

# ***Epidemische Algorithmen***

## ***Seminar Informationsmanagement im Web***

Vortrag von Konrad Rieck

Sommer 2003

Institut für Informatik – Fachbereich Mathematik und Informatik

Freie Universität Berlin

# Überblick

- ⑥ Motivation und Hintergrund
  - △ Replizierte Datenbanken
  - △ Synchronisationsverfahren
  
- ⑥ Epidemische Algorithmen
  - △ Epidemische Ausbreitung
  - △ Ereignisorientierte Algorithmen
  - △ Transaktionsorientierte Algorithmen
  
- ⑥ Schluss und Ausblick

# ***Motivation und Hintergrund***

# Motivation

- ⑥ Probleme von großen Diensten
  - △ Enorme Last für ein Rechnersystem
  - △ Ausfall von Rechnersystemen
  - △ Unzuverlässige Netzanbindungen
  
- ⑥ Lösungsmöglichkeit
  - △ Lastverteilung auf autonome Systeme
  - △ Datenverteilung mit Redundanz
  - △ Fehlertolerante Algorithmen

# Verteilte Datenbank

- ⑥ Autonome Rechnersysteme
  - △ Unabhängige Ressourcen
    - .. Netzanbindung
    - .. Stromversorgung
    - .. Hardware
  
- ⑥ Verbindung über Kommunikationsnetz
  
- ⑥ Verteilung des Datenbestandes
  - △ Fragmentierungsschema
  - △ Verteilungsschema

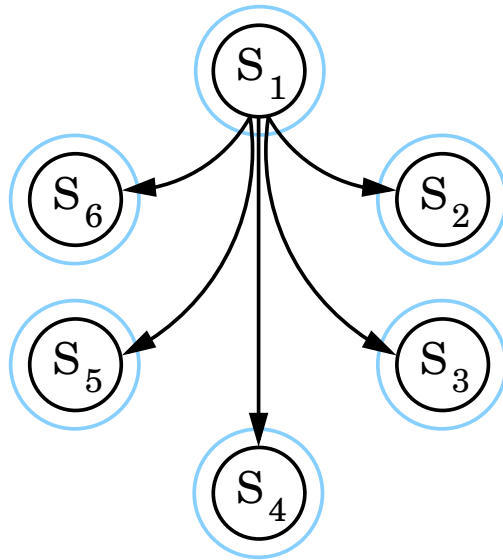
# Replizierte Datenbank

- ⑥ Spezialfall einer verteilten Datenbank
  - △ Keine Fragmentierung der Relationen
  - △ Verteilung mit maximaler Redundanz
- ⑥ Aufbau und Funktionsweise
  - △  $n$  Stationen  $S_1, S_2, \dots, S_n$
  - △ Jede Station verwaltet *Replikate*
  - △ Updates erfolgen lokal und dezentral
  - △ Synchronisation der Replikate

# Synchronisationsverfahren

- ⑥ Entscheidend für die Qualität der Datenbank
- ⑥ Kompromiss zwischen ...
  - △ ... Konsistenz der Replikate
  - △ ... Leistungsfähigkeit der replizierten Datenbank
- ⑥ Beispiele
  - △ Unicast oder Write-All Locking
  - △ Anti-Entropie Verfahren
  - △ Epidemische Algorithmen

# Unicast oder Write-All Locking

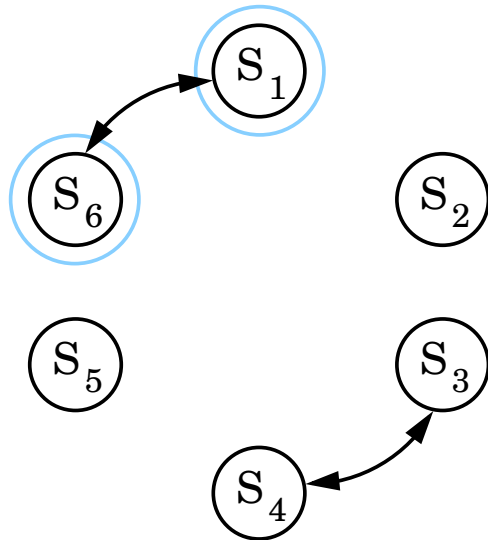


1. Lokales Update an  $S_1$
2. Kontakt und Sperrung
3. Transfer des Updates

- + Replikate immer konsistent
- Alle Stationen global bekannt und verfügbar
- Sperrung beeinträchtigt Leistungsfähigkeit



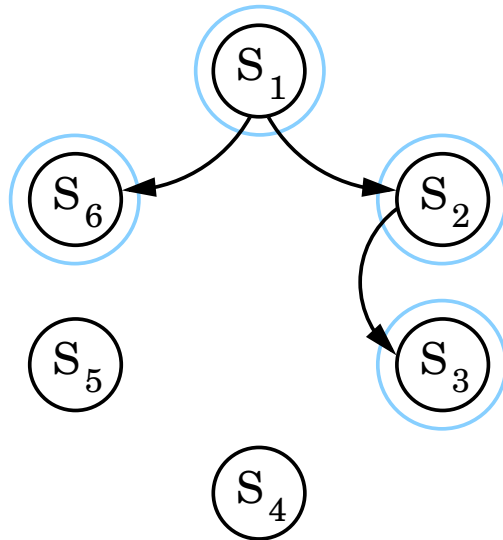
# Anti-Entropie Verfahren



1. Lokales Update an  $S_1$
2. Datenbankabgleich zufälliger Stationen
3. Update wird indirekt propagiert

- + Sperrungen nur an einzelnen Stationen
- + Nicht alle Stationen global bekannt oder verfügbar
- Replikate zeitweise inkonsistent
- Datenbankabgleich beeinträchtigt Leistungsfähigkeit

# Epidemische Algorithmen



1. Lokales Update an  $S_1$
2. Infektion einer Station
3. Epidemische Ausbreitung

- + Sperrungen nur an einzelnen Stationen
- + Nicht alle Stationen global bekannt oder verfügbar
- + Kein Datenbankabgleich, nur Transfer von Updates
- Replikate zeitweise inkonsistent

# ***Epidemische Algorithmen***

# Epidemische Ausbreitung (1)

- ⑥ Informationsausbreitung ähnlich ...
  - △ ... einer Epidemie
    - .. Virus springt von Wirt zu Wirt
    - .. Virus  $\hat{=}$  Nachricht über Update
    - .. Wirt  $\hat{=}$  Station
- ⑥ Lokales Update oder Nachricht über Update an  $S_i$
- ⑥ Nachricht über Update wird
  - △ ... gespeichert im lokalen Log  $L_i$
  - △ ... versendet an zufällige Station  $S_j$

## ***Epidemische Ausbreitung (2)***

- ⑥ Stationen in einem von drei Zuständen
  1. Empfänglicher Zustand
    - △ Zustand *vor* einer Epidemie.
  2. Infektiöser Zustand
    - △ Zustand nach lokalem Update oder nach Kontakt mit anderer infektiöser Station.
  3. Deaktivierter Zustand
    - △ Zustand ausgelöst mit Wahrscheinlichkeit  $\frac{1}{k}$  durch Kontakt mit bereits infizierter Station

# Ausbreitungsstärke

- ⑥ Ausbreitung hängt von Deaktivierungswahrscheinlichkeit  $\frac{1}{k}$  ab.

- ⑥ Anzahl  $s$  der nicht-infizierten Stationen

$$s = e^{-(k+1)(1-s)}$$

- ⑥  $s$  exponentiell abhängig von  $k$

- ⑥ Beispiel:  $k = 1 \Rightarrow s = 0.20$  und  $k = 2 \Rightarrow s = 0.06$

- △ Zu kleines  $k$ : Stationen werden nicht erreicht
- △ Zu großes  $k$ : Unnötige Kommunikationslast

- ⑥ Keine Garantie für vollständige Ausbreitung

# Ereignisorientierte Algorithmen

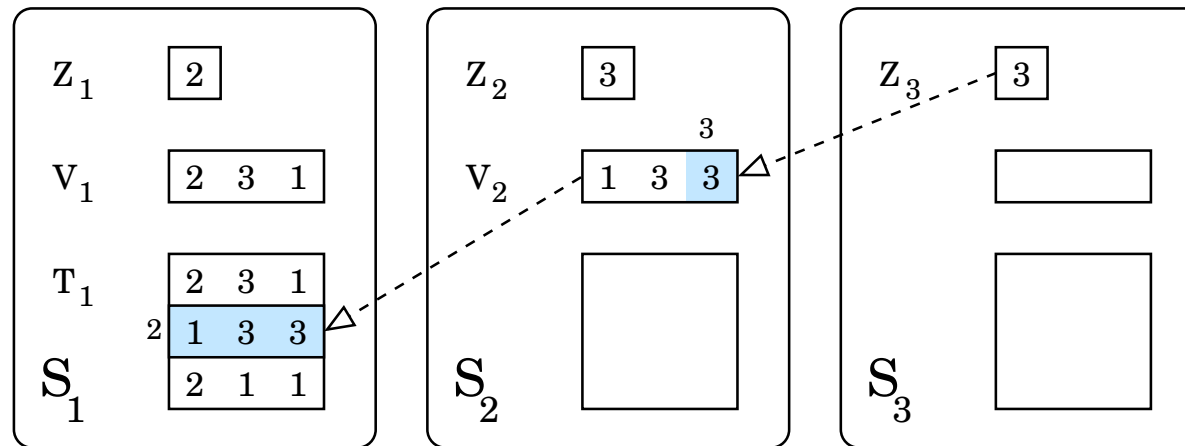
- ⑥ Problem der epidemischen Ausbreitung
  - △ Verlust des zeitlichen Zusammenhangs
    - ..  $S_i$  sendet Updates  $U_1, U_2, U_3$
    - ..  $S_j$  empfängt Updates  $U_3, U_2, U_1$
- ⑥ Lösungsmöglichkeit
  - △ Einführung eines Ereignismodells  $\langle E, \rightarrow \rangle$
  - △ Ereignismenge  $E$ , z.B.  $E = \{\text{READ}, \text{WRITE}\}$
  - △ „ $\rightarrow$ “ Relation, z.B.  $e \rightarrow f \equiv e$  vor  $f$  eingetreten.

# Realisierung der „ $\rightarrow$ “ Relation

- ⑥ „ $\rightarrow$ “ Relation verlangt Zeitmessung
- ⑥ Jede Station verwaltet eigene lokale Zeit
  - △ Lokale Zeit unabhängig von globaler Zeit
  - △ Änderung der lokalen Zeit ereignisorientiert
- ⑥ Lokale Zeit ermöglicht ...
  - △ ... Zeitstempel von Ereignissen  $TS(e)$
  - △ ...  $\forall e, f \in E : e \rightarrow f \Leftrightarrow TS(e) < TS(f)$
  - △ ... Speicherung von Zeiten (Vektoruhr  $V_i$ )
  - △ ... Wissen über Speicherung (Matrixuhr  $T_i$ )

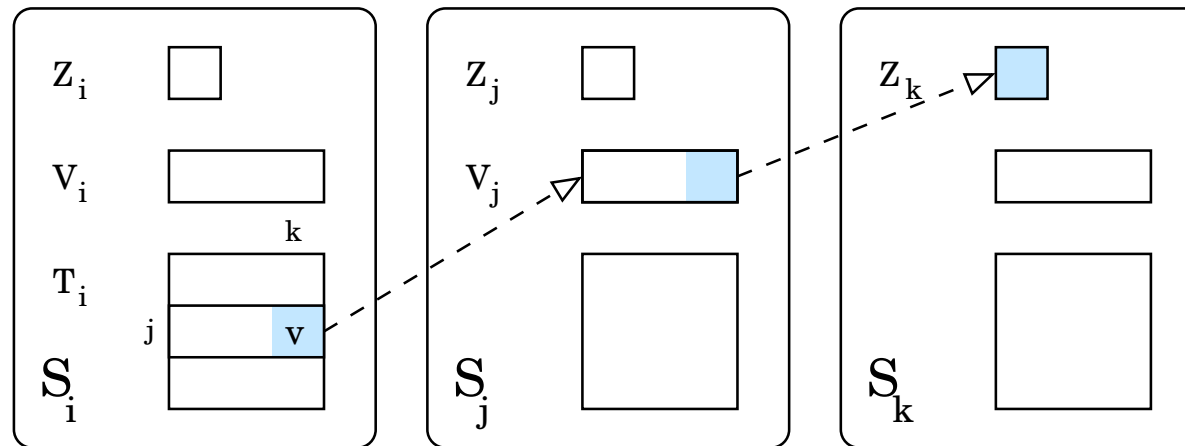


# Vektor- und Matrixuhren



- ⑥ Station  $S_i$  verwaltet unabhängige lokale Zeit  $Z_i$
- ⑥ Station  $S_i$  verwaltet Vektoruhr  $V_i = (Z_1, Z_2, Z_3)$
- ⑥ Station  $S_i$  verwaltet Matrixuhr  $T_i = (V_1, V_2, V_3)^T$
- ⑥ Es gilt  $\forall i, j, k : Z_k \geq V_j[k] \geq T_i[j, k]$

# Interpretation der Matrixuhr



- Station  $S_i$  kennt  $S_j$ 's Vektoruhr
- Station  $S_j$  verwaltet letzte Zeit von  $S_k$

$T_i[j, k] = v \Rightarrow S_i$  weiß, dass  $S_j$  alle Ereignisse von  $S_k$  bis zur Zeit  $v$  erhalten hat.

## Das HasRecvd Prädikat

- ⑥ Sei  $t$  ein Ereignis,  
 $S_k = \text{SITE}(t)$  die Station an der  $t$  aufgetreten ist,  
 $\text{TS}(t)$  die Vektorzeit als  $t$  aufgetreten ist und  
 $\text{TS}_k(t)$  die letzte bekannte Zeit von Station  $S_k$

$$\text{HASRECVD}(T_i, t, S_j) \equiv T_i[j, k] \geq \text{TS}_k(t)$$

- ⑥  $\text{HASRECVD}(T_i, t, S_j)$  ist wahr, wenn  $S_j$  das Ereignis  $t$  erhalten hat.
- ⑥  $\text{HASRECVD}(T_i, t, S_j)$  wird als *Timetable Eigenschaft* bezeichnet.

# Ein ereignisorientierter Algorithmus

- ⑥ Station  $S_i$  besitzt Log  $L_i$
- ⑥ Log  $L_i$  enthält Nachrichten über Updates
- ⑥ Station  $S_i$  wählt zufällig Station  $S_j$
- ⑥ Station  $S_i$  sendet alle Nachrichten  $t$  über Updates an  $S_j$  für die gilt  $\text{HASRECV}(T_i, t, S_j)$  ist falsch.
- ⑥ Station  $S_j$  erneuert Replikat
- ⑥ Station  $S_j$  trägt erhaltene Nachrichten in Log  $L_j$

# Transaktionsorientierte Algorithmen

## ⑥ Problem der ereignisorientierten Algorithmen

### △ Konflikte werden nicht berücksichtigt

- $S_i$  erhält Transaktion  $t_1 \equiv \text{READ}(x) \rightarrow \text{WRITE}(y)$
- $S_j$  erhält Transaktion  $t_2 \equiv \text{READ}(y) \rightarrow \text{WRITE}(x)$

## ⑥ Lösungsmöglichkeiten

### △ Pessimistischer epidemischer Algorithmus

- Optimistischer epidemischer Algorithmus
- Epidemischer Quorum-Algorithmus

# ***Pessimistischer Algorithmus***

Pessimistische Annahme: Konflikte häufig

- ⑥ Erweiterung des ereignisorientierten Algorithmus
- ⑥ Updates enthalten Transaktionen, nicht Ereignisse
- ⑥ Zu jeder Transaktion  $t$  werden übertragen ...
  - △ ... Zeitstempel  $TS(t)$  (Vektoruhr von Station)
  - △ ... Menge der gelesenen Objekte  $RS(t)$
  - △ ... Menge der geschriebenen Objekte  $WS(t)$

# Identifikation von Konflikten

- ⑥ Transaktionen  $t$  und  $t'$  im Konflikt, wenn ...
  1. Transaktionen greifen auf gleiche Objekte zu
    - △  $RS(t) \cap WS(t') \neq \emptyset$  *oder*
    - △  $WS(t) \cap WS(t') \neq \emptyset$  *oder*
    - △  $WS(t) \cap RS(t') \neq \emptyset$
  2. Transaktionen wurden gleichzeitig ausgeführt
    - △  $S_i$  erhielt zuerst  $t$ :  $TS_i(t) > TS_i(t')$  *und*
    - △  $S_j$  erhielt zuerst  $t'$ :  $TS_j(t') > TS_j(t)$
    - △ Inkompatibilität:  $TS(t) <> TS(t')$

# Das Aborted Prädikat

## ⑥ ABORTED Prädikat

$$\text{ABORTED}(t, S_j) \equiv \exists t' \in L_j : \text{TS}(t) \langle \rangle \text{TS}(t') \wedge \left( \begin{array}{l} \text{RS}(t) \cap \text{WS}(t') \neq \emptyset \quad \vee \\ \text{WS}(t) \cap \text{WS}(t') \neq \emptyset \quad \vee \\ \text{WS}(t) \cap \text{RS}(t') \neq \emptyset \end{array} \right)$$

- ⑥ Station  $S_j$  bricht Transaktion  $t$  und  $t'$  ab und sendet Abort-Nachricht, wenn  $\text{ABORTED}(t, S_j)$  wahr ist.



# Das Commit Prädikat

- ⑥ Station  $S_i$  führt Commit für Transaktion aus, wenn ...
  1. Alle Stationen das Update erhalten haben
  2. Keine Station eine Abort-Nachricht gesendet hat

- ⑥ COMMIT Prädikat

$$\text{COMMIT}(t, S_i) \equiv \forall j : \text{HASRECV}(T_i, t, S_j) \wedge \neg \text{ABORTED}(t, S_i)$$

- ⑥ Station  $S_i$  führt Commit für Transaktion  $t$  aus, wenn  $\text{COMMIT}(t, S_i)$  wahr ist.

# Verarbeitung einer Transaktion

- ⑥ Station  $S_i$  erhält Update der Transaktion  $t$
- ⑥ Station  $S_i$  setzt Schreibsperre auf  $WS(t)$
- ⑥ Station  $S_i$  wartet
- ⑥ Entweder  $COMMIT(t, S_i)$  wird wahr ...
  - △ Transaktion  $t$  wird übernommen und in  $L_i$  vermerkt
- ⑥ ... oder  $ABORTED(t, S_j)$  wird wahr
  - △ Transaktionen  $t$  und  $t'$  werden abgebrochen
  - △ Abort-Nachricht wird versendet
- ⑥ Sperre auf  $WS(t)$  wird gelöst

# Optimistischer Algorithmus

- ⑥ Beobachtung:
  - △ Pessimistischer Algorithmus langsam auf Grund längerer Sperrzeit

Optimistische Annahme: Konflikte selten

- ⑥ Variante des pessimistischen Algorithmus
- ⑥ Zu jeder Transaktion  $t$  wird gespeichert ...
  - △ Menge aller Transaktionen, die von  $t$  lesen  $\text{READFROM}(t)$ .

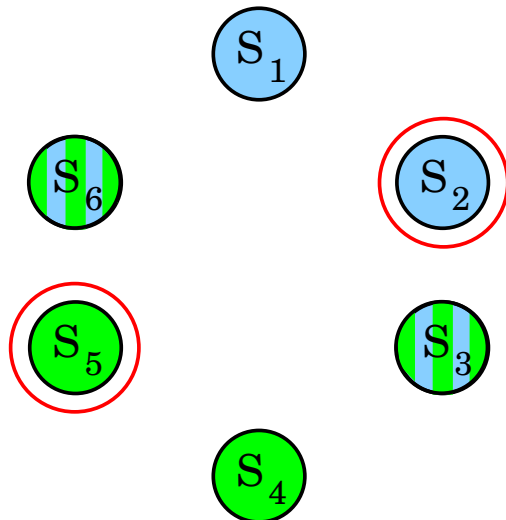
# Verarbeitung einer Transaktion

- ⑥ Station  $S_i$  erhält Update der Transaktion  $t$
- ⑥ Station  $S_i$  wartet
- ⑥ Entweder  $\text{COMMIT}(t, S_i)$  wird wahr ...
  - △ Transaktion  $t$  wird übernommen und in  $L_i$  vermerkt
- ⑥ ... oder  $\text{ABORTED}(t, S_j)$  wird wahr
  - △ Transaktionen  $t$  und  $t'$  werden abgebrochen
  - △ Transaktionen  $\text{READFROM}(t)$  und  $\text{READFROM}(t')$  werden abgebrochen. (Abbruch-Kaskade)
  - △ Abort-Nachricht wird versendet

# ***Epidemischer Quorum-Algorithmus***

- ⑥ Beobachtung:
  - △ Pessimistischer Algorithmus bricht bei einem Konflikt mindestens zwei Transaktionen ab.
- ⑥ Variante des pessimistischen Algorithmus
- ⑥ Verbesserung:
  - △ Votum welche Transaktionen abgebrochen bzw. übernommen werden soll.
  - △ Verwendung von Quoren und Antiquoren der Stationen

# Quoren und Antiquorum



1. Mehrheitsquorum 1
2. Mehrheitsquorum 2
3. Antiquorum

- ⑥ *Quoren* sind Teilmengen, die so gewählt werden, dass jede Schnittmenge von zwei Quoren nicht leer ist.
- ⑥ Ein *Antiquorum* ist jede Teilmenge, die ein Element aus jedem Quorum enthält.

# Verarbeitung einer Transaktion

- ⑥ Station  $S_i$  erhält Update der Transaktion  $t$
- ⑥ Station  $S_i$  setzt Schreibsperre auf  $WS(t)$
- ⑥ Stationen votieren bei Konflikt mit *Ja/Nein*
- 1. Station  $S_i$  erhält *Ja*-Quorum, so wird Commit ausgeführt
- 2. Station  $S_i$  erhält *Nein*-Antiquorum, so wird  $t$  abgebrochen und Abort-Nachricht versendet
- 3. Station  $S_i$  erhält Abort-Nachricht, so wird  $t$  abgebrochen
- ⑥ Sperre auf  $WS(t)$  wird gelöst

# ***Schluss und Ausblick***



# Schluss

- ⑥ Epidemische Algorithmen geeignet zur leistungsfähigen Synchronisation von replizierten Datenbanken
- ⑥ Epidemische Algorithmen können in unzuverlässigen Netzen eingesetzt werden
- ⑥ Synchronisation von Ereignissen und sogar Transaktionen möglich. Berücksichtigung von Konflikten
- ⑥ Verschiedene Algorithmen können entsprechend der Umgebung ausgewählt werden

# Vergleich der Algorithmen

- ⑥ Optimistischer Algorithmus (1) effizient in Umgebung mit wenigen Konflikten
- ⑥ Quorum-Algorithmus (2) effizienter als pessimistischer, da weniger Transaktionsabbrüche
- ⑥ Quorum-Algorithmus effizient in Umgebung mit häufigen Konflikten
- ⑥ Beide Algorithmen (1 & 2) effizienter als Unicast-Verfahren
- ⑥ Beide Algorithmen (1 & 2) geben keine Garantie auf vollständige Ausbreitung

# Ausblick

- ⑥ Epidemische Algorithmen könnten auch replizierte Fragmente synchronisieren
  - △ Mehr Ausfallsicherheit für verteilte Systeme mit fragmentierten Daten, z.B. P2P
- ⑥ Kombination aus Anti-Entropie und epidemischen Algorithmen garantiert vollständige Ausbreitung
  - △ Anti-Entropie Verfahren synchronisiert in großen Intervallen das System

***Fragen?***