



**HAL**  
open science

# Damaris: Leveraging Multicore Parallelism to Mask I/O Jitter

Matthieu Dorier, Gabriel Antoniu, Franck Cappello, Marc Snir, Leigh Orf

► **To cite this version:**

Matthieu Dorier, Gabriel Antoniu, Franck Cappello, Marc Snir, Leigh Orf. Damaris: Leveraging Multicore Parallelism to Mask I/O Jitter. [Research Report] RR-7706, 2011, pp.33. inria-00614597v2

**HAL Id: inria-00614597**

**<https://inria.hal.science/inria-00614597v2>**

Submitted on 6 Dec 2011 (v2), last revised 9 Apr 2012 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Damaris: Leveraging Multicore Parallelism to Mask I/O Jitter

Matthieu Dorier, Gabriel Antoniu,  
Franck Cappello, Marc Snir, Leigh Orf

**RESEARCH  
REPORT**

**N° 7706**

December 2011

Project-Teams KerData  
Joint INRIA/UIUC Laboratory  
for Petascale Computing





## Damaris: Leveraging Multicore Parallelism to Mask I/O Jitter

Matthieu Dorier<sup>\*</sup>, Gabriel Antoniu<sup>†</sup>,  
Franck Cappello<sup>‡</sup>, Marc Snir<sup>§</sup>, Leigh Orf<sup>¶</sup>

Project-Teams KerData  
Joint INRIA/UIUC Laboratory for Petascale Computing

Research Report n° 7706 — December 2011 — 30 pages

**Abstract:** With exascale computing on the horizon, the performance variability of I/O systems represents a key challenge in sustaining high performance. In many HPC applications, I/O is concurrently performed by all processes, which leads to I/O bursts. This causes resource contention and substantial variability of I/O performance, which significantly impacts the overall application performance. In this paper, we first explore the influence of user-configurable parameters and I/O approaches on write performance variability. We then propose a new approach, called Damaris, which leverages dedicated I/O cores on each multicore SMP node to efficiently perform asynchronous data processing and I/O. We evaluate our approach on two different platforms including the Kraken Cray XT5 supercomputer, with the CM1 atmospheric model, which is one of the target HPC applications for the Blue Waters project. By gathering data into large files while avoiding synchronization between cores, our solution brings several benefits: 1) it increases the sustained write throughput by a factor of almost 15; 2) it provides almost 70% overall application speedup on 9K cores; 3) it fully hides I/O-related costs; 4) it enables a 600% compression ratio without any additional overhead, leading to a major reduction of storage requirements.

**Key-words:** Exascale Computing, Multicore Architectures, I/O, Variability, Dedicated Cores

---

<sup>\*</sup> ENS Cachan Brittany, IRISA - Rennes, France. [matthieu.dorier@irisa.fr](mailto:matthieu.dorier@irisa.fr)

<sup>†</sup> INRIA Rennes Bretagne-Atlantique - France. [gabriel.antoniu@inria.fr](mailto:gabriel.antoniu@inria.fr)

<sup>‡</sup> INRIA Saclay - France, University of Illinois at Urbana-Champaign - IL, USA. [fci@lri.fr](mailto:fci@lri.fr)

<sup>§</sup> University of Illinois at Urbana-Champaign - IL, USA. [snir@illinois.edu](mailto:snir@illinois.edu)

<sup>¶</sup> Central Michigan University - MI, USA. [leigh.orf@cmich.edu](mailto:leigh.orf@cmich.edu)

**RESEARCH CENTRE  
RENNES – BRETAGNE ATLANTIQUE**

Campus universitaire de Beaulieu  
35042 Rennes Cedex

## Damaris: Exploitation du Parallelisme Multi-cœur pour Masquer la Variabilité des E/S

**Résumé :** Alors que l'ère de l'exascale se profile à l'horizon, la variabilité des performances des systèmes d'entrées-sorties (E/S) présente un défi majeur pour permettre d'atteindre des performances élevées et soutenues. Dans de nombreuses applications HPC, les E/S sont effectuées par tous les processus de manière concurrente. Cela engendre des pics d'utilisation des E/S, créant des conflits d'accès aux ressources et une variabilité substantielle dans les performances des E/S, impactant les performances globales de l'application de manière significative. Dans ce rapport, nous étudions l'influence de paramètres configurables par l'utilisateur sur la variabilité des E/S. Nous proposons une nouvelle approche, nommée Damaris, qui exploite des cœurs dédiés aux E/S sur chaque nœud SMP de manière à permettre des traitements et des mouvements de données efficaces et asynchrones. Nous évaluons notre approche sur deux plateformes, notamment sur le supercalculateur Kraken Cray XT5, avec l'application de modélisation atmosphérique CM1, une des applications visée par le projet Blue Waters. En rassemblant les données dans de plus gros fichiers tout en évitant la synchronisation entre les cœurs, notre solution apporte plusieurs avantages : 1) elle multiplie par 15 le débit atteint; 2) elle permet une accélération globale de l'application de presque 70% sur 9K cœurs; 3) elle cache tous coûts liés aux E/S; 4) elle permet d'atteindre un taux de compression de 600% sans surcoût, permettant ainsi une réduction importante de l'espace de stockage nécessaire.

**Mots-clés :** Exascale, Architectures Multi-cœurs, E/S, Variabilité, Cœurs Dédiés

## Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>4</b>  |
| <b>2</b> | <b>Background and related work</b>                               | <b>5</b>  |
| 2.1      | Understanding I/O jitter . . . . .                               | 5         |
| 2.2      | Approaches to I/O management in HPC simulations . . . . .        | 6         |
| <b>3</b> | <b>Impact of user-controllable parameters on I/O variability</b> | <b>7</b>  |
| 3.1      | Methodology . . . . .  | 7         |
| 3.2      | Results with the IOR benchmark . . . . .                         | 9         |
| 3.2.1    | Impact of the data size . . . . .                                | 9         |
| 3.2.2    | Impact of the striping policy . . . . .                          | 9         |
| 3.2.3    | Impact of the I/O approach . . . . .                             | 10        |
| 3.2.4    | Impact of the number of writers . . . . .                        | 12        |
| <b>4</b> | <b>The dedicated core approach</b>                               | <b>13</b> |
| 4.1      | Principle . . . . .  | 14        |
| 4.2      | Architecture . . . . .   | 14        |
| 4.3      | Key design choices . . . . .                                     | 15        |
| 4.4      | Client-side API . . . . .  | 16        |
| <b>5</b> | <b>Experimental results with CM1</b>                             | <b>18</b> |
| 5.1      | The CM1 application . . . . .                                    | 18        |
| 5.2      | Platforms and configuration . . . . .                            | 18        |
| 5.3      | Non-overlapping I/O approaches . . . . .                         | 19        |
| 5.3.1    | File-per-process method . . . . .                                | 19        |
| 5.3.2    | Using collective I/O and grouped I/O . . . . .                   | 20        |
| 5.4      | Our proposal: using dedicated I/O cores . . . . .                | 21        |
| <b>6</b> | <b>Discussion and related work</b>                               | <b>23</b> |
| 6.1      | Do all my cores are really needed? . . . . .                     | 23        |
| 6.2      | Positioning Damaris in the “I/O landscape” . . . . .             | 24        |
| <b>7</b> | <b>Conclusions</b>   | <b>25</b> |
| <b>A</b> | <b>Evaluating Damaris on a Power5 cluster</b>                    | <b>29</b> |
| A.1      | Platforms and configuration . . . . .                            | 29        |
| A.2      | Experimental results . . . . .                                   | 29        |

## 1 Introduction

As HPC resources approaching millions of cores become a reality, science and engineering codes invariably must be modified in order to efficiently exploit these resources. A growing challenge in maintaining high performance is the presence of high variability in the effective throughput of codes performing input/output (I/O) operations. A typical I/O approach in large-scale simulations consists of alternating computation phases and I/O phases. Often due to explicit barriers or communication phases, all processes perform I/O at the same time, causing network and file system contention. It is commonly observed that some processes exploit a large fraction of the available bandwidth and quickly terminate their I/O, then remain idle (typically from several seconds to several minutes) waiting for slower processes to complete their I/O. This jitter has been observed at 1000-way scale, where measured I/O performance can vary by several orders of magnitude between the fastest and slowest processes [32]. This phenomenon is exacerbated by the fact that HPC resources are typically used many concurrent I/O intensive jobs. This creates file system contention between jobs, further increases the variability from one I/O phase to another and leads to unpredictable overall run times.

While most studies address I/O performance in terms of aggregate throughput and try to improve this metric by optimizing different levels of the I/O stack from the file system to the simulation-side I/O library, few efforts have been made in addressing I/O jitter. Yet it has been shown [32] that this variability is highly correlated with I/O performance, and that statistical studies can greatly help addressing some performance bottlenecks. The origins of this variability can substantially differ due to multiple factors, including the platform, the underlying file system, the network, and the application I/O pattern. For instance, using a single metadata server in the Lustre file system [13] causes a bottleneck when following the file-per-process approach (described in Section 2.2), a problem that PVFS [8] or GPFS [28] are less likely to exhibit. In contrast, byte-range locking in GPFS or equivalent mechanisms in Lustre cause lock contentions when writing to shared files. To address this issue, elaborate algorithms at the MPI-IO level are used in order to maintain a high throughput [26].

The contribution of this paper is twofold. We first perform an in-depth experimental evaluation of I/O jitter for various I/O-intensive patterns exhibited by the Interleaved Or Random (IOR) benchmark [29]. These experiments are carried out on the Grid'5000 testbed [5] with PVFS as the underlying file system, and on the Kraken Cray XT5 machine [2] at the National Institute for Computational Sciences (NICS), which uses Lustre. As opposed to traditional studies based on statistical metrics (average, standard deviation), we propose a new graphical method to study this variability and show the impact of some user-controllable parameters on this variability. Grid'5000 was used because it offers a root access to reserved nodes, and allows to deploy and tune our own operating system and file system. Thus we experiment on a clean environment by reserving entire clusters in order not to interfere with other users. In contrast, the Kraken larger-scale production supercomputer (currently 11<sup>th</sup> in the Top500) is shared with other users.

Our second contribution aims to make one step further: we propose an approach that hides the I/O jitter exhibited by most widely used approaches to

I/O management in HPC simulations (described in Section 2). Based on the observation that a first level of contention occurs when all cores of a multicore SMP node try to access the network for intensive I/O at the same time, our new approach to I/O optimization, called Damaris (Dedicated Adaptable Middleware for Application Resources Inline Steering), leverages one core in each multicore SMP node to perform data processing and I/O asynchronously. This key design choice builds on the observation that it is often not efficient to use all cores for computation because of memory bandwidth limitations. Consequently, reserving one core for kernel tasks such as I/O management may help reducing jitter. Our middleware takes into account user-provided information related to the application behavior, the underlying I/O systems and the intended use of the output in order to perform “smart” I/O and data processing within SMP nodes.

The Damaris approach completely removes I/O jitter by performing scheduled asynchronous data processing and movement. We evaluate it with the Bryan Cloud Model, version 1 (CM1) [6] (one of the target applications for the Blue Waters [1] project) on Grid’5000 and on Kraken. We compare our approach to the classical file-per-process and collective-I/O approaches. By gathering data into large files while avoiding synchronization between cores, our solution increases the I/O throughput by a factor of 6 compared to standard approaches on Grid’5000, and by a factor of 15 on Kraken, hides all I/O-related costs, and enables a 600% compression ratio without any additional overhead.

This paper is organized as follows: Section 2 presents the background and motivations of our study. In Section 3 we run a set of benchmarks in order to characterize the correlation between the I/O jitter and a set of relevant parameters: the file striping policy, the I/O method employed (file-per-process or collective), the size of the data and the number of writers. Our new approach is presented in Section 4 together with an overview of its design and its API. We evaluate this approach with the CM1 atmospheric simulation running on 672 cores of the *paraphuie* cluster on Grid’5000, and on up to 9216 cores on Kraken. Section 6 finally presents a study of the theoretical and practical benefits of the approach, and positions it with respect to related works.

## 2 Background and related work

### 2.1 Understanding I/O jitter

Over the past several years, chip manufacturers have increasingly focused on multicore architectures, as the increase clock frequencies for individual processors has leveled off, primarily due to substantial increases in power consumption. These increasingly complex systems present new challenges to programmers, as approaches which work efficiently on simpler architectures often do not work well on these new systems. Specifically, high performance variability across individual components becomes more of an issue, and it can be very difficult to track the origin of performance weaknesses and bottlenecks. While most efforts today address performance issues and scalability for specific types of workloads and software or hardware components, few efforts are targeting the causes of performance variability. However, reducing this variability is critical, as it is an effective way to make more efficient use of these new computing platforms



through improved predictability of the behavior and of the execution run time of applications. In [30], four causes of jitter are pointed out:

1. Resource contention within multicore SMP nodes, caused by several cores accessing shared caches, main memory and network devices.
2. Communication, causing synchronization between processes that run within a same node or on separate nodes. In particular, access contention for the network causes collective algorithms to suffer from variability in point-to-point communications.
3. Kernel process scheduling, together with the jitter introduced by the operating system.
4. Cross-application contention, which constitutes a random variability coming from simultaneous access to shared components in the computing platform.

While issues 3 and 4 cannot be addressed by the end-users of a platform, issues 1 and 2 can be better handled by tuning large-scale applications in such way that they make a more efficient use of resources. As an example, parallel file systems represent a well-known bottleneck and a source of high variability [32]. While the performance of computation phases of HPC applications is usually stable and only suffer from a small jitter due to the operating system, the time taken by a process to write some data can vary by several orders of magnitude from one process to another and from one iteration to another. In [17], variability is expressed in terms of *interferences*, with the distinction between *internal interferences* caused by access contention between the processes of an application (issue 2), and *external interferences* due to sharing the access to the file system with other applications, possibly running on different clusters (issue 4). As a consequence, adaptive I/O algorithms have been proposed [17] to allow a more efficient and less variable access to the file system.

## 2.2 Approaches to I/O management in HPC simulations

Two main approaches are typically used for performing I/O in large-scale HPC simulations:

**The *file-per-process* approach** it consists in having each process write in a separate, relatively small file. Whereas this avoids synchronization between processes, parallel file systems are not well suited for this type of load when scaling to hundreds of thousands of files: special optimizations are then necessary [7]. File systems using a single metadata server, such as Lustre, suffer from a bottleneck: simultaneous creations of so many files are serialized, which leads to immense I/O variability. Moreover, reading such a huge number of files for post-processing and visualization becomes intractable.

**Using *collective I/O*** e.g. in MPI applications, all processes synchronize together to open a shared file, and each process writes particular regions of this file. This approach requires a tight coupling between MPI and the underlying file system [26]. Algorithms termed as “two-phase I/O” [3, 31] enable efficient

collective I/O implementations by aggregating requests and by adapting the write pattern to the file layout across multiple data servers [10]. When using a scientific data format such as pHDF5 [9] or pNetCDF [15], collective I/O avoids metadata redundancy as opposed to the file-per-process approach. However, it imposes additional process synchronization, leading to potential loss of efficiency in I/O operations. Moreover, pHDF5 and pNetCDF currently do not allow data compression.

It is usually possible to switch between these approaches when a scientific format is used on top of MPI; going from HDF5 to pHDF5 is a matter of adding a couple of lines of code, or simply changing the content of an XML file when using an adaptable layers such as ADIOS [18]. But users still have to find the best specific approach for their workload and choose the optimal parameters to achieve high performance and low variability. Moreover, the aforementioned approaches create periodic peak loads in the file system and suffer from contention at several levels. This first happens at the level of each multicore SMP node, as concurrent I/O requires all cores to access remote resources (networks, I/O servers) at the same time. Optimizations in collective-I/O implementations are provided to avoid this first level of contention; e.g. in ROMIO [11], data aggregation is performed to gather outputs from multiple cores to a subset of cores that interact with the file system by performing larger, contiguous writes.

### 3 Impact of user-controllable parameters on I/O variability

The influence of all relevant parameters involved in the I/O stack is not possible to ascertain in a single study: varying the network, the number of servers or clients, the file system together with its configuration, leads to too many degrees of freedom. Moreover, users can usually control only a restricted set of these parameters; we will therefore focus on such parameters only: the amount of data written at each I/O phase, the I/O approach used (collective write to a single shared file or individual files-per-process approach), the file striping policy inside the file system and the number of writing processes. The goal of this section is to emphasize their impact on I/O performance variability. We first present the methodology and metrics we use to compare experiments, then we experimentally characterize I/O variability on the Grid'5000 French testbed and on Kraken using the IOR benchmark. Since I/O performance variability is a well known problem on Lustre (in particular due to its single metadata server [23]), we have used PVFS on Grid'5000 to show that I/O variability also appears on other systems and at smaller scale.

#### 3.1 Methodology

The main problem when studying variability instead of performance involves choosing the right metric and the proper way to quickly get valuable insight into the I/O systems behavior in a large number of experiments. Choosing the write variance as a statistical metric for the variability is arguably not appropriate for three reasons: first, it highly depends on the underlying configuration; second, the values are hard to interpret and compare; finally, providing decent confident intervals for such a statistic would require thousands of runs for each tested

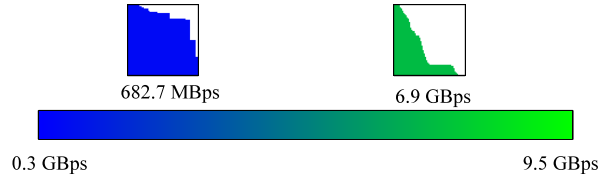


Figure 1: Visual representation of runs variability. The left run has a smaller throughput but less overall idle time than the right one.

configuration. Thus, we need new methods to investigate I/O variability at large scale.

### Visual representation of runs

We propose the following methodology to study the variability of I/O systems. We define a *run* as a single I/O phase of an application under a particular configuration. An *experiment* is a set of runs with the same configuration. For a particular run, we measure the time taken by each process to write its data. This time is reported on a small graph (that we call a *tile*), in which processes are sorted vertically by their write time, this time being represented by an horizontal line. We call the *variability* pattern the shape within a particular tile. As will be shown, for a given configuration we find recurrent patterns. An example of such a representation is shown in Figure 1

The throughput of a run, defined as the total amount of data divided by the time of the slowest process, is shown using a color scale for drawing these graphs, and also printed under each tile<sup>1</sup>. Compared to a computation of variance only, this method presents a more complete picture, enabling us to concurrently visualize different types of jitter including variability from a process to another, from a run to another and between experiments. The ratio between the white part and the colored part of a tile gives an overview of the time wasted by processes waiting for the slowest to complete the write. The shape within tiles provides some clues about the nature of the variability. The range of colors within a set of tiles provides a visual indication of the performance variability of a configuration for a given experiment, and also lets us compare experiments when the same color scale is used. In the following discussion, only a subset of representative runs are presented for each experiment.

### Platforms and tools considered

The first experiments have been conducted on two clusters in the Rennes site of the French Grid'5000. The *parapluie* cluster, featuring 40 nodes (2 AMD 1.7 GHz CPUs, 12 cores/CPU, 48 GB RAM), is used for running the IOR benchmark on a total of 576 cores (24 nodes). The latest version OrangeFS 2.8.3 of PVFS is deployed on 16 nodes of the *parapide* cluster (2 Intel 2.93 GHz CPUs, 4 cores/CPU, 24 GB RAM, 434 GB local disk), each node used both as I/O server and metadata server. In each cluster all nodes are communicate

<sup>1</sup>Note: this representation is color-based, black and white printing may lead to a loss of information.

through a 20G InfiniBand 4x QDR link connected to a common Voltaire switch. We use MPICH [22] with ROMIO [31] compiled against the PVFS library.

The last experiments aim to show the evolution of the jitter at large process counts. They are conducted on Kraken, the Cray XT5 machine running at NICS. Kraken has 9408 compute nodes running Linux Cray. Each node features 12 cores (two 2.6 GHz six-core AMD Opteron processors) and 16 GB of main memory, nodes are connected to each other through a Cray SeaStar2+ router. The Lustre file system available on Kraken is deployed on 336 OST (Object Storage Targets) running on 48 OSS (Object Storage Servers), and one MDS (MetaData Server). Kraken is shared by all users of the platforms as we are running experiments.

IOR [29] is a benchmark used to evaluate I/O optimizations in high performance I/O libraries and file systems. It provides backends for MPI-IO, POSIX and HDF5, and lets the user configure the amount of data written per client, the transfer size, the pattern within files, together with other parameters. We modified IOR in such a way that each process outputs the time spent writing; thus we have a process-level trace instead of simple average of aggregate throughput and standard deviation.

## 3.2 Results with the IOR benchmark

### 3.2.1 Impact of the data size

We first study the impact of the size of the data written by the processes on overall performance and performance variability. We perform a set of 5 experiments on Grid'5000 with a data size ranging from 4 MB to 64 MB per process using 576 writers. These data are written in a single shared file using collective I/O. Figure 2 represents 5 of these runs for each data size. We notice that for small sizes (4 MB) the behavior is that of a waiting queue, with a nearly uniform probability for a process to write in a given amount of time. However, the throughput is highly variable from one run to another, ranging from 4 to 8 GB/s. An interesting pattern is found when each process outputs 16 MB: the idle time decreases, the overall throughput is higher (around 7 GB/s) and more regular from a run to another. Moreover the variability pattern shows a more fair sharing of resources. Performance then decreases with 32 and 64 MB per process. When writing 64 MB per process, the throughput is almost divided by 7 compared to writing 16 MB.

IOR allows the user to configure the transfer size, i.e., the number of atomic write calls in a single phase. When this transfer size is set to 64 KB (which is equal to the MTU of our network and to the stripe size in the file system), we observed that writing 4, 8 or 16 MB per process leads to the same variability pattern and throughput (around 4 GB/s). Increasing the total data size to 32 or 64 MB shows however a decrease of performance and an increase of variability from one run to another. With 64 MB, a staircase pattern starts to appear, which might be due to the amount of data reaching cache capacity limits in the file system.

### 3.2.2 Impact of the striping policy

Like Lustre, PVFS allows the user to tune the file's striping policy across the I/O servers using `lfs setstripe` commands. The default distribution in PVFS

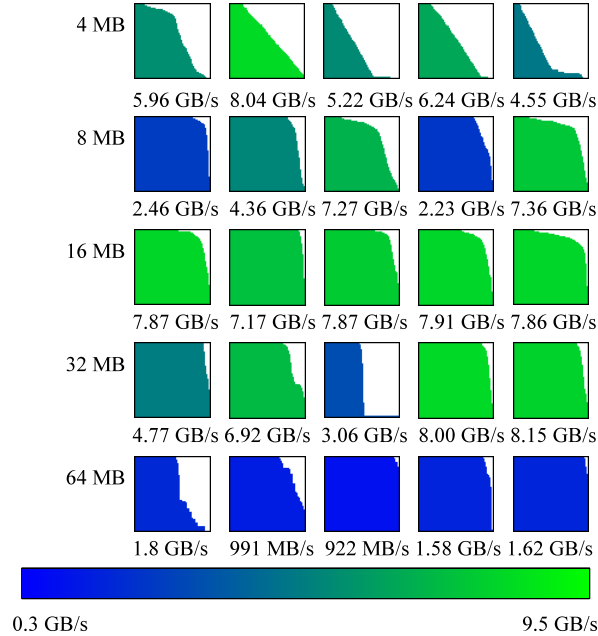


Figure 2: Impact of the data size on I/O performance variability. Experiments are conducted on 576 cores of Grid'5000, using Collective I/O.

consists in striping data in a round robin manner across the different I/O servers. The stripe size is one of the parameters that can be tuned. The PVFS developers informed us that given the optimizations included in PVFS and in MPICH, changing the default stripe size (64 KB) should not induce any change in aggregate throughput. We ran IOR three times, making each of the 576 processes writes 32 MB in a single write within a shared file using collective I/O. We observed that a similar maximum throughput of around 8 GB/s is achieved for both 64 KB, 4 MB and 32 MB for stripe size. However, the throughput was lower on average with a stripe size of 4 MB and 32 MB, and the variability patterns, which can be seen in Figure 3, shows more time wasted by processes waiting.

While the stripe size had little influence on collective I/O performance, this was not the case for the file-per-process approach. When writing more than 16 MB per process, the lack of synchronization between processes causes a huge access contention and some of the write operations to time out. The generally-accepted approach of increasing the stripe size or reducing the stripes count for the many-file approach is thus justified, but as we will see, this approach still suffers from high variability.

### 3.2.3 Impact of the I/O approach

We then compared the collective I/O approach with the file-per-process approach. 576 processes are again writing concurrently in PVFS. With a file size of 4 MB, 8 MB and 16 MB and a stripe size of 64 KB in PVFS, the file-per-process approach shows a much lower aggregate throughput than collective I/O.

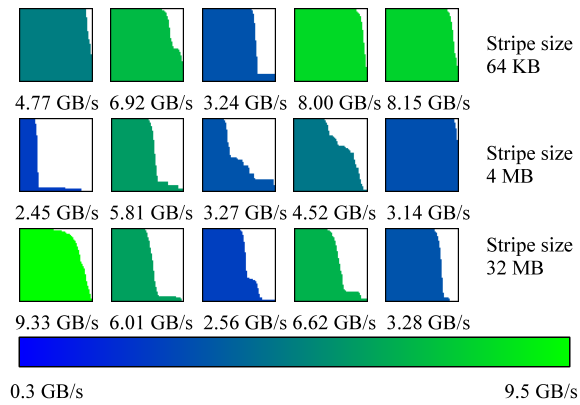


Figure 3: Throughput variability for three experiments with different stripe size. Experiments conducted on 576 cores of Grid’5000, using Collective I/O.

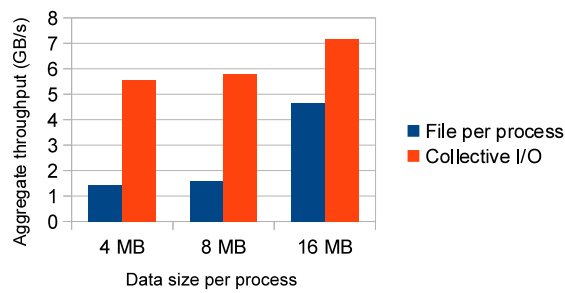


Figure 4: Aggregate throughput comparison between file-per-process and Collective I/O on 576 processes of Grid’5000.

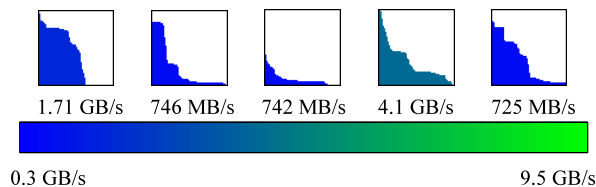


Figure 5: Throughput variability with the file-per-process approach on 576 cores of Grid’5000, each core writing 32 MB.

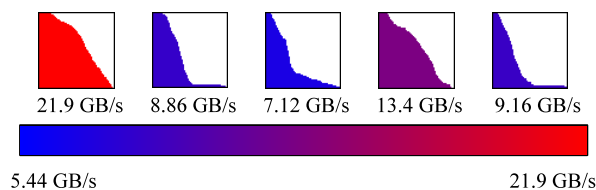


Figure 6: Throughput variability with the file-per-process approach on Grid’5000, writing 32 MB per process with a 32 MB stripe size.

Results (in terms of average aggregate throughput) are shown on Figure 4. Figure 5 presents the recurrent variability patterns of runs for 8 MB per process (runs writing 4 MB and 16 MB have similar patterns). We observe a staircase progression, which suggests communications timing out and processes retrying in series of bursts.

As previously mentioned, writing more data per client required to change the stripe size in PVFS. Yet the throughput obtained when writing 32 MB per process using a 32 MB stripe size and the file-per-process approach on Grid’5000 is found to be the highest over all experiments we did, with a peak throughput of more than 21 GB/sec. Thus Figure 6 shows the results with a different color scale in order not to be confused with other experiments. We notice that despite the high average and peak throughput, there are periods where a large number of processes are idle, and there is a high variance from a run to another.

### 3.2.4 Impact of the number of writers

On Kraken, we reproduced the experiment presented in Figure 5 of [23]: each process outputs 128 MB using a 32 MB transfer size and a file-per-process approach. Results are presented in Figure 7: the aggregate throughput increases with the number of processes, just as expected in [23]. An intuitive idea could have been that with larger process counts comes a higher variability from a run to another. This is however not the case. The explanation behind this is the law of large numbers: as the number of processes increases and since each process has to issue several transfers, we are more likely to observe a gaussian repartition of the time spent by a process writing its data. Indeed, with 8400 writers the variability patterns start to look like the cumulative distribution function of a gaussian. This effect has already been exemplified in [32].

In conclusion, the experiments presented above show that the performance of I/O phases is extremely variable, and that this variability, together with

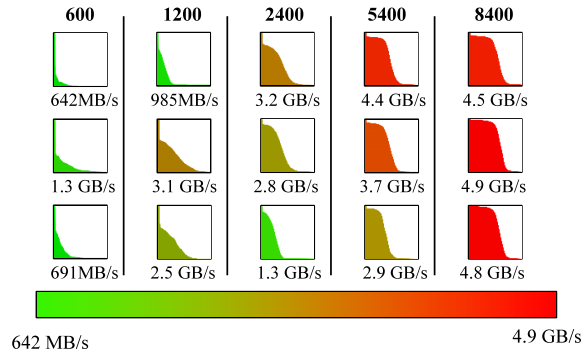


Figure 7: Throughput variability with the file-per-process approach on Kraken, with different number of writers. Each process writes 128 MB with a 32 MB transfer size.

the average and peak throughput achieved, highly depends on the data size, on the I/O approach and the number of writers. Experiments conducted on Grid’5000 with clusters fully reserved to our own experiments (including network interconnects) prove that a high variability can already be experienced with few concurrent accesses to the parallel file system. Whereas tuning the stripe size in PVFS has almost no influence when using collective I/O, it greatly helps to achieve a very high throughput with the file-per-process approach. We have successfully demonstrated that the chosen parameters and the I/O approach have a big impact on the performance and variability of I/O phases. The last experiment on Grid’5000 revealed that while a high throughput of at least 21 GB/s can be expected of our file system, poor configuration choices can lead to barely a few hundreds of MB/s.

The experiments conducted on Kraken with respect to scalability showed that in the presence of many independent writes, the variability pattern is becoming more predictable as we increase the number of writers. Still, we notice that a small set of processes manage to quickly terminate their I/O phase, which motivates further research of new I/O approaches that do not exhibit this variability.

## 4 The dedicated core approach

One conclusion of our study of I/O performance variability is that, in order to sustain a high throughput and a lower variability, it is preferable to write large data chunks while avoiding as much as possible access contentions at the level of the network interface and of the file system, for example by reducing the number of writers. Collective I/O reduces I/O variability, but introduces additional synchronizations that can limit the global I/O throughput at large process counts. On the other hand, the file-per-process approach can lead to a better throughput, but causes a high variability and complex output analysis. As the first level of contention occurs when several cores in a single SMP node try to access the same network interface, we propose to gather the I/O operations into one single core that will perform writes of larger data in each SMP node.



Moreover, this core will be dedicated to I/O (i.e. will not perform computation) in order to overlap writes with computation and avoid contention for accesses to the file system. We call this approach Damaris: this section introduces its design, implementation and API.

## 4.1 Principle

Damaris consists of a set of servers, each of which runs on a dedicated core in every SMP node used by the simulation. This server keeps data in a shared memory segment and performs post-processing, filtering and indexing in response to user-defined events sent either by the simulation or by external tools.

The buffering system running on these dedicated cores includes metadata information about incoming variables. In other words, clients do not write raw data but enriched datasets in a way similar to scientific data formats such as HDF5 or NetCDF. Thus dedicated cores have a full knowledge about incoming datasets and can perform “smart actions” on these data, such as writing important datasets in priority, performing compression, statistical studies, indexing, or any user-provided transformation.

## 4.2 Architecture

The architecture of Damaris is shown in Figure 8, representing a multicore SMP node in which one core is dedicated to Damaris, the other cores (only three represented here) being used by the simulation. As the number of cores per node increases, dedicating one core has a diminishing impact. Thus, our approach primarily targets SMP nodes featuring more than 8 cores.

As indicated above, the I/O pattern of large-scale applications is predictable, thus additional knowledge regarding the data that is expected to be written can be provided by the user. An important component in the design of Damaris is the external configuration file. This file (written by the user) provides a model of the data generated by the application. It follows an idea introduced by the ADIOS data format in [18], which allows the configuration of the I/O layers to be written into an external file, avoiding code recompilation when the user wants a different output. The two main goals of this external configuration are 1) to keep a high-level description of the datasets within the server, allowing higher-level data manipulations and 2) to avoid this description to be sent by clients through the shared memory (only data is sent together with the minimal descriptor that lets the server retrieve the full description in its metadata tree). The configuration file establishes a contract between clients and servers, defining the interactions between the two entities.

Communication between the computation cores and the dedicated cores is done through shared memory. A large memory buffer is created by the dedicated core at start time, with a size chosen by the user (thanks to the configuration file). Thus the user has a full control over the resources allocated to Damaris. When a client submits new data, it requests the allocation of a segment from this buffer, then copies its data using the returned pointer. By default, a mutex-based allocation algorithm is used to allow concurrent atomic allocations of segments by multiple clients. Even though this current allocation solution is already efficient, we are investigating possible algorithms leveraging *a priori* knowledge provided by the user to preallocate segments and avoid locking. As

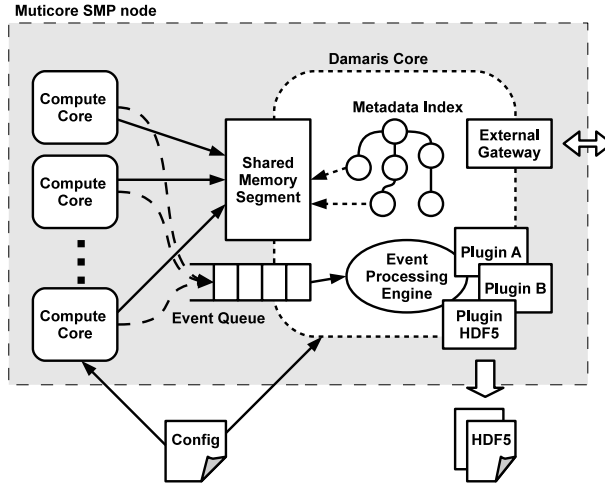


Figure 8: Design of the Damaris approach.

an example, if all clients are expected to write the same amount of data, the buffer can be split in as many parts as clients, allowing lock-free concurrent writes.

The event-queue is another shared component of the Damaris architecture. It is used by the clients either to inform the server that a write completed (*write-notification*), or to send *user-defined events*. The messages are pulled by an event processing engine (EPE). The configuration file provided by the user contains actions to be performed upon reception of write-notifications and user-defined events. If no action is defined for a given received event, the EPE just ignores it. Actions can be provided by the user through a plugin system. Such actions can prepare data for future analysis, for instance, or simply flush all in-memory data using a scientific I/O library.

All datasets written by the clients are uniquely characterized by a tuple  $\langle name, iteration, source, layout \rangle$ . *Iteration* gives the current step of the simulation, while *source* uniquely characterizes the client that has written the variable. The *layout* corresponds to a description of the structure of the data, for instance “3D array of 32 bits real values with extents  $n_x, n_y, n_z$ ”. For most simulations, this layout does not vary during runtime and can be provided also by the configuration file. Upon reception of a write-notification, the EPE will add an entry in a metadata structure associating the tuple with the received data (that stay in shared memory until actions are performed on them).

### 4.3 Key design choices

#### Behavior management and user-defined actions

The EPE can be enriched by plugins provided by the user. A plugin is a function (or a set of functions) embedded in a dynamic library that the actions manager will load and call in response to events sent by the application. The matching between events and expected reactions is provided by the external configuration file. Thus, it is easy for the user to define a precise behavior for Damaris by

simply changing the configuration file. Our intents when designing Damaris was to provide a very simple way for users to extend it and adapt it to the particular needs of their simulations.

### Event scope

It can be desirable to force a reaction only when a given subset of computation cores have sent an event. This can be done using the notion of event scope, which defines the subset of processes that has to perform a synchronized action, such as statistical data sharing for establishing a global diagnosis of the simulation. This synchronization involves distributed algorithms such as total order broadcast, and efficient scheduling schemes, that are beyond the scope of this paper.

### Interaction semantics

The use of a single message queue shared by all the cores and used both for user-defined events and write notifications ensures that a sequence of events or data sent by a compute core will be treated in the same order in the dedicated core. Yet, sequences of events coming from different compute cores may interleave. This semantics is important because the purpose of user-defined events is to force a behavior at a given execution point of the client, knowing what has already been written.

### Minimum-copy overhead

Another particularity of Damaris is the efficiency of interactions between clients and dedicated cores. At most a single copy from a local variable to the shared memory buffer is required to send data to the dedicated core. Damaris also includes the possibility to “write” data without actually making a copy: the simulation directly allocates its variables in the shared memory buffer. When the simulation finishes working on an array, it simply informs the dedicated core that the data can be considered as “written”. This is a strong point of Damaris that distinguishes it from other dedicated-process-based approaches [19, 16], in which multiple raw copies of the data have to be moved using either the network interface or kernel functionalities before actually reaching the I/O service.

### Persistency layer

Damaris interfaces with an I/O library such as HDF5 by using a custom persistency layer with HPF5 callback routines embedded in a plugin, as shown as an example in Figure 8.

## 4.4 Client-side API

Damaris is intended to be a generic, platform-independent, application-independent, easy-to-use tool. The current implementation is developed in C++ and uses the Boost library for interprocess communications and options parsing, and Xerces-C for XML configuration. It provides client-side interfaces for C, C++ and Fortran applications and requires only few minor changes in the application’s I/O routines, together with an external configuration file that describes the

data. The client-side interface is extremely simple and consists in four main functions (here in C):

- `dc_initialize("config.xml",core_id)` initializes the client by providing a configuration file and a client ID (usually the MPI rank). The configuration file will be used to retrieve knowledge about the expected data layouts.
- `dc_write("varname",step,data)` pushes some data into the Damaris' shared buffer and sends a message to Damaris notifying the incoming data. The dataset being sent is characterized by its name and the corresponding iteration of the application. All additional information such as the size of the data its layout (including its datatype) and additional descriptions are provided by the configuration file.
- `dc_signal("eventname",step)` sends a custom event to the server in order to force a behavior predefined in the configuration file.
- `dc_finalize()` frees all resources associated with the client. The server is also notified of client's disconnection and will not accept any other incoming data or event from this client.

Additional functions are available to allow direct access to an allocated portion of the shared buffer (`dc_alloc` and `dc_commit`), avoiding an extra copy from local memory to shared memory. Other functions allow the user to dynamically modify the internal configuration initially provided, which can be useful when writing variable-length arrays (which is the case in particle-based simulations, for example) or user-defined datatypes.

Below is an example of a Fortran program that makes use of Damaris to write a 3D array then send an event to the I/O core. The associated configuration file, which follows, describes the data that is expected to be received by the I/O core, and the action to perform upon reception of the event.

```

program example
  integer :: ierr, rank, step
  real, dimension(64,16,2) :: my_data
  call df_initialize("my_config.xml", rank, ierr)
  ...
  call df_write("my_variable", step, my_data, ierr)
  call df_signal("my_event", step, ierr)
  ...
  call df_finalize(ierr)
end program example

```

Associated XML configuration file:

```

<layout name="my_layout" type="real"
  dimensions="64,16,2" language="fortran" />
<variable name="my_variable" layout="my_layout" />
<event name="my_event" action="do_something"
  using="my_plugin.so" scope="local" />

```

## 5 Experimental results with CM1

We have evaluated our approach based on dedicated I/O cores against standard approaches with the CM1 atmospheric simulation, using our two platforms (Kraken and Grid'5000). Results on a third platform (BluePrint) were also presented in the first version of this research report. These results are in Appendix ??.

### 5.1 The CM1 application

CM1 [6] is used for atmospheric research and is suitable for modeling small-scale atmosphere phenomena such as thunderstorms and tornadoes. It follows a typical behavior of scientific simulations which alternate computation phases and I/O phases. The simulated domain is a fixed 3D array representing part of the atmosphere. The number of points along the  $x$ ,  $y$  and  $z$  axes is given by the parameters  $n_x$ ,  $n_y$  and  $n_z$ . Each point in this domain is characterized by a set of variables such as local *temperature* or *wind speed*.

CM1 is written in Fortran 95. Parallelization is done using MPI, by splitting the 3D array along a 2D grid. Each MPI rank is mapped into a pair  $(i, j)$  and simulates a  $n_{sx} * n_{sy} * n_z$  points subdomain. In the current release (r15), the I/O phase uses HDF5 to write one file per process. Alternatively, the latest development version allows the use of pHDF5. One of the advantages of using a file-per-process approach is that compression can be enabled, which is not the case with pHDF5. Using lossless gzip compression on the 3D arrays, we observed a compression ratio of 187%. When writing data for offline visualization, the floating point precision can also be reduced to 16 bits, leading to nearly 600% compression ratio when coupling with gzip. Therefore, compression can significantly reduce I/O time and storage space. However, enabling compression leads to an additional overhead, together with more variability both from a process to another, and from a write phase to another.

### 5.2 Platforms and configuration

- **On Kraken** we study the scalability of different approaches, including Damaris. Thus, the problem size vary from an experiment to another. When all cores are used by the simulation, each process handles a 44x44x200 points subdomain. Using Damaris, each non-dedicated core (11 out of 12) handles a 48x44x200 points subdomain, thus making the total problem size equivalent.
- **On Grid'5000** we run the development version of CM1 on 28 nodes, i.e. 672 cores. The total domain size is 1104x1120x200 points, so each core handles a 46x40x200 points subdomain with a standard approach, and a 48x40x200 points subdomain when one core out of 24 is used by Damaris.

CM1 takes an input file indicating whether each variable should be written or not, together with the frequency of outputs. We ran it with an output configuration and frequency corresponding to what can be expected by an atmospheric scientist from such a simulation.

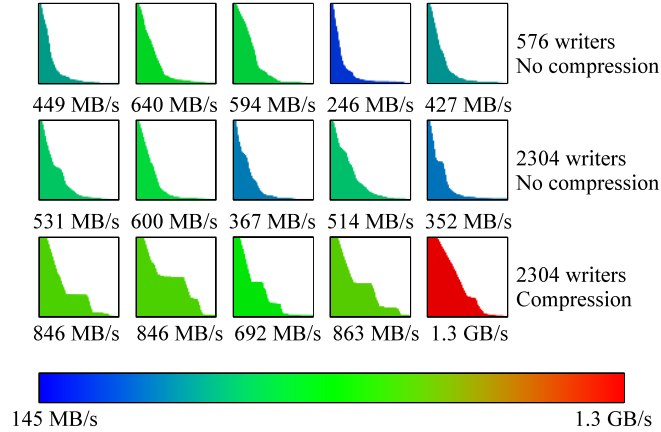


Figure 9: Variability patterns for three I/O phases of CM1 on Kraken, using different numbers of cores and enabling/disabling compression.

### 5.3 Non-overlapping I/O approaches

Our first investigations on optimizing I/O for CM1 started with standard approaches that do not overlap I/O with computation.

#### 5.3.1 File-per-process method

On Kraken, we run CM1 on 576 then 2304 cores. Each core writes 28.4 MB of data in an individual HDF5 file every 50 iterations of the model. Considering the I/O tips provided on the Kraken website [23], we set the stripe settings in Lustre so each file is handled by a single OST.

Figure 9 represents the variability patterns of five representative write phases using three different configurations: 576 or 2304 core, compression enabled<sup>2</sup> or not.

As we can see, a lot of time is wasted waiting for slow processes, both when using 576 and 2304 cores. The variability from a write phase to another is also clear. When enabling compression, we observe a recurrent pattern with different groups of processes. This pattern remains to be explained by investigating how HDF5 filters work and how the file size is correlated to the write time. We abandon the file-per-process approach at larger scale, due to a huge waste of performance and above all an extreme difficulty of post-processing the generated data (simply deleting all the files produced by CM1 required hours and heavy load on Lustre’s metadata server).

Note that, since CM1 uses HDF5 or pHDF5, we cannot have complete control over the data layout and on the number of atomic write calls to MPI subroutines. We thus observe that the achieved throughput is much lower than the performance measured in [23] and with the IOR benchmark in Section 3.

<sup>2</sup>The represented throughput corresponds to the apparent aggregate throughput, i.e. the total amount of raw uncompressed data divided by the time for the slowest process to complete the write phase.

On Grid5000 we run CM1 using 672 cores, writing a total of 15.8 GB uncompressed data (about 24 MB per process) at every 20 iterations. CM1 reported spending 4.22% of its time in I/O phases. Yet the fastest processes usually terminate their I/O in less than 1 sec, while the slowest take more than 25 sec, leading to an aggregate throughput of 695 MB/s on average and a standard deviation of 74.5 MB/s. Following the insights of Section 3, we increased the stripe size to 32 MB in PVFS. At such a “small” scale and using PVFS, the file-per-process approach remains a good strategy on Grid’5000 (the time spent in I/O phase is lower than the 5% usually accepted). Yet the I/O jitter is present and proves that a better solution can still be found.

A potential improvement in the file-per-process approach can be found by buffering several outputs using the `H5P_set_faple_core` feature of HDF5, which allows the datasets to be written in local memory and flushed only when closing the file. The study of this method is left as future work.

### 5.3.2 Using collective I/O and grouped I/O

The current development version of CM1 provides a pHDF5 backend that can gather MPI processes in groups. Each group creates and writes in a distinct file using collective-I/O. This way, the overall number of files can be reduced without the potential burden of a heavy synchronization of all processes at large process count. We thus tested CM1 writing either one shared file, or one file per node instead of one per core. Yet the use of pHDF5 functions forces a high level of synchronization between processes (several `MPI_Barriers` in write functions), which does not allow us to have a process-level view of the time spent in I/O and to use our jitter visualization method.

On Kraken, we experimented with 9216 cores writing a single shared file. Three I/O phases reported to last respectively 457 sec, 674 sec and 816 sec, i.e. showing a throughput of 569 MB/s, 385 MB/s and 318 MB/s. Considering that each computation phase lasts 220 sec, it is simply unacceptable to waste such an immense and unpredictable time performing I/O. By setting the stripe size to 32 MB instead of 1 MB in Lustre, the aggregate throughput went down to 90 MB/s, exemplifying the fact that bad choices of file system’s configuration can lead to extremely poor I/O performance. Similar results of both high variability and low throughput have been reported when using grouped-I/O to write one file per node. As these results don’t bring anything new, they won’t be commented in this paper.

On Grid’5000, when all processes write to a single shared file on PVFS, CM1 shows poor I/O performance (636 MB/s of aggregate throughput on average) and a relatively high variability, with a 101.8 MB standard deviation between distinct write phases. On this platform, the global shared file approach thus shows worse performance than the file-per-process approach, both in terms of average performance and performance variability. Figure 10 shows CPU usage on one of the nodes used by CM1 during the execution. Non-optimal use of computing resources can be noticed during I/O.

Writing one file per node using grouped I/O leads to a small improvement on Grid’5000: 750 MB/s on average. An expected result is the decrease of the associated standard deviation, from 101.8 MB to 36.8 MB. The improvement in standard deviation may not be significant, as the small number of runs does not lead to small enough confidence intervals. The output when using the grouped

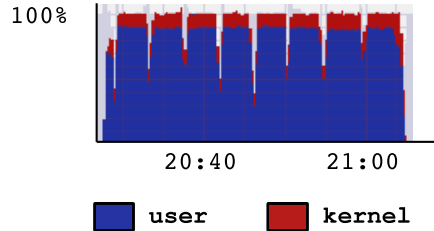


Figure 10: CPU usage measured on a node running CM1 with a collective I/O mechanism.

I/O produces 28 files instead of 672. But the compression capability is still lost due to the use of pHDF5. We initially thought such a grouped I/O approach might be advantageous only at extremely large scales, when both the file-per-process approach and the collective I/O approach show their limits, yet it didn't bring any improvement on Kraken compared to our previously presented results.

To conclude, collective-I/O achieves better performance than the file-per-process I/O on Grid'5000, but leads to unacceptably low performance on Kraken.

#### 5.4 Our proposal: using dedicated I/O cores

In this section, we evaluate the I/O performance and variability when using dedicated cores.

On Kraken, we ran CM1 with a Damaris I/O backend and experimented with 576, 2304 and 9216 cores. The configuration of CM1 and the expected output is similar to that of the two previous sections. One Damaris server is deployed on one of the 12 cores available on each SMP node, thus potentially reducing the computation power for the benefit of I/O efficiency. And the bet is widely worth it: the time spent in I/O phases from the point of view of the simulation is reduced to less than 1 sec, independently of the I/O phase or the scale. The jitter is now expressed in terms of tenths of seconds and cannot be distinguished with other communication- or memory-related jitters. In terms of global performance, the Figure 11 presents the total time spent by simulation performing 50 iterations and one write phase, with the different approaches considered. In every cases, the use of Damaris increases global performance, with a global improvement of nearly 70% on 9216 cores compared to the Collective-I/O approach.

From the point of view of the dedicated cores, each of them has to aggregate and write 322 MB of data, dividing the number of generated files by 12 and thus reducing the pressure on the metadata server. The average aggregate throughput is reported to be 1.8 GB/s for 576 cores (48 writers), and 6.1 GB/s for 2304 cores (192 writers), that is to say almost 15 times higher than the average throughput obtained with collective I/O on the same number of cores.

We enabled compression on the last 9216-cores experiment (768 writers), to achieve an apparent throughput of 4.1 GB/s and a compression ratio of 600%. The overhead and jitter induced by this compression is now completely hidden within the dedicated cores, and do not impact the running simulation.



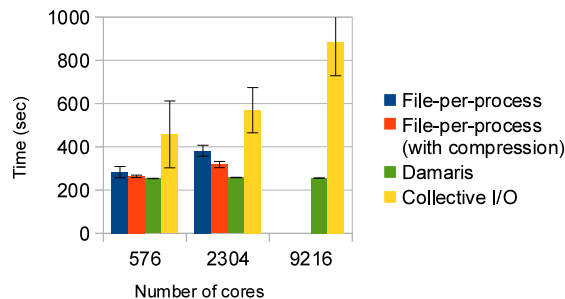


Figure 11: Comparison of performance (time to perform 50 iterations of computation and 1 write phase on Kraken) using the file-per-process, collective I/O and Damaris approaches. Standard deviations for 10 runs are provided. The file-per-process approach is abandoned when using 9216 cores because of the heavy pressure it puts on the file system and the difficulty to analyze the output.

Additionally, we implemented in Damaris the capability to schedule data movements. Instead of writing a file as soon as all the required variables have been written by the clients, the I/O cores compute an estimation of the computation time of an iteration (230 sec on Kraken). This time is then divided into as many slots as I/O cores, and each I/O core waits for its slot before actually writing. Evaluating this strategy on 2304 cores on Kraken, the aggregate throughput achieves 13.1 GB/s on average, that is to say 10 times higher than the best I/O throughput we achieved with the file-per-process approach at the same scale.

In all these experiments, the dedicated cores are only spending a small fraction of their time actually writing data (about 2%). Thus, the time spared can be used to perform data processing prior to storage. Investigating the potential use of this time is left as future work.

We did a similar experiment on Grid'5000, with one core out of 24 dedicated to I/O. CM1 writes every 20 iterations, that is to say about every 3 minutes of wallclock time. Here again Damaris aggregates data in one file per dedicated core, thus reducing the number of files from 672 to 28. Figure 12 shows the CPU usage in a node running CM1. The idle phases corresponding to I/O phases have disappeared and the CPU is constantly fully utilized. The I/O cores report to achieve a throughput of 4.32 GB/s, that is to say about 6 times higher than the throughput achieved with the collective. This throughput presents a standard deviation of 319 MB/s. But the most important thing to remember is that this variability is now completely hidden within a dedicated resources, and that the end-user doesn't have to care about it anymore. The CM1 simulation has thus recovered its original stability and predictability and dedicated cores are saving 75% of time that can be used for data processing.

We thus notice an improvement on both platforms together with a time savings that could be dedicated to data post-processing. CM1 also integrates a communication-intensive phase that periodically computes statistics in order to

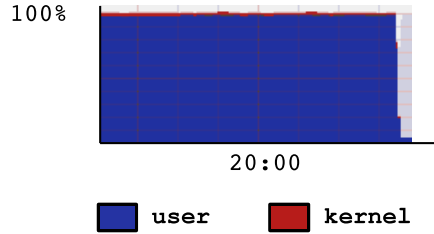


Figure 12: CPU usage measured on a node running CM1 with Damaris.

be able to stop the simulation if some values become physically irrelevant. This phase makes use of several calls to `MPI_Gather` and `MPI_Reduce` that cause a high overhead, reported to be more than 7% of walltime on Grid'5000. By moving this diagnosis into the dedicated cores, we could achieve a 7% additional speedup, to be added to the speedup already achieved by removing the I/O part.

## 6 Discussion and related work

In this section, we evaluate the benefits of our approach by computing mathematical bounds of effectiveness, and by comparing to experimental ascertainments. We then position our approach with respect to other related research.

### 6.1 Do all my cores are really needed?

Let us call  $W_{std}$  the time spent writing and  $C_{std}$  the computation time between two I/O phases with a standard approach,  $C_{ded}$  the computation time when the same workload is divided across one less core in each node. We here assume that the I/O time is completely hidden or at least negligible when using the dedicated core (which is experimentally verified) from the point of view of the simulation, and we call  $W_{ded}$  the time that the dedicated core spends writing. A theoretical performance benefit of our approach then occurs when

$$W_{std} + C_{std} > \max(C_{ded}, W_{ded})$$

Assuming an optimal parallelization of the program across  $N$  cores per node, and the worst case for Damaris where  $W_{ded} = N * W_{std}$ , we show that this inequality is true when the program spends at least  $p\%$  in I/O phase, with  $p = \frac{100}{N-1}$ . As an example, with 24 cores  $p = 4.35\%$ , which is already under the 5% usually admitted for the I/O phase of such applications. Thus assuming that the application effectively spends 5% of the time writing data, on a machine featuring more than 24 cores per node, it is more efficient to dedicate one core per node to hide the I/O phase.

However this is a simplified theoretical estimation. In practice, most HPC simulations doesn't exhibit a linear scalability. We have run CM1 with the exact same configuration and network topology on Grid'5000, dividing the workload across 24 cores per node first, then 23 cores per node. 200 iterations with 24 cores per node were done in 44'17" while only 41'46" were necessary to run 200

iterations with 23 cores per node. Similar ascertainment was made on Kraken, especially as we increase the total number of cores, using 11 cores per nodes instead of 12 already produces an improvement. The final statistics output by CM1 showed that all the phases that use all-to-all communications benefit from leaving one core out, thus actually reducing the computation time as memory contention is reduced. Such a behavior of multicore architectures is explained in [21] and motivates future studies regarding the appropriate number of cores that can be dedicated to services as Damaris does. Moreover, writing more data from a reduced number of cores has shown to be much more efficient than the worst case used in the theoretical bound.

## 6.2 Positioning Damaris in the “I/O landscape”

Through its capability of gathering data into larger buffers and files, Damaris can be compared to the ROMIO data aggregation feature [31]. Yet, data aggregation is performed synchronously in ROMIO, so all cores that do not perform actual writes in the file system must wait for the aggregator processes to complete their operations. Through space-partitioning, Damaris can perform data aggregation and potential transformations in an asynchronous manner and still use the idle time remaining in the dedicated cores.

Other efforts are focused on overlapping computation with I/O in order to reduce the impact of I/O latency on overall performance. Overlap techniques can be implemented directly within simulations [25], using asynchronous communications. Yet non-blocking file operations primitives are not part of the current MPI-2 standard. A quantification of the potential benefit of overlapping communication and computation is provided in [27], which demonstrates methodologies that can be extended to communications used by I/O. Our Damaris approach exploits the potential benefit of this overlap.

Other approaches leverage data-staging and caching mechanisms [24, 14], along with forwarding approaches [4] to achieve better I/O performance. Forwarding routines usually run on top of dedicated resources in the platform, which are not configurable by the end-user. Moreover, these dedicated resources are shared by all users, which leads to multi-application access contention and thus to jitter. However, the trend towards I/O delegate systems underscores the need for new I/O approaches. Our approach follows this idea but relies on dedicated I/O cores at the application level rather than hardware I/O-dedicated or forwarding nodes, with the advantage of letting the user configure its dedicated resources to best fit its needs.

The use of local memory to alleviate the load on the file system is not new. The Scalable Checkpoint/Restart by Moody et al. [20] already makes use of node-local storages to avoid the heavy load caused by periodic global checkpoints. Their work yet does not use dedicated resources or threads to handle or process data, and the checkpoints are not asynchronous. But it exemplified the fact that node-local memory can be used to delay or avoid file system’s interactions. When additional SSDs, Flash memory or disks are attached to compute nodes, both SCR and Damaris can leverage it to avoid consuming too much of RAM.

Some research efforts have focused on reserving computational resources as a bridge between the simulation and the file system or other back-ends such as visualization engines. In such approaches, I/O at the simulation level is replaced

by asynchronous communications with a middleware running on a separate set of computation nodes, where data is stored in local memory and processed prior to effective storage. PreDatA [33] is such an example: it performs in-transit data manipulation on a subset of compute nodes prior to storage, allowing more efficient I/O in the simulation and more simple data analytics, at the price of reserving dedicated computational resources. The communication between simulation nodes and PreDatA nodes is done through the DART [12] RDMA-based transport method, hidden behind the ADIOS interface which allows the system to any simulation that has an ADIOS I/O backend. However, given the high ratio between the number of nodes used by the simulation and the number of nodes used for data processing, the PreDatA middleware is forced to perform streaming data processing, while our approach using dedicated cores in the simulation nodes allows us to keep the data longer in memory or any node-local storage devices, and to smartly schedule all data operations and movements. Clearly some simulations would benefit from one approach or the other, depending on their needs in terms of memory, I/O throughput and computation performance, but also the two approaches – dedicating cores or dedicating nodes – are complementary and we could imagine a framework that make use of the two ideas.

Finally service-dedicated cores are also used in [16] and [19]. Yet in the first one, communication from clients to dedicated cores is done through a FUSE interface, which implies several copy of data and a loss of coupling between cores. Damaris, by using shared memory and a configuration file, allows minimum copy and efficient interactions between cores. In the second one, active buffers are handled by dedicated processes that can run on any node and interact with cores running the simulation through network. The implementation is done behind the MPI level, contrary to Damaris, which is designed to improve communications thanks to shared memory and to keep a high-level view of datasets. Moreover, performance evaluation in these two related works are focusing on throughput, while we show the benefits of our method with respect to both jitter and I/O performance.

## 7 Conclusions

Managing I/O variability in an efficient way can have a substantial impact on the ability to sustain a high performance on Petascale and Post-Petascale infrastructures. The high impact of I/O jitter on the overall HPC application performance has already been observed at smaller scales; understanding its behavior and proposing an efficient mechanism to reduce its effects is critical for preparing the advent of Post-Petascale machines. The contributions of this paper can be summarized as follows: 1) We propose a method to visualize some I/O variability patterns and understand the impact of user-controllable parameters on this variability. Through experiments on the Grid'5000 testbed and on Kraken supercomputer, we show that choosing the appropriate I/O approach and tuning these parameters can substantially change the I/O throughput by an order of magnitude. 2) After an in-depth discussion of the limitations of standard I/O approaches to both I/O throughput and I/O performance variability, we propose a new approach (called Damaris) which leverages dedicated I/O cores in multicore SMP nodes. This solution provides the capability to

better schedule data movement, but also to process the data prior to storage. Our implementation tested with the CM1 atmospheric model up to 9216 cores on Kraken proved to be effective: it completely hides the I/O costs, achieves a throughput 15 times higher than standard approaches and allows overhead-free data compression with up to 600% compression ratio.

Our future work will focus on several directions. We plan to quantify the optimal ratio between I/O cores and computation cores within a node for several classes of HPC simulations. We will also investigate ways to leverage spare time in the I/O cores. A very promising direction is to attempt a tight coupling between running simulations and visualization engines, enabling direct access to data by visualization engines (through the I/O cores) while the simulation is running. This could enable efficient inline visualization without blocking the simulation, thereby reducing the pressure on the file systems. Finally, we plan to explore coordination strategies of I/O cores in order to implement distributed I/O scheduling.

## Acknowledgments

This work was done in the framework of a collaboration between the KerData INRIA - ENS Cachan/Brittany team (Rennes, France) and the NCSA (Urbana-Champaign, USA) within the Joint INRIA-UIUC Laboratory for Petascale Computing. The experiments presented in this paper were carried out using the Grid'5000/ALADDIN-G5K experimental testbed, an initiative from the French Ministry of Research through the ACI GRID incentive action, INRIA, CNRS and RENATER and other contributing partners (see <http://www.grid5000.fr/> for details) and on Kraken at NICS (see <http://www.nics.tennessee.edu/>). The authors also acknowledge the PVFS developers, the HDF5 group, Kraken administrators and the Grid'5000 users, who helped properly configuring the tools and platforms used for this work. We acknowledge Robert Wilhelmson (NCSA, UIUC) for his insights on the CM1 application and Dave Semeraro (NCSA, UIUC) for our fruitful discussions on visualization/simulation coupling.

## References

- [1] BlueWaters project, <http://www.ncsa.illinois.edu/BlueWaters/>, viewed on December 2010.
- [2] Kraken Cray XT5 system, <http://www.nics.tennessee.edu/computing-resources/kraken>, viewed on September 2011.
- [3] Evaluation of Collective I/O Implementations on Parallel Architectures. *Journal of Parallel and Distributed Computing*, 61(8):1052 – 1076, 2001.
- [4] N. Ali, P. Carns, K. Iskra, D. Kimpe, S. Lang, R. Latham, R. Ross, L. Ward, and P. Sadayappan. Scalable I/O forwarding framework for high-performance computing systems. In *Cluster Computing and Workshops, 2009. CLUSTER '09. IEEE International Conference on*, pages 1–10, September 2009.

- 
- [5] R. Bolze, F. Cappello, E. Caron, M. Daydé, F. Desprez, E. Jeannot, Y. Jégou, S. Lanteri, J. Leduc, N. Melab, et al. Grid'5000: a large scale and highly reconfigurable experimental grid testbed. *International Journal of High Performance Computing Applications*, 20(4):481, 2006.
- [6] G. H. Bryan and J. M. Fritsch. A benchmark simulation for moist nonhydrostatic numerical models. *Monthly Weather Review*, 130(12):2917–2928, 2002.
- [7] P. Carns, S. Lang, R. Ross, M. Vilayannur, J. Kunkel, and T. Ludwig. Small-file access in parallel file systems. *Parallel and Distributed Processing Symposium, International*.
- [8] P. H. Carns, W. B. Ligon, III, R. B. Ross, and R. Thakur. Pvfs: a parallel file system for linux clusters. In *Proceedings of the 4th annual Linux Showcase & Conference - Volume 4*, page 28, 2000.
- [9] C. Chilan, M. Yang, A. Cheng, and L. Arber. Parallel I/O performance study with HDF5, a scientific data package, 2006.
- [10] A. Ching, A. Choudhary, W. keng Liao, R. Ross, and W. Gropp. Noncontiguous i/o through pvfs.
- [11] P. Dickens and J. Logan. Towards a high performance implementation of mpi-io on the lustre file system. *On the Move to Meaningful Internet Systems: OTM 2008*, pages 870–885, 2008.
- [12] C. Docan, M. Parashar, and S. Klasky. Enabling high-speed asynchronous data extraction and transfer using DART. *Concurrency and Computation: Practice and Experience*, 22(9).
- [13] S. Donovan, G. Huizenga, A. J. Hutton, C. C. Ross, M. K. Petersen, and P. Schwan. Lustre: Building a file system for 1000-node clusters, 2003.
- [14] F. Isaila, J. G. Blas, J. Carretero, R. Latham, and R. Ross. Design and evaluation of multiple level data staging for Blue Gene systems. *IEEE Transactions on Parallel and Distributed Systems*, 99(PrePrints), 2010.
- [15] J. Li, W. Liao, A. Choudhary, R. Ross, R. Thakur, W. Gropp, R. Latham, A. Siegel, B. Gallagher, and M. Zingale. Parallel netCDF: A high-performance scientific I/O interface. page 39. IEEE, 2006.
- [16] M. Li, S. Vazhkudai, A. Butt, F. Meng, X. Ma, Y. Kim, C. Engelmann, and G. Shipman. Functional partitioning to optimize end-to-end performance on many-core architectures. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–12. IEEE Computer Society, 2010.
- [17] J. Lofstead, F. Zheng, Q. Liu, S. Klasky, R. Oldfield, T. Kordenbrock, K. Schwan, and M. Wolf. Managing variability in the IO performance of petascale storage systems. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '10, pages 1–12, 2010.

- [18] J. F. Lofstead, S. Klasky, K. Schwan, N. Podhorszki, and C. Jin. Flexible IO and integration for scientific codes through the adaptable IO system (ADIOS). In *Proceedings of the 6th international workshop on Challenges of large applications in distributed environments, CLADE '08*, pages 15–24, 2008.
- [19] X. Ma, J. Lee, and M. Winslett. High-level buffering for hiding periodic output cost in scientific simulations. *IEEE Transactions on Parallel and Distributed Systems*, 17:193–204, 2006.
- [20] A. Moody, G. Bronevetsky, K. Mohror, and B. R. de Supinski. Design, modeling, and evaluation of a scalable multi-level checkpointing system. *SC Conference*.
- [21] S. Moore. Multicore is bad news for supercomputers. *Spectrum, IEEE*, 45(11):15, 2008.
- [22] Mpich2. <http://www.mcs.anl.gov/research/projects/mpich2/>.
- [23] NICS. <http://www.nics.tennessee.edu/io-tips>.
- [24] A. Nisar, W. keng Liao, and A. Choudhary. Scaling parallel I/O performance through I/O delegate and caching system. In *High Performance Computing, Networking, Storage and Analysis, 2008. SC 2008. International Conference for*, pages 1 – 12, November 2008.
- [25] C. M. Patrick, S. Son, and M. Kandemir. Comparative evaluation of overlap strategies with study of I/O overlap in MPI-IO. *SIGOPS Oper. Syst. Rev.*, 42:43–49, October 2008.
- [26] J.-P. Prost, R. Treumann, R. Hedges, B. Jia, and A. Koniges. MPI-IO/GPFS, an Optimized Implementation of MPI-IO on Top of GPFS.
- [27] J. C. Sancho, K. J. Barker, D. J. Kerbyson, and K. Davis. Quantifying the potential benefit of overlapping communication and computation in large-scale scientific applications. In *SC 2006 Conference, Proceedings of the ACM/IEEE*, page 17, november 2006.
- [28] F. Schmuck and R. Haskin. GPFS: A shared-disk file system for large computing clusters. In *Proceedings of the First USENIX Conference on File and Storage Technologies*, pages 231–244. Citeseer, 2002.
- [29] H. Shan and J. Shalf. Using IOR to Analyze the I/O performance for HPC Platforms. In *Cray User Group Conference 2007*, Seattle, WA, USA, 2007.
- [30] D. Skinner and W. Kramer. Understanding the causes of performance variability in HPC workloads.
- [31] R. Thakur, W. Gropp, and E. Lusk. Data Sieving and Collective I/O in ROMIO. *Frontiers of Massively Parallel Processing, Symposium on the*.
- [32] A. Uselton, M. Howison, N. Wright, D. Skinner, N. Keen, J. Shalf, K. Karavanic, and L. Oliker. Parallel i/o performance: From events to ensembles. In *Parallel Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, april 2010.

- [33] F. Zheng, H. Abbasi, C. Docan, J. Lofstead, Q. Liu, S. Klasky, M. Parashar, N. Podhorszki, K. Schwan, and M. Wolf. PreDatA – preparatory data analytics on peta-scale machines. In *Parallel Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, pages 1–12, April 2010.

## A Evaluating Damaris on a Power5 cluster

This section presents the results obtained on the BluePrint Power5 cluster at NCSA.

### A.1 Platforms and configuration

**BluePrint** is the Blue Waters interim system running at NCSA. It provides 120 nodes, each featuring 16 1.9GHz POWER5 cpus, and runs the AIX operating system. 64 GB of local memory is available on each node. GPFS is deployed on 2 dedicated nodes. CM1 was run on 64 nodes (1024 cores), with a 960x960x300 points domain. Each core handled a 30x30x300 points subdomain with the standard approach. When dedicating one core out of 16 on each node, computation cores handled a 24\*40\*300 points subdomain.

### A.2 Experimental results

We have evaluated the behavior the CM1’s file-per-process method on BluePrint with different output configurations, ranging from a couple of 2D arrays to a full backup of all 3D arrays. Compression was enabled, but the amount of data will be reported as logical data instead of physical bytes stored.

By mapping contiguous subdomains to cores located in the same node, we have also been able to implement a filter that aggregates datasets prior to actual I/O. The number of files is divided by 16, leading to 64 files created per I/O phase instead of 1024. With a dedicated core, CM1 reports spending less than 0.1 sec in the I/O phase, that is to say only the time required to perform a copy of the data in the shared buffer.

The results are presented in Figure 13. For each of the four experiments, CM1 runs during one hour and writes every 3 minutes. We plot the average values of the time spent by the slowest process in I/O phase, the fastest, and the mean. We observe that the slowest process is almost five times slower than the fastest.

We measured the time that dedicated cores spent writing depending on the amount of data. This time is reported on Figure 14 and compared to the time spent by the other cores to perform computations between two I/O phases. We observe that even with the largest amount of data, the I/O core stays idle 75% of the time, thus allowing more data post-processing in order to help further analysis.



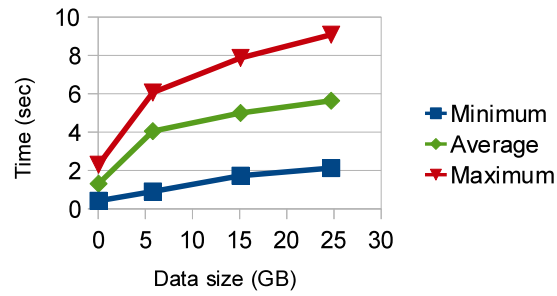


Figure 13: Maximum, minimum and average write time on BluePrint with different total output sizes, using the file-per-process approach, compression enabled.

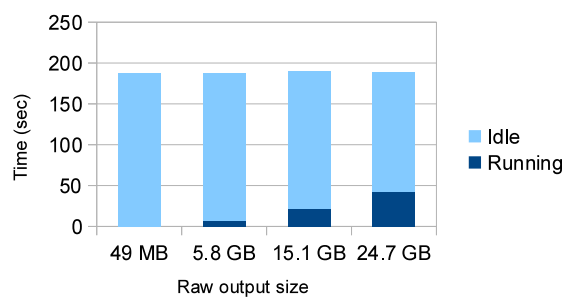


Figure 14: Write time and idle time in the I/O cores on BluePrint. Data size is the raw uncompressed data output from all dedicated cores together



**RESEARCH CENTRE  
RENNES – BRETAGNE ATLANTIQUE**

Campus universitaire de Beaulieu  
35042 Rennes Cedex

Publisher  
Inria  
Domaine de Volveau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399