

# Chapter 7: Lazy Replication

## Update Propagation Protocols (1)

- 1 Alternatives to update replicas:  
synchronous vs. asynchronous:
  - u synchronous – within transaction (eager replication),
  - u asynchronous – separate transaction to update replicas (lazy replication).
- 1 Synchronous updates of replicas: no scalability, in particular if each peer issues transactions. Each transaction must lock quorum or primary copy.

[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Update Propagation Protocols (2)

- 1 Furthermore:
  - u With synchronous updates, we must wait for the slowest one. **Illustration.**
  - u Failure of nodes – similar problem.

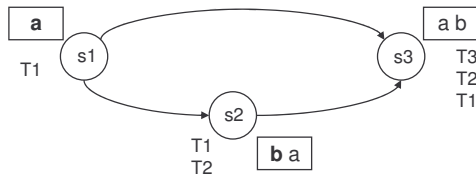
[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Update Propagation Protocols (3)

- 1 Lazy replication – underlying idea:
  - u If update has been successful on one node, it will eventually be successful on all other nodes as well. *Illustrated on subsequent slide.*
  - u I.e., transaction may already commit without having updated all replicas.

[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Lazy Replication & Primary-Copy Technique – Illustration



- 1 Three transactions:  $T_1$  on Site 1 updates a.  
 $T_2$  on Site 2 reads a and writes b.  
 $T_3$  on Site 3 reads a and b.
- 1 Lazy propagation of updates – possible sequence of executions:
  - u on Site 2:  $w_1[a]$   $r_2[a]$   $w_2[b]$
  - u on Site 3:  $w_2[b]$   $r_3[a]$   $r_3[b]$   $w_1[a]$



[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Update Propagation Protocols (4)

- 1 Eager vs. lazy:
  - u Lazy – performance tends to be better (but not by orders of magnitude, depending on the distribution scheme of the replicas),
  - u serializability typically not guaranteed with lazy.

[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Update Propagation Protocols (5)

- 1 Our topic in what follows: lazy protocols guaranteeing serializability.
  1. Protocols making assumptions regarding the arrangement of replicas.
  2. Hybrid protocol without such assumptions that is lazy whenever possible.

[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## System Model/Assumptions

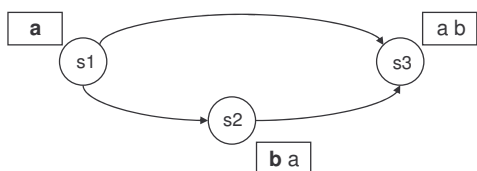
- 1 Each data object has a *primary site*.  
 $\Rightarrow$  *primary copy*,  
*secondary copies/replicas*.
- 1 Transaction has 'originating site'.
  - 1 Transaction can only modify data objects whose primary site = originating site; but may read any data object.
- 1 Nodes (sites) use 2PL.
- 1 Network is reliable; delivery of messages in FIFO order.
- 1 *Primary subtransaction*,  
*secondary subtransaction*.

[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Copy Graph

Copy Graph:

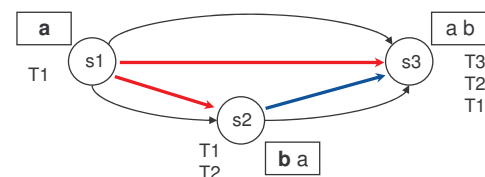
- 1 nodes  $\equiv$  sites.
- 1 Edge from  $s_i$  to  $s_j$  iff primary copy of a data object is on  $s_i$ , and  $\exists$  secondary copy on  $s_j$ .
- 1 Example: three sites, two data objects a and b.



- 1 *Backedges* := set of edges s.t. deletion of such edges makes the graph cycle-free.

[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## Example of Non-serializable Execution



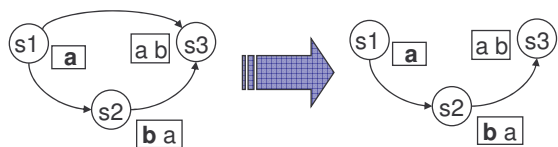
- 1 Three transactions:  $T_1$  on Site 1 updates a.  $T_2$  on Site 2 reads a and writes b.  $T_3$  on Site 3 reads a and b.
- 1 Lazy propagation of updates – possible sequence of executions:
  - u on Site 2:  $w_1[a]$   $r_2[a]$   $w_2[b]$
  - u on Site 3:  $w_2[b]$   $r_3[a]$   $r_3[b]$   $w_1[a]$



[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## DAG(WT) Protocol (1)

- 1 „DAG without Timestamps“
- 1 Prerequisite: copy graph is cycle-free.
- 1 Generate Tree T from copy graph:  $s_i$  child of  $s_j \Rightarrow s_i$  successor of  $s_j$  in T.
- 1 Example:



[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## DAG(WT) Protocol (2)

- 1 Node forwards update transactions to children in T.
- 1 Commit order
  - = order in which transactions arrive at node
  - = order in which transactions are forwarded to children.

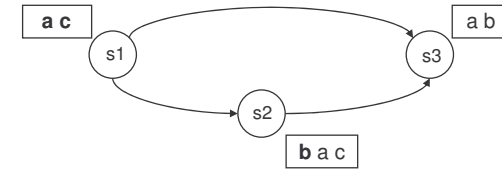
[Introduction](#)  
[DAG\(WT\)](#)  
[DAG\(T\)](#)  
[BackEdge](#)

## DAG(WT) Protocol (3)

- 1 Point that is still open:
  - u Local deadlocks feasible, i.e., transaction does not necessarily commit.
  - u Local deadlock feasible because of
    - conflict of T with transaction  $T_x$  that has not occurred at other sites. Illustration on subsequent slide.
    - Or interaction with transaction that has already occurred at other sites. No handshaking, only commit order is given.
  - u But local transaction must commit.
  - u Thus, victim selection policy should be fair, e.g.: last transaction.

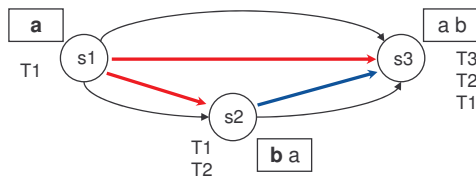
## DAG(WT) Protocol (4)

- 1 Point that is still open (cont.):
  - u Local deadlocks – illustration:



- $T_1: r_1(c) w_1(c) w_1(a)$ ,  $T_2: r_2(a) r_2(c) w_2(b)$
- Possible execution at s2:  
 $w_1(c) r_2(a) r_2(c) w_1(a)$

## Example of Non-serializable Execution



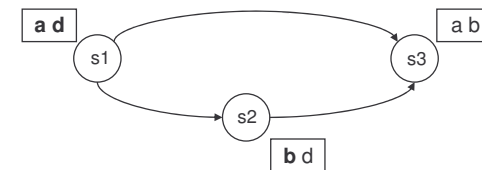
- 1 Three transactions:  $T_1$  on Site 1 updates a.  
 $T_2$  on Site 2 reads a and writes b.  
 $T_3$  on Site 3 reads a and b.
- 1 Lazy propagation of updates – possible sequence of executions:
  - u on Site 2:  $w_1[a] r_2[a] w_2[b]$
  - u on Site 3:  $w_2[b] r_3[a] r_3[b] w_1[a]$



Z

## DAG(WT) Protocol – Observation

- 1 Transaction is routed to nodes that are not relevant.



- 1 Delays.

## DAG(T) Protocol

- 1 „DAG with Timestamps“
- 1 Propagate update transactions along edges of copy graph.
- 1 Primary subtransactions have timestamp that specify execution order.
- 1 Outline of the following:
  - u Structure of timestamp (total order necessary),
  - u protocol itself.
- 1 Structure of timestamps is interesting, since their generation is decentralized.

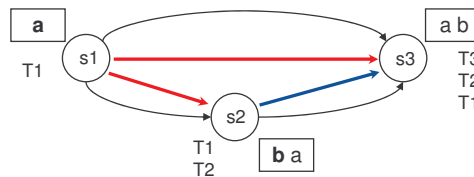
Introduction  
 DAG(WT)  
 DAG(T)  
 - Introduction  
 - Timestamps  
 - Protocol  
 BackEdge

## Timestamps (1)

- 1 In what follows:
  - u auxiliary notion *local timestamp counter*,
  - u auxiliary notion *timestamp of a site*,
  - u *timestamp of a transaction*.
- 1 Acyclicity results in total order of sites;  $s_i < s_{i+k}$
- 1 auxiliary construct – tuple corresponding to Node  $s_i$  :  $(s_i, LTS_i)$   
 $LTS = Local\ Timestamp\ Counter$ , counts the primary subtransactions that have committed there.

Introduction  
 DAG(WT)  
 DAG(T)  
 - Introduction  
 - Timestamps  
 - Protocol  
 BackEdge

## Local Timestamp Counters – Example



$(s1, 0), (s2, 0), (s3, 0)$   
 $(s1, 1)$   
 $(s2, 1)$

Introduction  
 DAG(WT)  
 DAG(T)  
 BackEdge

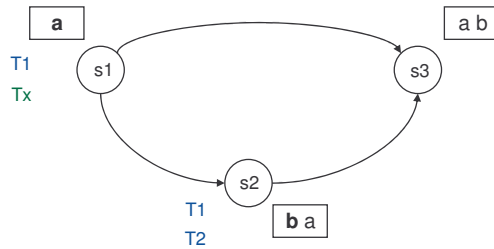
## Timestamps (2)

- 1 *Timestamp of a site  $s_i$* :
  - u vector of tuples,
  - u a tuple for  $s_i$  and further tuples for predecessors of  $s_i$  in copy graph.
  - u Important: tuple in vector ordered by sites.
- 1 Timestamp of a site reflects how many primary subtransactions and how many secondary subtransactions have committed there.

Introduction  
 DAG(WT)  
 DAG(T)  
 - Introduction  
 - Timestamps  
 - Protocol  
 BackEdge

## Timestamps (3)

Example:



Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

1	Site	Timestamp
	s <sub>1</sub>	<(s <sub>1</sub> , 0)>
	s <sub>2</sub>	<(s <sub>1</sub> , 0), (s <sub>2</sub> , 0)>
	...	...

2	Site	Timestamp
	s <sub>1</sub>	<(s <sub>1</sub> , 1)>
	s <sub>2</sub>	<(s <sub>1</sub> , 0), (s <sub>2</sub> , 0)>
	...	...

3	Site	Timestamp
	s <sub>1</sub>	<(s <sub>1</sub> , 1)>
	s <sub>2</sub>	<(s <sub>1</sub> , 1), (s <sub>2</sub> , 0)>
	...	...

4	Site	Timestamp
	s <sub>1</sub>	<(s <sub>1</sub> , 1)>
	s <sub>2</sub>	<(s <sub>1</sub> , 1), (s <sub>2</sub> , 1)>
	...	...

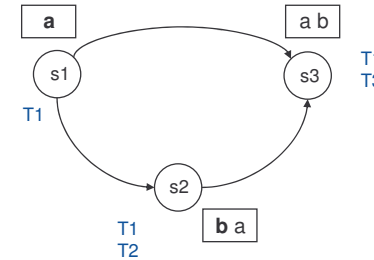
Klemens Böhm

Distributed Data Management: Lazy Replication - 21

## Timestamps (4)

- 1 *Timestamp*  $TS(T_i)$  of a transaction: timestamp of the site of the primary transaction immediately after the point of time of commit.

1 Example:



Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

T1: <(s<sub>1</sub>, 1)>  
T2: <(s<sub>1</sub>, 1), (s<sub>2</sub>, 1)>  
T3: <(s<sub>1</sub>, 1), (s<sub>3</sub>, 1)>

Klemens Böhm

Distributed Data Management: Lazy Replication - 22

## Timestamps (5)

- 1 Order of timestamps shall reflect commit order.
- 1 Lexicographic ordering < of timestamps:

$TS_1 < TS_2 \Leftrightarrow$

- u  $TS_1$  is prefix of  $TS_2$ , or
- u  $TS_1 = X(s_i, LTS_i) Y_1$ ,  $TS_2 = X(s_j, LTS_j) Y_2$ , and
  1.  $s_i > s_j$ , or
  2.  $s_i = s_j$ , and  $LTS_i < LTS_j$ .

Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

Klemens Böhm

Distributed Data Management: Lazy Replication - 23

## Timestamps – Explanation (1)

- 1 Order of timestamps shall reflect commit order.
- 1  $TS_1 < TS_2 \Leftrightarrow TS_1$  is prefix of  $TS_2$ .

Example: <(s<sub>1</sub>, 1)> describes state before <(s<sub>1</sub>, 1), (s<sub>2</sub>, 1)>.

1  $TS_1 < TS_2 \Leftrightarrow TS_1 = X(s_i, LTS_i) Y_1$ ,  $TS_2 = X(s_j, LTS_j) Y_2$ , and  $s_i = s_j$ , und  $LTS_i < LTS_j$ .  
Example: <(s<sub>1</sub>, 1)> describes state before <(s<sub>1</sub>, 2)>.

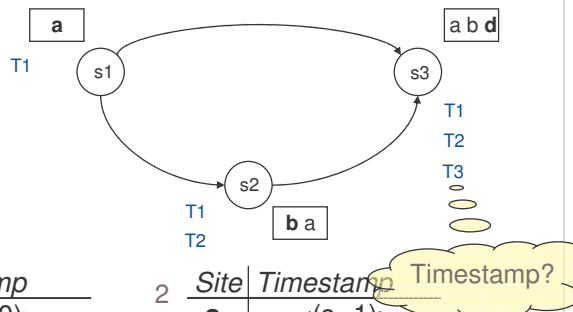
Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

Klemens Böhm

Distributed Data Management: Lazy Replication - 24

# Timestamps – Example, Continued (1)

Example:

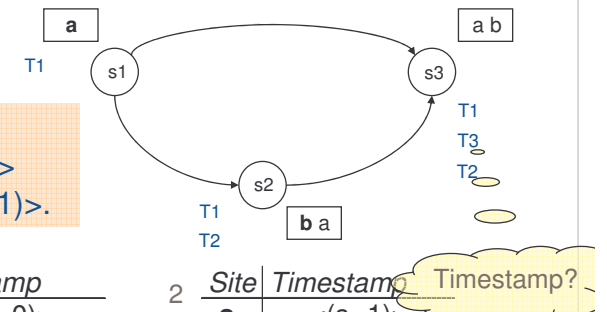


- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 0) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 0), (s_2, 0) \rangle$ |
| ...            | ...                                  |
- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 0), (s_2, 0) \rangle$ |
| ...            | ...                                  |
- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 1), (s_2, 0) \rangle$ |
| ...            | ...                                  |

- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 0), (s_2, 0) \rangle$ |
| ...            | ...                                  |
- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 1), (s_2, 1) \rangle$ |
| ...            | ...                                  |

# Timestamps – Example, Continued (2)

Example:



Hence,  
 $\langle (s_1, 1), (s_2, 1) \rangle$   
 $> \langle (s_1, 1), (s_3, 1) \rangle$ .

- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 0) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 0), (s_2, 0) \rangle$ |
| ...            | ...                                  |
- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 1), (s_2, 0) \rangle$ |
| ...            | ...                                  |
- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 1), (s_2, 0) \rangle$ |
| ...            | ...                                  |

- | Site           | Timestamp                            |
|----------------|--------------------------------------|
| s <sub>1</sub> | $\langle (s_1, 1) \rangle$           |
| s <sub>2</sub> | $\langle (s_1, 1), (s_2, 1) \rangle$ |
| ...            | ...                                  |

## Data Structures

Data structure for each node:

- timestamp vector of the node  
(= timestamp of the last committed secondary subtransaction + tuple of the node),
- waiting queues,  
one for each predecessor.

## Primary Transaction

Primary transaction commits at Node s<sub>i</sub>:

- increment LTS<sub>i</sub> (Local Timestamp Counter),
  - TS(T<sub>i</sub>) := TS(s<sub>i</sub>)
  - Send secondary (sub-)transaction to children of s<sub>i</sub>.
- 1 Timestamp of a site reflects how many primary subtransactions and how many secondary subtransactions have committed there.

## Secondary Transaction (1)

- 1 Assumption: only one secondary transaction at a time.
- 1 One waiting queue for each direct predecessor of the site in the copy graph.
- 1 Choose transaction from queues with minimal timestamp.
- 1 There must be at least one transaction in each queue before computing the min timestamp.
- 1 Idle nodes should commit dummy transactions from time to time.

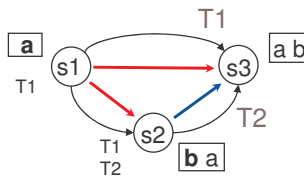
Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

## Secondary Transaction (2)

- 1 After commit:  $TS(s_i) := TS(T_i)(s_i, LTS_i)$
- 1 How do we know that next transaction in commit order is not being propagated through the network? Corresponding queue would be empty, furthermore: FIFO order.

Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

## Continuation of Example



- 1  $T_1$  has Timestamp  $(s_1, 1)$ ,
- 1  $T_2$  has Timestamp  $(s_1, 1), (s_2, 1)$ .  
(When  $T_1$  commits on  $s_2$ , the site timestamp is set to  $(s_1, 1), (s_2, 0)$ .)
- 1 Site  $s_3$ :  
timestamp of  $T_1$  is prefix of the one of  $T_2$ .  
 $\Rightarrow T_1$  is executed there before  $T_2$ .

Introduction  
DAG(WT)  
DAG(T)  
- Introduction  
- Timestamps  
- Protocol  
BackEdge

Z

## BackEdge Protocol – Motivation

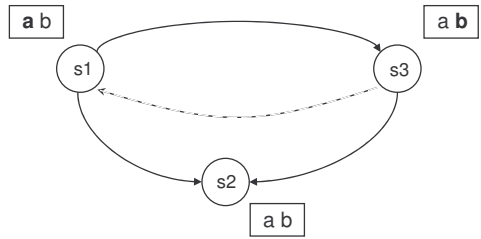
- 1 Copy graph must be acyclic for DAG(WT) and DAG(T) protocols.
- 1 Example:
  - u two sites  $s_1$  and  $s_2$ .
  - u  $s_1$  holds primary copy of  $a$  and copy of  $b$ ,  $s_2$  vice versa.
  - u  $T_1$  at Node  $s_1$  reads  $b$  and updates  $a$ ,  $T_2$  at Node  $s_2$  reads  $a$  and updates  $b$ .
  - u Both transactions execute concurrently and commit. **Illustration.**
  - u No serializability.

Introduction  
DAG(WT)  
DAG(T)  
BackEdge



## BackEdge Protocol – Overview

- Hybrid: for some replicas eager updates, for other ones lazy.
- Can be described both as extension of DAG(WT) and DAG(T).
- $G_{dag}$  – results from  $G$  by removing backedges.



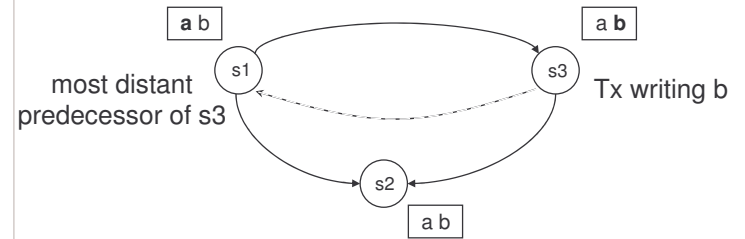
Klemens Böhm

Distributed Data Management: Lazy Replication - 33

Introduction  
DAG(WT)  
DAG(T)  
BackEdge

## BackEdge Protocol – Terminology (1)

- Backedge from  $s_i$  to  $s_j$ .  
 $\Rightarrow s_j$  is predecessor of  $s_i$  in  $G_{dag}$ .
- $T_i$  – primary subtransaction with node  $s_i$ .
- „Backedge subtransactions“: transactions  $S_1, \dots, S_j$  at predecessor  $s_{i1}, \dots, s_{ij}$  of  $s_i$  in  $T$ .  
 $s_{i1}$  „most distant“ of  $s_i$  etc.



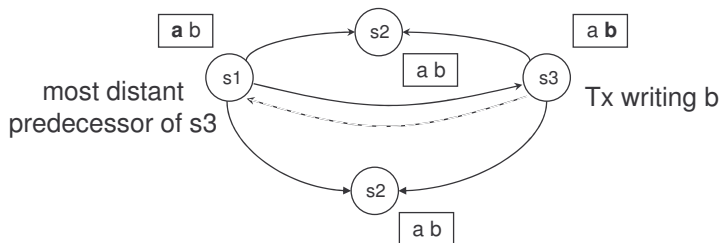
Klemens Böhm

Distributed Data Management: Lazy Replication - 34

Introduction  
DAG(WT)  
DAG(T)  
BackEdge

## BackEdge Protocol – Terminology (2)

- Illustration:



- How many predecessors does s3 have?

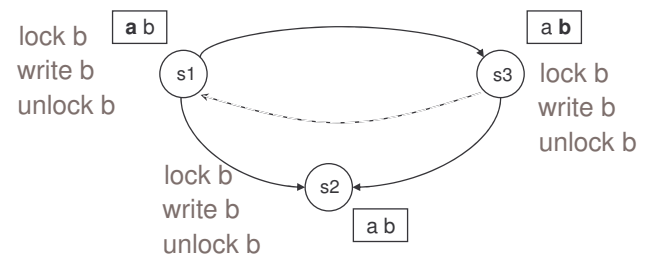
Klemens Böhm

Distributed Data Management: Lazy Replication - 35

Introduction  
DAG(WT)  
DAG(T)  
BackEdge

## BackEdge Protocol

- After execution of  $T_i$ ,  $S_1$  is sent to  $s_1$  – no commit, locks are not released.
- $S_2, \dots, S_j$  – as before, but without commit and without releasing locks.
- 2PC for  $T_i$ ,  $S_1, \dots, S_j$ .
- Remaining secondary subtransactions: lazy, as in DAG(T).



Klemens Böhm

Distributed Data Management: Lazy Replication - 36

Introduction  
DAG(WT)  
DAG(T)  
BackEdge

## BackEdge Protocol – Discussion

- 1 ,No magic' here.
- 1 Better than primary-copy schemes, according to simulations. Upto Factor 5, if there are many readers and few backedges.
- 1 Implementation on top of commercial DBMSs relatively easy.

## Literature

- 1 Yuri Breitbart et al.: Update Propagation Protocols for Replicated Databases. Proceedings SIGMOD'99.

## Potential Exam Questions

- 1 Illustrate that lazy replication schemes, when designed carelessly, may lead to inconsistencies.
- 1 Why is the topology of the copy graph important in the context of lazy replication?
- 1 Explain the different approaches from the lecture that ensure consistency with lazy replication.