

Probabilistic Communication and I/O Tracing with Deterministic Replay at Scale

Xing Wu¹ Karthik Vijayakumar¹ Frank Mueller¹ Xiaosong Ma^{1,2} Philip C. Roth²

¹ Department of Computer Science, North Carolina State University, Raleigh, NC 27695-7534

² Computer Science and Mathematics Division, Oak Ridge National Laboratory, Oak Ridge, TN 37831

Abstract

With today’s petascale supercomputers, applications often exhibit low efficiency, such as poor communication and I/O performance, that can be diagnosed by analysis tools. However, these tools either produce extremely large trace files that complicate performance analysis, or sacrifice accuracy to collect high-level statistical information using crude averaging.

This work contributes Scala-H-Trace, which features more aggressive trace compression than any previous approach, particularly for applications that do not show strict regularity in SPMD behavior. Scala-H-Trace uses histograms expressing the probabilistic distribution of arbitrary communication and I/O parameters to capture variations. Yet, where other tools fail to scale, Scala-H-Trace guarantees trace files of near constant size, even for variable communication and I/O patterns, producing trace files orders of magnitudes smaller than using prior approaches. We demonstrate the ability to collect traces of applications running on thousands of processors with the potential to scale well beyond this level. We further present the first approach to deterministically replay such probabilistic traces (a) without deadlocks and (b) in a manner closely resembling the original applications.

Our results show either near constant sized traces or only sub-linear increases in trace file sizes irrespective of the number of nodes utilized. Even with the aggressively compressed histogram-based traces, our replay times are within 12% to 15% of the runtime of original codes. Such concise traces resembling the behavior of production-style codes closely and our approach of deterministic replay of probabilistic traces are without precedence.

1. Introduction

As supercomputers progress in scale and capability toward exascale levels, characterization of communication and I/O behavior is becoming increasingly difficult due to system size and complexity. Today, many scientific applications execute in ten thousands of cores or more. Moreover, modern supercomputers are equipped with complex network interconnects to improve the speedup of parallel applications. Apart from the network complexity, different vendors employ different interconnect designs to improve the overall communication performance, thereby achieving better speedup. For example, the IBM Blue Gene family of supercomputers employs five different network interconnects [1]. Such interconnects mandate performance study of applications for efficient use of available resources. Even finding the most efficient task mapping to nodes has become difficult with complex, new system designs.

The large numbers of processors/cores, increased aggregate memory capacity, and optimized interconnects allow applications to grow not only in targeting larger problem sizes, but to explore more sophisticated communication models. Extreme-scale applications are often complex codes, integrating multiple software components exercising vastly different computation/communication models. Such codes are becoming more dynamic and diverging from strict, regular single program, multiple data (SPMD) behavior. Examples include multi-physics or coupled codes, where partitions of nodes implement different simulation models, work on separate datasets, or even conduct analytics tasks such as data reduction. Such applications exhibit multiple program, multiple data (MPMD) behavior as multiple nodes work on multiple sections of the program. *E.g.*, in climate simulations, some nodes simulate cli-

mate changes over land and while other nodes work on sea models. Hence, different modules, like land and sea, use different input data and algorithms resulting in different communication behavior within each module.

Several studies have investigated the communication and I/O characteristics of applications. They focus on three main classes: tracing tools, capable of capturing and recording all message events at the cost of high storage requirements; profiling tools, designed to provide low-overhead performance summaries trading off storage space for detail level; and communication and I/O kernels that eliminate computation and retain only application communication and I/O behavior. Although application kernels are designed to capture the exact application behavior, it is difficult to keep these kernels up-to-date since the applications constantly evolve over time. Application traces, in contrast, can be readily generated by simple instrumentation of an application, to keep up with a changing code base. This makes performance analysis via traces a preferred method to analyze parallel applications in practice.

The combination of job scale and application complexity, however, creates unique challenges for parallel tracing tools. On one end of the spectrum, traditional tracing tools (such as Vampir [10]) record all events sequentially for each parallel process. For large application runs on leadership-class supercomputers, this approach generates unmanageable trace file sizes, introducing prohibitive overheads, *e.g.*, for copying trace files from temporary to permanent storage, hitting the maximum storage limit, and even the need for a cluster plus another parallelized tool to perform trace analysis [2]. On the other end of the spectrum, tools that only report statistical information (such as mpiP [17]) may fail to deliver the level of detail needed in performance analysis or debugging.

On-the-fly trace compression [12, 15] provides lossless tracing *and* dramatically reduced trace file sizes, and it has recently been extended to conduct multi-level I/O tracing [19] in addition to capturing communication calls. However, effective compression builds on the homogeneous behavior across processes (inter-node compression) and repetitive behavior within a process (intra-node compression). With complex, irregular, or self-adjusting applications, such assumptions do not hold and compression fails due to mismatches between traced events.

In this work, we propose Scala-H-Trace, with a novel approach to collect concise traces for applications exhibiting *non-SPMD* behavior. In other words, while past approaches proved effective for the easier problems of tracing SPMD codes, this work focuses on the much harder problems of tracing non-SPMD codes. Scala-H-Trace is motivated by the tradeoff between exact details and manageability of trace file size. Although having exact details helps in root cause analysis, lossless tracing becomes increasingly unaffordable on ultra-scale machines. *Histograms* in Scala-H-Trace provide an opportunity to collect overall statistical details, *e.g.*, data send volume, which can be useful in studying network characteristics of the application. They provide an overall “big picture” of an application’s communication behavior. Scala-H-Trace employs histograms with multiple bins whose value ranges are dynamically adapted as trace data is recorded on-the-fly. In addition to representing a distribution, each bin also retains certain crude statistical information (min/max/avg/stddev), potentially useful for root cause analysis.

Scala-H-Trace also enables the user to set a *merge precision level* during trace collection. This precision level drives the com-

pression efficiency by collecting statistical information for varying traces in unique histogram bins. The trace precision is also ensured to fall within the user set value. If the trace precision falls below the specified precision, mismatching trace events are recorded without histogram-based compression, *i.e.*, traditional structural compression techniques are employed and may fail to provide concise traces in the absence of SPMD behavior of the code [12] but result in exact event recording. The size of such a trace file then becomes a function of the desired merge precision level, which can be tuned to obtain a manageable size while retaining trace artifacts suitable for performance analysis or even detect root causes in performance. At the same time, our unique approach to collect histogram-based statistical information captures the overall trend in communication and I/O behavior of applications executing on thousands of nodes. Our histogram-based approach also reduces the tracing overhead as the time taken to compress smaller histogram-based traces is considerably less than the traditional lossless traces.

While histogram-based tracing can effectively reduce trace data volumes, it creates several challenges for accurate replay of the traced events. To ensure the correctness of the captured trace and to reproduce the communication and I/O behavior, we also provide a novel replay facility. This replay tool reissues the recorded trace events without decompressing the compressed trace. If the compressed trace is lossless, sender-destination node information and communication volume are recorded precisely. Also, the causal ordering of the original application is preserved. For lossy, histogram-based traces, our tool employs a distributed, orchestrated and deterministic replay capability. Our goal in the replay of histogram-based traces is not to capture exact original events but rather the existence of a sequence of events with comparable timings and communication endpoints. Resulting information can be useful in identifying bottlenecks and also the communication pattern of a particular application.

Contributions: Our contributions are as follows:

- We provide the capability to record lossless and concise communication and I/O traces for non-SPMD programs.
- We create novel capabilities for more aggressive trace compression based on a precision level, selected by the user, that drives both compression efficiency and trace accuracy.
- We support a replay technique that reissues trace events without decompressing the original trace file.
- We provide a distributed approach to replay statistical traces that does not require back-channel communication to preserve causal event ordering for correctness.
- We ensure that replay is deterministic by (a) coordinating sender/receiver activity through receive reordering, (b) letting a node interpret the events of all other nodes and (c) ensuring that all nodes use the same random number sequence for probabilistic replay resulting in the same parameter and end-point choices for communication.
- We proved that deterministic replay after reordering is deadlock free [20].

We evaluated our approach with the Parallel Ocean Program (POP) and two benchmarks from NAS parallel benchmark suite. POP is both computational and I/O intensive and hence a representative application to evaluate our tool. Our results provide one to two orders of magnitude smaller trace files than any previous approach. We also evaluated our replay tool by reissuing histogram-based traced events. The replay time only deviated 12% to 15% from the original application’s time in most cases, even for most aggressively merged histogram-based traces.

2. Background

Scala-H-Trace is a novel design of a communication and I/O tracing tool that shares its methodology for representing the resulting trace on a single file (instead of one file per node), both otherwise relies on histogram-based data collection. While Scala-H-Trace was derived from the publicly available code of ScalaTrace/ScalaIOTrace [12, 15, 19], Scala-H-Trace provides entirely novel compression capabilities.

Scala[IO]Trace collects communication and optionally parallel I/O traces, using the MPI Profiling layer (PMPI) [9] through Umptire [18] to intercept MPI calls and to collect MPI traces. It features aggressive trace compression that generates a single, concise and lossless trace file from any large-scale parallel application run. It also preserves timing information in the compressed form along with the calling context of events being traced. In this paper, we develop Scala-H-Trace, which provides even more aggressive trace compression techniques to serve real-world scientific applications that do not show strict SPMD regularity.

Scala[IO]Trace performs two types of compression: *intra-node* and *inter-node*. The former exploits the repetitive nature of timestep simulation in parallel scientific applications. The latter exploits the homogeneity in behavior (SPMD) among different processes running the application. Intra-node compression is performed on-the-fly within a node. Inter-node compression is performed across nodes by forming a radix tree structure among all nodes and sending all intra-node compressed traces to respective parents in the radix tree. This results in a single compressed trace file capturing the entire application run across all nodes. The compression algorithm is discussed in detail elsewhere [12, 15]. Scala-H-Trace employs a different intra- and inter-node compression algorithm due to its reliance on histograms but still shares the reduction over a radix tree with Scala[IO]Trace.

2.1 Trace Compression

We briefly introduce several techniques used in Scala[IO]Trace to allow a later comparison with Scala-H-Trace. Repetitive events in different iteration of loops are collected as Regular Section Descriptors (RSD) [5] and power-RSDs capture RSD events in nested loops [8], both of which are represented in constant size. Consider the following code snippet:

```
for( i = 0; i < 10; i++ ) {
    for( j = 0; j < 100; j++ ) {
        compute1();
        MPI_Irecv(...); // Receive from left neighbor
        MPI_Isend(...); // Send to right neighbor
        MPI_Waitall(...);
    }
    MPI_Allreduce(...); // Collective reduction operation
}
```

Trace compression results in the following tuples: RSD1:{100, MPI_Irecv, MPI_Isend, MPI_Waitall} representing 100 iterations of MPI_Irecv, MPI_Isend and MPI_Waitall in the inner loop, PRSD1: {10, RSD1, MPI_Allreduce} denoting 10 iterations RSD1 followed by MPI_Allreduce in the outermost loop. The algorithm uses the calling context of events to match repetitive behavior. This ensures that identical MPI functions originating from different call paths of the application are not compressed together. Since trace events from different nodes are collected and merged in a single output trace file, the *task rank* information of nodes participating in an event is also compressed and encoded concisely in the compressed trace. This participant node information is represented in a tuple containing starting rank, total number of participants and an offset value separating ranks. Even multi-dimensional information is captured in this encoding format.

Apart from matching calling contexts, the compression algorithm matches function parameters and merges them along with

compressed events ensuring that no information is lost. In typical parallel applications, communication end-points differ across nodes as a result of communication with neighboring nodes. These varying end-points inhibit event compression. Scala[IO]Trace uses a unique location-independent encoding to represent communication end-points in events like `MPI_Send` and `MPI_Recv`. It also encodes MPI opaque pointers like `MPI_File` and `MPI_Comm`, which do not exhibit repetitive patterns, potentially inhibiting effective compression. There are special cases in which events with matching calling context can have non-matching function parameters. These non-matching function parameters are compressed using a vector representation, so that the particular event can be concisely represented in the trace.

2.2 Time Preservation

Another important feature of Scala[IO]Trace is the time preservation of captured traces. Instead of recording absolute timestamps, the tool records delta time of computation durations between adjacent communication calls. During RSD formation, instead of accumulating exact delta timestamps, statistical histogram bins are utilized to concisely represent timing details across the loop. These bins are comprised of statistical timing data (minimum, maximum, average and standard deviation). More details on collecting statistical timing information are provided elsewhere [15].

2.3 Timed Replay

Scala[IO]Trace not only supports scalable tracing, it also supports a scalable replay engine. Given a single, compressed trace file, the replay engine allows all I/O and communication calls to be reissued without trace decompression while preserving event ordering. The replay engine runs as an MPI job with the same number of tasks as its original application. It replays I/O and communication events in each node with their original parameters except for actual file content/message payloads. Instead, a random buffer of the same size as the original file/message buffer is used. Additionally, computation time on each node is simulated by a delay between traced events based on recorded delta time.

3. Inter-node Trace Compression

The SPMD nature of the scientific codes causes participants of a parallel application to produce similar per node traces. E.g., if we treat a trace as a sequence of MPI events, traces from different nodes tend to have similar subsequences that contains most of MPI events. In addition, loop structures captures by PRSDs in ScalaTrace facilitate compression as traces from different nodes tend to have similar PRSD nests. ScalaTrace originally required not just similar but rather *identical* patterns, i.e., it failed to fully exploit similarities for inter-node trace compression. More specifically, identical loop structures, i.e., PRSDs with identical length, iteration count, and MPI event sequence were required. While this approach works well with the perfect SPMD-style codes, it is subject to scalability problem when traces slightly diverge between nodes. For the example below, let T_i be traces where each letter in a trace “string” represents an MPI event and the pair of parentheses represent the loop structures. The coarse-grained trace matching algorithm may merge the per-node traces T_1 and T_2 to T_3 . Yet, an ideal compression would instead be trace T_4 .

$$T_1 : a(b(bcb)db)a \quad T_2 : a(b(beb)fb)a$$

$$T_3 : a(b(bcb)db)(b(beb)fb)a \quad T_4 : a(b(bceb)d)fb)a$$

Only if the inter-node trace matching algorithm does not miss the structural similarities can the probabilistic communication parameter compression (discussed below) be fully utilized. Hence, we have designed a novel, fine-grained event matching algorithm that recursively compares and merges the nested loop structures. Algorithm 1 outlines the recursive trace merging technique. This

algorithm traverses traces of two nodes, T_1 and T_2 , to identify the matching event pairs. Stand-alone events are compared by their MPI parameter values with the function `PARAM_MATCH`. If two events start structurally identical loop nests, i.e., loop nests with equal depth and equal iteration counts at each nest level, the function `MATCH_LOOP` is called. `MATCH_LOOP` then matches the loop bodies at each level starting from the innermost nest and recursively call itself if new matching loop heads are found. When a pair of matching events is identified, the preceding unmatched sequences are sequentially linked by `DO_MERGE`. Since we forward the cursors for both input sequences when a match is found, this algorithm, in practice, has a complexity of $O(n)$, where n is the length of the input traces given that two input traces are similar.

Algorithm 1 Recursive Trace Matching Algorithm

Precondition: T_1 and T_2 : input per node traces

Postcondition: T_1 and T_2 : recursively merged trace

```

1: procedure MATCH_TRACE( $T_1, T_2$ )
2:   for iter1  $\leftarrow T_1.head, T_1.tail$  do
3:     for iter2  $\leftarrow T_2.head, T_2.tail$  do
4:       if iter1 and iter2 start identical loop nests then
5:         MATCH_LOOP(iter1, iter2, depth_of_nest)
6:       else
7:         if PARAM_MATCH(iter1, iter2) then
8:           DO_MERGE(iter1, iter2)
9:         end if
10:      end if
11:    end for
12:  end for
13: end procedure

14: procedure MATCH_LOOP(loop1, loop2, depth)
15:   for iter1  $\leftarrow loop1.head, loop1.tail$  do
16:     for iter2  $\leftarrow loop2.head, loop2.tail$  do
17:       if iter1 == loop1.head && iter2 == loop2.head &&
18:         PARAM_MATCH(iter1, iter2) then
19:         DO_MERGE(iter1, iter2)
20:       end if
21:       if iter1 and iter2 are single events &&
22:         PARAM_MATCH(iter1, iter2) then
23:         DO_MERGE(iter1, iter2)
24:       end if
25:       if iter1 and iter2 start identical loop nests then
26:         MATCH_LOOP(iter1, iter2, depth_of_nest)
27:       end if
28:     end for
29:   end for
30:   if depth > 1 then
31:     MATCH_LOOP(iter1, iter2, depth-1)
32:   end if
33: end procedure

```

Algorithm 1 may still fail to generate the best inter-node compression because traversing two sequences with the double-nested loop structure does not guarantee identifying the longest common subsequence. As an example, consider T_1 and T_2 below. Algorithm 1 will return the sequence T_3 :

$$T_1 : abbbb \quad T_2 : bbbba \quad T_3 : bbbbabbbb$$

The matching event a is found before the longer subsequence $bbbb$. To solve this problem, we integrated a *Weighted Longest Common Subsequence* (WLCS) algorithm into Algorithm 1. WLCS is adapted from the classic *Longest Common Subsequence* (LCS) algorithm. Since the loop structures in the trace should be treated as a whole, we enhanced LCS such that the matching loop structures are evaluated with a weight that is equal to the length of their longest common subsequence. This addresses the above example to compress $bbbb$ first.

4. Histogram-Based Trace Collection

Noeth *et al.* [12] provide trace compression techniques resulting in an almost constant sized trace file or sublinear increases in trace file size with strong scaling (increasing number of nodes). Yet, these results only hold for SPMD-style benchmarks, not for production size applications with non-SPMD patterns. ScalalOTrace [19] provides mechanisms to collect both the communication and I/O traces for scientific applications like the Parallel Ocean Program (POP) [14]. But for some scientific applications, including POP, the inter-node compression technique fails to obtain a near-constant sized trace file with increasing number of nodes. Instead, we see a linear increase in the trace size due to non-SPMD style programming.

POP performs ocean simulation for multiple time steps. Each time step performs a set of computations and communications of an inner loop in multiple iterations. Due to different data-dependent convergence points in the computation across different timesteps, the number of inner loop iterations varies from timestep to timestep. Even though all MPI events originate from the same calling sequence (call stack), varying loop iteration counts in each timestep inhibit intra-node compression and thus negatively impact inter-node compression across all nodes. This behavior can also be observed in many Adaptive Mesh Refinement (AMR) applications in which the input set is dynamically rebalanced on a periodic basis.

To address these problems, we propose a novel method of tracing. We promote histogram-based trace information for a predefined user-tunable merge precision level to obtain higher compression rates of trace events — at the expense of accuracy. Consider the following 3 scenarios: (1) If the user sets the merge precision level to 100%, then only events with perfectly matching function parameters are merged. (2) If the user sets the merge precision level to 95%, then events with non-matching function parameters will be merged if and only if all pairs of parameters differ by no more than 5%. Should any pair of parameters exceed the 5% threshold, we fall back to lossless tracing. (3) If the user sets the merge precision level to 0%, then events with non-matching function parameters are also merged and the non-matching parameters are collected in histogram bins. Note that the function calling contexts always have to match for two events to be merged. Figure 1 explains the difference in the merge precision level and the precision level of the trace file. A merge precision level of 0% does not mean that the entire meaningful information is lost. Even at a 0% merge precision level, the statistical function parameters collected in histogram bins still capture the overall behavior of the application. Depending on user needs, the smallest traces with high application resemblance collected using a 0% merge precision level may be much more useful than unmanageably large trace files. We provide the option for users to decide on the tradeoff between the manageability of trace files vs. capturing the exact application behavior.

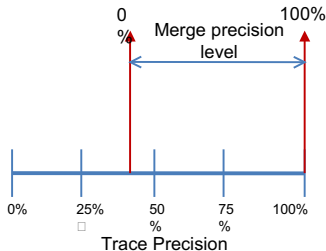


Figure 1. Trace precision vs Merge precision

Our approach uses histograms to collect probabilistic information on varying trace events and event parameters that otherwise inhibit trace compression. Histogram-based collection employs a technique to collect statistical information in dynamically balanced bins. The online balancing algorithm equalizes the number of items

per bin while adjusting their value range constraints. This ensures that the histogram captures outliers and other statistical distribution properties missed by simple aggregate statistical collection like maximum, minimum, average and variance. We also collect maximum and minimum participant rank information along with the frequency in bins so as to enable root cause detection, *e.g.*, due to load imbalance. Even with this lossy trace information, histograms help in providing more insight into the general characteristics of the traced application. Histogram details can be collected at various levels in the trace. The following explains what trace information is collected as histogram and discusses possible tradeoffs in collecting statistical information versus non-lossy information.

4.1 Intra-node Event Histogram

The loop iteration count denoted by PRSDs can be collected as a histogram. This enables better compression of repeating events in many scientific applications that otherwise would fail to compress due to data dependencies. Although the exact iteration count is lost in the final trace, the number of loop iterations directly depends on the computation, which, in turn, varies with different input sets. Hence, collecting statistical loop iteration counts only has a minor impact in capturing the communication behavior of the application. The main advantage of this approach is the ability to obtain a concise trace file by allowing a small percentage of lossy trace collection that otherwise would have resulted in a trace file of unmanageable size.

Consider the code snippet:

```
for( i = 0; i < 50; i++ ) {
  //Perform calculation till the result converges
  while(result > convergence_factor) {
    do_calculation();
    MPI_Irecv(...); //Receive from neighbors
    MPI_Send(...); //Send to neighbors
    MPI_Wait(...);
  }
}
```

In the above example, if the iteration count matches across time-steps, the resulting PRSD will be of form PRSD1: {50, RSD1}. Due to mismatching convergence points across different time-steps, the following sample events can occur:

```
RSD1: <39, MPI_Irecv, MPI_Send, MPI_Wait>
RSD2: <40, MPI_Irecv, MPI_Send, MPI_Wait>
RSD3: <39, MPI_Irecv, MPI_Send, MPI_Wait>
RSD4: <42, MPI_Irecv, MPI_Send, MPI_Wait> ... till RSD50
```

The expected PRSD is not formed due to mismatching RSDs across time steps. This leads to cascading compression failures across nodes. As a result, the trace file size increases linearly with the increase in participating nodes. Histogram-based trace collection ensures that the varying iteration count is captured in histogram bins. Hence, the resulting trace will have just one PRSD for the entire time-step calculation.

4.2 Inter-node Event Histogram

With inter-node event compression, compressed traces from different nodes are merged together. During this process, a radix tree structure is formed among all nodes. Child nodes send their respective intra-node compressed traces to their parents. A parent node performs compression of matching events from its child nodes. For each and every event from the parent node, a matching child event is searched. If there is a match, the parent event's participant list is updated with the rank of the child node and the child event is discarded. Other unmatched events will be reordered according to its dependency with other events.

In applications with non-SPMD behavior, loops created during intra-node compression can have matching events across nodes, but fail to compress across nodes due to a mismatch in the loop iteration count. This prevents the entire loop from being merged, increasing the trace file size linearly with the number of nodes.

As an example, consider the code snippet from the Section 4.1 again. Table 1 shows one such scenario in which computation dependent loop iterations fail to merge across nodes. By collecting loop iterations in histogram, all events merge successfully across nodes. Note that we enable merge only when all events inside the loop match perfectly. If events from two loop candidates do not match, then these loops are considered to represent *different* scenarios in the original application and are hence not merged.

Participants:0-3	Participants:4-7
Loop 50 times	Loop 51 times
MPI_Irecv from 0-3	MPI_Irecv from 4-7
MPI_Send to 0-3	MPI_Send to 4-7
MPI_Wait	MPI_Wait

Table 1. Varying Loop Iteration

4.3 Function Parameter Histogram

Apart from collecting loop iteration counts in histograms, MPI function parameters, such as Send/Recv volume, tag and sender/destination ranks, can also be recorded in histograms. The Send/Recv data volume is important to capture the network load due to the communication calls issued by the application. Send/Recv volumes often vary across different timesteps in AMR applications. This variation in volume inhibits compression of communication calls originating from the same call stack, thereby inhibiting compression across an entire loop due to a small deviation in the data volume parameters. There are other methods to collect exact volumes. One such method is to collect the volume information in a vector along with the rank information of participating nodes. But this results in a linear growth with the increase in number of participating nodes, which is non-scalable.

For applications that do not exhibit a regular communication pattern, it is impossible to compress repeating communication events originating from the same call stack with different sender/destination ranks. The approach of location-independent relative encoding of communication end-points provides a novel opportunity for event compression. But even this approach only succeeds in the case of applications with regular communication patterns. There are approaches in which the communication function call can be expressed as a PRSD but different end-points in different loop iterations have to be collected as a vector. Again, such an approach is not scalable for applications executed on thousands of nodes. An example of collecting function parameters as a vector is given below:

```

Loop iterations: 1 - 5
Participants: 0 - 9 (node ranks)
Event: MPI_Send
Data volume: 90 bytes [ranks: 0,1,4,5,8,9],
              92 bytes [ranks: 2,3,6,7]
Destination: relative rank 1 [ranks: 0,2,4,6,8],
              relative rank 9 [ranks: 1,3,5,7,9]

```

Assume MPI_Send is executed in a loop 5 times. With lossless trace collection, both data volume and destination will be recorded along with the rank information of the corresponding participant. The relative ranks shown above is location independent: 1 represents “the next right neighbor” and 9 represents “the next left neighbor”. This compression results in a more concise representation than its uncompressed equivalent, but it still suffers from increases in the trace size proportional to the number of nodes since no regularity for rank lists could be deduced.

Using histograms to collect relative end-points and data volume allows better compression of repeating events originating from the same call stack. For this example, histograms will record both destinations 1 and 9 in bins along with its frequency. In addition to binning communication end-points, we also collect relative ranks in a bitmap and encode it in the trace file. This provides information

on exact values that are missing from the histograms and aids post-mortem analysis tools. In the above example, an analysis tool may choose relative ranks of either 1 or 9 while relative ranks between 2 and 8 are excluded from the pseudo-random selection. We reiterate that we provide this lossy trace collection as a feature and the decision to use this feature is entirely upon the users. Users may choose to enable histogram-based tracing and configure the merge precision level in response to their application analysis needs, overheads and storage availability.

4.4 Histogram Construction

We have designed our system in a way to collect exact trace information as much as possible. Users can set a target merge precision level expressed as a percentage. Our compression algorithm attempts to match events originating from the same call stack and compresses events only if all function parameters match. Histogram collection is triggered only if there is a mismatch in function parameters or in the loop information. In such cases, the difference between two non-matching values is checked in terms of the user specified merge precision level. If the difference is within the target precision range, then events are merged and the non-matching parameters are recorded in a histogram from there on. If the difference falls out of the target precision range, then either event compression will fail or data is recorded in a vector as shown in the example in 4.3.

In our current implementation, the number of histogram bins is fixed at the start of the application run, but the value ranges in bins are dynamically adjusted. We provide an option to set an interval after which bins are adjusted. Two bins with the lowest frequencies are combined and the bin with maximum frequency is split into two bins. We further store auxiliary information in bins, such as minimum/maximum/average/variance, which are adjusted accordingly. Apart from per-bin statistical information, we also collect maximum/minimum values over the entire value range (all bins) and the node ranks associated with those. This provides outlier information and can be used in the replay studies and other performance analysis tools.

5. Deterministic Replay

While histogram-based trace collection is powerful in compressing irregular or dynamically changing events, the collected traces themselves create challenges for replaying and subsequent performance analysis. The core challenge of histogram-based replay is to ensure that events are issued in a deterministic manner across nodes and with coordinated parameter value selections for common communication end-points of sends and receives. Since Scala-H-Trace collects statistical values for communication volume, tags, and end-points, the conventional ScalaTrace replay design for lossless traces, which takes an independent, uncoordinated approach among nodes, can lead to potential deadlocks due to statistical uncertainty, or may fail to re-create the original communication or I/O pattern with reasonable proximity. The nature of histogram-based traces mandates a distributed, orchestrated replay with coordination among all participating nodes to ensure deterministic event sequences during replay for Scala-H-Trace. All nodes must read all trace events. They need to agree on a specific value selected from the statistical information found in the trace.

We have fundamentally redesigned the replay tool [12] to reissue MPI calls from lossless traces such that the trace data collected using histograms is honored during event replay. Our replay tool issues MPI calls using the compressed trace independent of the original application and without decompressing the trace. This tool verifies the correctness of the collected trace. It can also assist in the performance tuning of MPI communication and facilitates projections of network requirements for future procurements. Apart from replaying MPI calls, this basic replay framework can also be

extended to integrate with other performance analysis/tuning tools and it can be used to perform automated communication and network metric calculations.

Before we discuss the design of our new Scala-H-Trace replay tool, we first review the conventional design of replay for lossless traces in ScalaTrace [12]. For lossless traces, all participating nodes parse the trace file and *only act on events if the current node is a member of the participant list*. Then all nodes reissue MPI events one by one by identifying loops using the PRSD information and extracting individual MPI function parameters from the recorded trace. This replay tool also reads the delta time information from the trace and simulates the computation time by sleeping in place of computation. This simulates the exact communication and I/O behavior of the application in terms of interconnect characteristics, such as contention. The replay tool helps to verify the correctness of the trace. By design, it ensures absence of deadlocks if the original application did not have any deadlocks for a given trace. Replay also preserves the time taken in terms of the original application's runtime.

5.1 Scala-H-Trace Replay

With the histogram-based trace, the existing parallel replay functionality requires a complete overhaul to cope with statistical data instead of precise data. In our Scala-H-Trace approach, all participating nodes parse the entire trace file during replay. In contrast to ScalaTrace, *all nodes read and interpret all MPI events*. Such interpretation amounts to the selection of a random value following the histogram distribution of any recorded events, for each node in the trace. All nodes "know" the random values used by the other nodes. However, a given node only issues MPI calls if the current node is a member of the participant list in the recorded trace. The interpretation of histogram values for events that are not issued is crucial to provide efficient replay with histograms: It obviates a need to coordinate value selection across nodes and, hence, back-channel communication that might otherwise be required due to randomization, as discussed after the next paragraph.

During the random selection of replay parameters, end-points of MPI_Send/MPI_Isend events are selected. Upon encountering an MPI_Send, once a node identifies itself as a receiver, the receiver node issues a receive call (MPI_Irecv) instead of a MPI_Send. Hence, all receive communication events like MPI_Recv and MPI_Irecv are ignored. Since a particular receiver can also be a sender, only MPI_Irecv calls are internally issued followed by an internal MPI_Wait call when a node rank identifies itself as a receiver of a recorded MPI_Send event. Such internal MPI_Wait calls are issued last, after all ranks have been parsed and all MPI_Send/MPI_Irecv calls have been issued. Any MPI_Wait/Waitall calls in the original recorded trace are ignored.

The selection of a random value for histogram-recorded parameter for any event parameters, such as send/destination rank and data volume, requires that sending and receiving nodes make the same decision on matching end-points for a message exchanged between them. To ensure that sender and receiver nodes agree on their end-points for a message exchange, all nodes use the same random seed during initialization. Hence, all nodes agree on the random value upon each selection of a replay parameter within the range of 0 and total number of elements in the histogram. No coordination via communication is required as all nodes interpret all events in the same order, even if only a subset of (one or more) nodes actually issues an MPI call. Our randomization starts with a common seed across all nodes so that all histogram choices are deterministic. This alleviates communication overhead that would otherwise be required to coordinate sender/receiver selection from histograms. Instead, each node has the same random number sequence and interprets traces in the same manner albeit issuing only calls for events for their respective rank.

The selected random value is internally used to select an appropriate bin. The average value recorded for that bin is then chosen as a parameter for the MPI event. Histograms already record value distributions (iteration counts, send/recv ranks). By randomly selecting bins and values from bins respecting histogram frequencies, replay preserves fair value distributions. The following section discusses distributed coordination for random selection in more detail.

5.2 Challenges for Deterministic Replay: Point-to-Point Messages

The following code snippet is an example of climate simulation in which first 50 nodes work on land simulation and the next 50 nodes work on sea simulation. These simulations are performed in multiple time steps in which nodes perform calculations and communicate the result to surrounding neighbors. The destination nodes and communication volume can vary for land and sea simulation.

```
//Land simulation participants - Node 0 to 49
//Sea simulation participants - Node 50-99

int * resultbuf; //Buffer to hold results
for(timesteps from 0 to 100) {
    int destination[100]; //Array to hold dest. ranks
    int source[100]; //Array to hold source ranks
    do_calculations(resultbuf);
    if(simulation == land) {
        volume = 80 bytes;
        get_my_land_neighbors(destination, source);
    } else {
        volume = 90 bytes;
        get_my_sea_neighbors(destination, source);
    }
    for ( i = 0 to total_neighbors_count ) {
        MPI_Irecv (resultbuf,volume,source[i],...);
        MPI_Isend(resultbuf,volume,destination[i],...);
        MPI_Waitall (...); //Wait for Irecv
        MPI_Waitall (...); //Wait for Isend
    }
}
```

In the code above, all participating nodes perform calculations and communicate the results to corresponding neighbors. All MPI function calls originate from the same call stack but communication volume and source/destination endpoints vary across nodes. This results in perfectly compressed intra-node traces with the following events:

```
RSD1: {MPI_Irecv, MPI_Isend, MPI_Waitall, MPI_Waitall}
PRSD1: {total_neighbors_count, RSD1}
PRSD2: {total_timesteps, PRSD1}
```

Since communication volume and endpoints vary across nodes, the inter-node compression fails for the above section. With an appropriate user-specified merge precision level, communication volume and endpoints are collected in histogram bins during inter-node compression and the trace is compressed across all nodes. Hence, all nodes need to agree upon the message payload (data volume) and send/receive endpoints during replay. With such an agreement between nodes for the selection of a particular value for replay, potential deadlocks could occur.

For example, an original send/receive pair for sender (node 1) / receiver (node 2) might result in arbitrary selection of communication end-points without our distributed coordination scheme. In other words, the sender (node 10) may issue a message to node 20, both randomly selected on node 10, while node 20 interprets the send event as a message originating from node 13 and directed at node 23 per uncoordinated random selection. In such a case, node 10 would deadlock as the message is never received. Similarly, receives (or waits for completion of receives) may deadlock if no corresponding send is ever issued.

Our distributed, coordinated approach to randomized selection ensures that all nodes interpret the send event as originating from node 10 and being directed at node 20. While node 10 issues a send

(and node 20 a receive), all other nodes (13 and 23) will not issue any MPI call. The fact that the original event was a message from node 1 to 2 is probabilistically replayed as a message from one node (here: 10) to another (here: 20), *i.e.*, histograms do result in randomized end-points but retain the original number of messages for the example.

5.3 Challenges for Deterministic Replay: Collective Communication

With the coordinated replay approach, there are situations in which deadlocks can occur due to causal ordering of uncompressed traces. Consider the sample trace below with 8 nodes:

Participants:0-3	Participants:4-7
110-111: MPI_Irecv(2 iterations) 1st. iter: from 0-3, 2nd iter: from 4-7	130: MPI_Irecv from 0-3
112-113: MPI_Send(2 iterations) 1st. iter to 0-3, 2nd iter: to 4-7	131: MPI_Send to 0-3
114: MPI_Wait (2 counts)	132: MPI_Wait
115: MPI_Bcast	133: MPI_Bcast

Table 2. Uncompressed Trace

Both columns in Table 2 contain an uncompressed sequence of MPI events originating from the same call stack. Each MPI call is preceded by a sequence number as recorded by intra-node compression. The first set of nodes, 0 to 3, issues 2 sets of MPI_Irecv/Send/Wait calls followed by a MPI_Bcast. The second set of nodes, 4 to 7, issues only one MPI_Irecv/Send/Wait followed by MPI_Bcast. These events fail to compress due to mismatching Send/Recv counts across different sets of nodes. This results in the final trace with events 110-115 followed by events 130-133.

With the coordinated compression of Scala-H-Trace and its corresponding replay, the MPI_Send in sequence 112 will be issued and the corresponding MPI_Irecv will be issued internally by respective destination ranks as shown in the sample trace. (MPI_Irecv/Wait are ignored during histogram-based replay.) Next, MPI_Bcast will be issued by ranks 0 to 3. This will block ranks 0 to 3 until the corresponding MPI_Bcast (seq. 133) is issued by ranks 4-7. But before issuing the broadcast in sequence 133, nodes with ranks 4-7 issue the MPI_Send in sequence 131 with destination 0-3. Since nodes of ranks 0-3 are already blocked in MPI_Bcast (seq. 115), they cannot issue an corresponding internal MPI_Irecv, eventually leading to a situation in which nodes 4-7 cannot proceed to other events. This situation occurs frequently. In many scientific applications, two sets of nodes can execute different sections of a program leading to compression failure interspersed with collectives, such as barriers. This causal ordering of events in the trace can lead to deadlock when replayed using the above approach. We employ a novel design for the inter-node compression algorithm to forcibly merge collectives even if an entire PRSD loop of other events does not merge properly.

Inter-node compression attempts to match an entire sequence of events subject to the same PRSD loop across nodes. Even if there is a single mismatch, the entire sequence would conventionally not be merged but rather be written consecutively as shown in the sample trace above. We employ a novel design for inter-node compression to greedily merge any subset of events, *e.g.*, collectives inside a loop. We then rearrange other communication calls with collectives as synchronization points. This ensures that deadlocks cannot be introduced during the replay of MPI events.

A prove showing that our novel merge algorithm, which rearranges non-merging communication calls with a collective as a synchronization point, will not introduce deadlocks is provided in a technical report due to space constraints [20].

6. Experimental Results

We evaluated Scala-H-Trace in three aspects: (1) its effectiveness of trace file compression, and (2) its statistical trace replay feature and (3) its trace compression sensitivity to merge precision level settings. Experiments (1) and (2) utilize both the histogram compression approach and the WLCS-based recursive inter-node compression algorithm. Most of our experiments were conducted on Jaguar, the Cray XT4 system at ORNL. Each of compute node features a 2.1 GHz quad-core AMD Opteron 1354 processor and 8GB of DDR2 memory. The login nodes run a full-featured Linux version while the compute nodes run the Compute Node Linux microkernel. Due to unavailability of Jaguar in the final experimentation phase, the MG experiments were conducted on Jugene, an IBMBlue Gene/P system with 73,728 compute nodes and 294,912 cores, 2 GB memory per node, and the 3D torus and global tree interconnection networks.

We analyze the efficacy of Scala-H-Trace using a production-scale application, the Parallel Ocean Program (POP) [6], as the main challenge. The Parallel Ocean Program (POP) is an ocean circulation model developed at Los Alamos National Laboratory. Our experiments exercise a one degree grid resolution in which the problem size is 320x384 blocks and the individual block size is 5x6 resulting in a total of 4096 (64x64) blocks distributed to individual nodes. POP exhibits non-SPMD behavior, which leads to trace file size increases with the number of nodes for conventional trace tools, including ScalaTrace. POP is a large-scale application with challenging communication patterns. There five different dominant patterns equivalent to five micro-benchmarks, yet in combined complexity. Hence, this application provides an opportunity to show-case the effectiveness of histogram-based trace collection of Scala-H-Trace. We conducted experiments by varying the maximum number of blocks assigned to each node.

We further utilize the CG and MG benchmarks from the NAS parallel benchmark suite for inputs sizes C to study the efficacy of Scala-H-Trace for different types of application behavior. Both CG and MG mostly exhibit SPMD behavior but differ significantly in the communication pattern impacting the compression effectiveness during trace collection. These benchmarks are also selected from the NAS benchmarks as these two were the challenging cases for ScalaTrace’s lossless compression: Both were reported to result in sub-linear increases in the trace file size for ScalaTrace [12].

6.1 Trace Compression Effectiveness

We collected traces based on two different compression techniques. First, the original ScalaTrace is used, in which loop details and parameter values are captured losslessly and inter-node trace compression is performed with the coarse-grained matching scheme. Second, our novel histogram-based trace compression featuring Scala-H-Trace is used, in which trace information is collected in histograms for events and parameters that otherwise would not have compressed with the lossless trace compression, and inter-node compression is performed recursively. Trace file sizes are assessed under strong scaling, where we vary the number of nodes while keeping the overall problem size fixed. Lossless traces, obtained from ScalaTrace, are useful to identify exact details of the communication and I/O patterns exhibited by the application. Histogram-based traces are obtained from Scala-H-Trace, attempting to capture lossless information for trace events where feasible while non-matching events are recorded in histogram bins. We hypothesize that histograms suffice to capture the “big picture” of the application behavior and will assess this claim by accuracy of replay times relative to the original application. For applications exhibiting non-SPMD behavior, such as POP, histogram-based trace collection (Scala-H-Trace) collects concise traces, which could not otherwise be obtained with lossless trace compression (ScalaTrace).

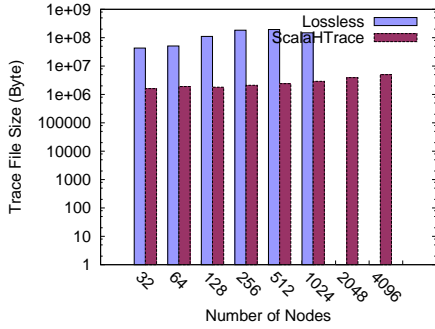


Figure 2. Parallel Ocean Program

Figure 2 depicts the trace file size for both lossless and histogram-based traces when varying the number of nodes. Note that the y axis is in log scale. Since POP exhibits non-SPMD behavior, we observe a linear increase in the trace file size in the case of lossless trace collection up to 256 nodes. The trace file size then stabilizes for 512 nodes and even declines for 1024 nodes. We identified that the timestep behavior becomes more regular at these levels, resulting in more effective inter-node compression. But we again observed an increasing trend in the case of 2048 nodes. For 2048 nodes and above, we could not even collect traces anymore as the trace file size was growing unmanageably fast and the time taken to merge hundreds of megabytes of per-node traces became prohibitive. With the histogram-based approach, there is a sub-linear increase in the trace file size. Moreover, histogram-based trace files are two orders of magnitude smaller than the lossless traces. This considerable reduction is obtained by aggressive compression of events and their associated function parameters in histograms. This clearly shows the efficacy of Scala-H-Trace to collect concise trace files even with applications exhibiting irregular behavior.

Figure 3 depicts trace file size for the CG benchmark. We observe an interesting trend in CG in which the trace file size for lossless traces is consistently 50% less than that of the histogram traces up to 1024 nodes, yet sizes match at 2048 node. Even though lossless traces are initially smaller than histogram traces, there is a consistent increase in the trace file size for the lossless case. In contrast, the size of histogram traces is almost constant with the increase in number of nodes. For lossless traces, non-matching function parameters for events with the same call stack are collected in vectors associated with a participant rank list. This representation is more concise than histograms for smaller number of nodes. With thousands of nodes, the vector-participant list pair for each event has increased in size to where it is at par with histogram traces. Unlike vector-participant lists, histogram representation is constant with the increase in number of nodes as the number of bins is fixed during the application run and even the outlier participant rank information is absorbed as constants in bins. It should also be noted that the trace file size for CG is in the order of hundreds of kilobyte. For larger applications with a similar communication behavior as CG yet with trace file sizes in hundreds of megabytes, such a linear (or even sub-linear) growth for lossless traces may simply not be scalable due to inter-node merge overheads, as discussed.

Figure 4 depicts the results for MG. MG exhibits a double nested 7-point stencil communication pattern in the 3D space. Due to the regular communication pattern and data-independent program behavior, compressing the MPI parameter values of MG works well for both lossless and histogram-based approaches. However, due to the slightly diverged per-node program behavior within a loop, the original inter-node compression algorithm of ScalaTrace failed to merge across communication groups. This caused trace sizes to increase linearly with the number of nodes.

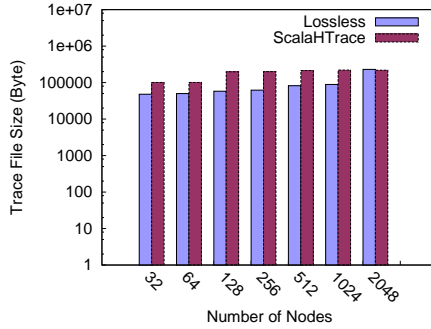


Figure 3. CG Benchmark

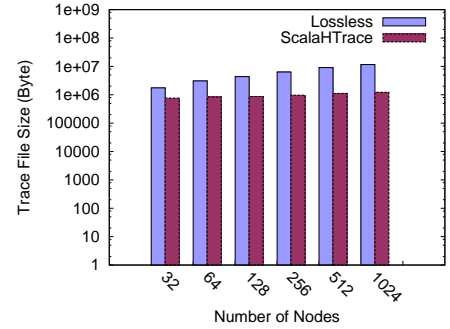


Figure 4. MG Benchmark

In contrast, with our novel fine-grained recursive approach, similar PRSDs are merged and the trace size grows sub-linearly, i.e., by a factor of two as the number of nodes is increased by a factor of 64.

6.2 Histogram-based Trace Replay

We studied the replay effectiveness of histogram-based traces by comparing the original application execution time with the time taken to replay the recorded events. We discuss the effectiveness of the distributed approach of replaying statistical traces. We also discuss the impact of trace compression on the replay behavior for histogram-based traces. We show that even with statistical histogram-based traces, replay can still be employed to check the correctness of a recorded trace and also to perform “what-if” analysis for system procurements.

Figure 5 depicts the replay time of histogram-based trace events compared to that of the application’s original execution time. The compressed traces are fully forced histogram trace events where any non-matching function parameters or loop iterations are collected as histograms. Even with these traces, we see that the replay time for traces collected for 32-512 number of nodes are within 5% of the original execution time (with the exception of replay time for 128 nodes). Replay time accuracy drops to 12% for 1024 and 2048 nodes. Due to our experiment with strong scaling for POP, the original execution time for both 1024 and 2048 nodes (30 seconds) is much lower than that for fewer nodes (>100 seconds) so that even small deviations in absolute values during replay increase the error percentage. We conjecture that such deviations are unrealistic as POP for this particular input does not scale beyond 512 nodes so that such short times are unrealistic. Similarly, this problem would not occur under weak scaling as runtimes would not decrease with larger number of nodes. Overall, we observe that the replay time is close to the original execution time, even for fully forced histograms, due to two reasons: (1) Since our histograms are dynamically balanced, a random value selected from histogram bins during replay falls within a commonly used value range in the original application run. (2) The inter-node compression algorithm effectively merged events across nodes so that communication calls are not split in the trace file.

We observe nearly 50% deviation in the case of 128 nodes for POP. To investigate the cause, we calculated the time spent by nodes in other communication calls and found that some nodes are engaged in more communication calls than the majority of nodes. This created load imbalances where the remaining nodes wait at collectives for nodes participating in larger number of events.

Figure 6 depicts the replay time for the CG benchmark. In the majority of cases, the replay time is with 10% to 15% of the original application runtime. Since the original execution time for CG is within 10 seconds for 1024 and 2048 nodes, even small changes in the absolute replay execution time increase the error percentage considerably. The replay time deviation can be attributed to the loop iterations recorded in histograms. Again, CG stops scaling at

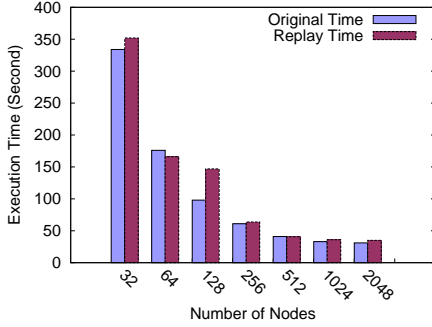


Figure 5. POP Replay

512 nodes for this input size so that larger application runs are unrealistic. Furthermore, if the random loop iteration selected from histograms is close to the maximum value, then all nodes participate in more communication calls than in the original application. This is a fundamental trade-off between accuracy and manageability of trace file sizes.

The replay time for MG benchmark is depicted in Figure 7. The averaged inaccuracy is 8.2% under strong scaling. We observe up to 34.2% inaccuracy for 2,048 nodes. This outlier is due to an excessively short runtime of 3.8s with an absolute error of just 1.3s. For 1,024 nodes, this decreases to 12.5% and for 512 to 5.3% and so on indicating that the problem is only due to excessively short runtimes. After discarding this outlier due to strong scaling limitations, the replay timing accuracy for MG is high. As discussed before, with the recursive inter-node trace compression, we are able to achieve a nearly constant trace sizes for MG even without the histogram-based probabilistic approach. Due to the elimination of the imprecision, the timing behavior of the trace replay highly resembles that of the original MG benchmark.

With the exception for MG, which fares equally well with histogram-based compression, replay for Scala-H-Trace generated traces with forced histograms results in runtimes that are within 12-15% of the original application for most cases. This result is interesting as forced histograms are equivalent to a 0% merge precision level, which is the most aggressive compression possible with Scala-H-Trace. More accurate replay may result from higher precision levels at the cost of slightly larger traces, as discussed next.

6.3 Trace Sensitivity Study

Finally, we study the effect of varying merge precision levels on trace file sizes. This experiment serves as an illustration for the benefits of user-specified merge precision levels as a means to steer compression, which should improve as precision decreases. Merge precision levels provide a tunable parameter to select target trace file size as required by operating environments or performance analysis experiments. Lossless traces may be desirable for exact analysis of application behavior and users with access to excessive storage may happily utilize this feature even if the trace file size becomes large. When desiring a more compact trace file and when inter-node merges become prohibitive for lossless traces, users can decrease merge precision level to target a desired trace file size and tracing overhead.

Figure 8 depicts the impact of verifying merge precision levels on the final trace file size. We fixed the number of nodes to 512 for POP and measured trace file sizes for varying merge precision levels. We observe that even with a small decrease in the merge precision from 100% to 95%, the trace size reduced by more than a factor of three. This significant reduction is due to merging events with varying numbers of loop iterations for the timestep in POP. With lossless traces, two different sets of loops with the exact same events fail to compress due to varying numbers of loop

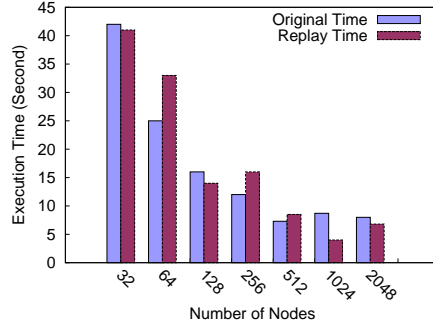


Figure 6. CG Replay

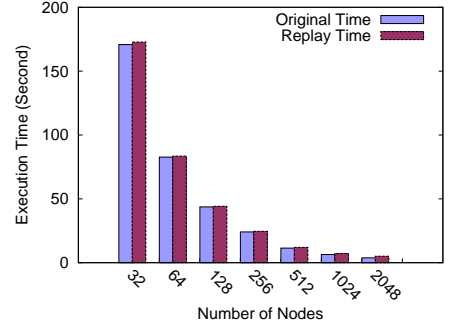


Figure 7. MG Replay

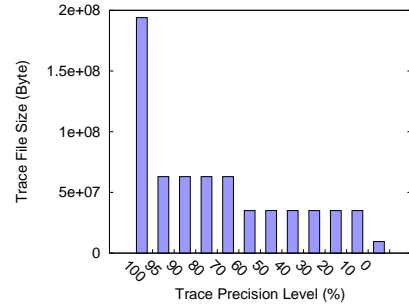


Figure 8. POP Trace Sensitivity for 512 nodes

iterations across the timestep. This variation is data dependent and induced by computation as explained in the section 4.1. The trace file size is constant up to a 70% merge precision level. At 60% precision, sizes drop again by almost 50%. This second reduction has been attributed to function parameters collected as histograms. Many events with varying function parameters are not combined under lossless tracing or result in vectors collected to represent varying parameters. Both contribute to the significant increase in the trace file size and prevent trace scalability with increasing numbers of nodes. Finally, another three-fold reduction in trace sizes is observed for forced histograms (0% merge precision level). At the 0% merge precision level, all non-matching values are represented as histograms, which results in the most concise trace possible with Scala-H-Trace. Overall, sensitivity experiments to merge precision levels show that small reductions in precision can significantly reduce the overall trace sizes. This particularly aids production-scale codes like POP, which otherwise cannot be feasibly traces without loss of information for thousands of nodes.

7. Related Work

There are several tools, such as TAU [16], Vampir [10], Paraver [13] and SCALASCA [4], that capture communication and/or I/O trace events using library instrumentation similar to Scala-H-Trace. But only few employ trace compression techniques to control the trace file size. Many of these tools depend on zlib for compression, which compresses blocks of data without preserving the structure of the trace, *i.e.*, post-processing/analysis only becomes feasible after decompression. This also increases the memory requirements, effectively rendering trace analysis infeasible on commodity desktops or laptops and sometimes even high-end workstations, depending on the uncompressed trace size. Unlike these techniques, ScalaTrace [12] compresses traces while preserving the trace structure in terms of order of events. As a result, post-processing/analysis can be performed without decompression. We utilize this concept of structure preserving compression in Scala-H-Trace. Yet while ScalaTrace and any of the aforementioned tracing tools record lossless traces with a subset or all event parameters, Scala-H-Trace establishes a different methodology. Parameters, event frequencies and

participant lists of nodes are recorded as histograms when lossless compressing cannot be established within a user-specified merge precision level. Employing statistical methods results in more concise traces even for non-SPMD programs at the expense of loss of information. Our replay tool uses an algorithm to issue events on-the-fly using the compressed traces, much like ScalaTrace. Yet recorded parameters are replayed in a probabilistic manner, which creates novel challenges that are met by our distributed approach to coordinate event replay across nodes.

The mpiP tool, a lightweight profiling library for MPI applications, collects statistical information about MPI functions [17]. It collects aggregate metrics like number of MPI events issued by the application and average execution times. This is useful to provide very high-level information on communication and I/O calls. Scala-H-Trace, in contrast, captures all events in traces and employs more sophisticated histogram bins only when the need arises for applications exhibiting non-SPMD behavior. Beside the histogram information, we also record outlier information associated with each bin to detect communication bottlenecks and to provide a “big picture” of communication and I/O events in applications.

Kluge *et al.* [7] employ pattern matching techniques similar to ours to capture POSIX I/O calls in parallel programs. Unlike our approach, they perform post-mortem pattern matching only after collecting the application traces. They read the collected trace and create an I/O dependency graph thereby preserving the event order to do pattern matching. Even though post mortem pattern matching reduces the trace volume, this approach limits its usefulness in memory constrained systems like the IBM BlueGene family. Without online compression, either the memory footprint increases by holding the recorded trace or trace events are frequently written to disk, which affects the application execution behavior. They also do not employ pattern matching across nodes so that they require a trace file per node. This limits their approach in that they struggle with applications utilizing thousands of nodes due to parallel file system constraints. Our approach is immune to such limitations as a single trace file captures the behavior of all nodes with statistical information on a per-event and per-parameter basis.

Gao *et al.* [3] developed an event trace compression technique that performs static analysis on the application binary and collects loops and functions as structures. Along with these structure, a path grammar is constructed on-the-fly. Path grammars are then utilized to encode paths taken during execution. These structures are compressed individually and stored. Even the iteration count is stored along with the compressed structure traces. This loosely resembles the RSD and PRSD technique used in related work [5, 8, 11, 15]. But unlike Gao *et al.*'s work, our tool does not require the construction of grammars for individual applications separately. Our work employs a generalized trace compression approach based on call path stacks and records parameters exploiting statistical means. It is sufficient to link the tool library along with the application to collect traces. This generalization also enables comparative trace studies between two different applications.

8. Conclusion

We presented the design and implementation of Scala-H-Trace, which provides novel capabilities for more aggressive trace compression than any previous approach. Scala-H-Trace utilizes histograms based on a user-specified merge precision level. It features a distributed approach to deterministically replay statistical histogram traces where events are reissued without decompressing the original trace file. Experimental results demonstrate the ability to obtain a single, near constant sized trace file, even for production-scale scientific applications such as POP with non-SPMD behavior. Results also show that replay time for traced events are within 12%-15% of the original application execution time in majority of the cases, even for the most aggressive “forced” histograms.

References

- [1] N. Adiga and et al. An overview of the BlueGene/L supercomputer. In *Supercomputing*, Nov. 2002.
- [2] H. Brunst, D. Kranzlmüller, and W. Nagel. Tools for Scalable Parallel Program Analysis - Vampir NG and DeWiz. *The International Series in Engineering and Computer Science, Distributed and Parallel Systems*, 777:92–102, 2005.
- [3] X. Gao, A. Snaveley, and L. Carter. Path grammar guided trace compression and trace approximation. *High-Performance Distributed Computing, International Symposium on*, 0:57–68, 2006.
- [4] M. Geimer, F. Wolf, B. J. N. Wylie, E. Abraham, D. Becker, and B. Mohr. The scalasca performance toolset architecture. In *International Workshop on Scalable Tools for High-End Computing*, June 2008.
- [5] P. Havlak and K. Kennedy. An implementation of interprocedural bounded regular section analysis. *IEEE Transactions on Parallel and Distributed Systems*, 2(3):350–360, July 1991.
- [6] P. W. Jones, P. H. Worley, Y. Yoshida, J. B. White, III, and J. Levesque. Practical performance portability in the parallel ocean program (pop): Research articles. *Concurr. Comput. : Pract. Exper.*, 17(10):1317–1327, 2005.
- [7] M. Kluge, A. Knüpfer, M. Müller, and W. E. Nagel. Pattern matching and i/o replay for posix i/o in parallel programs. In *Euro-Par '09: Proceedings of the 15th International Euro-Par Conference on Parallel Processing*, pages 45–56, Berlin, Heidelberg, 2009. Springer-Verlag.
- [8] J. Marathe and F. Mueller. Detecting memory performance bottlenecks via binary rewriting. In *Workshop on Binary Translation*, Sept. 2002.
- [9] MPI-2: Extensions to the message-passing interface. July 1997.
- [10] W. E. Nagel, A. Arnold, M. Weber, H. C. Hoppe, and K. Solchenbach. VAMPIR: Visualization and analysis of MPI resources. *Supercomputer*, 12(1):69–80, 1996.
- [11] M. Noeth, F. Mueller, M. Schulz, and B. R. de Supinski. Scalable compression and replay of communication traces in massively parallel environments. In *International Parallel and Distributed Processing Symposium*, Apr. 2007.
- [12] M. Noeth, F. Mueller, M. Schulz, and B. R. de Supinski. Scalatrace: Scalable compression and replay of communication traces in high performance computing. *Journal of Parallel Distributed Computing*, 69(8):969–710, Aug. 2009.
- [13] V. Pillet, J. Labarta, T. Cortes, and S. Girona. PARAVER: A tool to visualise and analyze parallel code. In *Proceedings of WoTUG-18: Transputer and occam Developments*, volume 44 of *Transputer and Occam Engineering*, pages 17–31, Apr. 1995.
- [14] The parallel ocean program (POP), 1996. <http://climate.lanl.gov/Models/POP/>.
- [15] P. Ratn, F. Mueller, B. R. de Supinski, and M. Schulz. Preserving time in large-scale communication traces. In *International Conference on Supercomputing*, pages 46–55, June 2008.
- [16] S. S. Shende and A. D. Malony. The tau parallel performance system. *Int. J. High Perform. Comput. Appl.*, 20(2):287–311, 2006.
- [17] J. Vetter and M. McCracken. Statistical scalability analysis of communication operations in distributed applications. In *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 2001.
- [18] J. S. Vetter and B. R. de Supinski. Dynamic software testing of mpi applications with umpire. In *Supercomputing*, page 51, 2000.
- [19] K. Vijayakumar, F. Mueller, X. Ma, and P. C. Roth. Scalable i/o tracing and analysis. In *PDSW '09: Proceedings of the 4th Annual Workshop on Petascale Data Storage*, pages 26–31, New York, NY, USA, 2009. ACM.
- [20] X. Wu, K. Vijayakumar, F. Mueller, X. Ma, and P. C. Roth. Probabilistic communication and i/o tracing with deterministic replay at scale. Technical Report TR 2011-6, Dept. of Computer Science, North Carolina State University, 2011.

Appendix

We next prove that our novel merge algorithm, which rearranges non-merging communication calls with a collective as a synchronization point, will not introduce deadlocks. We assume that the original application is deadlock free (which is reasonable since the event trace was collected from a terminating application) and provide the proof below:

THEOREM 1. *A replayed trace of a program with events reordered to synchronized collectives does not result in deadlock if the original trace was deadlock free.*

Proof: Follows from Lemmas 1-3. ■

LEMMA 1. *Lemma 1: A replayed trace of a program with only collectives will not deadlock.*

Proof: By construction of traces, all recorded participants engage in a collective during replay in the same order as recorded. Collectives are blocking. Hence, all participants complete a collective at (nearly) the same time (as collectives provide global/group synchronization). ■

LEMMA 2. *A replayed trace of a program with only point-to-point communication will not deadlock.*

Proof: Blocking/non-blocking sends are replayed in the same order as recorded. The corresponding receives are issues as non-blocking receives in the same order that the sends were issued. Once all non-blocking receives of a PRSD have been issued, waits are issued on all pending non-blocking receives. Hence, if the original trace did not deadlock, replayed point-to-point messages with identical receive ordering followed by wait ordering relative to sender ordering will not deadlock either. ■

LEMMA 3. *A replayed traces with mixed events of collectives and point-to-point messages will not deadlock.*

Proof: (a) Assume a trace with alternating phases of only point-to-point messages and only collectives. Since collectives provide a fence where all point-to-point messages are consumed prior to a collective, replaying such a trace is deadlock free for each region by Lemmas 1 and 2 and thereby also for the entire trace since each phase is causally independent.

(b) If a point-to-point message is crossing collectives (sent before but received after a collective across a pair of nodes that also participates in the collective), then the same send will be issued during replay before the corresponding collective, but the corresponding non-blocking receive will also be issued before the collective followed by a wait. Hence, if the application with its traced events was deadlock free, the replay will also be deadlock free.

Let $S = s_1, \dots, s_n$ be the set of event streams over n nodes where $s_i = e_1, \dots, e_{m_i}$ are ordered sequences of of point-to-point or collective communication events.

(c) By structural induction over (b): Let $s_i \in S$ and $s_k \in S$ be event streams with alternating phases of collectives and point-to-point messages or point-to-point messages crossing collectives as in proof step (a) or (b) of Lemma ??, respectively. Then let S^+ be the induced set of event streams of the traced application run with an additional point-to-point message sent from node i to node k denoted as $e_x^+ \rightarrow f_y^+$ where $e_x^+ \in s_i^+$, $f_y^+ \in s_k^+$ and $s_i^+ \in S^+$, $s_k^+ \in S^+$. Furthermore, let $\{e_x, e_{x+1}\} \subseteq s_i$ and $\{f_y, f_{y+1}\} \subseteq s_k$ be subsequences such that $\{e_x, e_x^+, e_{x+1}\} \subseteq s_i^+$ and $\{f_y, f_y^+, f_{y+1}\} \subseteq s_k^+$ for arbitrary $1 \leq x \leq m_i$ and $1 \leq y \leq m_k$. Furthermore, let S' be the set of event streams with reordered non-blocking receives in place of sends corresponding to S followed by waits that is deadlock free under replay.

The corresponding induced set of replayed event streams, S'^+ , is then $\{e_x, e_x^+, e_{x+1}\} \subseteq s_i'^+$ and $\{f_y, f_y^+, f_{y+1}\} \subseteq s_k'^+$ for $s_i'^+ \in S'^+$ and $s_k'^+ \in S'^+$. Since the application was deadlock free for S and S^+ and replay was deadlock free for S' , replay is also deadlock free for S'^+ since replay preserves ordering of event sequences $s_i'^+$ and $s_k'^+$ with respect to the send order of the application, i.e., a non-blocking receive (followed by a wait) is issued at the receiving node after $f_x \in s_k'^+$. ■

Notice that we fold the non-blocking receive and the corresponding wait into f_y^+ in the above notation to facilitate readability without loss of generality.