

Algorithm-Based Checkpoint-Free Fault Tolerance for Parallel Matrix Computations on Volatile Resources *

Zizhong Chen

University of Tennessee, Knoxville
zchen@cs.utk.edu

Jack J. Dongarra

University of Tennessee, Knoxville
and
Oak Ridge National Laboratory
dongarra@cs.utk.edu

ABSTRACT

As the desire of scientists to perform ever larger computations drives the size of today's high performance computers from hundreds, to thousands, and even tens of thousands of processors, node failures in these computers are becoming frequent events. Although checkpoint/rollback-recovery is the typical technique to tolerate such failures, it often introduces a considerable overhead, especially when applications modify a large amount of memory between checkpoints.

This paper presents an algorithm-based checkpoint-free fault tolerance approach in which, instead of taking checkpoints periodically, a coded global consistent state of the critical application data is maintained in memory by modifying applications to operate on encoded data. Although the applicability of this approach is not so general as the typical checkpoint/rollback-recovery approach, in parallel linear algebra computations where it usually works, because no periodical checkpoint or rollback-recovery is involved in this approach, partial node failures can often be tolerated with a surprisingly low overhead.

We show the practicality of this technique by applying it to the ScaLAPACK matrix-matrix multiplication kernel which is one of the most important kernels for ScaLAPACK library to achieve high performance and scalability. We address the practical numerical issue in this technique by proposing a class of numerically good real number erasure codes based on random matrices. Experimental results demonstrate that the proposed checkpoint-free approach is able to survive process failures with a very low performance overhead.

*This research was supported in part by the Los Alamos National Laboratory under Contract No. 03891-001-99 49 and the Applied Mathematical Sciences Research Program of the Office of Mathematical, Information, and Computational Sciences, U.S. Department of Energy under contract DE-AC05-00OR22725 with UT-Battelle, LLC.

Keywords

Matrix Computations, ScaLAPACK, Checkpoint, Fault Tolerance, Parallel and Distributed Systems.

1. INTRODUCTION

As the number of processors in today's high performance computers continues to grow, the mean-time-to-failure (MTTF) of these computers are becoming significantly shorter than the execution time of many current high performance computing applications. Even making generous assumptions on the reliability of a single processor or link, it is clear that as the processor count in high end clusters grows into the tens of thousands, the MTTF of these clusters will drop from a few years to a few hours, or less. The next generation DOE ASCI computers (IBM Blue Gene L) are being designed with 131,000 processors [1]. The failure of some nodes or links for such a large system is likely to be just a few minutes away [16]. In recent years, the trend of the high performance computing has been shifting from the expensive massively parallel computer systems to the clusters of commodity off-the-shelf systems [8]. While the commodity off-the-shelf cluster systems have excellent price-performance ratios, due to the low reliability of the off-the-shelf components in these systems, there is a growing concern with the fault tolerance issues in such system. The recently emerging computational grid environments [20] with dynamic resources have further exacerbated the problem. However, driven by the desire of scientists for ever higher levels of detail and accuracy in their simulations, many computational science programs are now being designed to run for days or even months. Therefore, the next generation computational science programs must be able to tolerate hardware failures.

Today's long running scientific applications typically deal with faults by checkpoint/restart approaches in which all process states of an application are saved into stable storage periodically. The advantage of this approach is that it is able to tolerate the failure of the whole system. However, in this approach, if one process fails, usually all surviving processes are aborted and the whole application is restarted from the last checkpoint. The major source of overhead in all stable-storage-based checkpoint systems is the time it takes to write checkpoints into stable storage [23]. The checkpoint of an application on a, say, ten-thousand-processor computer implies that all critical data for the application on all ten thousand processors have to be written into stable storage periodically, which may introduce an unacceptable amount

of overhead into the checkpointing system. The restart of such an application implies that all processes have to be recreated and all data for each process have to be re-read from stable storage into memory or re-generated by computation, which often brings a large amount of overhead into restart. It may also be very expensive or unrealistic for many large systems such as grids to provide the large amount of stable storage necessary to hold all process state of an application of thousands of processes.

In order to tolerate partial failures with reduced overhead, diskless checkpointing [23] has been proposed by Plank et. al. By eliminating stable storage from checkpointing and replacing it with memory and processor redundancy, diskless checkpointing removes the main source of overhead in checkpointing [23]. Diskless checkpointing has been shown to achieve a decent performance to tolerate single process failure in [21]. For applications which modify a small amount of memory between checkpoints, it is shown in [7] that even to tolerate multiple simultaneous process failures, the overhead introduced by diskless checkpointing is still negligible.

However, for applications, such as matrix-matrix multiplication, which modify a large amount of memory between checkpoints, due to the large checkpoint size, even diskless checkpointing still introduces a considerable overhead into applications. Firstly, a local in memory checkpoint has to be maintained in diskless checkpointing, which introduces a large amount of memory overhead and hurts the efficiency of applications. Secondly, the local checkpoint in diskless checkpointing has to be taken and encoded periodically, which introduce a considerable performance overhead into applications. Despite the checksum and reverse computation technique in [21] has reduced the memory overhead, the overhead to calculate the checkpoint encodings periodically does not change. Furthermore after failures, this technique increase the recovery overhead by reversing the computation.

In this paper, we present an algorithm-based checkpoint-free fault tolerance approach in which, instead of taking checkpoints periodically, a coded global consistent state of the critical application data is maintained in memory by modifying applications to operate on encoded data. Although this approach is not as generally applicable as typical checkpoint approaches, in parallel matrix computations where it usually works, because no periodical checkpoint and rollback-recovery are involved in this approach, fault tolerance for partial node failures can often be achieved with a surprisingly low overhead.

Despite the fact that there has been much research on algorithm-based fault tolerance [19] in which applications are modified to operate on encoded data to determine the correctness of some mathematical calculations on parallel platforms where failed processors produce incorrect calculations, to the best of our knowledge, this is the first time that applications are modified to operate on encoded data to maintain a global consistent state on parallel and distributed systems where failed processors stop working.

We show the practicality of this technique by applying it to the ScaLAPACK [2] matrix-matrix multiplication kernel which is one of the most important kernels for ScaLAPACK library to achieve high performance and scalability. We address the practical numerical issue in this technique by proposing a class of numerically good real number erasure codes based on random matrices. Experimental results

for matrix-matrix multiplication demonstrate that the proposed approach is able to survive a small number of process failures with a very low performance overhead.

Although the algorithm-based checkpoint-free fault tolerance approach presented in this paper is non-transparent and algorithm-dependent, it is meaningful in that

1. It can often achieve a surprisingly low overhead in parallel matrix computations where it usually works.
2. It is often possible to build it into frequently used numerical libraries such as ScaLAPACK to relieve the involvement of the application programmer.

The specific contributions this paper makes can be summarized as following

- **Algorithm-Based Checkpoint-Free Fault Tolerance:** We present an algorithm-based checkpoint-free fault tolerance approach in which, instead of taking checkpoint periodically, a coded global consistent state of the critical application data is maintained in memory by modifying applications to operate on encoded data. We show the practicality of this technique by applying it to the ScaLAPACK matrix-matrix multiplication kernel which is one of the most important kernels for ScaLAPACK to achieve high performance and scalability.
- **Numerically Good Real Number Erasure Codes:** We present a class of numerically good real number erasure codes based on random matrices which can be used to algorithm-based checkpoint-free fault tolerance technique to tolerate multiple simultaneous process failures. We prove our codes are numerically highly reliable.

The rest of this paper is organized as follows. Section 2 gives a brief review of the related work. Section 3 specifies the type of failures we focus on. Section 4 presents the basic idea of algorithm-based checkpoint-free fault tolerance. In Section 5, we present an example to demonstrate how algorithm-based checkpoint-free fault tolerance works in practice by applying this technique to the ScaLAPACK matrix-matrix multiplication kernel. In Section 6, we address the practical numerical issue in this technique by proposing a class of numerically good real number erasure codes. In Section 7, we evaluate the performance overhead of applying this technique to the ScaLAPACK matrix-matrix multiplication kernel. Section 8 compares algorithm-based checkpoint-free fault tolerance with existing works and discusses the limitations of this technique. Section 9 concludes the paper and discusses future work.

2. RELATED WORK

There is a rich literature on algorithm-based fault tolerance and fault tolerant computing in parallel and distributed systems. This section briefly reviews previous fault tolerance work related to our work.

2.1 Checkpoint/Restart

Checkpoint/restart is probably the most typical approach to tolerate failures in parallel and distributed systems. By writing checkpoints into stable storage periodically, the checkpoint/restart approach is able to tolerate the failure of the

whole system. Checkpoint could be performed either from a system level or from an application level. There is often a trade off between transparency and performance. Transparent system level approach usually introduces a higher performance overhead than non-transparent application level approach. The recent compiler based approach from [5] is a very promising approach in this class of approaches.

The advantage of this class of approaches is that it is very general and is able to tolerate the failure of the whole system.

The limitations of this class of approaches are that it generally needs stable storage to save a global consistent state periodically and that it aborts all survival processes even if only one of many processes failed.

2.2 Diskless Checkpointing

Diskless checkpointing [23] is a technique for checkpointing the state of a long running application on a distributed system without relying on stable storage. By eliminating stable storage from checkpointing and replacing it with memory and processor redundancy, diskless checkpointing removes the main source of overhead in checkpointing.

The advantages of this class of approaches are that it tolerates partial failures with reduced overhead and it does not rely on stable storage.

The limitation of this class of approaches is that it is not able to tolerate the failure of the whole system. And there is additional memory, processor, and network overhead that is absent in typical checkpoint/restart approach.

2.3 Checkpoint-Free Fault Tolerance

This class of approaches considers the specific characteristic of an application and designs fault tolerance schemes according to the specific characteristic of an application. In [16], Geist et. al. investigated the natural fault tolerance concept in which the application can finish the computation task even if a small amount of application data is lost, therefore, checkpoint can be avoided for this type of applications. In [4], a checkpoint-free scheme is given for iterative method. In [18], a checkpoint-free scheme is incorporated into a parallel direct search application.

The advantage of this class of approaches is that it is able to achieve very low overhead according to the specific characteristic of an application.

The limitation of this approach is that it is non-transparent and has to be designed according to the specific characteristic of an application

2.4 Algorithm based fault tolerance

Algorithm based fault tolerance [19] is a class of approaches which tolerant byzantine failures, in which failed processors continues to work but produce incorrect calculations. In this approach, applications are modified to operate on encoded data to determine the correctness of some mathematical calculations. This class of approaches can mainly be applied to applications performing linear algebra computations and usually achieves a very low overhead.

One of the most important characteristics of this research is that it assume a fail-continue model in which failed processors continues to work but produce incorrect calculations.

3. FAILURE MODEL

To define the problem we are targeting and clarify the differences with traditional algorithm-based fault tolerance, in this section, we specify the type of failures we are focusing on.

Assume the computing system consists of many nodes connected by network connections. Each node has its own memory and local disk. The communication between processes are assumed to be message passing. Assume the target application is optimized to run on a fixed number of processes.

We assume nodes in the computing system are volatile, which means a node may leave the computing system due to failure, or join the computing system after being repaired. Unlike in traditional algorithm-based fault tolerance which assumes a failed processor continues to work but produce incorrect results, in this paper, we assume a *fail-stop* failure model. That is the failure of a node will cause all processes on the failed nodes stop working. All data of the processes on the failed node is lost. The processes on survival nodes could not either send or receive any message from the processes on the failed node. Although there are many other type of failures exist, in this paper, we only consider this type of failures. This type of failure is common in today's large computing systems such as high-end clusters with thousands of nodes and computational grids with dynamic resources.

4. ALGORITHM-BASED CHECKPOINT-FREE FAULT TOLERANCE

In this section, we present the basic idea of algorithm-based checkpoint-free fault tolerance. We restrict our scope to the long running numerical computing applications only. As indicated in Section 5, this approach can mainly be applied to linear algebra computations on parallel and distributed systems.

4.1 Failure Detection and Location

It is assumed that fail-stop failures can be detected and located with the aid of the programming environment. Many current programming environments such as PVM [25], Globus [13], FT-MPI [10], and Open MPI [15] do provide this kind of failure detection and location capability. We assume the lost of partial processes in the message passing system does not cause the aborting of the survival processes and it is possible to replace the failed processes in the message passing system and continue the communication after the replacement. FT-MPI [10] is one such programming environments that support all these functionalities. In the rest of this section, we will mainly focus on how to recover the application.

4.2 Single Failure Recovery

Today's long running scientific programs typically deal with faults by checkpoint and rollback recovery in which all process states of an application are saved into certain storage periodically. If one process fails, the data on *all* processes has to be recovered from the last checkpoint. The checkpoint and rollback of an application on a, say, ten-thousand-processor computer implies that all critical data for the application on all ten thousand processors have to be saved into and recovered from some storage periodically, which may introduce an unacceptable amount of overhead

(both time and storage) into the checkpointing system. Considering that all data on all survival processes are still effective, it is interesting to ask: *is it possible to recover only the lost data on the failed process?*

Consider the simple case where there will be only one process failure. Before the failure actually occurs, we do not know which process will fail, therefore, a scheme to recover only the lost data on the failed process actually need to be able to recover data on *any* process. It seems difficult to be able to recover data on any process without saving all data on all processes somewhere. However, if we assume, at any time during the computation, the data on the i^{th} process P_i satisfies

$$P_1 + P_2 + \dots + P_{n-1} = P_n, \quad (1)$$

where n is the total number of process used for the computation. Then the lost data on *any* failed process would be able to be recovered from (1). Assume the j^{th} process failed, then the lost data P_j can be recovered from

$$P_j = P_n - (P_1 + \dots + P_{j-1} + P_{j+1} + \dots + P_{n-1})$$

In this very special case, we are lucky enough to be able to recover the lost data on *any* failed process without checkpoint due to the special *checksum relationship* (1). In practice, this kind of special relationship is by no means natural. However, it is natural to ask: *is it possible to design an application to maintain such a special checksum relationship on purpose?*

Assume the original application is designed to run on n processes. Let P_i denotes the data on the i^{th} computation process. The special checksum relationship above can actually be designed on purpose as follows

- Add another encoding process into the application. Assume the data on this encoding process is C . For numerical computations, P_i is often an array of floating-point numbers, therefore, at the beginning of the computation, we can create a checksum relationship among the data of all processes by initializing the data C on the encoding process as

$$P_1 + P_2 + \dots + P_n = C \quad (2)$$

- During the executing of the application, redesign the algorithm to operate both on the data of computation processes and on the data of encoding process in such a way that the checksum relationship (2) is always maintained.

The specially designed checksum relationship (2) actually establishes an equality between the data P_i on computation processes and the encoding data C on the encoding process. If any processor fails then the equality (2) becomes an equation with one unknown. Therefore, the data in the failed processor can be reconstructed through solving this equation.

4.3 Multiple Failure Recovery

The specially designed checksum relationship in the last sub-section can only survive one process failure. However, in today's high performance computers, there are usually more than one processors on each node. Hence, it is usual to run multiple processes on one node, which implies that the failure of one node often causes the lost of multiple processes. Furthermore, as the number of nodes increases, the

possibility of lost multiple nodes also increases. Therefore, it is often necessary to be able to survive multiple simultaneous process failures. In this section, we present a scheme which can be used to recover multiple simultaneous process failures.

Suppose there are n processes used for computation. Assume the data on the i -th computation process is P_i . In order to be able to reconstruct the lost data on m failed processes, another m processes are dedicated to hold m encodings (weighted checksums) of the computation data. At the beginning of the application, the weighted checksum C_j on the j th encoding process can be calculated from

$$\begin{cases} a_{11}P_1 + \dots + a_{1n}P_n & = C_1 \\ & \vdots \\ a_{m1}P_1 + \dots + a_{mn}P_n & = C_m, \end{cases} \quad (3)$$

where a_{ij} , $i = 1, 2, \dots, m$, $j = 1, 2, \dots, n$, is the weight we need to choose. During the executing of the application, the application need to be re-designed to operate both on the data of computation processes and on the data of encoding processes in such a way that the relationship (3) is always maintained.

We call the relationship (3) the *weighted checksum relationship*. We call $A = (a_{ij})_{mn}$ the *encoding matrix* for the weighted checksum relationship. The specially designed weighted checksum relationship (3) actually establishes m equalities between the data P_i on computation processes and the encoding data C_i on the encoding processes. If some processes fail, then the m equalities become a system of linear equations. Therefore, the lost data in the failed processes may be able to be reconstructed through solving the system of linear equations.

Suppose that k computation processes and $m - h$ encoding processes have failed, then there are $n - k$ computation processors and h encoding processes survive. If we look at the data on failed processors as unknowns, then (3) becomes m equations with $m - (h - k)$ unknowns.

If $k > h$, then there are less equations than unknowns. There is no unique solution for (3). The lost data on the failed processes can not be recovered.

However, if $k < h$, then there are more equations than unknowns. By appropriately choosing A , a unique solution for (3) can be guaranteed. Therefore, the lost data on the failed processes can be recovered by solving (3).

Without loss of generality, we assume: (1) the computational process j_1, j_2, \dots, j_k failed and the computational process $j_{k+1}, j_{k+2}, \dots, j_n$ survived; (2) the encoding process i_1, i_2, \dots, i_h survived and the encoding process $i_{h+1}, i_{h+2}, \dots, i_m$ failed. Then, in equation (3), P_{j_1}, \dots, P_{j_k} and $C_{i_{h+1}}, \dots, C_{i_m}$ become unknowns after the failure occurs. If we re-structure (3), we can get

$$\begin{cases} a_{i_1 j_1} P_{j_1} + \dots + a_{i_1 j_k} P_{j_k} & = C_{i_1} - \sum_{t=k+1}^n a_{i_1 j_t} P_{j_t} \\ & \vdots \\ a_{i_h j_1} P_{j_1} + \dots + a_{i_h j_k} P_{j_k} & = C_{i_h} - \sum_{t=k+1}^n a_{i_h j_t} P_{j_t} \end{cases} \quad (4)$$

and

$$\begin{cases} C_{i_{h+1}} & = a_{i_{h+1} 1} P_1 + \dots + a_{i_{h+1} n} P_n \\ & \vdots \\ C_{i_m} & = a_{i_m 1} P_1 + \dots + a_{i_m n} P_n. \end{cases} \quad (5)$$

Let A_r denote the coefficient matrix of the linear system (4). If A_r has full column rank, then P_{j_1}, \dots, P_{j_k} can be recovered by solving (4), and $C_{i_{h+1}}, \dots, C_{i_m}$ can be recovered by substituting P_{j_1}, \dots, P_{j_k} into (5).

Whether we can recover the lost data on the failed processes or not directly depends on whether A_r has full column rank or not. However, A_r in (4) can be any sub-matrix (including minor) of A depending on the distribution of the failed processors. If any square sub-matrix (including minor) of A is non-singular and there are no more than m process failed, then A_r can be guaranteed to have full column rank. Therefore, to be able to recover from any no more than m failures, the encoding matrix A has to satisfy: *any square sub-matrix (including minor) of A is non-singular.*

How can we find such kind of matrices? It is well known that some structured matrices such as Vandermonde matrix, Cauchy matrix, and DFT matrix satisfy any square sub-matrix (including minor) of the matrix is non-singular.

5. CHECKPOINT-FREE FAULT TOLERANCE FOR MATRIX MULTIPLICATION

As an example to demonstrate how the algorithm-based checkpoint-free fault tolerance works in practice, in this section, we apply this technique to the ScaLAPACK matrix-matrix multiplication kernel which is one of the most important kernels for ScaLAPACK to achieve high performance and scalability.

Actually, it is also possible to incorporate fault tolerance into many other ScaLAPACK routines through this approach. However, in this section, we will restrict our presentation to the matrix-matrix multiplication kernel. For the simplicity of presentation, in this section, we only discuss the case where there is only one process failure. However, just as described in the last section, it is straight to extend the result here to the multiple simultaneous process failures case by simply using the weighted checksum scheme in Section 4.3.

5.1 Two-Dimensional Block-Cyclic Distribution

It is well-known [2] that the layout of an application's data within the hierarchical memory of a concurrent computer is critical in determining the performance and scalability of the parallel code. By using two-dimensional block-cyclic data distribution [2], ScaLAPACK seeks to maintain load balance and reduce the frequency with which data must be transferred between processes.

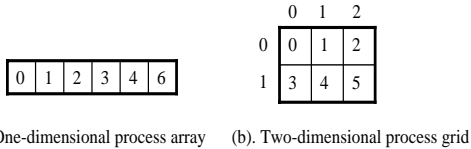


Figure 1: Process grid in ScaLAPACK

For reasons described above, ScaLAPACK organizes the one-dimensional process array representation of an abstract parallel computer into a two-dimensional rectangular process grid. Therefore, a process in ScaLAPACK can be referenced by its row and column coordinates within the grid. An example of such an organization is shown in Figure 1.

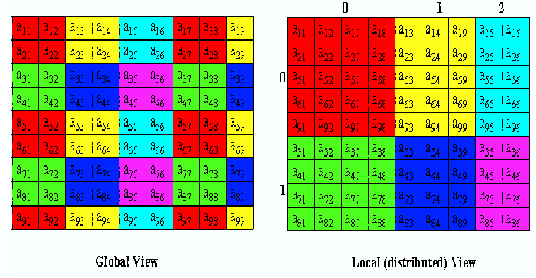


Figure 2: Two-dimensional block-cyclic matrix distribution

The two-dimensional block-cyclic data distribution scheme is a mapping of the global matrix onto the rectangular process grid. There are two pairs of parameters associated with the mapping. The first pair of parameters is (mb, nb) , where mb is the row block size and nb is the column block size. The second pair of parameters is (P, Q) , where P is the number of process rows in the process grid and Q is the number of process columns in the process grid. Given an element a_{ij} in the global matrix A , the process coordinate (p_i, q_i) that a_{ij} resides can be calculated by

$$\begin{cases} p_i = \lfloor \frac{i}{mb} \rfloor \bmod P, \\ q_i = \lfloor \frac{i}{nb} \rfloor \bmod Q, \end{cases}$$

The local coordinate (i_{p_i}, j_{q_i}) which a_{ij} resides in the process (p_i, q_i) can be calculated according to the following formula

$$\begin{cases} i_{p_i} = \lfloor \frac{i}{P} \rfloor \cdot mb + i \bmod mb, \\ j_{q_i} = \lfloor \frac{i}{Q} \rfloor \cdot nb + i \bmod nb, \end{cases}$$

Figure 2 is an example of mapping a 9 by 9 matrix onto a 2 by 3 process grid according two-dimensional block-cyclic data distribution with $mb = nb = 2$.

5.2 Encoding Two-Dimensional Block Cyclic Matrices

In this section, we will construct different encoding schemes which can be used to design checkpoint-free fault tolerant matrix computation algorithms in ScaLAPACK.

Assume a matrix M is originally distributed in a P by Q process grid according to the two dimensional block cyclic data distribution. For the convenience of presentation, assume the size of the local matrices in each process is the same. We will explain different coding schemes for the matrix M with the help of the example matrix in Figure 3. Figure 3 (a) shows the global view of an example matrix. After the matrix is mapped onto a 2 by 2 process grid with $mb = nb = 1$, the distributed view of this matrix is shown in Figure 3 (b).

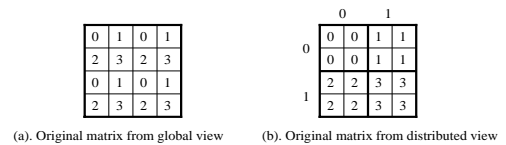


Figure 3: An example matrix

Suppose we want to tolerate a single process failure. We dedicate another $P+Q+1$ additional processes and organize the total $PQ+P+Q+1$ process as a $P+1$ by $Q+1$ process grid with the original matrix M distributed onto the first P rows and Q columns of the process grid.

The *distributed column checksum matrix* M^c of the matrix M is the original matrix M plus the part of data on the $(P+1)^{th}$ process row which can be obtained by adding all local matrices on the first P process rows. Figure 4 (b) shows the distributed view of the column checksum matrix of the example matrix from Figure 1. Figure 4 (a) is the global view of the column checksum matrix.

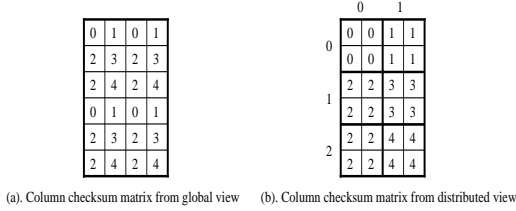


Figure 4: Distributed column checksum matrix of the example matrix

The *distributed row checksum matrix* M^r of the matrix M is the original matrix M plus the part of data on the $(Q+1)^{th}$ process columns which can be obtained by adding all local matrices on the first Q process columns. Figure 5 (b) shows the distributed view of the row checksum matrix of the example matrix from Figure 1. Figure 5 (a) is the global view of the row checksum matrix.

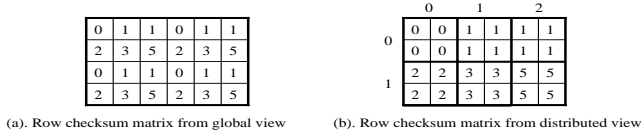


Figure 5: Distributed row checksum matrix of the original matrix

The *distributed full checksum matrix* M^f of the matrix M is the original matrix M , plus the part of data on the $(P+1)^{th}$ process row which can be obtained by adding all local matrices on the first P process rows, plus the part of data on the $(Q+1)^{th}$ process column which can be obtained by adding all local matrices on the first Q process columns. Figure 6 (b) shows the distributed view of the full checksum matrix of the example matrix from Figure 3. Figure 6 (a) is the global view of the full checksum matrix.

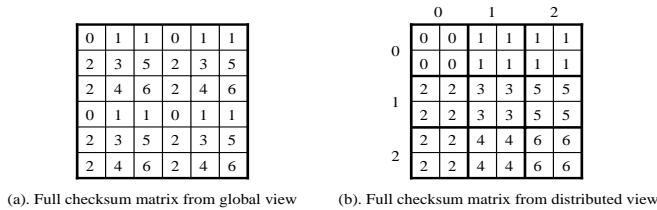


Figure 6: Distributed full checksum matrix of the original matrix

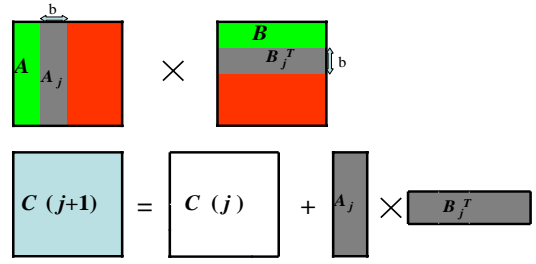


Figure 7: The j^{th} step of the matrix-matrix multiplication algorithm in ScaLAPACK

5.3 Scalable Universal Matrix Multiplication Algorithm

To achieve high performance, the matrix-matrix multiplication in ScaLAPACK uses a blocked outer product version of the matrix matrix multiplication algorithm. Let A_j denote the j^{th} column block of the matrix A and B_j^T denote the j^{th} row block of the matrix B . The following Figure 8 is the algorithm to perform the matrix matrix multiplication. Figure 7 shows the j^{th} step of the matrix matrix multiplication algorithm.

```

for j = 0, 1, ...
  row broadcast A_j;
  column broadcast B_j^T;
  C = C + A_j * B_j^T;
end

```

Figure 8: Scalable Universal Matrix-Matrix Multiplication Algorithm in ScaLAPACK

5.4 Maintaining Global Consistent States by Computation

Assume A , B and C are distributed matrices on a P by Q process grid with the first element of each matrix on process $(0,0)$. Let A^c , B^r and C^f denote the corresponding distributed checksum matrix. Let A_j^c denote the j^{th} column block of the matrix A^c and B_j^{rT} denote the j^{th} row block of the matrix B^r . We first prove the following fundamental theorem for matrix matrix multiplication with checksum matrices.

THEOREM 1. *Let $S_j = C^f + \sum_{k=0}^{j-1} A_k^c * B_k^{rT}$, then S_j is a distributed full checksum matrix.*

PROOF. It is straightforward that $A_k^c * B_k^{rT}$ is a distributed full checksum matrix and the sum of two distributed full checksum matrices is a distributed checksum matrix. S_j is the sum of j distributed full checksum matrices, therefore is a distributed full checksum matrix. \square

Theorem 1 tells us that at the end of each iteration of the matrix matrix multiplication algorithm with checksum matrices, the checksum relationship of all checksum matrices are still maintained. This tells us that a coded global consistent state of the critical application data is maintained in memory at the end of each iteration of the matrix matrix multiplication algorithm if we perform the computation with related checksum matrices.

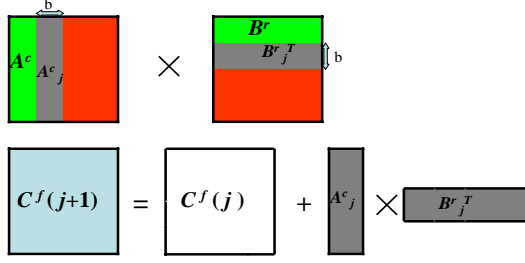


Figure 9: The j^{th} step of the fault tolerant matrix-matrix multiplication algorithm

However, in a distributed environment, different process may update their local data asynchronously. Therefore, if when some process has updated their local matrix and some process is still in the communication stage, a failure happens, then the relationship of the data in the distributed matrix will no be maintained and the data on all processes would not form a consistent state. But this could be solved by simply performing a synchronization before performing local update. Therefore, in the following algorithm in Figure 10, there will always be a coded global consistent state (i.e. the checksum relationship) of the matrix A^c , B^r and C^f in memory. Hence, a single process failure at any time during the matrix matrix multiplication would be able to recovered from the checksum relationship.

Despite in this algorithm, the only modification to the library routine is to perform a synchronization before local update. However the amount of modification necessary to maintain a consistent state is highly dependent on the characteristic of an algorithm. For example, in LU factorization, due to the damage of the linear relationship by the global row pivoting, one also needs to adjust the encodings appropriately when performing pivoting to maintain a consistent encoded state in memory.

```

construct checksum matrices  $A^c, B^r$ , and  $C^f$ ;
for  $j = 0, 1, \dots$ 
    row broadcast  $A_j^c$ ;
    column broadcast  $B_j^{rT}$ ;
    synchronize;
     $C^f = C^f + A_j^c * B_j^{rT}$ ;
end

```

Figure 10: A fault tolerant matrix-matrix multiplication algorithm

5.5 Overhead and Scalability Analysis

In this section, we analysis the overhead introduced by the algorithm-based checkpoint-free fault tolerance for matrix matrix multiplication.

For the simplicity of presentation, we assume all three matrices A, B , and C are square. Assume all three matrices are distributed onto a P by P process grid with m by m local matrices on each process. The size of the global matrices is Pm by Pm . Assume all elements in matrices are 8-byte double precision floating-point numbers. Assume every process has the same speed and disjoint pairs of processes can communicate without interfering each other. Assume it takes

$\alpha + \beta k$ seconds to transfer a message of k bytes regardless which processes are involved, where α is the latency of the communication and $\frac{1}{\beta}$ is the bandwidth of the communication. Assume a process can concurrently send a message to one partner and receive a message from a possibly different partner. Let γ denote the time it takes for a process to perform one floating-point arithmetic operation.

5.5.1 Time Complexity for Parallel Matrix Matrix Multiplication

Note that the sizes of all three global matrices A, B , and C are all Pm , therefore, the total number of floating-point arithmetic operations in the matrix matrix multiplication is $2P^3m^3$. There are P^2 process with each process execute the same number of floating-point arithmetic operations. Hence, the total number of floating-point arithmetic operations on each process is $2Pm^3$. Therefore, the time T_{matrix_comp} for the computation in matrix matrix multiplication is

$$T_{matrix_comp} = 2Pm^3\gamma.$$

In the parallel matrix matrix multiplication algorithm in Figure 8, the columns of A and the rows of B also need to broadcast to other column and row processes respectively. To broadcast one block columns of A using a simple *binary tree* broadcast algorithm, it takes $2(\alpha + 8bm\beta)\log_2 P$, where b is the row block size in the two dimensional block cyclic distribution. Therefore, the time T_{matrix_comm} for the communication in matrix matrix multiplication is

$$T_{matrix_comm} = 2\alpha\frac{Pm}{b}\log_2 P + 16\beta Pm^2\log_2 P.$$

Therefore, the total time to perform parallel matrix matrix multiplication is

$$\begin{aligned} T_{matrix_mult} &= T_{matrix_comp} + T_{matrix_comm} \\ &= 2Pm^3\gamma + 2\alpha\frac{Pm}{b}\log_2 P \\ &\quad + 16\beta Pm^2\log_2 P. \end{aligned} \quad (6)$$

5.5.2 Overhead for Calculating Encoding

To make matrix matrix multiplication fault tolerant, the first type of overhead introduced by the algorithm-based checkpoint-free fault tolerance technique is (1) constructing the distributed column checksum matrix A^c from A ; (2) constructing the distributed row checksum matrix B^r from B ; (3) constructing the distributed full checksum matrix C^f from C ;

The distributed checksum operation involved in constructing all these checksum matrices performs the summation of P local matrices from P processes and saves the result into the $(P + 1)^{\text{th}}$ process. Let T_{each_encode} denote the time for one checksum operation and T_{total_encode} denote the time for constructing all three checksum matrices A^c, B^r , and C^f , then

$$T_{total_encode} = 4T_{each_encode}$$

By using a *fractional tree* reduce style algorithm [24], the time complexity for one checksum operation can be expressed as

$$\begin{aligned} T_{each_encode} &= 8m^2\beta \left(1 + O\left(\left(\frac{\log_2 P}{m^2}\right)^{1/3}\right) \right) \\ &\quad + O(\alpha\log_2 P) + O(m^2\gamma) \end{aligned}$$

Therefore, the time complexity for constructing all three checksum matrices is

$$T_{total_encode} = 32m^2\beta \left(1 + O\left(\left(\frac{\log_2 P}{m^2}\right)^{1/3}\right) \right) + O(\alpha \log_2 P) + O(m^2\gamma). \quad (7)$$

In practice, unless the size of the local matrices m is very small or the size of the process grid P is extremely large, the total time for constructing all three checksum matrices is almost independent of the size of the process grid P .

The overhead (%) R_{total_encode} for constructing checksum matrices for matrix matrix multiplication is

$$R_{total_encode} = \frac{T_{total_encode}}{T_{matrix_mult}} = O\left(\frac{1}{Pm}\right) \quad (8)$$

From (8), we can conclude

1. If the size of the data on each process is fixed (m is fixed), then as the number of processes increases to infinite (that is $P \rightarrow \infty$), the overhead (%) for constructing the checksum matrices decreases to zero with a speed of $O(\frac{1}{P})$
2. If the number of processes is fixed (P is fixed), then as the size of the data on each process increases to infinite (that is $m \rightarrow \infty$) the overhead (%) for constructing the checksum matrices decreases to zero with a speed of $O(\frac{1}{m})$

5.5.3 Overhead for Performing Computations on Encoded Matrices

The fault tolerant matrix matrix multiplication algorithm in Figure 10 performs computations using checksum matrices which have larger size than the original matrices. However, the total number of processes devoted to computation also increases. A more careful analysis of the algorithm in Figure 10 indicates that the number of floating-point arithmetic operations on each process in the fault tolerant algorithm (Figure 10) is actually the same as that of the original non-fault tolerant algorithm (Figure 8).

As far as the communication is concerned, in the original algorithm (in Figure 8), the column (and row) blocks are broadcast to P processes. In the fault tolerant algorithms (in Figure 10), the column (and row) blocks now have to be broadcast to $P + 1$ processes.

Therefore, the total time to perform matrix matrix multiplication with checksum matrices is

$$T_{matrix_mult_checksum} = 2Pm^3\gamma + 2\alpha\frac{Pm}{b}\log_2(P+1) + 16\beta Pm^2\log_2(P+1).$$

Therefore, the overhead (time) to perform computations with checksum matrices is

$$T_{overhead_matrix_mult} = T_{matrix_mult_checksum} - T_{matrix_mult} = \left(2\alpha\frac{Pm}{b} + 16\beta Pm^2\right)\log_2\left(1 + \frac{1}{P}\right). \quad (9)$$

The overhead (%) $R_{overhead_matrix_mult}$ for performing computations with checksum matrices in fault tolerant matrix

matrix multiplication is

$$R_{overhead_matrix_mult} = \frac{T_{overhead_matrix_mult}}{T_{matrix_mult}} = O\left(\frac{1}{Pm}\right) \quad (10)$$

From (9), we can conclude that

1. If the size of the data on each process is fixed (m is fixed), then as the number of processes increases to infinite (that is $P \rightarrow \infty$), the overhead (%) for performing computations with checksum matrices decreases to zero with a speed of $O(\frac{1}{P})$
2. If the number of processes is fixed (P is fixed), then as the size of the data on each process increases to infinite (that is $m \rightarrow \infty$) the overhead (%) for performing computations with checksum matrices decrease to zero with a speed of $O(\frac{1}{m})$

5.5.4 Overhead for Recovery

The failure recovery contains two steps: (1) recover the programming environment; (2) recover the application data.

The overhead for recovering the programming environment depends on the specific programming environment. For FT-MPI [10] which we perform all our experiment on, it introduce a negligible overhead (refer Section 7).

The procedure to recover the three matrices A , B , and C is similar to calculating the checksum matrices. Except for matrix C , it can be recovered from either the row checksum or the column checksum relationship. Therefore, the overhead to recover data is

$$T_{recover_data} = 24m^2\beta \left(1 + O\left(\left(\frac{\log_2 P}{m^2}\right)^{1/3}\right) \right) + O(\alpha \log_2 P) + O(m^2\gamma) \quad (11)$$

In practice, unless the size of the local matrices m is very small or the size of the process grid P is extremely large, the total time for recover all three checksum matrices is almost independent of the size of the process grid P .

The overhead (%) $R_{recover_data}$ for constructing checksum matrices for matrix matrix multiplication is

$$R_{recover_data} = \frac{T_{recover_data}}{T_{matrix_mult}} = O\left(\frac{1}{Pm}\right) \quad (12)$$

6. PRACTICAL NUMERICAL ISSUE

The algorithm-based checkpoint-free fault tolerance presented in Section 4 involves solving system of linear equations to recover multiple simultaneous process failures. Therefore, in the practice of the algorithm-based checkpoint-free fault tolerance, the numerical issues involved in recovering multiple simultaneous process failures have to be addressed.

6.1 Numerical Stability of Real Number Codes

In Section 4, it has been derived that, to be able to recover from any no more than m failures, the encoding matrix A has to satisfy: any square sub-matrix (including minor) of A is non-singular. This requirement for the encoding matrix coincides with the properties for the generator matrices of real number Reed-Solomon style erasure correcting codes.

In fact, our weighted checksum encoding in Section 4.3 can be viewed as a version of the Reed-Solomon erasure coding scheme [22] in real number field. Therefore any generator matrix from real number Reed-Solomon style erasure codes can actually be used as the encoding matrix of algorithm-based checkpoint-free fault tolerance

In the existing real number or complex-number Reed-Solomon style erasure codes in literature, the generator matrices mainly include: Vandermonde matrix (Vander) [17], Vandermonde-like matrix for the Chebyshev polynomials (Chebvand) [3], Cauchy matrix (Cauchy), Discrete Cosine Transform matrix (DCT), Discrete Fourier Transform matrix (DFT) [12]. Theoretically, these generator matrices can all be used as the encoding matrix of the algorithm-based checkpoint-free fault tolerance scheme.

However, in computer floating point arithmetic where no computation is exact due to round-off errors, it is well known [14] that, in solving a linear system of equations, a condition number of 10^k for the coefficient matrix leads to a loss of accuracy of about k decimal digits in the solution. Therefore, in order to get a reasonably accurate recovery, the encoding matrix A actually has to satisfy *any square sub-matrix (including minor) of A is well-conditioned.*

The generator matrices from above real number or complex-number Reed-Solomon style erasure codes all contain ill-conditioned sub-matrices. Therefore, in these codes, when certain error patterns occur, an ill-conditioned linear system has to be solved to reconstruct an approximation of the original information, which can cause the loss of precision of possibly all digits in the recovered numbers.

6.2 Numerically Good Real Number Codes Based on Random Matrices

In this section, we will introduce a class of new codes that are able to reconstruct a very good approximation of the lost data with high probability regardless of the failure patterns processes. Our new codes are based on random matrices over real number field.

It is well-known [9] that Gaussian random matrices are well-conditioned. To estimate how well conditioned Gaussian random matrices are, we have proved the following Theorem:

THEOREM 2. *Let $G_{m \times n}$ be an $m \times n$ real random matrix whose elements are independent and identically distributed standard normal random variables, and let $\kappa_2(G_{m \times n})$ be the 2-norm condition number of $G_{m \times n}$. Then, for any $m \geq 2$, $n \geq 2$ and $x \geq |n - m| + 1$, $\kappa_2(G_{m \times n})$ satisfies*

$$\frac{\left(\frac{c}{x}\right)^{|n-m|+1}}{\sqrt{2\pi}} < P\left(\frac{\kappa_2(G_{m \times n})}{n/(|n-m|+1)} > x\right) < \frac{\left(\frac{c}{x}\right)^{|n-m|+1}}{\sqrt{2\pi}},$$

and

$$E(\ln \kappa_2(G_{m \times n})) < \ln \frac{n}{|n-m|+1} + 2.258,$$

where $0.245 \leq c \leq 2.000$ and $5.013 \leq C \leq 6.414$ are universal positive constants independent of m , n and x .

Due to the length of the proof for Theorem 1, we omit the proof here and refer interested readers to [6] for complete proof.

Note that any sub-matrix of a Gaussian random matrix is still a Gaussian random matrix. Therefore, a Gaussian

random matrix would satisfy any sub-matrix of the matrix is well-conditioned with high probability.

Theorem 2 can be used to estimate the accuracy of recovery in the weighted checksum scheme. For example, if an application uses 100,000 processes to perform computation and 20 processes to hold encodings, then the encoding matrix is a 20 by 100,000 Gaussian random matrix. If 10 processors fail concurrently, then the coefficient matrix A_r in the recovery algorithm is a 20 by 10 Gaussian random matrix. From Theorem 1, we can get

$$E(\log_{10} \kappa_2(A_r)) < 1.25$$

and

$$P(\kappa_2(A_r) > 100) < 3.1 \times 10^{-11}.$$

Therefore, on average, we will loss about one decimal digit in the recovered data and the probability to loss 2 digits is less than 3.1×10^{-11} .

7. EXPERIMENTAL EVALUATION

In this section, we experimentally evaluate the performance overhead of applying the algorithm-based checkpoint-free fault tolerance technique to the ScaLAPACK matrix-matrix multiplication kernel. We performed four sets of experiments to answer the following four questions:

1. What is the performance overhead of constructing checksum matrices?
2. What is the performance overhead of performing computations with checksum matrices?
3. What is the performance overhead of recovering FT-MPI programming environments?
4. What is the performance overhead of recovering checksum matrices ?

For each set of experiments, the size of the problems and the number of computation processes used are listed in Table 1.

Table 1: Experiment Configurations

Size of original matrix	12,800	19,200	25,600
Size of full checksum matrix	19,200	25,600	32,000
Process grid without FT	2 by 2	3 by 3	4 by 4
Process grid with FT	3 by 3	4 by 4	5 by 5

All experiments were performed on a cluster of 32 Pentium IV Xeon 2.4 GHz dual-processor nodes. Each node of the cluster has 2 GB of memory and runs the Linux operating system. The nodes are connected with a Gigabit Ethernet. The timer we used in all measurements is MPI_Wtime.

The programming environment we used is FT-MPI [10]. FT-MPI is a fault tolerant version of MPI that is able to provide basic system services to support fault survivable applications. FT-MPI implements the complete MPI-1.2 specification, some parts of the MPI-2 document and extends some of the semantics of MPI for allowing the application the possibility to survive process failures. FT-MPI can survive the failure of $n-1$ processes in a n -process job, and, if required, can re-spawn the failed processes. However, the application is still responsible for recovering the data structures and the data of the failed processes.

Although FT-MPI provides basic system services to support fault survivable applications, prevailing benchmarks show that the performance of FT-MPI is comparable [11] to the current state-of-the-art MPI implementations.

7.1 Overhead for Constructing Checksum Matrices

The first set of experiments is designed to evaluate the performance overhead of constructing checksum matrices. We keep the amount of data in each process fixed (that is the size of local matrices m fixed), and increase the size of the test matrices (hence the size of process grid).

Table 2: Time and overhead (%) for constructing checksum matrices

Size of original matrix	12,800	19,200	25,600
Exec. time for original matrix	442.9	695.0	989.8
Time for calculating encoding	38.0	40.8	43.2
Overhead (%) for encoding	8.6%	5.9%	4.4%

Table 2 reports the time for performing computations on original matrices and the time for constructing the three checksum matrices A^c , B^r , and C^f .

From Table 2, we can see that, as the size of the global matrices increases, the time for constructing checksum matrices increases only slightly. This is because, in the formula (7), when the size of process grid P is small, $32m^2\beta$ is the dominate factor in the time to constructing checksum matrices. Table 2 also indicates that the overhead (%) for constructing checksum matrices decreases as size of matrices increases, which is consistent with our theoretical formula (8) about the overhead for constructing checksum matrices.

7.2 Overhead for Performing Computations on Encoded Matrices

The algorithm-based checkpoint-free fault tolerance technique involve performing computations with checksum matrices, which introduces some overhead into the fault tolerance scheme. The purpose of this experiment is to evaluate the performance overhead of performing computations with checksum matrices.

Table 3: Time and overhead (%) for performing computations on encoded matrices

Size of original matrix	12,800	19,200	25,600
Size of full checksum matrix	19,200	25,600	32,000
Exec. time for original matrix	442.9	695.0	989.8
Exec. time for encoded matrix	462.6	716.4	1013.3
Increased time	19.7	21.4	23.5
Overhead (%)	4.4%	3.1%	2.4%

Table 3 reports the execution time for performing computations on original matrices and the execution time for performing computations on checksum matrices for different size of matrices.

Table 3 indicates the amount time increased for performing computations on checksum matrices increases slightly as the size of matrices increases. The reason for this increase is that, when perform computations with checksum matrices, column blocks of A^c (and row blocks of B^r) have to be broadcast to one more process. The dominate time for parallel matrix matrix multiplication is the time for computation which is the same for both fault tolerant algorithm and non-fault tolerant algorithm. Therefore, the amount time increased for fault tolerant algorithm increases only slightly as the size of matrices increases.

7.3 Overhead for Recovering FT-MPI Environment

The overhead for recovering programming environments depends on the specific programming environments. In this section, we evaluate the performance overhead of recovering FT-MPI environment.

Table 4: Time and overhead (%) for recovering FT-MPI environment

Size of original matrix	12,800	19,200	25,600
Exec. time for original matrix	442.9	695.0	989.8
Time for recovery FT-MPI	0.6	1.1	1.6
Overhead (%)	0.14%	0.16%	0.16%

Table 4 reports the time for recovering FT-MPI communication environment with single process failure. Table 4 indicates that the overhead for recovering FT-MPI is less than 0.2% which is negligible in practice.

7.4 Overhead for Recovering Application Data

The purpose of this set of experiments is to evaluate the performance overhead of recovering application data from single process failure.

Table 5: Time and overhead (%) for recovering replication data

Size of original matrix	12,800	19,200	25,600
Exec. time for original matrix	442.9	695.0	989.8
Time for recovery data	28.5	30.6	32.4
Overhead (%)	6.4%	4.4%	3.3%

Table 5 reports the time for recovering the three checksum matrices A^c , B^r , and C^f in the case of single process failure. Table 5 indicates that, as the size of the matrices increases, the time for recovering checksum matrices increases slightly and the overhead for recovering checksum matrices decreases, which again confirmed the theoretical results in Section 5.4.4.

8. DISCUSSION

This paper presented an algorithm-based checkpoint-free fault tolerance approach in which, instead of taking checkpoint periodically, a coded global consistent state of the crit-

ical application data is maintained in memory by modifying applications to operate on encoded data. Although the algorithm-based checkpoint-free fault tolerance in this paper share the same basic idea of modifying applications to operate on encoded data with the traditional the algorithm-based fault tolerance [19], they assume very different a failure model.

Compared with the typical checkpoint/restart approaches, the algorithm-based checkpoint-free fault tolerance in this paper can only tolerate partial process failures. It needs the support from programming environments to detect and locate failures. It requires the programming environments to be robust enough to survive node failures without suffering complete system failure. Both the overhead of and the additional effort to maintain a coded global consistent state of the critical application data in algorithm-based checkpoint-free fault tolerance is usually highly dependent on the specific characteristic of the application. Therefore, it is possible that the algorithm-based checkpoint-free fault tolerance approach introduces higher overhead than checkpoint approaches.

Unlike in typical checkpoint/restart approaches which involve periodical checkpoint and rollback-recovery, there is no checkpoint or rollback-recovery involved in this approach. Furthermore, in the algorithm-based checkpoint-free fault tolerance in this paper, whenever process failures occur, it is only necessary to recover the lost data on the failed processes. Therefore, for many applications, it is also possible for this approach to achieve a much lower fault tolerant overhead than typical checkpoint/restart approaches. As shown in Section 5 and 7, for matrix matrix multiplication, which is one of the most fundamental operations for computational science and engineering, as the size N of the matrix increases, the fault tolerance overhead decreases with the speed of $\frac{1}{N}$.

Although the algorithm-based checkpoint-free fault tolerance approach presented in this paper is non-transparent and algorithm-dependent, it is meaningful in that

1. It can often achieve a surprisingly low overhead in parallel matrix computations where it usually works.
2. It is often possible to build it into frequently used numerical libraries such as ScaLAPACK to relieve the involvement of the application programmer.

9. CONCLUSION AND FUTURE WORK

In this paper, we presented an algorithm-based checkpoint-free fault tolerance approach in which, instead of taking checkpoint periodically, a coded global consistent state of the critical application data is maintained in memory by modifying applications to operate on encoded data. Although the applicability of this approach is not so general as the typical checkpoint/rollback-recovery approach, in parallel matrix computations where it usually works, because no periodical checkpoint or rollback-recovery is involved in this approach, process failures can often be tolerated with a surprisingly low overhead.

We showed the practicality of this technique by applying it to the ScaLAPACK matrix-matrix multiplication kernel which is one of the most important kernels for ScaLAPACK library to achieve high performance and scalability. Experimental results demonstrated that the proposed checkpoint-

free approach is able to survive process failures with a very low performance overhead.

We addressed the practical numerical issue in this technique by proposing a class of numerically good real number erasure codes based on random matrices. We proved our codes are numerically highly reliable in recovering the lost data for multiple simultaneous process failures.

For the future, we plan to incorporate this fault tolerance technique into more ScaLAPACK library routines and more high performance computing applications. We would also like to evaluate this technique on systems with larger number of processors.

10. REFERENCES

- [1] N. R. Adiga and et al. An overview of the BlueGene/L supercomputer. In *Proceedings of the Supercomputing Conference (SC'2002), Baltimore MD, USA*, pages 1–22, 2002.
- [2] L. S. Blackford, J. Choi, A. Cleary, A. Petitet, R. C. Whaley, J. Demmel, I. Dhillon, K. Stanley, J. Dongarra, S. Hammarling, G. Henry, and D. Walker. ScaLAPACK: a portable linear algebra library for distributed memory computers - design issues and performance. In *Supercomputing '96: Proceedings of the 1996 ACM/IEEE conference on Supercomputing (CDROM)*, page 5, 1996.
- [3] D. L. Boley, R. P. Brent, G. H. Golub, and F. T. Luk. Algorithmic fault tolerance using the lanczos method. *SIAM Journal on Matrix Analysis and Applications*, 13:312–332, 1992.
- [4] G. Bosilca, Z. Chen, J. Langou, and J. Dongarra. Recovery patterns for iterative methods in a parallel unstable method. *Submitted to SIAM J. Scientific Computing*, 2004.
- [5] G. Bronevetsky, D. Marques, K. Pingali, and P. Stodghill. Automated application-level checkpointing of mpi programs. *ACM SIGPLAN PPOPP'03.*, 38(10):84–94, 2003.
- [6] Z. Chen and J. Dongarra. Condition numbers of gaussian random matrices. Technical Report ut-cs-04-539, University of Tennessee, Knoxville, Tennessee, USA, 2004.
- [7] Z. Chen, G. E. Fagg, E. Gabriel, J. Langou, T. Angskun, G. Bosilca, and J. Dongarra. Fault tolerant high performance computing by a coding approach. In *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2005, June 14-17, 2005, Chicago, IL, USA*. ACM, 2005.
- [8] J. Dongarra, H. Meuer, and E. Strohmaier. TOP500 Supercomputer Sites, 24th edition. In *Proceedings of the Supercomputing Conference (SC'1994), Pittsburgh PA, USA*. ACM, 1994.
- [9] A. Edelman. Eigenvalues and condition numbers of random matrices. *SIAM J. Matrix Anal. Appl.*, 9(4):543–560, 1988.
- [10] G. E. Fagg, E. Gabriel, G. Bosilca, T. Angskun, Z. Chen, J. Pjesivac-Grbovic, K. London, and J. J. Dongarra. Extending the MPI specification for process fault tolerance on high performance computing systems. In *Proceedings of the International Supercomputer Conference, Heidelberg, Germany*,

- 2004.
- [11] G. E. Fagg, E. Gabriel, Z. Chen, , T. Angskun, G. Bosilca, J. Pjesivac-Grbovic, and J. J. Dongarra. Process fault-tolerance: Semantics, design and applications for high performance computing. *Submitted to International Journal of High Performance Computing Applications*, 2004.
- [12] P. Ferreira and J. Vieira. Stable dft codes and frames. *IEEE Signal Processing Letters*, 10(2):50–53, 2003.
- [13] I. Foster and C. Kesselman. The globus toolkit. *The grid: blueprint for a new computing infrastructure*, pages 259–278, 1999.
- [14] G. H. Golub and C. F. Van Loan. *Matrix Computations*. The John Hopkins University Press, , 1989.
- [15] E. Gabriel, G. E. Fagg, G. Bosilca, T. Angskun, J. Dongarra, J. M. Squyres, V. Sahay, P. Kambadur, B. Barrett, A. Lumsdaine, R. H. Castain, D. J. Daniel, R. L. Graham, and T. S. Woodall. Open mpi: Goals, concept, and design of a next generation mpi implementation. In *PVM/MPI*, pages 97–104, 2004.
- [16] G. A. Geist and C. Engelmann. Development of naturally fault tolerant algorithms for computing on 100,000 processors. *Submitted to Journal of Parallel and Distributed Computing*, 2002.
- [17] C. N. Hadjicostis and G. C. Verghese. Coding approaches to fault tolerance in linear dynamic systems. *Submitted to IEEE Transactions on Information Theory*, 2004.
- [18] P. D. Hough, M. E. Goldsby, and and E. J. Walsh. Algorithm-dependent fault tolerance for distributed computing. Technical Report SAND2000-8219, Sandia Technical Report, 2000.
- [19] K.-H. Huang and J. A. Abraham. Algorithm-based fault tolerance for matrix operations,. *IEEE Transactions on Computers*, vol. C-33:518–528, 1984.
- [20] I. Foster and C. Kesselman. *The Grid: Blueprint for a New Computing Infrastructure*. Morgan Kaufman, San Francisco, 1999.
- [21] Y. Kim. *Fault Tolerant Matrix Operations for Parallel and Distributed Systems*. Ph.D. dissertation, University of Tennessee, Knoxville, June 1996.
- [22] J. S. Plank. A tutorial on Reed-Solomon coding for fault-tolerance in RAID-like systems. *Software – Practice & Experience*, 27(9):995–1012, September 1997.
- [23] J. S. Plank, K. Li, and M. A. Puening. Diskless checkpointing. *IEEE Trans. Parallel Distrib. Syst.*, 9(10):972–986, 1998.
- [24] P. Sanders and J. F. Sibeyn. A bandwidth latency tradeoff for broadcast and reduction. *Inf. Process. Lett.*, 86(1):33–38, 2003.
- [25] V. S. Sunderam. PVM: a framework for parallel distributed computing. *Concurrency: Pract. Exper.*, 2(4):315–339, 1990.