

A Randomized Parallel Sorting Algorithm With an Experimental Study

David R. Helman David A. Bader* Joseph JáJá†

Institute for Advanced Computer Studies &
Department of Electrical Engineering,
University of Maryland, College Park, MD 20742
{helman, dbader, joseph}@umiacs.umd.edu

August 15, 1996

Abstract

Previous schemes for sorting on general-purpose parallel machines have had to choose between poor load balancing and irregular communication or multiple rounds of all-to-all personalized communication. In this paper, we introduce a novel variation on sample sort which uses only two rounds of regular all-to-all personalized communication in a scheme that yields very good load balancing with virtually no overhead. Moreover, unlike previous variations, our algorithm efficiently handles the presence of duplicate values without the overhead of tagging each element with a unique identifier. This algorithm was implemented in *SPLIT-C* and run on a variety of platforms, including the Thinking Machines CM-5, the IBM SP-2, and the Cray Research T3D. We ran our code using widely different benchmarks to examine the dependence of our algorithm on the input distribution. Our experimental results illustrate the efficiency and scalability of our algorithm across different platforms. In fact, it seems to outperform all similar algorithms known to the authors on these platforms, and its performance is invariant over the set of input distributions unlike previous efficient algorithms. Our results also compare favorably with those reported for the simpler ranking problem posed by the NAS Integer Sorting (IS) Benchmark.

Keywords: Parallel Algorithms, Generalized Sorting, Integer Sorting, Sample Sort, Parallel Performance.

*The support by NASA Graduate Student Researcher Fellowship No. NGT-50951 is gratefully acknowledged.

†Supported in part by NSF grant No. CCR-9103135 and NSF HPCC/GCAG grant No. BIR-9318183.

1 Introduction

Sorting is arguably the most studied problem in computer science, both because of its intrinsic theoretical importance and its use in so many applications. Its significant requirements for interprocessor communication bandwidth and the irregular communication patterns that are typically generated have earned its inclusion in several parallel benchmarks such as NAS [7] and SPLASH [34]. Moreover, its practical importance has motivated the publication of a number of empirical studies seeking to identify the most efficient sorting routines. Yet, parallel sorting strategies have still generally fallen into one of two groups, each with its respective disadvantages. The first group, using the classification of Li and Sevcik [23], is the single-step algorithms, so named because data is moved exactly once between processors. Examples of this include sample sort [21, 10], parallel sorting by regular sampling [31, 24], and parallel sorting by overpartitioning [23]. The price paid by these single-step algorithms is an irregular communication scheme and difficulty with load balancing. The other group of sorting algorithms is the multi-step algorithms, which include bitonic sort [9], column sort [22], rotate sort [25], hyperquicksort [28], flashsort [29], B-flashsort [20], smoothsort [27], and Tridgell and Brent's sort [32]. Generally speaking, these algorithms accept multiple rounds of communication in return for better load balancing and, in some cases, regular communication.

In this paper, we present a novel variation on the sample sort algorithm [19] which addresses the limitations of previous implementations. We exchange the single step of irregular communication for two steps of regular communication. In return, we reduce the problem of poor load balancing because we are able to sustain a very high oversampling ratio at virtually no cost. Second, we efficiently accommodate the presence of duplicates without the overhead of tagging each element. And we obtain predictable, regular communication requirements which are essentially invariant with respect to the input distribution. Utilizing regular communication has become more important with the advent of message passing standards, such as MPI [26], which seek to guarantee the availability of very efficient (often machine specific) implementations of certain basic collective communication routines.

Our algorithm was implemented in a high-level language and run on a variety of platforms, including the Thinking Machines CM-5, the IBM SP-2, and the Cray Research T3D. We ran our code using a variety of benchmarks that we identified to examine the dependence of our algorithm on the input distribution. Our experimental results are consistent with the theoretical analysis and illustrate the scalability and efficiency of our algorithm across different platforms. In fact, it seems to outperform all similar algorithms known to the authors on these platforms, and its performance is indifferent to the set of input distributions unlike previous efficient algorithms.

The high-level language used in our studies is SPLIT-C [14], an extension of *C* for distributed memory machines. The algorithm makes use of MPI-like communication primitives but does not make any assumptions as to how these primitives are actually implemented. The basic data transport

is a **read** or **write** operation. The remote read and write typically have both blocking and non-blocking versions. Also, when reading or writing more than a single element, bulk data transports are provided with corresponding **bulk_read** and **bulk_write** primitives. Our collective communication primitives, described in detail in [6], are similar to those of the MPI [26], the IBM POWERparallel [8], and the Cray MPP systems [13] and, for example, include the following: **transpose**, **bcast**, **gather**, and **scatter**. Brief descriptions of these are as follows. The **transpose** primitive is an all-to-all personalized communication in which each processor has to send a unique block of data to every processor, and all the blocks are of the same size. The **bcast** primitive is used to copy a block of data from a single source to all the other processors. The primitives **gather** and **scatter** are companion primitives. **Scatter** divides a single array residing on a processor into equal-sized blocks, each of which is distributed to a unique processor, and **gather** coalesces these blocks back into a single array at a particular processor. See [3, 6, 4, 5] for algorithmic details, performance analyses, and empirical results for these communication primitives.

The organization of this paper is as follows. **Section 2** presents our computation model for analyzing parallel algorithms. **Section 3** describes in detail our improved sample sort algorithm. Finally, **Section 4** describes our data sets and the experimental performance of our sorting algorithm.

2 The Parallel Computation Model

We use a simple model to analyze the performance of our parallel algorithms. Our model is based on the fact that current hardware platforms can be viewed as a collection of powerful processors connected by a communication network that can be modeled as a complete graph on which communication is subject to the restrictions imposed by the latency and the bandwidth properties of the network. We view a parallel algorithm as a sequence of local computations interleaved with communication steps, where we allow computation and communication to overlap. We account for communication costs as follows.

Assuming no congestion, the transfer of a block consisting of m contiguous words between two processors takes $(\tau + \sigma m)$ time, where τ is the latency of the network and σ is the time per word at which a processor can inject or receive data from the network. Note that the bandwidth per processor is inversely proportional to σ . We assume that the bisection bandwidth is sufficiently high to support block permutation routing amongst the p processors at the rate of $\frac{1}{\sigma}$. In particular, for any subset of q processors, a block permutation amongst the q processors takes $(\tau + \sigma m)$ time, where m is the size of the largest block.

Using this cost model, we can evaluate the communication time $T_{comm}(n, p)$ of an algorithm as a function of the input size n , the number of processors p , and the parameters τ and σ . The coefficient of τ gives the total number of times collective communication primitives are used, and the coefficient

of σ gives the maximum total amount of data exchanged between a processor and the remaining processors.

This communication model is close to a number of similar models (e.g. [16, 33, 1]) that have recently appeared in the literature and seems to be well-suited for designing parallel algorithms on current high performance platforms.

We define the computation time T_{comp} as the maximum time it takes a processor to perform all the local computation steps. In general, the overall performance $T_{comp} + T_{comm}$ involves a tradeoff between T_{comp} and T_{comm} . In many cases, it is possible to minimize both T_{comp} and T_{comm} simultaneously, and sorting is such a case.

3 A New Sample Sort Algorithm

Consider the problem of sorting n elements equally distributed amongst p processors, where we assume without loss of generality that p divides n evenly. The idea behind sample sort is to find a set of $p - 1$ *splitters* to partition the n input elements into p groups indexed from 1 up to p such that every element in the i^{th} group is less than or equal to each of the elements in the $(i + 1)^{th}$ group, for $1 \leq i \leq p - 1$. Then the task of sorting each of the p groups can be turned over to the correspondingly indexed processor, after which the n elements will be arranged in sorted order. The efficiency of this algorithm obviously depends on how evenly we divide the input, and this in turn depends on how well we choose the *splitters*. One way to choose the *splitters* is by randomly sampling the input elements at each processor - hence the name **sample sort**.

Previous versions of sample sort [21, 10, 17, 15] have randomly chosen s *samples* from the $\frac{n}{p}$ elements at each processor, routed these ps *samples* to a single processor, sorted them at that processor, and then selected every s^{th} element as a *splitter*. Each processor P_i then performs a binary search on these *splitters* for each of its input values and then uses the results to route the values to the appropriate destination, after which local sorting is done to complete the sorting process. The first difficulty with this approach is the work involved in gathering and sorting the *samples*. A larger value of s results in better load balancing, but it also increases the overhead. The second difficulty is that no matter how the routing is scheduled, there exist inputs that give rise to large variations in the number of elements destined for different processors, and this in turn results in an inefficient use of the communication bandwidth. Moreover, such an irregular communication scheme cannot take advantage of the regular communication primitives proposed under the MPI standard [26]. The final difficulty with the original approach is that duplicate values are accommodated by tagging each item with a unique value [10]. This, of course, doubles the cost of both memory access and interprocessor communication.

In our solution, we incur no overhead in obtaining $\frac{n}{p^2}$ *samples* from each processor and in sorting these *samples* to identify the *splitters*. Because of this very high oversampling, we are able to replace

the irregular routing with exactly two calls to our **transpose** primitive, and, in addition, we are able to efficiently accommodate the presence of duplicates without resorting to tagging.

The pseudocode for our algorithm is as follows:

- **Step (1):** Each processor P_i ($1 \leq i \leq p$) randomly assigns each of its $\frac{n}{p}$ elements to one of p buckets. With high probability, no bucket will receive more than $c_1 \frac{n}{p^2}$ elements, where c_1 is a constant to be defined later.
- **Step (2):** Each processor P_i routes the contents of bucket j to processor P_j , for $(1 \leq i, j \leq p)$. Since with high probability no bucket will receive more than $c_1 \frac{n}{p^2}$ elements, this is equivalent to performing a **transpose** operation with block size $c_1 \frac{n}{p^2}$.
- **Step (3):** Each processor P_i sorts at most $\left(\alpha_1 \frac{n}{p} \leq c_1 \frac{n}{p}\right)$ values received in **Step (2)** using an appropriate sequential sorting algorithm. For integers we use the radix sort algorithm, whereas for floating point numbers we use the merge sort algorithm.
- **Step (4):** From its sorted list of $\left(\beta \frac{n}{p} \leq c_1 \frac{n}{p}\right)$ elements, processor P_1 selects each $\left(j\beta \frac{n}{p^2}\right)^{th}$ element as $\text{Splitter}[j]$, for $(1 \leq j \leq p - 1)$. By default, $\text{Splitter}[p]$ is the largest value allowed by the data type used. Additionally, for each $\text{Splitter}[j]$, binary search is used to determine the values $\text{Frac}_L[j]$ and $\text{Frac}_R[j]$, which are respectively the fractions of the total number of elements at processor P_1 with the same value as $\text{Splitter}[j - 1]$ and $\text{Splitter}[j]$ which also lie between index $\left((j - 1)\beta \frac{n}{p^2} + 1\right)$ and index $\left(j\beta \frac{n}{p^2}\right)$, inclusively.
- **Step (5):** Processor P_1 **broadcasts** the Splitter , Frac_L , and Frac_R arrays to the other $p - 1$ processors.
- **Step (6):** Each processor P_i uses binary search on its sorted local array to define for each of the p *splitters* a subsequence S_j . The subsequence associated with $\text{Splitter}[j]$ contains all those values which are greater than $\text{Splitter}[j - 1]$ and less than $\text{Splitter}[j]$, as well as $\text{Frac}_L[j]$ and $\text{Frac}_R[j]$ of the total number of elements in the local array with the same value as $\text{Splitter}[j - 1]$ and $\text{Splitter}[j]$, respectively.
- **Step (7):** Each processor P_i routes the subsequence associated with $\text{Splitter}[j]$ to processor P_j , for $(1 \leq i, j \leq p)$. Since with high probability no sequence will contain more than $c_2 \frac{n}{p^2}$ elements, where c_2 is a constant to be defined later, this is equivalent to performing a **transpose** operation with block size $c_2 \frac{n}{p^2}$.
- **Step (8):** Each processor P_i merges the p sorted subsequences received in **Step (7)** to produce the i^{th} column of the sorted array. Note that, with high probability, no processor has received more than $\alpha_2 \frac{n}{p}$ elements, where α_2 is a constant to be defined later.

We can establish the complexity of this algorithm with high probability - that is with probability

$\geq (1 - n^{-\epsilon})$ for some positive constant ϵ . But before doing this, we need to establish the results of the following lemmas.

Lemma 1: At the completion of **Step (1)**, the number of elements in each bucket is at most $c_1 \frac{n}{p^2}$ with high probability, for any $c_1 \geq 2$ and $p^2 \leq \frac{n}{3 \ln n}$.

Proof: The probability that exactly $c_1 \frac{n}{p^2}$ elements are placed in a particular bucket in **Step (1)** is given by the binomial distribution

$$b(s; r, q) = \binom{r}{s} q^s (1 - q)^{r-s}, \quad (1)$$

where $s = c_1 \frac{n}{p^2}$, $r = \frac{n}{p}$, and $q = \frac{1}{p}$. Using the following Chernoff bound [12] for estimating the tail of a binomial distribution

$$\sum_{s \geq (1+\epsilon)r q} b(s; r, q) \leq e^{-\frac{\epsilon^2 r q}{3}}, \quad (2)$$

the probability that a particular bucket will contain at least $c_1 \frac{n}{p^2}$ elements can be bounded by

$$e^{-(c_1-1)^2 \frac{n}{3p^2}}. \quad (3)$$

Hence, the probability that any of the p^2 buckets contains at least $c_1 \frac{n}{p^2}$ elements can be bounded by

$$p^2 e^{-(c_1-1)^2 \frac{n}{3p^2}} \quad (4)$$

and **Lemma 1** follows.

Lemma 2: At the completion of **Step (2)**, the total number of elements received by processor P_1 , which comprise the set of *samples* from which the *splitters* are chosen, is at most $\beta \frac{n}{p}$ with high probability, for any $\beta \geq 1$ and $p^2 \leq \frac{n}{3 \ln n}$.

Proof: The probability that processor P_1 receives exactly $\beta \frac{n}{p}$ elements is given by the binomial distribution $b\left(\beta \frac{n}{p}; n, \frac{1}{p}\right)$. Using the Chernoff bound for estimating the tail of a binomial distribution, the probability that processor P_1 receives at least $\beta \frac{n}{p}$ elements can be bounded by $e^{-(\beta-1)^2 \frac{n}{3p}}$ and **Lemma 2** follows.

Lemma 3: For each Splitter $[j]$, where $(1 \leq j \leq p)$, let SE_j and SS_j be respectively the sets of input elements and samples that are both equal in value to Splitter $[j]$, and let $|SS_j| \leq \lambda_j \frac{n}{p^2}$. Then, with high probability, no SE_j will contain more than $M_j \frac{n}{p}$ elements, where

$$M_j = \frac{(6\lambda_j + 1) + \sqrt{12\lambda_j + 1}}{6}. \quad (5)$$

Proof: The set of input elements $SE_j = \{x_{j_1}, x_{j_2}, \dots, x_{j_{M_j}}\}$ can have more than $M_j \frac{n}{p}$ members only if $\lambda_j \frac{n}{p^2}$ or less members are selected to be samples from the set $SE'_j = \{x_{j_1}, x_{j_2}, \dots, x_{j_{(M_j \frac{n}{p})}}\}$, which is the

set composed of the first $M_j \frac{n}{p}$ members in SE_j . However, since each element of SE'_j is independently chosen to be a sample with probability $\frac{1}{p}$, the probability of this event occurring is given by

$$\sum_{s \leq \lambda_j \frac{n}{p}} b\left(s; M_j \frac{n}{p}, \frac{1}{p}\right). \quad (6)$$

Using the following ‘‘Chernoff’’ type bound [18] for estimating the head of a binomial distribution

$$\sum_{s \leq \epsilon r q} b(s; r, q) \leq e^{-(1-\epsilon)^2 \frac{r q}{2}}, \quad (7)$$

where $s \leq \lambda_j \frac{n}{p^2}$, $r = M_j \frac{n}{p}$, and $q = \frac{1}{p}$, it follows that the probability that a set SE_j among the p sets of input elements has more than $M_j \frac{n}{p}$ is bounded by

$$\sum_{i=0}^{p-1} e^{-\left(1 - \frac{\lambda_j}{M_j}\right)^2 \frac{M_j n}{2p^2}}. \quad (8)$$

Using the fact that $p^2 \leq \frac{n}{3 \ln n}$, it is easy to show that the above sum can be bounded by $n^{-\epsilon}$, for some $\epsilon > 0$ and

$$M_j = \frac{(6\lambda_j + 1) + \sqrt{12\lambda_j + 1}}{6}. \quad (9)$$

The bound of **Lemma 3** will also hold if we include the subsets of elements and samples whose values fall strictly between two consecutive splitters.

Lemma 4: At the completion of **Step (7)**, the number of elements received by each processor is at most $\alpha_2 \frac{n}{p}$ with high probability, for any $\alpha_2 \geq 2.62$ ($\alpha_2 \geq 1.77$ without duplicates) and $p^2 \leq \frac{n}{3 \ln n}$.

Proof: Let Q be the set of input elements to be sorted by our algorithm, let R be the set of *samples* of **Step (4)** at processor P_1 with cardinality $\beta \frac{n}{p}$, and let S be the subset of R associated with $\text{Splitter}[j]$, which we define to be the samples in R with indices $\left((j-1)\left(\beta \frac{n}{p^2}\right) + 1\right)$ through $\left(j\beta \frac{n}{p^2}\right)$, inclusively. Let $Q_1 \frac{n}{p}$, $R_1 \frac{n}{p^2}$, and $S_1 \frac{n}{p^2}$ be respectively the number of elements in Q , R , and S with value equal to $\text{Splitter}[j-1]$, let $Q_2 \frac{n}{p}$, $R_2 \frac{n}{p^2}$, and $S_2 \frac{n}{p^2}$ be respectively the number of elements in Q , R , and S with values greater than $\text{Splitter}[j-1]$ but less than $\text{Splitter}[j]$, and let $Q_3 \frac{n}{p}$, $R_3 \frac{n}{p^2}$, and $S_3 \frac{n}{p^2}$ be respectively the number of elements in Q , R , and S with value equal to $\text{Splitter}[j]$.

According to **Step (6)** of our algorithm, processor P_j will receive

$$\left(\text{Frac}_L[j] \times Q_1\right) + Q_2 + \left(\text{Frac}_R[j] \times Q_3\right) \frac{n}{p} = \left(\frac{S_1}{R_1} Q_1 + Q_2 + \frac{S_3}{R_3} Q_3\right) \frac{n}{p} \quad (10)$$

elements. To compute the upper bound $\alpha_2 \frac{n}{p}$ on this expression, we first use **Lemma 3** to bound each $Q_i \frac{n}{p}$, giving us

$$\left(\frac{S_1}{R_1} \left(\frac{(6R_1 + 1) + \sqrt{12R_1 + 1}}{6}\right) + \left(\frac{(6S_2 + 1) + \sqrt{12S_2 + 1}}{6}\right) + \frac{S_3}{R_3} \left(\frac{(6R_3 + 1) + \sqrt{12R_3 + 1}}{6}\right)\right) \frac{n}{p} \quad (11)$$

Rearranging this expression, we get:

$$\begin{aligned} & \left(S_1 \left(1 + \frac{1}{6R_1} + \sqrt{\frac{1}{3R_1} + \frac{1}{36R_1^2}} \right) + \left(\frac{(6S_2 + 1) + \sqrt{12S_2 + 1}}{6} \right) \right. \\ & \quad \left. + S_3 \left(1 + \frac{1}{6R_3} + \sqrt{\frac{1}{3R_3} + \frac{1}{36R_3^2}} \right) \right) \frac{n}{p} \end{aligned} \quad (12)$$

Clearly, this expression is maximized for $R_1 = S_1$ and $R_3 = S_3$. Substituting these values and rearranging once again, we get:

$$\left(\left(\frac{(6S_1 + 1) + \sqrt{12S_1 + 1}}{6} \right) + \left(\frac{(6S_2 + 1) + \sqrt{12S_2 + 1}}{6} \right) + \left(\frac{(6S_3 + 1) + \sqrt{12S_3 + 1}}{6} \right) \right) \frac{n}{p} \quad (13)$$

Since $S_1 + S_2 + S_3 = \beta$, this expression is maximized for $S_1 = S_2 = S_3 = \frac{\beta}{3}$. Since **Lemma 2** guarantees that with high probability $\beta \geq 1$, **Lemma 4** follows with high probability for $\alpha_2 \geq 2.62$. Alternatively, if there are no duplicates, we can show that the bound follows with high probability for $\alpha_2 \geq 1.77$.

Lemma 5: If the set of input elements is arbitrarily partitioned into at most $2p$ subsets, each of size $X_i \frac{n}{p}$ ($1 \leq i \leq 2p$), with high probability at the conclusion of **Step (2)** no processor will receive more than $Y_i \frac{n}{p^2}$ elements from any particular subset, for $Y_i \geq (X_i + \sqrt{X_i})$ and $p^2 \leq \frac{n}{3 \ln n}$.

Proof: The probability that exactly $Y_i \frac{n}{p^2}$ elements are sent to a particular processor by the conclusion of **Step (2)** is given by the binomial distribution $b(Y_i \frac{n}{p^2}; X_i \frac{n}{p}, \frac{1}{p})$. Using the Chernoff bound for estimating the tail of a binomial distribution, the probability that from M possible subsets any processor will receive at least $Y_i \frac{n}{p^2}$ elements can be bounded by

$$\sum_{i=1}^M p e^{-\left(1 - \frac{Y_i}{X_i}\right)^2 \frac{X_i n}{3p^2}} \quad (14)$$

and **Lemma 5** follows for $M \leq 2p$.

Lemma 6: The number of elements exchanged by any two processors in **Step (7)** is at most $c_2 \frac{n}{p^2}$ with high probability, for any $c_2 \geq 5.42$ ($c_2 \geq 3.10$ without duplicates) and $p^2 \leq \frac{n}{3 \ln n}$.

Proof: Let U be the set of input elements to be sorted by our algorithm, let V be the set of elements held by intermediate processor P_i after **Step (2)**, and let W be the set of elements held by destination processor P_j after **Step (7)**. Let $U_1 \frac{n}{p}$, $V_1 \frac{n}{p^2}$, and $W_1 \frac{n}{p}$ be respectively the number of elements in U , V , and W with values equal to $\text{Splitter}[j - 1]$, let $U_2 \frac{n}{p}$, $V_2 \frac{n}{p^2}$, and $W_2 \frac{n}{p}$ be respectively the number of elements in U , V , and W with values greater than $\text{Splitter}[j - 1]$ but less than $\text{Splitter}[j]$, and let $U_3 \frac{n}{p}$, $V_3 \frac{n}{p^2}$, and $W_3 \frac{n}{p}$ be respectively the number of elements in U , V , and W with values equal to $\text{Splitter}[j]$.

According to **Step (6)** of our algorithm, intermediate processor P_i will send

$$\left((\text{Frac}_{\mathbf{L}}[j] \times V_1) + V_2 + (\text{Frac}_{\mathbf{R}}[j] \times V_3) \right) \frac{n}{p^2} \quad (15)$$

elements to processor P_j . To compute the upper bound $c_2 \frac{n}{p^2}$ on this expression, we first use **Lemma 5** to bound each V_k , giving us:

$$\left(\left(\text{Frac}_{\mathbf{L}}[j] \times (U_1 + \sqrt{U_1}) \right) + (U_2 + \sqrt{U_2}) + \left(\text{Frac}_{\mathbf{R}}[j] \times (U_3 + \sqrt{U_3}) \right) \right) \frac{n}{p^2} \quad (16)$$

Notice that since destination processor P_j receives respectively $\text{Frac}_{\mathbf{L}}[j]$ and $\text{Frac}_{\mathbf{R}}[j]$ of the elements at each intermediate processor with values equal to $\text{Splitter}[j-1]$ and $\text{Splitter}[j]$, it follows that $W_1 = \text{Frac}_{\mathbf{L}}[j] \times U_1$ and $W_3 = \text{Frac}_{\mathbf{R}}[j] \times U_3$. Hence, we can rewrite the expression above as

$$\left(\frac{W_1}{U_1} (U_1 + \sqrt{U_1}) + (U_2 + \sqrt{U_2}) + \frac{W_3}{U_3} (U_3 + \sqrt{U_3}) \right) \frac{n}{p^2} \quad (17)$$

Rearranging this expression, we get:

$$\left(W_1 \left(1 + \sqrt{\frac{1}{U_1}} \right) + (U_2 + \sqrt{U_2}) + W_3 \left(1 + \sqrt{\frac{1}{U_3}} \right) \right) \frac{n}{p^2} \quad (18)$$

Clearly, this expression is maximized for $U_1 = W_1$ and $U_3 = W_3$. Substituting these values and rearranging, we get:

$$\left(W_1 + \sqrt{W_1} + W_2 + \sqrt{W_2} + W_3 + \sqrt{W_3} \right) \frac{n}{p^2} \quad (19)$$

Since $W_1 + W_2 + W_3 = \alpha_2$, this expression is maximized for $W_1 = W_2 = W_3 = \frac{\alpha_2}{3}$. Since **Lemma 4** guarantees that with high probability $\alpha_2 \geq 2.62$, **Lemma 6** follows with high probability for $c_2 \geq 5.24$. Alternatively, if there are no duplicates, we can show that the bound follows with high probability for $c_2 \geq 3.10$.

With these bounds on the values of c_1 , α_2 , and c_2 , the analysis of our sample sort algorithm is as follows. **Steps (1), (3), (4), (6), and (8)** involve no communication and are dominated by the cost of the sequential sorting in **Step (3)** and the merging in **Step (8)**. Sorting integers using radix sort requires $O\left(\frac{n}{p}\right)$ time, whereas sorting floating point numbers using merge sort requires $O\left(\frac{n}{p} \log\left(\frac{n}{p}\right)\right)$ time. **Step (8)** requires $O\left(\frac{n}{p} \log p\right)$ time if we merge the sorted subsequences in a binary tree fashion. **Steps (2), (5), and (7)** call the communication primitives **transpose**, **bcast**, and **transpose**, respectively. The analysis of these primitives in [6] shows that with high probability these three steps require $T_{comm}(n, p) \leq \left(\tau + 2\frac{n}{p^2}(p-1)\sigma\right)$, $T_{comm}(n, p) \leq (\tau + 2(p-1)\sigma)$, and $T_{comm}(n, p) \leq \left(\tau + 5.24\frac{n}{p^2}(p-1)\sigma\right)$, respectively. Hence, with high probability, the overall complexity of our sample sort algorithm is given (for floating point numbers) by

$$\begin{aligned} T(n, p) &= T_{comp}(n, p) + T_{comm}(n, p) \\ &= O\left(\frac{n}{p} \log n + \tau + \frac{n}{p}\sigma\right) \end{aligned} \quad (20)$$

for $p^2 < \frac{n}{3 \ln n}$.

Clearly, our algorithm is asymptotically optimal with very small coefficients. But a theoretical comparison of our running time with previous sorting algorithms is difficult, since there is no consensus on how to model the cost of the irregular communication used by the most efficient algorithms. Hence, it is very important to perform an empirical evaluation of an algorithm using a wide variety of benchmarks, as we will do next.

4 Performance Evaluation

Our sample sort algorithm was implemented using SPLIT-C [14] and run on a variety of machines and processors, including the Cray Research T3D, the IBM SP-2-WN, and the Thinking Machines CM-5. For every platform, we tested our code on eight different benchmarks, each of which had both a 32-bit *integer* version (64-bit on the Cray T3D) and a 64-bit double precision floating point number (*double*) version.

4.1 Sorting Benchmarks

Our eight sorting benchmarks are defined as follows, in which n and p are assumed for simplicity to be powers of two and MAX_D , the maximum value allowed for *doubles*, is approximately 1.8×10^{308} .

1. **Uniform [U]**, a uniformly distributed random input, obtained by calling the C library random number generator *random()*. This function, which returns integers in the range 0 to $(2^{31} - 1)$, is seeded by each processor P_i with the value $(21 + 1001i)$. For the *double* data type, we “normalize” the *integer* benchmark values by first subtracting the value 2^{30} and then scaling the result by $(2^{-30} \times \text{MAX}_D)$.
2. **Gaussian [G]**, a Gaussian distributed random input, approximated by adding four calls to *random()* and then dividing the result by four. For the *double* data type, we normalize the *integer* benchmark values in the manner described for [U].
3. **Zero [Z]**, a zero entropy input, created by setting every value to a constant such as zero.
4. **Bucket Sorted [B]**, an input that is sorted into p buckets, obtained by setting the first $\frac{n}{p^2}$ elements at each processor to be random numbers between 0 and $(\frac{2^{31}}{p} - 1)$, the second $\frac{n}{p^2}$ elements at each processor to be random numbers between $\frac{2^{31}}{p}$ and $(\frac{2^{32}}{p} - 1)$, and so forth. For the *double* data type, we normalize the *integer* benchmark values in the manner described for [U].
5. **g -Group [g -G]**, an input created by first dividing the processors into groups of consecutive processors of size g , where g can be any integer which partitions p evenly. If we index these groups in

consecutive order from 1 up to $\frac{n}{g}$, then for group j we set the first $\frac{n}{pg}$ elements to be random numbers between $\left(\left(\left(\left(j-1\right)g + \frac{n}{2} - 1\right) \bmod p\right) + 1\right) \frac{2^{31}}{p}$ and $\left(\left(\left(\left(j-1\right)g + \frac{n}{2}\right) \bmod p\right) + 1\right) \frac{2^{31}}{p} - 1$, the second $\frac{n}{pg}$ elements at each processor to be random numbers between $\left(\left(\left(\left(j-1\right)g + \frac{n}{2}\right) \bmod p\right) + 1\right) \frac{2^{31}}{p}$ and $\left(\left(\left(\left(j-1\right)g + \frac{n}{2} + 1\right) \bmod p\right) + 1\right) \frac{2^{31}}{p} - 1$, and so forth. For the *double* data type, we normalize the *integer* benchmark values in the manner described for [U].

6. **Staggered [S]**, created as follows: if the processor index i is less than or equal to $\frac{n}{2}$, then we set all $\frac{n}{p}$ elements at that processor to be random numbers between $\left((2i-1) \frac{2^{31}}{p}\right)$ and $\left((2i) \frac{2^{31}}{p} - 1\right)$. Otherwise, we set all $\frac{n}{p}$ elements to be random numbers between $\left((2i-p-2) \frac{2^{31}}{p}\right)$ and $\left((2i-p-1) \frac{2^{31}}{p} - 1\right)$. For the *double* data type, we normalize the *integer* benchmark values in the manner described for [U].
7. **Deterministic Duplicates [DD]**, an input of duplicates in which we set all $\frac{n}{p}$ elements at each of the first $\frac{n}{2}$ processors to be $\log n$, all $\frac{n}{p}$ elements at each of the next $\frac{n}{4}$ processors to be $\log\left(\frac{n}{2}\right)$, and so forth. At processor P_p , we set the first $\frac{n}{2p}$ elements to be $\log\left(\frac{n}{p}\right)$, the next $\frac{n}{4p}$ elements to be $\log\left(\frac{n}{2p}\right)$, and so forth.
8. **Randomized Duplicates [RD]**, an input of duplicates in which each processor fills an array T with some constant number *range* (*range* is 32 for our work) of random values between 0 and $(range-1)$ whose sum is S . The first $\frac{T[1]}{S} \times \frac{n}{p}$ values of the input are then set to a random value between 0 and $(range-1)$, the next $\frac{T[2]}{S} \times \frac{n}{p}$ values of the input are then set to another random value between 0 and $(range-1)$, and so forth.

We selected these eight benchmarks for a variety of reasons. Previous researchers have used the **Uniform**, **Gaussian**, and **Zero** benchmarks, and so we too included them for purposes of comparison. But benchmarks should be designed to illicit the worst case behavior from an algorithm, and in this sense the **Uniform** benchmark is not appropriate. For example, for $n \gg p$, one would expect that the optimal choice of the *splitters* in the **Uniform** benchmark would be those which partition the range of possible values into equal intervals. Thus, algorithms which try to guess the *splitters* might perform misleadingly well on such an input. In this respect, the **Gaussian** benchmark is more telling. But we also wanted to find benchmarks which would evaluate the cost of irregular communication. Thus, we wanted to include benchmarks for which an algorithm which uses a single phase of routing would find contention difficult or even impossible to avoid. A naive approach to rearranging the data would perform poorly on the **Bucket Sorted** benchmark. Here, every processor would try to route data to the same processor at the same time, resulting in poor utilization of communication bandwidth. This problem might be avoided by an algorithm in which at each processor the elements are first grouped by destination and then routed according to the specifications of a sequence of p destination permutations. Perhaps the most straightforward way to do this is by iterating over

the possible communication strides. But such a strategy would perform poorly with the *g-Group* benchmark, for a suitably chosen value of g . In this case, using stride iteration, those processors which belong to a particular group all route data to the same subset of g destination processors. This subset of destinations is selected so that, when the g processors route to this subset, they choose the processors in exactly the same order, producing contention and possibly stalling. Alternatively, one can synchronize the processors after each permutation, but this in turn will reduce the communication bandwidth by a factor of $\frac{p}{g}$. In the worst case scenario, each processor needs to send data to a single processor a unique stride away. This is the case of the **Staggered** benchmark, and the result is a reduction of the communication bandwidth by a factor of p . Of course, one can correctly object that both the **g-Group** benchmark and the **Staggered** benchmark have been tailored to thwart a routing scheme which iterates over the possible strides, and that another sequences of permutations might be found which performs better. This is possible, but at the same time we are unaware of any single phase deterministic algorithm which could avoid an equivalent challenge. Finally, the **Deterministic Duplicates** and the **Randomized Duplicates** benchmarks were included to assess the performance of the algorithms in the presence of duplicate values.

4.2 Experimental Results

For each experiment, the input is evenly distributed amongst the processors. The output consists of the elements in non-descending order arranged amongst the processors so that the elements at each processor are in sorted order and no element at processor P_i is greater than any element at processor P_j , for all $i < j$.

Two variations were allowed in our experiments. First, radix sort was used to sequentially sort *integers*, whereas merge sort was used to sort double precision floating point numbers (*doubles*). Second, different implementations of the communication primitives were allowed for each machine. Wherever possible, we tried to use the vendor supplied implementations. In fact, IBM does provide all of our communication primitives as part of its machine specific Collective Communication Library (CCL) [8] and MPI. As one might expect, they were faster than the high level SPLIT-C implementation.

| Size | [U] | [G] | [2-G] | [4-G] | [B] | [S] | [Z] | [DD] | [RD] |
|-------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 256K | 0.019 | 0.019 | 0.020 | 0.020 | 0.020 | 0.020 | 0.016 | 0.016 | 0.018 |
| 1M | 0.068 | 0.068 | 0.070 | 0.070 | 0.070 | 0.069 | 0.054 | 0.054 | 0.058 |
| 4M | 0.261 | 0.257 | 0.264 | 0.264 | 0.263 | 0.264 | 0.202 | 0.226 | 0.213 |
| 16M | 1.02 | 1.01 | 1.02 | 1.02 | 1.02 | 1.02 | 0.814 | 0.831 | 0.826 |
| 64M | 4.03 | 4.00 | 4.00 | 3.99 | 4.03 | 4.00 | 3.21 | 3.20 | 3.27 |

Table I: Total execution time (in seconds) required to sort a variety of *integer* benchmarks on a 64-node Cray T3D.

| Size | [U] | [G] | [2-G] | [4-G] | [B] | [S] | [Z] | [DD] | [RD] |
|-------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 256K | 0.041 | 0.039 | 0.040 | 0.041 | 0.041 | 0.040 | 0.042 | 0.040 | 0.041 |
| 1M | 0.071 | 0.071 | 0.074 | 0.072 | 0.076 | 0.072 | 0.071 | 0.070 | 0.070 |
| 4M | 0.215 | 0.210 | 0.219 | 0.213 | 0.218 | 0.218 | 0.207 | 0.213 | 0.213 |
| 16M | 0.805 | 0.806 | 0.817 | 0.822 | 0.830 | 0.818 | 0.760 | 0.760 | 0.783 |
| 64M | 3.30 | 3.19 | 3.22 | 3.24 | 3.28 | 3.25 | 2.79 | 2.83 | 2.83 |

Table II: Total execution time (in seconds) required to sort a variety of *integer* benchmarks on a 64-node IBM SP-2-WN.

| Size | [U] | [G] | [2-G] | [4-G] | [B] | [S] | [Z] | [DD] | [RD] |
|-------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 256K | 0.022 | 0.022 | 0.023 | 0.023 | 0.023 | 0.022 | 0.021 | 0.021 | 0.021 |
| 1M | 0.089 | 0.089 | 0.088 | 0.089 | 0.090 | 0.088 | 0.082 | 0.082 | 0.083 |
| 4M | 0.366 | 0.366 | 0.364 | 0.366 | 0.364 | 0.362 | 0.344 | 0.344 | 0.341 |
| 16M | 1.55 | 1.55 | 1.50 | 1.54 | 1.53 | 1.52 | 1.45 | 1.46 | 1.47 |
| 64M | 6.63 | 6.54 | 6.46 | 6.44 | 6.46 | 6.52 | 6.23 | 6.25 | 6.24 |

Table III: Total execution time (in seconds) required to sort a variety of *double* benchmarks on a 64-node Cray T3D.

Tables I, II, III, and IV display the performance of our sample sort as a function of input distribution for a variety of input sizes. In each case, the performance is essentially independent of the input distribution. These tables present results obtained on a 64 node Cray T3D and a 64 node IBM SP-2; results obtained from the TMC CM-5 validate this claim as well. Because of this independence, the remainder of this section will only discuss the performance of our sample sort on the single benchmark [U].

The results in **Tables V and VI** together with their graphs in **Figure 1** examine the scalability of our sample sort as a function of machine size. Results are shown for the T3D, the SP-2-WN, and the CM-5. Bearing in mind that these graphs are log-log plots, they show that for a *fixed input size* n the execution time scales almost inversely with the number of processors p . While this is certainly the expectation of our analytical model for *doubles*, it might at first appear to exceed our prediction of an $O\left(\frac{n}{p} \log p\right)$ computational complexity for *integers*. However, the appearance of an inverse relationship is still quite reasonable when we note that, for values of p between 8 and 128, $\log p$ varies

| Size | [U] | [G] | [2-G] | [4-G] | [B] | [S] | [Z] | [DD] | [RD] |
|-------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 256K | 0.056 | 0.054 | 0.059 | 0.057 | 0.060 | 0.059 | 0.056 | 0.056 | 0.057 |
| 1M | 0.153 | 0.152 | 0.158 | 0.156 | 0.163 | 0.156 | 0.151 | 0.146 | 0.147 |
| 4M | 0.568 | 0.565 | 0.576 | 0.577 | 0.584 | 0.575 | 0.558 | 0.571 | 0.569 |
| 16M | 2.23 | 2.23 | 2.24 | 2.28 | 2.26 | 2.25 | 2.20 | 2.22 | 2.26 |
| 64M | 9.24 | 9.18 | 9.24 | 9.22 | 9.24 | 9.23 | 9.15 | 9.17 | 9.21 |

Table IV: Total execution time (in seconds) for required to sort a variety of *double* benchmarks on a 64-node IBM SP-2-WN.

| Sample Sorting of 8M Integers [U] | | | | | |
|-----------------------------------|----------------------|------|-------|-------|-------|
| Machine | Number of Processors | | | | |
| | 8 | 16 | 32 | 64 | 128 |
| CRAY T3D | 3.32 | 1.77 | 0.952 | 0.513 | 0.284 |
| IBM SP2-WN | 2.51 | 1.25 | 0.699 | 0.413 | 0.266 |
| TMC CM-5 | - | 7.43 | 3.72 | 1.73 | 0.813 |

Table V: Total execution time (in seconds) required to sort 8M *integers* on a variety of machines and processors using the [U] benchmark. A hyphen indicates that particular platform was unavailable to us.

| Sample Sorting of 8M Doubles [U] | | | | | |
|----------------------------------|----------------------|------|------|-------|-------|
| Machine | Number of Processors | | | | |
| | 8 | 16 | 32 | 64 | 128 |
| CRAY T3D | 5.48 | 2.78 | 1.44 | 0.747 | 0.392 |
| IBM SP2-WN | 7.96 | 4.02 | 2.10 | 1.15 | 0.635 |
| TMC CM-5 | - | - | 6.94 | 3.79 | 1.83 |

Table VI: Total execution time (in seconds) required to sort 8M *doubles* on a variety of machines and processors using the [U] benchmark. A hyphen indicates that particular platform was unavailable to us.

by only a factor of $\frac{7}{3}$. Moreover, this $O\left(\frac{n}{p} \log p\right)$ complexity is entirely due to the merging in **Step (8)**, and in practice, **Step (8)** never accounts for more than 30% of the observed execution time. Note that the complexity of **Step 8** could be reduced to $O\left(\frac{n}{p}\right)$ for *integers* using radix sort, but the resulting execution time would, in most cases, be slower.

The graphs in **Figure 2** examine the scalability of our sample sort as a function of problem size, for differing numbers of processors. They show that for a *fixed number of processors* there is an almost linear dependence between the execution time and the total number of elements n . While this is certainly the expectation of our analytic model for *integers*, it might at first appear to exceed our prediction of a $O\left(\frac{n}{p} \log n\right)$ computational complexity for floating point values. However, this appearance of a linear relationship is still quite reasonable when we consider that for the range of values shown $\log n$ differs by only a factor of 1.2.

Next, the graphs in **Figure 3** examine the relative costs of the eight steps in our sample sort. Results are shown for both a 64 node T3D and a 64 node SP-2-WN, using both the *integer* and the *double* versions of the [U] benchmark. Notice that for $n = 64M$ *integers*, the sequential sorting and merging performed in **Steps (3)** and **(8)** consume approximately 80% of the execution time on the T3D and approximately 70% of the execution time on the SP-2. By contrast, the two **transpose** operations in **Steps (2)** and **(7)** together consume only about 15% of the execution time on the T3D and about 25% of the execution time on the SP-2. The difference in the distribution between these two platforms is likely due in part to the fact that an *integer* is 64 bits on the T3D while only 32 bits on the SP-2. By contrast, *doubles* are 64 bits on both platforms. For $n = 64M$ *doubles*, the sequential sorting

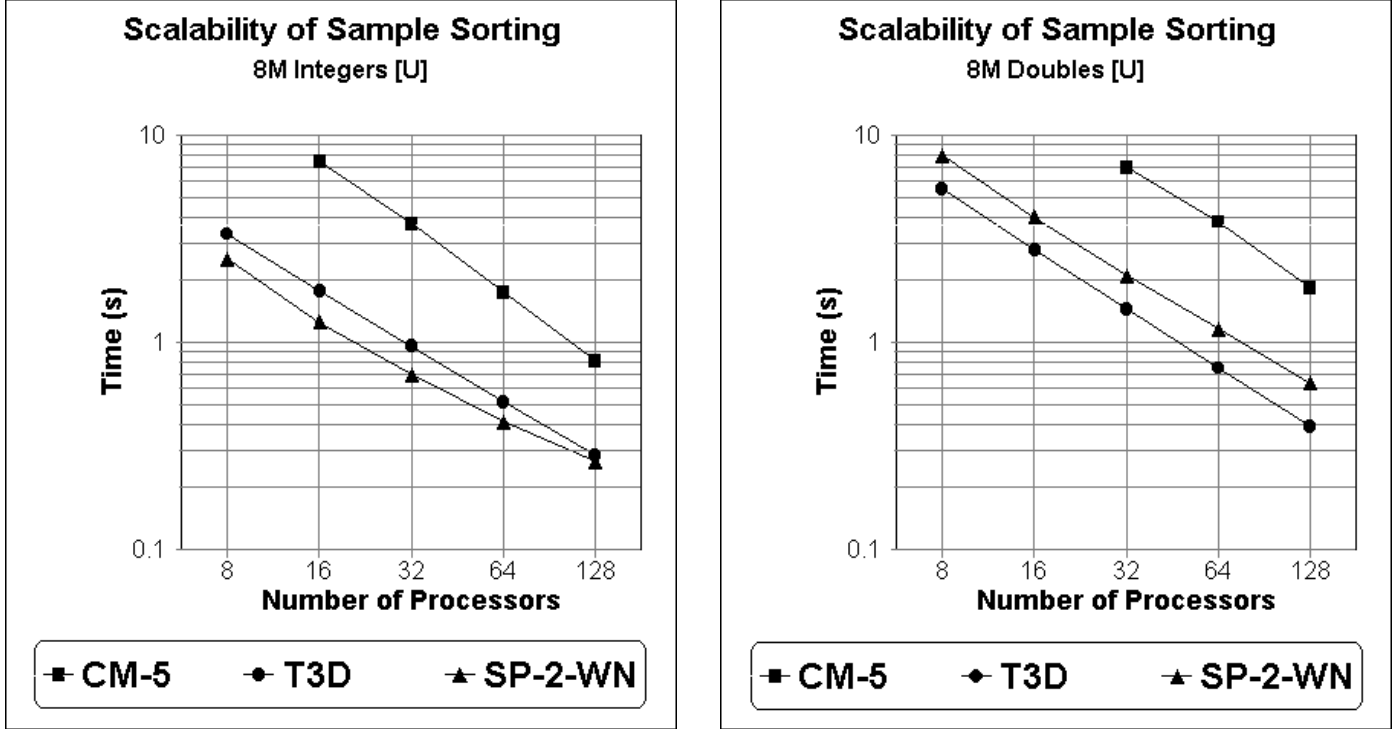


Figure 1: Scalability of sorting *integers* and *doubles* with respect to machine size.

and merging performed in **Steps (3)** and **(8)** consume approximately 80% of the execution time on both platforms, whereas the two **transpose** operations in **Steps (2)** and **(7)** together consume only about 15% of the execution time. Together, these results show that our algorithm is extremely efficient in its communication performance.

Finally, **Tables VII** and **VIII** show the experimentally derived expected value (E) and sample standard deviation (STD) of the coefficients c_1 , α_1 , c_2 , and α_2 used to describe the complexity of our algorithm in **Section 3**. The values in **Table VII** were obtained by analyzing data collected while sorting each of the *duplicate* benchmarks [DD] and [RD] 50 times on a 64-node Cray T3D. For each trial, the values recorded were the largest occurrence of each coefficient at any of the 64 processors. By contrast, the values in **Table VIII** were obtained by analyzing data collected while sorting each of the *unique* benchmarks [G], [B], [2-G], [4-G], and [S] 20 times. In every trial, a different seed was used for the random number generator, both to generate the benchmark where appropriate and to distribute the keys into bins as part of **Step (1)**. The experimentally derived expected values in **Table VII** for c_1 , α_1 , c_2 , and α_2 agree strongly with the theoretically derived bounds for *duplicate* keys of $c_1 \leq 2$, $\alpha_1 \leq c_1$, $c_2 \leq 5.24$, and $\alpha_2 \leq 2.62$ for $p^2 \leq \frac{n}{3 \ln n}$. Similarly, the experimentally derived expected values in **Table VIII** for c_1 , α_1 , c_2 , and α_2 agree strongly with the theoretically derived bounds for *unique* keys of $c_1 \leq 2$, $\alpha_1 \leq c_1$, $c_2 \leq 3.10$, and $\alpha_2 \leq 1.77$ for $p^2 \leq \frac{n}{3 \ln n}$. Note that expected values for c_2 and α_2 are actually less for *duplicate* values than for *unique* values, which is the opposite of what we might expect from the computed bounds. This might simply reflect our limited choice of

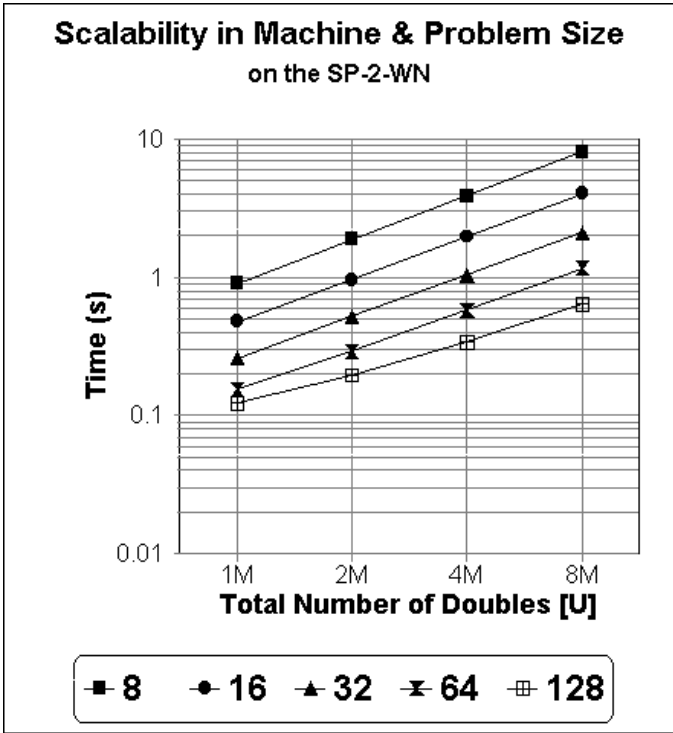
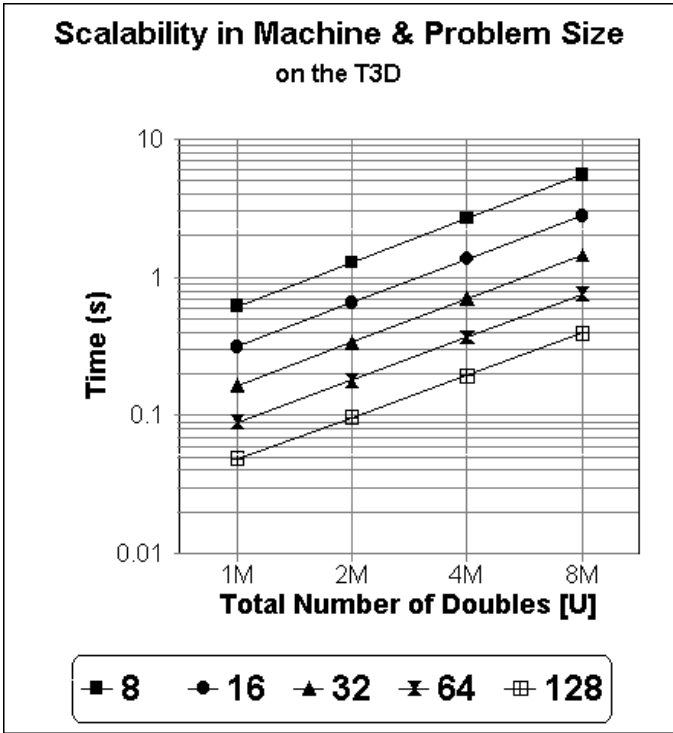
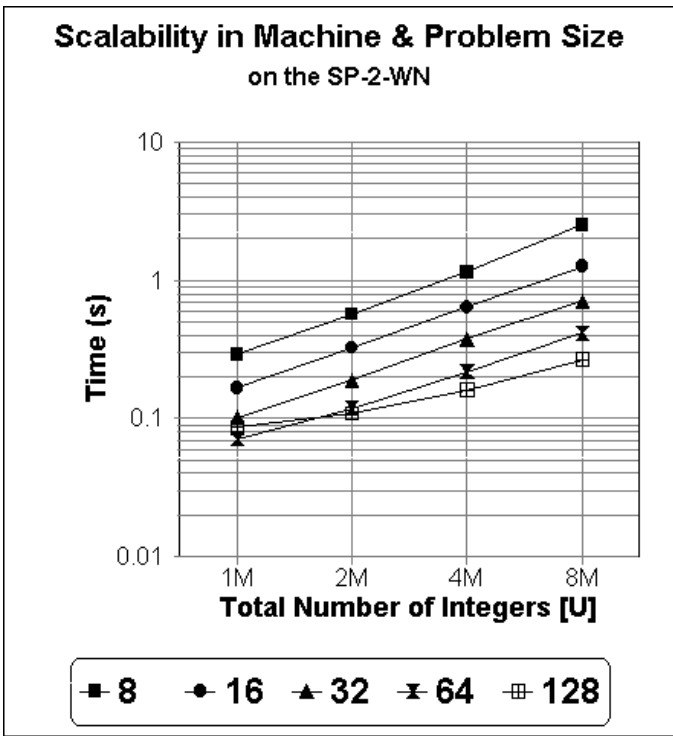
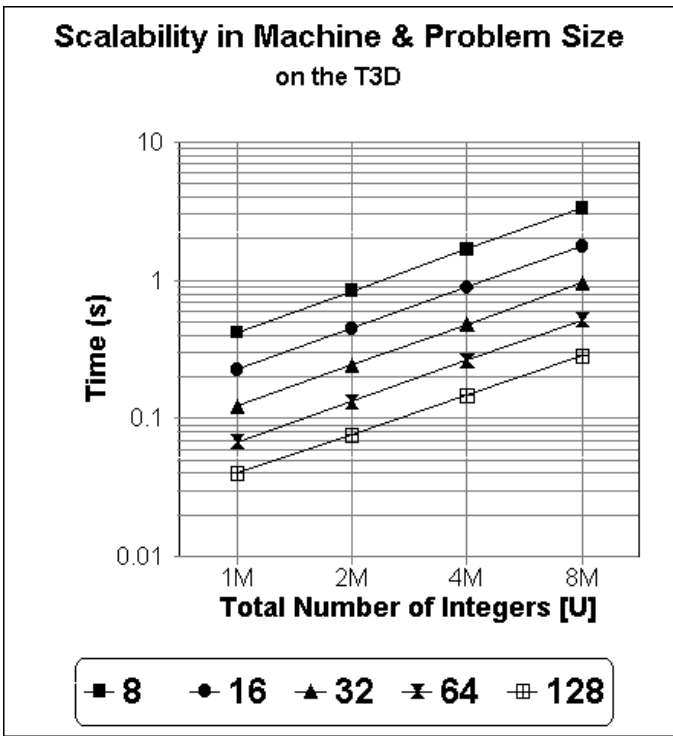


Figure 2: Scalability of sorting *integers* and *doubles* with respect to the problem size, for differing numbers of processors.

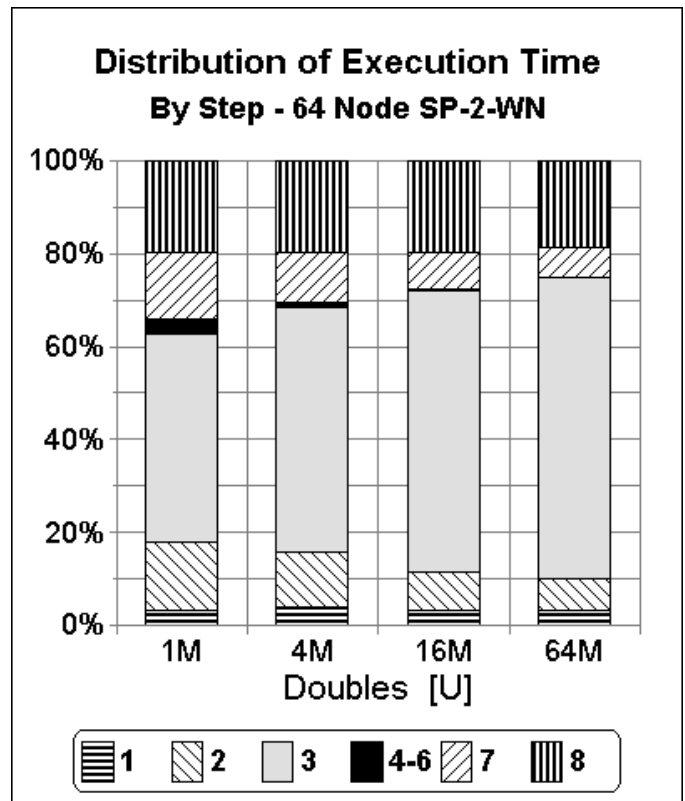
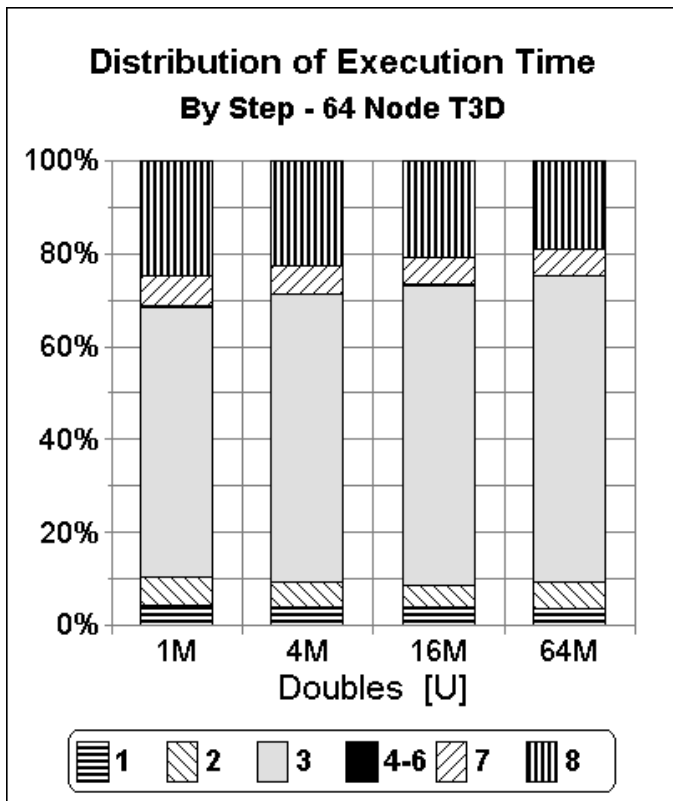
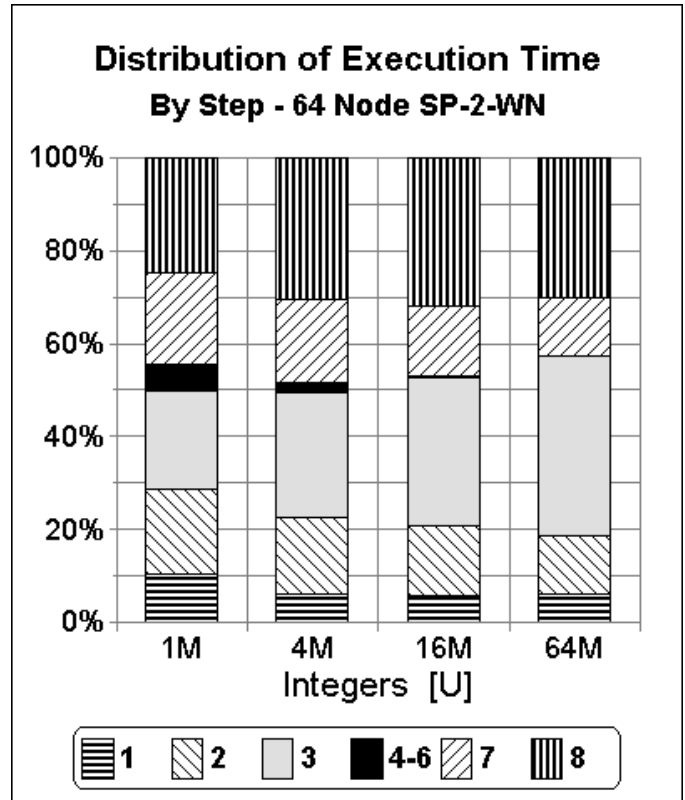
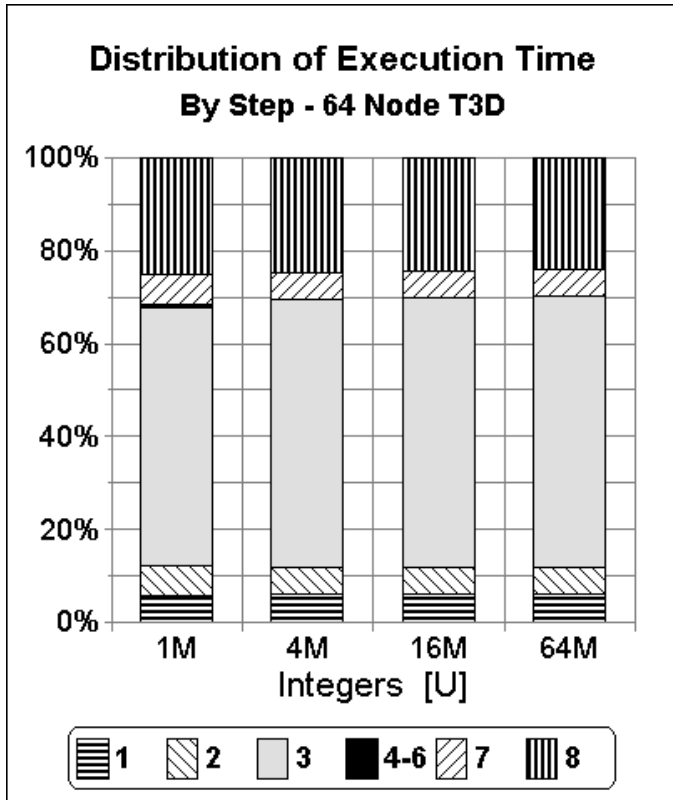


Figure 3: Distribution of execution time amongst the eight steps of sample sort. Times are obtained for both a 64 node T3D and a 64 node SP-2-WN using both the *integer* and the *double* versions of the [U] benchmark.

| keys/proc | $E(c_1)$ | $STD(c_1)$ | $E(\alpha_1)$ | $STD(\alpha_1)$ | $E(c_2)$ | $STD(c_2)$ | $E(\alpha_2)$ | $STD(\alpha_2)$ |
|-----------|----------|------------|---------------|-----------------|----------|------------|---------------|-----------------|
| 4K | 2.02 | 0.104 | 1.08 | 0.019 | 2.12 | 0.336 | 1.45 | 0.183 |
| 16K | 1.48 | 0.044 | 1.04 | 0.008 | 1.49 | 0.133 | 1.18 | 0.089 |
| 64K | 1.23 | 0.019 | 1.02 | 0.0003 | 1.24 | 0.062 | 1.09 | 0.044 |
| 256K | 1.11 | 0.009 | 1.01 | 0.002 | 1.12 | 0.026 | 1.04 | 0.020 |
| 1M | 1.06 | 0.005 | 1.00 | 0.001 | 1.06 | 0.015 | 1.02 | 0.012 |

Table VII: Statistical evaluation of the experimentally observed values of the algorithm coefficients on a 64 node T3D using the *duplicate* benchmarks.

| keys/proc | $E(c_1)$ | $STD(c_1)$ | $E(\alpha_1)$ | $STD(\alpha_1)$ | $E(c_2)$ | $STD(c_2)$ | $E(\alpha_2)$ | $STD(\alpha_2)$ |
|-----------|----------|------------|---------------|-----------------|----------|------------|---------------|-----------------|
| 4K | 2.02 | 0.091 | 1.08 | 0.017 | 2.64 | 0.935 | 1.55 | 0.181 |
| 16K | 1.48 | 0.044 | 1.04 | 0.007 | 1.65 | 0.236 | 1.25 | 0.074 |
| 64K | 1.23 | 0.021 | 1.02 | 0.0003 | 1.30 | 0.087 | 1.12 | 0.034 |
| 256K | 1.11 | 0.010 | 1.01 | 0.002 | 1.14 | 0.034 | 1.06 | 0.019 |
| 1M | 1.06 | 0.005 | 1.00 | 0.001 | 1.07 | 0.013 | 1.03 | 0.0108 |

Table VIII: Statistical evaluation of the experimentally observed values of the algorithm coefficients on a 64 node T3D using the *unique* benchmarks.

benchmarks, or it may suggest that the bounds computed for *duplicate* are looser than those computed for *unique* values.

4.3 Comparison with Previous Results

Despite the enormous theoretical interest in parallel sorting, we were able to locate relatively few empirical studies. Of these, only a few were done on machines which either were available to us for comparison or involved code which could be ported to these machines for comparison. In **Tables IX** and **X**, we compare the performance of our sample sort algorithm with two other sample sort algorithms. In all cases, the code was written in SPLIT-C. In the case of Alexandrov et al. [1], the times were determined by us directly on a 32 node CM-5 using code supplied by the authors which had been optimized for a Meiko CS-2. In the case of Dusseau [17], the times were obtained from the graphed results reported for a 64 node CM-5.

Finally, there are the results for the NAS Parallel Benchmark [30] for Integer Sorting (IS). The name of this benchmark is somewhat misleading. Instead of requiring that the integers be placed in sorted order as we do, the benchmark only requires that they be ranked without any reordering, which is a significantly simpler task. Specifically, the Class A Benchmark requires *ten* repeated rankings of a Gaussian distributed random input consisting of 2^{23} integers ranging in value from 0 to $(2^{19} - 1)$. The Class B Benchmark is similar, except that the input consists of 2^{25} integers ranging in value from 0 to $(2^{21} - 1)$. **Table XI** compares our results on these two variations of the NAS Benchmark with the best reported times for the CM-5, the T3D, and the SP-2-WN. We believe that our results, which were

| int./proc. | [U] | | [G] | | [2-G] | | [B] | | [S] | |
|------------|-------|-------|-------|-------|-------|------|-------|-------|-------|-----|
| | HBJ | AIS | HBJ | AIS | HBJ | AIS | HBJ | AIS | HBJ | AIS |
| 4K | 0.049 | 0.153 | 0.050 | 0.152 | 0.051 | 1.05 | 0.055 | 0.181 | 0.049 | † |
| 8K | 0.090 | 0.197 | 0.090 | 0.192 | 0.092 | 1.09 | 0.094 | 0.193 | 0.087 | † |
| 16K | 0.172 | 0.282 | 0.171 | 0.281 | 0.173 | 1.16 | 0.173 | 0.227 | 0.175 | † |
| 32K | 0.332 | 0.450 | 0.330 | 0.449 | 0.335 | 1.34 | 0.335 | 0.445 | 0.338 | † |
| 64K | 0.679 | 0.833 | 0.679 | 0.835 | 0.683 | 1.76 | 0.685 | 0.823 | 0.686 | † |
| 128K | 1.65 | 2.02 | 1.64 | 2.02 | 1.64 | 2.83 | 1.64 | 1.99 | 1.64 | † |
| 256K | 3.72 | 4.69 | 3.71 | 4.59 | 3.71 | 5.13 | 3.70 | 4.56 | 3.71 | † |
| 512K | 7.97 | 10.0 | 7.85 | 9.91 | 7.93 | 9.58 | 7.95 | 9.98 | 7.95 | † |

Table IX: Total execution time (in seconds) required to sort a variety of benchmarks and problem sizes, comparing our version of sample sort (**HBJ**) with that of Alexandrov et al. (**AIS**) on a 32-node CM-5.

†We were unable to run the (**AIS**) code on this input.

| int./proc. | [U] | | [B] | | [Z] | |
|------------|------|-----|------|-----|------|-----|
| | HBJ | DUS | HBJ | DUS | HBJ | DUS |
| 1M | 16.6 | 21 | 12.2 | 91 | 10.6 | 11 |

Table X: Time required per element (in microseconds) to sample sort 64M *integers*, comparing our results (**HBJ**) with those obtained from the graphed results reported by Dusseau (**DUS**) on a 64 node CM-5.

obtained using high-level, portable code, compare favorably with the other reported times, which were obtained by the vendors using machine-specific implementations and perhaps system modifications.

The only performance studies we are aware of on similar platforms for generalized sorting are those of Tridgell and Brent [32], who report the performance of their algorithm using a 32 node CM-5 on a uniformly distributed random input of signed integers, as described in **Table XII**.

5 Conclusion

In this paper, we introduced a novel variation on sample sort and conducted an experimental study of its performance on a number of platforms using widely different benchmarks. Our results illustrate the efficiency and scalability of our algorithm across the different platforms and appear to improve on all similar results known to the authors. Our results also compare favorably with those reported for the simpler ranking problem posed by the NAS Integer Sorting (IS) Benchmark.

We have also studied several variations on our algorithm which use differing strategies to ensure that every bucket in **Step (1)** receives an equal number of elements. The results obtained for these variations were very similar to those reported in this paper. On no platform did the improvements exceed approximately 5%, and in many instances they actually ran more slowly. We believe that a significant improvement of our algorithm would require the enhancement of the sequential sorting and merging in **Steps (3)** and **(8)**, and that there is little room for significant improvement in either the

| Comparison of NAS (IS) Benchmark Times | | | | | |
|--|----------------------|--------------------|----------|--------------------|----------|
| Machine | Number of Processors | Class A | | Class B | |
| | | Best Reported Time | Our Time | Best Reported Time | Our Time |
| CM-5 | 32 | 43.1 | 29.8 | NA | 124 |
| | 64 | 24.2 | 13.7 | NA | 66.4 |
| | 128 | 12.0 | 7.03 | NA | 33.0 |
| T3D | 16 | 7.07 | 12.3 | NA | 60.1 |
| | 32 | 3.89 | 6.82 | 16.57 | 29.3 |
| | 64 | 2.09 | 3.76 | 8.74 | 16.2 |
| | 128 | 1.05 | 2.12 | 4.56 | 8.85 |
| SP-2-WN | 16 | 2.65 | 10.3 | 10.60 | 46.6 |
| | 32 | 1.54 | 5.97 | 5.92 | 25.5 |
| | 64 | 0.89 | 3.68 | 3.41 | 13.6 |
| | 128 | 0.59 | 2.52 | 1.98 | 8.45 |

Table XI: Comparison of our execution time (in seconds) with the best reported times for the Class A and Class B NAS Parallel Benchmark for *integer* sorting. Note that while we actually place the integers in sorted order, the benchmark only requires that they be ranked without actually reordering.

| Problem Size | [U] | |
|--------------|------|------|
| | HBJ | TB |
| 8M | 4.57 | 5.48 |

Table XII: Total execution time (in seconds) required to sort 8M signed *integers*, comparing our results (**HBJ**) with those of Tridgell and Brent (**TB**) on a 32 node CM-5.

load balance or the communication efficiency.

6 Acknowledgements

We would like to thank Ronald Greenberg of UMCP’s Electrical Engineering Department for his valuable comments and encouragement.

We would also like to thank the CASTLE/SPLIT-C group at The University of California, Berkeley, especially for the help and encouragement from David Culler, Arvind Krishnamurthy, and Lok Tin Liu.

We acknowledge the use of the UMIACS 16-node IBM SP-2-TN2, which was provided by an IBM Shared University Research award and an NSF Academic Research Infrastructure Grant No. CDA9401151.

Arvind Krishnamurthy provided additional help with his port of SPLIT-C to the Cray Research T3D [2]. The Jet Propulsion Lab/Caltech 256-node Cray T3D Supercomputer used in this investigation was provided by funding from the NASA Offices of Mission to Planet Earth, Aeronautics, and Space Science. We also acknowledge William Carlson and Jesse Draper from the Center for Comput-

ing Science (formerly Supercomputing Research Center) for writing the parallel compiler AC (version 2.6) [11] on which the T3D port of SPLIT-C has been based.

We also thank the Numerical Aerodynamic Simulation Systems Division of NASA's Ames Research Center for use of their 160-node IBM SP-2-WN.

This work also utilized the CM-5 at National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, under grant number ASC960008N.

Please see <http://www.umiacs.umd.edu/research/EXPAR> for related work by the authors. All the code used in this paper is freely available to interested parties from our anonymous ftp site <ftp://ftp.umiacs.umd.edu/pub/EXPAR>. We encourage other researchers to compare with our results on similar inputs.

References

- [1] A. Alexandrov, M. Ionescu, K. Schauer, and C. Scheiman. LogGP: Incorporating Long Messages into the LogP Model - One step closer towards a realistic model for parallel computation. In *7th Annual ACM Symposium on Parallel Algorithms and Architectures*, pages 95–105, Santa Barbara, CA, July 1995.
- [2] R.H. Arpaci, D.E. Culler, A. Krishnamurthy, S.G. Steinberg, and K. Yelick. Empirical Evaluation of the CRAY-T3D: A Compiler Perspective. In ACM Press, editor, *Proceedings of the 22nd Annual International Symposium on Computer Architecture*, pages 320–331, Santa Margherita Ligure, Italy, June 1995.
- [3] D.A. Bader, D.R. Helman, and J. JáJá. Practical Parallel Algorithms for Personalized Communication and Integer Sorting. CS-TR-3548 and UMIACS-TR-95-101 Technical Report, UMIACS and Electrical Engineering, University of Maryland, College Park, MD, November 1995. To appear in *ACM Journal of Experimental Algorithmics*.
- [4] D.A. Bader and J. JáJá. Parallel Algorithms for Image Histogramming and Connected Components with an Experimental Study. Technical Report CS-TR-3384 and UMIACS-TR-94-133, UMIACS and Electrical Engineering, University of Maryland, College Park, MD, December 1994. To appear in *Journal of Parallel and Distributed Computing*.
- [5] D.A. Bader and J. JáJá. Parallel Algorithms for Image Histogramming and Connected Components with an Experimental Study. In *Fifth ACM SIGPLAN Symposium of Principles and Practice of Parallel Programming*, pages 123–133, Santa Barbara, CA, July 1995.
- [6] D.A. Bader and J. JáJá. Practical Parallel Algorithms for Dynamic Data Redistribution, Median Finding, and Selection. Technical Report CS-TR-3494 and UMIACS-TR-95-74, UMIACS and Electrical Engineering, University of Maryland, College Park, MD, July 1995. Presented at the 10th *International Parallel Processing Symposium*, pages 292-301, Honolulu, HI, April 15-19, 1996.
- [7] D. Bailey, E. Barszcz, J. Barton, D. Browning, R. Carter, L. Dagum, R. Fatoohi, S. Fineberg, P. Frederickson, T. Lasinski, R. Schreiber, H. Simon, V. Venkatakrisnan, and S. Weeratunga. The NAS Parallel Benchmarks. Technical Report RNR-94-007, Numerical Aerodynamic Simulation Facility, NASA Ames Research Center, Moffett Field, CA, March 1994.
- [8] V. Bala, J. Bruck, R. Cypher, P. Elustondo, A. Ho, C.-T. Ho, S. Kipnis, and M. Snir. CCL: A Portable and Tunable Collective Communication Library for Scalable Parallel Computers. *IEEE Transactions on Parallel and Distributed Systems*, 6:154–164, 1995.

- [9] K. Batcher. Sorting Networks and Their Applications. In *Proceedings of the AFIPS Spring Joint Computer Conference 32*, pages 307–314, Reston, VA, 1968.
- [10] G.E. Blelloch, C.E. Leiserson, B.M. Maggs, C.G. Plaxton, S.J. Smith, and M. Zagha. A Comparison of Sorting Algorithms for the Connection Machine CM-2. In *Proceedings of the ACM Symposium on Parallel Algorithms and Architectures*, pages 3–16, July 1991.
- [11] W.W. Carlson and J.M. Draper. AC for the T3D. Technical Report SRC-TR-95-141, Supercomputing Research Center, Bowie, MD, February 1995.
- [12] H. Chernoff. A Measure of Asymptotic Efficiency for Tests of a Hypothesis Based on the Sum of Observations. *Annals of Math. Stat.*, 23:493–509, 1952.
- [13] Cray Research, Inc. *SHMEM Technical Note for C*, October 1994. Revision 2.3.
- [14] D.E. Culler, A. Dusseau, S.C. Goldstein, A. Krishnamurthy, S. Lumetta, T. von Eicken, and K. Yelick. Parallel Programming in Split-C. In *Proceedings of Supercomputing '93*, pages 262–273, Portland, OR, November 1993.
- [15] D.E. Culler, A.C. Dusseau, R.P. Martin, and K.E. Schauser. Fast Parallel Sorting Under LogP: From Theory to Practice. In *Portability and Performance for Parallel Processing*, chapter 4, pages 71–98. John Wiley & Sons, 1993.
- [16] D.E. Culler, R.M. Karp, D.A. Patterson, A. Sahay, K.E. Schauser, E. Santos, R. Subramonian, and T. von Eicken. LogP: Towards a Realistic Model of Parallel Computation. In *Fourth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, May 1993.
- [17] A.C. Dusseau. Modeling Parallel Sorts with LogP on the CM-5. Technical Report UCB//CSD-94-829, Computer Science Division, University of California, Berkeley, 1994.
- [18] T. Hagerup and C. Rüb. A Guided Tour of Chernoff Bounds. *Information Processing Letters*, 33:305–308, 1990.
- [19] D.R. Helman, D.A. Bader, and J. JáJá. Parallel Algorithms for Personalized Communication and Sorting With an Experimental Study. In *Proceedings of the Eighth Annual ACM Symposium on Parallel Algorithms and Architectures*, pages 211–220, Padua, Italy, June 1996.
- [20] W.L. Hightower, J.F. Prins, and J.H. Reif. Implementations of Randomized Sorting on Large Parallel Machines. In *Proceedings of the 4th Symposium on Parallel Architectures and Algorithms*, pages 158–167, San Diego, CA, July 1992.
- [21] J.S. Huang and Y.C. Chow. Parallel Sorting and Data Partitioning by Sampling. In *Proceedings of the 7th Computer Software and Applications Conference*, pages 627–631, November 1983.
- [22] F.T. Leighton. Tight Bounds on the Complexity of Parallel Sorting. *IEEE Transactions on Computers*, C-34:344–354, 1985.

- [23] H. Li and K.C. Sevcik. Parallel Sorting by Overpartitioning. Technical Report CSRI-295, Computer Systems Research Institute, University of Toronto, Canada, April 1994.
- [24] X. Li, P. Lu, J. Schaeffer, J. Shillington, P.S. Wong, and H. Shi. On the Versatility of Parallel Sorting by Regular Sampling. *Parallel Computing*, 19:1079–1103, 1993.
- [25] J.M. Marberg and E. Gafni. Sorting in Constant Number of Row and Column Phases on a Mesh. *Algorithmica*, 3:561–572, 1988.
- [26] Message Passing Interface Forum. MPI: A Message-Passing Interface Standard. Technical report, University of Tennessee, Knoxville, TN, June 1995. Version 1.1.
- [27] C.G. Plaxton. Efficient Computation on Sparse Interconnection Networks. Technical Report STAN-CS-89-1283, Department of Computer Science, Stanford University, Stanford, CA, September 1989.
- [28] M.J. Quinn. Analysis and Benchmarking of Two Parallel Sorting Algorithms: Hyperquicksort and Quickmerge. *BIT*, 29:239–250, 1989.
- [29] J.H. Reif and L.G. Valiant. A Logarithmic Time Sort for Linear Sized Networks. *Journal of the ACM*, 34:60–76, 1987.
- [30] S. Saini and D.H. Bailey. NAS Parallel Benchmarks Results 12-95. Report NAS-95-021, Numerical Aerodynamic Simulation Facility, NASA Ames Research Center, Moffett Field, CA, December 1995.
- [31] H. Shi and J. Schaeffer. Parallel Sorting by Regular Sampling. *Journal of Parallel and Distributed Computing*, 14:361–372, 1992.
- [32] A. Tridgell and R.P. Brent. An Implementation of a General-Purpose Parallel Sorting Algorithm. Technical Report TR-CS-93-01, Computer Sciences Laboratory, Australian National University, Canberra, Australia, February 1993.
- [33] L.G. Valiant. A Bridging Model for Parallel Computation. *Communications of the ACM*, 33(8):103–111, 1990.
- [34] S.C. Woo, M. Ohara, E. Torrie, J.P. Singh, and A. Gupta. The SPLASH-2 Programs: Characterization and Methodological Considerations. In *Proceedings of the 22nd Annual International Symposium on Computer Architecture*, pages 24–36, June 1995.