

1-1-2007

Using Permuted States of Validated Simulation to Analyze Conflict Rates in Optimistic Replication

An-I Andy Wang
Florida State University

Geoffrey H. Kuenning
Harvey Mudd College

Peter Reiher
University of California - Los Angeles

Recommended Citation

An-I Andy Wang, Geoff Kuenning, and Peter Reiher. Using Permuted States of Validated Simulation to Analyze Conflict Rates in Optimistic Replication. *SCS Simulation: Transactions of the Society for Modeling and Simulation International*, 83(8):551-569, August 2007.

This Article is brought to you for free and open access by the HMC Faculty Scholarship at Scholarship @ Claremont. It has been accepted for inclusion in All HMC Faculty Publications and Research by an authorized administrator of Scholarship @ Claremont. For more information, please contact scholarship@cuc.claremont.edu.

Using Permuted States and Validated Simulation to Analyze Conflict Rates in Optimistic Replication

An-I A. Wang
Computer Science Department
Florida State University

Geoff H. Kuenning
Computer Science Department
Harvey Mudd College

Peter Reiher
Computer Science Department
University of California, Los Angeles

Keywords: permuted states, optimistic replication, conflict rates

Abstract

Optimistic replication provides high data availability in the presence of network outages. Although widely deployed, this relaxed consistency model introduces concurrent updates, whose behavior is poorly understood due to the vast state space.

This paper introduces the notion of permuted states to eliminate system states that are redundant and unreachable, which can constitute the majority of states (4069 out of 4096 for four replicas). With the aid of permuted states, we are for the first time able to construct analytical models beyond the two-replica case. By examining the analysis for 2 to 4 replicas, we can demystify the process of forming identical conflicts—the most common conflict type at high replication factors. Additionally, we have automated and optimized the generation of permuted states, which allows us to explore higher replication factors (up to 10 replicas) using hybrid techniques. It also allows us to validate our results with existing simulations based on actual replication mechanisms, which previously were analytically validated with only one pair of replicas.

Finally, we have discovered that update locality and bimodal access patterns are the primary factors contributing to the formation of identical conflicts.

1. INTRODUCTION

Optimistic replication is a tool to provide high data availability in the presence of network outages. As an example, optimistic replication allows multiple users to edit distant copies of the same data simultaneously, even without network connectivity. Data synchronization is achieved through a relaxed consistency model, which guarantees convergence and the correctness of data in the case of improper concurrent modifications, or conflicts. Common applications of optimistic replication include document sharing, banking, and reservation systems. Coda [13], Lotus Notes [6], Ficus [11], Oracle 7 [1], Bayou [15], Ingres, Microsoft Briefcase, and the Concurrent Version System are well-known research and commercial systems that use optimistic replication.

Although widely deployed, this relaxed consistency model introduces conflicts, whose behavior is not well understood. Empirical and simulation experience has shown evidence that conflicts occur infrequently at the level of aggregate statistics [7, 11, 17]. However, a theoretical result in the database literature suggests that the proliferation of conflicts will prevent optimistic replication from scaling [4]. To our knowledge, the current paper presents the first analytical modeling and understanding of conflict rates beyond two replicas, with analytical results validated by simulations built with actual optimistic mechanisms.

We have overcome several challenges to be able to characterize conflict rates analytically: (1) Defining the conflict rate in the analytical model. (2) Leveraging symmetries, permutations, and reachability to reduce the state space. (3)

Defining appropriate representations for automation and optimization of the state reduction process. (4) Studying the dominant conflict type (identical conflicts), via both analytical methods and simulation.

The following summarizes our major findings: (1) Conflicts are not directly captured by the system states. Since conflicts are detected when two replicas synchronize, a conflict occurs during a transition between two system states. Therefore, it is entirely possible for a system to be in a state with many conflicting data versions and without conflicts, as long as the system does not synchronize. (2) By exploiting the redundancy of states via permutation and removing unreachable states, we can reduce the state space by 2 to 6 orders of magnitude for as few as 6 replicas. (3) The temporal and spatial locality of updates interacts with optimistic mechanisms, resulting in significant changes in the number and types of conflicts.

In terms of methodology, our contributions include (1) a compact system-state representation that eliminates unimportant variation, (2) transition rules that can be used to automate analytical modeling at high replication factors, and (3) a hash-table heuristic for finding isomorphic states in this constrained problem domain.

2. BACKGROUND

In distributed environments, where component failure is the norm, replication provides high data availability by avoiding centralized, single-point failures. *Optimistic* replication further allows immediate access to any available replica of a data item, even during network outages. The tradeoff is permitting concurrent updates.

In many scenarios, this tradeoff is justifiable. First, for many applications, the majority of concurrent data modifications can proceed in parallel. With proper handling, the modifications can be later merged automatically or manually without data loss. Directories are an important example of this case. Independent file creations can be applied to two replicas of a directory and merged without causing problems [5]. Second, many applications (e.g., library database systems) can still provide meaningful service without immediate propagation of new updates.

Diverging data content requires a *reconciliation* process to bring replicas into synchronization, at some convenient time (e.g., when portable computers are temporarily connected to the network). Typically, reconciliation takes place between two replicas. Updates are tracked using either logging [13] or scanning [10]. Conflicts occur when different replicas of the same file are updated after the most recent reconciliation. Optimistic systems often provide extensible application-specific libraries to resolve the majority of conflicting updates automatically [9, 11]. The remaining conflicts require user intervention.

2.1 Definition of Conflicts

There are three common definitions for *conflicts*. The first is an update that conflicts with existing updates at any replica. This definition assumes oracle knowledge, which is not practical to measure in real systems.

The second definition is oriented toward the log-based reconciliation approach. At reconciliation time, both replicas replay logs of all updates since the last reconciliation between the same replica pair. Whenever two updates to different replicas of the same file are seen in the logs, a conflict is indicated.

The third definition is related to the scanning approach, in which a reconciliation-time scan detects updates and resolves conflicts. The difference from the second definition is that multiple updates are collapsed into one and will thus result in the report of only a single conflict. (In practice, most log-based systems optimize out multiple updates to save storage, which also causes conflicts to collapse. Thus, most real systems use the third definition.)

For the remainder of this paper, we will use the third definition. Without the loss of generality of our results, we will also assume bidirectional propagation of data at reconciliation time and deterministic resolution of conflicts.

2.2 Importance and Challenges of Analytical Modeling

Analytical modeling is important to understand the behavior of replication systems. In particular, a concise analytical form that can predict whether conflicts can be bounded under the worst scenario is invaluable for resource provisioning. An analytical model is also preferable for validating the correctness of simulations. Otherwise, even a trivial error in a simulation based on a single-point validation (i.e. two replicas) can go undetected and lead to very misleading conclusions for distributed systems.

Analytical modeling of the conflict rate is difficult for three reasons: First, the state space is exponential. By *state*, we mean the global system state. A state captures the relationship between any replica pair, so we can determine whether two replicas are the same or in conflict, and whether one replica has a more recent update or an older version of the data. Since conflicts are defined pairwise, two conflicting replicas might not be in conflict with a third. Therefore, each replica needs to track its update and conflict status relative to all other replicas. If each replica needs 2 states to indicate whether it is modified, a pairwise relationship needs 4 states. For R peer-to-peer replicas with $\frac{R(R-1)}{2}$ pairwise

relationships, we need $4^{\binom{R(R-1)}{2}}$ states. Even for 3 replicas, we need 64 states, which is prohibitive to track without automation.

Second, conflict resolution itself may lead to further conflicts, or *metaconflicts*. To illustrate, suppose we have many replicas of three conflicting data versions, X , Y , and Z . Pairwise synchronizations can result in versions XY , YZ , and XZ . Now, each meta-version is in conflict with three other versions. (e.g., XY is in conflict with Z , YZ , and XZ). Effectively, the original 3-way conflict has evolved into a 4-way conflict due to conflict resolutions. Therefore, the final conflict count is dependent on how data are propagated, in addition to the initial number of conflicting updates.

Third, a prior simulation study [16] suggests that a single class of metaconflicts—*identical conflicts*—accounts for the majority of conflicts at large replication factors. To illustrate,

suppose users K and L independently update separate replicas. K propagates her updates to user M ; L to N . When K and L reconcile, they detect a conflict and resolve it by generating an update to create version KL . Similarly, M and N detect a conflict and create MN . Now, when KL and MN reconcile, we have an identical conflict, since the content of KL and MN is the same. Unfortunately, the base case for identical conflicts involves 4 replicas, or 4,096 states, making it difficult to characterize their causes using existing analytical methods.

3. PERMUTED STATES

How do we visualize the problem so that the number of states is tractable, at least for the 4-replica case? Although 4 replicas may sound limiting, the analysis for the 4-replica case contains 4,092 more states than the 2-replica case, which is a giant leap in the analytical sense. As we will see, the 4-replica case also introduces critically important behaviors (such as conflict-resolution loops) that are characteristic of much larger systems, and do not appear in smaller cases.

Our solution is to transform the problem and analyze it in the domain of combinatorics. We use an event-based model in which time is measured in terms of “interesting” system events (updates and reconciliations).

Figure 1 illustrates the system states for two replicas, with λ as the probability of having an update at either replica, and μ as the probability having a pair-wise reconciliation process, respectively, as the next system event. We use a Poisson interarrival model. At each state, the outbound update probabilities sum to λ , and the outbound reconciliation probabilities sum to μ . Finally, the sum of outbound λ and μ at each state is 1.

This analysis assumes uniform update and reconciliation probabilities across all replicas, which is necessary to make the analysis tractable. However, the resulting model can be used to cross-validate simulations at moderate replication factors. (Previously, simulations are analytically validated for only two replicas.) The validated simulations can then be reconfigured to account for nonuniform access patterns to explore higher replication factors.

Each replica is represented by a dot. In the starting state (shaded), two replicas are identical, represented via a horizontal line connecting the two. If reconciliation occurs, the replicas remain identical, so the starting state transitions back to itself.

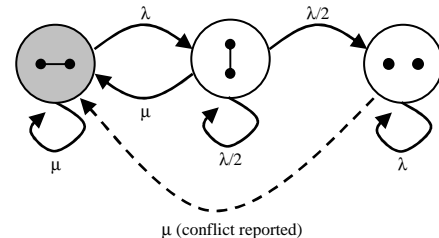


Figure 1: The state-transition diagram for two replicas. The shaded circle is the starting state, where replicas are identical. In the rightmost state is the conflict state, where replicas are in conflict. The solid lines mark various transitions among states, and the dashed line marks a conflict-reporting transition. Note that a conflict is not reported until the system transitions back to the starting state as the result of reconciling between two conflicting replicas.

If one of the replicas is updated, we move to the middle state, where the update-receiving replica *dominates* the *subordinate* one. This relationship is represented by a non-horizontal line, where the upper replica dominates the lower one. Note that regardless of which is updated, we are guaranteed to transition from the starting state to the middle state. By decoupling the state of the system from the labeling of individual replicas, each state effectively captures all *isomorphic* system states resulted from permuting the replica identifications. We refer to this type of state representation as *permuted states*.

In the case of reconciliation between a dominating replica and its subordinate, the content of the dominating replica will replace that of the subordinate, and then both replicas will be marked as identical (transition back to the starting state). An update to the dominating replica will not change its dominance over the subordinate replica. However, an update to the subordinate replica breaks its subordinate relationship to its dominating replica, and the system enters the rightmost state (conflict).

Conflicting replicas (dots) are not connected by lines. An update to either of the conflicting replicas will leave both in conflict. However, a reconciliation between two conflicting replicas will lead to identical replicas (the starting state or the convergence state), with a reported conflict.

Note that a system can be in a state with conflicting replicas without reporting conflicts, since conflicts are detected only at reconciliation time. Therefore, the *conflict rate* used in this paper, or the probability of having conflicts due to a system event (either update or reconciliation), is computed by obtaining the equilibrium probability of a state that contains replicas in conflict, multiplied by the probability of traversing its conflict-resolving transition.

3.1 Analysis for Two Replicas—the Base Case

With the state diagram in Figure 1, we can assign probability variables p_0 (leftmost) to p_2 (rightmost) to each state. When the system is in equilibrium, the outbound transition flow at each state should be equal to the inbound flow, resulting in a system of linear equations ((1), (2), and (3)). Also, the sum of probability at each state should be 1 (4).

$$\lambda p_0 = \mu p_1 + \mu p_2 \quad (1)$$

$$\left(\frac{\lambda}{2} + \mu\right) p_1 = \lambda p_0 \quad (2)$$

$$\mu p_2 = \frac{\lambda}{2} p_1 \quad (3)$$

$$p_0 + p_1 + p_2 = 1 \quad (4)$$

$$p_0 = \frac{\lambda\mu + 2\mu^2}{(\lambda + 2\mu)(\lambda + \mu)} \quad (5)$$

$$p_1 = \frac{2\lambda\mu}{(\lambda + 2\mu)(\lambda + \mu)} \quad (6)$$

$$p_2 = \frac{\lambda^2}{(\lambda + 2\mu)(\lambda + \mu)} \quad (7)$$

$$p_{\text{conflict}} = \mu p_2 = \frac{\lambda^2 \mu}{(\lambda + 2\mu)(\lambda + \mu)} \quad (8)$$

Solving for the probability of the conflicting state (7), the probability of reporting a conflict p_{conflict} (8) can be computed by multiplying (7) by μ , the probability of taking the transition that resolves conflicting replicas.

As expected, the conflict rate has an intimate dependency on both the update arrival rate and reconciliation rate, even for this simple two-replica case. Figure 2 shows the percentage contribution of each state as a function of λ/μ , superimposed with the conflict-rate curve. (We will postpone discussing the simulation validation.) As λ/μ approaches 0, the probability of convergence (p_0) approaches 1, and the probability of being in the conflict state (p_2) approaches 0. As λ/μ increases asymptotically, the probability of convergence (p_0) approaches 0, and the probability of being in the conflict state (p_2) approaches 1.

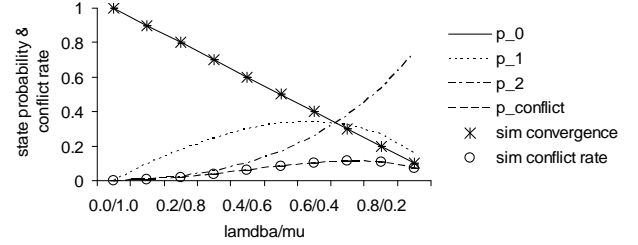


Figure 2: Percentage contributions of states for two replicas, superimposed with analytical conflict rate curves and simulation validation data points. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events. Since the confidence intervals are less than 1% of the mean and clutter the simulation points, we have removed them for the clarity of presentation.

For the conflict rate (8), as μ approaches 1 (λ/μ approaches 0), the high frequency of reconciliation will bring the conflict rate (p_{conflict}) to 0. Intriguingly, as μ approaches 0, the lack of opportunities to report conflicts with reconciliation processes will also bring the conflict rate down to 0. This finding is consistent with prior findings [18]. Thus, a system can spend most of its time updating two conflicting replicas, but only one conflict is reported per reconciliation process. Also, through this exhaustive range of ratios between the update and reconciliation rates, we can see that the conflict rate can actually be bounded under this system setting as a fraction of the total update and reconciliation events (11%) by solving equation (8) for its maximum. This is an important insight for capacity planning under optimistic replication.

3.2 Validation for Two Replicas

To validate the analytical results based on the use of permuted states, we compare them with results obtained from a version-vector-based simulation of optimistic replication, similar to the one used in a prior study [16]. Briefly, each replica keeps a local “version vector” of update counters for all replicas. A replica increments its local counter (a version vector element) whenever it performs an update. At reconciliation time, two replicas compare their version vectors. If every counter of replica X is greater than or equal to the corresponding counter of the replica Y , X dominates Y . If X dominates Y , and if Y dominates X , X and Y are equal. If neither X dominates Y , nor Y dominates X , we have a conflict. When reconciling, a subordinate replica copies the version vector from the dominant one. To merge conflicting version vectors, each counter is set to the greater of the corresponding version vector elements. The counter of the conflict-resolving replica is incremented by one, indicating that a new version is generated as a

result of resolving conflicts. The simulation includes only one replicated item. We follow the methodology presented in [16].

All simulation results are presented at the 90% confidence level. We assume that all updates and reconciliations take place instantly.

Figure 2 also shows the validation results based on simulation. The results match well with the model based on permuted states. For the first time, simulation results for optimistic replication have been cross-validated with an analytical model, although with only two replicas.

3.3 Generalization to Three Replicas

The three-replica case tests whether we can reduce the original 64 states down to a more manageable number, and whether the notation and manipulation rules invented for the two-replica case can be generalized.

3.3.1 State-Transition Diagram for Three Replicas

Figure 3 shows the transition diagram for three replicas.

Surprisingly, the three-replica case can be completely captured with only 8 permuted states. The starting state (state 0), which is also the convergence state, is easily generalized from the two-replica case. However, since an update can be applied to any one of the replicas, the update-receiving replica dominates the remaining two identical ones (state 1). At this point, an update to one of the two identical subordinate replicas will first break off the update-receiving replica from the replica that dominates it, and second will make the update-receiving replica dominate its original identical partner (state 2).

At state 2, it is interesting to note that the two dominating replicas are in conflict, while dominating the same replica. A reconciliation between the two dominating replicas will lead to the report of a conflict and transition to state 3. An update to the subordinate replica will break the update-receiving replica from all its dominating replicas and reach state 4.

At state 3, reconciling between any dominating version and the subordinate version will lead to the convergence state (state 0). An update to any dominating version will lead to state 6. An update to the subordinate replica will break off the update-receiving replica from its dominating replicas and reach state 5.

At state 4, an update to any replica will leave all three replicas in conflict. Reconciling any pair of replicas will lead to the report of a conflict and a transition to state 5.

At state 5, an update to one of the identical replicas will lead to state 7; reconciling between one of the identical replicas and the replica in conflict will lead to the report of a conflict and the generation of a new version that dominates the replica not involved in reconciliation (state 3).

At state 6, the dominance relationship is transitive. An update to the top dominating replica results in a self-transition. An update to the middle dominating replica breaks its relationship with the top dominating replica and leads to state 2, while preserving its dominating relationship to the subordinate replica. An update to the subordinate replica will break its relationship to both dominating replicas and lead to state 7. Reconciling the top dominating replica with either other replica will result in state 3. Reconciling between the bottom two replicas will result in state 1.

At state 7, an update to the subordinate replica will lead to state 4. Reconciling between the dominating and subordinate replicas will result in state 5. Reconciling the dominating replica with the conflicting replica will result in state 3. Reconciling the subordinate replica with the conflicting replica is equivalent of making an update to the subordinate replica (state 4) and then reconciling it with the conflicting replica. The result is state 5.

Overall, the three-replica case demonstrates the richness of behaviors in optimistic replication. Interestingly, not all conceivable states are possible. For example, it is not possible to have one replica dominating over two replicas in conflict. Using

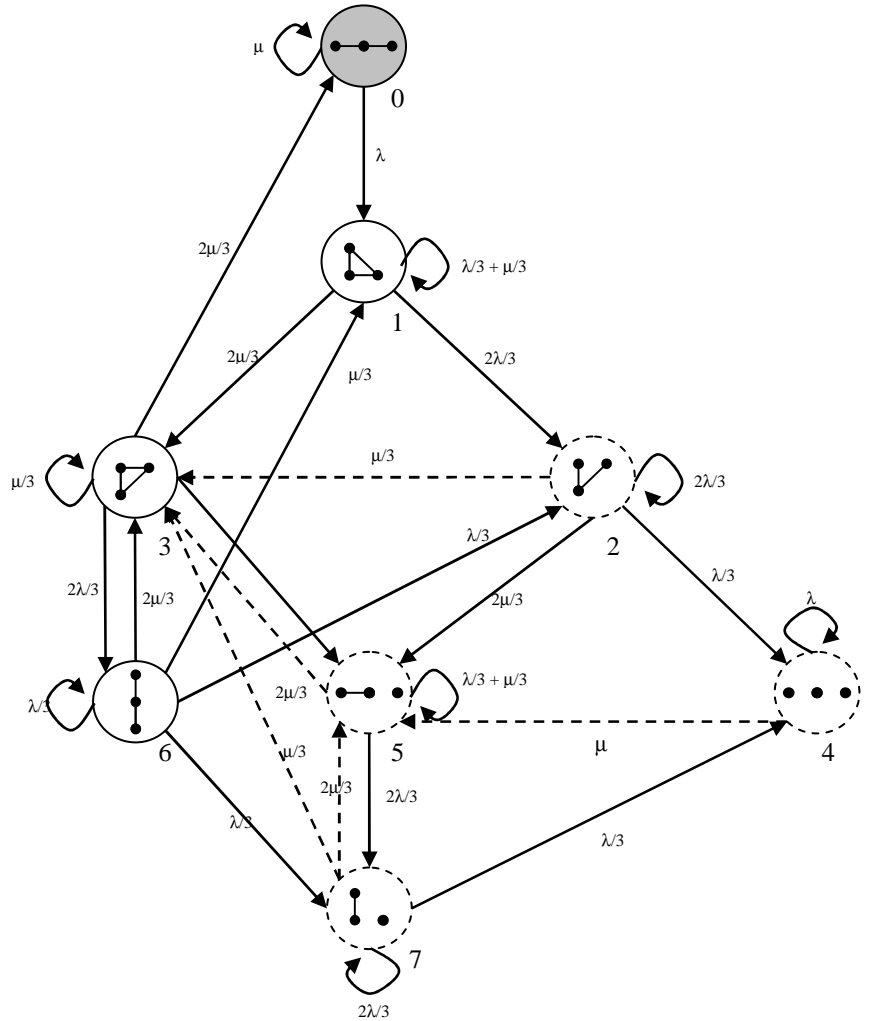


Figure 3: The state-transition diagram for three replicas. The shaded state is the starting state, where replicas are identical. States with dashed outlines contain conflicting replicas. The solid lines mark various transitions among states, and the dashed lines mark conflict reporting transitions.

permuted states for analysis eliminates both isomorphic states and unreachable states.

3.3.2 Analyses and Validation for Three Replicas

Similar to the two-replica analysis, we set up a system of equations based on the state-transition diagram in Figure 3. The exact equations used are listed in Appendix A. Note that the conflict rate is the sum of the products of the probability of each conflict-originating state and its outbound transition probability of conflict-resolving edges (9).

$$P_{\text{conflict}} = \frac{\mu}{3} p_2 + \mu p_4 + \frac{2\mu}{3} p_5 + \frac{2\mu}{3} p_7 \quad (9)$$

After solving the system of equations with MathCAD, the following equations highlight our findings:

$$p_0 = \frac{4\mu^3}{(3\lambda + 2\mu)(\lambda + 2\mu)(\lambda + \mu)} \quad (10)$$

$$p_4 = \frac{2\lambda^3}{(2\lambda + 3\mu)(\lambda + 2\mu)(\lambda + \mu)} \quad (11)$$

$$p_{\text{conflict}} = \frac{2\lambda^2\mu(3\lambda^2 + 11\lambda\mu + 9\mu^2)}{(2\lambda + 3\mu)(3\lambda + 2\mu)(\lambda + 2\mu)(\lambda + \mu)} \quad (12)$$

One immediate surprise from these resulting equations is that the order of complexity is smaller than expected. For nine equations and eight unknowns, we would expect the resulting equations to have exponents of seven to eight; on the other hand, we have four to five, suggesting that our permuted state representation may be further compacted. For example, a state with a single inbound transition can be merged with the state that makes the inbound transition. Based on the equations listed in Appendix A, states p_0 and p_6 can be directly substituted with state p_3 , and the overall system can be characterized with six states.

Figure 4 shows the probability contributions of each state. Similar to the two-replica case, as λ/μ approaches 0, the probability of convergence (p_0) approaches 1. As λ/μ increases asymptotically, the probability of being in the fully divergent state (p_4) approaches 1, where all replicas are conflicting versions. Also, in the case of a 0.9/0.1 λ/μ ratio, the most probable states (p_4 , p_5 , and p_7) contain conflicting versions.

Figure 5 shows the percentage contribution for the all-identical and all-conflict states and the conflict-rate curve. The remaining states have been removed for clarity. Compared to the two-replica case, the additional replica has increased the conflict-rate potential from 11% to 17%. Also, the peak of the curve shifts left, meaning that fewer updates are needed to cause more conflicts, with a fixed number of reconciliation processes.

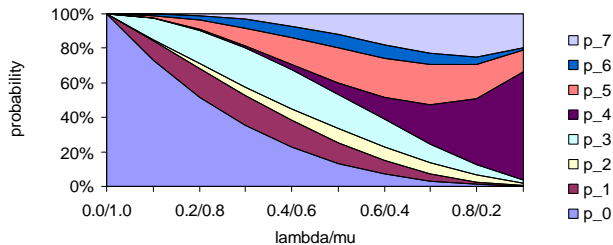


Figure 4: Percentage contribution of states in the three-replica case.

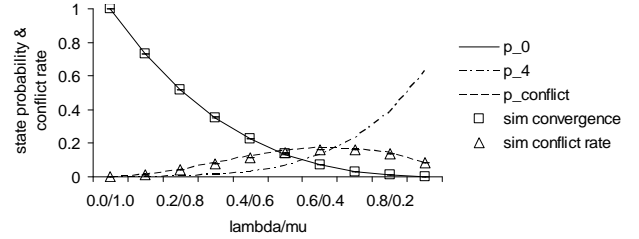


Figure 5: Percentage contributions of states for three replicas, superimposed with analytical conflict-rate curves and simulation validation data points. The simulation was repeated 5 times with different random seeds, each with 100,000 update and reconciliation events (i.e., 100K events in total). Since the confidence intervals are less than 1% of the mean and clutter the simulation points, we have removed them for the clarity of presentation.

In terms of the simulation validation, the data points once again confirm the validity of modeling optimistic replication via permuted states.

3.4 The Base Case for Identical Conflicts

The four-replica case is the base case for identical conflicts. This exploration helps us to gain insights on how identical conflicts are formed and how to predict their proliferation. Rather than burdening the reader with the same analyses as for two and three replicas, this section will highlight only the results of interest.

Surprisingly, the four-replica case contains only 27 permuted states, which is valuable for analyzing systems like Oceanstore [8], where the core writable replicas have only a replication factor of four. (We omit the state equations, since they are lengthy and not very informative.) Five of the states have only a single inbound transition, so the 27 states can potentially be compacted down to 22. Unfortunately, MathCAD did not find a closed form for these equations, so we used the Microsoft Excel solver to generate Figures 6 and 7, with 10,000 iterations to approximate each data point. The simulation results match well with the curves derived from analytical equations. (Note that the peak conflict rate has increased and shifted slightly to the left; this effect is even more pronounced in Figure 15, later in the paper.)

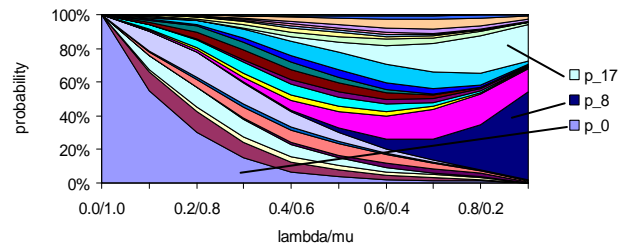


Figure 6: Percentage contribution of states in the four-replica case.

As in the three-replica case (Figure 4), a few states have dominating contributions when the system operates with extreme λ/μ ratios, namely those close to the convergence state p_0 (states with many lines interconnecting dots that represent replicas), and those close to the divergent state p_8 (states with few lines

interconnecting the dots that represent replicas). This finding prompts the question of whether it is possible to trim the state space for any number of replicas down to a characteristic subset of states. We will leave the answer as future work.

With the aid of permuted states, for the first time, we can understand and enumerate the cases where identical conflicts are generated. The formation of an identical conflict goes through a setup sequence, as shown in Figure 8. Incidentally, this sequence is in the example we used in Section 2.2. During this setup phase, the system first enters the state with two pairs of identical replicas, replicas 1 and 2, and replicas 3 and 4. A reconciliation process between one replica from each pair (e.g., replicas 1 and 4) will result in a new pair of identical replicas, dominating the replicas 2 and 3, which remain in conflict with each other. (A unique version number is generated for each conflict-resolved version.) When replicas 2 and 3 reconcile, they will create a new version of data during the conflict reconciliation and break away from the dominating replicas 1 and 4.

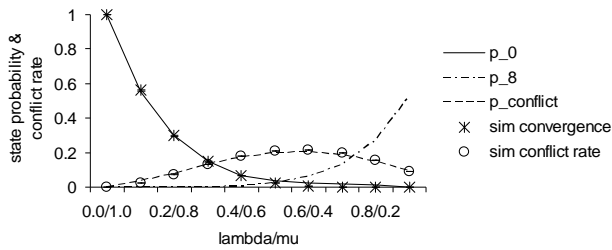


Figure 7: Percentage contributions of states for four replicas, superimposed with analytical conflict rate curves and simulation validation data points. The analytical numbers were approximated by the Microsoft Excel solver with 10,000 iterations per data point. p_0 is the convergence state, and p_8 is the fully divergent state. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events. Since the confidence intervals are less than 1% of the mean and clutter the simulation points, we have removed them for the clarity of presentation.



| | Replica 1 | | Replica 2 | | Replica 3 | | Replica 4 | |
|-------------|-----------|-------------|-----------|-------------|-----------|-------------|-----------|-------------|
| Event | Con-tent | Ver | Con-tent | Ver | Con-tent | Ver | Con-tent | Ver |
| | X | 1000 | X | 1000 | Y | 0010 | Y | 0010 |
| Recon(1, 4) | XY | 2010 | X | 1000 | Y | 0010 | XY | 2010 |
| Recon(2, 3) | XY | 2010 | XY | 1110 | XY | 1110 | XY | 2010 |

Figure 8: A setup sequence of events for identical conflicts in the four-replica case. The states being reconciled are in boldface.

Figure 9 shows the states involved in forming identical conflicts. For simplicity, not all outbound transitions are shown. Basically, after the setup sequence, as long as one replica from each version pair remains, they can reconcile and form identical conflicts. The other two replicas can be in a number of states resulting from updates and reconciliations.

Take the first right-looping branch from Figure 9 as an illustrative example. Figure 10 starts with two sets of identical

pairs, continued from Figure 8. Based on version vectors, replicas 1 and 4 are identical; replicas 2 and 3 are identical. If replicas 1 and 2 reconcile, we will have an identical conflict based on their version information; however, the content is the same. Replicas 3 and 4 can reconcile and form another identical conflict.

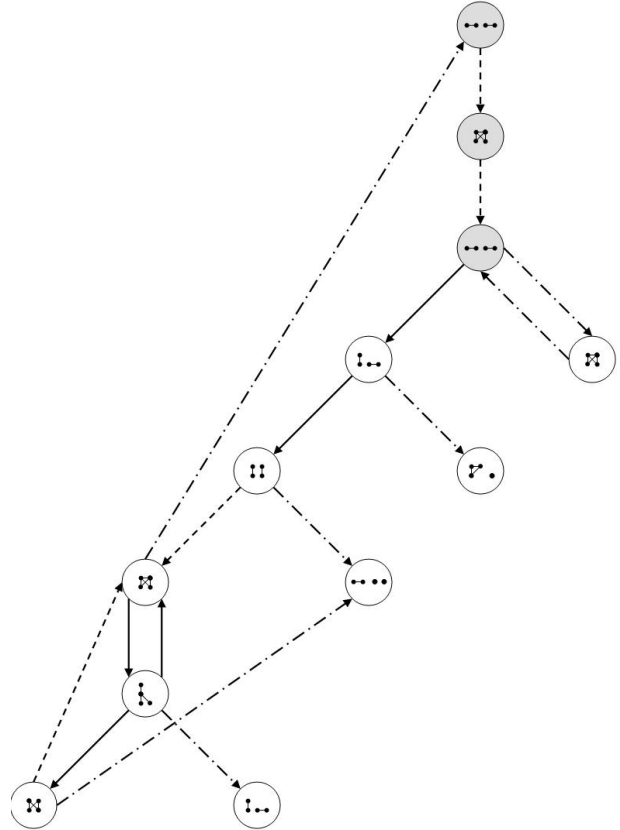


Figure 9: A subset of four-replica states and transitions that illustrates ways to form identical conflicts. The shaded states are the setup sequence. The solid lines mark various transitions among states; the dashed lines mark conflict transitions, and the dotted dashed lines mark identical-conflict transitions.



| | Replica 1 | | Replica 2 | | Replica 3 | | Replica 4 | |
|-------------|-----------|-------------|-----------|-------------|-----------|-------------|-----------|-------------|
| Event | Con-tent | Ver | Con-tent | Ver | Con-tent | Ver | Con-tent | Ver |
| | XY | 2010 | XY | 1110 | XY | 1110 | XY | 2010 |
| Recon(1, 2) | XY | 3110 | XY | 3110 | XY | 2010 | XY | 1100 |
| Recon(3, 4) | XY | 3110 | XY | 3110 | XY | 2120 | XY | 2120 |

Figure 10: The first right looping branch from Figure 9.

3.5 Comments on Identical Conflicts

During the process of analyzing the base case for identical conflicts, we made two disturbing observations. (1) For the four-replica scenario, we can see that identical conflicts can potentially be self-inducing due to looping behavior. Based on a prior simulation study on large-scale optimistic replication systems [16],

identical conflicts constitute most of the conflicts, and this looping behavior may be a contributing factor. (2) Identical conflicts are defined as a function of data content, in addition to the system states. In Figure 9, the same states can generate either regular or identical conflicts. However, tracking the data content and system states is analytically prohibitive even for few replicas, since we can no longer compact states effectively.

4. AUTOMATIC GENERATION OF PERMUTED STATES

Although permuted states enable a significant step forward in understanding replicated systems, there is still a limit to the size of the analysis. In exploring higher replication factors, we thus turn from pure analysis to an analytically validated simulation. However, such validation requires an automated way to generate analytical solutions for higher replication factors. Although the number of states still grows rapidly, the ability to validate a simulation up to 10 replicas can cover common replication deployment scenarios and give confidence in the accuracy of even larger simulations. Even without closed-form solutions to equations, using random traversals of the states with appropriate probabilities for updates and reconciliations can be a good sanity check for simulation results.

Four steps are involved in automating state generation: (1) Map the diagram state representation into a data structure; (2) define rules to transition among states; (3) remove isomorphic states resulting from permuted labeling of replicas; and (4) traverse states with the given transition probabilities.

4.1 Data Structure Representation

The first simplification in our automation is an assumption of global knowledge, which allows us to directly translate diagram states into graph-based representations. To illustrate, identical replicas 1 and 2 have the state $(=,=)$. The first row belongs to

replica 1. The equality symbols indicate that replica 1 is identical to itself and replica 2. The second row belongs to replica 2, showing that replica 2 is identical to replica 1 and itself.


If replicas 1 and 2 are in conflict, we have state $(*,=)$. The

asterisks show that replica 1 is in conflict with replica 2 (vector 1), and replica 2 is in conflict with replica 1 (vector 2).

If replica 1 dominates replica 2, we have the state $(=, >)$. The

greater-than sign shows that replica 1 dominates replica 2 (vector 1), and the less-than sign shows that replica 2 is subordinate to replica 1 (vector 2).

In this simple 2-replica scenario, vectors 1 and 2 appear to contain redundant information. However, as the number of replicas increases, we need each replica to track its relationship to others to capture the full complexity of system states. For example, with

, replicas 1 (top left), 2 (top right), and 3 (the bottom) can be represented with the state $(=, *, >)$, $(*, =, >)$, $(<, <, =)$


4.2 State Transition Rules

The state transition rules are also translated from the state diagram. An update to a dominating replica does not change the dominance of the replica. An update to a subordinate replica breaks its relationship with its dominating replicas. An update to a number of identical replicas makes one of the replicas dominating over all other replicas.

Reconciliation rules fall into three categories: (1) Reconciling two identical replicas: No actions are needed. (2) Reconciling a dominant and a subordinate replica: The subordinate replica first receives an update from the dominating one (with update rules applied, meaning that if a subordinate is dominated by two replicas, it has to break off from both dominating replicas first) and copies over the vector from the dominating replica. (3) Reconciling conflicting replicas: Each replica first receives an update to form a new data version (with similar update rules applied). For each vector element, if one of the replicas dominates a third replica not involved in reconciliation, both replicas are set to dominate the third. If one replica is identical to a third, both reconciling replicas are set to be identical to each other, dominating the third. Note that since conflict resolution is equivalent to first applying chosen updates to the conflicting replicas so as to make them equal, and then reconciling them, they cannot be subordinate to any other replica after reconciliation.

4.3 Isomorphic State Reduction

As we analyze the system at the level of permuted states, the labeling of replicas becomes irrelevant. For example, in the two-replica case, we make no distinctions between replica 1 dominating over 2 $(=, >)$, and replica 2 dominating over 1 $(=, <)$. In the

three-replica case of , we make no distinctions among $(=, *, >)$, $(=, >, *)$, and $(=, <, <)$, $(*, =, >)$, $(<, =, <)$, and $(>, =, *)$, $(<, <, =)$, $(*, >, =)$, $(>, *, =)$

This problem of compacting isomorphic states into a permuted state is similar to the problem of finding isomorphic graphs. Unlike subgraph isomorphism, the isomorphic graph problem is neither NP-complete, nor a P-problem [14]. However, known algorithms for constant-bounded vertex degree are $O(n^4)$ [2], which is still too computationally intensive to scale well.

An intriguing observation is that one state can be turned into another by swapping two corresponding rows and columns, with the invariant that the diagonal entries are always '='. For example, swapping the first and third rows of $(=, *, >)$ leads to $(<, <, =)$; the

first and third columns are then swapped to obtain a isomorphic state, $(=, <, <)$. This raises the possibility that we might find a

reduced representation of the state that captures these variations due to permutations. If so, we can look up the isomorphic states in a dictionary. We currently use a hash table for this purpose.

To construct the reduced representation, we need to consider the following constraints: (1) We need a function that is commutative for row elements and column elements, so that any

swapping between two rows and columns results in the same value. (2) We need to break the diagonal symmetry of matrices by using two different functions for rows and columns, or our reduction will not be as effective in eliminating unintended collisions (false isomorphisms). False compaction of states will confound the probability of reaching distinct states and increase the error of our analyses. (3) We need to account for all matrix elements.

The ‘=’, ‘>’, ‘<’, and ‘*’ symbols are first mapped to four 32-bit numbers (chosen randomly at design time). The current hash function is the sum of values from three functions (integer overflow is ignored): (1) the product of the sums of each row element, (2) the sum of the products of each column element, and (3) the sum of all elements. The observed false compaction rate is around 0.6% for 6 replicas, found by comparing the number of obtained states with that obtained by brute-force permutation of matrix rows and columns. The false compaction rate is expected to become lower as the number of replicas increases, because two valid matrices first have to follow the transitivity constraints (if replica A dominates B, and if replica B dominates C, then replica A dominates C) and be hashed to the same value.

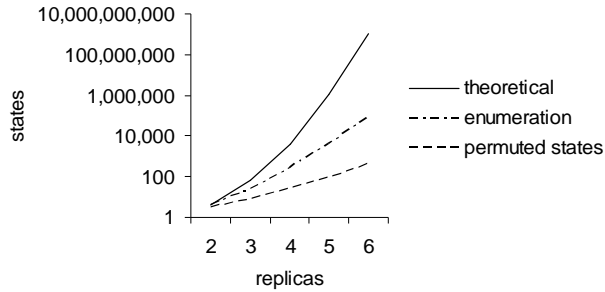


Figure 11: Comparison of number of system states obtained by theoretical limits (4 states for each replica pair), brute-force enumeration of reachable states, and permuted states.

Figure 11 compares the effectiveness of the permuted-state approach to the number of states obtained by theoretical limits (4 states for each replica pair) and brute-force enumeration of reachable states. Note that the y axis uses a logarithmic scale. Using permuted states for as few as 6 replicas can reduce the state space by 2 to 6 orders of magnitude.

4.4 Validation

After building the permuted states, all we need to do is construct the state diagram through systematic state enumeration, and randomly traverse the graph with specified update and reconciliation arrival rates. For $\lambda = \mu$, we validated our automated analytical method (which contains 488,013 permuted states) against a simulation with up to 10 replicas (Figure 12). Intuitively, λ should be much greater than μ in real systems. However, based on trace analyses [16], thanks to write-back caching and the work cycle (2-day weekends and 8-hour working days), the average λ and μ are not that far apart.

In the past, simulations of optimistic replication have been validated only for 2 replicas. Now we have two independent implementations of optimistic models that cross-validate well even at 10 replicas. We are thus more confident in using simulation to explore optimistic systems at higher replication factors.

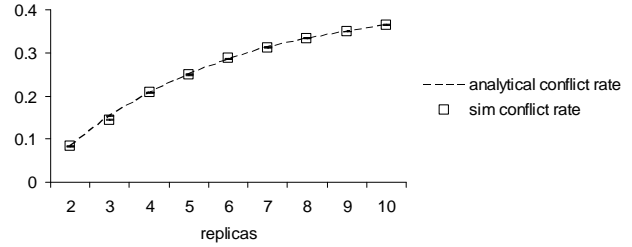


Figure 12: Comparison between automated analytical modeling based on permuted states and simulation with counter-based version vectors. The analytical numbers were obtained with 100,000 random transitions with $\lambda = \mu$. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events. Since the confidence intervals are less than 1% of the mean and clutter the simulation points, we have removed them for the clarity of presentation.

5. IDENTICAL CONFLICTS REVISITED

A prior trace-based simulation study has shown that identical conflicts account for a significant fraction of conflicts at high replication factors (around 50 replicas) [16]. However, the simulation validated by our analytical model shows that identical conflicts are relatively rare events compared to the overall conflicts (Figure 13).

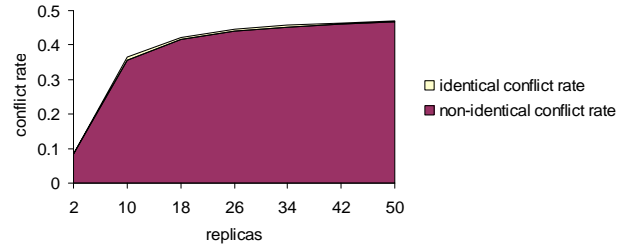


Figure 13: A stack graph of non-identical and identical conflict rates at different replication factors. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events, with $\lambda = \mu$. The identical-conflict curve is just above the non-identical one. The difference is too small to be seen.

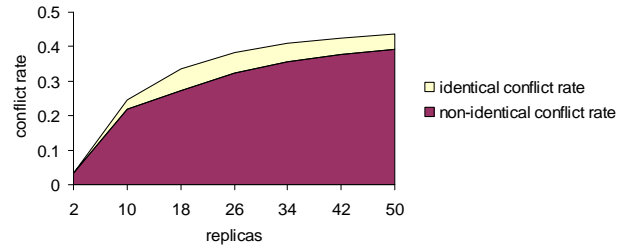


Figure 14: A stack graph of non-identical and identical conflict rates at different replication factors, with 90% of the updates going to 10% of the replicas. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events, with $\lambda = \mu$.

The prior study also suggested that identical conflicts are caused by access locality—most updates are applied to a subset of replicas. To test this possibility, we adjusted our simulation to have 90% of the updates take place at only 10% of the replicas.

Figure 14 shows that the impact of access locality is clearly visible, but the effect is within 10% of the total. To make sure that we were within the plausible range of parameter settings, we examined conflict and identical conflict rates as a function of λ/μ , with 90% of the updates going to 10% of the replicas.

Figure 15 shows the decomposition of identical and non-identical conflicts with a wide range of λ/μ ratios. The non-identical conflict rate is low for both high and low λ/μ ratios, due to either the lack of updates to create diverging versions, or the lack of reconciliation to detect conflicts. When $\lambda = \mu$, the identical conflict rate is expected to be low. However, as the λ/μ ratio decreases, the identical conflict rate accounts for most conflicts.

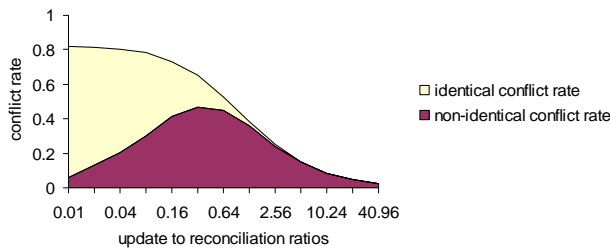


Figure 15: A stack graph of non-identical and identical conflict rates for 50 replicas and various λ/μ ratios (in log scale), with 90% of updates going to 10% of all replicas. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events.

From the viewpoint of traces, updates occur more frequently than reconciliations on average, to amortize the cost of reconciliation over time. Therefore, we hypothesize that a high average λ/μ ratio with accompanying high identical conflict rate is induced by a bimodal traffic pattern, where updates arrive in bursts to maintain an upper range of λ/μ , while the system is running a lengthy background reconciliation most of the time. This temporal locality of updates is also sensible from the viewpoint of a single replicated file, which is likely to be updated intensively over short durations, with periods of no updates but many reconciliations.

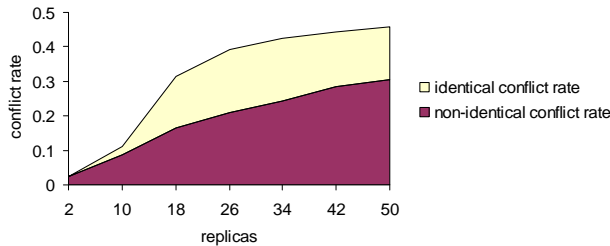


Figure 16: A stack graph of non-identical and identical conflict rates under different replication factors, with 90% of updates to 10% of all replicas and with a bimodal traffic pattern where the λ/μ ratio is 8 for one-third of the time to reflect 8-hour working days, and 0.08 to reflect non-working hours. The simulation was repeated 5 times with different random seeds, each with a total of 100,000 update and reconciliation events.

We model this behavior of update bursts with an irregular square wave, to reflect weekly activities. The function parameters are extracted from the same trace used in [16]. The function consists of five eight-hour working sessions, with $\lambda/\mu = 8$, each followed by a 16-hour “off” session dominated by reconciliation, with $\lambda/\mu = 0.08$. The five eight-hour days are followed by two 24-hour reconciliation periods. A working hour has an average of 3 updates, and the remaining hours have an average of one reconciliation per hour. The overall average λ/μ is about 0.86, computed based on the aggregate number of updates and reconciliations. Access locality still applies.

Figure 16 shows a much more drastic decomposition of conflict rate compared to Figure 14. Clearly, bimodal access patterns have a large impact on identical conflicts. Intriguingly, based on Figure 15, an average λ/μ of 0.86 should produce relatively few identical conflicts. However, beyond 10 replicas, identical conflicts can account for up to 48% of the total. This suggests that load generators based on aggregate mean arrival rates of updates and reconciliations are not suitable for studying optimistic replication. Also, given that our square function is a crude approximation of a trace, a direct trace-based simulation would be expected to have more identical conflicts, which is consistent with prior findings [16].

Intriguingly, the introduction of access and temporal localities to the workload affects the overall conflict statistics (beyond 10 replicas) very little, but the internal composition changes dramatically.

6. RELATED WORK

The permuted-state approach suggests that it might be profitable to revisit existing approaches to evaluating optimistic replication. One early approach, by Golding [3], was to measure the mean time to converge with R conflicting replicas. No updates were allowed during the convergence process.

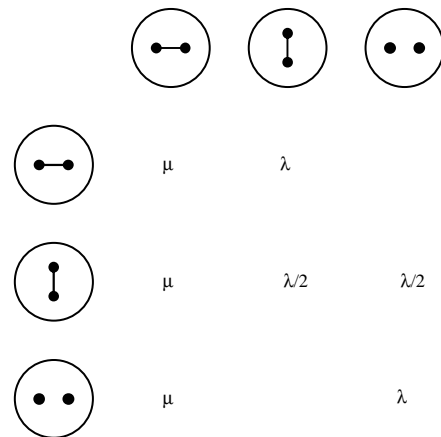


Figure 17: The state-transition table for two replicas. The rows represent the *from* states, and the columns the *to* states.

Golding’s approach is equivalent to making λ/μ be zero. Based on Figures 4 and 6, for 3 and 4 replicas optimistic replication under this setting is dominated by only a few states. Many states can be overlooked because they are only reachable through a mixture of update and reconciliation events. On the other hand, we have also observed that an optimistic replication

system typically operates with alternating extremes of λ/μ . Therefore, Golding’s approach does not reflect the case where reconciliation events dominate the system. With only the state transitions, one might conclude that Golding’s approach misses 33% of system states for 2 replicas (Figure 17), since it starts with the fully divergent states where all replicas are in conflict and transitions directly back to the starting state. The intermediate state is not exercised at all. For three replicas, Golding’s approach misses 38% of states; 56% for four replicas. In practice, however, many missed states are not heavily exercised, since alternating between extreme λ/μ ratios is the norm.

Gray et al. [4] studied replication under a database workload, with relatively uniform access patterns to all replicated items. Their results suggest that the conflict rate grows at a rate that is prohibitive for scaling of optimistic replication. However, Gray’s analytical model assumes an access pattern that is not applicable in environments where update locality is the norm. Also, due to the strong correlation between the usage model and the working day, our traffic pattern is bimodal. Our model cycles through extreme update-to-reconciliation ratios.

Kistler and Satyanarayanan [7] have conducted an empirical study of disconnected operation in the Coda file system, showing a low likelihood of concurrent updates [9]. A study of the Ficus file system [11] showed that optimistic replication used in an office environment achieved an extremely low conflict rate after the automation of conflict resolution for many applications and after removing identical conflicts. The study reported many identical conflicts, but their relationships to the system parameters and their implications on scaling were not explained. Neither the Coda nor the Ficus experience has examined the relationship between the update-to-reconciliation ratio and the formation of identical conflicts, which constitute the majority of conflicts.

There are relatively few studies that use both simulation and analytical methods to investigate the causes of conflicts and identical conflicts. Through a trace-driven simulation, one paper observed the inverse relationship between update locality and conflict and identical conflict rates [16]. However, our study has further investigated the effect of a bimodal traffic pattern on optimistic replicated systems. Another paper analytically characterized the conflict rate, but the results were limited to two replicas and not generalizable [18]. Our use of permuted states can capture the combinatorial growth of states at a small scale, which is representative of most replication scenarios.

There have also been other studies that examined the service quality of optimistic replication [12, 19]. However, the behavior of the conflict-rate curve was not deeply explored in these studies.

7. FUTURE WORK

Through this exploration of optimistic replication with both analytical and simulation approaches, we have begun to gain more mature intuition about system behavior. Although the state space of optimistic replication is large, we believe that a system can be reasonably characterized with fewer than 200 states, to capture all major aggregate statistics. We intend to design, implement, analyze, and automate algorithms to extract the top contributing states. The aggregate statistics obtained from a *trimmed* state diagram will be compared with the full state diagram for verification. If successful, we can use this simplified model to provide system feedback, prediction, and tuning at runtime.

Since traffic characteristics can significantly influence the fraction of conflicts that are identical, a fruitful area for future research would be to construct traffic filters that can shape the decomposition of conflicts. Ideally, we want no conflicts. If that is not possible, we want most conflicts to be identical, for easy resolution. (Note that for modeling purposes, an automated resolver such as those described in [9] and [11] would cause non-identical conflicts to behave as if they were identical ones.)

8. LESSONS & CONCLUSION

The results presented in this paper capture several iterations of *experimentation* with analytical methods, and many findings are not obvious in retrospect. We originally made naïve attempts to cluster states with conflicts into small sets of super states to simplify the computation. However, the results were similar to variable substitutions in complex equations. Although the resulting state-transition diagram had fewer states, the complexity of the equations remained unchanged.

After we discovered that conflicts occur as transitions, not states, we tried to insert *probing* states into each conflict-generating transition, in the hope that the equilibrium probability of the probing state would capture the conflict probability. Unfortunately, the probing states significantly distorted the results, making it difficult to compute conflict probabilities.

We have described methods to represent, automate, and optimize permuted states, which has enabled us to use analytical methods to explore the 4-replica base case of identical conflicts, and automate the analytical investigation up to 10 replicas. All results have been independently confirmed by a simulation based on version vectors. As a consequence, we discovered that update locality and bimodal access patterns are the primary factors that influence the fraction of identical conflicts.

The analysis of problems with exponential state spaces is always challenging. By introducing the concept of permuted states, we have developed a new technique that makes the base-case analysis of complex replicated systems tractable. As a result, we have been able to characterize and quantify important system behaviors that have previously been unrecognized or poorly understood.

Appendix A: State-Equilibrium Equations for Three Replicas

$$\lambda p_0 = \frac{2\mu}{3} p_3$$

$$\frac{2(\lambda + \mu)}{3} p_1 = \lambda p_0 + \frac{\mu}{3} p_6$$

$$\left(\frac{\lambda}{3} + \mu\right) p_2 = \frac{2\lambda}{3} p_1 + \frac{\lambda}{3} p_6$$

$$\left(\lambda + \frac{2\mu}{3}\right) p_3 = \frac{2\mu}{3} p_1 + \frac{\mu}{3} p_2 + \frac{2\mu}{3} p_5 + \frac{2\mu}{3} p_6 + \frac{\mu}{3} p_7$$

$$\mu p_4 = \frac{\lambda}{3} p_2 + \frac{\lambda}{3} p_7$$

$$\frac{2(\lambda + \mu)}{3} p_5 = \frac{2\mu}{3} p_2 + \frac{\lambda}{3} p_3 + \mu p_4 + \frac{2\mu}{3} p_7$$

$$\left(\frac{2\lambda}{3} + \mu\right) p_6 = \frac{2\lambda}{3} p_3$$

$$\left(\frac{\lambda}{3} + \mu\right) p_7 = \frac{2\lambda}{3} p_5 + \frac{\lambda}{3} p_6$$

$$p_0 + p_1 + p_2 + p_3 + p_4 + p_5 + p_6 + p_7 = 1$$

REFERENCES

- [1] Daniels D, Doo LB, Downing A, Elsbernd C, Hallmark G, Jain S, Jenkins B, Lim P, Smith G, Souder B, Stamos J. Oracle's Symmetric Replication Technology and Implications for Application Design. *Proceedings of SIGMOD Conference*, p. 467, 1994.
- [2] Foggia P, Sansone, C, Vento M. A Performance Comparison of Five Algorithms for Graph Isomorphism. *Proceedings of the 3rd Workshop on Graph-based Representations in Pattern Recognition*, 2001.
- [3] Golding RA. Weak-Consistency Group Communication and Membership. Ph.D. Dissertation, Department of Computer Science, University of California, Santa Cruz, 1992.
- [4] Gray J, Helland P, O'Neil P, Shasha D. The Dangers of Replication and a Solution. *Proc. of the 1996 ACM SIGMOD Conference*, pp.173-182, 1996.
- [5] Guy R, Popek G, Page TW. Consistency Algorithms for Optimistic Replication. *Proceedings of the 1st International Conference on Network Protocols, IEEE*, October 1993.
- [6] Kawell LJ, Beckhardt S, Halvorsen T, Ozzie R, Greif I. Replicated Document Management in a Group Communication System. Groupware: *Software for Computer-Supported Cooperative Work*, IEEE Computer Society Press, pp. 226-235, 1992.
- [7] Kistler JJ, Satyanarayanan M. Disconnected Operation in the Coda File System. *ACM Transactions on Computer Systems*, 10(1), February 1992.
- [8] Kubiawicz J, Bindel D, Chen Y, Czerwinski S, Eaton P, Geels D, Gummadi R, Rhea S, Weatherspoon H, Weimer W, Wells C, and Zhao B. *Proceedings of the Ninth international Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2000)*, November 2000.
- [9] Kumar P, Satyanarayanan M. Flexible and Safe Resolution of File Conflicts. *Proceedings of the 1995 USENIX Technical Conference*, pp. 95-106, January 1995.
- [10] Page T, Guy R, Heidemann J, Ratner D, Reiher P, Goel A, Kuenning G, Popek G. Perspectives on Optimistically Replicated, Peer-to-Peer Filing. *Software—Practice and Experience*, December 1997.
- [11] Reiher P, Heidemann J, Ratner D, Skinner G, Popek G. Resolving File Conflicts in the Ficus File System. *Proceedings of USENIX Conference*, pp. 183-195, June 1994.
- [12] Rowstron AIT, Lawrence N, Bishop CM. Probabilistic Modeling of Replica Divergence. *Proc. of the 8th IEEE Workshop on Hot Topics in Operating Systems*, May 2001.
- [13] Satyanarayanan M. Coda: A Highly Available File System for a Disconnected Workstation Environment. *Proceedings of the 2nd Workshop on Workstation Operating Systems*, September 1989.
- [14] Skiena S, Graph Isomorphism. *Implementing Discrete Mathematics: Combinatorics and Graph Theory with Mathematica*. Reading, Massachusetts, Addison-Wesley, pp. 181-187, 1990.
- [15] Terry DB, Theimer MM, Petersen K, Demers AJ, Spreitzer MJ, Hauser CH. Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System. *Proceedings of the 15th ACM Symposium on Operating Systems Principle*, December 1995.
- [16] Wang AIA. A Simulation Evaluation for Optimistically Replicated Filing Environments. Master's Thesis. Computer Science Department, University of California, Los Angeles, 1998.
- [17] Wang AIA, Peter Reiher, and Rajive Bagrodia. A Simulation Evaluation of Optimistically Replicated Filing in Mobile Environments. *Proceedings of the 18th IEEE International Performance, Computing, and Communication Conference (IPCCC)*, February 1999.
- [18] Wang AIA, Peter Reiher, Rajive Bagrodia, and Geoffrey Kuenning. Understanding the Behavior of the Conflict-Rate Metric in Optimistic Peer Replication. *Proceedings of the 5th IEEE International Workshop on Mobility in Databases and Distributed Systems (MDDS)*, Aix-en-Provence, France, September 2002.
- [19] Yu H, Vahdat A. Design and Evaluation of a Continuous Consistency Model for Replicated Servers. *Proc. of the 4th Symposium on Operating Systems Design and Implementation*, October 2000.

An-I A. Wang is an assistant professor of computer science at Florida State University. He received his Ph.D. and M.S. in computer science from UCLA in 2003 and 1998, and his B.A. in computer science from UC Berkeley in 1995. His research interests include file systems, optimistic peer-to-peer replication, performance evaluation, ad hoc network routing, operating systems, and distributed systems.

Geoff H. Kuenning is an associate professor of computer science at Harvey Mudd College. He received his Ph.D. in computer science from UCLA in 1997, and his B.S. and M.S. in computer science from Michigan State University in 1973 and 1974. From 1974 to 1989, he worked in the areas of operating systems and embedded systems. His research interests include file systems, performance analysis, and computer system security.

Peter Reiher is an adjunct associate professor of computer science at UCLA. He received his Ph.D. and his M.S. in computer science from UCLA in 1987 and 1984, respectively. He received his B.S. in electrical engineering from the University of Notre Dame in 1979. Dr. Reiher's research interests include active networks, advanced operating systems, parallel discrete event simulation, and security for distributed systems.