

to be combined and communicated are partitioned into d (the dimension of the cube) equal parts, where one part is exchanged in each direction. This kind of approach to increasing the utilization of the communication network is discussed in [4], Chapter 21. The effect is to reduce the communication overhead by multiplying the β term in the time complexities of Versions 2 and 3 by a factor d , which is then carried through to reduce the execution time of the hybrid algorithms as well. This is beyond the scope of this paper.

Acknowledgements

The author would like to thank Dr. Israel Nelken and Dr. Bernard L. Buzen for their helpful comments.

References

- [1] E. Anderson, A. Benzon, J. Dongarra, S. Moulton, S. Oosterhout, and R. van de Geijn. LAPACK for distributed memory architectures. *Parallel and Distributed Processing for Scientific Computing, Fifth Symposium*, 1991.
- [2] J.J. Dongarra and R.A. van de Geijn. Reduction of communication overhead on distributed memory architectures. *LAPACK: Linear Algebraic Programming*. University of Tennessee, 1991.
- [3] T.H. Dunigan. Performance of the Cray T3E. Technical Report 11491, Oak Ridge National Laboratory, 1997.
- [4] G. Fox, M. Johnson, G. Lyzenga, D. A. Ratner, and A. C. Stevens. *Concurrent Processors*, volume 1. Academic Press, 1988.
- [5] C. Moler, J. Little, and R. J. Hanson. *Matrix Computations*, 2nd edition. John Wiley & Sons, 1979.

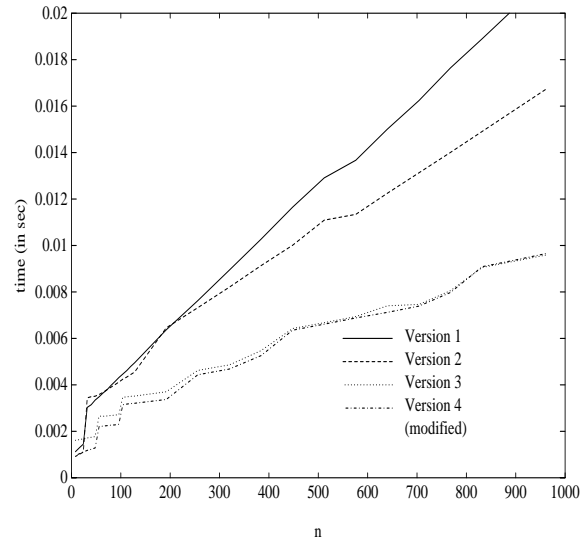
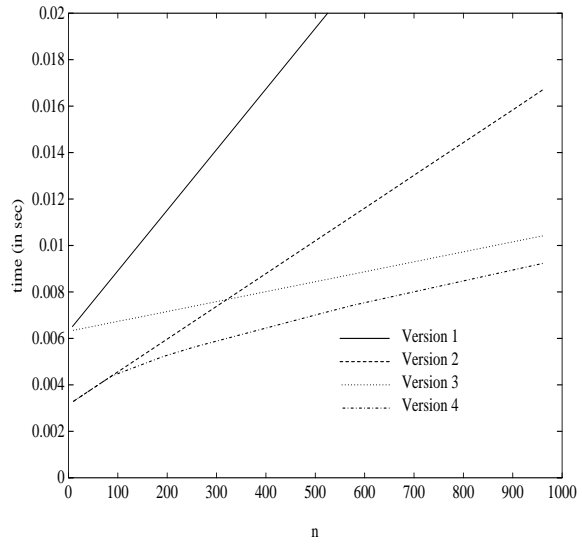


Figure 8: Observed time as a function of vector length n on 64 nodes. In this version, the length of vectors being communicated has not been adjusted, and the algorithm has been modified as described in Section 6, with $\alpha = 480\mu\text{sec}$, $\beta = 2\mu\text{sec}$, and $\gamma = .35\mu\text{sec}$.

communicated is less than 25. While there is a dramatic drop in execution time for Version 2 when $n \leq 25$, Version 3 is particularly aided by the reduction in startup time, since the number of vectors being communicated are halved at each step. For example, on a cube of dimensions 8, if the vector is of length 512, the vectors communicated during the last two steps are of length 8, respectively, thereby reducing the execution time. This demonstrates the importance of the modification in Section 4. The following modification appears to be an optimal choice as described in Section 4. The number of vectors communicated at each stage until the vector length is less than 25, is the same as the number of vectors thereafter, is used. The experimental method shows significant improvement.

7 Conclusion

The implementation of the open source software is available at <https://github.com/robertmiller/parallel-matrix-vector-multiplication>.



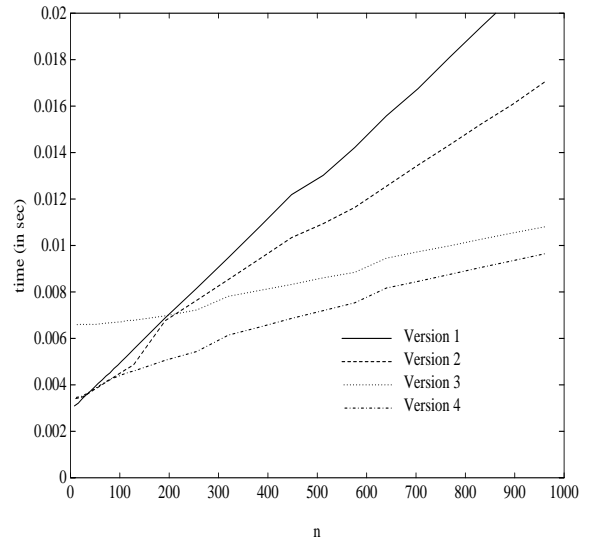
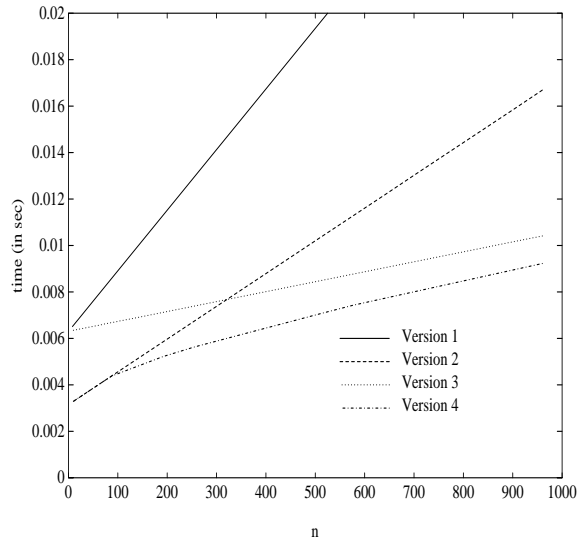


Figure 7: Predicted (*left*) and observed (*right*) time as a function of vector length n on 64 nodes when $\alpha = 525\mu\text{sec}$, $\beta = 2\mu\text{sec}$, and $\gamma = .35\mu\text{sec}$.

6 Experiments on the Intel iPSC/860

To test our theoretical results, we implemented the various global combine operations on the Intel iPSC/860. Our experiments centered around a specific global combine operation, the summation of single precision floating point vectors.

The Intel iPSC/860 is a commercial parallel processor that consists of 64 processors connected in a hypercube topology. Although this machine is somewhat different from the one in Section 2, it can be programmed in such a way that all assumptions in Section 2 are satisfied.

In [3] it is shown that the cost for communicating a floating point number on the

iPSC/860 is roughly given by Assumption 5, with $\alpha = 136\mu\text{sec}$.

The cost for communicating k single precision floating point numbers is $\gamma = .35\mu\text{sec}$.

For $k \leq 25$ or less are communicated, the communication cost is $\beta = 2.0\mu\text{sec}$.

To avoid this complication, we padded all communication with zeros.

For larger k , we padded all communication with zeros.

There is some overhead in the

general bookkeeping, yielding an overhead of $\beta = 2.0\mu\text{sec}$.

The first series of experiments

approaches described in [3]

communication cost

to

```

hybridCOMB2R(n, x, y, d)
begin
  if n <  $\frac{2\alpha}{(d-1)(\beta+\gamma)+\gamma}$  then
    if bit(i,d-1)=1 then
      send (n, x, ngbr(i,d-1))
    else
      recv (n, y, ngbr(i,d-1))
      combine (n, x, y)
      if d-1 > 0 call hybridCOMB2R(n, x, y, d-1)
  else
    let x0 = x[0,...,n/2-1], x1 = x[n/2,...,n]
    if bit(i,d-1)=0 then
      send(n/2, x1, ngbr(i,d-1))
      recv(n/2, y, ngbr(i,d-1))
      combine(n/2, x0, y)
      if d-1>0 call hybridCOMB2R(n/2, x0, y, d-1)
      recv(n/2, x1, ngbr(i,d-1))
    else
      send(n/2, x0, ngbr(i,d-1))
      recv(n/2, y, ngbr(i,d-1))
      combine(n/2, x1, y)
      if d-1>0 call hybridCOMB2R(n/2, x1, y, d-1)
      send(n/2, x1, ngbr(i,d-1))
end

```

Figure 6: Optimal hybrid global combine-to-root routine

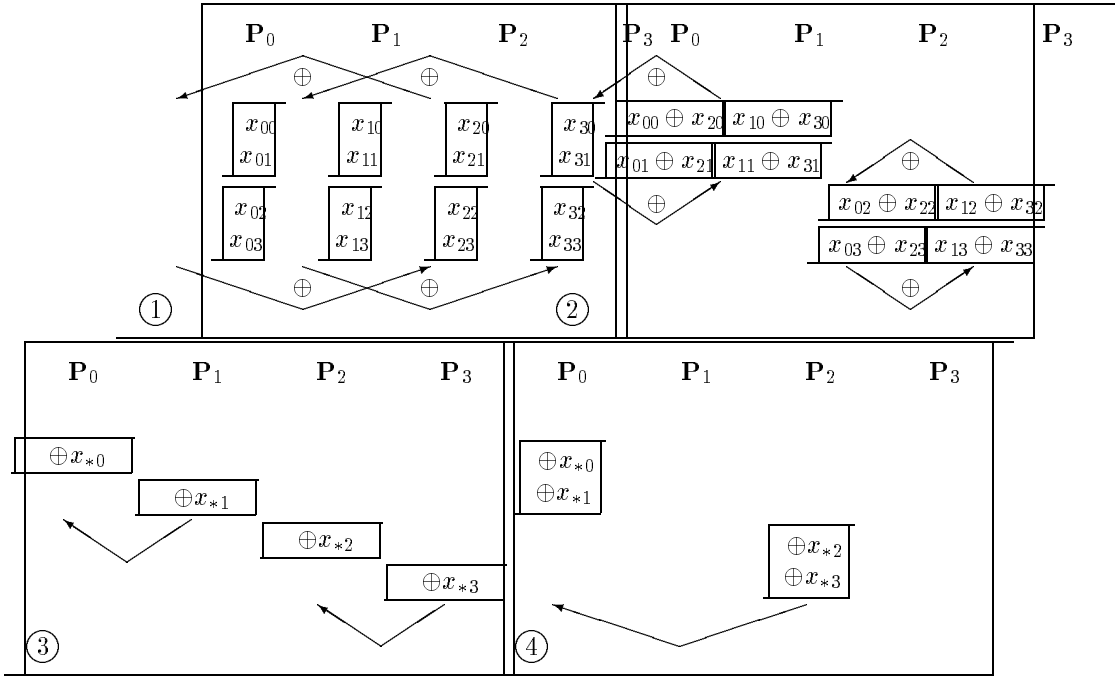


Figure 5: Second approach to combine-to-root on 4 nodes

5 Combine-to-Root

Some algorithms require the result of the global combine to only be known at the root. Without loss of generality, we can take this node to be P_0 .

The first approach pro


```

hybridCOMB(n, x, y, d, S)
begin
  if  $S_{d-1} = 0$ a then
    send (n, x, ngr(i,d-1))
    recv (n, y, ngr(i,d-1))
    combine (n, x, y)
    if d-1 > 0 call hybridCOMB(n, x, y, d-1, S)
  else
    let x0 = x[0,...,n/2-1], x1 = x[n/2,...,n]
    if bit(i,d-1)=0 then
      send(n/2, x1, ngr(i,d-1))
      recv(n/2, y, ngr(i,d-1))
      combine(n/2, x0, y)
      if d-1>0 then call hybridCOMB(n/2, x0, y, d-1, S)
      send(n/2, x0, ngr(i,d-1))
      recv(n/2, x1, ngr(i,d-1))
    else
      send(n/2, x0, ngr(i,d-1))
      recv(n/2, y, ngr(i,d-1))
      combine(n/2, x1, y)
      if d-1>0 call hybridCOMB(n/2, x1, y, d-1, S)
      send(n/2, x1, ngr(i,d-1))
      recv(n/2, x0, ngr(i,d-1))
  end
end

```

^aIn Section 4 it will be shown that an optimal hybrid strategy can be obtained by deleting S from the calling sequence and replacing this condition by

$$n < 2\alpha / ((d-1)(\beta + \gamma) + \gamma)$$

Figure 4: Hybrid global combine routine

Since $T(S, 2^{-(d-j-1)}n, j) > T(S, 2^{-(d-j)}n, j) \geq 0$ and $T(S, n, d) \leq T(R, n, d)$, we conclude that

$$\alpha + 2^{-(d-j-1)}n(\beta + \gamma) < 2\alpha + 2^{-(d-j)}n(2\beta + \gamma)$$

and hence

$$n < 2^{d-j} \frac{\alpha}{\gamma} < 2^{d-k} \frac{\alpha}{\gamma} \leq 2^{d-k} \frac{\alpha}{k(\beta + \gamma) + \gamma}$$

which contradicts the definition of k .

Case 2b: $j \leq k$. Then $S = (S_0, \dots, S_k, 1, \dots, 1)$ and

$$T(S, n, d) = \sum_{i=1}^{d-k-1} (2\alpha + 2^{-k}n(2\beta + \gamma))$$

However,

$$T(\bar{S}, n, d) = \sum_{i=1}^{d-k-1} (2\alpha + 2^{-k}n(2\beta + \gamma))$$

and, by Case 1 above,

$$T(\bar{S}, n, d) > T(S, n, d)$$

which a contradiction.

Proof: The first four results follow immediately from the definition of $T(S, n, d)$. This follows from the observation that the time of the first j steps of the process is the same as the time for the last $d - j$ steps. Hence, when

combines of the partial result on subcubes of dimension $d - 1$. This suggests a family of strategies that combine Versions 2 and 3.

We can generate 2^d separate strategies by considering all possible combinations of the d coordinates. This indicates that for the d -dimensional case, there are 2^d different strategies that can be generated by combining the partial results from the $(d-1)$ -dimensional case.

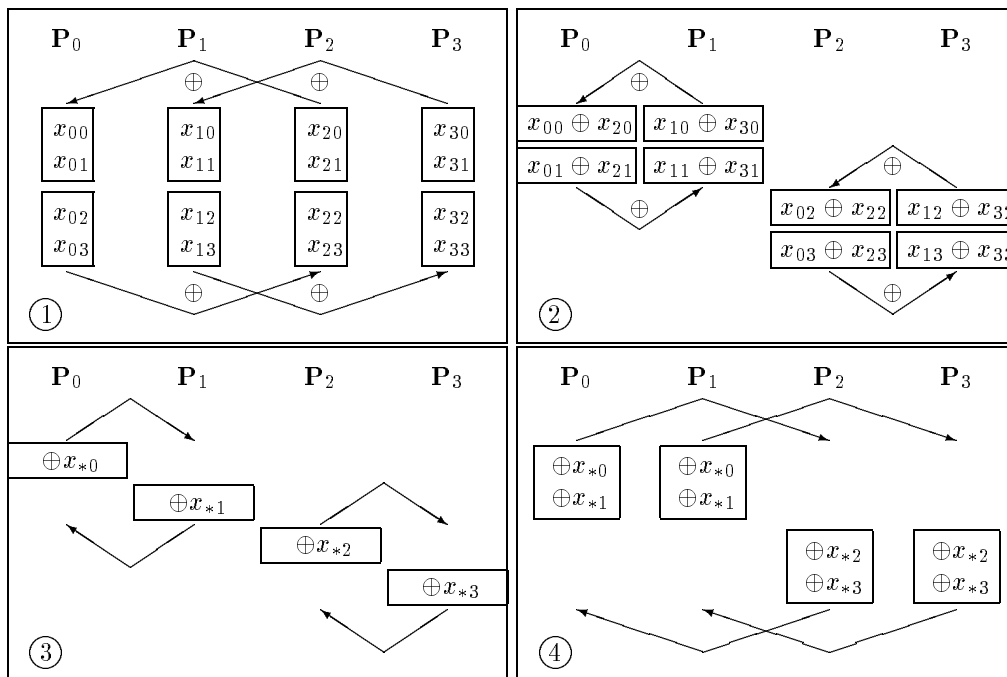


Figure 3: Version 3 on 4 nodes

not sent. This process proceeds for direction $d-2, \dots, 0$, where the size of the vector communicated and combined is halved at each step, resulting in a single vector that results from

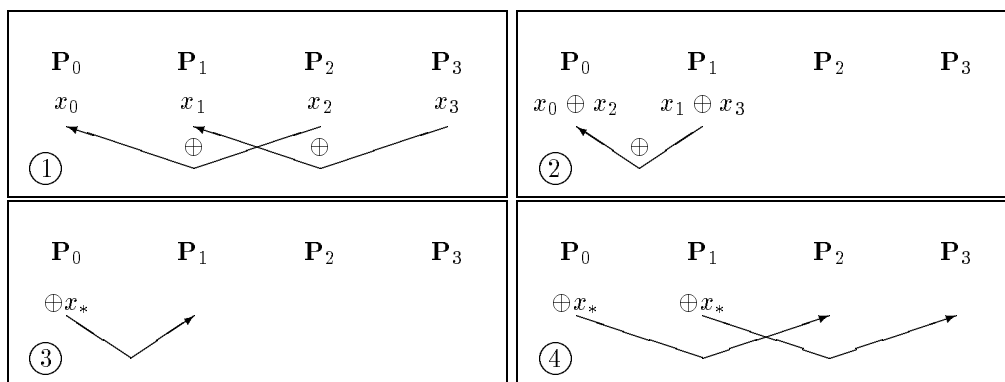


Figure 1: Version 1 on 4 nodes

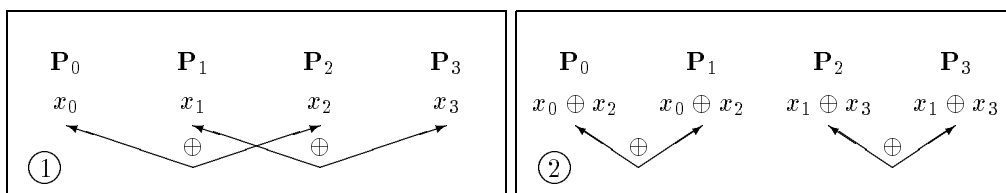


Figure 2: Version 2 on 4 nodes

3.2 Version 2

A second approach, described in [4], Section 14-5.4, starts by
 contents of vector x_i with its neighbor in di
 with the content

2 Assumptions

Target architectures for our algorithm are distributed memory hypercube multiprocessors including Multiple Instruction Multiple Data (MIMD) multiprocessors such as NCUBE2 and Transputer based architectures. For our theoretical

LAPACK Working Note 29

On Global Combine Operations*

Robert A. van de Geijn

Department of Computer Science
University of Tennessee
Knoxville, Tennessee 37996-1301

and

Department of Computer Sciences
University of Texas at Austin
Austin, Texas 78712

na.vandegeijn@na-net.ornl.gov

Abstract

We discuss a hybrid strategy for implementing global combine operations on distributed memory MIMD multi-computers. A theoretical analysis is given and results from its implementation on the Intel iPSC/860 are reported.

1 Introduction

In this paper, we address the implementation of the combine operation when vectors of data to be combined are distributed among the processors (nodes) of a MIMD hybrid computer. Several solutions to this problem have appeared in the literature, including one that combines two of the most common techniques. Our