# Scalable and Fault Tolerant Failure Detection and Consensus*

Amogh Katti, Giuseppe Di Fatta
School of Systems Engineering
University of Reading
Reading, RG6 6AY, UK
{a.p.katti, G.DiFatta}@reading.ac.uk

Thomas Naughton,
Christian Engelmann
Computer Science and Mathematics Division
Oak Ridge National Laboratory
Oak Ridge, TN 37831-6173, USA
{naughtont, engelmannc}@ornl.gov

## ABSTRACT

Future extreme-scale high-performance computing systems will be required to work under frequent component failures. The MPI Forum's User Level Failure Mitigation proposal has introduced an operation, *MPI_Comm_shrink*, to synchronize the alive processes on the list of failed processes, so that applications can continue to execute even in the presence of failures by adopting algorithm-based fault tolerance techniques. This *MPI_Comm_shrink* operation requires a fault tolerant failure detection and consensus algorithm. This paper presents and compares two novel failure detection and consensus algorithms. The proposed algorithms are based on Gossip protocols and are inherently fault-tolerant and scalable. The proposed algorithms were implemented and tested using the Extreme-scale Simulator. The results show that in both algorithms the number of Gossip cycles to achieve global consensus scales logarithmically with system size. The second algorithm also shows better scalability in terms of memory and network bandwidth usage and a perfect synchronization in achieving global consensus.

## Categories and Subject Descriptors

D.4.4 [**Software**]: Communications Management—*message sending, network communication*; D.4.8 [**Software**]: Performance—*measurements, simulation*; D.4.5 [**Reliability**]: Fault-tolerance

---

## Keywords

Fault-tolerant MPI, User-level failure mitigation, Failure detection, Consensus, Gossip protocols

## 1. INTRODUCTION

Resilience is a critical challenge as high-performance computing (HPC) systems continue to increase component counts, individual component reliability decreases (due to shrinking process technology [1] and near-threshold voltage operation [12]), and software complexity increases [17]. In spite of frequent faults, errors and failures, parallel application correctness and execution efficiency are essential to ensure the success of extreme-scale HPC systems.

As cost constraints limit resilience mitigation in hardware, a cooperative approach between hardware and software is needed to efficiently mitigate faults, errors, and failures at the appropriate layer. However, application-level checkpoint/restart has been the dominant HPC fault tolerance method for decades. Every detectable uncorrectable error results in a complete application abort and restart from previously saved checkpoint state, even if the error could have been more efficiently handled by the application, e.g., using forward error correction with erasure codes [11].

Algorithm-based fault tolerance (ABFT) [9, 11, 13] may be able to deal with loss of application state, such as caused by a failing compute process or node, more efficiently through reconfiguration and adaptation without the need for a more drastic recovery measure, such as a global rollback. The employed fault tolerance techniques by the application may include error correction using data redundancy or encoding, and re-execution using local checkpoints.

The Message Passing Interface (MPI) is the dominant parallel programming interface for facilitating communication between compute processes in HPC. Despite its popularity, MPI is not fault tolerant. Recent efforts in MPI fault tolerance focused on user-level failure mitigation (ULFM) [2]. The proposed ULFM extensions to the MPI standard enable applications to be notified of MPI process failures, to create a new MPI communicator object that excludes known failed MPI processes from further communication using the *MPI_Comm_shrink()* operation, and to achieve a uniform agreement on a value among the non-failed MPI processes using the *MPI_Comm_agree()* operation.

This paper particularly focuses on the implementation of the proposed *MPI_Comm_shrink()* operation and its need to perform an agreement on the group of failed MPI processes among the non-failed MPI processes, even while MPI

process failures occur. More precisely, the new MPI communicator object created by the collective *MPL_Comm_shrink()* operation contains a consistent group of MPI processes at every participating MPI process that excludes at least every failed MPI process that has previously raised a failure notification to the application. MPI processes that fail during the *MPL_Comm_shrink()* operation may be excluded as well, but consistently at every participating MPI process.

According to the ULFM proposal, only fail-stop process failures are considered; when a process fails, it stops communicating with the rest of the processes. The method of failure detection is not defined. However, the ULFM proposal specifies that an operation involving a failed MPI process must always complete in a finite amount of time. If an operation does not involve a failed MPI process, it must not raise a MPI process failure exception. This provides implementers with different options for failure detection. For example, a correct MPI implementation may provide failure detection only for MPI processes involved in an ongoing operation and may postpone detection of other failures until necessary. The knowledge about detected failed MPI processes is local and only constructed as globally consistent state in the form of a new MPI communicator object using the *MPL_Comm_shrink()* operation, thus requiring a fault-tolerant consensus algorithm to uniformly agree on the group of failed MPI processes.

This paper investigates the use of Gossip-based protocols to detect failures and disseminate this information in a consistent manner to support the *MPL_Comm_shrink()* operation. The proposed consensus algorithm can also be used to support the *MPL_Comm_agree()* operation.

Gossip-based (or Epidemic) protocols are a robust and scalable communication paradigm to disseminate information in a large-scale distributed environment using randomized communication. They have the advantage of inherent robustness and scalability with respect to global communication schemes based on deterministic communication patterns. Applications and services based on Gossip-based protocols for large- and extreme-scale systems have been proposed in many fields of distributed computing. Recently, Gossip-based protocols have also been applied in the context of HPC [18, 19].

In this paper, two Gossip based failure detection and consensus algorithms using randomized pinging have been developed and tested by means of simulations. The first algorithm stores the system view in a matrix at each process to facilitate consensus detection. It detects failures before and during the execution of the algorithm and is hence completely fault tolerant. The second algorithm detects consensus on the failed processes using a heuristic method based on a list of failed processes, thus increasing the memory scalability with respect to the first algorithm. It also transfers less Gossip data as only the list of failed processes is sent and hence consumes negligible network bandwidth. For both algorithms, the number of Gossip cycles to detect consensus scales logarithmically with the system size.

The paper is structured as follows. Section 2 details the proposed Gossip-based failure detection and consensus algorithms, where Section 2.1 focuses on failure detection using stochastic pinging, Section 2.2 on achieving consensus using global knowledge, and Section 2.3 on achieving consensus using heuristics. Section 3 presents experimental results and Section 4 discusses related work. Section 5 concludes the paper with a summary and a discussion of future work.

## 2. GOSSIP-BASED FAILURE DETECTION AND CONSENSUS

The *MPL_Comm_shrink()* operation must implement a consensus algorithm that achieves agreement on the set of failed MPI processes that have previously raised a failure notification, i.e., on the group of known failed MPI processes at each participating MPI process at the start of the operation. MPI process failures that occur during the operation will eventually be detected and corresponding failure notifications will be raised during successive MPI communication operations. To avoid coarse-grain iterative agreement on the group of failed MPI processes with successive calls to *MPL_Comm_shrink()* by an application, the consensus algorithm may additionally include an agreement on the group of MPI processes that fail during the *MPL_Comm_shrink()* operation, i.e., during the consensus algorithm.

In both cases, the consensus algorithm needs to be fault tolerant, i.e., deal with already known or newly detected MPI process failures. Also, in both cases, a complete failure detector is implemented that detects fail-stop MPI process failures by combining failure detection and consensus. Every MPI process independently detects MPI process failures. Consensus on the failed MPI processes is then achieved by aggregating these MPI process failure detections with the help of a consensus algorithm.

A straightforward way to detect MPI process failures during the *MPL_Comm_shrink()* operation is pinging, wherein a process asks another whether it is alive. A reply indicates a positive response, while a failure notification from the underlying MPI runtime marks a negative response. The failure detector in the MPI runtime can be based on a simple communication timeout. In an alternative method, a process periodically sends a heartbeat message to let another process know that it is alive. The receiving process monitors the incoming heartbeat messages and marks a process as failed upon a failure notification from the underlying MPI runtime. The failure detector can be based on a simple communication timeout for the periodic heartbeat. The work presented in this paper is based on pinging for failure detection during the *MPL_Comm_shrink()* operation, using the already existing MPI runtime reporting mechanism for process failures.

The consensus during the *MPL_Comm_shrink()* operation involves all fault free processes agreeing on the group of failed processes. In general, a solution to the consensus problem exists only in certain environments [20]. For instance, consensus is not possible in completely asynchronous environments [7]. However, in an asynchronous environment, failures can still be detected with completeness and accuracy, leading to a uniform view of the system at each process [5], although, group membership may not be agreed upon [4]. In the context of MPI, the proposed ULFM extensions are based on realistic assumptions, such as fail-stop, no recovery, a synchronous model, and only short periods with exceptionally high MPI process failure rates.

This paper proposes two scalable failure detection and consensus algorithms for *MPL_Comm_shrink()* based on Uniform Gossiping. The Gossip messages are implicitly used to implement stochastic failure detection. Consensus is achieved and detected by maintaining the system state in a matrix or by maintaining a failed process list at each process. Gossip

messages carry failure information with them and disseminate known failures at exponential speed. When a process $p$ sends a Gossip message to process $q$, this process $q$ comes to know about the failures that were directly detected by $p$, thus detecting failures indirectly. This Gossip message from $p$ not only contains the failures directly detected by $p$, but also indirectly detected failures through received Gossip messages. These indirect detections of process $p$ are propagated as well to process $q$, resulting in exponential information dissemination. When the failure information of a process is disseminated to all fault free processes, consensus on its failure is reached in a logarithmic number of Gossip cycles (as shown in the experiments in Section 3).

The proposed failure detection and consensus algorithms work under the following assumptions.

- Processes are assumed to be connected by a reliable communication medium.
- A synchronous system model is assumed, i.e., a non-failed MPI process responds to a message within a known, finite amount of time.
- The fail-stop model is assumed, i.e., a failed MPI process stops communicating.
- Faults are assumed to be permanent, i.e., a failed MPI process does not recover.
- A process once detected as failed is detected to have failed by all the processes eventually.
- Periods of system stability are assumed, i.e., MPI process failures during the consensus algorithm will, at some point, stop for a long enough period to reach consensus.

## 2.1 Failure Detection using Stochastic Pinging

In this Section, we discuss the MPI process failure detection feature of a *MPI_Comm_shrink()* operation using stochastic pinging as part of a Gossip-based protocol. Every process independently detects failures by pinging a random process periodically. During a Gossip cycle of length $T_{Gossip}$ units, a process $i$ randomly selects a process $j$ to ping according to a uniform probability distribution function. If process $j$ replies by the end of the current Gossip cycle, then process $i$ finds it to be alive; failed otherwise. Figure 1 shows the pseudocode for the algorithm.

| At process $i$ |
| --- |
| **At each cycle (every $T_{Gossip}$ time units):** |
| 1    $j = getRandomProcess()$ |
| 2    send a ping message to $j$ |
| **At event: received a ping message from $j$:** |
| 3    send a reply message to $j$ |
| **At event: timeout without receiving a reply from $j$:** |
| 4    mark $j$ to have failed |

**Figure 1: Failure detection using pinging**

During a cycle of Uniform Gossiping, the probability of a process being selected as destination of zero, one or more ping messages follows a binomial distribution. Ultimately any failed process is quickly detected by one or more of the non-failed processes, thus initiating the epidemic exponential propagation of information to achieve consensus. The adopted Gossip-based approach can tolerate moderate-to-low message loss rates and delays as it is intrinsically fault tolerant.

## 2.2 Consensus using Global Knowledge

In this Section, we discuss achieving consensus on the set of failed MPI processes during the *MPI_Comm_shrink()* operation by maintaining global knowledge at each MPI process. As shown in the algorithm of Figure 2 each process $p$ detects failures by pinging random processes and maintains a fault matrix $F_p$ to store the status of all processes as believed by it and also by all the other processes. An entry $F_p[i,j]$ in a process's fault matrix indicates the status of process $j$ as detected by process $i$ (1 if detected to have failed; 0 otherwise). The algorithm can be divided into four logical tasks according to the functions performed: (1) initialization, (2) failure detection, (3) fault matrix update, and (4) check for consensus.

| At process $p$ |
| --- |
| **Require: Fault Matrix $F_p[j,k]$ where** $0 \le j,k < n$ |
| $F_p[j,k]$ - status of process $k$ as detected by process $j$ |
| **Initialisation:** |
| //all processes are assumed to be alive at start |
| 1    $for(j = 0, j < n, j + +)$ |
| 2      $for(k = 0, k < n, k + +)$ |
| 3        $F_p[j,k] = 0$ |
| 4      $endfor$ |
| 5    $endfor$ |
| **At each cycle (every $T_{Gossip}$ time units):** |
| //failure detection using pinging |
| 6    $q = getRandomProcess()$ |
| 7    send a message of type ping to $q$ piggybacking $F_p$ |
| 8    set timeout event $E_q = < T_{Gossip}, q >$ for receiving reply message from $q$ |
| 9    $for(k = 0, k < n, k + +)$ //check for consensus on k |
| 10     $temp = 0$ |
| 11     $for(j = 0, j < n, j + +)$ |
| 12       $if(Fp[j,k]||Fp[p,j])$ |
| 13         $temp = temp + 1$ |
| 14       $endif$ |
| 15     $endfor$ |
| //consensus is reached when all fault free |
| //processes have detected the failed process |
| 16     $if(temp == n)$ |
| 17       $consensus\_reached(k)$ |
| 18     $endif$ |
| 19   $endfor$ |
| **At event: received a message from $r$ piggybacked** **with $F_r$:** |
| 20   $if(message\_type == ping)$ |
| 21     send a message of type reply to $r$ piggybacking $F_p$ |
| 22   $endif$ |
| //merge the fault matrices |
| 23   $for(k = 0, k < n, k + +)$ |
| 24     $for(j = 0, j < n, j + +)$ |
| //propagation of remote failure detections |
| 25       $if(j \neq p)$ |
| 26         $F_p[j,k] = F_p[j,k]||F_r[j,k]$ |
| 27       $else$ |
| //indirect local failure detection |
| 28         $F_p[p,k] = F_p[p,k]||F_r[r,k]$ |
| 29       $endif$ |
| 30     $endfor$ |
| 31   $endfor$ |
| **At event: timeout $E_q$ and no reply message received** **from $q$:** |
| //mark $q$ to have failed (direct failure detection) |
| 32   $F_p[p,q] = 1$ |

**Figure 2: Failure detection and consensus by maintaining global knowledge (algorithm 1)**

Every process initializes with the assumption that every

other process in the system is alive and no other process has yet detected any failures (lines 1-5 of Figure 2).

To detect failures, every $T_{Gossip}$ time units a process $p$ randomly selects a process $q$ and sends a ping message to it piggybacking the fault matrix $F_p$ (lines 6-7). Sending the entire matrix facilitates the propagation of not only process $p$'s detections but also other processes' detections known to process $p$, thus exponentially propagating information. A timeout event is then set to receive a reply from $q$ during the current Gossip cycle (line 8). An asynchronous reply follows upon reception of this ping message (lines 20-22). At the expiry of the current Gossip cycle if a reply message from $q$ has not been received by $p$, $q$ is detected (directly) to have failed (line 32).

Upon reception of a Gossip (ping or reply) message at $p$ from $r$, the local fault matrix $F_p$ is updated by performing an $OR$ operation of the corresponding elements in $F_r$ except the row $p$ (lines 23-31 except line 28). The row $p$ in $F_p$ is updated to include the detections of process $r$ (row $r$ in $F_r$) and thus performing indirect failure detections at $p$ (line 28).

Finally, to check if consensus has been reached on the failure of a process $k$ at $p$, a logical $OR$ operation is performed between the corresponding elements of the $k^{th}$ column of the fault matrix and its $p^{th}$ row. Consensus is reached when all fault free processes have detected the faulty one (lines 9-19).

In this algorithm processes have to maintain local knowledge of the entire system state, with $O(n^2)$ memory requirement, in order to check for consensus. Network bandwidth consumption is also high due to the transfer of this local system state as part of the Gossip. Designing scalable consensus algorithms, which achieve consensus without maintaining global knowledge at each process and with moderate network utilization, is a challenging task. Previous work [15] has shown how to detect convergence in epidemic aggregation by means of heuristic methods, such as comparing local state with a random sample of remote states or running multiple protocol instances locally. In the next section a similar heuristic method for consensus detection with a more efficient memory and bandwidth utilization is presented.

## 2.3 Efficient Heuristic Consensus

Maintaining a matrix of size $n^2$ elements, where $n$ is the number of processes in the system, and sending it as part of the Gossiping consumes a lot of memory and network bandwidth. Storing only the failed processes in a list and sending the same while Gossiping can avoid this. The algorithm in Figure 3 uses a list instead of a matrix and detects consensus using an heuristic method.

Each process maintains a fault list $L_p$ to store the failed processes known to it. An entry in this list is a 2-tuple $< r, ccnt >$, where $r$ is the rank of the failed process and $ccnt$ is the consensus count associated with it. $ccnt$ is the length of an unbroken sequence of pinged processes that have $r$ in their failed process list. The algorithm can be divided into four logical tasks according to the functions performed: (1) initialization, (2) failure detection, (3) updating the fault list and (4) checking for consensus.

Every process starts with the assumption that every other process in the system is alive and hence has its fault list empty (line 1 of Figure 3).

To detect failures, every $T_{Gossip}$ time units a process $p$ randomly selects a process $q$ and sends a ping message to it, piggybacking the fault list $F_p$ (lines 2-3). A timeout event is

| At process $p$ |
|---|
| **Require:** Fault list $L_p = \{< r, ccnt >, ...\}$ |
| **Initialisation:** |
| 1    $L_p \leftarrow \{\}$ |
| **At each cycle (every $T_{Gossip}$ time units):** |
|      //failure detection using pinging |
| 2      $q = getRandomProcess()$ |
| 3      send a message of type ping to $q$ piggybacking $L_p$ |
| 4      set timeout event $E_q = < T_{Gossip}, q >$ for receiving reply message from $q$: |
|      //check for consensus |
| 5      foreach entry $< r, ccnt >$ in $L_p$ |
| 6        $if((curr\_cycle > log(n))$ && $(ccnt \geq MIN\_CCNT))$ |
| 7          $consensus\_reached(r)$ |
| 8        $endif$ |
| 9      end foreach |
| **At event: received a message from $r$ piggybacked with $L_r$:** |
| 10    $if(message\_type == ping)$ |
| 11      send a message of type reply to $r$ piggybacking $L_p$ |
| 12    $endif$ |
|      //merge the contents of the two lists and update ccnt of |
|      //each entry (indirect failure detection) |
| 13    $L_p \leftarrow merge(L_p, L_r)$ |
| **At event: timeout $E_q$ and no reply message received from $q$:** |
|      //mark $q$ to have failed and add it to the list (indirect |
|      //failure detection) |
| 14    $L_p \leftarrow L_p \cup \{< q, 0 >\}$ |

**Figure 3: Efficient failure detection and consensus (algorithm 2)**

then set to receive a reply from $q$ during the current Gossip cycle (line 4). An asynchronous reply follows upon reception of this ping message (lines 10-12). At the expiry of the current Gossip cycle, if a reply message from $q$ has not been received by $p$, $q$ is detected to have failed and is added to the fault list with $ccnt$ initialized to 0 (line 14).

Upon reception of a Gossip (ping or reply) message at $p$ from $r$, the remote fault list $F_r$ is merged with the local fault list $F_p$ (line 13). This merge operation includes incrementing the $ccnt$ of an element $< r, ccnt >$ in the local list if $r$ is also present in the received fault list. The $ccnt$ of the element is reset if the received fault list does not contain the element $< r, ccnt >$. If the received fault list contains an element $< r, ccnt >$ not present in the local fault list, it is added to the local fault list with $ccnt$ set to 0.

Finally, to check if consensus has been reached on the failure of process $r$, the value of the associated $ccnt$ is checked. The gossiping is performed for at least $log(n)$ number of cycles to allow the spreading of the failure detection information in the system. When a minimum value of the consensus counter ($MIN\_CCNT$) is reached, consensus on the failure of process $r$ is detected (lines 5-9).

In algorithm 1, when a process $p$ has detected consensus on $k$ it is certain that all processes have detected the failure of $k$, although they may not have detected consensus yet. Whereas in algorithm 2 when a process $p$ has detected consensus on $k$, some processes may have not detected the failure of $k$ yet. The initial phase of $log(n)$ cycles before detecting consensus is used to avoid false consensus detections and a small value of the threshold $MIN\_CCNT$ can be adopted.

## 3. EXPERIMENTAL RESULTS

The algorithm for consensus using global knowledge and

the algorithm for consensus using the heuristic approach were both implemented in the form of MPI applications, using basic MPI point-to-point communication primitives. The fault matrix of algorithm 1 is implemented as an integer matrix, whereas a failed process identifier in algorithm 2 is an integer. Fault injection was simulated by excluding a process from further communication.

The experiments were performed on a Linux cluster computer with one head node and 16 compute nodes. The head node has two AMD Opteron 4386 3.1 GHz processors with eight cores per processor and 64 GB RAM. The compute nodes have one Intel Xeon E3-1220 3.1GHz processor with four cores per processor and 16 GB RAM. The entire system has a total of 80 compute cores. The nodes are connected by Gigabit Ethernet. The system is running the Ubuntu 12.04 LTS operating system and Open MPI 1.6.5.

Experiments were executed using the Extreme-scale Simulator (xSim 0.5) [6, 14] atop the Linux cluster to evaluate the algorithms at significantly larger scale than the available physical system. xSim is a performance investigation toolkit that permits running MPI applications in a controlled environment with a large number of concurrent execution threads, while observing application performance and resilience in a simulated extreme-scale system. Using a lightweight parallel discrete event simulation, xSim executes an MPI application on a much smaller system in a highly oversubscribed fashion with a virtual wall clock time, such that performance data can be extracted based on a processor and a network model. xSim is designed like a traditional MPI performance tool, as an interposition library that sits between the MPI application and the MPI library, using the MPI profiling interface. In previous experiments, it has been run up to 134,217,728 communicating MPI ranks using a 960-core Linux cluster.

The simulator was deployed on the Linux cluster computer by associating one simulator MPI process per physical processor core. Within each simulator MPI process, a number of concurrent execution threads are executed, each representing an individual MPI process that is located on a processor core within a simulated HPC system. The execution timing of these simulated MPI processes is based on a processor model with a 1-to-1 performance match to the physical AMD processor core the simulator is running on and a network interconnect model with a basic star topology, 1 $\mu$s link latency, and infinite bandwidth. The model parameters were set to nominal values as the experiments investigate the number of Gossip cycles required to reach consensus.

## 3.1 Consensus using Global Knowledge

The first part of the experiments investigated the characteristics of the algorithm for consensus using global knowledge. The maximum number of Gossip cycles in the experiments was set to $5log(n)$, where $n$ is the number of simulated MPI processes. This is large enough to allow all the alive processes to reach consensus on the injected failures as the information dissemination speed of Gossip-based protocols is exponential. The Gossip cycle length for a given system size was set to allow the matrix merge operations to complete within the cycle plus the round trip time between the farthest processes in the system. This is necessary as matrix merge operations consume substantial amount of the cycle time. Experiments were carried out to test the algorithm's

scalability and its fault tolerance property.

Failures were injected into randomly chosen simulated MPI processes. In the first set of experiments, failures were injected right before the failure detection and consensus algorithm. In the second set of experiments, failures were injected during its execution to test the fault tolerance property of the algorithm. In all cases, the processes reach consensus on the injected failure(s) at different cycles. Hence, the cycle number of the last process reaching consensus is recorded.

Figure 4 shows the relation between the number of cycles taken to reach consensus and the system size for a single failure injected before the algorithm. It is evident that the number of cycles to reach consensus varies logarithmically with the system size. Figure 5 shows the exponential spreading of failure detection information at a particular process for the injected failure. Both figures demonstrate the logarithmic complexity of the algorithm. Figure 6 shows the distribution of the cycle number at which different processes reach consensus. In Figure 7, multiple (four) failures were injected before the algorithm and their effect on consensus time was observed. It took only 1 or 2 cycles more than in the single failure case (Figure 4). In Figure 8, multiple (four) failures were injected into randomly chosen processes at random cycles during the algorithm execution. The number of cycles needed to reach consensus increased slightly.
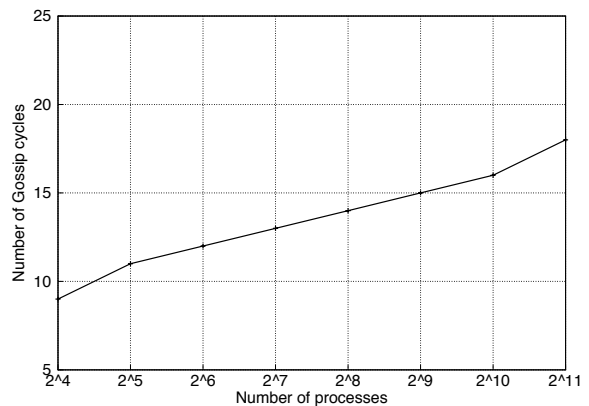


**Figure 4: Number of cycles to achieve global consensus after a single failure injection (algorithm 1)**

## 3.2 Efficient Heuristic Consensus

The second part of the experiments investigated the characteristics of the algorithm for consensus using the heuristic approach. The maximum number of Gossip cycles to run in the experiments was again set to $5log(n)$, where $n$ is the number of simulated MPI processes. It is large enough to allow all the alive processes to reach consensus on the injected failures. The length of each Gossip cycle was set to allow the merge operation to complete within the cycle plus the round trip time between the farthest processes in the system. The $MIN\_CCNT$ was set to 3.

Figure 9 shows the relation between the number of cycles taken to reach consensus on a single injected failure and the system size. It is evident that the number of cycles
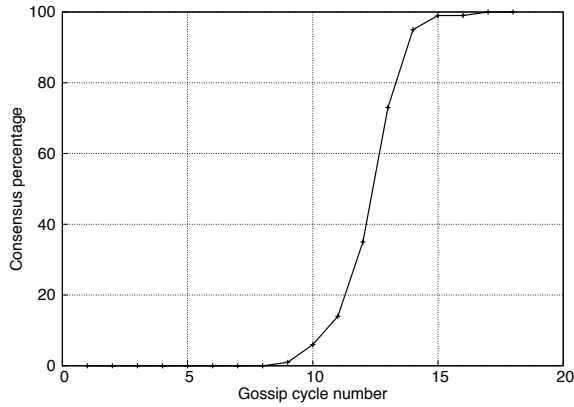
**Figure 5: Local consensus progress at a process after a single failure injection for system size of 2048 (algorithm 1)**
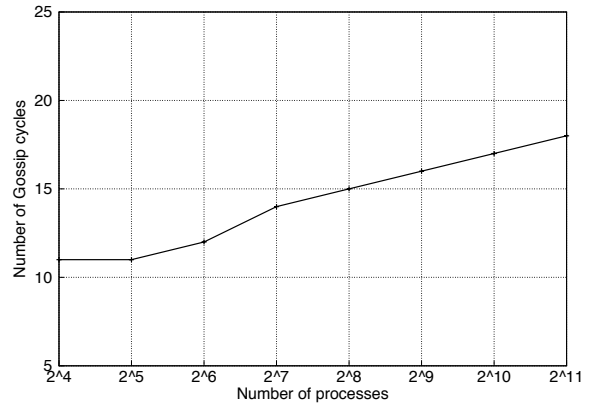


**Figure 7: Number of cycles to achieve global consensus after multiple (4) failures, which were injected before algorithm execution (algorithm 1)**
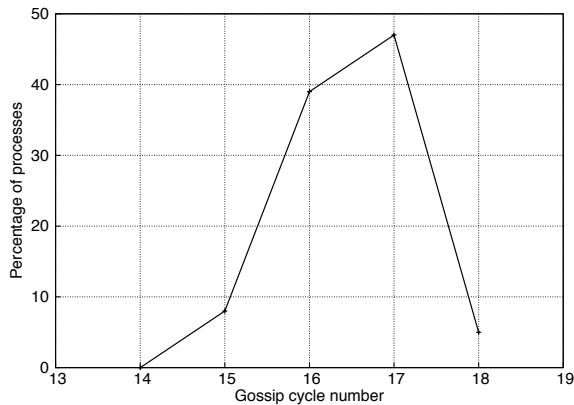


**Figure 6: Consensus detection spread for a system size of 2048 (algorithm 1)**
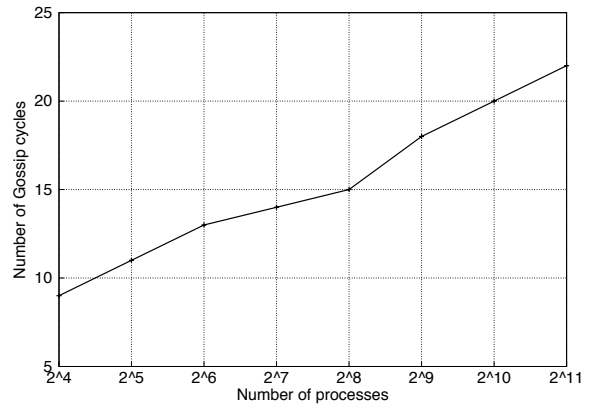


**Figure 8: Number of cycles to achieve global consensus with multiple (4) failures, which were injected during algorithm execution (algorithm 1)**

to reach consensus varies logarithmically with the system size as expected. Moreover, in this case no variance in the cycle number at which different processes reach consensus was observed, thus achieving perfect synchronization.

The algorithm is scalable in terms of memory as it needs to store only the list of failed processes at each process. The two algorithms were compared for their bandwidth utilization. The amount of data exchanged between simulated MPI processes is reported by xSim. Figure 10 shows the bandwidth consumed per process at increasing system sizes. It can be observed that the heuristic-based algorithm transfers a negligible amount of data and is significantly more efficient than the algorithm using global knowledge.

## 4. RELATED WORK

This section discusses failure detection and consensus algorithms that have been considered for HPC and compares them with the proposed algorithms.

### 4.1 Failure Detection Algorithms

A failure detection algorithm, for fail-stop type failures, using heartbeat messages was proposed in [16]. Every process maintains a log, called Gossip list, that contains a number (called heartbeat value) for each member process to represent its aliveness. Every $T_{Gossip}$ time units the process increments its own heartbeat value in the list and gossips the list to a randomly chosen process. When a process receives a remote Gossip list, it updates its Gossip list by setting the heartbeat value for each process to the maximum value in the two lists. The Gossip list is monitored continuously and if the heartbeat value for a process has not risen for $T_{cleanup}$ time units, it is suspected to have failed. There is, however, chance of an alive process to be falsely suspected as the algorithm is based on randomized Gossiping.

In the proposed algorithms, there is direct failure detection without passing through the suspicion phase of the above Gossip-based failure detection. Also, since suspicions
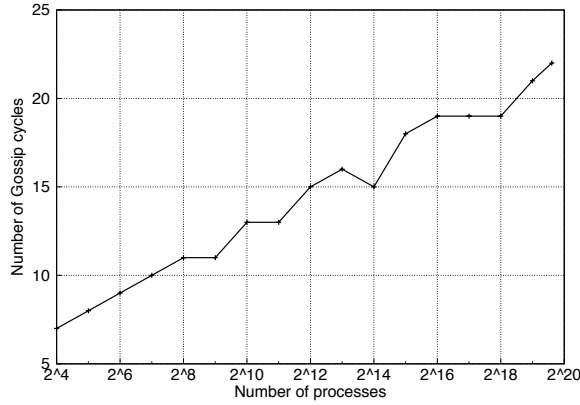
**Figure 9: Number of cycles to achieve global consensus after a single failure injection (algorithm 2)**
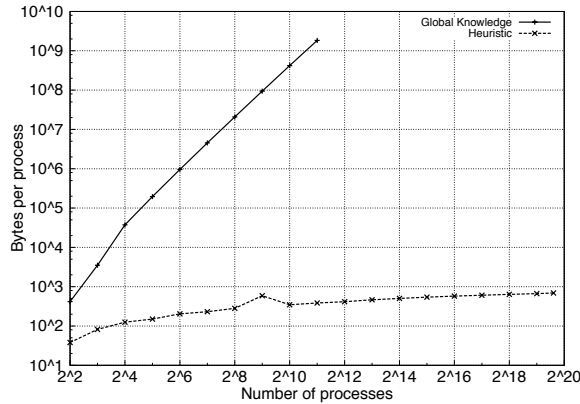


**Figure 10: Total bandwidth utilization of the consensus algorithms with a single failure injection**

become detections after a majority vote, the assumption that no more than one third of processes do not fail during a single Gossip iteration is relaxed in the proposed approach.

Another algorithm for failure detection was given in [8] and is based on randomized pinging. A process $p$ randomly pings another process $q$. $q$ is found to be alive if a reply is received in time. Otherwise the $p$ asks $k$ randomly chosen peers to ping $q$ as well. If no peer receives a reply, $p$ detects $q$ to have failed. This approach to failure detection takes network link failures into account along with process failures.

## 4.2  Consensus Algorithms

Both centralized, i.e., using a coordinator, and completely distributed consensus algorithms are available. Fault tolerant versions of the two-phase and three-phase consensus algorithms are discussed in the following. Distributed consensus algorithms based on Gossiping, which is inherently fault tolerant, are also discussed.

### 4.2.1  Coordinator-based Approaches

*Over a static tree structure.*

A two-phase consensus algorithm to aggregate locally known failures is given in [10]. Failures are presumed to be detected prior to invoking the algorithm. It is a fault tolerant log-scaling two-phase algorithm running over a tree topology. It is based on reliable gather and broadcast operations, both built on a multi-level, tree structure communication topology. The coordinator is at the root of the tree and makes the final decision based on the votes by the participants organized in the form of parent and children below it. Votes from the participants at the leaves and intermediate levels of the tree are gathered by the coordinator at the root going through the intermediate parent nodes and the decision from the coordinator is broadcast in the reverse direction. In the presence of faults, the parent of a failed participant recursively adopts its children. A child, upon detection of the failure of its parent, queries its grandparent to know how it should continue in the algorithm. If it has not voted yet, it will participate in the voting phase. If it has voted but the parent failed before propagating the vote up, it will cast the vote again; if the parent propagated the vote, it will participate in the decision phase. A termination detection algorithm is invoked upon failure of the coordinator. The parent directly below the coordinator decides the termination status: success with failed process list if it received decision from the coordinator, or abort if it has not voted yet. The participant, thus making the decision, propagates it below the tree. However, the termination detection will not be of any help if the coordinator failed without propagating its decision after the voting phase. The tree is rebalanced with the new alive processes between invocations of the algorithm and thus the approach involves maintenance of the communication structure, which is not needed by the algorithms proposed in this paper.

Failures that happen during the execution of the two-phase consensus are not included in the final list and a coordinator failure aborts the algorithm. Hence, it is not a completely fault tolerant consensus algorithm. In contrast, the approach presented in this paper is completely fault tolerant, as it allows for failures during the consensus.

*Over a dynamic tree structure.*

The approach in [3] assumes that a process knows of a few failed processes (a failure detector is again assumed to be readily available) to contribute towards the global list of failed processes. It is built around a reliable broadcast algorithm that dynamically constructs the broadcast tree and a three-phase consensus algorithm.

In the *BALLOTING* phase, the root generates a ballot (which is a sequence number to differentiate between iterations) and broadcasts it, including the known list of failed processes. The child, upon receipt of the broadcast message, checks whether it has any new failed processes known to it but not in the ballot it received. It sends a *REJECT* message piggybacked with the *ACK* message including the new failed processes if any; it accepts the ballot otherwise and sends an *ACCEPT* message piggybacked with the *ACK* message. The root starts the next phase if it receives *ACCEPT* message piggybacked with the *ACK* message from all its children and if any child rejects the ballot it updates its set of failed processes and tries again. In the second phase, the root broadcasts the *AGREE* message with the ballot. Now the participants know that the ballot has been agreed

upon by everyone and they agree to the ballot for the second phase. Then the root starts the third phase by broadcasting the *COMMIT* message. The participants upon receipt of the *COMMIT* message commit to the ballot.

A failure of the root is checked by every process and when a process detects that all processes with ranks lower than itself have failed it appoints itself as the new root. The new root restarts the algorithm from whatever state it is in. Note that if the root fails when it is in the *BALLOT-ING* state, the new root has to start all over again, wasting all the iterations performed so far. Failures of participants are handled by repeating phase one of the algorithm with a new sequence number included in the ballot. When a broadcast message belonging to an old iteration arrives at a process, which is not in the *BALLOTING* state, a *NAK* with *AGREE_FORCED* is forwarded to the root to clean up the old circulating broadcast messages. This message is also used by the root to start Phase 2 of the algorithm (with the assumption that process failures will subside and cease).

It was tested by injecting failures into randomly chosen processes before and during the execution. The processes that fail during the operation of the algorithm may or may not be included into the final list of failed processes. Moreover, every failure that happens while the algorithm is running requires the algorithm to start all over again by rebuilding the communication structure. The algorithms proposed in this paper do not require any communication structure.

### 4.2.2 Gossip-based Approaches

In Gossip-based approaches, failure detection is performed as explained in Section 2.1 and consensus is (in combination with failure detection) also implemented using Gossiping [16], and hence completely fault tolerant. Each process maintains a suspicion matrix $S$ to store the status of processes as detected by all the processes. An entry $S_{i,j}$ in a process's suspicion matrix indicates the status of process $j$ as detected by process $i$ (1 if suspected to have failed 0 otherwise). It also maintains a fault vector $F$ to store the processes' status as decided by all processes. An entry $F_i$ in this vector is 1 if the majority of the processes suspect process $i$ to have failed. The suspicion matrix is sent to other processes as part of Gossiping. When a process receives this suspicion matrix, it merges it with its own suspicion matrix. A process updates its fault vector by examining the suspicions of all processes. If the majority of the processes have suspected a process to have failed, it decides that the process has indeed failed and updates its fault vector to reflect this. Consensus is reached when all processes have detected a process to have failed. Upon detection of consensus a process broadcasts consensus message to all the live processes.

Because every process needs to maintain a suspicion matrix of $O(n^2)$, where $n$ is the system size, this algorithm is memory intensive and does not scale well. For increasing process numbers (beyond 48) the consensus time was found to increase exponentially. The first algorithm proposed in this paper also uses a matrix at each process, but the second algorithm uses only a list at each process. The algorithms given in this paper scale logarithmically with system size.

The proposed approach in this paper bypasses the failure suspicion phase based on distributed diagnosis. Since suspicions become detections after a majority vote, the assumption that no more than one third of processes do not fail during single Gossip iteration is relaxed in this approach.

Experiments in the state of the art HPC failure detection and consensus literature, have featured not more than a few thousand processes, whereas the proposed algorithm (algorithm 2) scaled to hundreds of thousand processes on a small cluster computer.

## 5. CONCLUSION AND FUTURE WORK

Failure detection and consensus for a fault-tolerant MPI enable HPC applications to adopt algorithm-based fault tolerance techniques to cope with MPI process failures more efficiently. Centralized methods for failure detection and consensus are based on a coordinator and do not scale well to very large and extreme-scale systems. Completely distributed algorithms based on Gossiping that were previously proposed in the literature consume an inordinate amount of time, memory and network bandwidth.

In this work two novel failure detection and consensus algorithms that use randomized pinging were presented. The first approach is based on global knowledge: each process maintains a local view of the entire system state to achieve consensus on failed processes. A Gossip protocol is used to detect failures and to exponentially propagate them in the system until the local views converge. The second algorithm does not rely on global knowledge and adopts a heuristic method to achieve consensus on failures.

Both algorithms were implemented as MPI applications and tested using the Extreme-scale Simulator. The results confirm their expected scalability and fault tolerance properties. In both algorithms, the number of Gossip cycles to achieve consensus on failures scales logarithmically with the system size. The second algorithm has significantly lower memory and bandwidth utilization and has shown to be able to achieve a perfect consensus synchronization as well.

The first algorithm can be implemented with boolean matrices at each process to increase scalability. The second algorithm's memory scalability can be further improved by maintaining the status of processes in a bit vector, if the number of failures in the system is high. It would also be interesting to investigate an efficient algorithm with a different heuristic approach for detecting consensus asynchronously and without the guaranteed initial propagation phase: this would allow detecting consensus on failures that happen both before and during the execution of the algorithm. Investigating processes entering alive state from faulty state is also interesting. This would avoid false positives and also allow process recovery.

Further future work in this area focuses on implementing *MPI_Comm_shrink()* and *MPI_Comm_agree()* with different approaches (static and dynamic tree, as well as, the different Gossip-based variants) and compare them using the Extreme-scale Simulator with architectural models of future-generation HPC systems.

## 6. ACKNOWLEDGEMENTS

## 7. REFERENCES

[1] R. C. Baumann. Radiation-induced soft errors in advanced semiconductor technologies. *IEEE Transactions on Device and Materials Reliability (TDMR)*, 5(3):305–316, 2005.

[2] W. Bland, G. Bosilca, A. Bouteiller, T. Herault, and J. Dongarra. A proposal for user-level failure mitigation in the mpi-3 standard. *Department of Electrical Engineering and Computer Science, University of Tennessee*, 2012.

[3] D. Buntinas. Scalable distributed consensus to support mpi fault tolerance. In *Parallel & Distributed Processing Symposium (IPDPS), 2012 IEEE 26th International*, pages 1240–1249. IEEE, 2012.

[4] T. D. Chandra, V. Hadzilacos, S. Toueg, and B. Charron-Bost. On the impossibility of group membership. In *Proceedings of the fifteenth annual ACM symposium on Principles of distributed computing*, pages 322–330. ACM, 1996.

[5] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM (JACM)*, 43(2):225–267, 1996.

[6] C. Engelmann. Scaling to a million cores and beyond: Using light-weight simulation to understand the challenges ahead on the road to exascale. *Future Generation Computer Systems (FGCS)*, 30(0):59–65, Jan. 2014.

[7] M. J. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM (JACM)*, 32(2):374–382, 1985.

[8] I. Gupta, T. D. Chandra, and G. S. Goldszmidt. On scalable and efficient distributed failure detectors. In *Proceedings of the twentieth annual ACM symposium on Principles of distributed computing*, pages 170–179. ACM, 2001.

[9] K.-H. Huang and J. A. Abraham. Algorithm-based fault tolerance for matrix operations. *IEEE Transactions on Computers (TC)*, C-33(6):518–528, 1984.

[10] J. Hursey, T. Naughton, G. Vallee, and R. L. Graham. A log-scaling fault tolerant agreement algorithm for a fault tolerant mpi. In *Recent Advances in the Message Passing Interface*, pages 255–263. Springer, 2011.

[11] L. Kaplan, P. Briggs, M. Ohlrich, and W. Leslie. Resilience to various failures for read-mostly in-memory data structures. In $26^{th}$ *IEEE International Parallel and Distributed Processing Symposium (IPDPS) 2012: Workshops*, pages 1572–1580, 2012.

[12] H. Kaul, M. Anders, S. Hsu, A. Agarwal, R. Krishnamurthy, and S. Borkar. Near-threshold voltage (NTV) design: Opportunities and challenges. In *Proceedings of the 49th Annual Design Automation Conference*, DAC '12, pages 1153–1158, New York, NY, USA, 2012. ACM.

[13] H. Ltaief, E. Gabriel, and M. Garbey. Fault tolerant algorithms for heat transfer problems. *Journal of Parallel and Distributed Computing (JPDC)*, 68(5):663–677, 2008.

[14] T. Naughton, C. Engelmann, G. Vallée, and S. Böhm. Supporting the development of resilient message passing applications using simulation. In *Proceedings of the $22^{nd}$ Euromicro International Conference on Parallel, Distributed, and network-based Processing (PDP) 2014*, pages 271–278, Turin, Italy, Feb. 12-14, 2014. IEEE Computer Society, Los Alamitos, CA, USA. Acceptance rate 32.6% (73/224).

[15] P. Poonpakdee, N. G. Orhon, and G. Di Fatta. Convergence detection in epidemic aggregation. In *Euro-Par 2013: Parallel Processing Workshops*, pages 292–300. Springer, 2014.

[16] S. Ranganathan, A. D. George, R. W. Todd, and M. C. Chidester. Gossip-style failure detection and distributed consensus for scalable heterogeneous clusters. *Cluster Computing*, 4(3):197–209, 2001.

[17] M. Snir, R. W. Wisniewski, J. A. Abraham, S. V. Adve, S. Bagchi, P. Balaji, J. Belak, P. Bose, F. Cappello, B. Carlson, A. A. Chien, P. Coteus, N. A. Debardeleben, P. Diniz, C. Engelmann, M. Erez, S. Fazzari, A. Geist, R. Gupta, F. Johnson, S. Krishnamoorthy, S. Leyffer, D. Liberty, S. Mitra, T. Munson, R. Schreiber, J. Stearley, and E. V. Hensbergen. Addressing failures in exascale computing. *International Journal of High Performance Computing Applications (IJHPCA)*, 28(2):127–171, May 2014.

[18] P. Soltero, P. Bridges, D. Arnold, and M. Lang. A gossip-based approach to exascale system services. In *Proceedings of the 3rd International Workshop on Runtime and Operating Systems for Supercomputers*, page 3. ACM, 2013.

[19] H. Straková, G. Niederbrucker, and W. N. Gansterer. Fault tolerance properties of gossip-based distributed orthogonal iteration methods. *Procedia Computer Science*, 18:189–198, 2013.

[20] J. Turek and D. Shasha. The many faces of consensus in distributed systems. *Computer*, 25(6):8–17, 1992.