

# Evaluation of HPC applications' Memory Resource Consumption via Active Measurement

Marc Casas, Greg Bronevetsky

## Abstract—

As the number of compute cores per chip continues to rise faster than the total amount of available memory, applications will become increasingly starved for memory storage capacity and bandwidth, making the problem of performance optimization even more critical. Also, understanding and optimizing the usage of an increasing number of hierarchical memory levels and complex cache management policies is becoming a very hard task. We propose a methodology for measuring and modeling the performance of hierarchical memories in terms of the application's utilization of the key memory resources: capacity of a given memory level and bandwidth between two levels. This is done by actively interfering with the application's use of these resources. The application's sensitivity to reduced resource availability is measured by observing the effect of interference on application performance. The resulting resource-oriented model of performance both greatly simplifies application performance analysis and makes it possible to predict an application's performance when running with various resource constraints. This is useful to predict performance for future memory-constrained architectures. This paper applies the proposed methodology to 6 important and well known High Performance Computing (HPC) codes to show the strength and the potential of analysis based on resource-oriented measurements.

**Index Terms**—Multi-core architectures, Memory Hierarchy, Performance Analysis.



## 1 INTRODUCTION

Modern computing systems achieve high performance by, among other things, the efficient usage to the memory hierarchy, which combines small amounts of expensive and fast memories with larger, cheaper and slower components. This hierarchical design provides a good balance between cost, performance and storage capacity. The main drawback that hierarchical designs have is their programmability, as it is very difficult to tune real codes to fully exploit the potential that hierarchical memories have. Significant research has been devoted to develop cache-friendly algorithms [13], [8] and performance analysis tools to understand applications' memory usage [30], [24], [18]. However, the goal of easy-to-use memory optimization techniques is still far from reach.

Modern architectural designs provide increasing improvements in computation capability while maintaining a constant power dissipation density by increasing the number of cores on each chip. Since the power and cost efficiency of memory designs are not improving at the same rate, the amount of memory per compute core is dropping [20]. This is especially true for High Performance Computing (HPC) systems, where hard limits on power costs will mean that next-generation Exascale systems may provide one or two

orders of magnitude less memory capacity and bandwidth per core than today's systems [20]. These limitations will force application designers to fundamentally rethink how their algorithms utilize the memory system and will make effective memory optimization methodologies critical for maintaining application performance on future systems.

Ensuring that applications optimally use the memory hierarchy or restructuring algorithms to leverage hierarchies that are deeper (more levels) and thinner (fewer resources per core) requires a detailed analysis of how an application uses memory. Although there exists a wide range of tools to help with this task, they have key limitations. Simulation-based tools such as *cachegrind* [25] and *gem5* [2] can analyze the application's behavior in great detail and can predict the performance of any collection of applications running on any hardware configuration. However, such tools run hundreds or thousands times slower than native execution and cannot simulate the commercial architectures on which almost all applications run because simulator developers have no access to their proprietary details. These limitations have motivated work on tools based on monitoring hardware performance counters. These tools report metrics such cache miss rates or instructions per cycle for various code regions [30], conduct complex analyses of such counter data [18], [3], [4] or connect them to other aspects of the application, such as data structures [24]. Although these tools are efficient and precisely capture the state of the hardware and how it is utilized by the application, this information is not actionable in most cases. First, the metrics reported by these tools are so low-level that they can only be interpreted by the most hardware-savvy developers. Further, this information is not useful for predicting how the application may behave in alternate scenarios, such as if its available resources are reduced by the execution of

---

*The research leading to these results has received funding from the European Research Council under the European Union's 7th FP (FP/2007-2013) / ERC GA n. 321253. M. Casas is supported by the Secretary for Universities and Research of the Ministry of Economy and Knowledge of the Government of Catalonia and the Co-fund programme of the Marie Curie Actions of the 7th R&D Framework Programme of the European Union (Contract 2013 BP\_B 00243)*

- M. Casas is with the Barcelona Supercomputing Center, Technical University of Catalonia, Barcelona  
E-mail: [marc.casas@bsc.es](mailto:marc.casas@bsc.es)
- G. Bronevetsky is with Google

other software or because the application runs on a new platform. The limitations of today’s techniques motivate the need for a new approach that combines the predictive capability of simulation-based tools with the high performance of counter-based methods for real commercial proprietary hardware.

This paper presents a new performance analysis technique that addresses this need by capturing the application’s effective use of the storage capacity of different levels of the memory hierarchy as well as the bandwidth between adjacent levels. Our approach models various memory components as resources and measures how much of each resource the application uses *from the application’s own perspective*. To the application a given amount of a resource is “used” if not having this amount will degrade the application’s performance. This is in contrast to the hardware-centric perspective that considers “use” as any hardware action that utilizes the resource, even if it has no effect on performance. For instance, while from the hardware perspective cache storage capacity is “used” when live data is stored in it, to the application it is “used” only if the data is part of the application’s active working set. This paper specifically focuses on measuring storage capacity in caches that are shared by multiple cores and the bandwidth between them and higher levels of the memory hierarchy. In addition to measuring use, we also quantify the application’s sensitivity to being provided less of the resource than it needs to run without suffering performance degradations. This predicts how well the application would run in scenarios where less of the resource is provided, such as the memory hierarchy of a future system (e.g. a node of an Exascale system).

We measure the application’s use of memory resources via the proactive methodology illustrated in Figure 1. While the application runs, it uses a given fraction of each resource at a given level of a memory hierarchy (denoted “level X cache”). If this level X cache is shared among multiple cores it is possible to measure this use by running on another core an interference thread that utilizes a known amount of cache capacity or bandwidth, where the use of a separate core limits the thread’s effects to just the chosen resource. The algorithm increases the amount of resource used until the main application’s performance is observed to degrade. The difference between the total amount of the resource and the amount used by the interference thread at that point is the amount actively used by the application. If the application’s performance is not sensitive to the interference then it either primarily uses a higher level of the memory hierarchy (little use of level X) or doesn’t fit in level X and is thus not sensitive to reductions in its capabilities. The two cases can be differentiated by observing the application’s miss rates for level X cache.

Overall, this paper’s contributions are:

- A methodology to actively measure the application’s resource use in terms of the effect of availability of this resource on its performance [5], [6].
- A novel interference-based mechanism to simulate a reduction in available memory storage and bandwidth on a given real hardware architecture.
- A validation technique for quantifying the real effects of our interference mechanisms on the application

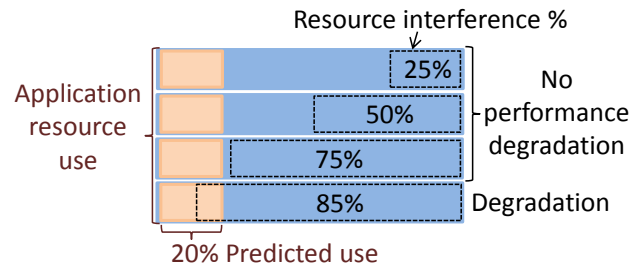


Fig. 1. The application’s resource use is measured by interfering with increasing fractions of the resource until this affects application performance

and bounding the effects of a given source on interference on unrelated resources.

- A method to predict how the application’s performance will degrade on alternative, less capable memory hierarchies.

Section 2 presents our interference measurement methodology. This approach is validated in Section 3 using microbenchmarks with well-characterized memory behavior. Section 4 describes in detail the HPC workloads used in this paper. Section 5 shows how our techniques can be applied to HPC applications to measure their cache storage and bandwidth requirements. Section 6 provides three examples of non-trivial insights that are provided by active measurement techniques.

## 2 MEASUREMENT METHODOLOGY

Each level of the memory hierarchy provides two resources: storage capacity and communication bandwidth to the (larger and slower) level below it. Our Active Measurement methodology measures the application’s use of these resources by comparing its performance when running on a given memory system to its performance when less of a given memory system resource is available due to the execution of a special interference thread. Specifically, the techniques presented in this paper focus on portions of the memory hierarchy that are shared among multiple cores. This restriction enables more precise measurements by making it easier to isolate the effects of the interference thread to analyze just the specific resource it is designed to target because the interference threads run on different cores that share the resource. This section details the design of these threads and Section 3 experimentally validates that each thread is effective at using up the resource that it targets and uses few other resources.

Our experiments focus on the following architecture: 2-socket nodes with 8-core Intel Xeon E5-2670 processors. The L1 and L2 caches are private to each core, while the L3 cache is shared among all the cores on a socket. The L3 on the Intel Xeon E5-2670 is 20MB in size and this architecture is thus denoted Xeon20MB (details in Table 1).

### 2.1 Memory Bandwidth Interference

The pseudo-code of the bandwidth interference thread, denoted BWThr is shown in Figure 2. This code attempts to transfer as much data as possible between one memory hierarchy level and the next by issuing a large number of

	Cache	Capacity	Line Size	Associativity
Private	L1 I	32KB	64 bytes	8-way
	L1 D	32KB	64 bytes	8-way
	L2	256KB	64 bytes	8-way
Shared	L3	20MB	64 bytes	20-way

TABLE 1

8-core Intel Xeon E5-2670 memory hierarchy

```

long long int* buf_0 =
    malloc(sizeof(long long int)*bufSize);
...
long long int* buf_numBufs =
    malloc(sizeof(long long int)*bufSize);

for(int i=0; 1; i++) {
    buf_0[identity(largePrime*i)%bufSize]++;
    ...
    buf_numBufs[identity(largePrime*i)%bufSize]++;
}

```

Fig. 2. Pseudo-code of the bandwidth interference thread BWThr

memory accesses that will miss in the first level. We induce frequent cache misses by allocating a buffer and iterating it with a stride that is coprime with the buffer size. The use of a large number keeps the number of iterations between adjacent accesses to the same location large, the constant stride makes it possible for the hardware prefetcher to help use up more bandwidth and the coprimality of the buffer size and the stride ensures that there are no aliasing effects, that is, enforces that the interference thread’s working set size is actually the whole memory buffer. Further, to ensure that the compiler cannot perform any optimizations based on the simple access pattern, the computation of the strided index  $largePrime \cdot i$  is wrapped inside a call to an `identity` function that is located in a different file and thus not available at compile time. One side-effect of this is that the compiler can no longer transform the loop to issue the maximum number of simultaneous memory accesses that the underlying hardware can support. We overcome this problem by simultaneously performing this procedure for many buffers at the same time (our experiments use 44, which we discovered to be sufficient), maximizing the concurrent memory traffic. Other approaches based on threads that make extensive use of the memory system have been used to study memory bandwidth and latency of IBM systems [26].

## 2.2 Cache Storage Interference

Figure 3 shows the pseudo-code of the storage interference thread, denoted CSthr. It allocates a buffer of a given size and then randomly touches elements in this buffer. If the buffer fits inside the high-level private caches of CSthr’s processor this iteration has little effect on the cache shared by it and the application. This is because once the data is fetched into this cache during the initial iterations of the inner loop, there will be no additional cache misses. However, once the array grows larger than the private caches, the random order of the memory accesses ensures that many

```

int* buf = malloc(sizeof(int)*bufSize);
while(1) buf[random_position]++;

```

Fig. 3. Pseudo-code of the storage interference thread CSthr

of the buffer accesses will miss in the private cache and will always hit in the lower-level shared cache. A random memory access pattern ensures a more intense perturbation than linear access patterns because the probability of consecutively accessing two addresses of the same cache line is very low, meaning that almost every access misses in the L1 and L2 and hits in the L3. Further, the use of random access ensures that the hardware pre-fetcher will not recognize the access pattern and thus will not fetch in additional addresses outside the target buffer. Because CSthr spends all of its time passing over the buffer, the application threads have little chance to use the cache space assigned to the buffer before CSthr accesses this line again and pulls this resource away from the application. The design of CSthr ensures that it predictably utilizes a fixed fraction of the target shared cache and prevents the application from making any productive use of it.

## 3 VALIDATION

In this section we evaluate the amount of storage and bandwidth resources the interference threads utilize and validate our results. We also demonstrate that each interference thread only utilizes its target resource, meaning that they affect application behavior orthogonally.

### 3.1 Memory Bandwidth Interference

The bandwidth used by BWThr is computed based on the number of L3 cache misses it incurs, which is obtained by reading the hardware performance counter available in the Xeon20MB machine. Each miss causes a full cache line to be transferred from main memory to the L3 of its core. Since BWThr also updates the array positions it access, its bandwidth usage is computed as:

$$BW = 2 \cdot \frac{\text{cache\_line\_size} \cdot \#\text{cache\_misses}}{\text{ExecutionTime}} \quad (1)$$

Our measurements indicate that using a 520KB buffer a single BWThr utilizes 5.6GB/s per core in Xeon20MB. Since Xeon20MB provides around 40GB/s of bandwidth between the L3 cache and memory according to the STREAM benchmark [23], 8 BWThr running on 8 different cores would consume the whole available bandwidth.

### 3.2 Cache Storage Interference

The CSthr utilizes cache storage by repeatedly accessing a range of memory addresses, attempting to deny the application their use. However, because there is a time window between adjacent touches by CSthr of the same cache location it is possible for the application to bring its own data into this location and make productive use of it before it is evicted by CSthr. As such, the total amount of storage utilized by CSthr cannot be computed directly and must be computed based on its effects on representative applications. As we describe in Section 3.2.3, we create several synthetic benchmarks with different well-known memory access patterns based on probability distributions. We then use the knowledge of a each benchmark’s distribution to derive the L3 cache miss rate that the benchmark would observe when running on a fully associative cache of a given

```

int* buf = malloc(sizeof(int)*bufSize);
for (int i=0; i<N_ACCESSES; i++) {
    int value = buf[X()];
    // Some computation involving value
}

```

Fig. 4. Probabilistic Memory Access Algorithm.  $X()$  is a random variable that has a probability distribution function  $f$  associated.

size. Given the L3 miss rate observed when each benchmark runs concurrently with CSThr we can then compute the effective cache size that is available to the benchmark at a given level of CSThr interference. Our experiments show that in most cases CSThr has a consistent effect for all the benchmarks and identify the narrow range of applications and interference levels for which the effects of CSThr cannot be accurately quantified.

### 3.2.1 Benchmark Design

Figure 4 shows the skeleton of the synthetic benchmarks we use to validate CSThr. It loops over a buffer  $N\_ACCESS$  times and in each iteration reads the value at a buffer index chosen randomly from some probability distribution and performs some number of computations on this location. Every time that a new index is generated, an access to the memory must be performed to get the data. In our experiments we created several different variants of this benchmark for a range of probability distributions and different degrees of memory access patterns. Table 2 lists all the distributions considered and represents both a wide range of access patterns as well as degrees of spatial locality (depends on the standard deviation, which is varied widely). The computation shown in the skeleton displayed by Figure 4 involves  $value$ , which is stored in a register. Then, the average time between two memory accesses can be varied by setting the computation performed after each read to be 1, 10 or 100 integer additions. By doing that, we consider not only the memory access pattern itself, but also several degrees of memory access frequency.

The next step in our analysis is to derive the expected cache miss rate for each benchmark given the amount of cache storage that is available. This is computed via the following formula, which uses the probability that a given randomly sampled index is in the cache to compute the Expected Hit Rate (EHR):

$$EHR = \sum_{i \in \text{buffer}} P(i \text{ is accessed}) \cdot P(i \in \text{cache}) \quad (2)$$

$P(i \text{ is accessed})$  is equal to the probability mass function  $f(i)$  of the distribution because we have designed a specific set of benchmarks with this property. The probability of  $i$  being in the cache, if the cache capacity is smaller than the buffer size, is equal to the probability that  $i$  has been accessed in any of the  $N$  previous memory accesses, being  $N$  the maximum number of memory addresses that can be stored in the cache, that is, the total cache capacity. Thus,  $P(i \in \text{cache})$  is equal to the product  $f(i) \cdot Cache\_capacity$ .

Pattern Name	Statistical Distribution	Distribution Parameters	Standard Deviation
Norm 4	Normal	$\mu=n/2 \sigma=n/4$	$n/4$
Norm 6	Normal	$\mu=n/2 \sigma=n/6$	$n/6$
Norm 8	Normal	$\mu=n/2 \sigma=n/8$	$n/8$
Exp 4	Exponential	$\lambda=4/n$	$n/4$
Exp 6	Exponential	$\lambda=6/n$	$n/6$
Exp 8	Exponential	$\lambda=8/n$	$n/8$
Laplace 4	Laplace	$\mu=n/2 \sigma=n/4$	$n/4$
Laplace 6	Laplace	$\mu=n/2 \sigma=n/6$	$n/6$
Logistics 4	Logistics	$\mu=n/2 s=n/4$	$n^2\pi^2/48$
Logistics 6	Logistics	$\mu=n/2 s=n/6$	$n^2\pi^2/108$
Logistics 8	Logistics	$\mu=n/2 s=n/8$	$n^2\pi^2/192$
Tri 1	Triangular	$a=0 b=0.4n c=n$	$n^2/18$
Tri 2	Triangular	$a=0 b=0.6n c=n$	$n^2/18$
Tri 3	Triangular	$a=0 b=0.8n c=n$	$n^2/18$
Uni	Uniform	$a=0 b=n$	$n^2/12$

TABLE 2

Memory access patterns considered.  $n$  is the size of the buffer.

Putting everything together:

$$\begin{aligned}
EHR &= \sum_{i \in \text{buffer}} f(i) \cdot f(i) \cdot Cache\_capacity \quad (3) \\
&= Cache\_capacity \cdot \sum_{i \in \text{buffer}} f(i)^2 \quad (4)
\end{aligned}$$

This formula makes several assumptions:

- The size of the buffer must be larger than the size of the cache. This assumption is satisfied just by using sufficiently large buffers.
- The formula applies to steady state execution, after the algorithm has warmed up the cache by loading its contents based on its probability distribution. We satisfy this assumption by setting  $N\_ACCESS$  to be much larger than the buffer sizes.
- It assumes that the cache is fully associative, which is not true in general. As such, despite the fact that the model is in general accurate, as it is demonstrated in section 3.2.2, it slightly underestimates the number of cache misses because set-associative miss ratios are in general larger than fully-associative ones.
- No considerations are made in terms of the impact of the cache line size or, in other words, it assumes that each cache line size is able to store just one element of the buffer. That is not true in case of real hardware but this simplification does not really impact model's accuracy since  $Cache\_Storage \gg Cache\_line\_size$ .

The dependence of the cache miss rate on cache size has been studied for generic applications in the past by using statistical models empirically derived [15]. Our model offers more insight, as it is not empirical, but it requires a precise description of the application's memory access pattern, which is not always available.

### 3.2.2 Validation of the Probabilistic Model

The variety of distributions used in our synthetic benchmarks induce a wide range of memory access patterns and L3 cache miss rates from below 10% to above 80%. The probability distributions with larger standard deviations

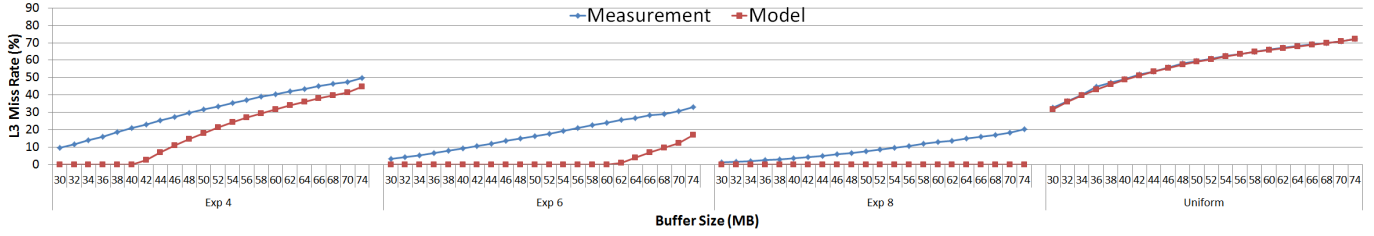


Fig. 5. Model Evaluation. Experiments with memory accesses patterns exp4, exp6, exp8 and Uniform.

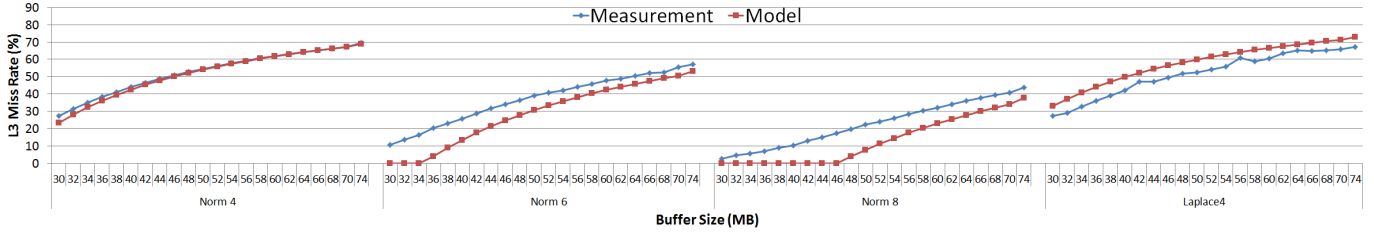


Fig. 6. Model Evaluation. Experiments with memory accesses patterns Norm4, Norm6, Norm8 and Laplace4.

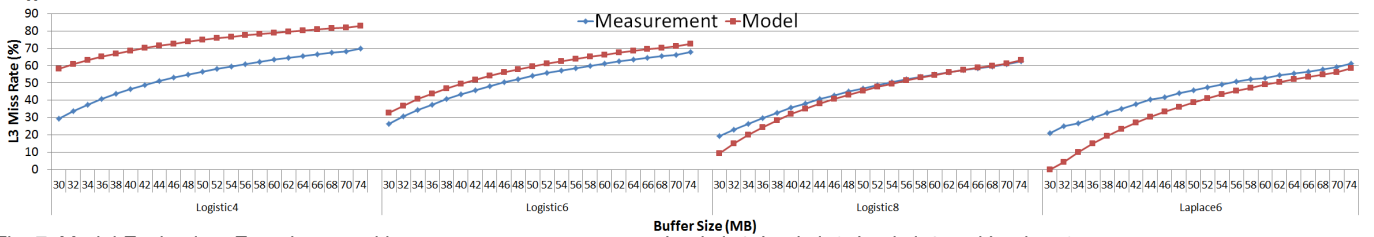


Fig. 7. Model Evaluation. Experiments with memory accesses patterns Logistic4, Logistic6, Logistic8 and Laplace6.

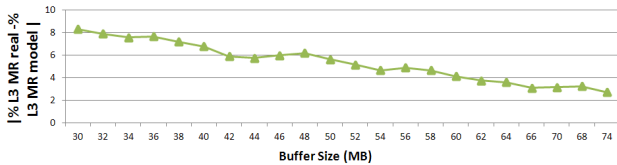


Fig. 8. Model Evaluation.

access with a higher probability wider sets of data, which decreases the chances of accessing two elements of the same cache line in a short period of time. This fact induces higher miss rates due to worse spatial memory locality. Further, cache miss rates rise as the buffer size increases since more memory is available for selection.

Equation 4 predicts the benchmark's miss rate based on its available storage capacity. We validated the equation by running the synthetic benchmarks defined in table 2 on Xeon20MB and comparing the predicted miss rates on the 20MB of L3 cache known to be available with the real cache miss rates measured by hardware counters.

In Figure 5 we show results considering patterns exp4, exp6, exp8 and Uniform considering buffer sizes from 30 up to 74 MB. The results show the impact of model's assumptions concerning cache associativity, as the real miss rates are significantly underestimated when they are below 25%. As the buffer size is increased the benchmarks' locality is reduced and model's assumptions have less impact. We can also observe the impact of the statistical distributions' standard deviations in the accuracy of the model. Those distributions with highest deviations (exp4 and uniform) are

better predicted since their memory access patterns are less impacted by model's assumptions.

Figures 6 and 7 show similar conclusions. Predictions improve as the buffer sizes increase as model's assumptions become less critical. Also, the model provides the most accurate predictions for those patterns defined by distributions with the largest deviations. The behavior of the three triangular distributions defined in Table 2 is very close to the one shown by the uniform distribution. The graphs concerning the three triangular distribution patterns have been omitted as they show the same trend as the uniform distribution pattern.

Figure 8 shows the absolute differences between these two numbers, averaged over all the distributions in Table 2. The average absolute distance between the measured cache miss rates and the predicted ones is always less than 10% and the average plus one standard deviation is 15% or less. Error is higher for small buffer sizes because the model assumes that caches are fully associative. This under-predicts the real cache miss rate when the cache is not heavily used and there are few cache misses. However, as the cache becomes fully utilized by larger buffer sizes and most memory accesses become misses (above 50% miss rate), the details of cache associativity become unimportant and the model's error drops to under 5%. In summary, the data shows that the simple analytic model of the behavior of the synthetic benchmarks is accurate in general, especially for large buffer sizes. Exhaustive studies on the relationship between cache associativity and miss rate have been done [16] and are consistent with the accuracy of our analytical model.

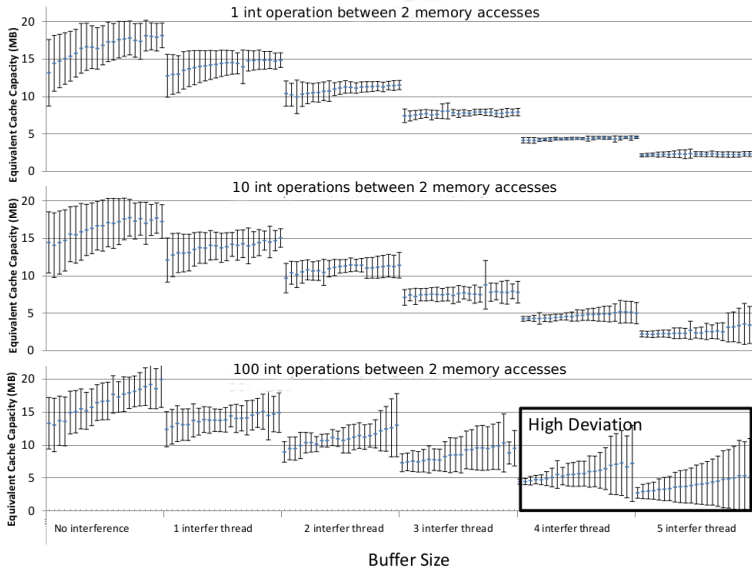


Fig. 9. Evaluation of the cache capacity interference. In the top, the benchmarks compute 1 integer sum between two consecutive memory accesses. In the middle, they perform 10 integer sums and at the bottom 100 integer sums are computed between memory accesses.

### 3.2.3 Measuring Cache Storage Use

Having shown that Equation 4 accurately predicts the L3 miss rate of the synthetic benchmarks based on the available cache storage capacity we can now use it to predict the effective storage available to these benchmarks when CSThr interferes with a portion of it. This is done by running experiments where CSThr interferes with a given synthetic benchmark’s use of cache storage, and measuring the resulting L3 miss rate. We then invert the formula in Equation 4 and given the observed miss rate compute the effective amount of available cache storage.

We conducted this evaluation on the Xeon20MB architecture, with 0 to 5 CSThrs, each using 4MB buffers. The synthetic benchmarks were parameterized with 10 different probability distributions (Table 2), 3 different degrees of memory access frequency (1, 10 and 100 integer additions per load) and 22 different buffer sizes from 30MB to 74MB for a total of 660 different configurations. The L3 cache miss rates of these variants range from less than 10% to more than 80%, with a similar variation on L3 cache memory access frequencies. This breadth of coverage makes our validation representative of a wide range of real applications, as it is pointed out in [31], where it is shown that L3 miss rates in SPEC 2006 applications range from 20% to 100% in commodity hardware, which is equivalent to the range covered by our 660 benchmark configurations.

The goal of this experiment is to determine whether the effects of different numbers of CSThrs are consistent across this wide range of memory access patterns. Figure 9 shows the results of our experiment. The charts from top to bottom show results with different degrees of memory access frequency, least frequency on top and most on the bottom. The charts from left to right show results with different numbers of CSThrs, with no interference on the left and then 1 through 5 CSThrs on the right.

Each chart shows on the y axis the amount of cache storage that is available to the benchmarks, as computed

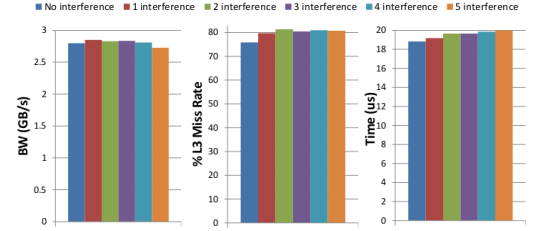


Fig. 10. Bandwidth from main memory to L3, L3 cache miss rate and time to do  $10^7$  iterations over its main loop of the BWThr when running concurrently with 0 and 5 CSThrs.

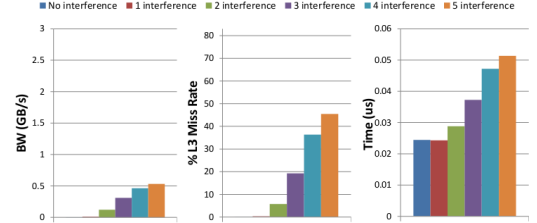


Fig. 11. Bandwidth from main memory to L3, L3 cache miss rate and time to perform a read, an arithmetic addition and a write of the CSThrs when running concurrently with 0 and 5 BWThrs.

by our formula and the x-axis shows the size of the buffer used by the benchmark (concrete buffer sizes omitted to improve readability and are the same as in Figure 8). In addition to the average predicted storage (across all the 10 probability distributions), denoted by the thick horizontal line, the charts show the region that includes the average plus and minus the standard deviation. The length of these intervals represents the dispersion of the measurements across all the distributions and is smaller for more consistent predictions of cache storage capacity. As such, the lower the length is, the more consistent is the estimation of the equivalent cache capacity.

The data displayed in the column tagged as “No Interference” shows results computed without any CSThr. As such, the predicted cache capacity should match the real cache capacity of the machine, which happens more clearly when buffer sizes close to 74MB are used. The buffer sizes close to 30MB show lower values because, as explained in Section 3.2.1, the model under-estimates the cache miss rates due to its assumption of fully associative caches. Thus, given the higher miss rates measured in real experiments, the inverse of Equation 4 under-estimates the available amount of cache storage. The error of the predictions drops as the buffer size is increased, as seen in Figure 8, until the predictions reach the correct capacity of 20MB.

The outcome of the experiments with no interference is the same for all degrees of memory access frequency (1, 10 or 100 operations between loads) since in these experiments the benchmarks do not compete for the cache. When 1 CSThr is used, the model predicts an effective cache capacity of approximately 15MB for all degrees of memory access frequency and 12MB with 2 CSThrs. When we use 3, 4, or 5 CSThrs the effective cache capacity is approximately 7, 5, and 2.5MB respectively.

One interesting phenomenon is that as the frequency of memory accesses decreases, the standard deviation of the predictions raises. Although this effect is weak for 1

CSThr, it becomes stronger as more CSThrs are added. Our experiments show that for low access frequencies the effects of CSThrs are erratic to a point where with 5 CSThrs interfere with either as much as the entire cache or as little as just half of it.

Overall, this validation quantifies the accuracy of the validation methodology and identifies the properties of applications for which it has high error: high degrees of interference and memory access frequency. Fortunately 100 arithmetic operations per load represent a very low memory access frequency. In particular, in scientific applications, which are typically more regular than most other application domains, 1-10 operations per load is considered the most common range. Our experiments thus show that our validation methodology is accurate for most real-world workloads.

### 3.3 Orthogonality of Cache Storage and Bandwidth Interference

Sections 3.1 and 3.2 have demonstrated that BWThr and CSThrs consistently utilize a given amount of cache bandwidth and storage, respectively. However, since both techniques use the memory hierarchy there is a risk that they utilize additional resources besides the ones they target. If so, they would interfere with the application in multiple ways, making each measurement reflect the use of a complex combination of system resources that has little intuitive meaning to the application developer. It is thus necessary to establish under which circumstances these threads have orthogonal effects: that each thread type almost exclusively affects the resource it targets and no other. By running the interference threads on separate cores we have ensured that no resources private to the application's cores are used. In this section we quantify the degree to which BWThr utilize cache storage and CSThrs utilize cache bandwidth.

To measure this we executed the BWThr (520KB buffers) and CSThrs (4MB buffers) simultaneously on different cores of the same Xeon20MB socket to measure each others' resource use. Figure 10 shows the effect on the execution of a single BWThr of concurrently running between 0 and 5 CSThrs, measured in terms of (i) the amount of bandwidth from the main memory to the L3 cache effectively used by the BWThr, (ii) the measured L3 cache miss rate of the BWThr, and (iii) the total time required to iterate  $10^7$  times over the BWThr's main loop, shown in Figure 2. The data shows that the BWThr behaves the same regardless of the number of CSThrs that are running concurrently with it. That implies that we can run up to 5 CSThrs without significantly impacting the memory bandwidth.

Figure 11 shows the opposite experiment: 1 CSThr running concurrently with between 0 and 5 BWThr. The bandwidth plot shows the bandwidth from main memory to L3 consumed by the CSThr and the execution time plot shows the average time the CSThr takes to perform a read, an arithmetic addition and a write operation. The data shows that a single BWThr has no impact on the CSThr's performance and 2 BWThr have a small effect. However, the CSThr is impacted significantly by the execution of 3, 4 and 5 BWThr, which implies that 3 or more BWThr utilize significant amounts of cache capacity. This induces L3

misses in the CSThr, which cause it to slow down and use more bandwidth. As discussed in Section 2.1, each BWThr utilizes 5.6GB/s of total memory bandwidth, which means that up to 11.2GB/s can be stolen without impacting cache capacity. Since the total memory bandwidth measured in the Xeon20MB architecture is near 40GB/s, that means we can impact on 28% of it while keeping the independence of the interference threads. Another important measurement provided in the left hand side of Figure 11 is the bandwidth utilization of the CSThr from main memory to L3. A single CSThr without additional interference utilizes very little memory bandwidth. These experiments show the design choices under which the orthogonality of CSThr and BWThr is guaranteed.

These results demonstrate that our measurements are indeed highly focused for a significant part of their dynamic range. Further, when their effects are orthogonal they can be used to represent the application's overall memory behavior as a simple 2-dimensional linear space where BWThr and CSThr identify the basis vectors: the application's utilization of storage and bandwidth. This projection enables application developers and architects to reason about an application's memory use in the same terms as if they were using a cache simulator without paying the big burden in terms of computing time that architectural simulation has.

## 4 APPLICATION SET

Our experiments focus on the following applications:

- AMG [11] - An implementation of the Algebraic Multi Grid Solver by using the Hypre library.
- FFTW [12] - Fast Fourier Transform library that uses hierarchical composition of multiple FFT algorithms, applied to perform a 2D transform of a matrix.
- Lulesh [19] - The Livermore Unstructured Lagrangian Explicit Shock Hydrodynamics simulation that is a materials science proxy application, executed on a cube domain.
- MCB [14] - A continuous energy Monte Carlo Burn-up Simulation Code for studying nuclear waste transmutation systems.
- MILC [1] - The MIMD Lattice Computation, a Quantum Chromodynamics simulation with lattice size.
- VPFIT [21] - A structure sensitive crystal plasticity simulation code.

The set of applications is representative of the typical workloads that run in HPC infrastructures. AMG carries out several iterations of an iterative solver over the same linear system at different levels of granularity, which means that it behaves like a CPU intensive benchmark when it operates over a dense representation of the system and like a communication and memory bound application when it performs solver iterations over a sparse representation of the system. Thus, AMG runs will display very different phases. FFTW and VPFIT applications contain expensive all-to-all communications. The difference between these two applications is that VPFIT performs expensive computation between two communication phases while FFTW does not. As such, VPFIT has some flexibility to overlap communication and computation while FFTW has much less. Lulesh is a typical

finite difference method code with local communication phases interleaved by intensive computation phases. MCB is a monte carlo simulation code, which means that it does not have much communication and, therefore, its usage of the interconnecting network is expected to be low. Finally, MILC spends most of its time running the conjugate gradient solver, which means that most of its communications involve point to point communications with the neighbors and global reductions once in a while.

## 5 PARALLEL APPLICATION STUDIES

Having presented and validated our basic measurement methodology we now show how it can be applied to parallel applications to gain insight into how they utilize the memory hierarchy. Our evaluation focuses on two performance aspects of the applications presented in section 4.

- Measuring the amount of L3 storage and L3↔Main Memory bandwidth used by these applications, and
- Characterizing their sensitivity to being provided less of either of these resources.

In each experiment we increase the amount of storage or bandwidth utilized by the BWThrs and CSThrs to observe both the point where the application’s execution time degraded (indicated the amount of resource used by the application) as well as the relationship between resource availability and performance degradation. Importantly, the slowdown of each process is stochastic, with individual instructions on different processors affected very differently by interference. This non-deterministic slowdown of instructions introduces noise into the application’s execution, which is a well-known source of slowdown for parallel applications [27], [17].

### 5.1 Experimental Setup

We run our parallel experiments on a cluster of NUMA nodes where each one of them has two Xeon20MB sockets. Nodes have 32GB of RAM and are connected via Infini-Band QDR (QLogic) interconnect (40Gb/sec bandwidth). Our measurements are performed by running the benchmarks on several nodes while allocating some cores to the application and executing BWThrs or CSThrs on some or all of the remaining cores. We run MCB on 24 MPI processes, Lulesh on 64 and AMG, MILC, VPFIT and FFTW on 72. For each application, we consider several mappings: 1, 2, 3, 4 and 6 MPI processes per Xeon20MB socket in case of MCB, MILC, VPFIT, AMG and FFTW. In case of Lulesh we map 1, 2, 4 and 8 MPI processes per socket due to some application specific restrictions. In bandwidth experiments we run 1 or 2 BWThrs on dedicated cores with a buffer size of 520KB each to interfere with up to 28% of the total memory bandwidth. In storage experiments we run between 1 and 5 CSThrs with a buffer size of 4MB each to utilize upto 87% of the total L3 cache capacity (17.5MB out of 20MB).

### 5.2 MCB

The top graphs of Figure 12 show in detail the performance degradations we measure for several different mappings of the processes of MCB to compute nodes when the input set

size is 20,000 particles. The x-axis corresponds to different numbers of CSThrs or BWThrs running on the available cores. The y-axis shows MCB’s execution time in different mapping and interference levels. Note that since different mappings leave different numbers of cores available, not all combinations of mapping and interference can be run. The bottom graphs show the performance degradations measured when running MCB on problems with 20,000 to 260,000 particles. They correspond to MCB runs on 24 Xeon20MB processors (12 nodes), with 1 MCB process and 7 available cores per processor. The graphs on the left focus on storage interference and the graphs on the right present bandwidth interference results.

In the top-left graph of Figure 12 we can see the consistency of the performance degradations across all the considered MPI mappings: the more processes mapped on each processor, the less cache capacity is available for each process and thus the same performance degradation is induced with fewer CSThrs. This representation suggests a simple way to calculate the average cache capacity utilization of MCB processes. For each process mapping, we consider the experiments with no performance degradation and pick the one that has the most CSThrs. We then consider the experiments with performance degradation and pick the one with the fewest CSThrs. Our prior analysis has determined that 1, 2, 3, 4, and 5 CSThrs with a 4MB buffer size leave 15, 12, 7, 4 and 3MB of cache capacity available to the application. We use this information to compute for each of the above configurations the ratio  $\frac{\text{Available\_cache\_capacity}}{\#\text{processes}}$  to obtain the upper and lower bound on the amount of storage available to each application process. For MCB running on 20,000 particles (top-left graph of Figure 12) each MCB process needs between 3.75 and 5MB of L3 capacity when 4 processes are mapped on each processor, 3.5-6MB are required when 2 run on each processor and 4-7MB for 1 process per processor. Performing the same analysis for the memory bandwidth we calculate that each MCB process’ bandwidth utilization is 8.6-10.0GB/s and 9.6-11.5GB/s when we map 4 and 3 processes per socket.

The bottom-left graph of Figure 12 shows MCB’s performance degradation as 1 through 5 CSThrs are executed on the available cores, which is equivalent to an L3 cache size of 15MB, 12MB, 7MB, 4MB and 3MB, respectively, according to the measurements described in Section 3.2.3. The data shows that when MCB simulates 20,000 to 260,000 particles, there is little performance degradation with one, two or three CSThrs and significant degradation of 20-25% with four or five CSThrs. This means that on this input range each MCB process uses between 4MB and 7MB of the L3 cache. The fact that even when given 2.5MB MCB’s performance only suffers by less than 30% indicates that this application would not perform much worse even if the L3 cache was not available on this architecture.

The bottom-right graph of Figure 12 shows MCB’s performance degradation when one or two BWThrs are executed to reduce the available bandwidth (as discussed above, running more than 2 BWThrs also uses up cache storage). The impact on MCB’s performance grows as the number of particles increases from 20,000 to 90,000 because its communication and thus miss rate grows with increasing workloads. Above 90,000 particles the impact of bandwidth



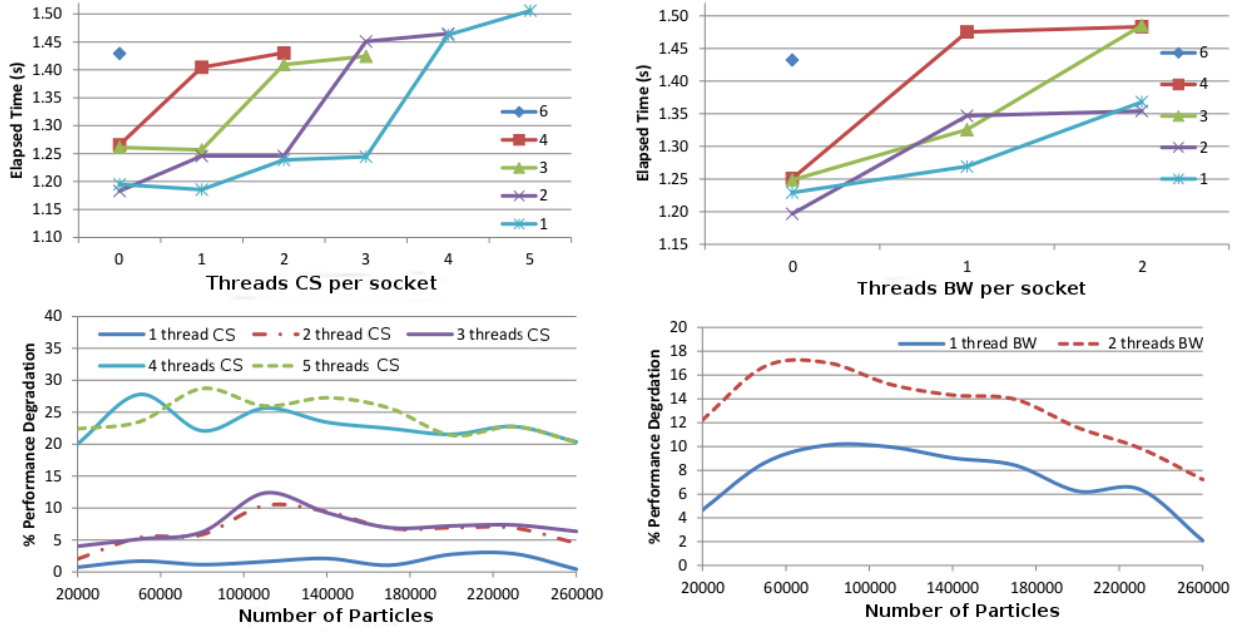


Fig. 12. Performance Degradations of MCB on 24 MPI tasks. The two top figures show results obtained considering several MPI mappings and using a 20,000 particles domain. The two figures at the bottom show results obtained when we map 1 MPI tasks per processor. Numbers of particles between 20,000 and 260,000 are considered.

interference drops because the application spends more time computing and less time communicating, which reduces the pressure on the memory buses.

### 5.3 Lulesh

Figure 13 shows the performance degradation of Lulesh [19] as it runs with 64 MPI processes. The physical domains simulated by Lulesh are cubes of sizes from  $22 \times 22 \times 22$  to  $36 \times 36 \times 36$  units (the size of one dimension is reported on the x axis), using the same format as Figure 12. The graphs in the top show performance degradations we measure running Lulesh on  $22 \times 22 \times 22$  cube across different process mappings. Experiments with 4 processes per processor show that Lulesh overflows the L3 cache when any number of CSThrs run, meaning that for all inputs each Lulesh process uses more than 3.5MB of storage (15MB that one CSThr leaves available, divided by 4 Lulesh processes per processor).

The bottom-left graph of Figure 13 shows the performance degradation of Lulesh with increasing number of CSThrs when mapping 1 process per socket. When domain size is  $32 \times 32 \times 32$  or smaller degradation is less than 5% for 1 and 2 CSThrs but more than 10% for 5 CSThrs, indicating that each Lulesh process uses between 2.5MB and 10MB of cache storage. For larger cubes Lulesh overflows the L3 cache with any amount of storage interference, suggesting that for these input sizes Lulesh processes use more than 15MB of cache each. The bottom-right graph of Figure 13 shows that performance degradation is larger than 10% when domain sizes are 32 and 36 with one or two BWThrs. This is consistent with the fact the L3 cache memory is not large enough for sizes bigger than 30 so the application needs the memory bus regularly to fetch data into the cache memory.

### 5.4 MIMD Lattice Computation (MILC)

Experiments obtained considering the MILC code are shown in Figure 14. The two top graphs are obtained using a lattice size size of  $n_x=16$ ,  $n_y=32$ ,  $n_z=32$ ,  $n_t=36$  as input. The top-left graph shows how MILC suffers no degradation when one MPI process is mapped per socket and less than 4 CSThrs threads co-run with it. Moderate performance degradation is experimented by MILC when 1 MPI process shares the socket resources with 4 CSThrs. The degradation becomes significant when 5 cache interference are active. This behavior indicates that each MILC MPI process needs between 5 and 7.5 MB of L3 shared cache when the input size is  $x=32$ ,  $n_y=16$ ,  $n_z=32$ ,  $n_t=36$ . In terms of memory bandwidth, the results are shown in the top right hand side of Figure 14. When just 1 or 2 MPI processes are mapped per socket, there is no impact. However, when more than 2 MPI processes are mapped per socket, there is a significant performance impact when just one BWThrs is executed. That means that 3 MPI processes already consume the whole memory bandwidth capacity of Xeon20MB processors, 40 GB/s. In the two bottom graph we show results obtained by mapping just 1 MPI process per socket and considering several lattice sizes. The  $n_x$  and  $n_y$  sizes are specified in the x-axis of the graph while  $n_z$  and  $n_t$  are 32 and 36 respectively. Larger lattice sizes show more sensitivity to shared memory resources. In particular, when the considered lattice size is equal to  $n_x=32$ ,  $n_y=32$ ,  $n_z=32$ ,  $n_t=36$ , we observe a significant sensitivity to shared L3 cache capacity (45% degradation when 5 CSThrs co-run with the main application) and memory bandwidth (12% degradation when 2 BWThrs are considered).

### 5.5 Crystal Viscoplasticity Proxy Application (VPFFT)

Results concerning the Crystal Viscoplasticity Proxy Application (VPFFT) are shown in Figure 15. The considered

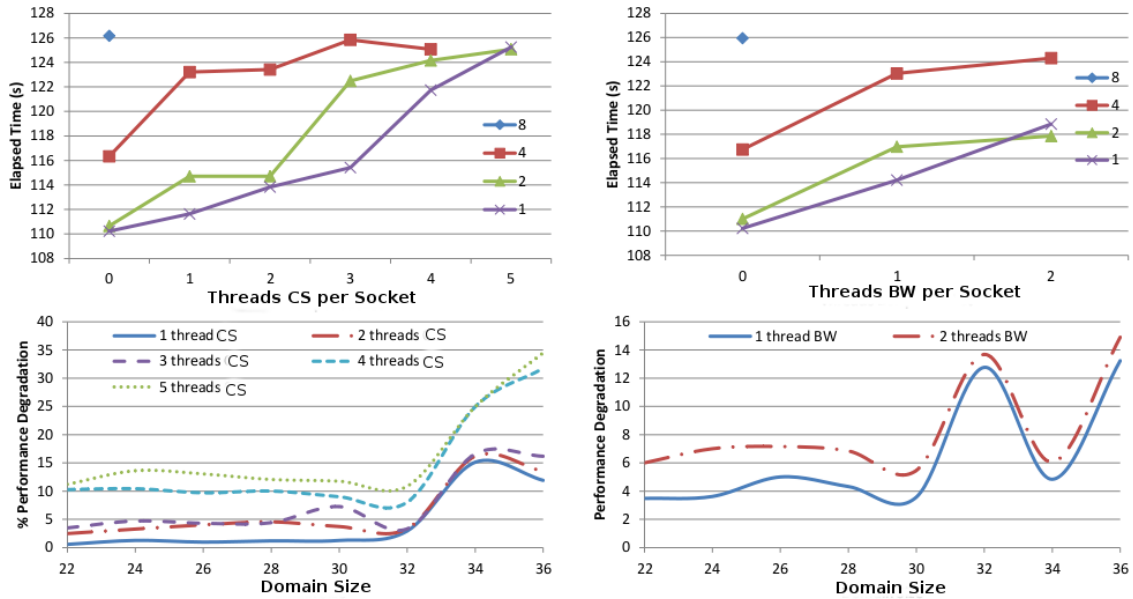


Fig. 13. Performance Degradations of Lulesh on 64 MPI tasks. The two top figures show results obtained considering several MPI mappings and using a  $22 \times 22 \times 22$  domain. The two figures at the bottom show results obtained when we map 1 MPI process per processor. The domains considered are between  $22 \times 22 \times 22$  and  $36 \times 36 \times 36$ .

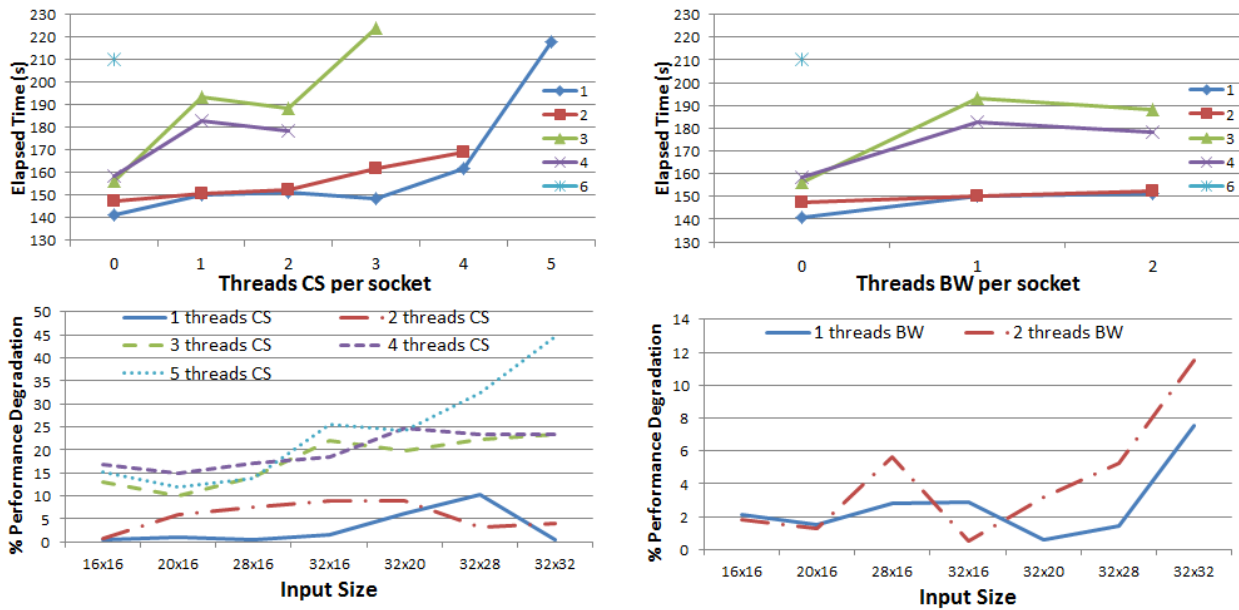


Fig. 14. Performance Degradations of MILC on 72 MPI tasks. The two top figures consider an input domain of  $n_x=16, n_y=32, n_z=32, n_t=36$ . The two graph at the bottom are obtained running a single MPI process per socket and considering several lattice sizes, from  $n_x=16, n_y=16, n_z=32$  and  $n_t=36$  to  $n_x=32, n_y=32, n_z=32$  and  $n_t=36$ .

input domain is a  $4 \times 4 \times 4$  cube and the maximum number of iterations is 300. The two graphs at the bottom show results obtained considering 5 different mappings, from 1 MPI process per socket up to 6. The impact of the interference threads, either if they are CSThrs or BWThrs, is negligible except when there are 3 or more MPI process in the same socket. In such scenarios, the application suffers significant performance slowdown when run together with 2 or more CSThrs. Since 2 CSThrs steal 10MB of the Xeon20MB L3 cache, we can infer that 3 MPI process need between 10 and 15MB to run at maximum performance, that is, each MPI process requires between 3.3 and 5MB of L3 cache storage.

A very interesting property of the VPFft application is that it gets some benefit from mapping at the same socket several MPI processes. As we can see in the top-left figure, there is a very important improvement if 3 MPI processes are mapped in the same socket instead of 1 or 2. Such improvement comes from that fact that it is much cheaper to exchange data between two cores of the same socket than two cores belonging to different sockets. However, the benefits of such approach are not that important when we map more than 3 MPI processes per socket. As said before, one MPI process needs between 3.3 and 5MB of L3 capacity which means that when we map more than 3 processes we are close to

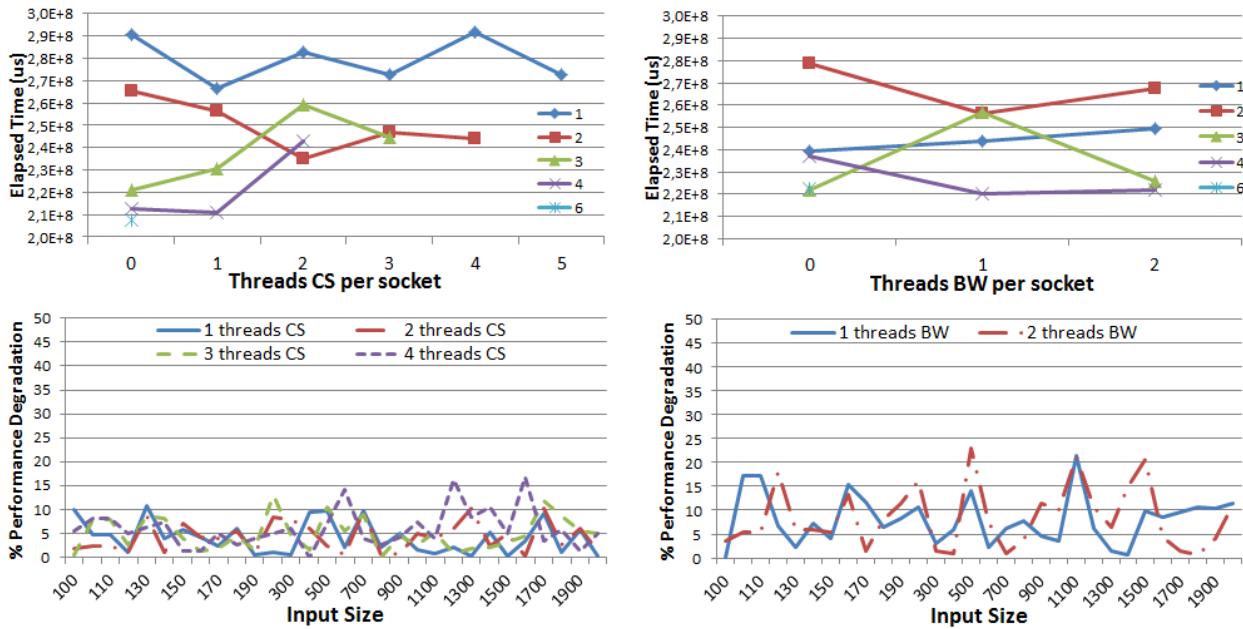


Fig. 15. Performance Degradations of VPFFT on 72 MPI tasks considering a cubic input domain. The two graphs in the top consider several mappings of MPI processes per socket. The two in the bottom consider several input sizes when mapping a single MPI process to a socket.

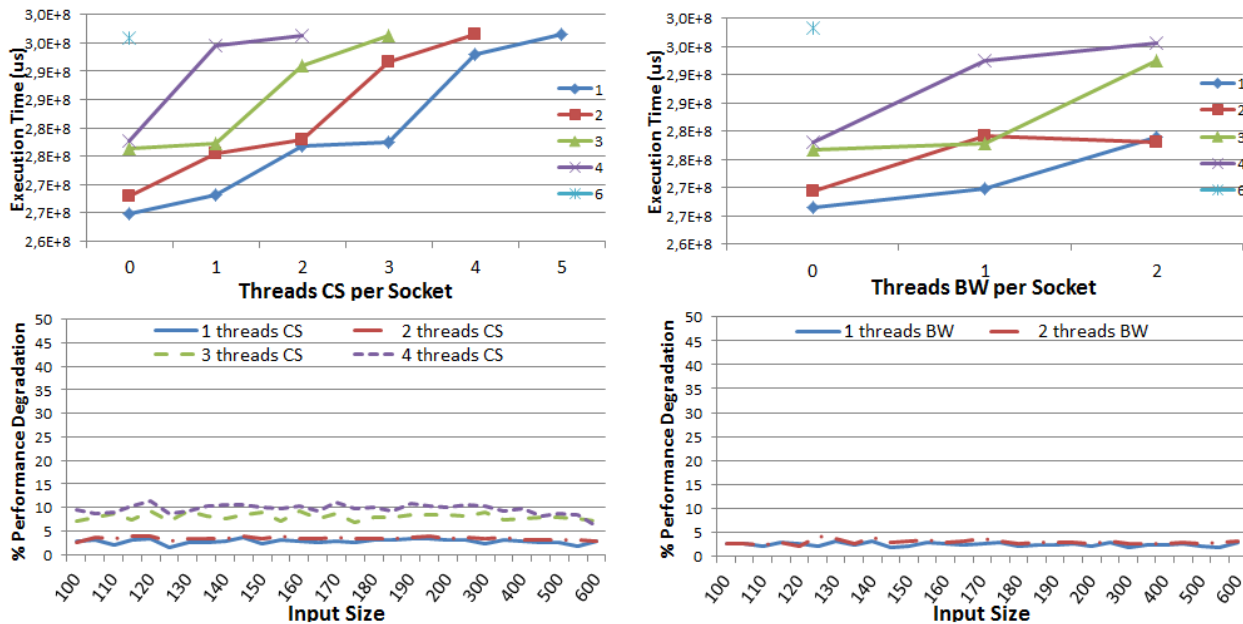


Fig. 16. Performance Degradations of AMG run on 72 MPI tasks and applied to a 3D Laplace problem on a cube. In the top we see results mapping different numbers of MPI processes to each socket and considering a  $160 \times 160 \times 160$  cube. In the two graphs displayed at the bottom we map one single MPI process per socket and show results considering different cube sizes. The x-axis display the size of the x dimension of the cube. The sizes of dimensions y and z are both 160.

the 20MB total capacity. Such lack of cache storage reduces the benefits of faster communications.

## 5.6 Algebraic Multi-Grid (AMG) Solver

The impact of shared memory resources on the AMG solver has also been evaluated. The results are shown in Figure 16. The two top graphs display results considering a 3D Laplace problem on a  $160 \times 160 \times 160$  cube. When we map a single MPI process to each socket, the significant performance degradation takes place when 4 CSThrs are active, which means that each MPI process needs between 5

and 7.5MB of the Xeon20MB L3 cache. In terms of memory bandwidth, significant performance degradations take place when 3 MPI processes are executed with 2 BWThrs in the same socket, meaning that the 3 MPI ranks need between 28.8 and 34.4 GB/s of memory bandwidth (each BWThrs consumes 5.6 GB/s, as we show in section 3). The two graphs at the bottom of Figure 16 show results mapping one MPI processes per socket and considering cube sizes from  $100 \times 160 \times 160$  up to  $600 \times 160 \times 160$ . As we can see in the figures, the considered problem sizes do not impact the application sensitivity to shared memory resources. In

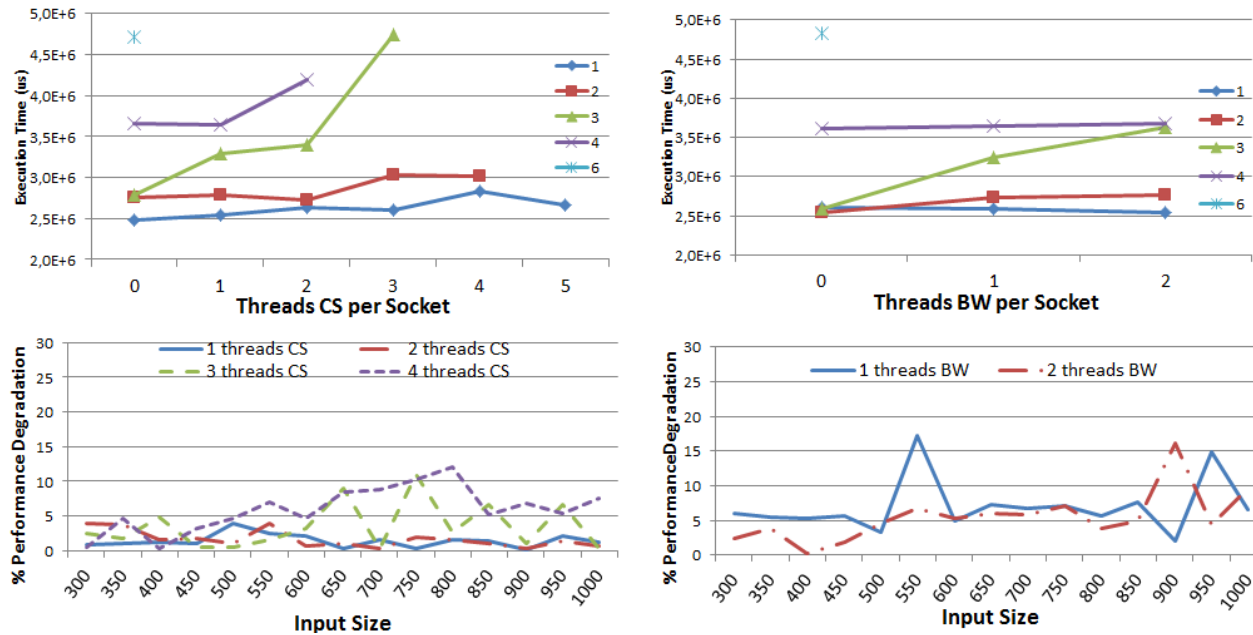


Fig. 17. Performance Degradations of suffered by FFTW performed over a 2D matrix and run on 72 MPI tasks. In the top we see results obtained from a 600x600 matrix and considering several mappings. The two graphs at the bottom show results mapping a single MPI process per socket and considering matrix sized from 300x300 up to 1000x1000.

any case, the measured performance degradations are never larger than 10%.

## 5.7 Fast Fourier Transform (FFTW)

Results concerning a 2D Fast Fourier Transformed implemented using the FFTW library [12] are shown in Figure 17. The two graphs in the top display results obtained by using a 600x600 matrix as input set. The top left graph shows results considering several mappings, from 1 to 6 MPI processes per socket, and different cache capacity interference degrees. The top right graph shows similar experiments considering the memory bandwidth interference threads. The experiments regarding the CSThrs interference show that when 3 MPI processes run together with a single interference thread there is already a significant performance perturbation. That means that the 3 MPI processes need between 15 and 20MB, which is between 5 and 6.6MB per MPI process. The experiments regarding the memory bandwidth, displayed at the top right graph of Figure 17, show that there is a significant performance slowdown when 3 MPI process are run together with a single BWThrs. This behavior indicates that 3 MPI process consume between 37.2 and 40 GB/s. The two figures at the bottom of Figure 17 show results concerning different input matrices, from 300x300 up to 1000x1000. The left figure shows results obtained by using different numbers of CSThrs. We can see that performance slowdowns become more significant once the input set sizes are increased, specially when considering matrices larger than 600x600. On the bottom right we see the results obtained when the application runs together with the BWThrs. The behavior seems less dependent of the input set size although there is some spiky behavior due the natural irregular behavior of MPI collectives that the FFT code has.

## 6 ACTIVE MEASUREMENTS USEFULNESS

Active measurement techniques are able to emulate application's behavior in computational environments with restricted memory resources, which allows a wide range of possible analysis from the parallel programmer's perspective. In this section, we describe three possible analysis based on the data shown in section 5.

**Finding the Optimal Mapping:** By putting a single MPI process per socket and progressively filling the empty cores with CSThrs or BWThrs, we can realize the amount of shared memory resources each MPI tasks needs to avoid performance penalties. The user could blindly try several possible mappings and select the best one in terms of performance, but that would not provide any insight. In contrast, our method provides the reason why an application experiments performance degradations when too many MPI processes are mapped in the same socket. For example, FFTW experiments a severe performance degradation when 6 MPI processes are mapped in the same socket (from 2.5 seconds of execution time when a single MPI process is mapped per socket to 4.5 seconds when 6 are mapped), as we can see in the two top graphs of Figure 17. Does this degradation come from the exhaustion of L3 cache capacity or is the lack of available memory bandwidth what brings such large degradations? By inspecting the left hand side graph in the top of Figure 17, we can see that such huge performance degradation appears when there is not enough cache capacity. The right hand side graph shows the performance degradation due to the lack of memory bandwidth, which is not that intense. Other tools like the Structural Simulation Toolkit (SST) [29] may give a similar insights but they require a significant simulation effort.

**Applications Behavior on Reduced Memory Resources:** Applications' requirements of shared memory resources may vary depending on the input data. While Performance

Tools based on sampling and visualization of hardware counters, like Paraver [28] or HPCtoolkit [24], are very good in providing an analysis of program performance, they cannot provide significant insight about the performance impact of cache storage reductions or input set size changes. Simulators are very expensive even when they use sophisticated statistical techniques [34]. In contrast, active measurement methods provide useful hints. For example, at the bottom of Figure 13 we display how the performance degradation that lulesh suffers under restricted cache capacity environments increases from 13% to 34% as the input cube-domain sizes grow from 22x22x22 to 36x36x36. Depending on the size of the problem, computational scientists may decide which machine should be used to run their simulations by taking into account analysis derived from active measurement techniques.

**Applications co-schedule:** Computational resources are scarce and potential co-scheduling opportunities are rarely explored because the potential performance degradation they may bring to parallel applications are unknown. However, active measurement techniques make possible to estimate such degradation and provide some hints to computational scientists on whether or not it is convenient to put more than one parallel workload in the same socket. Despite the fact that we do not conduct any co-schedule experiment in this paper, active measurements can certainly provide useful insights to do it. Also, the information derived from active measurements can be used to improve the scheduling of parallel workloads by feeding the OS or the runtime system with information regarding applications' resource needs.

## 7 RELATED WORK

Some approaches [32] and [7] use interference threads to steal available cache storage from a targeted application. However, they do not control the working set size of their cache-stealing routines or validate the exact resources they utilize.

Eklov et al's work [9], [10] is a valuable contribution. Like us, Eklov et al develop interference workloads that utilize cache storage and bandwidth. However, their work falls short in a few critical areas. First, our methodology validates the independence between the BWThr and the cache interference (CSThr), while Eklov et al do not consider the possible impact of the Bandwidth Bandit on the cache storage capacity. As such, the performance slowdowns reported by Eklov et al can be misleading since the Bandwidth Bandit can erase some cache lines and thus increase applications cache miss rate. The resulting performance slowdowns are thus difficult to interpret since they are caused by interference in multiple resources. Second, while BWThr is much simpler than Eklov et al's Bandwidth Pirate, experiments show that it is at least as effective at reducing available memory bandwidth. As reported in the top of page 6 of [10], the Bandwidth Pirate method can steal up to 4.6 GB/s (43% of total system bandwidth) with an unknown impact on cache storage. Our method steals up to 28% of total bandwidth without significantly impacting cache storage. Third, our work proposes a method to precisely evaluate the effective reduction in cache capacity due to

cache interference while Eklov et al use of a simple heuristic to estimate the point where their cache storage interference is inaccurate, what causes them to assign low accuracy to experiments where 66%-75% of cache capacity is stolen. Our more accurate validation methodology allows us to ensure accurate results when interfering with up to 87% of cache storage capacity.

In [22] a "bubble" kernel with tunable capacity and memory activity is presented. While the bubble is very useful on analyzing applications' performance degradation due to generic memory activity, it is not able to decompose such degradation into several factors, as our work does. As such, using the bubble we would be able to predict the performance interference between co-located applications, but not the particular resource that is exhausted. In [33] the approach is improved to provide accurate QoS control and maximized server utilization via interference measurement.

Finally, techniques to accurately measure the parameters of hardware components have been proposed [36], [35]. The authors argue that existing micro-benchmarks are inadequate, and present novel micro-benchmarks for determining the parameters of all levels of the memory hierarchy, including registers, all cache levels and the translation look-aside buffer. The experimental results show that the tool successfully determines memory hierarchy parameters on many current platforms.

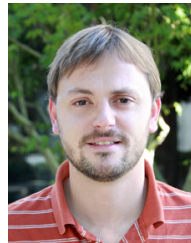
## 8 CONCLUSION

It is becoming critical to quantify how applications use the memory hierarchy to enable developers to identify performance bottlenecks. We present and validate the Active Measurement methodology, which quantifies an application's utilization of the memory hierarchy, specifically the storage and bandwidth of shared caches, and predicts the application's performance when the required memory resources are not available. Our approach is a significant improvement over the prior work. It provides information that today can only be derived by simulators but it is much faster than current simulation-based techniques and unlike simulation makes predictions for any architecture the application may run on with no information about its proprietary internal details. Further, it is much more actionable than performance counter analysis techniques because it can predict application performance when different amounts of key resources are available.

## REFERENCES

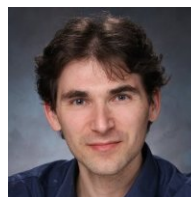
- [1] G. Bauer, S. Gottlieb, and T. Hoefler. Performance Modeling and Comparative Analysis of the MILC Lattice QCD Application `su3_rmd`. In *CCGRID*, pages 652–659, 2012.
- [2] A. Butko, R. Garibotti, L. Ost, and G. Sassatelli. Accuracy evaluation of GEM5 simulator system. *IEEE International Workshop on Reconfigurable Communication-centric Systems-on-Chip (ReCoSoC)*, July 2012.
- [3] M. Casas, R. Badia, and J. Labarta. Automatic analysis of speedup of MPI applications. In *Proceedings of the 22Nd Annual International Conference on Supercomputing, ICS '08*, pages 349–358, New York, NY, USA, 2008. ACM.
- [4] M. Casas, R. Badia, and J. Labarta. Prediction of behavior of MPI applications. In *Cluster Computing, 2008 IEEE International Conference on*, pages 242–251, Sept 2008.

- [5] M. Casas and G. Bronevetsky. Active measurement of memory resource consumption. In *Proceedings of the 2014 IEEE 28th International Parallel and Distributed Processing Symposium, IPDPS '14*, pages 995–1004, 2014.
- [6] M. Casas and G. Bronevetsky. Active measurement of the impact of network switch utilization on application performance. In *Proceedings of the 2014 IEEE 28th International Parallel and Distributed Processing Symposium, IPDPS '14*, pages 165–174, 2014.
- [7] X. Chen, C. Xu, R. P. Dick, and Z. M. Mao. Performance and power modeling in a multi-programmed multi-core environment. *Proc. of the Design Automation Conference (DAC)*, 2010.
- [8] T. M. Chilimbi, M. D. Hill, and J. R. Larus. Cache-conscious structure layout. *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 1–12, 1999.
- [9] D. Eklov, N. Nikoleris, D. Black-Schaffer, and E. Hagersten. Cache pirating: Measuring the curse of the shared cache. *Proc. of the International Conference on Parallel Processing (ICPP)*, 2011.
- [10] D. Eklov, N. Nikoleris, D. Black-Schaffer, and E. Hagersten. Bandwidth bandit: quantitative characterization of memory contention. In *Proceedings of the 21st international conference on Parallel architectures and compilation techniques, PACT '12*, pages 457–458, New York, NY, USA, 2012. ACM.
- [11] R. D. Falgout and U. M. Yang. hypre: a library of high performance preconditioners. In *Preconditioners, Lecture Notes in Computer Science*, pages 632–641, 2002.
- [12] M. Frigo and S. G. Johnson. The Design and Implementation of FFTW3. *Proceedings of the IEEE*, 93(2):216–231, 2005. Special issue on “Program Generation, Optimization, and Platform Adaptation”.
- [13] M. Frigo, C. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. *Proc. of the IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 285–297, 1999.
- [14] N. Gentile and B. Miller. Monte carlo benchmark. <http://www.osti.gov/estsc/details.jsp?rcdid=4793>.
- [15] A. Hartstein, V. Srinivasan, T. R. Puzak, and P. G. Emma. On the nature of cache miss behavior: Is it 2? *J. Instruction-Level Parallelism*, 10, 2008.
- [16] M. D. Hill and A. J. Smith. Evaluating associativity in cpu caches. *IEEE Trans. Comput.*, 38(12):1612–1630, Dec. 1989.
- [17] T. Hoefler, T. Schneider, and A. Lumsdaine. Characterizing the influence of system noise on large-scale applications by simulation. *Proceedings of the 2010 ACM/IEEE conference on Supercomputing (SC)*, 2010.
- [18] K. A. Huck, A. D. Malony, S. Shende, and A. Morris. Scalable, automated performance analysis with tau and perexplorer. *PARCO*, pages 629–636, 2007.
- [19] I. Karlin, A. Bhatele, J. Keasler, B. L. Chamberlain, J. Cohen, Z. DeVito, R. Haque, D. Laney, E. Luke, F. Wang, D. Richards, M. Schulz, and C. Still. Exploring traditional and emerging parallel programming models using a proxy application. In *27th IEEE International Parallel & Distributed Processing Symposium (IEEE IPDPS 2013)*, Boston, USA, May 2013.
- [20] P. Kogge. Exascale computing study: Technology challenges in achieving exascale systems. Technical report, DARPA IPTO, 2008.
- [21] R. A. Lebensohn, A. K. Kanjarla, and P. Eisenlohr. An elastoviscoplastic formulation based on fast fourier transforms for the prediction of micromechanical fields in polycrystalline materials. *International Journal of Plasticity*, 3233(0):59 – 69, 2012.
- [22] J. Mars, L. Tang, R. Hundt, K. Skadron, and M. L. Soffa. Bubble-up: increasing utilization in modern warehouse scale computers via sensible co-locations. In *Proceedings of the 44th Annual IEEE/ACM International Symposium on Microarchitecture, MICRO-44 '11*, pages 248–259, New York, NY, USA, 2011. ACM.
- [23] J. D. McCalpin. Memory bandwidth and machine balance in current high performance computers. *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter*, pages 19–25, 1995.
- [24] J. Mellor-Crummey, R. Fowler, and D. Whalley. Tools for application-oriented performance tuning. *Proc. of the 15th International Conference on Supercomputing*, pages 154–165, 2001.
- [25] N. Nethercote, R. Walsh, and J. Fitzhardinge. Building workload characterization tools with valgrind, October 2006.
- [26] D. M. Pase. Linpack HPL performance on IBM eserver 326 and xseries 336 servers. *IBM Technical Report*, 2005.
- [27] F. Petrini, D. J. Kerbyson, and S. Pakin. The case of the missing supercomputer performance. *Proceedings of the 2003 ACM/IEEE conference on Supercomputing (SC)*, 2003.
- [28] V. Pillet, J. Labarta, T. Cortes, and S. Girona. Paraver: A tool to visualize and analyze parallel code. In *Proceedings of WoTUG-18: Transputer and occam Developments*, volume 44, pages 17–31. mar, 1995.
- [29] A. F. Rodrigues, R. C. Murphy, P. Kogge, and K. D. Underwood. The structural simulation toolkit: Exploring novel architectures. In *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing, SC '06*, New York, NY, USA, 2006. ACM.
- [30] M. Schulz, J. Galarowicz, D. Maghrak, W. Hachfeld, D. Montoya, and S. Cranford. Open — speedshop: An open source infrastructure for parallel performance analysis. *Scientific Programming*, pages 105 – 121, 2008.
- [31] K. Singh, M. Bhadauria, and S. A. McKee. Real time power estimation and thread scheduling via performance counters. *SIGARCH Comput. Archit. News*, 37(2):46–55, July 2009.
- [32] C. Xu, X. Chen, R. P. Dick, and Z. M. Mao. Cache contention and application performance prediction for multi-core systems. *Proc. of the Intl. Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2010.
- [33] H. Yang, A. Breslow, J. Mars, and L. Tang. Bubble-flux: precise online qos management for increased utilization in warehouse scale computers. *SIGARCH Comput. Archit. News*, 41(3):607–618, June 2013.
- [34] J. J. Yi, D. J. Lilja, and D. M. Hawkins. Improving computer architecture simulation methodology by adding statistical rigor. *Computers, IEEE Transactions on*, 54(11):1360–1373, 2005.
- [35] K. Yotov, K. Pingali, and P. Stodghill. Automatic measurement of memory hierarchy parameters. *Proceedings of the 2005 ACM SIGMETRICS international conference*, pages 181–192, 2005.
- [36] K. Yotov, K. Pingali, and P. Stodghill. X-ray: A tool for automatic measurement of hardware parameters. *Proceedings of the 2nd Second International Conference on the Quantitative Evaluation of Systems*, pages 168–177, 2005.



**Marc Casas** Marc Casas is a senior researcher at the Barcelona Supercomputing Center (BSC). He received his BS and MS degrees in mathematics from the Technical University of Catalonia (UPC), and his PhD degree in 2010 from the Computer Architecture Department of UPC. He was a postdoctoral research scholar at the Lawrence Livermore National Laboratory (LLNL) from 2010 until 2013 working on algorithmic-based fault tolerance and active measurement methods based on software interference. He has

received several awards, like a Marie Curie Fellowship, and some of his papers have been awarded in conferences like Euro-Par or SC. His current research interests are focused on all the aspects involving the interaction between the parallel programming model, the runtime system and the hardware.



**Greg Bronevetsky** Greg Bronevetsky received his PhD degree in Computer Science from Cornell University in 2006. He was then awarded the Lawrence post-doctoral fellowship at Lawrence Livermore National Laboratory (LLNL) and became a full staff member in 2009. He received the DOE Early Career Award in 2009 and the Presidential Early Career award for Scientists and Engineers (PECASE) in 2010. Dr Bronevetsky joined Google in 2015. His research focuses on understanding and managing the behavior of

computer software and hardware, including static analysis of software, modeling application performance, root cause analysis of failures and software-level resilience.