

# Interprocessor Communication with Limited Memory

Ali Pinar, Bruce Hendrickson

Ali Pinar is with the Computational Research Division of Lawrence Berkeley Lab, Berkeley, CA 94720-8139. E-mail: apinar@lbl.gov. Bruce Hendrickson is with the Discrete Algorithms and Math Sandia National Laboratories, Albuquerque, NM 87185-1110. E-mail: bah@cs.sandia.gov

This work was funded by the Applied Mathematical Sciences program, U.S. Department of Energy, Office of Energy Research and performed at Sandia, a multiprogram laboratory operated by Sandia Corporation, a Lockheed-Martin Company, for the U.S. DOE under contract number DE-AC-94AL85000. The first author is also supported by the Director, Office of Science, Division of Mathematical, Information, and Computational Sciences of the U.S. Department of Energy under contract DE-AC03-76SF00098.

### Abstract

Many parallel applications require periodic redistribution of workloads and associated data. In a distributed memory computer, this redistribution can be difficult if limited memory is available for receiving messages. We propose a model for optimizing the exchange of messages under such circumstances which we call the *minimum phase remapping problem*. We first show that the problem is NP-Complete, and then analyze several methodologies for addressing it. First, we show how the problem can be phrased as an instance of multi-commodity flow. Next, we study a continuous approximation to the problem. We show that this continuous approximation has a solution which requires at most two more phases than the optimal discrete solution, but the question of how to consistently obtain a good discrete solution from the continuous problem remains open. We also devise a simple and practical approximation algorithm for the problem with a bound of 1.5 times the optimal number of phases. We also present an empirical study of variations of our algorithms which indicate that our approaches are quite practical.

### Keywords

Interprocessor communication, Dynamic load balancing, Data migration, Scheduling, NP-Completeness, Approximation algorithms.

## I. INTRODUCTION

In many parallel computations, the workload needs to be periodically redistributed among the processors. When computational work varies over time, the tasks must be re-assigned to keep the workload balanced. On a distributed memory computer, this generally requires that data structures associated with the computations be transferred between processors. Many examples of this phenomena occur in scientific computing, including adaptive mesh refinement, particle simulations with short- or long-range forces, state-dependent physics models, and multi-physics or multi-phase simulations. For many such simulations, the limiting resource is not computation but memory. Important examples include differential equation solvers using adaptive meshes, and large scale particle simulations. Scientists commonly choose to use the minimum number of processors upon which a particular simulation can fit, or they choose a problem size which fills the memory of a particular number of processors. However, the dynamic memory requirements of such applications make such targeting difficult.

A number of algorithms and software tools have been developed to repartition the work among processors (see for example, [1], [2] and references therein). However, the mechanics

of actually moving large amounts of data has received much less attention. When the processors have sufficient memory, the simplest way to transmit the data is quite effective. Each processor can execute the following steps.

- (1) Allocate space for my incoming data.
- (2) Post an asynchronous receive for my incoming data.
- (3) Barrier.
- (4) Send all my outgoing data.
- (5) Free up space consumed by my outgoing data.
- (6) Wait for all my incoming data to arrive.

The barrier in step (3) ensures that no messages arrive until the processor is ready to receive them, so no buffering is needed.

Unfortunately, this protocol can fail when memory is limited. It requires a processor to have sufficient memory to hold both the outgoing and the incoming data simultaneously, since incoming messages can arrive before space for outgoing data is freed. An alternative way to view this issue is that for a period of time, the data being transferred consumes space on both the sending and receiving processors. A protocol that alleviates this problem is desirable for three reasons. First, since many scientific calculations are memory limited, reserving space for this communication operation limits the size of the calculations that can be performed. Second, the amount of memory required by this protocol is unpredictable because the data remapping requirements depend upon the computation. Thus, one is forced to set aside a conservative amount of space, which is likely to be wasteful. Third, a general purpose tool for dynamic load balancing should be robust in the presence of limited memory. On the hopefully infrequent occasions when memory limitations prohibit direct transfers, a good tool should apply an alternative strategy instead of exiting. The construction of such a tool (Zoltan [3]) inspired our interest in this problem. The desire to impact Zoltan has a number of implications for this research. We are interested in algorithms with attractive practical performance, not just asymptotic bounds on worst-case behavior. We also need algorithms that are straightforward and efficient to implement.

To address the deficiencies of the standard protocol, we propose a simple modification. Instead of sending all of the data at once, we will send it in phases. After each phase,

processors can free up the memory of the data they have sent. That memory is now available for the next communication phase. Each phase adds an overhead due to the latency costs, thus it is important to limit the total number of phases, which is an approximation to minimize total remapping time. The results in §VII will confirm that the overall communication time is strongly effected by the number of phases.

More formally, consider a set of  $P$  processors. The amount of data that needs to be communicated between processors is a *transfer request*. We will assume that the request is *feasible*, that is the end result of satisfying the transfer request does not violate any processor's memory constraints. We will let  $T_{ij}$  denote the total volume of data to be transferred from processor  $i$  to processor  $j$ . We now wish to perform the requested transfer in a sequence of phases. Let  $t_{ij}^l$  denote the volume of data transfer from processor  $i$  to processor  $j$  in phase  $l$ , and let  $A_i^l$  be the memory available to processor  $i$  at the beginning of phase  $l$ . We will also use  $R_i^l$  and  $S_i^l$  to denote the total volume of data received and sent by processor  $i$  in phase  $l$  (i.e.,  $R_i^l = \sum_{j=1}^P t_{ji}^l$  and  $S_i^l = \sum_{j=1}^P t_{ij}^l$ ). At each step the finite-memory constraint requires that  $R_i^l \leq A_i^l$  for  $i = 1, 2, \dots, P$ . The available memory after each phase can be computed as  $A_i^{l+1} = A_i^l + S_i^l - R_i^l$ . Our objective is to find a schedule of transfers that obeys the memory constraint and satisfies the transfer request in a minimal number of phases. We will call this the *minimum phase remapping problem*.

In §II, we show that the problem of determining whether a given transfer can be completed in a specified number of phases is NP-Complete. The remainder of this paper focuses on formulations and approximation algorithms that could be used in practice. In §III we present a reduction of our problem to multi-commodity flow. We present a continuous relaxation of the problem in §IV, and a practical approximation algorithm in §V. Although the emphasis of this paper is a theoretical analysis, we discuss some practical issues associated with the application of our techniques in §VI, and results of empirical studies in §VII. Earlier versions of some portions of this paper have appeared in conference proceedings [4], [5].

Despite its practical importance, efficient data transfers has not been well-studied. Some standard collective communication operations can be implemented in ways that limit memory usage, but the general problem we consider seems to be new. Cypher and Konstan-

tinidou designed memory efficient message passing protocols [6]. However, their work addressed exchange of tokens as opposed to variable sized messages, and they did not explicitly consider the effect of finite memory in the processors. Their work conceptually divides a process into communication and application processes. Communication processes receive unit-size messages and copy them to application processes. It is assumed that application processes have enough memory, and the goal is to limit the memory requirement of the communication processes. Very recently, Hall *et al.* studied the problem for large storage systems [7]. They assume the total time for subdividing a message and sending subdivided messages is about the same as sending the entire message. So, they study the problem on unit-size messages. This assumption might hold for huge volumes of data to be transferred, as in the case of reorganizing a database, however it is far from being true for dynamic load balancing applications due to very high message setup times. Their model also restricts each processor to send and receive only one message in a phase.

## II. COMPLEXITY

In this section we show that determining whether a given transfer can be completed in a specified number of phases is NP-Complete. Our proof uses a reduction from the Hamiltonian Circuit problem, which is known to be NP-Complete [8]. Recall that a Hamiltonian Circuit is a cycle in a graph that visits each vertex once. Given an instance of the Hamiltonian Circuit problem, the basic idea of our reduction is to construct an instance of the data transfer problem in which there is but a single unit of usable memory. This unit is a *token* that is passed between processors, and possession of the token allows a processor to receive data in the next phase. A solution to the data remapping problem occurs if and only if the token can be passed in a cycle among all the processors, which implies the existence of a Hamiltonian Circuit. To see how this can be done, consider the Hamiltonian Circuit problem posed in the left portion of Fig. 1. From this instance, we construct the data remapping problem in the right portion of the figure. The data remapping problem contains the original graph as its *core* (represented in the figure with dark lines) after replacing vertices with processors and replacing edges with unit-volume data transfers. It also contains a *chain* of processors to the left. The bottom processor in this chain has free memory that will percolate upwards with each phase, finally allowing

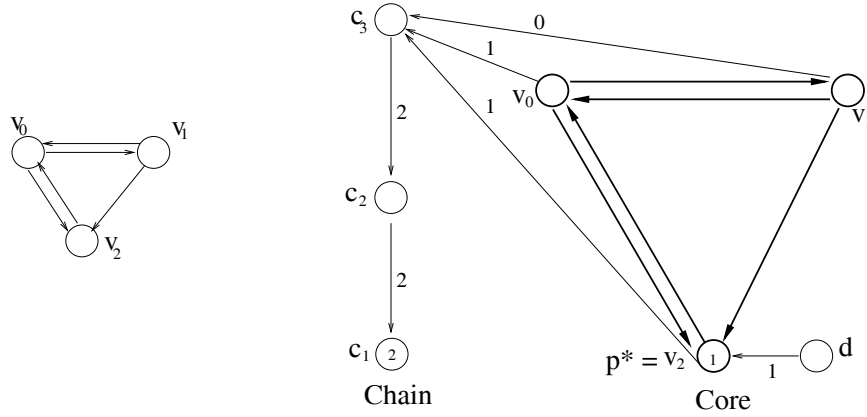


Fig. 1. Construction for NP-Completeness proof.

all the data transfers to be completed. Given a graph  $G = (V, E)$ , we construct a data remapping problem as follows.

- For each vertex  $v_i$  in  $V$  there is a processor  $p_i$ . We refer to these processors as core processors. Each edge of  $E$  is a unit-volume transfer.
- Add a chain of  $|V|$  processors  $c_1, \dots, c_{|V|}$  along with transfer requests from  $c_{i+1}$  to  $c_i$  for  $i = 1, 2, \dots, |V| - 1$ , each with volume  $|E| - |V|$ .
- Add a transfer request from each core processor  $p_i$  to the top of the chain  $c_{|V|}$ . This transfer has volume equal to one less than the in-degree of  $v_i$  in  $G$ .
- Add a dummy processor  $d$  and a unit-weight transfer connecting  $d$  to an arbitrary processor  $p^*$  in the core.
- Give  $|E| - |V|$  units of free memory to processor  $c_1$  and 1 unit of free memory to  $p^*$ . No other processor has free memory.

Consider what happens as the data remapping occurs. In the first phase,  $c_2$  will send its data to  $c_1$ , moving the free memory one step up in the chain. After  $|V| - 1$  phases, this free memory will have arrived at  $c_{|V|}$ , the top of the chain. Meanwhile the token, which started at  $p^*$ , will have meandered about, enabling some data to be transferred. In phase  $|V|$ , processor  $c_{|V|}$  has enough memory to receive all it needs from the core processors. During this phase, the token can take one more step. The messages sent to  $c_{|V|}$  free up memory in the core processors. Specifically, at the end of phase  $|V|$ , each core processor  $p_i$  has  $\text{indegree}(p_i) - 1$  units of free memory (one processor has an additional unit due to the token). In phase  $|V| + 1$ , core processor  $p_i$  can now receive all it needs, minus 1. The

transfers to  $p_i$  can be completed in this phase if and only if one of the data transfers to  $p_i$  has previously been handled by the token. But the only way for the token to visit all the core processors in  $|V|$  phases is to complete a Hamiltonian Circuit of the core graph. Note that the token must end up where it started, at processor  $p^*$ , to enable the transfer from  $d$  to occur at phase  $|V| + 1$ . This argument leads to the following result.

*Theorem 1:* Determining whether an instance of the data remapping problem can complete in a specified number of phases is NP-Complete.

*Proof:* Given an instance of the Hamiltonian Circuit problem  $G = (V, E)$ , construct a data remapping problem as described above. As sketched above, the data remapping problem finishes in  $|V| + 1$  phases if the core graph has a Hamiltonian Circuit. If the core graph does not have a Hamiltonian Circuit, then one of its processors will not have been visited by the token by the end of phase  $|V|$ . That unvisited processor,  $p_i$ , still needs to receive  $\text{indegree}(v_i)$  data, but has only  $\text{indegree}(v_i) - 1$  units of available memory, so the data transfers cannot complete in  $|V| + 1$  phases. The construction of the data remapping problem is polynomial, so we can conclude that the data remapping problem is NP-Hard. A given solution can be verified in polynomial time, thus the problem is in NP. ■

### III. MULTI-COMMODITY FLOW FORMULATION

In this section, we present a multi-commodity flow (MCF) formulation to determine whether a given transfer can complete in a specified number of phases [9]. Once we can solve the decision problem, the number of phases in an optimal solution can be determined using parametric search. This formulation enables use of MCF technology to solve the minimum phase data remapping problem optimally. This might be helpful for three reasons. First, some MCF problems can be solved relatively quickly, despite their intractability in the general case. Second, the continuous version of the MCF problem can be solved in polynomial time and the solution can be used as a heuristic for the integer problem. Finally, MCF solvers will find an optimal solution if runtime is not an issue.

In our MCF formulation, each processor corresponds to a commodity. Let  $P$  be the number of processors and  $L$  be the number of phases. We must decide if a remapping can complete in  $L$  phases. As depicted in Fig. 2, our MCF graph contains a sequence of components, one for each phase. Each component allows for the communication that

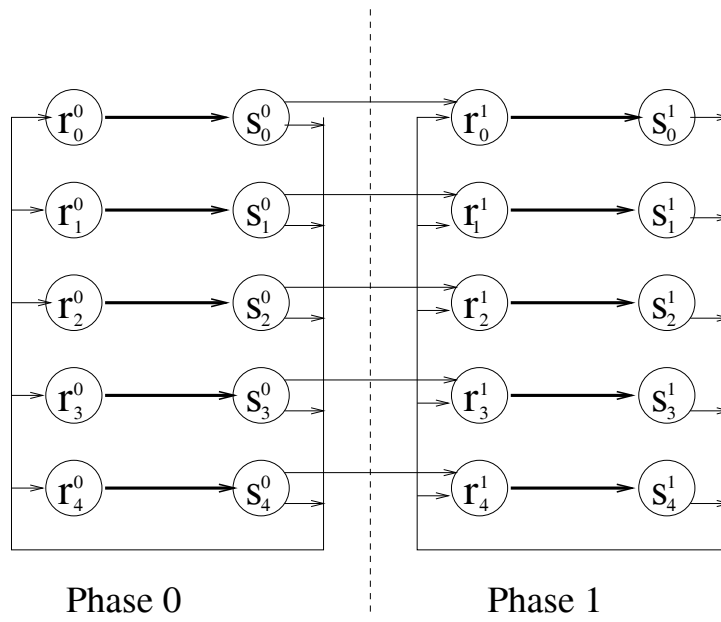


Fig. 2. MCF graph for 5 processors and 2 phases.

occurs in the corresponding phase.

The MCF graph  $G = (V, E)$  has  $2PL$  vertices. Each processor is represented by  $2L$  vertices: two processors (one sender and one receiver) at each phase. We will use  $r_i^l$  and  $s_i^l$  to denote receiver and sender, respectively, for processor  $i$  in phase  $l$ . A sender vertex of the first phase is the source of a commodity with volume equal to the total volume of the data originally stored by this processor. A receiver vertex in the last phase is a destination for a set of commodities that corresponds to data to be stored by this processor after remapping is complete.

In the MCF graph, there is an edge from  $r_i^l$  to  $s_i^l$  for  $l = 1, \dots, L$  and  $i = 1, \dots, P$ . The capacity of an edge is equal to the total memory on the respective processor. There are also edges from each sender vertex  $s_i^l$  to all other receiver vertices  $r_j^l$  in the same phase to enable data exchange between any pair of processors in a phase. These edges have infinite capacities.

With this construction, all processors first receive the data in a phase and then send their messages. This corresponds to first allocating space for the data to be received and then sending the outgoing data. The edges from receivers to senders within a phase guarantee that there is sufficient space to allocate memory for the incoming data before



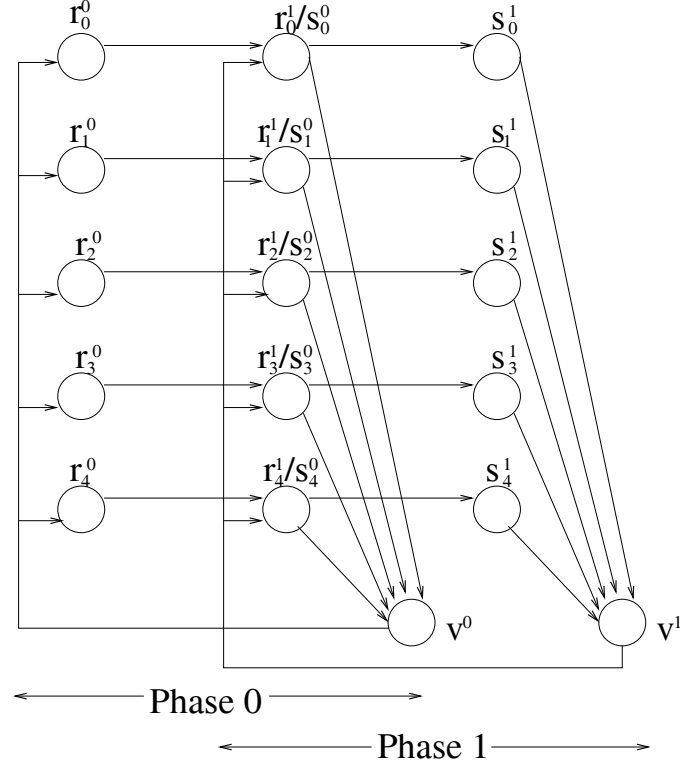


Fig. 3. MCF graph after reduction.

releasing the space for the data being shipped out, so that the memory constraints are guaranteed to be satisfied.

Finally, there is an edge (with infinite capacity) from each sender  $s_i^l$  to the receiver in the next phase  $r_i^l$  for  $l = 1, \dots, L - 1$ . The flow on these edges corresponds to data that is already in the memory of a processor at the beginning of a phase. The graph for  $P = 5$  and  $L = 2$  is depicted in Fig. 2.

*Theorem 2:* There exists a solution to the remapping problem if and only if there exists a solution to the MCF formulation.

*Proof:* We can replace a data transfer from processor  $i$  to processor  $j$  in phase  $l$  with flow on edge  $(s_i^l, r_j^l)$  of equal volume. As argued above, memory constraints on the processors are satisfied if and only if the capacity constraints on the edges are satisfied in  $G$ . Thus, the feasibility of one solution implies the feasibility of the other. ■

In this formulation the number of commodities is equal to the number of processors, and the graph has  $2PL$  vertices and  $P^2L$  edges. The number of vertices and edges can

be reduced for a more efficient formulation. First, we can replace the crossbar between senders and receivers in a phase  $l$  with a vertex  $v^l$  and edges from all senders of phase  $l$  to  $v^l$  and edges from  $v^l$  to all receivers of phase  $l$ . Second, we can merge the senders of phase  $l$  with receivers of phase  $l+1$ . The graph after these reductions is depicted in Fig. 3. This improved formulation has  $PL + L + P$  vertices and  $(3L + 1)P$  edges.

#### IV. CONTINUOUS RELAXATION

Although the multi-commodity flow formulation from §III provides a methodology for solving instances of the minimum phase remapping problem, runtime can still be exponential in the problem size. In this section, we describe an efficient solution for an approximation to the remapping problem. In the approximation, integral constraints on the volume of data transfers are relaxed to allow continuous values. Naturally, the volume of transfer between two processors in a phase must be an integer. But integer solutions near the continuous ones can be used as heuristics. Note that the unit of data transfer is only a byte, whereas the volume of data being transferred is often on the order of megabytes. So, conversion from a continuous solution to an integer solution will often be a small perturbation, and so heuristics based upon this idea may be generally effective. However, bad cases for this heuristic exist, as discussed at the end of this section.

As defined in the introduction,  $T_{ij}$  denotes the total volume of data to be communicated from processor  $i$  to processor  $j$ , and  $t_{ij}^l$  denotes the volume of data transferred from processor  $i$  to processor  $j$  in phase  $l$ . The memory available to processor  $i$  at the beginning of phase  $l$  is denoted by  $A_i^l$ . We also use  $R_i$  and  $S_i$  to denote the total volume of data received and sent by processor  $i$  during remapping.

Let  $T$  be the total volume of data to be transferred, and  $M$  be the total volume of available memory in the system, then  $L = \lceil \frac{T}{M} \rceil$  is a lower bound on the number of phases. We will divide each message into  $L$  equal pieces, i.e.,  $t_{ij}^0 = t_{ij}^1 = \dots = t_{ij}^{L-1} = \frac{T_{ij}}{L}$  and send a piece at each phase. If the memory constraints are satisfied, then the data transfers will complete in precisely  $L$  phases. However, there is no guarantee that memory constraints will not be violated. As a solution to this problem, we will use preprocessing and postprocessing phases to ensure feasibility of the phases in between.

*Lemma 1:* If the following conditions are satisfied, the continuous version of the remap-

ping problem can be completed in  $L = \lceil \frac{T}{M} \rceil$  phases.

(i)  $S_i = R_i$  for all processors.

(ii)  $A_i^0 \geq \frac{R_i}{L}$ .

*Proof:* At each phase processor  $i$  will receive  $\frac{R_i}{L}$  units of data. By the second condition, each processor has sufficient memory for the first phase. By the first condition, each processor ships out  $\frac{S_i}{L} = \frac{R_i}{L}$  units of data at each phase, which frees up sufficient memory for the next phase. ■

*Lemma 2:* A solution for a continuous version of the data remapping problem for transfer request  $\mathcal{R}$  can be performed via the following three steps.

1. one preprocessing phase.
2. a new transfer request  $\mathcal{R}'$  where  $S_i = R_i$  and  $A_i^0 \geq \frac{R_i}{L}$ .
3. one postprocessing phase.

*Proof:* In the preprocessing phase we will reorganize the data to satisfy conditions (i) and (ii) from Lemma 1, and define a new mapping of the data. After the new mapping is complete, a single postprocessing phase will be sufficient to get all of the data to the correct processor.

In the preprocessing step, all processors  $i$  with  $R_i < S_i$  will transfer some of their outgoing data to processors  $j$  in which  $R_j > S_j$ , so that in subsequent phases  $R_i = S_i$ . Note that if the transfer request is feasible, then  $R_j - S_j \leq A_j^0$ . Thus, this rearrangement can be completed in a single phase.

Next, as second part of the preprocessing step, all processors  $i$  with  $A_i < \frac{R_i}{L}$  will transfer some of their outgoing data to processors  $j$  with  $A_j > \frac{R_j}{L}$ . To avoid disturbing the first property, the sending processors will also pass equal amounts of receiving assignment. Once again, this step can be completed in one phase, since by construction, the receiving processors have sufficient space.

Notice that, the actual data being transferred is irrelevant—we are just trying to balance the numbers—so a send and receive operation can cancel each other. This enables merging of the two steps above into one phase.

After the new transfer request  $\mathcal{R}'$  is realized, we must correct for the transfer of receiving assignments. This correction is the purpose of the postprocessing phase. Under the

transfer of receiving assignments, each processor is either a sender or a receiver of such assignments. So, during postprocessing, each processor will either receive or send data, but not both. Since the initial remapping is feasible, each processor has enough memory for the data to be received, thus the postprocessing can be completed in one phase. ■

The complexity of constructing the solution for the preprocessing phase is linear in the number of processors. To see this, divide the processors into two lists: those with  $R_i < S_i$  and those  $R_j > S_j$ . Now step through the lists together, transferring sending responsibility from a processor in the  $i$  list to one in the  $j$  list. Each transfer balances  $R_i$  and  $S_i$  for a processor in one of the lists. The same can be applied to balance initial available memories. Notice that the preprocessing step uniquely describes the postprocessing phase, and remapping for  $\mathcal{R}'$  is straightforward.

*Theorem 3:* Given a transfer request  $\mathcal{R}$ , the continuous version of the data remapping problem can be completed in  $\lceil \frac{T}{M} \rceil + 2$  phases.

*Proof:* By Lemma 2,  $\mathcal{R}$  can be completed by pre- and postprocessing steps, along with a transfer request  $\mathcal{R}'$  satisfying conditions of Lemma 1. Notice that the total volume of data to be transferred,  $T'$  in  $\mathcal{R}'$  is no greater than  $T$  in  $\mathcal{R}$ , and the total available memory in the system does not change:  $M = M'$ . Hence by Lemma 1,  $\mathcal{R}'$  can be completed in  $\lceil \frac{T'}{M} \rceil \leq \lceil \frac{T}{M} \rceil$  phases. Together with one preprocessing and one postprocessing phases, remapping can be completed in  $\lceil \frac{T}{M} \rceil + 2$  phases. ■

It is worth noting that a good solution for this continuous approximation may not yield a good solution for the true discrete problem. For instance, consider the example depicted in Fig. 4.



Fig. 4. Catastrophic instance for continuous relaxation.

This example consists of two groups of processors, with no communication between the groups, and there is only one unit of available memory. Available memory must be pos-

sessed by each component in turn, and this requires temporarily moving some data from one component to the other to transfer the free memory, as will be discussed in more detail in the next section. In the preprocessing step described in the proof of Lemma 2, this available memory will be divided into two groups of processors, but the fractional transfers that follow give no insight into the correct way to orchestrate the data transfers for this instance. Specifically, in the continuous solution all processors are identical, so no information is gleaned about the necessity of working on components in turn.

## V. EFFICIENT APPROXIMATION ALGORITHMS

In this section, we describe the basics of a family of efficient algorithms that provide solutions, in which the number of phases is at most 1.5 times that of an optimal solution. The algorithm is motivated by some simple observations. First, the maximum amount of data that can be transferred in a phase is equal to the total amount of free memory in the parallel machine. Let  $M$  be the total available memory in the parallel machine, and let  $T$  be the total volume of data to be moved. Note that  $M$  doesn't change between phases.

*Lemma 3:* The minimum number of phases in a solution is  $\lceil \frac{T}{M} \rceil$ .

This bound can only be achieved only if available memory is used to receive messages at each phase. Thus, free memory is wasted if it resides on a processor that has no data to receive. Our algorithm works by redistributing free memory to processors that can use it. Equivalently, data is *parked* on a processor with free memory, which it cannot use, to free up memory on processors that can use it. We will park only data that must be transferred eventually.

### A. Parking

Parking aims to utilize memory that would otherwise be wasted. Consider a processor that has received all its data and still has available memory. This memory cannot be utilized in subsequent phases, which decreases the total memory usable for communication, thus potentially increasing the number of phases. Instead, another processor can temporarily move some of its data to this processor to free up space for messages. An example is illustrated in Fig. 5. In this simple example, the top two processors want to exchange 100 units of data, but each has only one unit of available memory. A simplis-

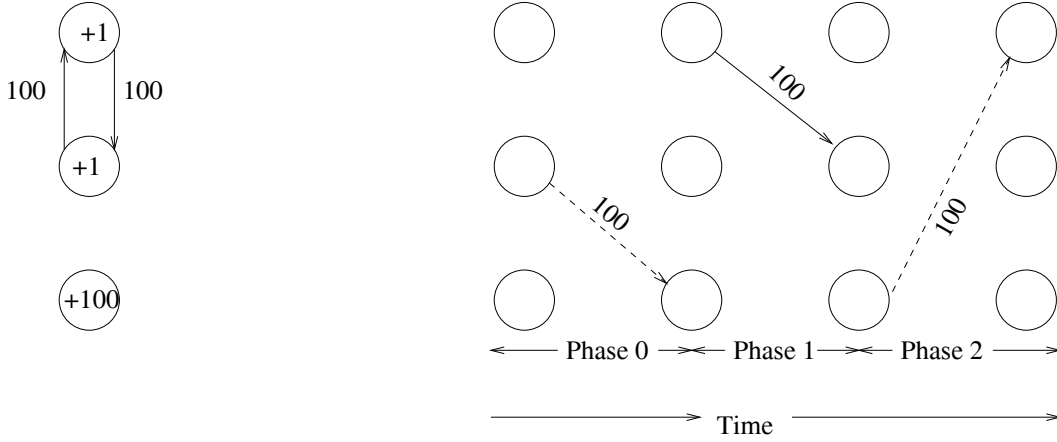


Fig. 5. Example of the utility of parking.

tic approach will require 100 phases. However, the third processor has 100 units of free memory. By *parking* data on this third processor (i.e. transferring free memory to another processor), the number of phases can be reduced to three.

More formally, if a processor has  $k$  units of data left to receive and  $m$  units of free memory, then it has *parking space* of  $\max(0, m - k)$  units. A processor has data to park if the incoming data exceeds available memory, and the quantity of this parkable data is  $\max(0, k - m)$  units. The parkable data consists of data that eventually must be sent to another processor. Note that if the transfer request is feasible, then a processor must send out  $\max(0, k - m)$  units. Any processor that has parking space can store parkable data from another processor, thereby maximizing the amount of usable free memory. This parked data merely takes an extra step on the way to its final destination. Exploiting this observation will allow us to construct an approximation algorithm.

In our algorithm, we merely store data in a parking space, and then forward it to its correct destination when the destination processor has available memory. Note that it is inconsequential which processor owns the parked data. In other words, parking spaces are indistinguishable. What potentially effects performance is which processors shunt their data to a parking space.

*Lemma 4:* It is sufficient to park data at most once to obtain an optimal solution.

*Proof:* Assume there is a solution that parks some data  $D$  twice. Let  $p_1$  and  $p_2$  be the first and second processors on which  $D$  is parked. After data is moved from  $p_1$  to  $p_2$ , if

no other processor uses available memory at  $p_1$ , then there was never a need to move data to  $p_2$ . If another processor  $p_i$ , parks data on  $p_1$ , then we can rearrange the data movement with  $D$  staying in  $p_1$ , and  $p_i$  parking on  $p_2$ , due to indistinguishability of parking spaces. ■

It is worth noting that parking is not just a heuristic but a requirement in some cases. Consider the example in Fig. 5, modified so that there is no available memory in the top two processors. In this case, the transfer request is still feasible, but realizing the remapping requires parking.

### B. An Approximation Algorithm

In this section, we describe an algorithm that obtains a solution with at most 1.5 times the optimal number of phases. The algorithm is quite generic and allows for a number of possible enhancements.

#### *Algorithm 1:*

- A processor receives as much data as it can in each phase (i.e., if a processor has available memory at the end of a phase, then this processor does not have any more data to receive).
- If the transfer request cannot be completed in the next phase, then park as much data as possible (i.e., park the minimum of the total parkable data and the total parking space).

Note that many details about the algorithm are unspecified: If I have more incoming data than free memory, which messages should I receive in the current phase? If several processors want to park data, but limited parking spaces are available, which should succeed? We will show below that with any answers to these questions, the resulting algorithm generates a solution with no more than 1.5 times the optimal number of phases. Intelligent answers to these questions could be used to devise algorithms with better practical (or perhaps theoretical) performance.

*Lemma 5:* The total volume of data transferred by Algorithm 1 is at most  $\lceil \frac{3T}{2} \rceil$ .

*Proof:* Let  $T_p$  be the volume of data transferred through parking, and let  $T_d$  be the data transferred directly. Data is transferred either directly or through parking, thus  $T = T_p + T_d$ .

It is enough to park data once due to Lemma 4, thus parked data is moved twice, and

the total volume of data moved is  $2T_p + T_d = T + T_p$ . Recall that by definition parkable data, data is parked only if it will help receiving more data in the next phase. Thus each parked unit of data enables at least one direct transfer, the algorithm guarantees that  $T_p \leq T_d$ . Thus at most half of  $T$  can be transferred through parking, i.e.,  $T_p \leq \frac{T}{2}$ , and the total volume of data moved is  $T + T_p \leq T + \frac{T}{2} = \frac{3T}{2}$ . ■

*Theorem 4:* Algorithm 1 constructs a solution with at most  $\lceil \frac{3T}{2M} \rceil + 1$  phases.

*Proof:* The algorithm makes use of all  $M$  units of available memory until the amount of parkable data is less than the amount of parking space. It then completes in at most two additional phases, one in which some data is parked, and a final phase in which each processor has enough memory to receive all its messages. By Lemma 5 we know that the total volume of data transferred in the algorithm is at most  $\lceil \frac{3T}{2} \rceil$ . With  $M$  units of transfer in all, but the last two phases, the process can be completed in at most  $\lceil \frac{3T}{2M} \rceil + 2$  phases.

We will now decrease the bound to  $\lceil \frac{3T}{2M} \rceil + 1$ . Let  $l$  be the number of phases for the algorithm to complete the data remapping process. The total volume of data transferred is  $(l-2)M + x$ , where  $x$  is the volume of data transferred in the last two phases:  $1 < x \leq 2M$ . From Lemma 5 we know that

$$\lceil \frac{(l-2)M + x}{M} \rceil \leq \lceil \frac{3T}{2M} \rceil.$$

But simple algebra reveals that

$$l - 1 \leq \lceil \frac{(l-2)M + x}{M} \rceil.$$

Combining these inequalities

$$l - 1 \leq \lceil \frac{3T}{2M} \rceil,$$

and the result follows. ■

Combined with Lemma 3, Theorem 4 shows that Algorithm 1 is a  $3/2$  approximation algorithm for the minimum phase remapping problem. Without a better lower bound, this value of  $3/2$  is tight as illustrated by the example in Fig. 6.

This example consists of an odd number of processors  $P$ . All but one of them are organized in pairs which exchange a single unit of data. Only the unpaired processor has a single



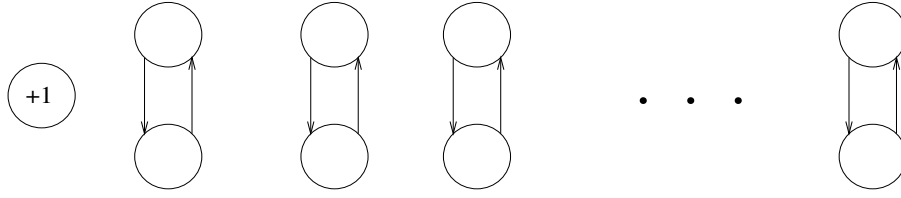


Fig. 6. Example to show the tightness of the 1.5 bound.

unit of available memory. The total volume of data to be moved is  $T = P - 1$ . The only way for a pair to exchange their data is first to park a unit elsewhere, so a total of  $\frac{P-1}{2}$  units of parking are needed. Hence, the total volume of data transferred is  $\frac{3(P-1)}{2} = \frac{3T}{2}$ , and the number of phases is  $\frac{3T}{2M}$ , since  $M = 1$ .

## VI. PRACTICAL CONSIDERATIONS

In this section we consider some more practical aspects of the minimum phase remapping problem, setting the stage for experimental results in the subsequent section. First we discuss practical aspects of algorithms for scheduling the communication phases, and efficient implementations of heuristics based on the templates of the preceding sections. Then we consider how to implement parking, and finally we address some practical limitations of our model and how to address them.

### A. Distributed and Centralized Heuristics

Algorithms for minimizing the number of phases can be either centralized or distributed. Centralized heuristics reduce the problem onto one processor, construct a solution on this processor, and then broadcast the solution to all processors. Reducing the whole problem onto a single processor enables the use of sophisticated algorithms to construct optimal or near optimal solutions. However, these methods invariably suffer from the gather and scatter overhead due to collecting the information on one processor and broadcasting the solution. Moreover, the remaining processors wait idle while one processor is constructing a schedule, which hinders scalability.

An alternative is a distributed method, where processors collectively produce a solution using local information. A distributed approach is more scalable, but the lack of global information is likely to lead to lower quality solutions. We present experiments that

**While** I have data to send or receive **do**

Apportion available memory to senders.

Post an asynchronous receive for each allotment.

Inform senders their allotments and collect my allotments.

**for** each sending allotment I am given

Construct and send message.

Free up the memory for data sent.

Fig. 7. Template of a distributed algorithm for memory-constrained interprocessor communication

compare the performance of the two approaches in the next section. Here, we discuss how to implement a distributed scheduling algorithm efficiently. Fig.7 presents a distributed algorithm for interprocessor communication with limited memory. Notice that there are no global operations in this algorithm; processors communicate only with their neighbors – processors they need to exchange data with. This locality is important for scalability.

The number of phases is determined by how available memory is apportioned to the processors, which is not specified in Fig. 7. Below, we describe five heuristics that rely only on local information. No extra communication is required to make apportionment decisions. Four of our heuristics are non-preemptive, that is we chose a processor and grant it as much memory as it requires or is available, and repeat the process while there is available memory. We experimented with the following four criteria for choosing a sender.

- *First-fit* (FF): Choose the first processor on the list of senders.
- *Max-first* (MF): Choose the processor that has the most to send.
- *Random* (RD): Choose a processor uniformly at random.
- *Fortune Wheel* (FW): Choose a processor randomly, but the probability of a processor being chosen is directly proportional to how much it has to send.

We also have a preemptive heuristic, where we grant allotments so that the maximum remaining volume to be sent by any sender is minimized. For example if processor 1 has 5 units whereas processor 2 has 4 units to send, we grant 5 units of available memory as 3 units to processor 1 and 2 units to processor 2, which leaves 3 units to send for

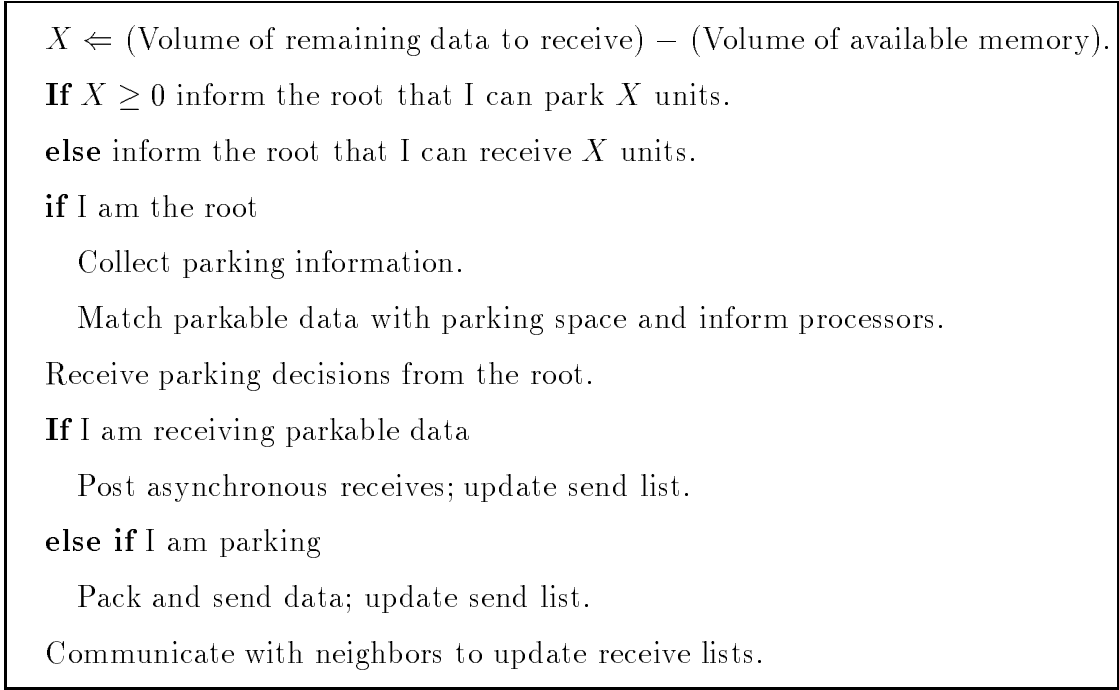


Fig. 8. Template for parking

each processor. We refer to this heuristic as the *Min-max* (MM) heuristic. We compare practical performances of these heuristics in the experimental results section.

### B. Parking

Parking helps to decrease the number of phases by utilizing available memory on processors that do not have more to receive. However, the parking operation itself can be expensive for several reasons. First, the total volume of communication increases, since data being parked are moved twice. More importantly, parking is a global operation: any processor needing to receive can park data on any other processor with available memory, and so some kind of global oversight and control is necessary.

Implementation of a parking algorithm requires three steps: making parking decisions, moving parked data, and updating send/receive information. An algorithm that follows these steps is presented in Fig. 8. This algorithm uses a root processor to make the parking decisions. Moving data is straightforward, since processors already know to whom to send and from whom to receive. Finally, we can update receive lists by communicating with the neighbor processors. As can be seen, parking is not a simple operation, and as will be

shown in §VII, its overhead can not always be compensated for by decrease in the number of phases.

### *C. Limitations of the Model*

The preceding sections of this paper have concentrated on the amount of available memory, ignoring how that memory is organized. However, the adaptive applications that require periodic data redistribution typically make extensive use of dynamic memory management. As a consequence, the available memory may not be contiguous, which limits the immediate practical utility of our techniques. In the following discussion, we assume that the message passing software requires outgoing/incoming messages to be read from/to a single, contiguous block of memory. This is typical of current implementations of MPI.

The complexity of memory layout adds several complications to our simple theoretical model. First, outgoing data will not generally be contiguous, so space to buffer it on the sending processor will be needed. Second, if available memory is fragmented, then incoming messages can only be as large as the fragments, not the entire free memory. Once data has arrived, it may need to be copied into non-contiguous portions of an application data structure. Finally, our analysis assumed that total available memory at the end of a phase could be updated simply by the difference between the total amounts of outgoing and incoming data. Although this may be true, it says nothing about fragment sizes, which are the true, practical limit on useful memory for message passing.

Despite these caveats, our model and algorithms can sometimes be applied without modifications. Consider the case where the application controls memory and packs all data contiguously. Available memory will then consist of a single, contiguous block. If the application data structures allow for the memory to be reorganized to make outgoing message data contiguous, then no further buffering of outgoing data will be necessary. Once messages have arrived, memory can be locally reorganized to place them into the application data structures in the appropriate way, while keeping the available memory contiguous. In this way, the objections raised above are circumvented and our model and techniques are directly applicable.

However, most applications do not allow for such straightforward memory management

```

While data exchange is not complete do
    Determine my receive buffer and send buffer sizes.
    Communicate my send buffer size to processors I have data for.
    While there is available memory and data to receive do
        Choose a set of senders who will send to me.
        Choose message sizes that senders will send to me.
        Post asynchronous receives.
        Inform senders of their allotments.
    for each sending allotment I am given
        Construct and send message.
        Free up the memory for data sent.

```

Fig. 9. Template of a robust algorithm for memory-constrained interprocessor communication.

techniques. For such applications, our theoretical model provides guidance, but not robust theoretical guarantees of performance. In Fig. 9 we sketch an algorithm that performs robustly, but without the certainty of bounds on the number of phases. The basic idea is to have each processor determine its available memory at the beginning of each phase and to restrict the sizes of incoming messages to ensure that this memory is not exceeded. No assumptions are necessary about the amount of available memory in subsequent phases.

## VII. EXPERIMENTAL RESULTS

We have implemented the algorithms discussed in §VI, and performed a series of empirical studies on the LBL/NERSC Alvarez Cluster, an 80 node Pentium III Myrinet cluster with 2 processors in each node. All methods were implemented in C, and use MPI for message passing.

The problem instances in our experiments come from the data distribution associated with load balancing the operation a particular preconditioner known as overlapped Schwartz domain decomposition. The communication involved can be described in a fairly simple way. Initially, a graph is partitioned in such a way that each processor has about the same number of vertices. Now imagine that each processor needs to perform an op-

eration that depends on the number of vertices it owns plus the number of vertices on other processors that are adjacent to ones it owns. These “overlapped subdomains” will not generally be well balanced. But with a modest adjustment to the original partition, the overlapped domains can be made balanced (see [10] for details). The communication operation we are using in our experiments is that of transforming the original partition into a balanced overlapped partition.

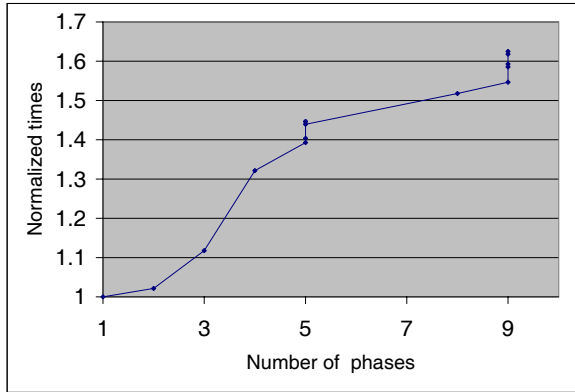
This communication is typical of those encountered in dynamic load balancing. The communication pattern is sparse and unstructured. Note also that different initial partitions lead to different communication patterns, so it is easy to generate a set of different, but related remapping instances. In all our experiments we used the matrix `ocean`, which has 143,437 rows and 819,186 nonzeros, with 8 double-precision numbers for each mesh node.

We report our results on four sets of experiments. The objective of the first set of experiments is to determine whether or not minimizing the number of phases truly is a valid metric for reducing overall communication cost. Our second experiment studies the relative merits of distributed and centralized approaches as discussed in §VI-A. Third, we study the performance of several different distributed heuristics, including those introduced in §VI-A. Finally, we consider the practical merits of parking.

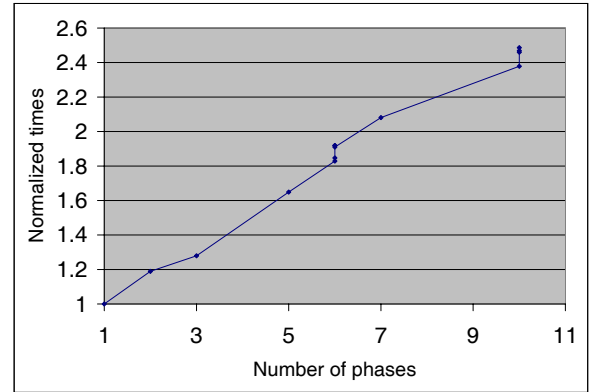
- Experiment 1: Model Validation

On the first set of experiments, we investigated the relation between the number of phases and the communication time. For scheduling methods, we implemented the *Random*, *Fortune Wheel* and *First Fit* techniques from §VI-A. Keeping the communication pattern the same, we reduced the amount of memory the algorithms were allowed to use, which increased the number of phases, but kept the volume of data transfer unchanged.

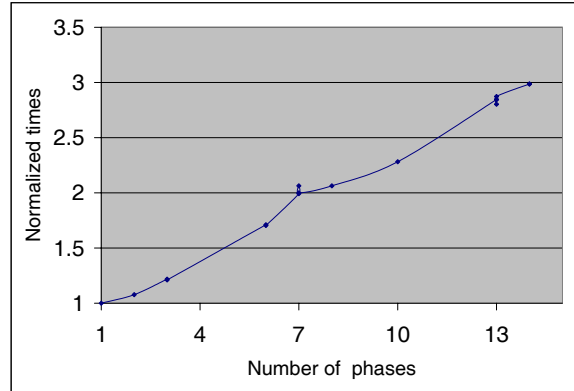
The results of these experiments are presented in Fig. 10. In this figure, the horizontal axis correspond to the number of phases, and the vertical axis correspond to the reorganization time normalized with respect to reorganization time without memory limitations, i.e., communication in a single phase. Each point in these graphs corresponds to a single instance in our experiments. It is important to note that decision making part in this heuristics take almost no time, thus the runtimes are essentially only the communication



32 Processors



64 Processors



128 Processors

Fig. 10. Correlation between number of phases and communication times.

times. The results show that runtimes grow roughly linearly with increasing number of phases, and that the the cost of each phase is a non-negligible fraction of the single phase communication cost (particularly for the larger numbers of processors). Thus, minimizing the number of phases is a valid objective for scheduling.

- Experiment 2: Centralized Versus Distributed Methods

In Fig. 11 we compare centralized and distributed methods. In this figure, times for the centralized methods are broken down to four components: “gather” is the time to collect problem information on one processor, “solve” is the time that one processor spends constructing a schedule, and “scatter” is the time to inform all processors of the schedule. We

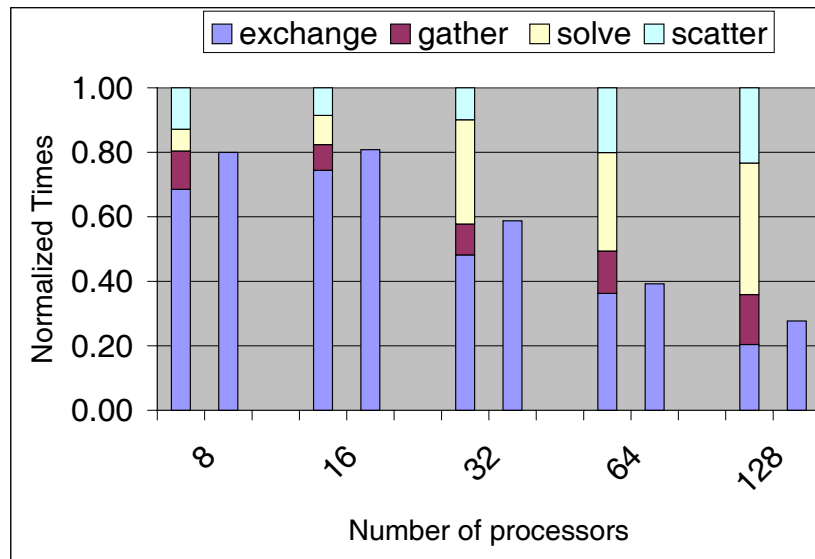


Fig. 11. Centralized and Distributed Methods

experimented with the centralized and distributed implementations of the five heuristics described in §VI-A on 8 problem instances. To take maximal advantage of the availability of all the information on a single processor, the centralized methods used parking to minimize the number of phases. The distributed methods in this experiment did not use parking.

For comparison purposes, times are normalized with respect to the total runtime of the centralized implementation. We have not observed significant differences in the relative performances of the different methods, so we present only the averages of normalized times in Fig. 11. As seen in the figure, the overhead for the centralized methods is quite high due to the gather and scatter time, and the serial algorithmic bottleneck. However, the centralized method is able to produce modestly better schedules. Our data indicates that distributed methods are generally more efficient. However, if the schedule is to be reused then perhaps the significant overhead induced by the centralized method could be justified.

- Experiment 3: Different Distributed Methods

We have implemented and compared the five distributed heuristics introduced in §VI-A, as well as the continuous relaxation method (CR) described in §IV. Since we are interested



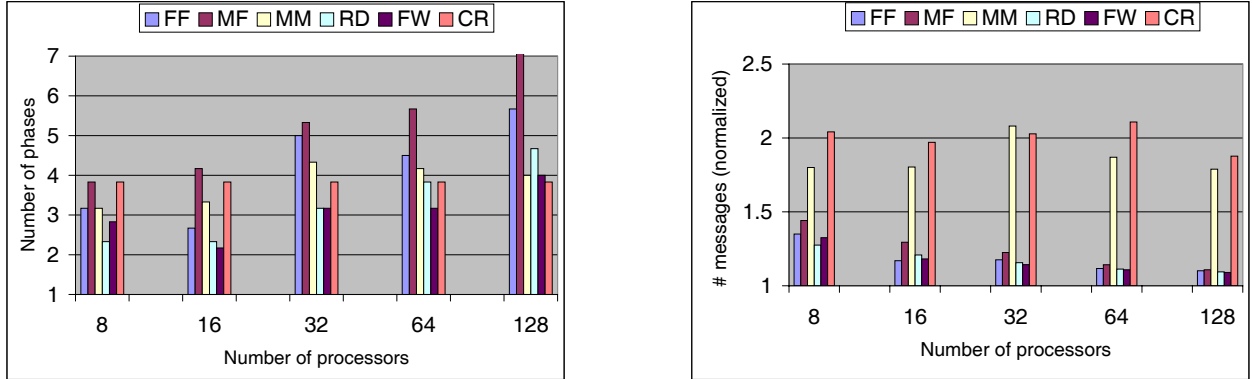


Fig. 12. Comparison of scheduling heuristics. (a) Number of phases (b) Number of messages

in a general purpose tool, our principle concern is with the total time consumed by a collection of invocations. For this reason, we ran the algorithms on 10 problem instances, counting the number of phases in the communication schedules they produce. The results of the average number of phases for each method are displayed in Fig. 12(a). It can be seen that randomized techniques (RD and FW) have the best performance, followed by the min-max technique (MM). The continuous relaxation method suffers from its pessimistic nature, since it always requires 2 extra phases, even for very simple problems. The first-fit heuristic (FF) suffers from the fact that communication lists are often sorted, and so all processors will want to communicate with the lowest numbered processors first. The poor performance of max-first (MF) can be attributed to natural communication patterns of dynamic load balancing. Dynamic balancing shifts load from overloaded processors to underloaded processors, thus some overloaded processors send a lot of data, but receive (almost) none. The max-first technique is likely to give priority to such processors, and so lead to poor distribution of the communication load.

Although the number of phases is our primary metric, the runtime of each phase is also important. The number of messages in a phase can effect performance due to large message setup times. Thus, schedules with fewer number of messages are preferred. Fig. 12(b) compares the heuristics in terms of the number of messages. As before, numbers are computed as the average of 10 different instances, but here the results are normalized by the number of communicating pairs of processors (i.e., the number of messages without

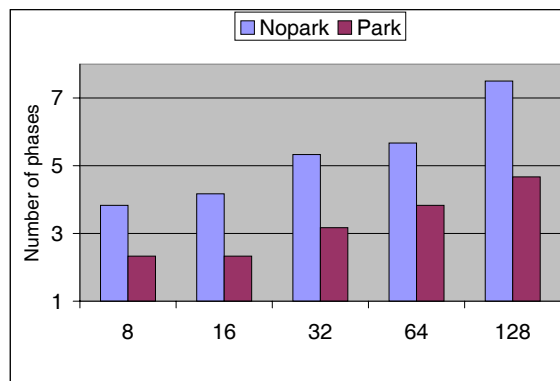


Fig. 13. Effect of parking on the number of phases

memory limitations). As expected, the number of messages is very high for the min-max and continuous relaxation heuristics. Performances of the other heuristics are very close. Recall that the continuous relaxation technique reorganizes data in pre and postprocessing phases, thus the total volume of communication is higher than what is required. We observed that in practice the additional data movement for this method is very high – 52%, 46% and 44% for 32, 64 and 128 processors, respectively. Overall, we can conclude that randomized techniques perform better than the other heuristics.

- Experiment 4: Parking

As discussed in §V, parking utilizes memory that would otherwise be wasted, potentially decreasing the number of phases. Our experiments as presented in Fig. 13, show that in practice the decrease in number of phases due to parking can be very significant. In many cases, the number of phases can be halved by using parking.

We have already discussed the overhead due to parking in §VI-B. In our implementation, parking requires a centralized algorithm, and so bears the costs associated with gathering data and broadcasting an answer as discussed above. In addition, parking increases the total volume of communication since data are routed through an intermediate processor instead of being sent directly to their destination. Our experiments showed that this is not a significant problem. The total volume of communication increases by only an average of 1% for 32 processors, 3% for 64 processors, and 3% for 128 processors, which can easily be compensated for by the significant decrease in the number of phases. The important

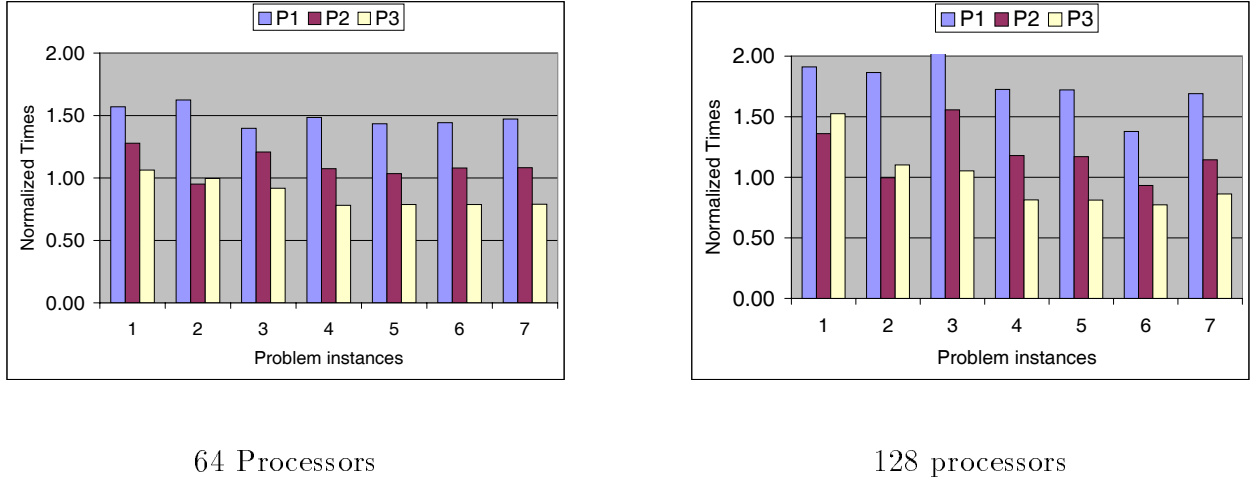


Fig. 14. Effect of parking on the number of phases

problem is how to limit the cost of the parking operation itself.

To retain the benefits, but reduce the cost, we tried using parking less frequently to decrease its cost, but still be able limit the number of phases. Fig. 14 presents the results for these experiments. In this figure, P1 corresponds to parking at each phase, and P2 (P3) corresponds to parking at every second (third) phase. The total available memory decreases from left to right for the problem instances. The times are normalized with respect to the times of solutions without parking. The results show that parking can decrease overall data reorganization times, especially when memory is very limited.

## VIII. CONCLUSION

We have addressed the problem of actually migrating data in the case of limited memory. When migrating data in an adaptive computation, it can happen that processors do not have enough memory to allocate space for their incoming data before they can release the space consumed by their outgoing data. In this case, the remapping operation must be decomposed into phases so that processors free up memory for the data they shipped out at end of a phase, making it available for the incoming data in the next phase.

In this paper, we studied how to complete the remapping operation in the minimum number of phases, the problem we call minimum phase remapping. We showed that the problem of determining whether a given transfer can be completed in a specified number of

phases is NP-Complete. A reduction of the minimum phase remapping problem to multi-commodity flow was presented. We showed how a continuous relaxation of the problem admitted a simple solution with two more phases than that of an optimal solution, but it might be difficult to obtain a good discrete solution from this continuous one. We also devised a practical approximation algorithm with a bound of 1.5 times the optimal solution.

We then conducted a series of empirical studies to validate the relevance of the number of phases as a metric of communication cost, and to study the performance of different versions of our algorithms. These experiments indicate that our methods can be implemented efficiently in parallel, and that multi-phase remapping can be performed inexpensively.

Our interest in this problem arises from our collaborative efforts to build general purpose libraries to support complicated parallel applications. Specifically, we are interested in providing additional robustness to parallel libraries. The functionality described here is now being added into Zoltan, a public-domain dynamic load balancing tool [3].

## REFERENCES

- [1] G. Cybenko, "Dynamic load balancing for distributed memory multiprocessors," *J. Parallel Distrib. Comput.*, vol. 7, pp. 279–301, 1989.
- [2] B. Hendrickson and K. Devine, "Dynamic load balancing in computational mechanics," *Comp. Meth. Appl. Mech. Eng.*, vol. 184, no. 2–4, pp. 485–500, 2000.
- [3] Karen D. Devine, Bruce Hendrickson, Erik G. Boman, Matthew M. St. John, and Courtenay Vaughan, "Zoltan: A dynamic load-balancing library for parallel applications – user's guide," Tech. Rep. SAND99-1377, Sandia National Laboratories, Albuquerque, NM, 1999, <http://www.cs.sandia.gov/Zoltan/>.
- [4] A. Pinar and B. Hendrickson, "Interprocessor communication with memory constraints," in *Proc. 12th ACM Symp. Parallel Algorithms and Architectures*, July 2000, pp. 39–45.
- [5] A. Pinar and B. Hendrickson, "Communication support for adaptive computation," in *Proc. 10th SIAM Conf. Parallel Processing for Scientific Computing*, March 2001.
- [6] R. Cypher and S. Konstantinidou, "Bounds on the efficiency of message-passing protocols for parallel computers," *SIAM J. Comput.*, vol. 25, no. 5, pp. 1082–1104, 1996.
- [7] Joseph Hall, Jason Hartline, Anna Karlin, Jared Saia, and John Wilkes, "On algorithms for efficient data migration," in *Proc. Symposium on Discrete Algorithms*, 2001.
- [8] R. M. Karp, "Reducibility among combinatorial problems," in *Complexity of Computer Computations*, R. E. Miller and J. W. Thatcher, Eds., pp. 85–103. Plenum Press, New York, NY, 1972.
- [9] Ravindra K. Ahuja, Thomas L. Magnanti, and James B. Orlin, *Network Flows: Theory, Algorithms and Applications*, Prentice Hall, Englewood Cliffs, NJ, 1993.
- [10] A. Pinar and B. Hendrickson, "Partitioning for complex objectives," in *Proc. 15th Intl. Parallel & Distrib. Proc. Symp.* 2001, IEEE.