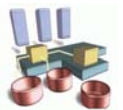
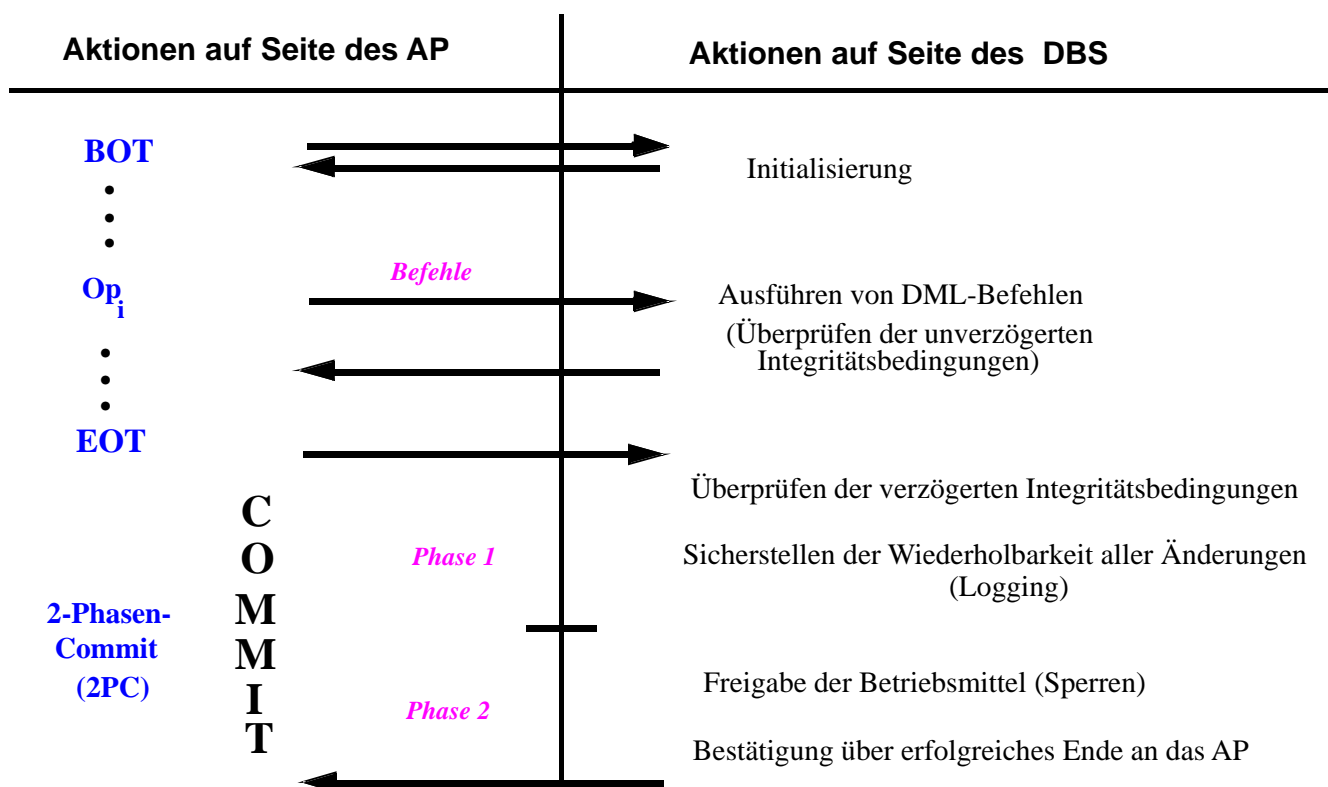


# 6. Verteilte Transaktionsverwaltung

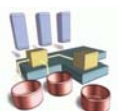
- Einführung
- Synchronisation
  - Verfahrensüberblick, Sperrverfahren
  - Deadlock-Behandlung
    - Timeout
    - Deadlock-Vermeidung: Wait/Die-, WoundWait-Verfahren
    - globale Deadlock-Erkennung
- Commit-Protokolle: Anforderungen
- Basis-Protokoll: 2-Phasen-Commit
  - Ablauf
  - Fehlerbehandlung
  - Lineares 2PC
  - Hierarchisches 2PC, XA-Protokoll
- 1-Phasen-Commit
- 3-Phasen-Commit



## Transaktionsverarbeitung in zentralisierten DBS

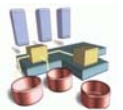


Weiterarbeit im Anwendungsprogramm(AP)



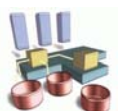
# Transaktionsverwaltung in VDBS

- ACID-Eigenschaften von Transaktionen auch im verteilten Fall zu garantieren: Atomarität, Konsistenz, Isolation, Dauerhaftigkeit
- Logging und Recovery
  - wesentliche Neuerung: globales Commit-Protokoll
  - Robustheit gegenüber partiellen Fehlern, insbesondere Kommunikationsfehlern (Netzwerkpartitionierungen u. ä.)
- Integritätssicherung
  - Integritätsbedingungen auf fragmentierten Relationen
  - Überwachung zB im Rahmen eines erweiterten Commit-Protokolls



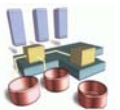
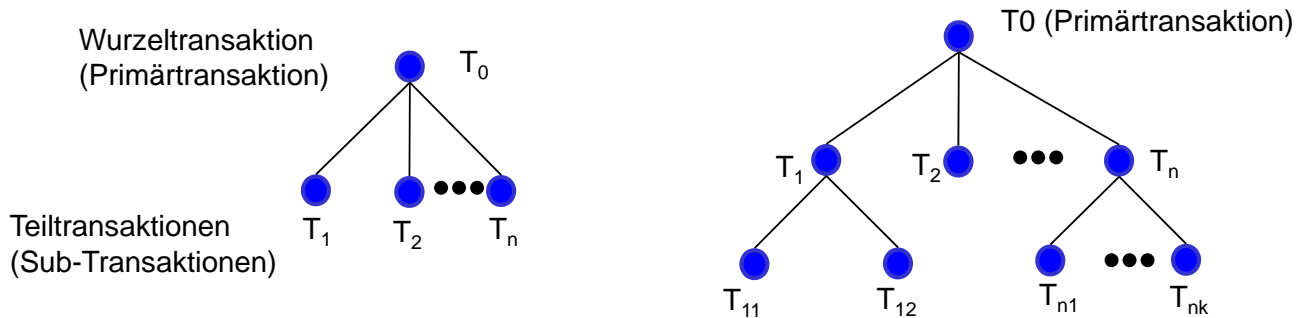
## Transaktionsverwaltung in VDBS (2)

- Synchronisation
  - Wahrung der globalen Serialisierbarkeit
  - rechnerübergreifende Abhängigkeiten (globale Deadlocks u. ä.)
- Replikation
  - Wahrung von Replikationstransparenz
  - Optimierung von Leistung und Verfügbarkeit



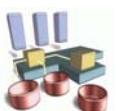
# Transaktionsstruktur

- Transaktionsaufbau: BOT,  $OP_1$ ,  $OP_2$ , ...,  $OP_n$ , COMMIT
- Transaktionsbaum: Kontrollstruktur in einer verteilten Umgebung
  - Transaktionsbaum repräsentiert Aufrufbeziehungen
  - 1-stufige oder geschachtelte Teiltransaktionen
  - isoliertes Rücksetzen von Sub-Transaktionen wird i.a. nicht unterstützt:  
Abbruch einer Teiltransaktion führt zum Abbruch der Gesamttransaktion



## Synchronisation

- Mehrbenutzerbetrieb führt ohne Synchronisation zu **Anomalien**:
  - Lost Updates
  - Dirty Reads
  - Non-repeatable Reads
  - Phantome
- Korrektheitskriterium der (**globalen**) **Serialisierbarkeit**:  
*gleichzeitige (und verteilte) Ausführung mehrerer Transaktionen ist äquivalent zu wenigstens einer seriellen Ausführung derselben Transaktionen*



# Synchronisation (2)

## ■ Sperrprotokolle garantieren Serialisierbarkeit

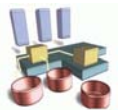
- vor jedem Objektzugriff muß Sperre mit ausreichendem Modus angefordert werden
- gesetzte Sperren anderer Transaktionen sind zu beachten
- am Transaktionsende (2. Commit-Phase) werden alle Sperren freigegeben

## ■ Einfachster Ansatz: RX - Sperrverfahren

		aktueller Sperrmodus		
		NL	R	X
angeforderter Sperrmodus	R	+	+	-
	X	+	-	-

NL: no lock,  
R: read lock;  
X: eXclusive lock

- + Sperre wird gewährt
- Sperrkonflikt



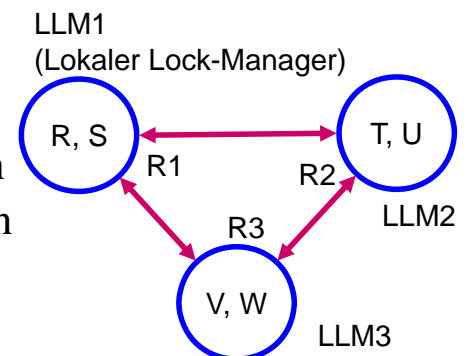
# Synchronisation in VDBS

## ■ Zentrale Sperrverfahren inakzeptabel

- Knotenautonomie, Kommunikationsaufwand

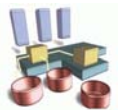
## ■ Verteilte Sperrverfahren

- jeder Rechner verwaltet Sperren für lokale Daten
- keine eigene Nachrichten für Sperranforderungen
- Sperrfreigabe innerhalb des Commit-Protokolls
- am weitesten verbreiteter Ansatz für VDBS



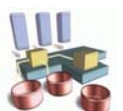
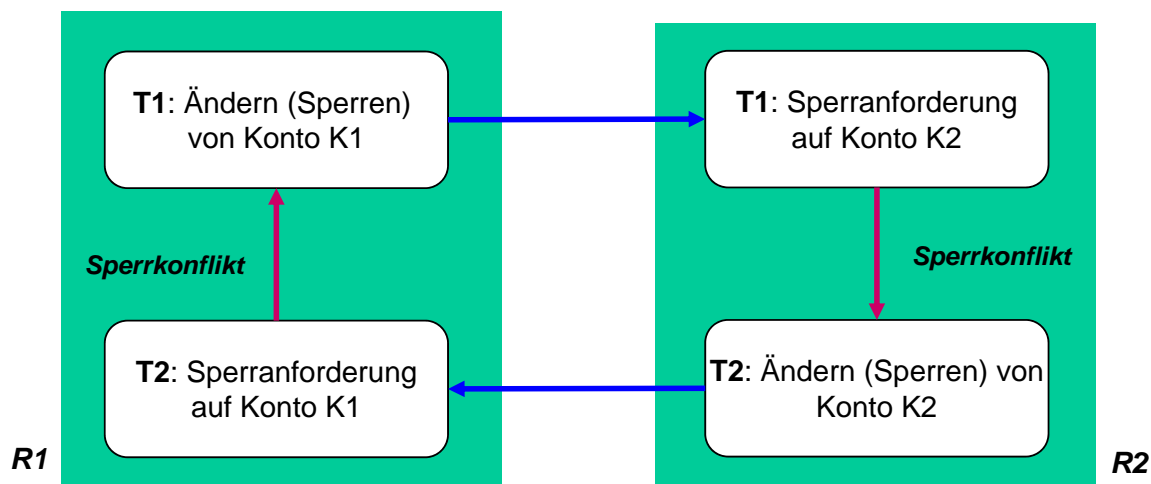
# Synchronisation in VDBS (2)

- Zeitmarkenverfahren (timestamp ordering)
  - Transaktionen erhalten eindeutige Zeitmarken bei BOT
  - alle Zugriffe müssen in Reihenfolge der BOT-Zeitmarke erfolgen
  - keine Deadlocks, jedoch viele Rücksetzungen
- Optimistische Synchronisation
  - keine Sperren, sondern Konflikttest am Transaktionsende (Validierung)
  - Problem zahlreicher Rücksetzungen und des „Verhungerns“ von Transaktionen



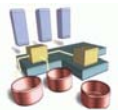
## Deadlock-Behandlung

- Sperrverfahren erfordern Deadlock-Behandlung
- Deadlock: zyklische Wartebeziehung zwischen zwei oder mehr Transaktionen
  - Deadlock-Behandlung in DBS erfordert Rücksetzung von Transaktionen
- Beispiel (Überweisungen zwischen Konten K1 und K2)



# Deadlock-Behandlung: Timeout-Verfahren

- Festsetzen einer maximalen Wartezeit auf eine Sperre (Timeout)
- Nach Überschreiten des Timeout wird wartende Transaktion zurückgesetzt
- Vorteile
  - jeder Deadlock wird irgendwann aufgelöst
  - geringer Implementierungsaufwand
  - keine zusätzliche Nachrichten für Deadlock-Behandlung
  - kann auch bei heterogenen DBS genutzt werden
- Nachteile
  - unnötige Rücksetzungen (Erreichen des Timeout auch ohne Deadlock)
  - Wahl des Timeouts (viele Rücksetzungen vs. später Auflösung von Deadlocks)
- Nutzung in den meisten VDBS bzw. verteilten Transaktionssystemen



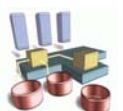
## Deadlock-Behandlung: Weitere Lösungsmöglichkeiten

### Deadlock-Vermeidung (Avoidance)

- Zuweisung einer eindeutigen *Transaktionszeitmarke* bei BOT
- Entscheidung bei einem Sperrkonflikt, ob Warten zugelassen wird oder eine Transaktionsrücksetzung erfolgt
- kein Wartezyklus (Deadlock) möglich, falls stets nur ältere Transaktionen auf jüngere warten dürfen (Bsp.: Wait/Die) bzw. nur jüngere auf ältere Transaktion warten (Bsp.: Wound/Wait-Verfahren)
- Behandlung globaler Deadlocks ohne Kommunikation!

### Deadlock-Erkennung (Detection)

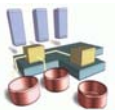
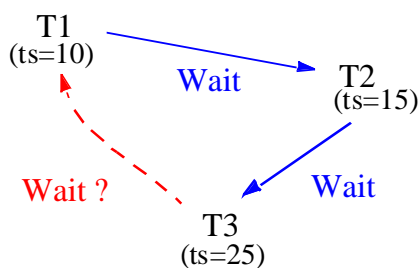
- Führen eines globalen Wartegraphen (-> Kommunikationsbedarf)
- Zyklensuche zur Erkennung von Deadlocks
- potentiell geringste Anzahl von Rücksetzungen



# Deadlock-Vermeidung: Wait/Die

- Zuweisung einer eindeutigen *Transaktionszeitmarke*  $ts(T)$  bei Beginn jeder Transaktion  $T$
- WAIT/DIE-Verfahren
  - anfordernde Transaktion  $T_i$  wird zurückgesetzt, falls sie in Sperrkonflikt mit älterer Transaktion gerät
  - $T_i$  wartet bei einem Sperrkonflikt, falls sie älter ist als Sperrbesitzer

```
 $T_i$  fordert Sperre, Konflikt mit  $T_j$ :  
if  $ts(T_i) < ts(T_j)$  { $T_i$  älter als  $T_j$ }  
then WAIT ( $T_i$ )  
else Rollback ( $T_i$ ) {'Die'}
```



## Wait/Die (2)

- Kein Zyklus möglich, da nur ältere auf jüngere Transaktionen warten
- Bewertung
  - Behandlung globaler Deadlocks ohne Kommunikation
  - unnötige Rücksetzungen (ohne Vorliegen eines Deadlocks) möglich

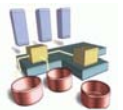


# Deadlock-Vermeidung: Wound/Wait

- preemptiver Ansatz: Sperrbesitzer wird zurückgesetzt, wenn er bei einem Sperrkonflikt jünger als anfordernde Transaktion ist
- kein Zyklus möglich, da nur jüngere auf ältere Transaktionen warten

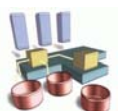
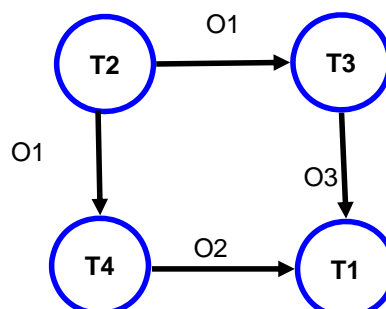
```
Ti fordert Sperre, Konflikt mit Tj:  
if ts (Ti) < ts (Tj) {Ti älter als Tj}  
then ABORT (Tj) {'Wound'}  
else Wait (Ti)
```

- Bewertung ähnlich zu Wait/Die
  - noch stärkere Bevorzugung älterer Transaktionen positiv zu bewerten



# Deadlock-Erkennung

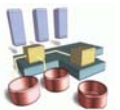
- Explizites Führen eines *Wartegraphen* (wait-for graph)
  - Knoten: laufende Transaktionen
  - gerichtete Kanten: Wartebeziehungen aufgrund von Sperrkonflikten
  - Zyklus kennzeichnet Deadlock
- Zyklensuche zur Erkennung von Verklemmungen
  - bei jedem Sperrkonflikt bzw.
  - verzögert (z. B. über Timeout gesteuert)
- Deadlock-Auflösung durch Zurücksetzen einer oder mehrerer Transaktionen, deren Wegfall Zyklus auflöst
  - z. B. Verursacher des Deadlocks
  - "billigste" Opfer





# Deadlock-Erkennung in Verteilten DBS

- Nachrichtenaustausch zur Erstellung des Wartegraphen
- Zentrale Deadlock-Erkennung
  - ausgezeichnete Knoten verwaltet globalen Wartegraphen
  - hoher Kommunikationsaufwand (Übertragung aller neuen / wegfallenden Wartebeziehungen)
  - Einschränkung der Knotenautonomie
- Verteilte Deadlock-Erkennung
  - verteilte Verwaltung eines globalen Wartegraphen
  - korrektes Verfahren schwierig zu realisieren:
    - Nachrichtenverzögerungen
    - Empfangs- ≠ Sendereihenfolge
    - Rechnerausfälle
  - oftmals
    - doppelte Erkennung von Deadlocks
    - "falsche" Deadlocks

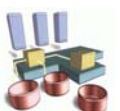


## Verteilte Deadlock-Erkennung in R\*

- "Deadlock Detector" pro Rechner, der periodisch lokalen Wartegraphen auf Zyklen durchsucht
  - spezieller Knoten "EXTERNAL" im Wartegraph zur Darstellung von Wartebeziehungen zu Sub-Transaktionen auf anderen Rechnern
  - Zyklus mit EXT-Knoten kennzeichnet potentiellen globalen Deadlock
- Weitergabe der Zyklusinformation an andere Rechner, um globalen Deadlock ggf. zu erkennen
- Ausgangsbeispiel

R1:

R2:



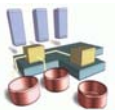
# Deadlock-Erkennung in R\* (2)

## ■ Kooperation mit anderen Rechnern:

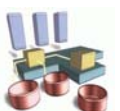
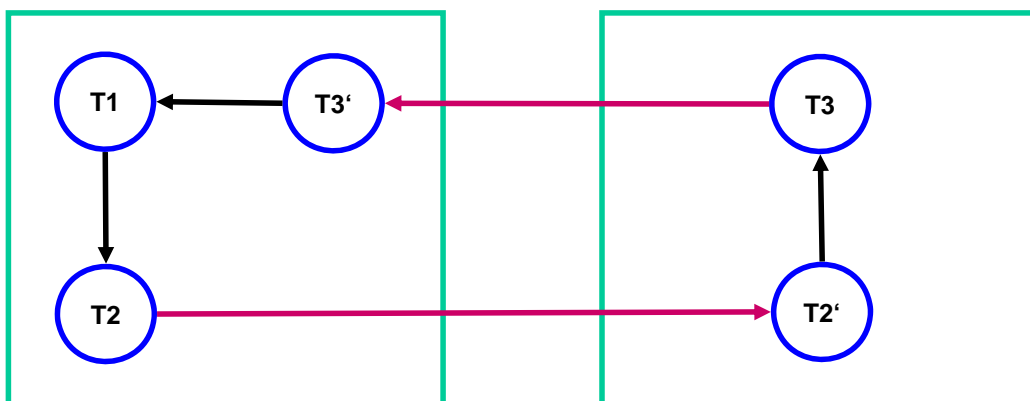
1. Empfange Deadlock-Information anderer Rechner
2. Erweitere damit lokalen Wartegraphen
3. löse lokale/vollständige Zyklen durch Bestimmung und Rücksetzung eines "Opfers" auf
4. für globale Zyklen sende lineare Darstellung:  
 $EXT \rightarrow T_1 \rightarrow T_2 \rightarrow \dots \rightarrow T_n \rightarrow EXT$   
an Rechner, auf den  $T_n$  wartet (falls  $T_1 > T_n$ )

## ■ Bewertung

- **max.  $N(N-1)/2$  Nachrichten** zur Erkennung eines globalen Deadlocks bei  $N$  beteiligten Rechnern
- Erkennung falscher Deadlocks möglich



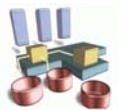
## Beispiel



# Commit-Protokolle

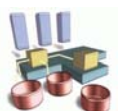
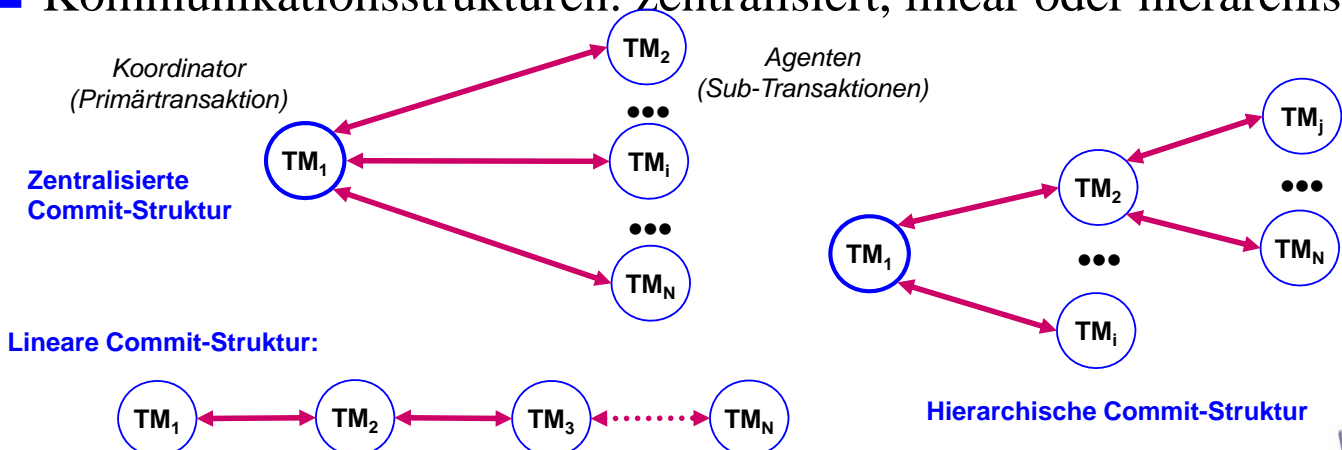
- Sicherstellen der **Atomizität** verteilter Transaktionen durch rechnerüber-greifendes Mehrphasen-Commit-Protokoll
- Anforderungen an geeignetes Commit-Protokoll:
  - Korrektheit
  - Geringer Aufwand (#Nachrichten, #Log-Writes)
  - Geringe Antwortzeitverlängerung
  - Robustheit gegenüber Rechnerausfällen und Kommunikationsfehlern
  - Knotenautonomie: jeder an einer verteilten Transaktionsausführung beteiligte Rechner soll möglichst lange das Recht des einseitigen Transaktionsabbruchs (*unilateral abort*) haben

"Nicht-Fehler-Fall" ist zu optimieren

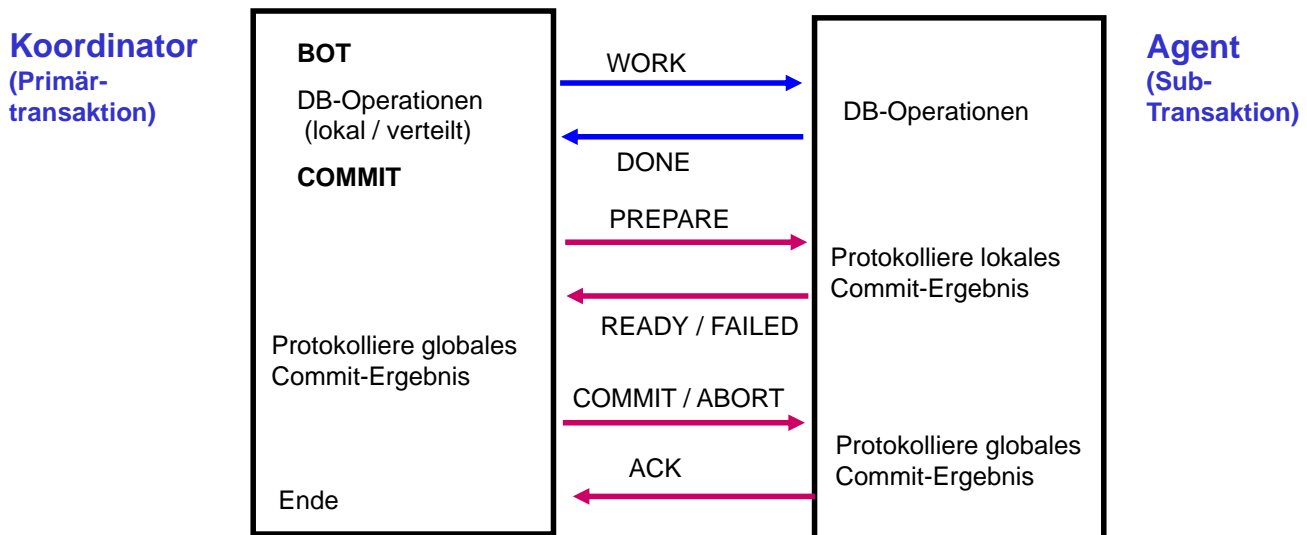


## Commit-Protokolle (2)

- Ausführung des Commit-Protokolls erfolgt durch *Transaktions-Manager (TM)* an jedem Knoten (1 Koordinator + N-1 Agenten)
- Wesentliche Alternativen
  - 2-Phasen-Commit (zentral, linear, hierarchisch)
  - 1-Phasen-Commit
  - 3-Phasen-Commit
- Kommunikationsstrukturen: zentralisiert, linear oder hierarchisch



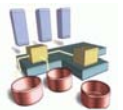
# Zentralisiertes 2-Phasen-Commit (N=2)



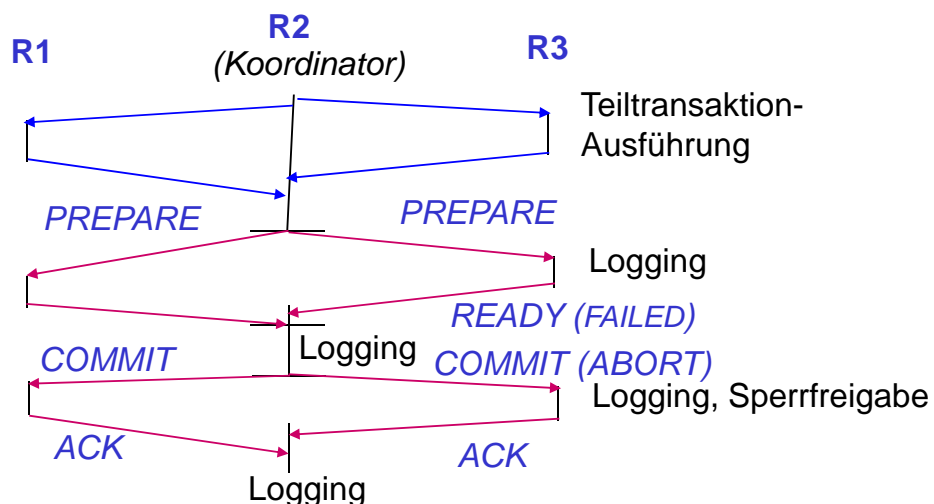
## ■ Aufwand

- Erfolgreicher Ausgang: 4 Nachrichten, 4 Log-Writes
- ABORT-Nachrichten gehen nur an Teiltransaktionen, die nicht mit FAILED gestimmt haben

## ■ Kernproblem Koordinatorausfall => ggf. lange Blockierung



# Zentralisiertes 2-Phasen-Commit (N=3)

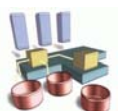


## ■ Basisverfahren:

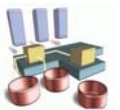
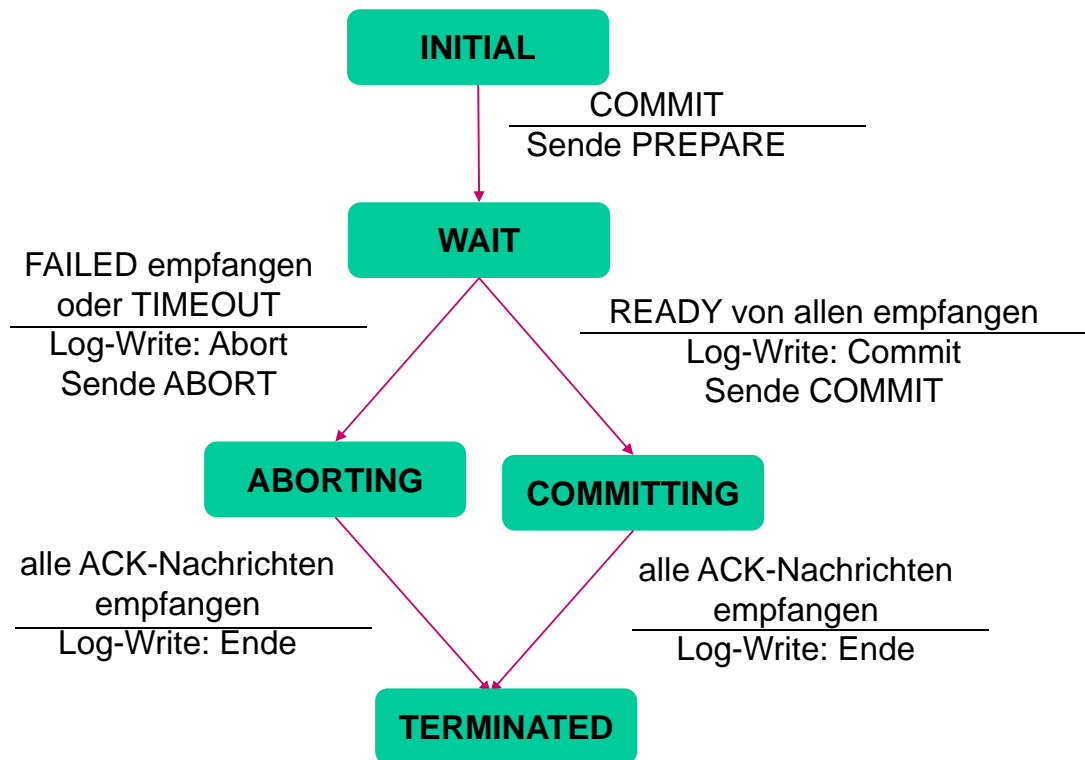
- 4 (N-1) Nachrichten (N = Anzahl der beteiligten Knoten)
- 2 N Log-Writes

## ■ Optimierung für lesende Sub-Transaktionen (Anzahl M)

- 4 (N-1) - 2M Nachrichten für M < N, 2 (N-1) für M=N
- 2 N - M Log-Writes



# Zustandsgraph Koordinator (2PC)



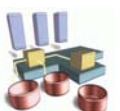
## 2PC: Fehlerbehandlung

### ■ Timeout-Bedingungen Koordinator:

- WAIT => setze Transaktion zurück; verschicke ABORT-Nachr.
- ABORTING, COMMITTING => vermerke Agenten, für die ACK noch aussteht

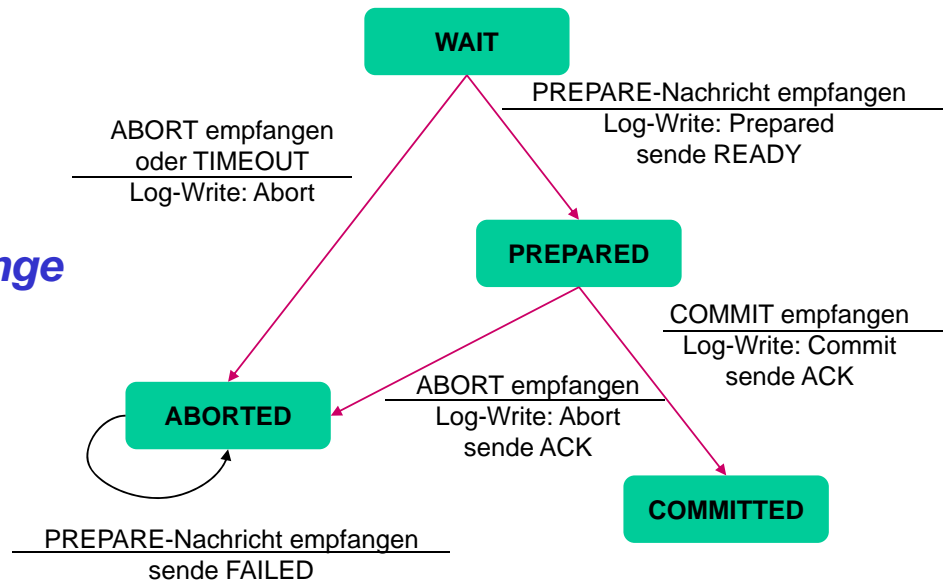
### ■ Ausfall des Koordinatorknotens:

- Log-Zustand TERMINATED:
  - UNDO bzw. REDO-Recovery, je nach Transaktionsausgang
  - keine "offene" Teiltransaktionen möglich
- Log-Zustand ABORTING:
  - UNDO-Recovery
  - ABORT-Nachricht an Rechner, von denen ACK noch aussteht
- Log-Zustand COMMITTING:
  - REDO-Recovery
  - COMMIT-Nachricht an Rechner, von denen ACK noch aussteht
- Sonst: UNDO-Recovery



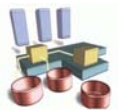
## 2PC: Fehlerbehandlung (2)

### Zustandsübergänge Agent



### ■ Timeout-Bedingungen für Agenten:

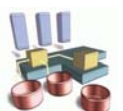
- *WAIT* => setze Teiltransaktion zurück (unilateral ABORT)
- *PREPARED* => *erfrage Transaktionsausgang bei Koordinator (bzw. anderen Rechnern)*



## 2PC: Fehlerbehandlung (3)

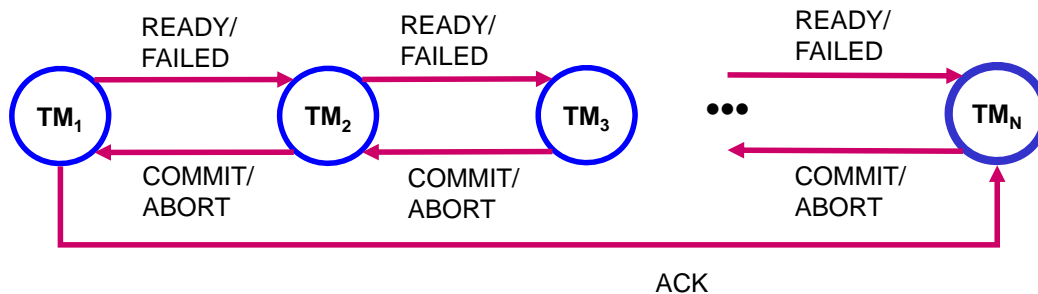
### ■ Rechnerausfall für Agenten:

- Log-Zustand COMMITTED: REDO-Recovery
- Log-Zustand ABORTED bzw. kein 2PC-Log-Satz vorhanden: UNDO-Recovery
- *Log-Zustand PREPARED: Anfrage an Koordinator-Knoten, wie Transaktion beendet wurde (Kordinator hält Information, da noch kein ACK erfolgte)*

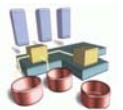
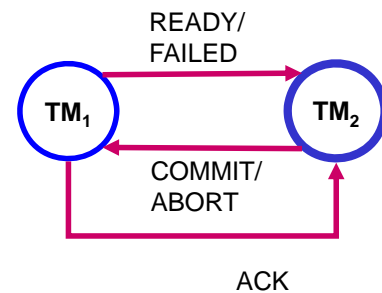


# Lineares 2PC

- Sequentielle Commit-Behandlung, dafür Halbierung der Nachrichtenanzahl:  $(N-1) + N = 2N-1$
- Transfer der Commit-Koordinierung: Koordinatorrolle geht auf letzten Agenten über ("*Last Agent*"-Optimierung)

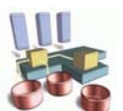
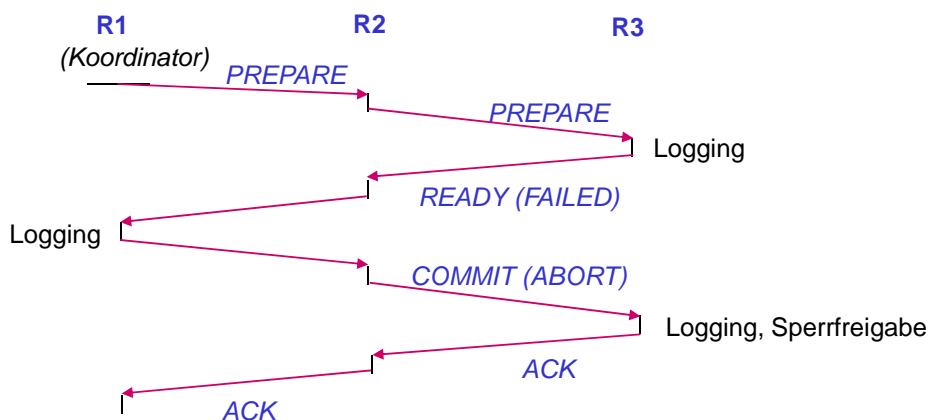


- Besonders vorteilhaft für  $N=2$ : **3 Nachrichten**



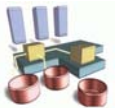
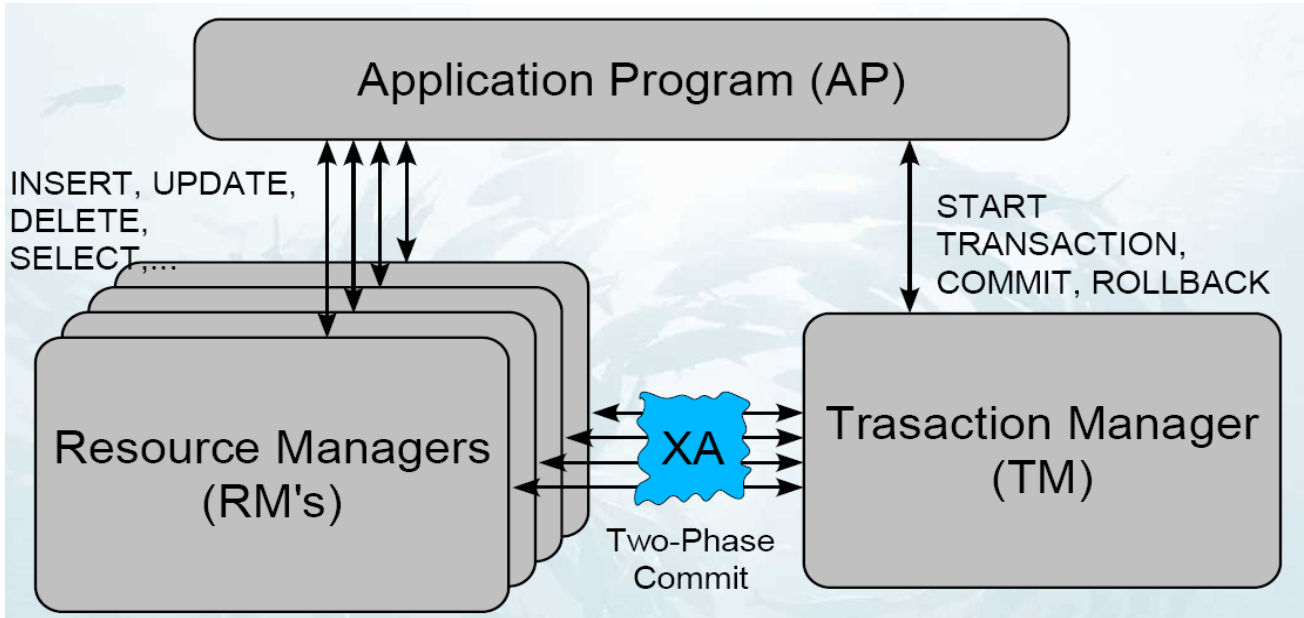
# Hierarchisches 2PC

- allgemeineres Ausführungsmodell mit beliebiger Schachtelungstiefe
  - Nutzung im XA-Standard für verteilte Transaktionen
  - Antwortzeiterhöhung steigt mit Schachtelungstiefe
- Bekanntester Vertreter: *Presumed-Abort-Protokoll*
  - Optimierung für ABORT: keine ACK-Nachrichten und kein synchrones Logging
    - wenn keine Angaben im Log gefunden werden, wird per Default ABORT angenommen
  - Optimierung für lesende Teiltransaktionen:
    - kein Logging, Sperrfreigabe in Phase 1
    - Kommunikation für zweite Phase wird umgangen

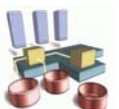
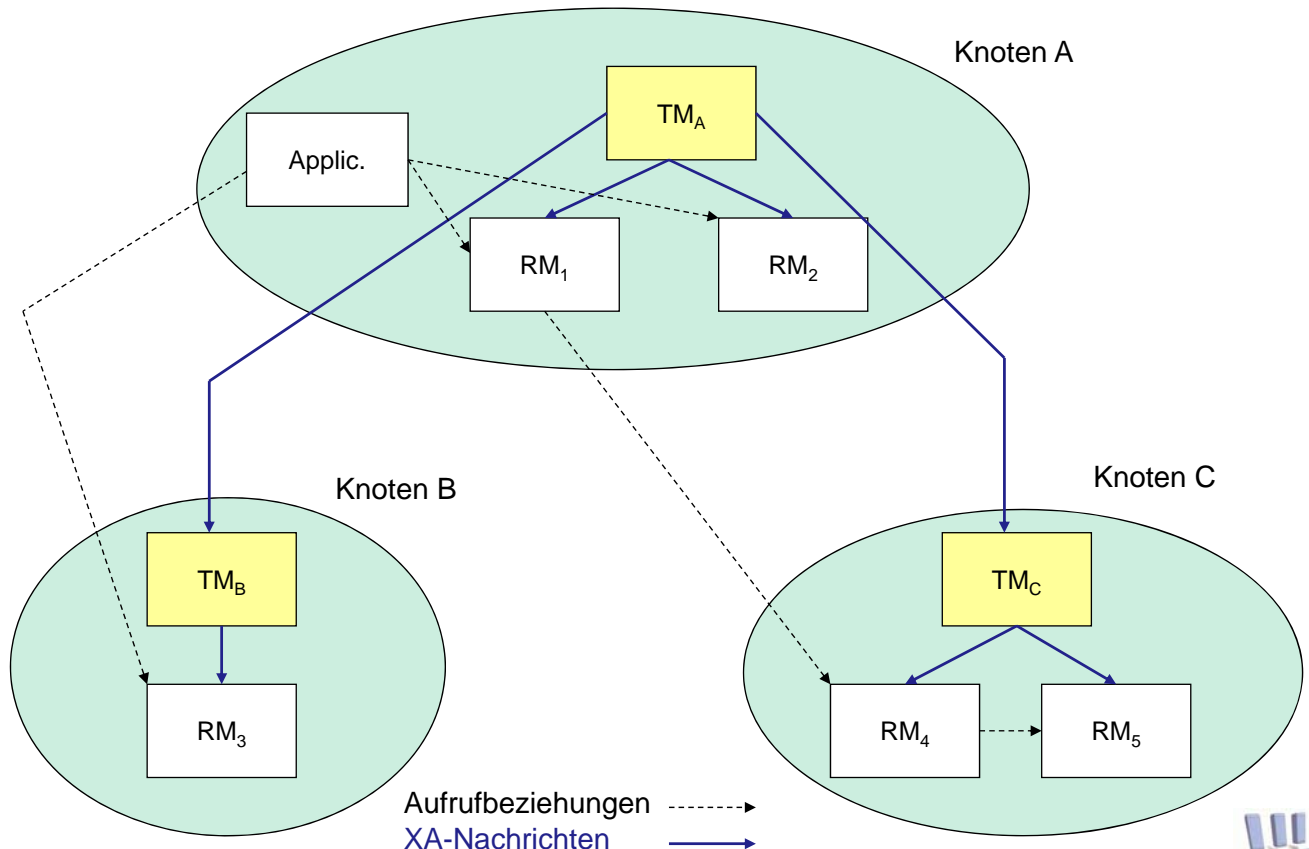


# X/OPEN Modell für verteilte Transaktionen

- TM: Teil von Betriebssystem oder Applikations-Server
- DBS wichtiger RM-Typ
  - XA-Unterstützung durch Oracle, DB2, SQLServer, MySQL, ...

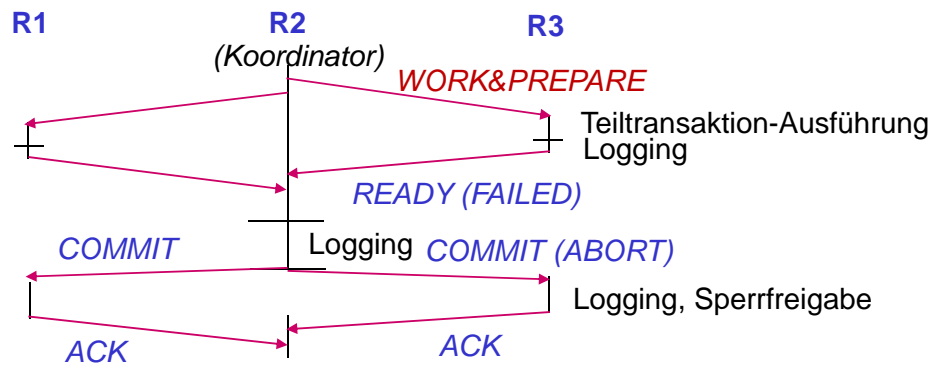


## Verteilte Transaktionsausführung

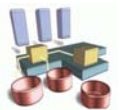




# 1-Phasen-Commit

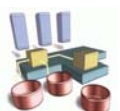
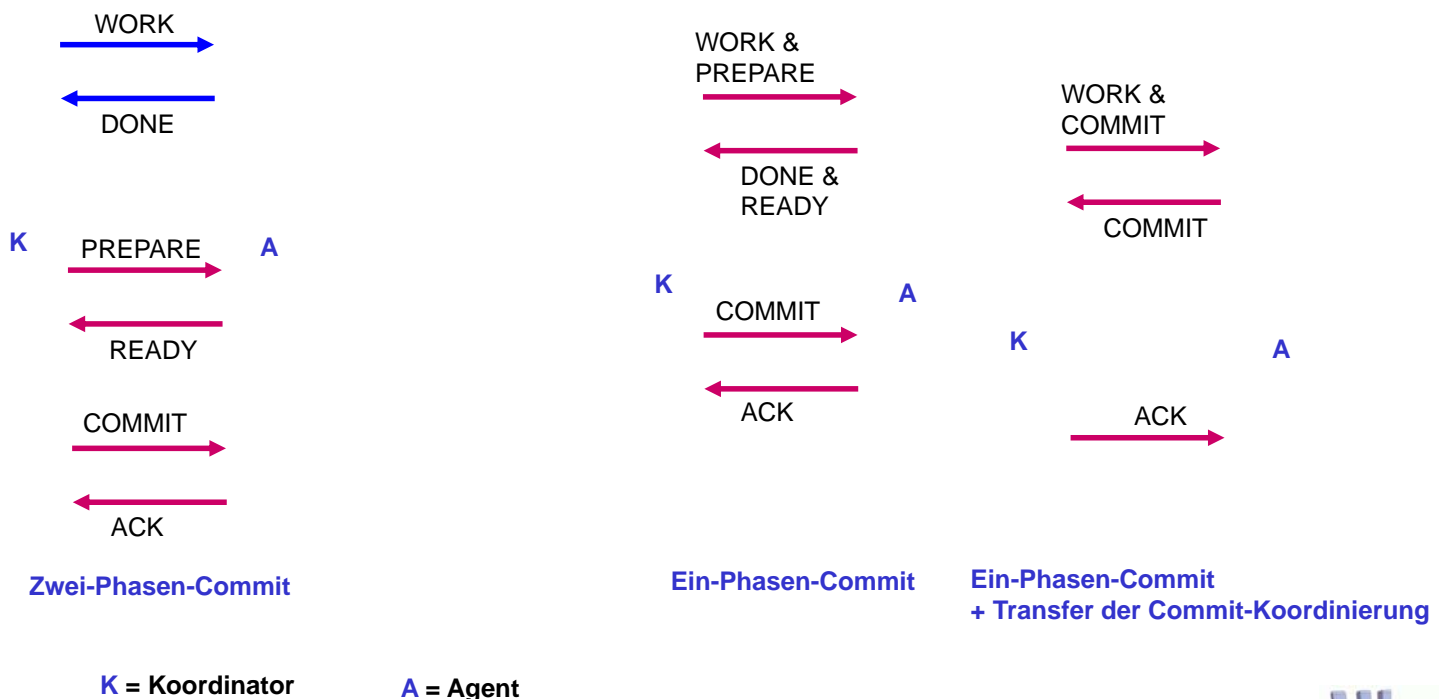


- Teil-Transaktionen sichern ihre Änderungen bereits vor Rückgabe der Ergebnisse an Primärtransaktion
  - nach lokalem Commit am Koordinator steht Erfolg der Transaktion fest
- Einsparung von 2 Nachrichten pro Agent: 2 (N-1) Nachrichten
  - Besonders vorteilhaft für kurze (verteilte) Transaktionen
- Nachteile
  - starke Abhängigkeit vom Koordinator durch frühzeitigen Verzicht auf Unilateral Abort
  - mehrfaches Logging (READY) pro Transaktion und Knoten möglich



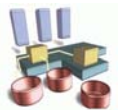
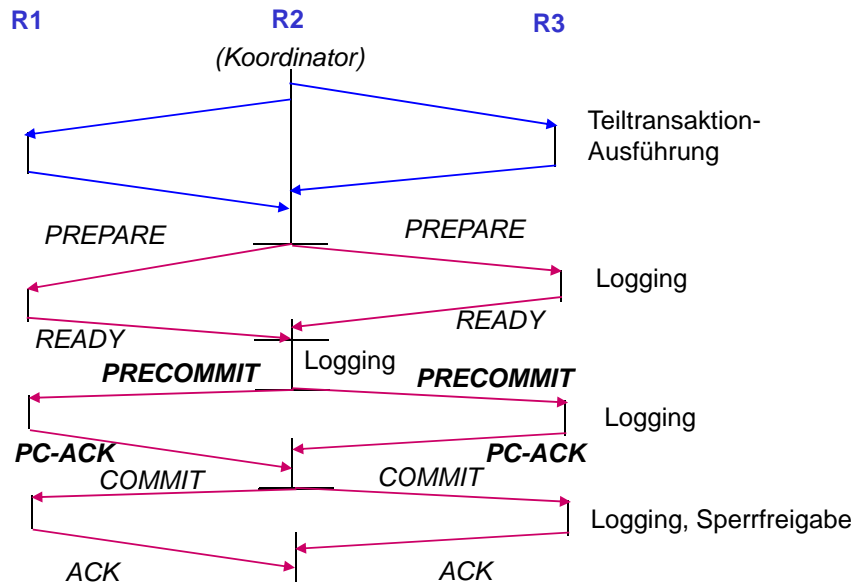
# 1-Phasen-Commit (2)

- Für N=2 kann weitere Nachricht durch **Transfer der Commit-Koordinierung** eingespart werden



# 3-Phasen-Commit

- Nicht-blockierendes Verfahren
- Annahmen:
  - keine Netzwerkpartitionierung
  - höchstens  $K < N$  Rechner fallen gleichzeitig aus



## 3-Phasen-Commit (2)

- ABORT-Behandlung wie im 2PC
- neue Zwischenphase (Zustand PRECOMMIT) falls alle Agenten Phase 1 mit READY abschließen:
  - Koordinator teilt PRECOMMIT allen Teiltransaktionen mit
  - nach Eingang von  $K$  Quittungen (PC-ACK) erfolgt COMMIT-Entscheidung
  - erst jetzt ist Überleben der Transaktion gewährleistet
- **Koordinatorausfall:** Wahl eines neuen Koordinators
  - Erfragen des Transaktionszustandes noch nicht abschließend bearbeiteter Transaktionen bei überlebenden Rechnern:
  - Commit / Abort (bzw. keine Information): Mitteilung des Transaktionsausgangs
  - Precommit bei wenigstens einem überlebenden Rechner: Commit-Protokoll wird von neuem Koordinator mit Verschicken von Precommit-Nachrichten fortgeführt
- löst Blockierungsproblem im Prepared-Zustand
  - lag negative Koordinator-Entscheidung vor: kein Precommit möglich
  - positive Koordinator-Entscheidung: wenigstens 1 Knoten muss Precommit-Zustand haben
  - Precommit-Zustand beim Koordinator: positive oder negative Entscheidung möglich

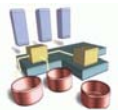


# Nachrichtenbedarf verteilter Commit-Protokolle

N: Anzahl beteiligter Knoten

M: Anzahl lesender Sub-Transaktionen

	Allgemein	Beispiel 1 (N=2, M=0)	Beispiel 2 (N=10, M=5)
1-Phasen-Commit	$2*(N-1)$	2	18
Lineares 2PC	$2*N-1$	3	19
zentralisiertes/hierarchisches 2PC	$4*(N-1)-2M$	4	26
3-Phasen-Commit	$6*(N-1)-4M$	6	34



## Zusammenfassung

- ACID-Eigenschaften für verteilte Transaktionen
- Synchronisation
  - Sicherstellung der globalen Serialisierbarkeit
  - verteilte Sperrverfahren vorzuziehen (geringer Kommunikationsaufwand, weniger Rücksetzungen als mit Zeitstempel- und optimistischen Verfahren)
- Globale Deadlock-Behandlung
  - einfachste Lösung: Timeout
  - Deadlock-Vermeidung (z.B. Wound/Wait) vermeidet Kommunikation, führt jedoch zu unnötigen Rücksetzungen
  - verteilte Deadlock-Erkennung aufwändig: reduziert aber Rücksetzungen
- Verteilte Commit-Protokolle
  - Sicherstellung der Atomarität und Dauerhaftigkeit bei verteilten Änderungen
  - Standardverfahren: hierarchisches 2PC
  - Varianten mit verbesserter Leistungsfähigkeit oder Verfügbarkeit (1PC, 3PC ...)
  - relativ hoher Aufwand

