

# Predicting Performance on SMPs. A Case Study: The SGI Power Challenge

Nancy M. Amato<sup>†</sup>   Jack Perdue<sup>†</sup>   Andrea Pietracaprina<sup>‡</sup>   Geppino Pucci<sup>‡</sup>   Mark Mathis<sup>†</sup>

<sup>†</sup>Department of Computer Science  
Texas A&M University, College Station, TX, USA.  
{amato,jkp2866,mmathis}@cs.tamu.edu

<sup>‡</sup>Dipartimento di Elettronica e Informatica  
Università di Padova, Italy.  
{andrea,geppo}@artemide.dei.unipd.it

## Abstract

*We study the issue of performance prediction on the SGI-Power Challenge, a typical SMP. On such a platform, the cost of memory accesses depends on their locality and on contention among processors. By running a carefully designed suite of microbenchmarks, we provide quantitative evidence that memory hierarchy effects impact performance far more substantially than other phenomena related to contention. We also fit three cost functions based on variants of the BSP model, which do not account for the hierarchy, and a newly defined function  $F$  expressed in terms of hardware counters, which captures both memory hierarchy and contention effects. We test the accuracy of all the functions on both synthetic and application benchmarks showing that, unlike the other functions,  $F$  achieves an excellent level of accuracy in all cases. Although hardware counters are only available at run-time, we give evidence that function  $F$  can still be employed as a prediction tool by extrapolating values of the counters from pilot runs on small input sizes.*

## 1 Introduction

Despite the vast body of ingenious parallel algorithmic techniques developed over the last two decades, the widespread use of parallel computers is still hampered by the difficulty of exploiting their massive computational potential to an extent that warrants their large cost. Indeed, it has often been noted that theoretically efficient algorithms

exhibit poor performance when implemented on real machines. Very often, this is due to the inadequacies of the cost functions employed to predict performance, which do not properly account for – or totally disregard – aspects of the machine that have a major impact on performance.

Although much progress has been made, the development of adequate tools for accurately predicting performance on real machines remains one of the most challenging problems in parallel processing. We believe that further progress towards this goal requires a tighter coupling of cost models to architectures than has been previously employed.

In this work, we study the issue of performance predictability on *Symmetric MultiProcessors (SMPs)*. These widespread parallel platforms are composed of powerful off-the-shelf microprocessors interacting through a distributed shared-memory via a communication medium, such as a bus. In such a system, the cost of an access to a shared datum may vary dramatically: from a few cycles if the data is in first-level cache (L1), to tens of cycles for second-level (L2) cache, to hundreds of cycles if the data must be accessed from main memory. The cost may be even greater in the presence of high contention among the processors for the bus or memory banks, or of (false) data sharing.

**Our Contribution** In this paper, we study the relative impact on performance of hierarchy and contention phenomena on an SGI-Power Challenge (SGI-PC), which is a typical representative of the class of SMPs. More specifically, we designed a suite of synthetic microbenchmarks which exercise different usages of the hierarchy under a set of controlled scenarios obtained by varying the level and type of contention among the processors. Based on the access times measured through the microbenchmarks, we infer parameter values for a set of linear cost functions inspired by some variants of the popular *Bulk-Synchronous Parallel (BSP)* model [14], and for a newly defined function describing the memory hierarchy usage of a program which relies on the MIPS R10000 hardware counters. The “BSP-like” func-

---

\*This research was supported in part by NATO CRG 961243 “Bulk Synchronous Computational Geometry,” and by NCSA grant CCR970010N. The work at Texas A&M was also supported by the NSF CAREER award CCR-9624315 and grants IRI-9619850, ACI-9872126, EIA-9805823, EIA-9810937, by DOE ASCI ASAP (Level 2 Program) grant B347886, and by the Texas Higher Education Coordinating Board grant ARP-036327-017. Perdue and Mathis supported in part by Dept. of Education Graduate Fellowships. The work at Padova was also supported by MURST of Italy under project “Algorithms for Large Data Sets: Science and Engineering.”

tions account for bus and bank contention and for data sharing but disregard hierarchy effects, while the counter-based function attempts to encompass all of these phenomena.

We test the accuracy of the cost functions on the whole set of microbenchmarks, and on application benchmarks, namely, the NAS suite of parallel benchmarks and three sorting algorithms. Our tests show that the function based on hardware counters, which accounts for the crucial impact of the memory hierarchy on performance, achieves an excellent level of accuracy in all cases, with times less than a factor of 2 away from actual times on average, and less than a factor of 3 in the worst case. In contrast, the cost functions inspired by the *BSP*-like models, which disregard hierarchy effects, provide performance predictions that can be more than two orders of magnitude away from actual times. Our study provides quantitative evidence that the memory hierarchy is the primary factor that affects performance on SMPs, while the impact of the other phenomena, although noticeable, is considerably less crucial.

It has to be remarked that while the *BSP*-like cost functions are easily computed *a priori* by code inspection, the counter-based function involves quantities (such as the number of accesses at the various levels of the memory hierarchy) that are easily computable only at run-time. However, we claim that the latter function can still be used as a prediction tool whenever accurate guesses of such quantities can be inferred from the code or extrapolated from pilot runs on small input sizes. In order to validate our claim, we provide examples of performance prediction for large sorting instances based on extrapolation of the relevant counters.

The counter-based function may also prove useful in the design of software systems, compilers, or large applications, to profile the memory hierarchy usage of critical portions of their code. Additionally, given a history of performance for an algorithm with varying input data sets on an architecture, a run-time system could use past performance data to predict the potential impact of rescheduling the order in which processes are provided system resources.

**Previous Work** The issue of performance prediction of parallel software has received considerable attention over the last decade, often within the context of the more general quest for a bridging model of parallel computation, i.e., one that balances conflicting requirements such as simplicity, accuracy and generality. One of the most popular attempts at defining a bridging model is Valiant's *BSP* model [14]. The *BSP* is a bulk-synchronous model where computation is organized as a sequence of *supersteps* separated by barrier synchronizations, and processors operate asynchronously within each superstep. The large body of work this model has generated has demonstrated its suitability for the development of portable software (see e.g., [9]).

Although the original *BSP* is meant to model message-passing architectures, two *BSP* variants specifically tailored to shared-memory systems have been recently developed. Namely, the *Queuing Shared Memory (QSM)* [8] and the  $(d, x)$ -*BSP* [4], which both embody some aspects of memory contention. In particular, *QSM*'s cost function includes a parameter that accounts for the maximum number of concurrent accesses to the same memory location, while  $(d, x)$ -*BSP*'s cost function accounts for memory bank contention (parameters  $d$  and  $x$  represent, respectively, bank delay and the ratio of the number of banks to processors). The set of *BSP*-like cost functions considered in this paper include the *QSM* and the  $(d, x)$ -*BSP*, and a third function, inspired by the *Extended BSP (EBSP)* model [10], which extends the *BSP* to account for unbalanced communication. Although *EBSP* was meant to model message-passing systems, we obtain the cost function by reinterpreting its original definition for a shared-memory machine.

In [1] the *Parallel Memory Hierarchy (PMH)* model is introduced which uses a single mechanism to model both interprocessor communication and the memory hierarchy in a parallel computer through a tree-structured view of the machine's organization. Although the model encompasses parameters which characterize the performance at each level of the tree, it does not provide a global cost function that can be used to predict program performance.

Finally, hardware counters are extensively used to profile sequential and parallel code. Example uses of the counters on the SGI-PC can be found in [16].

## 2 Hardware and Software Platforms

The SGI-PC configuration we used consists of eight R10000 194 MHz processors, each provided with a 32 KB on-chip instruction cache, a 32 KB on-chip level-1 (L1) data cache, and a 1 MB off-chip unified (instructions and data) level-2 (L2) cache. Cache line size is 32B (8 words) for L1 and 128B (32 words) for L2. Both L1 and L2 are two-way set associative. An 8-way interleaved, 2 GB main memory distributed across 8 banks is accessed by the processors through a 1.2GB/s shared-bus using a cache-coherent protocol [12].

All our experiments on the SGI-PC were coded according to an SPMD bulk-synchronous programming style [14, 8], where all processors execute the same program consisting of a sequence of supersteps separated by barriers. In a superstep, each processor performs a number of *memory accesses* (load or store instructions) on words which may reside either in the processor's L1/L2 caches or in main memory, and a number of *local operations* on data held in registers. Barriers were implemented using the SGI native `m_sync()` primitive. In this work, we are mainly interested in predicting the cost of memory accesses and we will

not deal with local operations.

The running time of each superstep was measured by mapping the CPU cycle counter to memory (`syssgi()` and `mmap()`), reading that value as each superstep started and ended, and using a scaling factor (provided by `syssgi()`) to convert clock cycles to microseconds. Also, in each superstep we monitored loads/stores issued, and L1/L2 cache misses/writebacks at each individual processor by means of some hardware counters provided by the R10000 design [13]. We verified experimentally that the counters are non-intrusive. (The employed counters are described in Section 4.2.)

### 3 Experimental Test Bed

In this section, we describe a suite of simple microbenchmarks whose purpose is to measure the cost of accessing the SGI-PC memory system under a variety of scenarios. The suite was designed with the intention of ascertaining the relative impact of the following key phenomena on access time:

**Locality:** the level of the hierarchy where the accesses take place (i.e., either L1/L2 caches or main memory);

**Bus Contention:** the volume of bus traffic generated by the accesses;

**Bank Contention:** the number of accesses directed to the same memory bank;

**Coherence:** the different coherency activities triggered by the accesses.

In the generic microbenchmark, a number of *active processors* perform a sequence of accesses (either a load or a store sequence) to (possibly coincident) sub-arrays of a large shared array. The sequence of accesses is iterated several times to filter out noise in the measurements and cold-start effects. In order to exercise different combinations of the four phenomena illustrated above, we instantiate this generic microbenchmark by varying the number of active processors, which affects bus contention, and by varying access stride, size and base address of the sub-arrays, which affects locality, bank contention and coherence. More specifically, in a microbenchmark, the stride, the sub-array size, and the number of accesses for each active processor are the same, while suitable base addresses are chosen so that processors work either exclusively on distinct sub-arrays or concurrently on the same sub-array. Finally, all the accesses are evenly distributed among a number of *target banks*, so that each target bank serves the same number of requests.

We designed microbenchmarks with 1, 2, 4, 8 processors active, strides ranging from 1 to 256 words, three

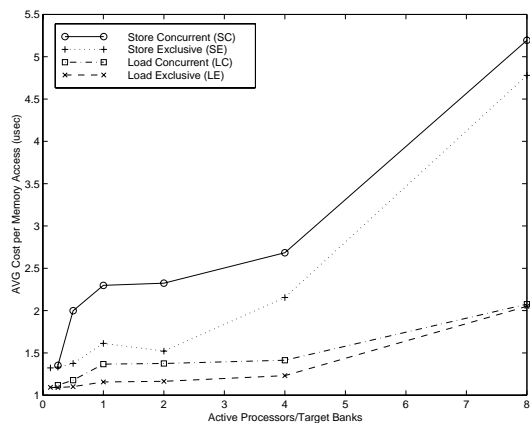
subarray sizes, namely, `size(L1)`, `size(L2)`, and  $2 \times \text{size}(L2)$ , and 0, 1, 2, 4, 8 target banks. All combinations of the above parameters have been exercised once for loads and once for stores, and by having processors operate once exclusively and once concurrently on sub-arrays. We refer to a microbenchmark using the mnemonic code `px.sty.sz.w.bz`, where  $x$  denotes the number of active processors,  $y$  the stride,  $w$  the subarray size ( $w = 1$  for `size(L1)`,  $w = 2$  for `size(L2)`, and  $w = 3$  for  $2 \times \text{size}(L2)$ ), and  $z$  the number of target banks. Moreover, a two letter suffix to the mnemonic code identifies the type of accesses (L for load or S for store) and the sharing of sub-arrays (E for exclusive or C for concurrent accesses). (See [2] for a full description of the microbenchmarks.)

A few examples should help clarify how the parameters that characterize our microbenchmarks can be set to exercise different memory usages. Consider, for instance, microbenchmark `px.st1.sz1.b0.LE`, where  $x$  active processors load consecutive words from distinct sub-arrays of size `size(L1)`. After the first iteration, all data reside in the processors' L1 caches, hence all future loads will take place in L1. Similarly, in microbenchmark `px.st8.sz2.b0.LE`, after the first iteration every load performed by an active processor will result in one L1 miss and one L2 hit. Finally, in microbenchmark `px.st256/i.sz3.bi.LE`, every load is served by main memory, and all active processors distribute their accesses evenly among  $i$  banks.

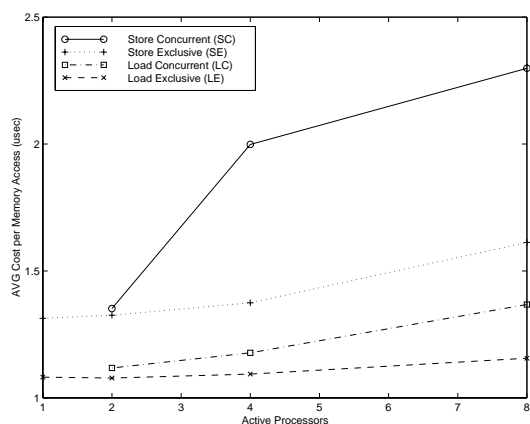
The two graphs in Figure 1 show the impact of bank and bus contention, respectively, on the access time of a small sample of microbenchmarks. Specifically, in Figure 1(a) access times are plotted, for both loads and stores and for both the exclusive and the concurrent scenarios, as a function of the ratio of active processors to target banks. Clearly, in our microbenchmarks such a ratio is directly proportional to bank contention. Figure 1(b) shows a similar plot, as a function of the number of active processors, for a subset of microbenchmarks where the number of target banks is fixed equal to the number of active processors. That is, the access times are shown as a function of bus contention for fixed bank contention. We note that for loads the impact of both bus and bank contention is somewhat minor, while for stores it is more significant, especially when combined with coherency traffic (concurrent scenario). However, in any case neither of these phenomena affects access times by more than a factor of 4, while, as we will see in the following sections, access times may vary by more than two orders of magnitude due to memory hierarchy effects.

### 4 Predicting Performance

In this section, we introduce a number of cost functions which can be used to predict the running time of a super-



(a)



(b)

**Figure 1.** Access times measured from a sample of microbenchmarks, plotted as a function of (a) ratio of active processors to target banks; and (b) number of active processors.

step on the SGI-PC. As mentioned before, we will focus on the contribution of memory accesses to the running time. In the first subsection, we present three functions inspired by the *BSP* variants mentioned in the introduction. These functions base their predictions on quantities that are computable, to a large extent, by inspecting the (assembly) code. The code allows one only to distinguish between operations on data held in registers, which we regard as local operations, and accesses to the rest of the memory system, without explicitly distinguishing between accesses to caches or main memory. Consequently, the three *BSP*-like functions considered here treat all such accesses in the same way, and cannot account for memory hierarchy effects. In the second subsection, a new function is defined which explicitly accounts for the memory hierarchy and is expressed in terms of the values of some MIPS R10000 hardware counters.

## 4.1 *BSP*-like cost functions

We define three cost functions based, respectively, on the *QSM*, *EBSP*, and  $(d, x)$ -*BSP* models. The functions, which are described below, are used to predict the running time of a superstep. In all of the functions  $H$  represents the maximum number of loads or stores issued by a processor in the superstep. In the following, we will use the term *access* to refer to either a load or a store operation.

**Function *QSM*** Function *QSM* is defined as

$$\max\{g1_{QSM} \cdot H, g2_{QSM} \cdot K\},$$

where  $K$  is the maximum number of accesses performed on the same word by all processors,  $g1_{QSM}$  is the cost per access experienced by a processor, and  $g2_{QSM}$  is the cost per access relative to a single word. According to the model, the second term is expected to dominate when there is high contention at a cell.

Since the function does not distinguish among accesses to different levels of the hierarchy, or among loads and stores, which generally have different costs, one cannot provide unique values for the parameters but, rather, intervals of possible values. On the SGI-PC, we employed the suite of microbenchmarks described in Section 3 to obtain minimum and maximum values for each parameter involved in the function by considering best case and worst case scenarios with respect to phenomena not captured by the parameter.

For  $g1_{QSM}$ , the minimum value<sup>1</sup> (0.0083) resulted from the microbenchmark performing loads from L1 (p1.st1.sz1.b0.LE), while the maximum (5.35) resulted from the microbenchmark performing stores to main memory with maximum bank and bus contention, and maximum coherency traffic (p8.st256.sz3.b1.SC). For  $g2_{QSM}$ , the minimum value (0.001) resulted from the microbenchmark with all processors repeatedly loading the same word (a variant of p8.st1.sz1.b0.LC), while the maximum (0.67) resulted from the microbenchmark with all processors repeatedly storing the same word (a variant of p8.st256.sz3.b1.SC).

Hence, we can define an “optimistic” version ( $QSM_{\min}$ ) and a “pessimistic” version ( $QSM_{\max}$ ) of the *QSM* function based, respectively, on the minimum and maximum values of  $g1_{QSM}$  and  $g2_{QSM}$ . The two versions are:

$$\begin{aligned} QSM_{\min} &= \max\{0.0083 \cdot H, 0.001 \cdot K\}, \\ QSM_{\max} &= \max\{5.35 \cdot H, 0.67 \cdot K\}. \end{aligned}$$

**Function *EBSP*** Function *EBSP* is defined as

$$\max\{g1_{EBSP} \cdot H, g2_{EBSP} \cdot M/p\},$$

<sup>1</sup>All parameter values are expressed in  $\mu\text{sec}$  per access.

where  $M$  is the total number of accesses performed by all  $p$  processors,  $g1_{EBSP}$  is the cost per access experienced by a processor when no other processor accesses memory, and  $g2_{EBSP}$  is the cost per access experienced by a processor when all other processors perform approximately the same number of accesses, thus resulting in potentially high congestion. Clearly, we expect  $g1_{EBSP} \leq g2_{EBSP}$ . According to the *EBSP* philosophy, the running time of a superstep characterized by an unbalanced access pattern, i.e., one with  $H \gg M/p$ , is determined by the time taken by the processor performing the largest number of accesses, as if it worked in isolation, hence the term  $g1_{EBSP} \cdot H$  in the function dominates. However, when  $H \simeq M/p$  the high traffic may slow the processors down, and the term  $g2_{EBSP} \cdot M/p$  dominates.

As before, we determine minimum and maximum values of the parameters. For  $g1_{EBSP}$ , the minimum value (0.0083) resulted from loads from L1 (p1.st1.sz1.b0.LE), while the maximum (1.31) resulted from the microbenchmark with one processor accessing main memory with maximum bank contention (p1.st256.sz3.b1.SE). For  $g2_{EBSP}$ , the minimum value (0.0083) resulted from the microbenchmark with all processors accessing their L1 caches (p8.st1.sz1.b0.LE), while the maximum (5.35) resulted from the microbenchmark with all processors performing the same number of accesses to main memory with maximum bank and bus contention and maximum coherency traffic (p8.st256.sz3.b1.SC).

The corresponding “optimistic” and “pessimistic” versions of function *EBSP* are:

$$\begin{aligned} EBSP_{\min} &= \max\{0.0083 \cdot H, 0.0083 \cdot M/p\}, \\ EBSP_{\max} &= \max\{1.31 \cdot H, 5.35 \cdot M/p\}. \end{aligned}$$

**Function *DXBSP*** Function *DXBSP* is defined as

$$\max\{g1_{DXBSP} \cdot H, g2_{DXBSP} \cdot M_b\},$$

where  $M_b$  is the total number of accesses directed to the same bank,  $g1_{DXBSP}$  is the cost per access experienced by a processor when bank contention is low, and  $g2_{DXBSP}$  is the cost per access at a bank experienced in the presence of high bank contention. According to the  $(d, x)$ -*BSP* model, the second term is expected to dominate only when many accesses hit the same bank, which becomes a bottleneck.

For  $g1_{DXBSP}$ , the minimum value (0.0083) resulted from one processor performing loads from L1 (p1.st1.sz1.b0.LE), while the maximum (3.40) resulted from the microbenchmark where all processors access main memory with moderate bank contention but high bus contention and coherency traffic (p8.st32.sz3.b8.SC). For  $g2_{DXBSP}$ , the minimum value (0.26) and the maximum value (0.33) resulted from microbenchmarks with all processors performing, respectively, loads and stores

on distinct words in the same bank (p8.st256.sz3.b1.LE and p8.st256.sz3.b1.SE). Here, we chose microbenchmarks without concurrent accesses to make sure that the running time was indeed dominated by bank delay and not by other factors (e.g., coherency activities).

The corresponding “optimistic” and “pessimistic” versions of function *DXBSP* are:

$$\begin{aligned} DXBSP_{\min} &= \max\{0.0083 \cdot H, 0.26 \cdot M_b\}, \\ DXBSP_{\max} &= \max\{3.40 \cdot H, 0.33 \cdot M_b\}. \end{aligned}$$

## 4.2 A cost function based on hardware counters

As will be demonstrated by the validations reported in the next section, the accuracy of all of the above *BSP*-like functions is strongly limited by the fact that they disregard the memory hierarchy, which is the main reason for the high variance in the parameter values, and, consequently, for the large gap between the optimistic and pessimistic versions of these functions.

We define a new cost function  $F$  that is based on the MIPS R10000 hardware counters shown in Table 1. The counters provide a detailed account of the memory hierarchy usage. Function  $F$  is defined under the assumption that the running time of a superstep is determined by one of the following factors: (1) the accesses issued by some processor at the various levels of the hierarchy, (2) the traffic on the bus caused by accesses to main memory, or (3) bank contention caused by accesses targeting the same bank. To reflect this assumption,  $F$  takes the maximum of three functions  $F1$ ,  $F2$  and  $F3$  defined as follows. Let  $p_i$  and  $b_j$  denote, respectively, the  $i$ -th processor and the  $j$ -th memory bank, with  $0 \leq i < p$  and  $0 \leq j < q$ , where  $q$  denotes the number of memory banks. Then:

$$\begin{aligned} F1 &= \max_{0 \leq i < p} (g1_{F1} \cdot (\text{LD}(p_i) + \text{ST}(p_i)) + \\ &\quad + g2_{F1} \cdot \text{L1M}(p_i) + g3_{F1} \cdot \text{L2M}(p_i) + \\ &\quad + g4_{F1} \cdot \text{L1W}(p_i) + g5_{F1} \cdot \text{L2W}(p_i)) \\ F2 &= g1_{F2} \cdot \sum_{0 \leq i < p} \text{L2M}(p_i) + g2_{F2} \cdot \sum_{0 \leq i < p} \text{L2W}(p_i) \\ F3 &= \max_{0 \leq j < q} (g1_{F3} \cdot \text{LM}(b_j) + g2_{F3} \cdot \text{WB}(b_j)), \end{aligned}$$

where counters  $\text{LD}(p_i)$ ,  $\text{ST}(p_i)$ ,  $\text{L1M}(p_i)$ ,  $\text{L2M}(p_i)$  and  $\text{L2W}(p_i)$  denote the values of the respective counters of processor  $p_i$ , while  $\text{LM}(b_j)$  and  $\text{WB}(b_j)$  denote, respectively, the total number of L2 misses served by bank  $b_j$  and the quadwords written back from L2 to bank  $b_j$ , and they must be inferred from counters L2M and L2W and from code inspection.

$F1$  accounts for the usage of the memory hierarchy by individual processors. Specifically,  $g1_{F1}$  reflects the cost of accessing L1 and is multiplied by the total number of loads

MIPS R10000 Hardware Counters	
LD (C0E2)	Loads issued
ST (C0E3)	Stores issued
L1M (C1E9)	L1 misses
L2M (C1EA)	L2 misses
L1W (C1E6)	L1 lines written back from L1 to L2
L2W (C0E7)	Quadword writebacks from L2 to RAM

**Table 1.** MIPS R10000 hardware counters used for  $F$

$F_i$	$g1_{F_i}$	$g2_{F_i}$	$g3_{F_i}$	$g4_{F_i}$	$g5_{F_i}$
$F1$	0.0088	0.055	0.97	0.013	0.026
$F2$	0.15	0.02			
$F3$	0.26	0.051			

**Table 2.** Parameters measured for  $F1$ ,  $F2$  and  $F3$ .

and stores made by a processor, since all of them act eventually on L1;  $g2_{F1}$  and  $g4_{F1}$  account for data movement costs (respectively misses and writebacks) between L1 and L2; analogously,  $g3_{F1}$  and  $g5_{F1}$  account for the cost of data movements between L2 and main memory.  $F2$  accounts for bus contention and should dominate when most of the processors are active and each active processor issues many requests for data in main memory. Specifically,  $g1_{F2}$  and  $g2_{F2}$  reflect the delay introduced by the bus for an L2 miss and quadword writeback, respectively. Finally,  $F3$  accounts for bank contention and is expected to dominate when many active processors issue many access requests served by the same bank. Parameters  $g1_{F3}$  and  $g2_{F3}$  reflect the time taken by a bank to serve an L2 miss and quadword writeback, respectively.

Table 2 shows the values of the parameters obtained by fitting each function (through least squares fitting) on a set of microbenchmarks where that function is expected to dominate. We used microbenchmarks p1.st1.sz1.b0.\*E, p1.st8.sz2.b0.\*E, and p1.st32.sz3.b8.\*E for  $F1$ , microbenchmarks p8.st256.sz3.b8.\*E for  $F2$ , and microbenchmarks p8.st256.sz3.b1.\*E for  $F3$ , where \* stands for both L and S.

Although  $F$  cannot be immediately computed by simply inspecting the code, it may still be used to predict the performance of an application in situations where accurate estimates of the relevant counters can be inferred *a priori* or extrapolated from pilot runs on small input sizes. An example of such a use is given in Section 5.4. Moreover, the function could prove useful in the design of software systems, compilers, or large applications, to profile the memory hierarchy usage of critical portions of their code.

## 5 Validations

In this section we investigate the quality of predictions for the cost functions introduced in Section 4 by checking their accuracy over a set of synthetic access patterns and over a number of real applications, namely, three bulk-synchronous implementations of parallel sorting and the NAS Parallel Benchmarks [6, 7]. Specifically, we determined measured and predicted times (indicated by  $T_M$  and  $T_P$ , respectively) and calculated the prediction error as

$$ERR = \frac{\max\{T_M, T_P\}}{\min\{T_M, T_P\}},$$

which indicates how much smaller or larger the predicted time is with respect to the measured time.

### 5.1 Synthetic Access Patterns

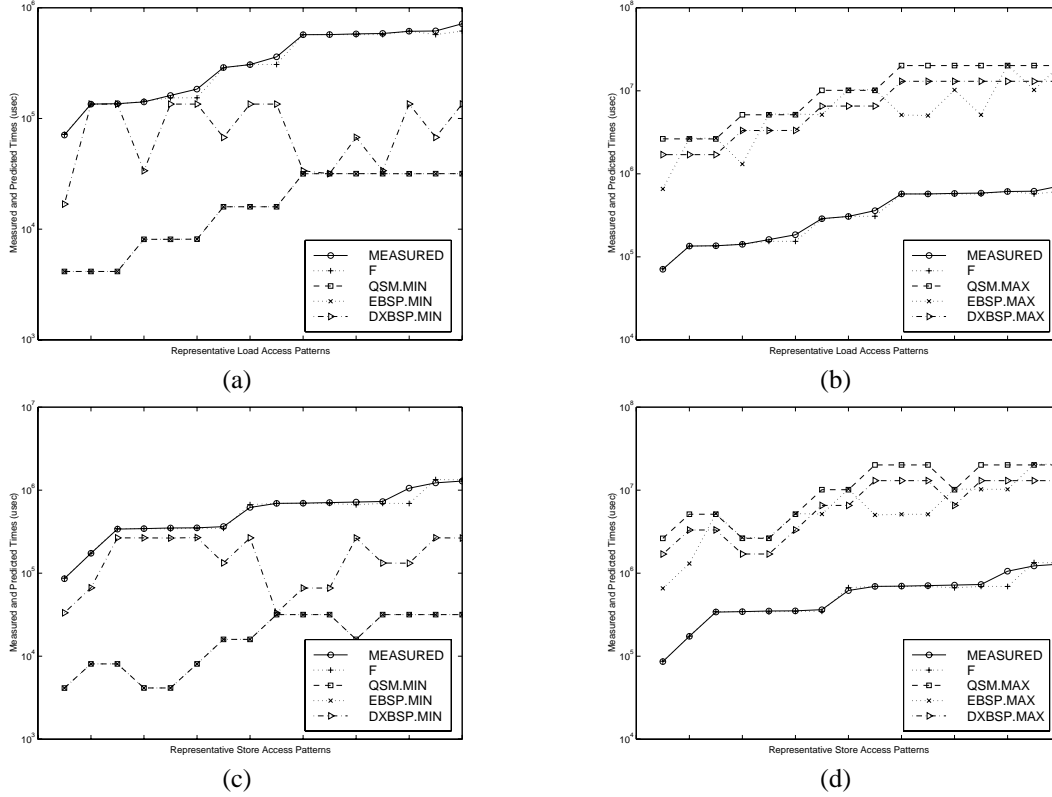
Synthetic *access patterns* were obtained by running the original set of microbenchmarks under a variety of scenarios featuring different phenomena that could have an impact on access time. More specifically, along with the already discussed Load/Store–Exclusive/Concurrent combinations, we introduced other variants where processors access data which are present as *clean*, *shared* or *dirty* in other processors’ caches, so as to include patterns exercising different aspects of the coherency protocol. Overall, we obtained a *Validation Suite* (VS, for short) of 412 different access patterns.

Table 3 reports, for all functions, the average and maximum values of ERR obtained over the entire suite VS, while Figure 2 shows plots of measured and predicted running times of a subset of VS derived from the microbenchmarks employed for Figure 1.(a).

The results of the validations highlight the fact that disregarding hierarchy effects when evaluating performance has a significant negative impact on predictive accuracy, while modeling other architectural aspects, such as bus and bank contention, or concurrency to the same cell, yields a rather modest payoff in achieving higher accuracy. Moreover, the

Function		AVG ERR	MAX ERR
QSM	MIN	24.15	88.02
	MAX	53.85	636.79
EBSP	MIN	24.15	88.02
	MAX	27.29	648.35
DXBSP	MIN	6.36	31.84
	MAX	34.8	411.46
$F$		1.19	1.91

**Table 3.** Cost function errors for VS.



**Figure 2.** Measured and predicted times for a set of 16 access patterns consisting of loads (Plots (a) and (b)) and stores (Plots (c) and (d)), derived from the microbenchmarks employed in Figure 1.(a). The measured time and  $F$  are shown on all plots. Plots (a) and (c) show the optimistic versions of the  $BSP$ -like functions, while Plots (b) and (d) their pessimistic versions.

fact that function  $F$  exhibits high accuracy on all the experiments in VS suggests that phenomena that were disregarded when defining the function (such as some types of coherency overheads) have only a minor impact on performance. Finally, we note that although quantifying bank contention, as required to apply functions  $F3$  (hence  $F$ ) and  $DXBSP$ , may be hard to do in practical situations, a number of measurements, not reported in this extended abstract, show that  $F$  remains quite accurate even if only the maximum between  $F1$  and  $F2$  is considered.

## 5.2 Sorting Programs

Our first set of applications consists of three sorting algorithms: *samplesort* [15], *columnsort* [11], and a parallel version of *radixsort* [3]. These algorithms were chosen because they are well-understood parallel algorithms that exhibit a variety of communication patterns. We coded all algorithms in a bulk-synchronous fashion, with no special effort made to optimize the implementations, since our goal was to test the accuracy of the cost functions rather than to develop efficient algorithms.

Each algorithm was run on a wide range of input sizes;

namely,  $n/p = 10^5 \cdot i$ , for  $1 \leq i \leq 10$ , where  $n$  is the total number of keys to be sorted. For each run, we applied the functions to every superstep. For the most part, applying the functions was straightforward. The only difficulty was posed by functions  $DXBSP$  and  $F3$ , which require knowledge of the maximum number of L2 misses and writebacks targeting a particular memory bank. While this information is known for the synthetic access patterns of VS, it cannot be readily obtained for real applications. However, for the purposes of our validations, we have used estimates which assume that all such accesses are equally distributed amongst the eight banks. This is a reasonable assumption for the sorting programs, whose access patterns tend to be balanced amongst the 8 banks.

A summary of our results is contained in Table 4 where we report, for each superstep, the maximum value of ERR over all runs (the average values of ERR are similar and can be found in [2]). Since the sorting algorithms exhibit a high degree of locality, we would expect the optimistic versions of the  $BSP$ -like functions to perform much better than their pessimistic counterparts, and indeed this is the case (errors are not shown for  $EBSP_{\min}$  and  $DXBSP_{\min}$  because they

Sorting Programs – Maximum Prediction Errors						
Sort	SSStep	QSM		EBSP	DXBSP	F
		MIN	MAX	MAX	MAX	
Radix	SS1	2.12	400.14	320.48	258.55	1.41
	SS2	3.35	505.89	298.79	326.88	1.86
	SS3	3.08	473.30	295.88	305.83	1.83
	SS4	2.72	321.56	302.11	207.78	1.39
Sample	SS1	2.62	339.75	196.86	219.53	1.44
	SS2	2.17	320.95	252.25	207.38	1.15
	SS3	1.98	404.26	195.76	261.21	1.30
	SS4	2.89	287.72	247.31	185.91	1.11
	SS5	2.58	361.36	327.08	233.49	1.26
Column	SS1	3.44	268.13	205.23	173.25	1.06
	SS2	2.46	268.13	264.49	173.25	2.05
	SS3	2.88	230.37	228.11	148.85	1.88
	SS4	2.61	245.56	247.10	158.67	2.09
	SS5	1.36	484.93	280.03	313.34	1.16

**Table 4.** Sorting algorithms: accuracy of the cost functions on individual supersteps.

are almost identical to the errors for  $QSM_{\min}$ ). Although the difference is not as dramatic as for the synthetic applications,  $F$  is still seen to be significantly more accurate than any of the  $BSP$ -like functions. This indicates that disregarding hierarchy effects results in a noticeable lack of accuracy even for regular programs.

### 5.3 NAS Parallel Benchmarks

We tested our functions on the *NAS Parallel Benchmarks (NPB)* [6, 7]. NPB is a set of 8 programs, derived from computational fluid dynamics (CFD) applications, that is designed to evaluate the performance of parallel machines. The **CG** kernel solves an unstructured sparse linear system by the conjugate gradient method. **BT** is an embarrassingly parallel kernel that generates pairs of Gaussian random deviates and tabulates the number of pairs in successive square annuli. **FT** is a fairly standard implementation of a 3D FFT PDE. **IS** is an integer sorting program; the keys are generated in an initial sequential portion. The **LU** solver application is a diagonal pipelining computation that results in a large number of small messages. **MG** is a simple 3D multi-grid benchmark. The **SP** and **BT** applications both solve three sets of uncoupled systems of equations using a multi-partition scheme [5] which provides good load balance and uses coarse-grained communication.

The benchmarks are MPI-based source-code implementations that are intended to run ‘as is’. As we did not wish to alter the benchmarks, we decided to treat each program as a single ‘superstep’ for the purposes of evaluating the cost functions, that is, measurements (times and counters) were performed external to the benchmark execution using the `SGI Perfex` utility. This seemed a reasonable compromise as the NAS benchmarks have only a few barrier synchronizations, which are very fast on the SGI-PC, and, moreover, with the exception of **LU** exchange small numbers of

NAS Parallel Benchmarks (NPB) – Prediction Errors						
NPB	MIN	QSM		EBSP	DXBSP	F
		MAX	MAX	MAX	MAX	
<b>CG</b>	2.46	258.31	210.10	166.91	1.46	
<b>EP</b>	2.42	262.53	252.32	169.64	1.02	
<b>FT</b>	2.05	309.40	245.64	199.92	1.63	
<b>IS</b>	1.57	404.81	354.47	261.57	1.39	
<b>LU</b>	2.15	295.01	236.80	190.62	1.32	
<b>MG</b>	1.57	403.48	289.11	260.71	1.73	
<b>BT</b>	2.77	229.68	189.08	148.41	1.05	
<b>SP</b>	2.13	298.69	194.21	193.00	1.05	

**Table 5.** NAS Parallel Benchmarks: accuracy of the cost functions on the individual benchmarks.

large messages, hence the overhead introduced by the MPI message-passing routines is rather limited. Finally, two of the benchmarks, **BT** and **SP** required a square number of processors to run. In these cases, we used a nine-processor configuration of the machine but applied the functions derived for the eight-processor configuration, assuming that parameter values would not change significantly. Our assumption is confirmed by the very low errors obtained by  $F$  on these latter benchmarks.

Table 5 reports the prediction errors (ERR) incurred by the functions on each NPB (errors are not shown for  $EBSP_{\min}$  and  $DXBSP_{\min}$  because they were practically identical to the errors for  $QSM_{\min}$ ). Here, the distinction between average and maximum errors does not make sense, since only one input size was used for each benchmark. (For the  $DXBSP$  and the  $F3$  functions we made the same approximation for bank contention as was done for the sorting applications.) Again, it can be seen that the  $F$  function performs significantly better than the other functions, with almost perfect predictions for the **EP**, **BT**, and **SP** benchmarks, and discrepancies of less than 75% in all cases. As with the sorting programs, the optimistic versions of the  $BSP$ -like functions perform much better than their pessimistic counterparts, which can be attributed to the high locality and regularity exhibited by the benchmarks. However, even the former are significantly worse than the  $F$  function, in almost all cases, with errors up to 180%. Again, the reasonable level of accuracy attained by the optimistic versions of the  $BSP$ -like cost functions is to be attributed to the high locality and regularity exhibited by the benchmarks.

### 5.4 Extrapolating Performance

One of the advantages of the  $BSP$ -like functions over the counter-based function  $F$ , is that, to a large extent, the programmer can easily determine the input values for the function (e.g.,  $H$  or  $M$ ). However, as we have seen, these functions may not provide meaningful predictions as they all fail to account for hierarchy effects.



Sort	SSStep	Errors for $F$ Measured Counts		Errors for $F$ Estimated Counts	
		AVG	MAX	AVG	MAX
Radix	SS1	1.22	1.32	1.01	1.09
	SS4	1.11	1.16	1.13	1.16
Sample	SS2	1.05	1.09	1.20	1.21
	SS4	1.06	1.11	1.03	1.04
	SS5	1.13	1.17	1.24	1.26
Column	SS1	1.12	2.49	1.17	1.72
	SS2	1.80	1.89	2.02	2.12
	SS3	1.66	1.69	1.84	1.87
	SS4	1.78	1.83	2.05	2.06
	SS5	1.16	1.17	1.88	1.90

**Table 6.** Sorting algorithms: comparison of  $F$ 's accuracy with measured vs estimated counters, over selected sorting supersteps and large input sizes.

While the counter-based function exhibits excellent accuracy, it seems that one should actually run the program to obtain the required counts, which would annihilate its potential as a performance predictor. However, there are many cases where counter values can be guessed in advance with reasonable confidence, and then plugged into  $F$  to obtain accurate predictions. In fact, we claim that in some cases meaningful estimates for the counters can be derived by extrapolating values for large problem sizes from pilot runs of the program on small input sets.

To investigate the above claim, we developed least-squares fits for each of the counters used in  $F$  for those supersteps in our three sorting algorithms that had significant communication. The input size  $n$  of the sorting instance was used as the independent variable. For each counter, we obtained the fits on small input sizes ( $n/p = 10^5 \cdot i$ , for  $1 \leq i \leq 5$ ), and then used the fits to forecast the counter values for large input sizes ( $n/p = 10^5 \cdot i$ , for  $5 < i \leq 10$ ). These estimated counter values were then plugged in  $F$  to predict the execution times for the larger runs.

The results of this study are summarized in Table 6. Interestingly, but somewhat accidentally, in some cases the predictions obtained with the estimated counter values were actually slightly better than those obtained with the measured counter values (e.g., Column sort's superstep 1, Radix sort's superstep 1, and Sample sort's superstep 4). More importantly, however, it can be seen that in *all* cases, the level of accuracy of  $F$  using the extrapolated counter values was not significantly worse than what was obtained with the actual counter values. These preliminary results indicate that at least in some situations a hardware counter-based function does indeed have potential as an *a priori* predictor of performance.

## References

- [1] B. Alpern, L. Carter, and E. Feig. Uniform memory hierarchies. In *Proc. of the 31st IEEE FOCS*, pages 600–608, 1990.
- [2] N. M. Amato, J. Perdue, A. Pietracaprina, G. Pucci, and M. Mathis. Predicting performance on SMPs. a case study: The SGI Power Challenge. Technical Report 99-020, Dept. of Computer Science, Texas A&M University, 1999.
- [3] G. E. Blelloch, C. E. Leiserson, B. M. Maggs, C. G. Plaxton, S. J. Smith, and M. Zagha. An experimental analysis of parallel sorting algorithms. *Theory of Computing Systems*, 31(2):135–167, 1998.
- [4] G.E. Blelloch, P.B. Gibbons, Y. Mattias, and M. Zagha. Accounting for memory bank contention and delay in high-bandwidth multiprocessors. *IEEE Trans. Par. Dist. Sys.*, 8(9):943–958, 1997.
- [5] J. Bruno and P. R. Cappello. Implementing the beam and warming method on the hypercube. In *Proceedings of 3rd Conference on Hypercube Concurrent Computers and Applications*, Pasadena, CA, jan 1988.
- [6] D. Bailey *et al.* The NAS parallel benchmarks. *Int. J. Supercomputer Appl.*, 5(3):63–73, 1991.
- [7] D. Bailey *et al.* The NAS parallel benchmarks 2.0. Technical Report NAS-95-020, NASA Ames Research Center, <http://www.nas.nasa.gov/npb/>, 1995.
- [8] P. B. Gibbons, Y. Matias, and V. Ramachandran. Can a shared-memory model serve as a bridging-model for parallel computation? In *Proc. ACM Symp. Par. Alg. Arch. (SPAA)*, pages 72–83, 1997.
- [9] M. W. Goudreau, K. Lang, S. B. Rao, T. Suel, and T. Tsantilas. Portable and efficient parallel computing using the BSP model. *IEEE Trans. Comput.*, 48(7):670–689, 1999.
- [10] B. H. H. Juurlink and H. A. G. Wijshoff. A quantitative comparison of parallel computation models. In *Proc. ACM Symp. Par. Alg. Arch. (SPAA)*, pages 13–24, 1996.
- [11] T. Leighton. Tight bounds on the complexity of parallel sorting. *IEEE Trans. Comput.*, c-34(4):344–354, 1985.
- [12] Silicon Graphics Corporation 1995. *SGI Power Challenge: User's Guide*, 1995.
- [13] Silicon Graphics Corporation 1997. *Definition of MIPS R10000 Performance Counters*, 1997. <http://www.sgi.com/processors/r10k/performance.html>.
- [14] L. Valiant. Bridging model for parallel computation. *Comm. ACM*, 33(8):103–111, 1990.
- [15] Y. Won and S. Sahni. A balanced bin sort for hypercube multiprocessors. *J. Supercomputing*, 2:435–448, 1988.
- [16] M. Zagha, B. Larson, S. Turner, and M. Itzkowitz. Performance analysis using the MIPS R10000 performance counters. In *Proc. Supercomputing*, 1996.