- [21] R. Tewari. Robustness in Replicated Databases. PhD thesis, Rutgers University, 1990.
- [22] R. Tewari and N.R. Adam. Distributed file allocation with consistency constraints. In Proceedings of the 12th IEEE International Conference on Distributed Computing Systems, June 1992.
- [23] R. H. Thomas. A majority consensus approach to concurrency control for multiple copy databases. ACM Transactions on Database Systems, 4(2):180-209, June 1979.
- [24] J. D. Ullman. Principles of Database and Knowledge-Base Systems. Computer Science Press, 1988.

- [6] D. Davcev and W. A. Burkhard. Consistency and recovery control for replicated files. In Proc. 10th ACM Symposium on Operating System Principles, pages 87-96, 1985.
- B. Gavish and H. Pirkul. Computer and database location in distributed computer systems. IEEE Transactions on Computers, C-35(7):583-590, July 1986.
- [8] D. K. Gifford. Weighted voting for replicated data. In 7th ACM SIGOPS Symposium on Operating Systems Principles, pages 150-159, December 1979.
- [9] C. L. Huang and V. O. K. Li. Regeneration-based multiversion dynamic voting scheme for replicated database systems. In 6th IEEE Conference on Data Engineering, pages 370-377, 1990.
- [10] S. Jajodia and D. Mutchler. Dynamic voting. In ACM SIGMOD Conference, pages 227-238, 1987.
- [11] S. Jajodia and D. Mutchler. Integrating static and dynamic voting protocols to enhance file availability. In 4th IEEE International Conference on Data Engineering, pages 144-153, 1988.
- [12] S. Jajodia and D. Mutchler. A pessimistic consistency control algorithm for replicated files which achieves high availability. *IEEE Transactions on Software Engineering*, 15(1):39-46, Jan 1989.
- [13] S. Jajodia and D. Mutchler. Dynamic voting algorithms for maintaining consistency of a replicated database. ACM Transactions on Database Systems, 15(2), June 1990.
- [14] D. E. Long and J. L. Carroll. Regeneration protocols for replicated objects. In 5th IEEE Conference on Data Engineering, pages 538-545, 1989.
- [15] D. E. Long and J. L. Carroll. The reliability of regeneration-based replica control protocols. In 9th IEEE Conference on Distributed Computing Systems, pages 465-473, 1989.
- [16] S. Mahmoud and J. S. Riordan. Optimal allocation of resources in distributed information networks. ACM Transactions on Database Systems, 1(1):66-78, March 1976.
- [17] H. L. Morgan and K. D. Levin. Optimal program and data locations in computer networks. Communications of ACM, 20(5):315-322, May 1977.
- [18] J. F. Paris. Voting with witnesses: A consistency scheme for replicated files. In 6th IEEE Conference on Distributed Computing Systems, pages 606-612, 1986.
- [19] C. Pu, J. D. Noe, and A. Proudfoot. Regeneration of replicated objects: A technique and its Eden implementation. IEEE Transactions on Software Engineering, 14(7):936-945, July 1988.
- [20] R. V. Renesse and A. S. Tanenbaum. Voting with ghosts. In 8th IEEE Conference on Distributed Computing Systems, pages 456-462, 1988.

ensures that the last copy is lost only when the last site of the current majority partition fails. This leads to improved availability as compared to the pure dynamic voting algorithm. An availability analysis of the RVC algorithm for the case of site failures confirms that RVC provides higher availability than previously proposed algorithms.

Future research in this area needs to address the issue of determining the vector GT(O) using as inputs the communication costs, a statistical profile of the database and availability/reliability considerations. Another issue that needs further consideration is the determination of the number of virtual copies of each data object that should be created as well as their location. Finally, an important question is: at which site should the data object be regenerated when the associated threshold GT(O) has been crossed. We are currently investigating these issues.

Acknowledgments

We acknowledge constructive comments and suggestions by anonymous referees that have helped to improve the quality of the paper. Thanks also to Dr. Bharat Bhargava of Purdue University and Dr. David Rozenshtein of Rutgers University for their careful reading of the paper. Financial support from the GSM Research Resources Committee, Rutgers University, is acknowledged by Nabil Adam and Rajiv Tewari. Summer research support from Temple University for 1990 is acknowledged by Rajiv Tewari. A preliminary version of the algorithm appears in the 11th International Conference on Distributed Computing Systems, May 1991, Arlington, TX.

References

- N.R. Adam and R. Tewari. Regeneration with virtual copies for replicated databases. In Proceedings of the 11th IEEE International Conference on Distributed Computing Systems, pages 429-436, May 1991.
- [2] P. A. Bernstein and N. Goodman. An algorithm for concurrency control and recovery in replicated distributed databases. ACM Transactions on Database Systems, 9(4):596-615, December 1984.
- [3] P. A. Bernstein, V. Hadzilacos, and N. Goodman. Concurrency Control and Recovery in Database Systems. Addison-Wesley Publishing Co, 1987.
- [4] B. Bhargava and Z. Ruan. Site recovery in replicated distributed database systems. In 6th IEEE Conference on Distributed Computing Systems, pages 621-627, May 1986.
- [5] R. G. Casey. Allocation of copies of a file in an information network. In AFIPS Conference Proceedings, pages 617-625, 1972.

Figure 8: Availability Analysis of the RVC Algorithm

Figure 7: State Transition Diagram for RVC Algorithm

with double ovals are the states where regeneration occurs because of the generation threshold of GT = 2. At each of these states one virtual copy is converted to a real copy.

For illustrative purposes Figure 7 is the state transition diagram for the case of 4 real copies and 2 virtual copies, with GT = 2. This is the most directly comparable case with previous studies. The state transition diagram is drawn in two parts. The first part of the diagram indicates state transitions when real copies fail and the second part indicates state transitions when virtual copies fail.

To follow the construction of Figure 7, it is necessary to observe that the state where all sites are connected is (4,4,0,2), which implies that four out of four sites that took part in the last update are operational; zero sites are not working; and the system has two virtual copies. Starting from the initial state (4,4,0,2), we can obtain the upper part of Figure 7 by letting the real copies fail and recover. The lower part of Figure 7 is obtained by letting the virtual copies fail and recover. The two parts taken together model the complete system. A total of 26 states are required to model all possible state transitions.

Using this state transition diagram, we write down the corresponding flow balance equations for the 26 states. This system of simultaneous equations together with the augmented equation that all probabilities sum to one, is then solved². The results of the availability analysis are plotted in Figure 8. The availability numbers for majority voting and dynamic linear algorithms have been taken from [13] for comparison purposes, after validating the results for the majority voting algorithm.

The results of the availability analysis indicate that the performance of the RVC algorithm for the case of site failures is superior to the Dynamic-Linear and Majority Voting algorithms. The trade off for obtaining this improved availability is in terms of higher costs for maintaining the virtual copies. The improved performance of the RVC algorithm results from selective regeneration using a generation threshold, which is economical over wide area networks. The improvement in availability obtained by our algorithm easily offsets the cost of maintaining the virtual copies.

8 Conclusion and Future Research

In this paper, we propose an algorithm that utilizes regeneration in distributed systems to improve the availability of replicated data. The proposed algorithm has an advantage over a pure regeneration strategy or a pure dynamic voting strategy, since it automatically maintains system availability using the generation threshold GT(O), so that manual intervention is avoided for all but total failure. Thus, it is self-adaptive to changes in system configuration by regenerating data objects at other sites when there is a danger of losing the last real copy of the data object. On the other hand when sites recover, some real copies may be converted back to virtual copies in an effort to obtain the optimal initial allocation. Our algorithm

²Maple V, a symbolic manipulation system from the University of Waterloo, was used for symbolic calculations.

so SG(H) must have an edge connecting T_i and T_k ; either T_i must precede T_k $(T_i \longrightarrow T_k)$, or T_k must precede T_i $(T_k \longrightarrow T_i)$. Therefore, SG(H) induces a write order. 2

Read Order : Let H be a RVC history. SG(H) induces a read order for H.

Suppose T_j reads x from T_i , T_k writes x ($i \neq k$ and $j \neq k$), and T_i precedes T_k ($T_i \longrightarrow T_k$). By property RVC_2 , T_j reads a read quorum of x; and by RVC_1 , T_k writes a write quorum of x. Since any two read and write quorums must intersect, there is a copy of x, say X_M , such that $r_j[x_M]$ and $w_k[x_M]$ are in H. Since these are conflicting operations, either $r_j[x_M] < w_k[x_M]$ or $w_k[x_M] < r_j[x_M]$. In the former case, $T_j \longrightarrow T_k$ is in SG(H), and we have proved an induced read order. If $w_k[x_M] < r_j[x_M]$, then we can show a contradiction in terms of the version number. Therefore, we have proved that $T_j \longrightarrow T_k$, so SG(H) induces a read order for H.

The two properties read order and write order, imply that SG(H) is an RDSG of H. Also, by condition RVC_4 , SG(H) is acyclic. Hence, H is 1SR 2

7 Availability Analysis of the RVC Algorithm

The performance of the RVC algorithm is analyzed using a Markov model for the case of site failures. Even though the RVC algorithm can handle network partitioning, it is difficult to perform a general analysis for all possible topologies of the network. Previous studies [10, 13, 15] have used the Markov analysis technique, hence we can compare our results to the results obtained by them for comparable voting based algorithms.

The assumptions inherent in a Markov analysis are:

- 1. The communication links do not fail. Only sites are subject to failures and repairs.
- 2. Failures at various sites occur according to a Poisson process with failure rate p. Repairs at various sites also occur according to a Poisson process with repair rate q.
- 3. Updates occur instantaneously, hence communication delays can be ignored.
- 4. Updates are more frequent than failures or repair. Hence, after any failure or repair, an update always arrives at a functioning site.

These assumptions are required only for the performance analysis of the algorithm, and not for the proper functioning of the algorithm itself. A state transition diagram is depicted in Figure 7, and this state transition diagram depicts the state using a vector notation, state vector = (X, Y, Z, W), where X is the number of sites out of Y that are up, Y is the update sites cardinality, Z is the number of other sites (not taking part in updates) that are up, and W is the number of sites with virtual copies. The states depicted

Every RVC history has the following properties:

 RVC_1 : If T_i writes x, then H will contain $w_i[x_{A_1}], \ldots, w_i[x_{A_n}]$ for some write quorum $Q_w(x) = \{x_{A_1}, \ldots, x_{A_n}\}$ of x's current physical copies.

This is the write rule for the dynamic voting with virtual copies algorithm, and conforms to procedure WRITE discussed before.

 RVC_2 : If T_i reads x, then H contains $r_j[x_{A_1}], \ldots, r_j[x_{A_n}]$ for some read quorum $Q_r(x) = \{x_{A_1}, \ldots, x_{A_n}\}$ of x. Then H contains $rr_j[x_{A_k}]$ for some x_{A_k} in $Q_r(x)$ where $VN(x_{A_k}) = max\{VN(x_i) \mid x_i \in Q_r(x)\}$. Here rr_j refers to the physical copy read for item j, since after the collection of the read quorum, the real read will be done from one of the current physical copies, and not from the virtual copies.

 RVC_2 says that each transaction that reads data item x, reads a read quorum of x and selects from that read quorum, a real copy with the maximum version number. Since the TM cannot determine the real read until it knows the version numbers of all copies in $Q_r(x)$, all read of those copies must precede $rr_i[x_A]$.

 RVC_3 : Every $r_j[x_A]$ follows at least one $w_i[x_A], i \neq j$.

This condition requires each copy of X (whether real or virtual) to be initialized before it can be included in a read quorum. Without loss of generality, we can assume that the initial values of the version numbers for all copies of all data items to be initialized to zero.

 RVC_4 : SG(H) is acyclic.

This condition says that the underlying concurrency control scheduler uses a correct concurrency control algorithm. Thus the RVC algorithm can work with any correct CC algorithm like the distributed 2PL, distributed timestamping, or any distributed optimistic technique. RVC is not restrictive with regard to the concurrency control algorithm, unlike the available copies algorithm that assumes strict 2PL scheduling.

Theorem 1: Every RVC history is 1SR.

Proof:

We shall prove that every RVC history H is 1SR, by proving that SG(H) is an RDSG of H. Then, since SG(H) is known to be acyclic by RVC_4 , H will be 1SR. First, we prove that H satisfies the following properties:

Write Order : Let H be a RVC history. SG(H) induces a write order for H.

We can prove this by the following reasoning. Let T_i and T_k write x. Since all write quorums of a data object intersect, there exists a real copy x_A that T_i and T_k both write. These writes on x_A conflict,

This example traces our algorithm for one complete cycle starting from a non-partitioned state and continuing through several site failures. The recovery of sites is further traced until all sites have recovered. Notice that due to the value of generation threshold chosen (GT(O) = 2), we end the cycle with two virtual copies at sites D and E, just as we had started with two virtual copies at sites D and E. This illustrates the self adaptive nature of the algorithm.

6 Correctness of the RVC Algorithm

Dynamic voting as proposed by Jajodia and Mutchler [10] has been proved correct. The read, write, failure and recover mechanisms are different for our algorithm than they are for plain dynamic voting, and for this reason we cannot assume that the proof of dynamic voting would directly apply here.

Our objective is to show that the replicated data (RD) histories produced by the RVC algorithm are onecopy serializable (1SR). We will utilize a modified serializability graph (SG) in the proof. The serializability graph is a graph that is used in proving the correctness of concurrency control algorithms [24]. We also assume that a software module called the transaction manager (TM) is running at each site that handles conflict serialization, and we further assume that the underlying concurrency control protocol used by the TM is two phase locking (2PL). These assumptions are standard practice on all major database systems [3].

The SG models the fact that two transactions, that have conflicting requests for the same data item, must be synchronized, even if they do not access the same copy of that data item. This is consistent with the one-copy view of the database held by the user. A node n_i precedes node n_j , denoted by $n_i \ll n_j$, in a directed acyclic graph if there is a path from n_i to n_j . Given a replicated data history H, a replicated data serialization graph (RDSG) is an augmented graph of SG(H) (i.e. SG(H) with possibly other edges added) such that the following conditions hold:

- 1. If T_i and T_j write data item x, then either $T_i \ll T_j$ or $T_j \ll T_i$.
- 2. If T_j reads-from T_i , T_k writes some copy of x $(k \neq i, k \neq j)$, and $T_i \ll T_k$, then $T_j \ll T_k$.

A RDSG satisfying condition (1), is said to induce a *write order* for H. If it satisfies condition (2), the RDSG is said to induce a *read order* for H. Therefore, RDSG(H) is an extension of SG(H), that induces a read order and a write order for H. With these definitions, we can state the following important property from [3]:

Bernstein et al[1987]: Let H be an RD history. If H has an acyclic RDSG, then H is 1SR.

Let us define a RVC (Regeneration with Virtual Copies) history as an RD history that models the execution of the RVC algorithm. We now state our theorem:

Since the generation threshold (=2) has been crossed, D' is upgraded to a real copy by copying information from site C. Now, if a partition occurs, C and D are both partitions containing one real copy. The tie for majority partition will be broken in favor of the lexicographically largest copy, in this case, C. Thus, site C will be the majority partition. We have thus shown that read and write operations are always performed in a majority partition. A majority partition can always be obtained down to the last copy. At any state, for any possible partitioning, the topological information about the network carried in the U_i vector enables the collection of votes to form a majority partition.

We now continue the example to illustrate site recovery and integration.

State VI (After site B recovers)

Suppose, site B recovers (and runs the function *recover*) and finds that it can communicate with C, D. Then site B determines that it is in a majority partition, and it proceeds to integrate by copying the object's information from C or D, giving:

The threshold (GT = 2) is reached, hence the lexicographically lowest copy (i.e. at D) is converted into a virtual copy.

State VII (After site A recovers)

Site A recovers (and an update arrives) and runs the function *recover*, described later, for each data object owned by site A. The state of the system after site A recovers is:

State VIII (After site E recovers)

After site E recovers, all sites can communicate with each other and all sites integrate with each other, forming a non-partitioned network. The state of the system at this point is:

State III (After site E fails and one update takes place) The state of the system after site E's failure is:

The virtual copy at site E is not included in the write quorum, and the update vectors of all the operational sites reflect the set of sites that participated in the last update to data object O. The version number of E remains 8, indicating that it has become out-of-date. It should be noted that if site E had failed, and there was no update at that point in time, the state information of all the other sites would not have been updated instantaneously merely on failure of E. The state information would be updated at the next update request to the data object O.

State IV (After site B fails and one update takes place) The state of the system after site B's failure is:

Now consider some hypothetical partitioning, say A/CD'. CD' would be the majority partition, since it has 2 out of three current copies (real and virtual). For partitioning AC/D', AC would be the majority partition. Hence for any possible partitioning of the functional sites, we can find a majority partition, thereby guaranteeing mutual exclusion. If there are ties in the voting process, the ties are resolved by using lexicographical ordering rules as pointed out in the Figure 2

State V (After Site A fails and one update takes place) The state of the system after site A's failure is:



5 An Example of the RVC Algorithm

For the purpose of illustrating the detailed operation of the RVC algorithm, consider a five site network with the value of the generation threshold for the given data object O, GT(O) assumed to be set to 2 for illustration purposes. We have left the procedure for determining the GT(O) unspecified at this point, but one technique for determining it is to set GT(O) to the number of copies of data object O to attain a minimum tolerable availability in the system. This can be obtained as a result of a file allocation algorithm that incorporates availability constraints. Details of such an algorithm can be found in [22].

State I (initial state)

The initial state of the system has version numbers equals to six and, as depicted by the update vectors, all five sites that have a copy of the data object are in communication with each other.

State II (after two updates)

This state depicts the system after two updates have been successfully propagated. The update vectors show that each site can communicate with all other sites.

Figure 4: The DO_WRITE Module

procedure DO_WRITE(0: object)
begin
if majority $(O, \text{ site}, C)$ then
for each data object $i \in S$
if i is a real copy
perform the write
update state information
else if i is a virtual copy
update state information;
$\operatorname{commit}(C)$
else
$\operatorname{abort}(C)$
end; {DO_WRITE}

Figure 5: The DO_FAILURE Module

<pre>procedure DO_FAILURE(i: site_id);</pre>
begin
for each data object $O \in i$
if (O is a real copy) and $(GT(O) < k)$ then
determine the site j ,
to regenerate data object O
upgrade a virtual copy of O to
real copy at site j
end; {DO_FAILURE}

Figure 2: The Function Majority function majority(O: object, i: site_id, var C: set_of_sites): boolean; begin $C \longleftarrow$ set of sites that communicate with this site for each $i \in C$ begin read v_i, U_i end $\mathsf{let} \; v_{\textit{max}} = \textit{max}\{v_i \; | \; i \in C\}$ for each $i \in C$ begin if $v_i = v_{max}$ then include i in set Send let $U = U_i$ for some $i \in S$ if $(\operatorname{card}(S) > 1/2 \operatorname{card}(U))$ or $(\operatorname{card}(S) = 1/2 \operatorname{card}(U) \text{ and } \max_{i \in U}(v_i) \in S)$ then majority := true else majority := false end; {majority}

associated with each data object consists of:

$$v_{O} = integer$$
 $U_{A}^{O} = (A, B, C, D, E)$

where, v_O indicates the version number of object O, and U_A^O indicates the update site vector of object O at site A. This vector contains a list of the sites that participated in the last update to object O.

I Iguie 0. The Do_RebitD Module
procedure DO_READ(0: object)
50811
if majority $(O, \text{ site}, C)$ then
select any real copy $i\in S$
read from <i>i</i>
$\operatorname{commit}(C)$
else
$\operatorname{abort}(C)$
end; {DO_READ}

Figure 3: The DO_READ Module

- Step 3: When a site detects the failure of another site, it runs the DO_FAILURE module. For each data object O belonging to the failed site¹, if the copy at the failed site was a real copy and the generation threshold GT(O) has been crossed, convert one virtual copy to a real copy. The site at which the data object is to be regenerated is determined based on communication cost minimization considerations, node utilization considerations or network reliability considerations depending on the organizational priorities. This module ensures that no data object loses its last real copy, unless there are no more sites to regenerate on.
- Step 4: When a site recovers, the DO_RECOVER module is executed by that site. This module first determines whether the site is in the current majority partition with respect to *each* data object stored at that site. For each data object O that has a current majority partition it determines if the GT(O) of that data object has been exceeded. If so, real copies of the data object are converted to virtual copies in accordance with the vector T(O, j) of initial file allocations. Otherwise, the data object O is made current and its state information is updated. The objective of this module is to try to move towards the optimal file allocation stored in the vector T(O, j) for each data object O.

An important consideration for our algorithm is that there may be a high communication cost if a real copy of a large object needs to be updated during recovery at a distant site. This is handled in two ways in our algorithm. First, regeneration is performed selectively according to predefined rules. Selective regeneration is described more fully in Figures 3 through 6. Second, if there is a need during selective regeneration to create a real copy of a large data object at a remote site, this can be done by using the *differential file* concept. Only the *diffs* between the version number at the remote site and subsequent updates need to be sent to the remote site, where local processing is performed to bring the data object up to date. The technique of sending diffs saves communication costs, specially in the case of large data objects over long distances. Hence, selective regeneration is a preventive approach to minimize incurring large communication costs, whereas differential updates minimize communication costs when an update is absolutely necessary.

Descriptive code for the RVC algorithm is provided in Figures 2 through 6. Figure 2 describes the *majority* function, which determines whether a given site i belongs to a majority partition with respect to a data object O. The other modules in the algorithm use function *majority* to initiate the collection of a majority of current sites in the present partition.

The notation used in the description is as follows. Suppose the distributed system consists of sites $(A, B, C, \ldots,)$. Each data object could be replicated at one or more of the sites. If a data object O has a real copy at site A, then we denote it by A_O . A virtual copy at site A is denoted by A'_O . The state information

¹ This information is maintained in the allocation table T(O, j). The table T(O, j) is itself fully replicated.

data objects (i.e. data objects requiring large storage). In such environments the RVC algorithm can work in conjunction with a file allocation algorithm to optimize system performance and availability.

3 Design Objectives of the RVC Algorithm

The design objectives for the algorithm are:

- 1. The algorithm should be applicable to geographically distributed architecture.
- 2. The algorithm should be able to handle Network partitioning as well as site failure.
- 3. A read quorum should require only one current real copy of a data object. The other virtual copies comprising the quorum could have state information, but no associated data.
- 4. A write quorum should require a minimum of one current real copy of a data item. The other copies could be virtual copies.
- 5. Virtual copies can be upgraded to real copies according to a prespecified protocol, whenever the number of real copies drops below a specified generation threshold.
- 6. The algorithm should be truly distributed. That is, each partition should be able to decide autonomously whether it is a majority partition based on the state information associated with the data objects in that partition.

4 Description of the RVC Algorithm

We first present the algorithm's major steps followed by a detailed description of the algorithm.

- Step 1: When a read request for a data object *O* is received at a given site, the DO_READ module is executed by that site, which then runs the function *majority* to check if it is in the current majority partition with respect to the data object *O*. If it is, the read request is satisfied by a real copy in that partition. Otherwise, the request is turned down. A quorum made up of at least one real copy will be allowed.
- Step 2: When a write request for a data object O is received at a given site, that site will execute the DO_WRITE module. If the site determines that it is in a current majority partition, the update will be propagated to each site that has a real copy and the associated state information will be updated. For those sites that have virtual copies, only the state information will be updated.

The motivation for our algorithm can be described with reference to Figure 1. We consider two cases: site failures, and network partitioning. For illustrative purposes, we consider a network with seven sites, numbered one through seven. Assume that the file allocation algorithm determines that a copy of the file should be located at sites 1, 3, 5 and 7.

Consider the case of site failures which is depicted in the first diagram of Figure 1 and suppose that sites start failing in the order of 1, 3, 5 and 7. Even though there has been no partitioning, and sites 2, 4 and 6 are operational, the availability of the file went down to zero (since not even one copy of the file is accessible). At this point any voting based algorithm will not allow operations to continue in any partition. Our proposed algorithm, which is referred to as Regeneration with Virtual Copies (RVC) and is described in section 4 avoids this problem by keeping virtual copies at sites 2, 4, and 6 and having these virtual copies participate in the voting process. Thus, as soon as sites start to fail, data object copies will be selectively regenerated at other sites. In the context of systems subject to site failures only, we observe that the last copy to fail will be a real copy, and a copy is available as long as at least one site in the network is operational.

In the case of network partitioning, depicted in the second diagram of Figure 1, we start with the initial configuration having data object copies at sites 1, 3, 5 and 7. Suppose that the network partitions into sites 1 and 3 in one partition and sites 5 and 7 in the second partition. The dynamic voting algorithm with linearly ordered copies (where ties are broken in favor of the lexicographically lowest numbered copy) will allow updates and reads in the partition containing sites 1 and 3. If another partition occurs, resulting in having sites 1 and 3 in different partitions, processing will be allowed only in the partition containing sites 1 and 4. The danger here is that if site 1 fails, dynamic voting class of algorithms will allow processing in none of the partitions. It is this shortcoming of the voting class of algorithms that we wish to overcome through our proposal.

The RVC algorithm allows processing to continue by regenerating the virtual copy kept at site 4 to a real copy. Thus, unlike currently available dynamic voting schemes, even if site 1 fails, site 4 will, under the RVC algorithm, still allow operations to continue.

It should be noted that regeneration will lead to sub-optimal allocation of files, since the new allocation obtained by regeneration will not be the same as the optimal solution determined by the file allocation algorithm. However, this sub-optimal mode of operation is preferable to having no data object copies available at all. As sites recover, and/or communication links are repaired, the system can be restored to its original "optimal" configuration.

To summarize, our work combines the advantages of the dynamic voting approach with the regeneration approach to achieve a selective regeneration policy that satisfies consistency constraints. Our algorithm would be effective in such applications as military tracking files, version management of graphics files in CAD/CAM systems, and other image tracking applications involving infrequent updates to relatively large

Figure 1: Motivation for RVC Algorithm

Initial allocation of data objects to nodes of the DCS can be performed by a file allocation algorithm such as the algorithms proposed in [5, 17, 16, 7, 21]. We have recently proposed an algorithm[22] to calculate the optimal number of copies of a data object, given communication costs between sites, storage costs and node and site reliabilities. This algorithm gives the overall system reliability or availability for a specified number of copies of a data object, hence it can be used to calculate GT(O) by specifying the minimum reliability or availability desired by the system administrator or organizational policies.

GT(O), therefore represents the optimal number of file copies to be maintained in the DCS for each data object. We maintain an allocation table T(O, j), whose row indices represent data objects and column indices represent the sites at which data object O has been initially allocated. The table T(O, j) can be a zero-one matrix, where an element T_{Oj} is set to 1 if site j has a copy of the data object O, and 0 otherwise. The allocation table T(O, j) is itself replicated at all sites, since it has to be available to determine a majority partition for each read and write operation.

Dynamic voting as proposed in [6, 10, 12] requires real copies of data objects at all sites at which the data object is replicated. We propose to integrate *Virtual Copies* with dynamic voting. Virtual copies contain only state information and participate in the voting process both for read and write quorums, whereas real copies have state information as well as associated data of the object. We also impose the constraint that each quorum must contain at least one real copy of the data object.

Regeneration in our protocol refers to converting a virtual data object copy to a real data object copy, whenever the generation threshold has been crossed. This ensures that the last copy to fail is a real copy. Regeneration will also work to convert real copies to virtual copies, when the need arises, e.g. when sites are recovering, resulting in a surfeit of real copies. This will result in a self-regulating (and self-adaptive) system that monitors the number of real and virtual copies of all data objects in the system, and seeks to maintain a predefined level of the number of real and virtual copies in the system. called regeneration, which is quite similar to file migration. Regeneration in distributed computing systems has been recently suggested in [19, 15, 14, 9, 1] as a mechanism for maintaining consistency of replicated data in distributed computing systems. Distributed databases and distributed file service mechanisms are examples of distributed computing systems.

The algorithm proposed in [19], is an example of a regeneration algorithm that implements a replicated directory system. Such a system allows the selection of arbitrary objects to be replicated, the choice of the number of replicas of each object, and the placement of copies on machines (that are assumed to have independent failure modes). The replication level is restored by automatically replacing lost copies (due to node crashes) on other active sites. The algorithm uses a read one write all strategy (ROWA); if some copies are found to be inaccessible, new replicas are created to replace them. Some limitations of this algorithm are: it is not applicable to network partitioning; it is proposed for distributed systems on a local area network, but not for distributed systems on a wide area network; and it requires additional storage overhead due to replicated directories that are used to point to current copies.

Any mutual consistency protocol that implements mutual exclusion, including voting and dynamic voting involves data transfer over the network. This data transfer is necessary since all copies that participate in an update to any data object have to be made current. If the database is distributed over a wide area network, this will involve large data transfers over long distances. This is contrary to the philosophy of distributed database design, which advocates minimizing data transfer over the network as one of its primary objective. This problem is further compounded if updates are numerous. At every update, large data transfer will be required in order to keep the copies of each data object synchronized. Hence, a regeneration approach in wide area distributed computing systems has to minimize data transfer by regenerating only when absolutely necessary. In this paper we propose a consistency control algorithm that selectively utilizes regeneration, and also accommodates the case of network partitioning. Our algorithm provides greater availability than previous voting based algorithms by incorporating selective regeneration.

The rest of paper is organized as follows. In the next section we present a motivation for the proposed algorithm, followed by the design objectives of the algorithm. A description of the algorithm in terms of the DO_READ, DO_WRITE, DO_FAILURE and DO_RECOVER modules is presented in section 4. An example in the section 5 illustrates the working of our algorithm. We present a proof of correctness of the algorithm in section 6 and a stochastic availability analysis in section 7.

2 Motivation for the Proposed Algorithm

In wide area network environments, the key to using regeneration is to use it selectively. Thus, we propose a policy of selective regeneration through a mechanism that is referred to as generation threshold (GT(O)). one, and only one partition. Such a partition is referred to as the *majority partition*. Consistency control algorithms ensure that user requests are processed in such a manner that mutual consistency of all data objects is preserved when site failures and/or network partitioning occur.

Consistency control algorithms can be classified according to whether they can handle site failures only, or both site failures and network partitioning. Two representative algorithms that can handle site failures only, are ROWAA (Read-One-Write-All-Available) [4] and Available Copies[2]. These algorithms function by reading *any* available copy of the data object (preferably local) and writing into all available copies. Under these algorithms a user query can be always satisfied as long as at least *one* copy of the desired data object is available. Updates of a data object, in this case, will always be satisfied as long as the last copy of a data object is not lost.

Voting algorithms for preserving mutual consistency of replicated data by mutual exclusion have been proposed by Thomas[23] and Gifford[8]. Voting type of algorithms can handle network partitioning in addition to site failures. This is achieved by a process of quorum collection which lets a network partition decide autonomously whether it constitutes a majority of current copies of a data object. If so, query and update operations would be allowed to proceed. The voting algorithms proposed in [23, 8] utilize a static quorum, which results in a limited database availability.

The dynamic voting algorithm proposed by Jajodia and Mutchler [10] and later refined in [11, 12], uses a version number, VN_i , of a copy, which counts the number of successful updates to the data object. The current version number is the maximum of the version numbers of all copies of a data object. A copy is current if its version number equals the current version number of the replicated data object. A partition is defined as a majority partition if it contains a majority of the current copies. Associated with each copy at a site *i* is another integer called the update site cardinality, SC_i . The update site cardinality is used in determining whether the current partition is a majority partition or not. SC_i maintains the number of copies of the data object that participated in the last update, hence each partition can autonomously decide whether they have a majority of copies with respect to the version number and the site cardinality. Based upon the values of VN_i and SC_i , certain rules are defined for reading and updating data objects. Rules for merging of data objects are also specified. A recent modification of this algorithm can be found in [13].

Other algorithms based on voting are: Voting with Witnesses[18] and Voting with Ghosts[20]. The voting with witnesses algorithm proposes data object copies called *witnesses* that can attest to the state of a data object by maintaining state information in the form of the version number. Witnesses can take part in the collection of read quorums. The voting with *ghosts* algorithm proposes data object copies that again carry only state information, but take part in write quorums. The objective of these algorithms is to improve the availability of the DCS.

A different approach to maintaining replicated databases has been adopted by proponents of a technique

Regeneration with Virtual Copies for Distributed Computing Systems

Nabil R. Adam Rutgers University Newark, NJ 07102 adam@adam.rutgers.edu Rajiv Tewari Temple University Philadelphia, PA 19122 tewari@cis.temple.edu

May 1, 1991 Revised July 1, 1994

Abstract

We consider the consistency control problem for replicated data in a distributed computing system (DCS) and propose a new algorithm to dynamically regenerate copies of data objects in response to node failures and network partitioning in the system. The DCS is assumed to have strict consistency constraints for data object copies. The new algorithm combines the advantages of voting based algorithms and regeneration mechanisms to maintain mutual consistency of replicated data objects in the case of node failures and network partitioning. Our algorithm extends the feasibility of regeneration to DCS on wide area networks, and is able to satisfy user queries as long as there is one current partition in the system. A stochastic availability analysis of our algorithm shows that it provides improved availability as compared to previously proposed dynamic voting algorithms.

1 Introduction

In a distributed computing environment, two types of failures may occur: the processor at a given site may fail (referred to as site failure), and communication between two sites may fail (referred to as communication link failure). When a site fails, processing at that site stops and the contents of the volatile storage are destroyed. Communication links may fail due to such reasons as noise in the link, or temporary link malfunction.

Link failures may result in *network partitioning*, isolating the network into two (or more) connected components, such that nodes within a given component are able to communicate with one another but not with nodes in other components. If we model the distributed computing system by a network where nodes represent sites and arcs represent links, then partitioning divides the operational sites into two or more components. Each component is referred to as a *partition*. Since these components cannot communicate, mutual consistency of replicated data can be preserved only if user requests are allowed to be processed in