* in das Fragment *TheolProfs* (= Löschen aus *PhiloProfs*, Einfügen in *TheolProfs*)
- Ändern der abgeleiteten Fragmentierung von *Vorlesungen* (= Einfügen der von Sokrates gehaltenen Vorlesungen in *TheolVorls*, Löschen der von ihm gehaltenen Vorlesungen aus *PhiloVorls*)

28

Allokationstransparenz



Benutzer müssen Fragmentierung kennen, aber nicht den „Aufenthaltort“ eines Fragments

Beispielanfrage:

```
select Gehalt
from ProfVerw
where Name = ‚Sokrates‘;
```

29

Allokationstransparenz (Forts.)



unter Umständen muss Originalrelation rekonstruiert werden

Beispiel:

Verwaltung möchte wissen, wieviel die C4-Professoren der Theologie insgesamt verdienen

da Fragmentierungstransparenz fehlt muss die Anfrage folgendermaßen formuliert werden:

```
select sum (Gehalt)
from ProfVerw, TheolProfs
where ProfVerw.PersNr = TheolProfs.PersNr and
      Rang = ‚C4‘;
```

30

Lokale Schema-Transparenz



Der Benutzer muss auch noch den Rechner kennen, auf dem ein Fragment liegt.

Beispielanfrage:

```
select Name
from TheolProfs at STheol
where Rang = ‚C3‘;
```

31

Lokale Schema-Transparenz (Forts.)



Ist überhaupt Transparenz gegeben?

Lokale Schema-Transparenz setzt voraus, dass alle Rechner dasselbe Datenmodell und dieselbe Anfragesprache verwenden.

⇒ vorherige Anfrage kann somit analog auch an Station *S_{Philo}* ausgeführt werden

Dies ist nicht möglich bei Kopplung unterschiedlicher DBMS.

Verwendung grundsätzlich verschiedener Datenmodelle auf lokalen DBMS nennt man „Multi-Database-Systems“ (oft unumgänglich in „realer“ Welt).

32

Anfrageübersetzung und Anfrageoptimierung TUM

Voraussetzung: Fragmentierungstransparenz

Aufgabe des Anfrageübersetzers: Generierung eines Anfrageauswertungsplans auf den Fragmenten

Aufgabe des Anfrageoptimierers: Generierung eines möglichst effizienten Auswertungsplanes
 → abhängig von der Allokation der Fragmente auf den verschiedenen Stationen des Rechnernetzes

33

Anfragebearbeitung bei horizontaler Fragmentierung TUM

Übersetzung einer SQL-Anfrage auf dem globalen Schema in eine äquivalente Anfrage auf den Fragmenten benötigt 2 Schritte:

1. Rekonstruktion aller in der Anfrage vorkommenden globalen Relationen aus den Fragmenten, in die sie während der Fragmentierungsphase zerlegt wurden. Hierfür erhält man einen algebraischen Ausdruck.
2. Kombination des Rekonstruktionsausdrucks mit dem algebraischen Anfrageausdruck, der sich aus der Übersetzung der SQL-Anfrage ergibt.

34

Beispiel TUM

```

select Titel
from Vorlesungen, Profs
where gelesenVon = PersNr and
        Rang = ,C4';
    
```

Der entstandene algebraische Ausdruck heißt **kanonische Form** der Anfrage:

35

Algebraische Äquivalenzen TUM

Für eine effizientere Abarbeitung der Anfrage benutzt der Anfrageoptimierer die folgende Eigenschaft:

$$(R_1 \cup R_2) \bowtie_p (S_1 \cup S_2) = (R_1 \bowtie_p S_1) \cup (R_1 \bowtie_p S_2) \cup (R_2 \bowtie_p S_1) \cup (R_2 \bowtie_p S_2)$$

Die Verallgemeinerung auf n horizontale Fragmente R_1, \dots, R_n von R und m Fragmente S_1, \dots, S_m von S ergibt:

$$(R_1 \cup \dots \cup R_n) \bowtie_p (S_1 \cup \dots \cup S_m) = \bigcup_{1 \leq i \leq n} \bigcup_{1 \leq j \leq m} (R_i \bowtie_p S_j)$$

Falls gilt: $S_i = S \bowtie_p R_i$ mit $S = S_1 \cup \dots \cup S_n$
 Dann gilt immer:

$$R \bowtie_p S = \bigcup_{1 \leq i \leq m} (R_i \bowtie_p S_i)$$

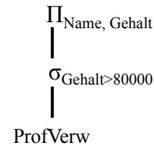
36

Optimierung bei vertikaler Fragmentierung



Für unser Beispiel gilt:

Alle notwendigen Informationen sind in *ProfVerw* enthalten \Rightarrow der Teil mit Vereinigung und Join kann „abgeschnitten“ werden. Das ergibt den folgenden optimierten Auswertungsplan:



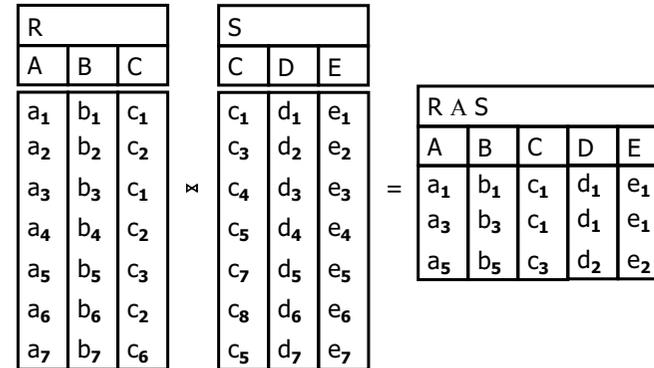
Beispiel für eine schlecht zu optimierende Anfrage:
(Attribut *Rang* fehlt in *ProfVerw*)

```

select Name, Gehalt, Rang
from Professoren
where Gehalt > 80000;
    
```

41

Der natürliche Verbund zweier Relationen R und S



42

Join-Auswertung in VDBMS



spielt kritischere Rolle als in zentralisierten Datenbanken

Problem: Argumente eines Joins zweier Relationen können auf unterschiedlichen Stationen des VDBMS liegen

2 Möglichkeiten: Join-Auswertung mit und ohne Filterung

43

Join-Auswertung



Betrachtung des allgemeinsten Falles:

äußere Argumentrelation *R* ist auf Station *St_R* gespeichert

innere Argumentrelation *S* ist dem Knoten *St_S* zugeordnet

Ergebnis der Joinberechnung wird auf einem dritten Knoten *St_{Result}* benötigt

44

Join-Auswertung ohne Filterung



Nested-Loops

Transfer einer Argumentrelation

Transfer beider Argumentrelationen

45

Nested-Loops



Iteration durch die äußere Relation R mittels Laufvariable r und Anforderung des/der zu jedem Tupel r passenden Tupel $s \in S$ mit $r.C = s.C$ (über Kommunikationsnetz bei St_g)

Diese Vorgehensweise benötigt pro Tupel aus R eine Anforderung und eine passende Tupelmenge aus S (welche bei vielen Anforderungen leer sein könnte)

⇒ es werden $2 * |R|$ Nachrichten benötigt

46

Der natürliche Verbund zweier Relationen R und S



R		
A	B	C
a ₁	b ₁	c ₁
a ₂	b ₂	c ₂
a ₃	b ₃	c ₁
a ₄	b ₄	c ₂
a ₅	b ₅	c ₃
a ₆	b ₆	c ₂
a ₇	b ₇	c ₆

S		
C	D	E
c ₁	d ₁	e ₁
c ₃	d ₂	e ₂
c ₄	d ₃	e ₃
c ₅	d ₄	e ₄
c ₇	d ₅	e ₅
c ₈	d ₆	e ₆
c ₅	d ₇	e ₇

47

Transfer einer Argumentrelation



1. vollständiger Transfer einer Argumentrelation (z.B. R) zum Knoten der anderen Argumentrelation
2. Ausnutzung eines möglicherweise auf $S.C$ existierenden Indexes

48

Der natürliche Verbund zweier Relationen R und S



R		
A	B	C
a ₁	b ₁	c ₁
a ₂	b ₂	c ₂
a ₃	b ₃	c ₁
a ₄	b ₄	c ₂
a ₅	b ₅	c ₃
a ₆	b ₆	c ₂
a ₇	b ₇	c ₆

S		
C	D	E
c ₁	d ₁	e ₁
c ₃	d ₂	e ₂
c ₄	d ₃	e ₃
c ₅	d ₄	e ₄
c ₇	d ₅	e ₅
c ₈	d ₆	e ₆
c ₅	d ₇	e ₇

49

Transfer beider Argumentrelationen



1. Transfer beider Argumentrelationen zum Rechner St_{Result}
2. Berechnung des Ergebnisses auf dem Knoten St_{Result} mittels
 - a) Merge-Join (bei vorliegender Sortierung)
 - oder
 - b) Hash-Join (bei fehlender Sortierung)
 - evtl. Verlust der vorliegenden Indexe für die Join-Berechnung
 - kein Verlust der Sortierung der Argumentrelation(en)

50

Join-Auswertung mit Filterung



Verwendung des Semi-Join-Operators zur Filterung

Schlüsselidee: nur Transfer von Tupeln, die passenden Join-Partner haben

Benutzung der folgenden algebraischen Eigenschaften:

- $R \bowtie S = R \bowtie (R \bowtie S)$
- $R \bowtie S = \Pi_C(R) \bowtie S$

51

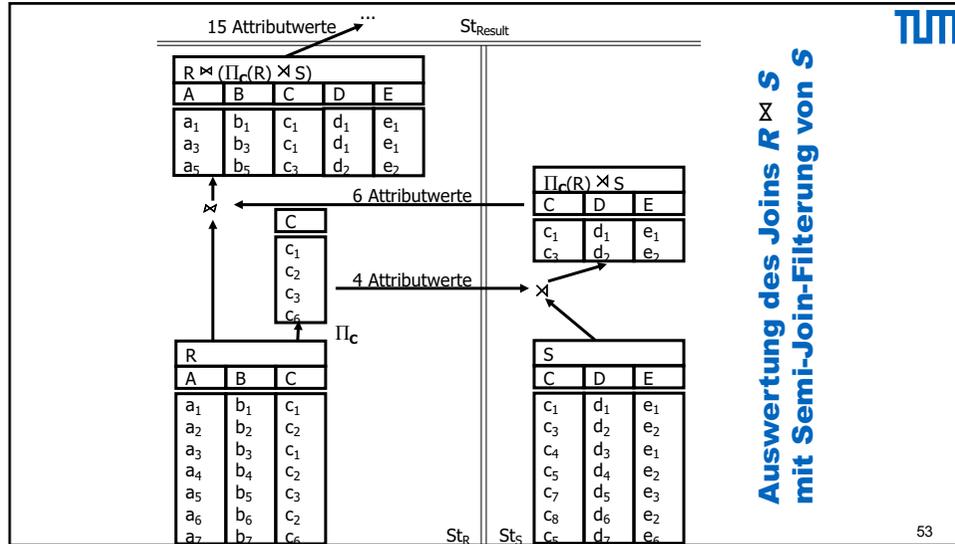
Join-Auswertung mit Filterung (Beispiel, Filterung der Relation S)



1. Transfer der unterschiedlichen C-Werte von $R (= \Pi_C(R))$ nach St_S
2. Auswertung des Semi-Joins $R \bowtie S = \Pi_C(R) \bowtie S$ auf St_S und Transfer nach St_R
3. Auswertung des Joins auf St_R , der nur diese transferierten Ergebnistupel des Semi-Joins braucht

Transferkosten werden nur reduziert, wenn gilt:
 $\|\Pi_C(R)\| + \|R \bowtie S\| < \|S\|$
 mit $\|R\| =$ Größe (in Byte) einer Relation

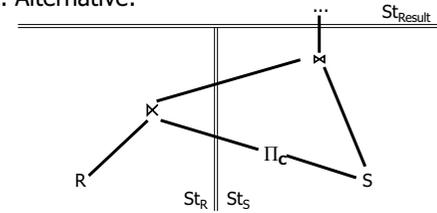
52



53

Alternative Auswertungspläne

1. Alternative:



2. Alternative:

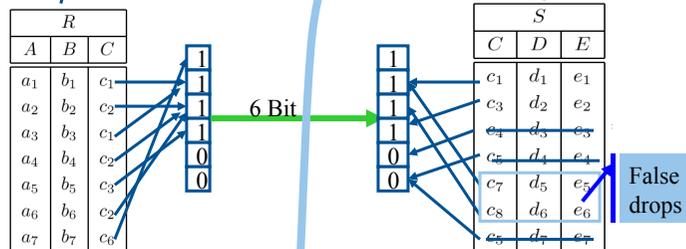
$$(R \bowtie \Pi_C(S)) \bowtie (\Pi_C(R) \bowtie S)$$

54

Join mit Hashfilter (Bloom-Filter)

15 Attribute

12 Attribute (4 Tupel inkl. 2 False Drops)



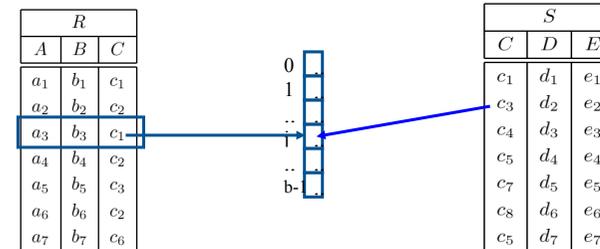
55

Join mit Hashfilter

(False Drop Abschätzung)

Wahrscheinlichkeit, dass ein bestimmtes Bit j gesetzt ist

- W. dass ein bestimmtes $r \in R$ das Bit setzt: $1/b$
- W. dass kein $r \in R$ das Bit setzt: $(1-1/b)^{|R|}$
- W. dass ein $r \in R$ das Bit gesetzt hat: $1 - (1-1/b)^{|R|}$



56

Join mit Hashfilter



(False Drop Abschätzung)

W. dass irgendein $r \in R$ ein bestimmtes Bit gesetzt hat: $1 - (1-1/b)^{|R|}$

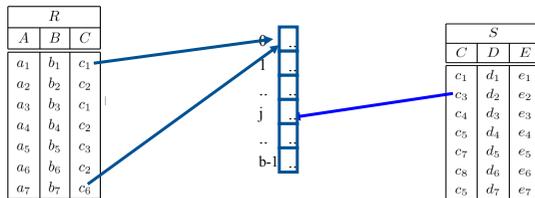
Wieviele Bits sind gesetzt?

- $b * [1 - (1-1/b)^{|R|}]$

Mehrere $r \in R$ können dasselbe Bit setzen

Approximation: alle $r \in R$ setzen unterschiedliche Bits

- W. dass ein bestimmtes Bit j gesetzt ist: $|R| / b$
- $b \gg |R|$



57

Join mit Hashfilter



(False Drop Abschätzung)

W. dass irgendein $r \in R$ ein bestimmtes Bit gesetzt hat:

- $1 - (1-1/b)^{|R|}$

W. dass ein bestimmtes $s \in S$ ausgewählt wird:

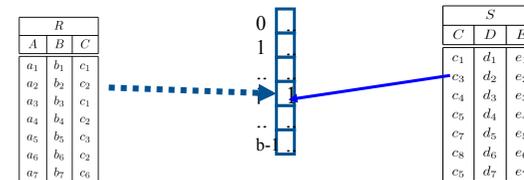
- $1 - (1-1/b)^{|R|}$

Wieviele $s \in S$ werden ausgewählt?

- $|S| * [1 - (1-1/b)^{|R|}]$

Approximation: alle r setzen unterschiedliche Bits

- W. dass ein bestimmtes Bit j gesetzt ist: $|R| / b$
- $|S| * (|R|/b)$ Elemente aus S werden ausgewählt



58

Parameter für die Kosten eines Auswertungsplan



Kardinalitäten von Argumentrelationen

Selektivitäten von Joins und Selektionen

Transferkosten für Datenkommunikation (Verbindungsaufbau + von Datenvolumen abhängiger Anteil für Transfer)

Auslastung der einzelnen VDBMS-Stationen

Effektive Anfrageoptimierung muss auf Basis eines Kostenmodells durchgeführt werden und soll mehrere Alternativen für unterschiedliche Auslastungen des VDBMS erzeugen.

59

Transaktionskontrolle in VDBMS



Transaktionen können sich bei VDBMS über mehrere Rechnerknoten erstrecken

⇒ Recovery:

- ♦ Redo: Wenn eine Station nach einem Fehler wieder anläuft, müssen alle Änderungen einmal abgeschlossener Transaktionen - seien sie lokal auf dieser Station oder global über mehrere Stationen ausgeführt worden - auf den an dieser Station abgelegten Daten wiederhergestellt werden.
- ♦ Undo: Die Änderungen noch nicht abgeschlossener lokaler und globaler Transaktionen müssen auf den an der abgestürzten Station vorliegenden Daten rückgängig gemacht werden.

60

EOT-Behandlung

Die EOT (End-of-Transaction)-Behandlung von globalen Transaktionen stellt in VDBMS ein Problem dar.

Eine globale Transaktion muss atomar beendet werden, d.h. entweder

- **commit:** globale Transaktion wird an allen (relevanten) lokalen Stationen festgeschrieben
- oder
- **abort:** globale Transaktion wird gar nicht festgeschrieben

⇒ Problem in verteilter Umgebung, da die Stationen eines VDBMS unabhängig voneinander „abstürzen“ können

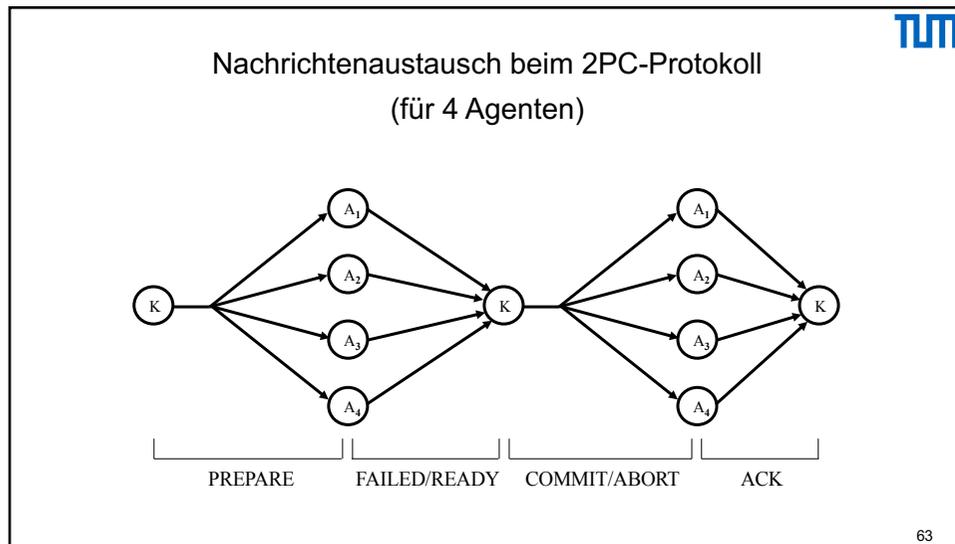
61

Problemlösung: Zweiphasen-Commit-Protokoll

→ gewährleistet die Atomarität der EOT-Behandlung

das 2PC-Verfahren wird von sog. Koordinator K überwacht und gewährleistet, dass die n Agenten (=Stationen im VDBMS) A_1, \dots, A_n , die an einer Transaktion beteiligt waren, entweder alle von Transaktion T geänderten Daten festschreiben oder alle Änderungen von T rückgängig machen

62



Ablauf der EOT-Behandlung beim 2PC-Protokoll

K schickt allen Agenten eine **PREPARE**-Nachricht, um herauszufinden, ob sie Transaktionen festschreiben können

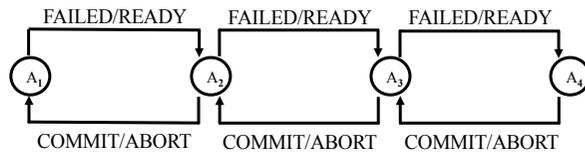
jeder Agent A_i empfängt **PREPARE**-Nachricht und schickt eine von zwei möglichen Nachrichten an K :

- **READY**, falls A_i in der Lage ist, die Transaktion T lokal festzuschreiben
- **FAILED**, falls A_i kein **commit** durchführen kann (wegen Fehler, Inkonsistenz etc.)

- hat K von **allen** n Agenten A_1, \dots, A_n ein **READY** erhalten, kann K ein **COMMIT** an alle Agenten schicken mit der Aufforderung, die Änderungen von T lokal festzuschreiben; antwortet einer der Agenten mit **FAILED** od. gar nicht innerhalb einer bestimmten Zeit (*timeout*), schickt K ein **ABORT** an alle Agenten und diese machen die Änderungen der Transaktion rückgängig
- haben die Agenten ihre lokale EOT-Behandlung abgeschlossen, schicken sie eine **ACK**-Nachricht (=acknowledgement, dt. Bestätigung) an den Koordinator

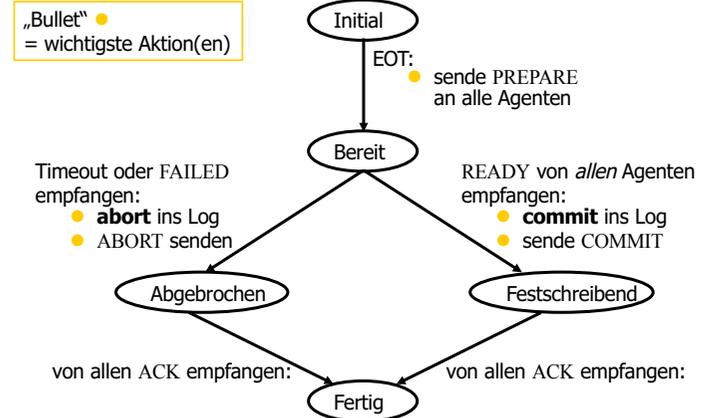
64

Lineare Organisationsform beim 2PC-Protokoll



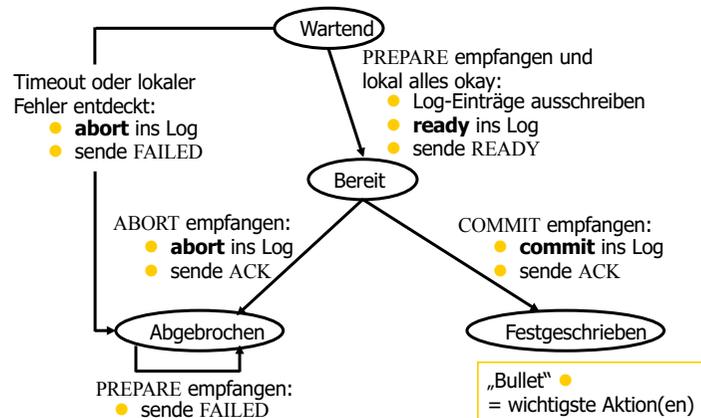
65

Zustandsübergang beim 2PC-Protokoll: Koordinator



66

Zustandsübergang beim 2PC-Protokoll: Agent



67

Fehlersituationen des 2PC-Protokolls



- Absturz eines Koordinators
- Absturz eines Agenten
- verlorengegangene Nachrichten

68

Absturz eines Koordinators



- Absturz vor dem Senden einer COMMIT-Nachricht → Rückgängigmachen der Transaktion durch Versenden einer ABORT-Nachricht
- Absturz nachdem Agenten ein READY mitgeteilt haben → **Blockierung der Agenten**
 - ⇒ Hauptproblem des 2PC-Protokolls beim Absturz des Koordinators, da dadurch die Verfügbarkeit des Agenten bezüglich andere globaler und lokaler Transaktionen drastisch eingeschränkt ist

Um Blockierung von Agenten zu verhindern, wurde ein **Dreiphasen-Commit-Protokoll** konzipiert, das aber in der Praxis zu aufwendig ist (VDBMS benutzen das 2PC-Protokoll).

69

Absturz eines Agenten



- antwortet ein Agent innerhalb eines Timeout-Intervalls nicht auf die PREPARE-Nachricht, gilt der Agent als abgestürzt; der Koordinator bricht die Transaktion ab und schickt eine ABORT-Nachricht an alle Agenten
- „abgestürzter“ Agent schaut beim Wiederanlauf in seine Log-Datei:
 - kein **ready**-Eintrag bzgl. Transaktion T → Agent führt ein abort durch und teilt dies dem Koordinator mit (*FAILED*-Nachricht)
 - **ready**-Eintrag aber kein **commit**-Eintrag → Agent fragt Koordinator, was aus Transaktion T geworden ist; Koordinator teilt *COMMIT* oder *ABORT* mit, was beim Agenten zu einem Redo oder Undo der Transaktion führt
 - **commit**-Eintrag vorhanden → Agent weiß ohne Nachfragen, dass ein (lokales) Redo der Transaktion nötig ist

70

Verlorengegangene Nachrichten



- PREPARE-Nachricht des Koordinators an einen Agenten geht verloren **oder**
- READY-(oder FAILED-)Nachricht eines Agenten geht verloren
 - nach Timeout-Intervall geht Koordinator davon aus, dass betreffender Agent nicht funktionsfähig ist und sendet ABORT-Nachricht an alle Agenten (Transaktion gescheitert)
- Agent erhält im Zustand **Bereit** keine Nachricht vom Koordinator
 - Agent ist blockiert, bis COMMIT- oder ABORT-Nachricht vom Koordinator kommt, da Agent nicht selbst entscheiden kann (deshalb schickt Agent eine „Erinnerung“ an den Koordinator)

71

Mehrbenutzersynchronisation in VDBMS



Serialisierbarkeit

Zwei-Phasen-Sperrprotokoll in VDBMS

- lokale Sperrverwaltung an jeder Station
- globale Sperrverwaltung

72

Serialisierbarkeit



Lokale Serialisierbarkeit an jeder der an den Transaktionen beteiligten Stationen reicht nicht aus. Deshalb muß man bei der Mehrbenutzersynchronisation auf **globaler Serialisierbarkeit** bestehen.

Beispiel (lokal serialisierbare Historien):

S ₁			S ₂		
Schritt	T ₁	T ₂	Schritt	T ₁	T ₂
1.	r(A)				
2.		w(A)			
			3.		w(B)
			4.	r(B)	

T₁ ↔ T₂

73

Lokale Sperrverwaltung



globale Transaktion muss vor Zugriff/Modifikation eines Datums A, das auf Station S liegt, eine Sperre vom Sperrverwalter der Station S erwerben

Verträglichkeit der angeforderten Sperre mit bereits existierenden Sperrern kann lokal entschieden werden → favorisiert lokale Transaktionen, da diese nur mit ihrem lokalen Sperrverwalter kommunizieren müssen

74

Globale Sperrverwaltung



= alle Transaktionen fordern alle Sperrern an einer einzigen, ausgezeichneten Station an.

Nachteile:

zentraler Sperrverwalter kann zum Engpass des VDBMS werden, besonders bei einem Absturz der Sperrverwalter-Station („rien ne vas plus“)

Verletzung der lokalen Autonomie der Stationen, da auch lokale Transaktionen ihre Sperrern bei der zentralisierten Sperrverwaltung anfordern müssen

→ zentrale Sperrverwaltung i.a. **nicht** akzeptabel

75

Deadlocks in VDBMS



Erkennung von Deadlocks (Verklemmungen)

- zentralisierte Deadlock-Erkennung
- dezentrale (verteilte) Deadlock-Erkennung

- Vermeidung von Deadlocks

76

„Verteilter“ Deadlock



S ₁			S ₂		
Schritt	T ₁	T ₂	Schritt	T ₁	T ₂
0.	BOT				
1.	lockS(A)				
2.	r(A)				
6.		lockX(A) ~~~~~	3.		BOT
			4.		lockX(B)
			5.		w(B)
			7.	lockS(B) ~~~~~	

77

Timeout



betreffende Transaktion wird zurückgesetzt und erneut gestartet
→ einfach zu realisieren

Problem: richtige Wahl des Timeout-Intervalls:

- zu lang → schlechte Ausnutzung der Systemressourcen
- zu kurz → Deadlock-Erkennung, wo gar keine Verklemmung vorliegt

78

Zentralisierte Deadlock-Erkennung



Stationen melden lokal vorliegende Wartebeziehungen an neutralen Knoten, der daraus globalen Wartegraphen aufbaut (Zyklus im Graphen → Deadlock)

→ sichere Lösung

Nachteile:

- hoher Aufwand (viele Nachrichten)
- Entstehung von Phantom-Deadlocks (=nicht-existierende Deadlocks) durch „Überholen“ von Nachrichten im Kommunikationssystem

79

Dezentrale Deadlock-Erkennung



lokale Wartegraphen an den einzelnen Stationen →

Erkennen von lokalen Deadlocks

Erkennung globaler Deadlocks:

- jeder lokale Wartegraph hat einen Knoten *External*, der stationenübergreifenden Wartebeziehungen zu externen Subtransaktionen modelliert
- Zuordnung jeder Transition zu einem Heimatknoten, von wo aus *externe Subtransaktionen* auf anderen Stationen initiiert werden

Die Kante $External \rightarrow T_i$
wird für jede „von außen“ kommende Transaktion T_i eingeführt.

Die Kante $T_j \rightarrow External$
wird für jede von außen kommende Transaktion T_j dieser Station eingeführt, falls T_j „nach außen“ geht.

80

Beispiel:

S_1 Heimatknoten von T_1 , S_2 Heimatknoten von T_2

Wartegraphen:

S_1 : $External \rightarrow T_2 \rightarrow T_1 \rightarrow External$

S_2 : $External \rightarrow T_1 \rightarrow T_2 \rightarrow External$

↓

S_2 : $External \leftrightarrow T_1 \leftrightarrow T_2 \leftrightarrow External$

$T_1 \rightarrow T_2 \rightarrow T_1$
 $T_2 \rightarrow T_1 \rightarrow T_2$

Zur Reduzierung des Nachrichtenaufkommens wird der Pfad
 $External \rightarrow T_1' \rightarrow T_2' \rightarrow \dots \rightarrow T_n' \rightarrow External$
nur weitergereicht, wenn T_1' einen kleineren Identifikator als
 T_n' hat (= path pushing).

TUM

81

Deadlock-Vermeidung

TUM

optimistische Mehrbenutzersynchronisation:
nach Abschluss der Transaktionsbearbeitung wird Validierung durchgeführt

Zeitstempel-basierende Synchronisation:
Zuordnung eines Lese-/Schreib-Stempels zu jedem Datum
entscheidet, ob beabsichtigte Operation durchgeführt werden kann ohne
Serialisierbarkeit zu verletzen oder ob Transaktion abgebrochen wird (**abort**)

82

Deadlockvermeidung bei Sperrbasierte Synchronisation

TUM

wound/wait:

- nur jüngere Transaktionen warten auf ältere;
- fordert ältere Transaktion Sperre an, die mit der von der jüngeren Transaktion gehaltenen nicht verträglich ist, wird jüngere Transaktion abgebrochen

wait/die:

- nur ältere Transaktionen warten auf jüngere;
- fordert jüngere Transaktion Sperre an, die mit der von der älteren Transaktion gehaltenen nicht kompatibel ist, wird jüngere Transaktion abgebrochen

83

Voraussetzungen für Deadlockvermeidungsverfahren

TUM

Vergabe global eindeutiger Zeitstempel als
Transaktionsidentifikatoren

lokale Zeit Stations-ID

- lokale Uhren müssen hinreichend genau
aufeinander abgestimmt sein

84

Synchronisation bei replizierten Daten

ROWA: Read One, Write All**Problem:**

Zu einem Datum A gibt es mehrere Kopien A_1, A_2, \dots, A_n , die auf unterschiedlichen Stationen liegen.

Eine Lesetransaktion erfordert nur eine Kopie, bei Änderungstransaktionen müssen aber alle bestehenden Kopien geändert werden.

⇒ hohe Laufzeit und Verfügbarkeitsprobleme



85

Quorum-Consensus Verfahren

Ausgleich der Leistungsfähigkeit zwischen Lese- und Änderungstransaktionen

→ teilweise Verlagerung des Overheads von den Änderungs- zu den Lesetransaktionen indem den Kopien A_i eines replizierten Datums A individuelle Gewichte zugeordnet werden

- *Lesequorum* $Q_r(A)$
- *Schreibquorum* $Q_w(A)$

Folgende Bedingungen müssen gelten:

1. $Q_w(A) + Q_w(A) > W(A)$
2. $Q_r(A) + Q_w(A) > W(A)$



86

Beispiel

Station (S_i)	Kopie (A_i)	Gewicht (w_i)
S_1	A_1	3
S_2	A_2	1
S_3	A_3	2
S_4	A_4	2

$$W(A) = \sum_{i=1}^4 w_i(A) = 8$$

$$Q_r(A) = 4$$

$$Q_w(A) = 5$$



87

Zustände

a) vor dem Schreiben eines Schreibquorums

Station	Kopie	Gewicht	Wert	Versions#
S_1	A_1	3	1000	1
S_2	A_2	1	1000	1
S_3	A_3	2	1000	1
S_4	A_4	2	1000	1

b) nach dem Schreiben eines Schreibquorums

Station	Kopie	Gewicht	Wert	Versions#
S_1	A_1	3	1100	2
S_2	A_2	1	1000	1
S_3	A_3	2	1100	2
S_4	A_4	2	1000	1



88

Peer to Peer-Informationssysteme



Seti@Home

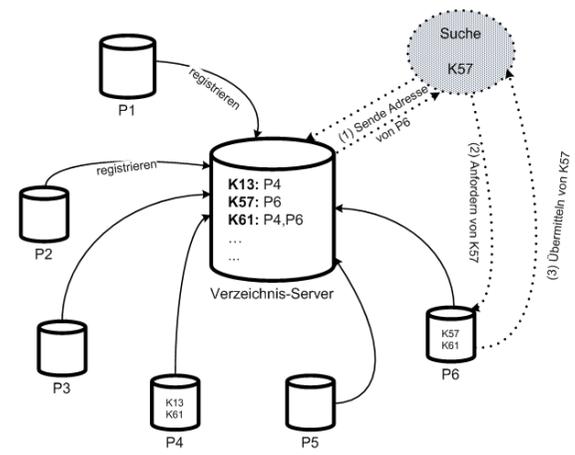
- P2P number crunching

Napster

- P2P file sharing / Informationsmanagement

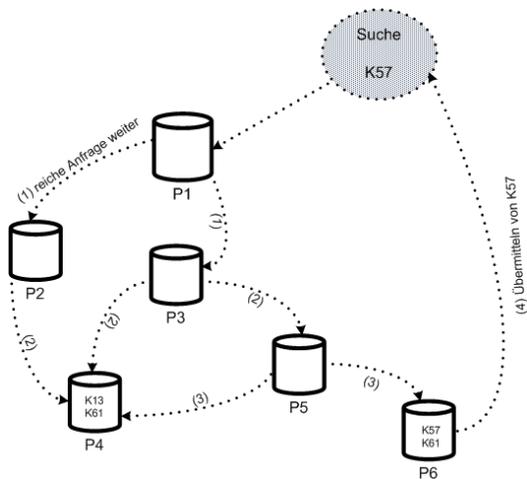
89

Napster-Architektur



90

Gnutella-Architektur



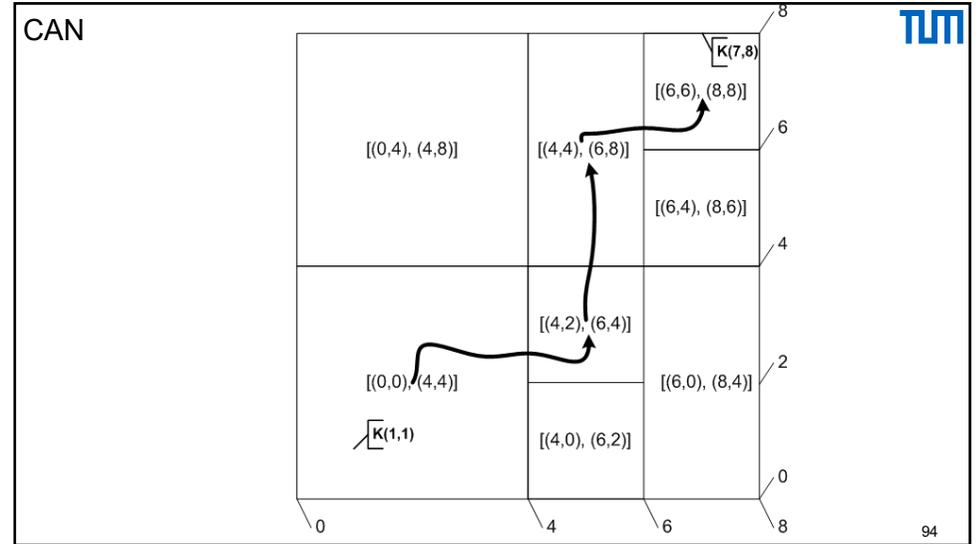
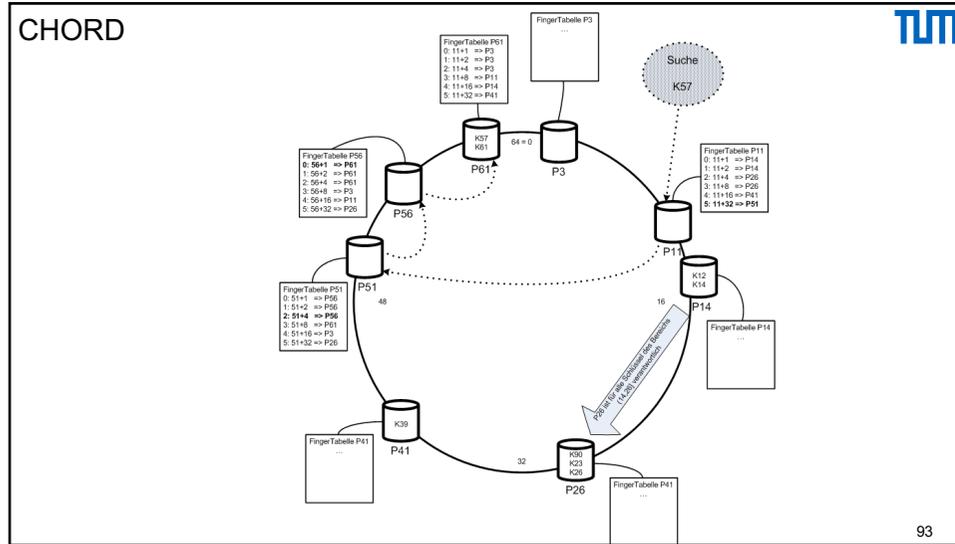
91

DHT: Distributed Hash Table



- Basieren auf „consistent hashing“
- Vollständige Dezentralisierung der Kontrolle
- Dennoch zielgerichtete Suche

92



No-SQL Datenbanken

Internet-scale Skalierbarkeit

CAP-Theorem: nur 2 von 3 Wünschen erfüllbar

- Konsistenz (Consistency)
- Zuverlässigkeit/Verfügbarkeit (Availability)
- Partitionierungs-Toleranz

No-SQL Datenbanksysteme verteilen die Last innerhalb eines Clusters/Netzwerks

- Dabei kommen oft DHT-Techniken zum Einsatz

95

Schnittstelle der No-SQL Datenbanken

Insert(k,v)

Lookup(k)

Delete(k)

Extrem einfach → effizient

Aber: wer macht denn die Joins/Selektionen/...

- → das Anwendungsprogramm

96

Konsistenzmodell: CAP



Relaxiertes Konsistenzmodell

- Replizierte Daten haben nicht alle den neuesten Zustand
 - Vermeidung des (teuren) Zwei-Phasen-Commit-Protokolls
- Transaktionen könnten veraltete Daten zu lesen bekommen
- Eventual Consistency
 - Würde man das System anhalten, würden alle Kopien irgendwann (also eventually) in denselben Zustand übergehen
- Read your Writes-Garantie
 - Tx leist auf jeden Fall ihre eigenen Änderungen
- Monotonic Read-Garantie
 - Tx würde beim wiederholten Lesen keinen älteren Zustand als den vorher mal sichtbaren lesen

97

Systeme



MongoDB
Cassandra
Dynamo
BigTable
Hstore
SimpleDB
S3
CockroachDB

98