

Adding the Time Dimension to Majority Voting Strategies*

Hüseyin Aysan, Sasikumar Punnekkat, and Radu Dobrin
Mälardalen Real-Time Research Centre, Mälardalen University, Västerås, Sweden
{huseyin.aysan, sasikumar.punnekkat, radu.dobrin}@mdh.se

Abstract

Real-time applications typically have to satisfy high dependability requirements and require fault tolerance in both value and time domains. A widely used approach to ensure fault tolerance in dependable systems is the N-modular redundancy (NMR) which typically uses a majority voting mechanism. However, NMR primarily focuses on producing the correct value, without taking into account the time dimension. In this paper, we propose a new approach, Voting on Time and Value (VTV), applicable to real-time systems, which extends the modular redundancy approach by explicitly considering both value and timing failures, such that correct value is produced at correct time, under specified assumptions. We illustrate the proposed approach by an algorithm applicable for triple modular redundancy (TMR).

1. Introduction

Most real-time applications typically have to satisfy high dependability requirements due to their interactions and possible impacts on the environment. Ensuring dependable performance of such systems typically involves both fault prevention and fault tolerance approaches in their design. Usage of redundancy is the key for achieving fault tolerance and it has been employed successfully in the physical, temporal, information and analytical domains of a large number of critical applications. Static techniques such as N-modular redundancy (NMR) have been used in safety and mission critical applications, most often in the well-known form of triple-modular redundancy (TMR), where three nodes are used for replication [9]. The key attraction of this approach lies in its low overhead and fault masking abilities, without the need for backward recovery. The disadvantages include the cost of redundancy and single point failure mode of the voter. Traditionally, voters are constructed as simple electronic circuits so that a very high reliability can be achieved. Usage of triplicated voters has been

employed to take care of the single-point failure mode in case of highly critical systems [8]. Surveys and taxonomies on several voting strategies have been presented [7, 5].

Replicated nodes' output delivery times can vary due to several factors, such as clock drifts, node failures, processing and scheduling variations at node level, as well as communication delays. Most of the existing voting strategies, however, focus solely on masking value failures by assuming that the system is tightly synchronized, as presented in [6]. On the other hand, loosely synchronized systems may be an attractive alternative due to, e.g., low overheads, requiring, however, specifically designed asynchronous voting algorithms to compensate for the timing variations.

A simple approach towards tolerating both value and timing failures in a replica using the NMR approach could be adding time stamps to the replica outputs. Then, majority voting on time stamp values could detect possible timing anomalies of the nodes, under the unrealistic assumptions that the communication is ideal and nodes never halt. Moreover, this approach is unable to mask late timing failures.

Shin and Dolter [11] proposed two voting techniques applicable to real-time systems, relaxing the tight synchronization requirements, viz., *Quorum Majority Voting (QMV)* and *Compare Majority Voting (CMV)*. QMV performs majority voting among the received values as soon as $2n+1$ out of $3n+1$ replicas deliver their outputs to the voter, thus, guaranteeing detection of majority of non-faulty values even in the case n replicas fail. CMV masks failures of n out of $2n+1$ replicas as in basic majority voting. The main difference is that in CMV the output is delivered as soon as a majority consisting of identical values has been received, i.e., without waiting for the rest of the replicas. Both QMV and CMV provide outputs within a bounded time interval, as long as the assumptions regarding the maximum number of failures hold. However, QMV and CMV are unable to detect assumption violations in the time domain.

In this paper, we propose a novel approach, *Voting on Time and Value (VTV)*, which performs majority voting in both time and value domains. Our approach enhances the fault tolerance abilities of NMR by restricting the replica outputs to be both correct in value, and delivered within a

*This work was partially supported by the Swedish Foundation for Strategic Research via the strategic research centre PROGRESS.

specified admissible time interval, under specified assumptions. Furthermore, our approach is able to detect assumption violations in time domain.

The rest of the paper is organized as follows: In Section 2 we present the system model and the assumptions used in this paper. Section 3 describes our approach, illustrates it by an instantiation to a system using triplicated nodes. We conclude the paper in Section 4 outlining the on-going and future work.

2. System Model

In this paper, we assume a distributed real-time system, where each critical node is replicated for fault tolerance, and replica outputs are voted to ensure correctness in both value and time. For the sake of readability, in the rest of the paper, we denote the i^{th} replica of a node N by N_i . The output delivered by N_i , is specified by two domain parameters, viz., value and time [1, 10, 3]:

$$\text{Specified output for } N_i = \langle v_i^*, t_i^*, \Delta_v, \Delta_t \rangle$$

where v_i^* is the correct value, t_i^* is the correct time point when the output should be delivered, $[v_i^* - \Delta_v, v_i^* + \Delta_v]$ is the admissible value range and $[t_i^* - \Delta_t, t_i^* + \Delta_t]$ is the admissible time interval for output delivery as per the real-time system specifications.

An output delivered by N_i is denoted as:

$$\text{Delivered output from } N_i = \langle v_i, t_i \rangle$$

where v_i is the value and t_i is the time point at which the value was delivered.

We define the output generated by replica N_i as *incorrect in value domain* if $v_i < v_i^* - \Delta_v$ or $v_i > v_i^* + \Delta_v$, and *incorrect in time domain* if $t_i < t_i^* - \Delta_t$ (*early timing failure*), or if $t_i > t_i^* + \Delta_t$ (*late timing failure*).

Assumptions: Our approach relies on the following set of assumptions (to a large extent based on [4]):

1. non-faulty nodes produce values within a specified admissible range and within a specified time interval after each computation block
2. replica outputs with incorrect values do not form majority
3. incorrectly timed replica outputs do not form majority
4. a maximum permissible *drift* δ from the global time is specified and ensured by infrequent synchronization (which is significantly less costly than tight synchronization)
5. the voter does not fail.

3. Voting on Time and Value (VTV)

In this section we present our novel voting strategy that explicitly considers failures in both time and value domains. As a consequence of assumption 5, in the worst case, the maximum deviation between any two replica outputs is 2δ . Hence, in VTV approach, agreement in the time domain is reached when a majority of replicas deliver their outputs within this derived time interval of 2δ (referred to as feasible window henceforth). If a node has n replicas, then at least $m = \lceil \frac{n+1}{2} \rceil$ outputs from these replicas need to match for establishing majority. The number of groups with m sequential replica outputs within n replica outputs is $n-m+1$. Since the majority in time domain can be formed by any of these groups, a separate feasible window needs to be initiated upon receiving each of first $n-m+1$ replica outputs. We keep track of the feasible windows by using simple countdown timers. Once an agreement in time domain is obtained, then values are voted. If an agreement in value domain is not obtained for a particular feasible window, the process continues with subsequent feasible windows, until a majority in time and an agreement in value can be formed, or an assumption violation is detected.

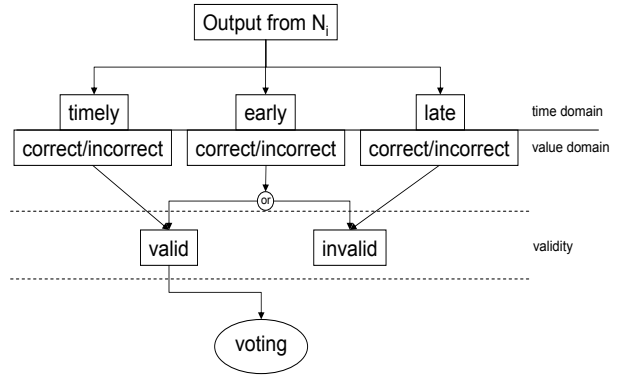


Figure 1. Replica output flow through voter

Depending on the real-time application characteristics, a value produced by a node may be considered *valid* or *invalid* for the purpose of voting, in case it is produced early. An illustration of replica output flow through the voter is given in Figure 1. An issue is the choice of the set of *valid* values to be used in the voting mechanism, i.e., *all* received values vs. *all timely* received values. We illustrate this voting dilemma by using the scenario described in Figure 2. Let us assume, e.g., an airbag control system where a sensor is replicated in five different nodes and produces one out of two values periodically, e.g., value a in case of a collision detection and value b otherwise. If a collision is detected at a time $t \leq t_1$ let us assume that the airbag has to inflate within a time interval $[t_{start}, t_{end}]$, where $t_2 < t_{start} \leq t_3$

and $t_5 \leq t_{end}$. In our example, the first two values are detected as *early* and the last three are identified as *timely*. However, in this case, an early value has to be taken into consideration in the voting since an early collision detection is still a valid output with respect to the value domain. Thus, the output has to be voted upon receiving the last value at time t_5 , among *all* values, i.e., $a, a, a, b,$ and b , resulting in an output a at time $(t_5 + \epsilon)$ (where ϵ is the time required for the voting and is assumed to be negligible in this paper for simplifying the presentation).

On the other hand, let us assume that the same Figure 2 illustrating an altitude measurement sensor in an airplane, replicated by five nodes to read and output the altitude periodically to the voter, where data freshness may be a more desirable aspect. As the correct window of time for the output is the same as described in the previous example, the only relevant values to be taken into consideration by the voter are $a, b,$ and b corresponding to the time points $t_3, t_4,$ and t_5 respectively. Hence, the output produced at time $(t_5 + \epsilon)$ is b .

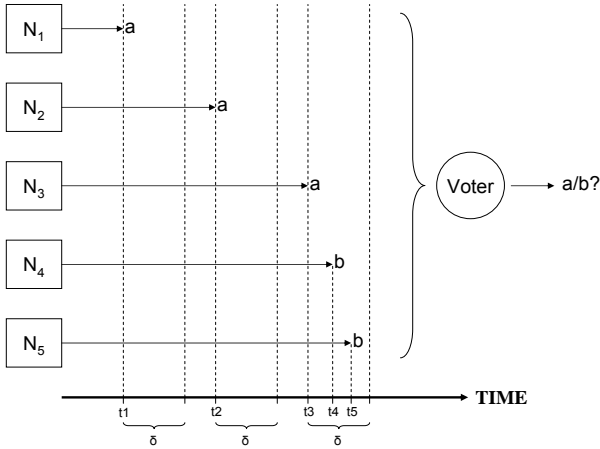


Figure 2. Voting dilemma

Upon finding a feasible window, if a majority in value domain is obtained with all the values received so far, the voter delivers the majority value without waiting for the rest of the replicas. Otherwise, either a majority in value domain, receipt of all replica outputs, or the end of the feasible window is waited for, whichever comes first. If a majority in value domain is obtained while waiting, it is delivered as the correct output. The decision on whether the early generated replica outputs are involved in value voting or not results in two cases at this point:

Case 1 *Early and timely outputs are considered valid.* If the end of the feasible window is reached with a majority among the received values, it is delivered as correct output.

Case 2 *Only timely outputs are considered valid.* If the end of the feasible window is reached with a majority among the timely received values, it is delivered as correct output.

If the end of the feasible window is reached without an agreement in value domain, the process continues with a subsequent feasible window. If the last feasible window is reached, or all replica outputs are received without reaching an agreement on the values, disagreement is signalled to the rest of the system.

3.1 VTV in TMR

In this section, we present an instantiation of our approach to triple modular redundancy which can tolerate single node failures in value domain, time domain or both (Algorithm 1). In this example, we assume early timing failures as *invalid* for the purpose of voting. However, the validity of such values can be easily tuned in the algorithm.

Majority in time domain is achieved if at least two values are delivered to the voter within a time interval less than or equal to 2δ , since this is the maximum deviation in time among all the values as long as there is no failure. Majority in value domain is formed if at least two of the *timely* outputs have the same value.

The algorithm signals disagreement in case majority condition is not satisfied in any of the domains, thus enabling a fail-safe or fail-stop behavior of the system.

The replicated nodes' output values are stored in local variables V1, V2 and V3. Values are assigned to these variables in the order of receiving inputs from the nodes (i.e., the first received value is stored in V1, the second one in V2 and the last one in V3). Two countdown timers, C1 and C2, initially set to 2δ , are used to keep track of feasible windows in order to identify majority in time domain.

The algorithm waits for the first node output to be delivered and then starts C1. It continues by waiting for the second node output and starts C2 upon its arrival. If both values have arrived before C1 expires, and have matching values, the voter will output the correct value. Otherwise we have two cases:

Case 1 C1 has not reached zero, but the values V1 and V2 do not match. In this case, the algorithm waits for V3 until C1 reaches zero. If the third value arrives before C1 reaches zero and matches either V1 or V2, the algorithm outputs the matching value since all values are timely and there is an agreed value. In case of assumption violation, i.e., there exists no replica output pair matching in value domain, the algorithm signals *disagreement*. If the third value does not arrive before C1 reaches zero, the algorithm waits for V3 until C2 reaches zero. If V3 is received and matches V2 before

Algorithm 1: VTV

```
input :  $v_1, v_2, v_3 = NULL$ 
output:  $v_{out}$  or indication of disagreement
/* Inputs are ordered wrt reception */
/* Voting in value domain is performed
   among timely received values */
1  $C_1, C_2 \leftarrow 2\delta$ ; // countdown timers
2 while  $v_1 = NULL$  do wait;
3 start  $C_1$ ;
4 while  $v_2 = NULL$  do wait;
5 start  $C_2$ ;
6 if  $C_1 > 0$  then
7   if  $v_1 = v_2$  then
8     output  $v_1$ ;
9   else
10    while  $C_1 > 0$  and  $v_3 = NULL$  do wait;
11    if  $C_1 > 0$  and ( $v_3 = v_1$  or  $v_3 = v_2$ ) then
12      output  $v_3$ ;
13    else if  $v_3 \neq NULL$  then
14      signal disagreement;
15    else
16      while  $C_2 > 0$  and  $v_3 = NULL$  do wait ;
17      if  $v_3 = v_2$  then
18        output  $v_3$ ;
19      else
20        signal disagreement;
21      end
22    end
23  end
24 else if  $C_2 > 0$  then
25   while  $C_2 > 0$  and  $v_3 = NULL$  do wait;
26   if  $v_3 = v_2$  then
27     output  $v_3$ ;
28   else
29     signal disagreement;
30   end
31 else
32   signal disagreement;
33 end
```

C_2 reaches zero, the algorithm outputs the matching value. Otherwise the algorithm signals *disagreement*.

Case 2 C_1 has reached zero. In this case, V_1 is considered invalid, and the algorithm waits for V_3 until C_2 reaches zero, as only a match between V_2 and V_3 may result in an agreement. If the values do not match or V_3 has not been received at all, the algorithm signals *disagreement*.

4. Conclusions

In this paper we have presented a new voting strategy called Voting on Value and Time (VTV) for redundant systems, to explicitly consider both value and timing failures

for achieving fault tolerance in real-time applications. Under specified failure assumptions, our method is capable of producing the correct output as well as identifying the correct window of time in which the output has to be delivered.

We have presented an algorithm for the particular case where one output is replicated in three different nodes, and illustrated the basic idea on how we perform the voting in both value and time domain.

Our ongoing research indicates that VTV, when used in the general case to mask arbitrary number of value and timing failures, is cost-effective in comparison with the number of nodes required by majority voting in NMR. The main reason is that, in our approach, a non-faulty node can be successfully used to mask both a value and a timing failure in the voting procedure.

References

- [1] A. Avizienis, J. Laprie, and B. Randell. Fundamental concepts of dependability. *Research Report N01145, LAAS-CNRS*, April 2001.
- [2] D. Blough and G. Sullivan. A comparison of voting strategies for fault-tolerant distributed systems. *Proceedings of the Ninth Symposium on Reliable Distributed Systems*, pages 136–145, 1990.
- [3] A. Bondavalli and L. Simoncini. Failure classification with respect to detection. *Proceedings of 2nd IEEE Workshop on Future Trends in Distributed Computin*, pages 47–53, 1990.
- [4] P. Ezhilchelvan, J.-M. Helary, and M. Raynal. Building responsive tmr-based servers in presence of timing constraints. *Object-Oriented Real-Time Distributed Computing, 2005. ISORC 2005. Eighth IEEE International Symposium on*, pages 267–274, 2005.
- [5] F. D. Giandomenico and L. Strigini. Adjudicators for diverse-redundant components. *Proceedings of the Ninth Symposium on Reliable Distributed Systems*, pages 114–123, 1990.
- [6] H. Kopetz. Fault containment and error detection in the time-triggered architecture. *Autonomous Decentralized Systems, 2003. ISADS 2003. The Sixth International Symposium on*, pages 139–146, 2003.
- [7] G. Latif-Shabgahi and a. S. B. J.M. Bass. A taxonomy for software voting algorithms used in safety-critical systems. *IEEE Transactions on Reliability*, 53(3):319–328, 2004.
- [8] R. E. Lyons and W. Vanderkulk. The use of triple-modular redundancy to improve computer reliability. *Journal of Research and Development*, 6:200–209, 1962.
- [9] J. V. Neuman. Probabilistic logics and the synthesis of reliable organisms from unreliable components. *Automata Studies*, pages 43–98, 1956.
- [10] D. Powell. Failure mode assumptions and assumption coverage. *Proceedings of 22nd International Symposium on Fault-Tolerant Computing*, pages 386–395, 1992.
- [11] K. Shin and J. Dolter. Alternative majority-voting methods for real-time computing systems. *IEEE Transactions on Reliability*, 38(1):58–64, 1989.