

# Transactions for Distributed Wikis on Structured Overlays<sup>\*</sup>

Stefan Plantikow, Alexander Reinefeld, and Florian Schintke

Zuse Institute Berlin

**Abstract.** We present a transaction processing scheme for structured overlay networks and use it to develop a distributed Wiki application based on a relational data model. The Wiki supports rich metadata and additional indexes for navigation purposes.

Ensuring consistency and durability requires handling of node failures. We mask such failures by providing high availability of nodes by constructing the overlay from replicated state machines (*cell model*). Atomicity is realized using two phase commit with additional support for failure detection and restoration of the transaction manager. The developed transaction processing scheme provides the application with a mixture of pessimistic, hybrid optimistic and multiversioning concurrency control techniques to minimize the impact of replication on latency and optimize for read operations. We present pseudocode of the relevant Wiki functions and evaluate the different concurrency control techniques in terms of message complexity.

*Keywords.* Distributed transactions, content management systems, structured overlay networks, consistency, concurrency control.

## 1 Introduction

*Structured overlay networks (SONs)* provide a scalable and efficient means for storing and retrieving data in distributed environments without central control. Unfortunately, in their most basic implementation, SONs do not offer any guarantees on the ordering of concurrently executed operations.

Transaction processing provides concurrently executing clients with a single, consistent view of a shared database. This is done by bundling client operations in a transaction and executing them as if there was a global, serial transaction execution order. Enabling structured overlays to provide transaction processing support is a sensible next step for building *consistent* decentralized, self-managing storage services.

We propose a transactional system for an Internet-distributed content management system built on a structured overlay. Our emphasis is on supporting

---

<sup>\*</sup> This work was partially supported by the EU projects SELFMAN and CoreGRID.

transactions in dynamic decentralized systems where nodes may fail at a relatively high rate. The chosen approach provides clients with different concurrency control options to minimize latency.

The article is structured as follows: Section 2 describes a general model for distributed transaction processing in SONs. The main problem addressed is masking the unreliability of nodes. Section 3 presents our transaction processing scheme focusing on concurrency control. This scheme is extended to the relational model and exemplified using the distributed Wiki in Section 4. Finally, in Section 5, we evaluate the different proposed transaction processing techniques in terms of message complexity.

## 2 Transactions on Structured Overlays

Transaction processing guarantees the four ACID properties: *Atomicity* (either all or no data operations are executed), *consistency* (transaction processing never corrupts the database state), *isolation* (data operations of concurrently executing transactions do not interfere with each other), *durability* (results of successful transactions survive system crashes). Isolation and consistency together are called *concurrency control*, while *database recovery* refers to atomicity and durability.

*Page model.* We only consider transactions in the *page model* [1]: The database is a set of uniquely addressable, single objects. Valid elementary operations are reading and writing of objects, and transaction abort and commit. The model does not support predicate locking. Therefore, phantoms can occur and consistent aggregation queries are not supported. The page model can naturally be applied to SONs. Objects are stored under their identifier using the overlay's policy for data placement.

### 2.1 Distributed Transaction Processing

Distributed transaction processing guarantees the ACID properties in scenarios where clients access multiple databases or different parts of the same database located on different nodes. Access to local databases is controlled by *resource manager (RM)* processes at each participating node. Additionally, for each active transaction, one node takes the role of the *transaction manager (TM)*. The TM coordinates with the involved RMs to execute a transaction on behalf of the client. The TM also plays an important role during the execution of the distributed atomic commit protocol.

Distributed transaction processing in a SON requires distribution of resource and transaction management. The initiating peer can act as TM. For resource management, it is necessary to minimize the communication overhead between RM and storing node. Therefore, in the following, we assume that each peer of the overlay performs resource management for all objects in its keyspace partition.

## 2.2 The Cell Model for Masking Churn

Distributing the resource management over all peers puts tight restrictions on messages delivered under transaction control. Such messages may only be delivered to nodes that are currently responsible for the data. This property is known as *lookup consistency*. Without lookup consistency, a node might erroneously grant a lock on a data item or deliver outdated data. It is an open question how lookup consistency can be guaranteed efficiently in the presence of frequent and unexpected node failures (*churn*). Some authors (e.g. [2]) have suggested protocols that ensure consistent lookup if properly executed by *all* joining and leaving nodes. Yet large-scale overlays are subject to considerable amounts of churn [3] and therefore correct transaction processing requires masking it.

*Cell model.* Instead of constructing the overlay network using single nodes, we propose to build the overlay out of *cells*. Each cell is a dynamically sized group of physical nodes [4] that constitute a *replicated state machine (RSM, [5])*. Cells utilize the chosen RSM algorithm to provide replicated, atomic operations and high availability. This can be exploited to

- mask churn and therefore guarantee lookup consistency,
- provide stable storage for transactional durability,
- ensure data consistency using atomic operations,
- minimize overhead for routing to replicas (cell nodes form a clique).

For the underlying nodes, we assume the *crash-stop* failure model. This model is common for SONS because it is usually unknown whether a disconnected node will rejoin again later. We do not cover the distribution of physical nodes on cells, nor do we consider Byzantine failures. We assume that cells never fail unexpectedly and always execute the overlay algorithm orderly. If too many cell nodes fail, the cell destroys itself by executing the overlay’s leave protocol. The data items are re-distributed among neighboring cells. For simplification, we also assume that the key-space partition associated to each cell does not change during transaction execution.

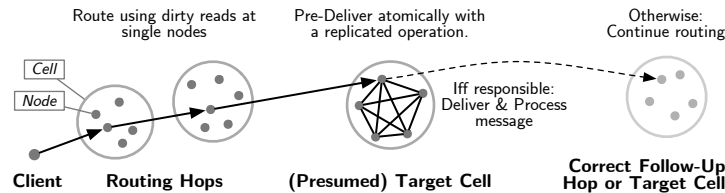


Fig. 1. Cell routing using dirty reads.

*Cell routing.* The execution of replicated operations within the cells comes at a considerable cost: RSMs are implemented using some form of atomic broadcast which in turn depends on an implementation of a consensus protocol. Yet, modern consensus algorithms like *Fast Paxos* [6] require at least  $N(\lfloor 2N/3 \rfloor + 1)$  messages for consensus between  $N$  nodes. While this cost is hardly avoidable for consistent replication, it is as well unacceptable for regular message routing. Hence we propose to use dirty reads (i.e. to read the state of one arbitrary node). When the state of a node and its cell are temporarily out of sync, routing errors may occur. To handle this, the presumed target cell will pre-deliver the message using a replicated operation (Fig. 1). Pre-delivery first checks, whether the presumed target cell is currently responsible for the message. If that is the case, the message is delivered and processed regularly. Otherwise, message routing is restarted from the node that initiated pre-delivery.

Replicated operations will only be executed at a node after all of its predecessor operations have been finished. Therefore, at pre-delivery time, the presumed target cell either actually is responsible for the message or a previously executed replicated operation has changed the cell's routing table consistently such that the *correct* follow-up routing hop or target cell for the message is known. A message reaches its destination under the assumption that cell routing table changes are sufficiently rare, and intermediate hops do not fail.

### 3 Concurrency Control and Atomic Commit in SONs

In this section, we present appropriate concurrency control and atomic commit techniques for overlays based on the cell model.

*Atomic Operations.* Using RSMs by definition [5] allows the execution of atomic and totally ordered operations. This already suffices to implement transaction processing, e.g. by using *pessimistic two phase locking (2PL)* and an additional distributed atomic commit protocol. However, each replicated operation is expensive, thus any efficient transaction processing scheme for cell-structured overlays should aim at minimizing the number of replicated operations.

*Optimistic concurrency control (OCC).* OCC executes transactions against a *local* working copy (working phase). This copy is validated just before the transaction is committed (validation phase). The transaction is aborted if conflicts are detected during validation. As every node has (a possibly temporarily deviating) local copy of its cell's shared state, OCC is a prime candidate for reducing the number of replicated operations by executing the transaction against single nodes of each involved cell.

#### 3.1 Hybrid Optimistic Concurrency Control

Plain OCC has the drawback that long-running transactions using objects which are frequently accessed by short-running transactions may suffer starvation due

to consecutive validation failures. This is addressed by *hybrid optimistic concurrency control (HOCC, [7])* under the assumption of *access invariance*, i.e. repeated executions of the same transaction have identical read and write sets.

HOCC works by executing *strong two phase locking (SS2PL)* for the transaction’s read and write sets at the beginning of the validation phase. In case of a validation failure, locks are kept and the transaction logic is re-executed. Now, access invariance ensures that this second execution cannot fail because all necessary locks are already held by the transaction. However, it is required that optimistically read values do not influence the result of the re-execution phase. Otherwise, consistency may be violated.

The use of SS2PL adds the benefit that no distributed deadlock detection is necessary if a global validation order between transactions is established. A possible technique for this has been described by Agrawal et. al [8]: Each cell  $v$  maintains a strictly monotonic increasing timestamp  $t_v$  for the largest, validated transaction. Before starting the validation, the transaction manager suggests a validation time stamp  $t > t_v$  to all involved cells. After each such cell has acknowledged that  $t > t_v$  and updated  $t_v$  to  $t$ , the validation phase is started. Otherwise the algorithm is repeated. Gruber [9] optimized this approach by including the current  $t_v$  in every control message.

### 3.2 Distributed Atomic Commit

*Distributed atomic commit (DBAC)* requires consensus between all transaction participants on the transaction’s termination state (committed or aborted). If DBAC is not guaranteed, the ACID properties are violated.

We propose a blocking DBAC protocol that uses cells to treat TM failures by replicating transaction termination state. Every transaction is associated with a unique identifier (TXID). The overlay cell corresponding to that TXID is used to store a *commit record* holding the termination state and the address of the TM node (an arbitrary, single node of the TXID cell). If no failures occur, regular *two-phase atomic commit (2PC)* is executed. Additionally, after all prepared-messages have been received and before the final commit messages are sent, the TM first writes the commit record. If the record is already set to abort, the TM aborts the transaction. If RMs suspect a TM failure, they read the record to either determine the termination state or initiate transaction abort. Optionally, RMs can restore the TM by selecting a new node and updating the record appropriately. Other RMs will notice this when they reread the modified record.

### 3.3 Read-only Transactions

In many application scenarios simple read-only transactions are much more common than update transactions. Therefore we optimize and extend our transaction processing scheme for read-only transactions by applying techniques similar to *read-only multiversioning (ROMV, [10])*.

All data items are versioned using unique timestamps generated from each node's loosely synchronized clock and globally unique identifier. Additionally, we maintain a *current version* for each data item. This version is accessed and locked exclusively by HOCC transactions as described above and implicitly associated with the cell's maximum validation timestamp  $t_v$ . The current version decouples ROMV and HOCC.

Our approach moves newly created versions to the future such that they never interfere with read operations from ongoing read-only transactions. This avoids the cost associated with distributed atomic commit for read-only transactions but necessitates it to execute reads as replicated operations. Read-only transactions are associated with their start time. Every read operation is executed as a replicated operation using the standard multiversioning rule [11]: The result is the oldest version which is younger than the transaction start time. If this version is the current version, the maximum validation timestamp  $t_v$  will be updated. This may block the read operation until a currently running validation is finished. Update transactions create new versions of all written objects using  $t > t_v$  during atomic commit.

## 4 Algorithms for a Distributed Wiki

In this section, we describe the basic algorithms of a distributed content management system built on a structured overlay with transaction support.

### 4.1 Mapping the Relational Model

So far we only considered uniquely addressable, uniform objects. In practice, many applications use more complex, relational data structures. This raises the question of how multiple relations with possibly multiple attributes can be stored in a single structured overlay. To address this, first, we assume that the overlay supports range queries [12, 13] over a finite number of index dimensions.

Storing multiple attributes requires mapping them on index dimensions. As the number of available dimensions is limited, it is necessary to partition the attributes into disjoint groups and map these groups instead. The partition must be chosen in such a way that fast primary-key based access is still possible. Depending on their group membership, attributes are either primary, index, or non-indexed data attributes. Multiple relations can be modeled by prefixing primary keys with a unique relation identifier.

### 4.2 Notation

Table 1 contains an overview of the pseudocode syntax from [14]. Relations are represented as sets of tuples and written in CAPITALS. Relation tuples are addressed by using values for the primary attributes in the fixed order given by the relation. For reasons of readability, tuple components are addressed using unique labels (Such labels can easily be converted to positional indexes). Range queries are expressed using labels and marked with a "?".

**Table 1.** Pseudocode notation

Syntax	Description
<b>Procedure</b> Proc ( $arg_1, arg_2, \dots, arg_n$ )	Procedure declaration
<b>Function</b> Fun ( $arg_1, arg_2, \dots, arg_n$ )	Function declaration
<b>begin ... commit (abort) transaction</b>	Transaction boundaries
ADDRESS <sup>ZIB</sup> "	Read tuple from relation
ADDRESS <sup>ZIB</sup> " $\leftarrow$ ("Takustr. 7", "Berlin")	Write tuple to relation
$\Pi_{attr_1, \dots, attr_n}(M) = \{\pi_{attr_1, \dots, attr_n}(t) \mid t \in M\}$	Projection
$\forall t \in \text{tuple set} : \text{RELATION} \stackrel{+}{\leftarrow} t$ bzw. $\stackrel{-}{\leftarrow} t$	Bulk insert and delete
DHT <sup>?</sup> <sub>key<sub>1</sub>="a", key<sub>2</sub></sub> or DHT <sup>?</sup> <sub>key<sub>1</sub>="a", key<sub>2</sub>=*</sub>	Range query with wildcard

### 4.3 Wiki

A *Wiki* is a content management system that embraces the principle of minimizing access barriers for non-expert users. Wikis like [www.wikipedia.org](http://www.wikipedia.org) comprise millions of pages written in a simplified, human-readable markup syntax. Each page has a unique name which is used for hyperlinking to other Wiki pages. All pages can be read and edited by any user, which may result in many concurrent modification requests for hotspot pages. This makes Wikis a perfect test-case for our distributed transaction algorithm.

Modern Wikis provide a host of additional features, particularly to simplify navigation. In this paper we exemplarily consider backlinks (a list of all the other pages linking to a page) and recent changes (a list of recent modifications of all Wiki pages). We model our Wiki using the following two relations:

Relation	Primary attributes	Index attributes	Data attributes
CONTENT	<i>pageName</i>	<i>ctime</i> (change time)	<i>content</i>
BACKLINKS	<i>referencing</i> (page), <i>referenced</i> (page)	-	-

Wiki operations use transactions to maintain global consistency invariants:

- CONTENT always contains the current content for all pages,
- BACKLINKS contains proper backlinks for all pages contained in CONTENT,
- users cannot modify pages whose content they have never seen (explained below).

The function WikiRead (Alg. 4.1) delivers the content of a page and all backlinks pointing to it. This requires a single read for the content and a range query to obtain the backlinks. Both operations can be executed in parallel.

The function WikiWrite (Alg. 4.2) is more complex because conflicting writes by multiple users must be resolved. This can be done by serializing the write requests using locks or request queues. If conflicts are detected during (atomic) writes by comparing last read and current content, the write operation is aborted. Users may then manually merge their changes and retry. This approach is similar to the compare-and-swap instructions used in modern microprocessors and to the concurrency control in version control systems.<sup>1</sup> We realize the compare-and-swap in WikiWrite by using transactions for our distributed Wiki. First, we precompute which backlinks should be inserted and deleted. Then, we compare the current and old page content and abort if they differ. Otherwise all updates are performed by writing the new page content and modifying BACKLINKS. The update operations again can be performed in parallel.

---

**Algorithm 4.1** WikiRead: Read page content

---

```

1: function WikiRead (pageName)
2:   begin transaction read-only
3:     content  $\leftarrow \pi_{\text{content}}(\text{CONTENT}_{\text{pageName}})$ 
4:     backlinks  $\leftarrow \Pi_{\text{referenced}}(\text{BACKLINKS}_{\text{referencing}=\text{pageName}, \text{referenced}}^?)$ 
5:   commit transaction
6:   return content, backlinks
7: end function

```

---



---

**Algorithm 4.2** WikiWrite: Write new page content and update backlinks

---

```

1: procedure WikiWrite (pageName, contentold, contentnew)
2:   refsold  $\leftarrow \text{Refs}(\text{content}_{\text{old}})$ 
3:   refsnew  $\leftarrow \text{Refs}(\text{content}_{\text{new}})$ 
4:   refsdel  $\leftarrow \text{refs}_{\text{old}} \setminus \text{refs}_{\text{new}}$  — precalculation
5:   refsadd  $\leftarrow \text{refs}_{\text{new}} \setminus \text{refs}_{\text{old}}$ 
6:   txStartTime  $\leftarrow \text{CurrentTimeUTC}()$ 
7:   begin transaction
8:     if  $\pi_{\text{content}}(\text{CONTENT}_{\text{pageName}}) = \text{content}_{\text{old}}$  then
9:        $\text{CONTENT}_{\text{pageName}} = (\text{txStartTime}, \text{content}_{\text{new}})$ 
10:       $\forall t \in \{(ref, \text{pageName}) \mid ref \in \text{refs}_{\text{add}}\} : \text{BACKLINKS} \stackrel{\perp}{\leftarrow} t$ 
11:       $\forall t \in \{(ref, \text{pageName}) \mid ref \in \text{refs}_{\text{del}}\} : \text{BACKLINKS} \stackrel{\leftarrow}{\leftarrow} t$ 
12:     else
13:       abort transaction
14:     end if
15:   commit transaction
16: end procedure

```

---

<sup>1</sup> Most version control systems provide heuristics (e.g. merging of different versions) for automatic conflict resolution that could be used for the Wiki as well.



---

**Algorithm 4.3** SetPageMetadata: Write page metadata attributes

---

**Require:**  $changeEnv$  environment describing changes to be made

- 1: **procedure** SetPageMetadata ( $pageName$ ,  $content_{old}$ ,  $changeEnv$ )
- 2:   **begin transaction**
- 3:     **if**  $\pi_{content}(CONTENT_{pageName}) = content_{old}$  **then**
- 4:        $\forall (anAttrName \leftarrow anAttrValue) \in changeEnv :$
- 5:          $METADATA_{pageName, anAttrName} \leftarrow anAttrValue$
- 6:     **else**
- 7:       **abort transaction**
- 8:     **end if**
- 9:   **commit transaction**
- 10: **end procedure**

---

The list of recently changed pages can be generated by issuing a simple range query inside a transaction and sorting the results appropriately.<sup>2</sup>

#### 4.4 Wiki with Metadata

Often it is necessary to store additional metadata with each page (e.g. page author, category). To support this, we add a third relation `METADATA` with primary key attributes  $pageName$  and  $attrName$  and data attribute  $attrValue$ . Alternatively we could also add metadata attributes to `CONTENT`. But this would not be scalable as current overlays only provide a limited number of index dimensions.

Modifying page metadata (Alg. 4.3) requires verifying that the page has not been changed by some other transaction. Otherwise new metadata could be associated wrongly to a page (This is similar to storing wrong backlinks). For reading page metadata, a simple range query suffices [14].

## 5 Evaluation

The presented algorithms for ensuring consistency mainly require the atomicity property while only few restrictions are placed on the serial execution order of operations. Thus in theory, a high degree of concurrency is possible. This is especially interesting for range queries like `RecentChanges` which can utilize the overlay's capabilities to multicast to many nodes in parallel.

Table 2 shows the communication overhead of various concurrency control schemes. We compare the different schemes using an example transaction that consists of  $k$  serial steps. Each step executes data operations in parallel on  $N$  cells (one operation per cell).

For every scheme, we distinguish the number and type of operations necessary to carry out the transaction:  $U$  is a simple unreplicated operation,  $R$  is a replicated operation, and  $L$  is a lookup (routing) operation. The cost is split into one-time (initial and DBAC) overhead, the cost per step, and the total cost.

---

<sup>2</sup> The complete range query is:  $\{CONTENT_{pageName=*, ctime=*}^{\leftarrow} \}_{\# < resultLimit}$

**Table 2.** Comparison of concurrency control methods

Transaction type	One-time overhead for $N$ cells	Ops per step on $N$ cells in parallel	Total for $k$ serial steps
(1) Atomic Write	$1L$	$1R$	$1L + 1R$ , because $k, N = 1$
(2) Read-Only Trans.	$NL$	$NR$	$NL + kNR$
(3) Pess. 2PL + 2PC	$NL + 2NR$	$NR$	$NL + (k + 1)NR$
(4) Hyb. Opt. + 2PC	$NL + 2NR$	$NU$	$NL + (k - 1)NU + 2NR$
(5) Hyb. Opt. + 2PC + Validation Error	$NL + 3NR$	$2NU$	$NL + (2k - 2)NU + 3NR$

(2) to (4) use the 2PC variant described in 3.2. For our evaluation, we assume that no failures occur during the commit.

Totals include DBAC costs and take the possible combined sending of messages into account (e.g. combining last write operation with validate and prepare). The evaluated concurrency control schemes are:

- (1) a simple, replicated operation on a single cell,
- (2) a read-only multiversioning transaction (Sec. 3.3),
- (3) a pessimistic 2PL transaction,
- (4) a HOCC (Sec. 3.1) transaction without validation failure, and
- (5) a HOCC transaction with validation failure and re-execution of transaction logic.

HOCC reduces the number of necessary replicated operations for  $k > 1$ . For  $k = 1$  and a transaction on a single cell, ACID is already provided by using a RSM and no DBAC is necessary. For  $k = 1$  and a transaction over multiple cells, HOCC degenerates into 2PL: the data operations on the different cells are combined with validate-and-prepare messages and executed using single replicated operations.

Read-only transactions use more replicated operations but save the DBAC costs of HOCC. This makes them well-suited for quick, parallel reads. But long running read transactions might be better off using HOCC if the performance gained by optimism outweighs DBAC overhead and validation failure chance.

Using cells yields an additional benefit. If replication was performed above the overlay layer, additional routing costs of  $(r - 1)N$  lookup messages would be necessary ( $r$  is the number of replicas).

## 6 Related Work

Mesaros et. al describe a transaction processing scheme for overlays based on 2PL [15]. Lock conflicts are resolved by giving higher priority to older transactions and forcing the losing transaction into the 2PL shrinking phase. Transactions are executed by forming a dynamic multicast group consisting of all

participating nodes. The article does not address issues of lookup consistency and replication.

OceanStore [16] uses a two-tier approach for multiversioning-based replication. On the first layer, a small set of replicas forms a primary ring. On the second layer, additional replicas cache object versions. Replicas are located using the Tapestry overlay network. Primary ring replicas use a Byzantine agreement protocol to serially execute atomic operations.

Etna [17] is a system for executing atomic read and write operations in a Chord-like overlay network. Operations are serialized using a primary copy and replicated over  $k$  successors using a consensus algorithm.

Both articles do not describe how full transaction processing can be built on top of atomic operations. For OceanStore, multiversioning [11] is proposed [16]. The inherent cost of replicated transaction execution is handled using the caching tier. However, this comes at the price of reduced consistency.

As an alternative to our solution for atomic commitment, Moser et al. [18] describe a non-blocking approach based on Paxos commit. Their solution treats the set of all replicas of all accessed items as a whole and fixes this set at commit time. They suggest the use of symmetric replication [19] to achieve availability. Instead of using RSMs inside cells, encoding schemes like Reed-Solomon codes could be used, as proposed by Litwin et al. [20] to ensure proper availability.

## 7 Summary

We presented a transaction processing scheme suitable for a distributed Wiki application on a structured overlay network. While previous work on overlay transactions has not addressed node unreliability, we identified this as a key requirement for consistency and proposed the cell model as a possible solution.

The developed transaction processing scheme provides applications with a mixture of concurrency control techniques to minimize the required communication effort. We showed core algorithms for the Wiki that utilize overlay transaction handling support and evaluated different concurrency control techniques in terms of message complexity.

## References

1. Gray, J.: The transaction concept: Virtues and limitations. In: Proceedings of the 7th Intl. Conf. on Very Large Databases. (1981) 144–154
2. Ghodsi, A.: Distributed k-Ary System: Algorithms for Distributed Hash Tables. PhD thesis, KTH Stockholm (2006)
3. Li, J., Stribling, J., Gil, T.M., Morris, R., Kaashoek, M.F.: Comparing the performance of distributed hash tables under churn. In: IPTPS '04. (February 2004)
4. Schiper, A.: Dynamic group communication. *Distributed Computing* **18**(5) (2006) 359–374
5. Schneider, F.B.: The state machine approach: A tutorial. Technical Report TR 86-800, Dept. of Comp. Sci., Cornell University (December 1986)

6. Lamport, L.: Fast paxos. Technical Report MSR-TR-2005-112, Microsoft Research (January 2006)
7. Thomasian, A.: Distributed optimistic concurrency control methods for high-performance transaction processing. *IEEE Transactions on Knowledge and Data Engineering* **10**(1) (January/February 1998) 173–189
8. Agrawal, D., Bernstein, A.J., Gupta, P., Sengupta, S.: Distributed optimistic concurrency control with reduced rollback. *Distributed Computing* (2) (1987) 45–59
9. Gruber, R.E.: Optimistic Concurrency Control for Nested Distributed Transactions. PhD thesis, Massachusetts Institute of Technology (June 1989)
10. Mohan, C., Pirahesh, H., Lorie, R.: Efficient and flexible methods for transient versioning of records to avoid locking by read-only transactions. In: *SIGMOD '92: Proceedings of the 1992 ACM SIGMOD Intl. Conf. on Management of data*, New York, NY, USA, ACM Press (1992) 124–133
11. Reed, D.P.: Naming and synchronization in a decentralized computer system, PhD thesis. Technical Report MIT-LCS-TR-205, MIT (September 1978)
12. Schütt, T., Schintke, F., Reinefeld, A.: Structured Overlay without Consistent Hashing: Empirical Results. In: *Proceedings of the Sixth Workshop on Global and Peer-to-Peer Computing (GP2PC'06)*. (May 2006)
13. Andrzejak, A., Xu, Z.: Scalable, efficient range queries for grid information services. In: *2nd IEEE Intl. Conf. on Peer-to-Peer Computing (P2P2002)*. (2002)
14. Plantikow, S.: Transactions for distributed wikis on structured overlay networks (in German). Diploma thesis, Humboldt-Universität zu Berlin (April 2007)
15. Mesaros, V., Collet, R., Glynn, K., Roy, P.V.: A transactional system for structured overlay networks. Technical Report RR2005-01, Université catholique de Louvain (UCL) (March 2005)
16. Rhea, S., Eaton, P., Geels, D., Weatherspoon, H., Zhao, B., Kubiawicz, J.: Pond: The OceanStore Prototype. In: *Proceedings of the 2nd USENIX Conf. on File and Storage Technologies*. (2003) 1–14
17. Muthitachoen, A., Gilbert, S., Morris, R.: Etna: A fault-tolerant algorithm for atomic mutable DHT data. Technical Report MIT-CSAIL-TR-2005-044 and MIT-LCS-TR-993, CSAIL, MIT (2005)
18. Moser, M., Haridi, S.: Atomic commitment in transactional DHTs. In: *First CoreGRID European Network of Excellence Symposium*. (2007)
19. Ghodsi, A., Alima, L.O., Haridi, S.: Symmetric replication for structured peer-to-peer systems. In: *The 3rd Intl. Workshop on Databases, Information Systems and peer-to-Peer Computing*. (2005)
20. Litwin, W., Schwarz, T.: LH\*RS: a high-availability scalable distributed data structure using Reed Solomon Codes. In: *SIGMOD '00: Proceedings of the 2000 ACM SIGMOD Intl. Conf. on Management of Data*, ACM Press (2000) 237–248