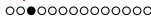


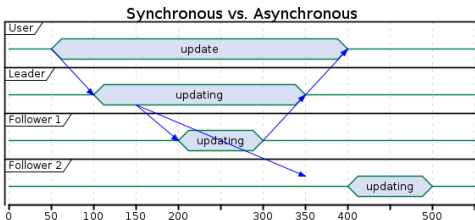
Wiederholung

- Probleme mit Uhren
- Synchronisation von Uhren
- Logische Zeit
 - Lamport
 - Vektor
- Mutex
- Wahlalgorithmen



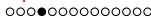
Replikation

Synchronous vs. Asynchronous



- Replikation zu Follower 1 ist synchron.
- Replikation zu Follower 2 is asynchron.

Was sind die Vor- und Nachteile von (a)synchroner Replikation?



Synchronous vs. Asynchronous 2

Vorteil synchroner Replikation

- Bestätigte writes wurden bereits repliziert.
- Absturz des Leader ohne Datenverlust

Nachteil synchroner Replikation

- Nicht verfügbar bei Absturz von Followern

Kombination synchroner und asynchron üblich

Beispiel:

- 1 synchroner Follower
- Absturz: Neuen synchronen Follower wählen

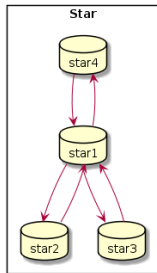
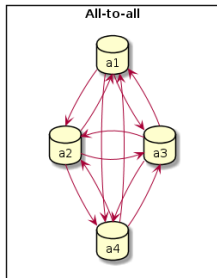
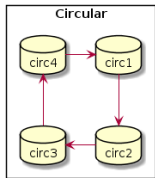
Replication Logs: Write Ahead Log (WAL)

- SQL Datenbanken verwalten WAL.
- An WAL wird nur angefügt (append-only)
- Einträge in WAL werden von den Followern angewandt

Problem: Struktur des WAL üblicherweise Implementierungsdetail

- Verschiedene Versionen einer DB gleichzeitig?
- Schließt Zero-Downtime-Updates aus
- Logische Logs oft von der Storage Engine entkoppelt

Topologien



Kennzahlen

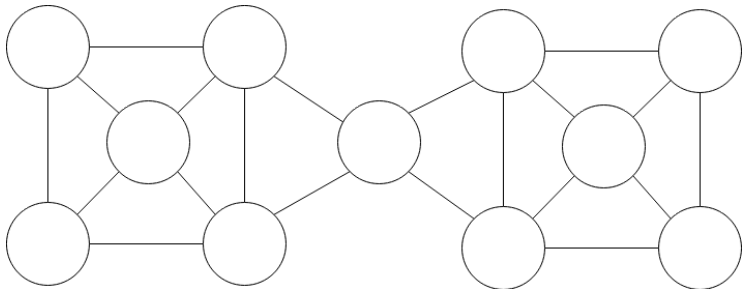
Durchmesser Der längste kürzeste Pfad.

Bisektionsbandbreite Kanten zu löschen zum Zerlegen in zwei ähnlich große Teile



Übung Toplogien

Bestimme Durchmesser, Bisektionsbandbreite und Knoten/Kantenkonnectivität für folgende Toplogie:





Leaderless Replication

- Verbreitet durch Amazons Dynamo DB
- Auch Riak, Cassandra, Voldemort
- Writes auf jedem Knoten
- Meist „Quorum“ Reads und Writes.



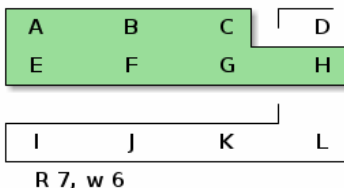
○○○○○○○○○○○○○○○●○○○○○ ○○○○



○○○○○○○○○○○○○ ○○○○○○○○○○○○○○ ○

Replikation

Quorum: Write-Write-Konflikte vermeiden



- Wenn $w \leq \frac{n}{2}$ können 2 Nutzer widersprüchliche Daten schreiben.
- Beim Lesen erkennbar, da $r > n - w$
- write-write Konflikt oder stale data

CAP (Consistency, Availability, Partition Tolerance)

CAPs Consistency (Konsistenz)

CAP meint mit Consistency Linearizability.

- Jede Operation atomar
- Jeder Prozess sieht die gleiche Ordnung von Operationen.
- Diese Ordnung entspricht der realen Ordnung von Operationen.

CAP (Consistency, Availability, Partition Tolerance)

CAPs Availability (Verfügbarkeit)

- Jede Anfrage an einen verfügbaren Knoten liefert eine Antwort
- Auch bei Netzwerkpartitionen.

CAP (Consistency, Availability, Partition Tolerance)

CAPs Partition (Teilung)

Eine Netzwerkpartition teilt das System in 2 oder mehr Seiten. Die einzelnen Seiten können nicht miteinander kommunizieren.

Kann man entscheiden, dass keine Netzwerkpartitionen auftreten?

CAP (Consistency, Availability, Partition Tolerance)

CAP - Bewertung

- Linearizability **OR** Total Availability.
- CAP Theorem wird oft im Marketingsprech verwendet, um Datenbanken zu charakterisieren.
- Konsistenzmodelle neben Linearizability, verschiedene Grade von Verfügbarkeit

Welche Kombinationen sind möglich?

Availability

- Total Available / High Available
- Sticky Available
- Unavailable

Literatur: Highly Available Transactions: Virtues and Limitations
Bailis et al. (2013).

Total Available / High Available

- Antwort erhält, wer **einen** korrekten (nicht versagenden) Server kontaktieren kann
- Auch bei Netzwerkpartitionen zwischen Servern

Sticky Available

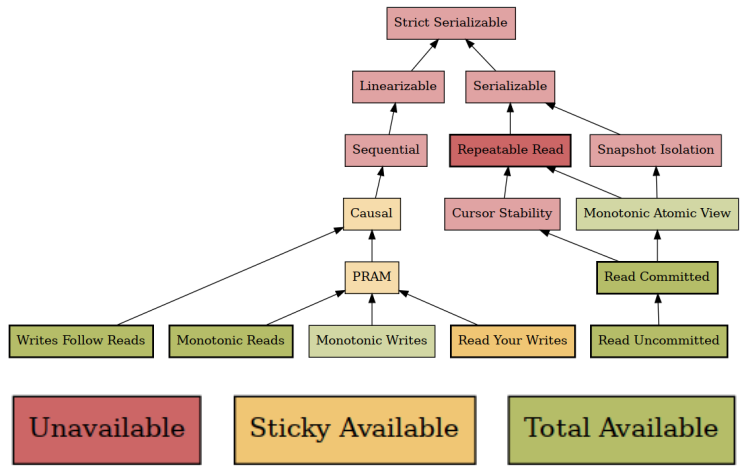
- Antwort erhält, wer einen Server kontaktieren kann, der den gesamten, **dem Nutzer bekannten Zustand** beinhaltet
- Auch bei Netzwerkpartitionen zwischen Servern

Sticky Available - Beispiel

- Daten auf mehrere Server repliziert
- Jede Replika enthält alle Daten
- Nutzer kontaktiert immer denselben Server
- \Rightarrow Sticky Available

Consistency

Consistency

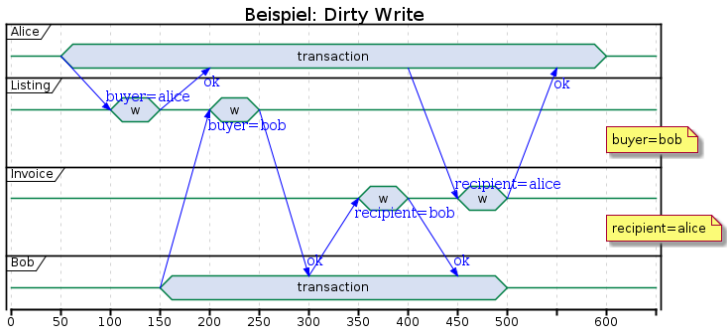


Isolation

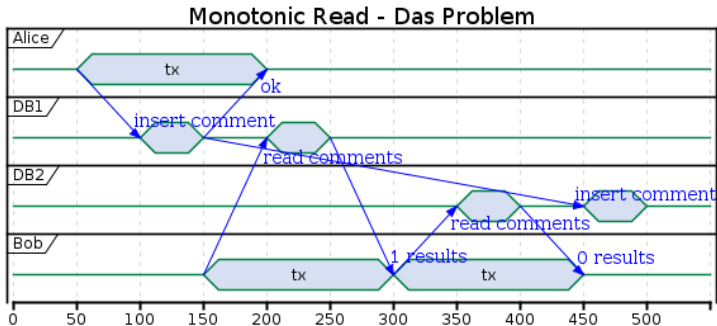
Read Uncommitted

Verhindert dirty writes:

Beim Schreiben in eine DB werden nur Daten überschrieben, die bereits committed wurden.

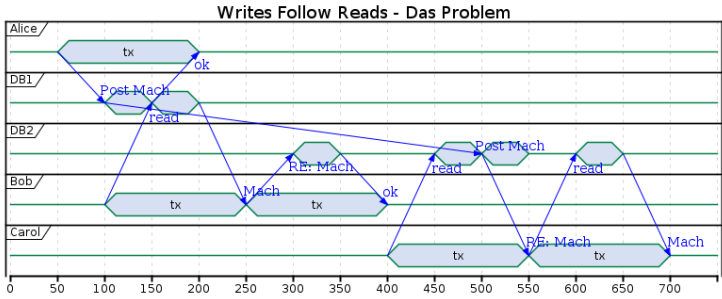


Monotonic Reads/Writes



Monotonic Reads Sobald ein Wert von einem Client gelesen wurde, wird dieser oder ein späterer Wert bei folgenden reads gelesen.

Writes Follow Reads



CALM: Monotonicity

Formale Definition:

Ein Programm P ist monoton, wenn für alle input sets S, T gilt:
wenn $S \subseteq T$ dann $P(S) \subseteq P(T)$

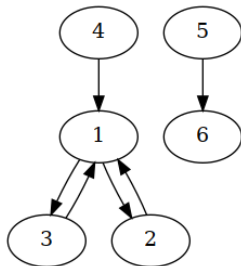
Beispiel: Set addition

- $S = \text{add}(a), \text{add}(b)$
- $T = \text{add}(a), \text{add}(b), \text{add}(c)$
- $S \subseteq T = \text{true}$
- $P(S) = \{a, b\}$
- $P(T) = \{a, b, c\}$
- $P(S) \subseteq P(T) = \text{true}$

CALM Theorem

Distributed Deadlock Detection - Berechnung

Jede Maschine übermittelt ihre Kanten.



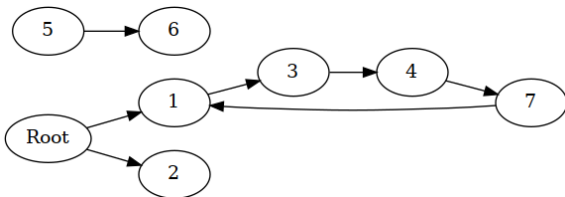
Können zusätzliche Kanten dazu führen Deadlocks auflösen?
 → Nein. Zusätzliche Kanten führen nur zu evtl. zusätzlichen Deadlocks. → monoton!



CALM Theorem

Distributed Garbage Collection - Berechnung

Jede Maschine übermittelt ihre Kanten.



Können weitere Kanten dazu führen unser Ergebnis ändern?

→ Ja. Eine Kante von 1 zu 5 würde dazu führen, dass 5 und 6 nicht collected werden können → nicht monoton! → Wir benötigen Koordination

Composability

Wenn die Funktionen f und g monoton sind, ist auch $f(g(x))$ monoton.

→ Wir können monotone Programme aus monotonen Operationen konstruieren.

Beispiel: Shopping Cart

Können wir den Shopping Cart einer Webanwendung monoton gestalten?

Idee: Wir modellieren den Inhalt des Shopping Cart als Set.

Das Hinzufügen eines Artikels ist damit monoton und wir benötigen keine Koordination.

Aber wie können wir einen Artikel aus dem Shopping Cart entfernen?

Idee: Wir verwalten ein Add-Set und ein Remove-Set.

Shopping Cart: Bewertung

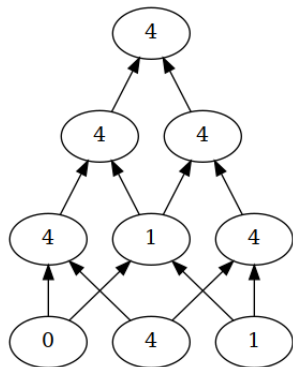
Die Verwaltung des Shopping Carts kann ohne Koordination erfolgen.

Allerdings benötigen wir Koordination, sobald wir den Einkauf tätigen.

- > Wir müssen sicherstellen, dass alle Änderungen des Shopping Carts gesehen wurden.
- > Aber wir konnten die nötige Koordination im System reduzieren.

CRDTs

Beispiel: Integer, max



- assoziativ: $\max(\max(1,2),3) = \max(1,\max(2,3))$
- kommutativ: $\max(1,2) = \max(2,1)$
- idempotent: $\max(1,2) = \max(\max(1,2),2)$

Operation Based Replication

- Das System repliziert Operationen anstatt State.
- Benötigt einen zuverlässigen broadcast channel, der die vom CRDT bestimmte Ordnung berücksichtigt (Beispiel: kausal).
- Operationen die nach der Ordnung gleichzeitig (concurrent) stattfinden, müssen kommutativ sein.

Beispiel:

- $+7, -5 = -5, +7$

Operation Based vs. State Based

- State Based:
 - Simpler, da keine Ordnung der Nachrichten benötigt wird und jede Änderung lokal betrachtet werden kann.
 - Nachrichten müssen irgendwann ankommen, aber Reihenfolge ist egal.
 - Überträgt immer gesamten State.
 - Keine Gruppenzugehörigkeit nötig.
- Operation Based:
 - Komplexer, sämtliche Nachrichten betrachtet.
 - Überträgt nur die Operationen. („diff“)
 - Benötigt Gruppenzugehörigkeit.

G-Set (Grow only)

```
state = Set()
def add(element):
    state.add(element)

def merge(other_state):
    state.union(other_state)

def contains(element):
    return element in state
```


Einschränkung 2P-Set?

- einmal hinzugefügtes Element kann nie wieder hinzugefügt werden.

Weitere CRDTs

- PN-Set: Counter für jedes Element, Wert des Counter entscheidet über Set-Zugehörigkeit.
- 2P2P Graph: Je ein 2P Set für Knoten und Kanten.
- Verschiedene Implementierungen von Registern.
- Datentypen für kollaborative Textbearbeitung.

→ [Conflict-free_replicated_data_type#Known_CRDTs](#)

Zusammenfassung

- CRDTs können genutzt werden, um Koordination zu vermeiden oder zumindest einzuschränken.
- CRDTs benötigen eine Form der Garbage Collection, um performant zu bleiben.
 - Garbage Collection benötigt wiederum Koordination.
- CRDTs werden in Verteilten Systemen eingesetzt: Riak, Redis, Dynamo

Ich wünsche Ihnen unkoordinierten Erfolg!



Verweise I

- Adya, A. and Liskov, B. (1999). Weak consistency: A generalized theory and optimistic implementations for distributed transactions.
- Bailis, P., Davidson, A., Fekete, A., Ghodsi, A., Hellerstein, J. M., and Stoica, I. (2013). Highly available transactions: Virtues and limitations. *Proc. VLDB Endow.*, 7(3):181–192.
- Shapiro, M., Preguiça, N., Baquero, C., and Zawirski, M. (2011). A comprehensive study of convergent and commutative replicated data types.

Bilder: