

Performance Enhancement Strategies for Multi-Block Overset Grid CFD Applications

M. Jahed Djomehri^a, Rupak Biswas^{b,*}

^a*CSC, NASA Ames Research Center, Moffett Field, CA 94035, USA*

^b*NAS Division, NASA Ames Research Center, Moffett Field, CA 94035, USA*

Abstract

The overset grid methodology has significantly reduced time-to-solution of high-fidelity computational fluid dynamics (CFD) simulations about complex aerospace configurations. The solution process resolves the geometrical complexity of the problem domain by using separately generated but overlapping structured discretization grids that periodically exchange information through interpolation. However, high performance computations of such large-scale realistic applications must be handled efficiently on state-of-the-art parallel supercomputers. This paper analyzes the effects of various performance enhancement strategies on the parallel efficiency of an overset grid Navier-Stokes CFD application running on an SGI Origin2000 machine. Specifically, the role of asynchronous communication, grid splitting, and grid grouping strategies are presented and discussed. Details of a sophisticated graph partitioning technique for grid grouping are also provided. Results indicate that performance depends critically on the level of latency hiding and the quality of load balancing across the processors.

Key words: parallel performance, load balancing, multi-block applications, computational fluid dynamics

1 Introduction

The overset grid methodology [2] for high-fidelity computational fluid dynamics (CFD) simulations about complex aerospace configurations falls into the general class of Schwartz decomposition methods [15]. The solution process resolves the geometrical complexity of the problem domain by generating and using overlapping multi-block structured discretization grids. This overset

* Corresponding author. Tel.: +1-650-604-4411; Fax: +1-650-604-3957.

Email addresses: djomehri@nas.nasa.gov (M. Jahed Djomehri),
Rupak.Biswas@nasa.gov (Rupak Biswas).

approach typically employs a Chimera interpolation technique [18] to periodically update and exchange inter-grid boundary information.

However, to reduce time-to-solution, high performance computations of such large-scale realistic applications must be handled efficiently on state-of-the-art parallel supercomputers. Several pertinent papers describing numerical results and parallel implementations have been presented over the years at many conferences [1,7,14,17]. Those that are relevant to performance analysis are relatively outdated; their experiments were conducted on machines with low compute power using small-scale test problems [6,12,13,21]. This paper is the first attempt to report a detailed parallel performance evaluation of a high-fidelity multi-block overset CFD production code running large-scale complex-geometry applications on large numbers of processors. A preliminary version can be found in [3].

Various parallel programming paradigms have been developed for both distributed and shared memory systems. Currently, the most popular paradigms are message passing, shared memory programming, and their hybrid combination. Widely used scientific programs suitable for most modern architectures are implemented using a message passing paradigm, such as MPI, mainly for portability reasons. Fortunately, the overset grid method can readily employ MPI to exploit its coarse-grained parallelism as well as communicate information between distributed overlapping grids.

The parallel efficiency of the overset approach depends primarily upon the proper distribution of the computational workload and the minimization of the communication overhead among the processors. For most practical computational problems, optimal load balancing to minimize processor idle time is a challenging task. Overset applications with tens of millions of grid points may consist of many overlapping grids. Smart clustering of individual grids (also known as blocks or zones) into groups should therefore not only consider the total number of “weighted” grid points (described in Section 3.3), but also the size and connectivity of the inter-grid data. Major challenges during the grouping process may arise due to the wide variation in block sizes and the disparity in the number of inter-grid boundary points. Note also that for large processor sets, the overhead associated with boundary data exchange may adversely affect parallel performance.

This paper analyzes the effects of various performance enhancement techniques on the parallel efficiency of an overset grid Navier-Stokes CFD application called OVERFLOW. Specifically, the role of asynchronous communication, grid splitting, and grid grouping strategies are presented and discussed. First, we study the effect of synchronous and asynchronous communication via MPI. The asynchronous exchange is an attempt to relax the communication schedule in order to hide latency. Second, the splitting of large blocks as a means of controlling the computational load is analyzed. This is particularly important for scalability, where the same grid system must be retained for executing on different numbers of processors. Finally, two grid clustering techniques are examined: one based on a naive bin-packing approach and the

other using a more sophisticated graph partitioning method. All our experiments are conducted on an SGI Origin2000 machine using three test cases that simulate complex rotorcraft vortex dynamics and consist of between 63 million and 78 million grid points. Results indicate that performance depends critically on the level of latency hiding and the quality of load balancing across the processors.

The remainder of this paper is organized as follows. Section 2 provides a brief description of the OVERFLOW application. The performance enhancement techniques of grid splitting, asynchronous communication, and grid grouping are described in Section 3. Parallel performance results are presented and critically analyzed in Section 4. Finally, Section 5 concludes the paper with a summary and some key observations.

2 Numerical Methodology

In this section, we provide a brief overview of the overset grid CFD application called OVERFLOW, including the basics of its solution process, grid connectivity, and message-passing parallelization model.

2.1 Solution Process

The high-fidelity overset grid application, called OVERFLOW [2], owes its popularity within the aerodynamics community due to its ability to handle complex designs consisting of multiple geometric components, where individual body-fitted grids can be constructed easily about each component. The grids are either attached to the aerodynamics configuration (near-body), or are detached (off-body). The union of near- and off-body grids covers the entire computational domain (see Fig. 1(a) for a simple schematic).

OVERFLOW uses a Reynolds-averaged Navier-Stokes solver, augmented with a number of turbulence models. In this work, a special version of the code,

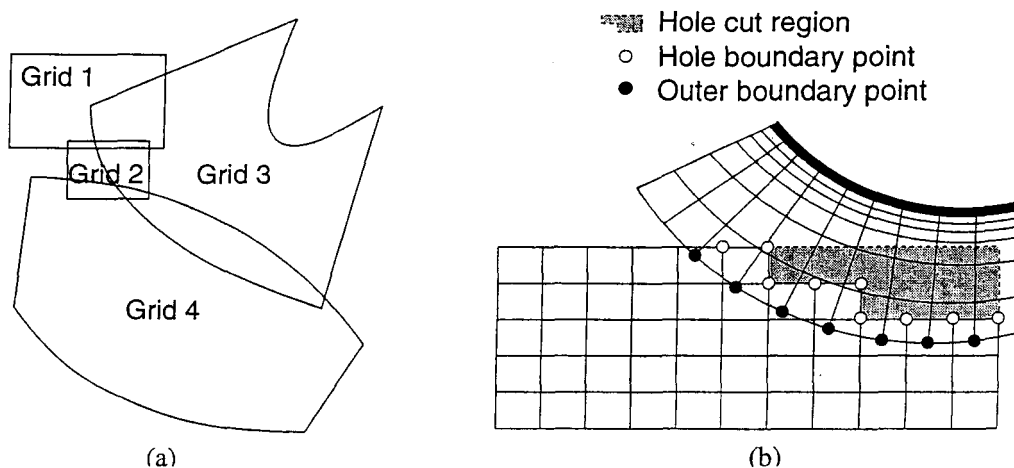


Fig. 1. (a) Overset grid schematic; (b) hole and outer inter-grid boundary points.

named OVERFLOW-D [9,10], is used. Unlike the original version which is primarily meant for fixed-body (static) grid systems, OVERFLOW-D is explicitly designed to simplify the modeling of components in relative motion (dynamic grid systems). For example, in typical rotary-wing problems, the near-field is modeled with one or more grids around the moving rotor blades. The code then automatically generates Cartesian “background” or “wake” grids, called bricks, that encompass these curvilinear near-body grids. At each time step, the flowfield equations are solved independently on each zone in a sequential manner. Overlapping boundary points or inter-grid data are updated from previous solutions prior to the start of the current time step using a Chimera interpolation procedure [18]. The code uses finite differences in space, with a variety of implicit/explicit time stepping.

2.2 Grid Connectivity

The Chimera interpolation procedure [18] determines the proper connectivity of the individual grids. To update inter-grid boundary data, the scheme has to process two types of boundary points: “hole” and “outer” boundary points (see Fig. 1(b)). Holes are cut in grids which intersect solid surfaces, such as when a portion of an overset grid lies inside a physical body. The hole boundary points are on the surfaces of these cuts. All other inter-grid boundary points are classified as outer. Adjacent grids are expected to have at least a one-cell, or a single fringe, overlap to ensure the continuity of the solutions; for higher-order accuracy and to retain certain physical features in the solution, a double fringe overlap is sometimes used [19]. A program named Domain Connectivity Function (DCF) [11] computes the inter-grid donor points that have to be supplied to other grids. The DCF procedure is incorporated into the OVERFLOW-D code and fully coupled with the flow solver. For dynamic grid systems, DCF has to be invoked at every time step to create new holes and inter-grid boundary data.

2.3 MPI Parallelization Model

The parallel version of the OVERFLOW-D application has been developed around the multi-block feature of the sequential code, which offers a natural coarse-grained parallelism [21]. The main computational logic at the top level of the sequential code consists of a “time-loop”, a “grid-loop”, and a “subiteration-loop”. The last two loops are nested within the time-loop. Within the grid-loop, solutions are obtained on the individual grids with imposed boundary conditions, where the Chimera interpolation procedure successively updates inter-grid boundaries after computing the numerical solution on each grid. Convergence of the solution process is accelerated by the subiteration-loop. Upon completion of the grid-loop, the solution is automatically advanced to the next time step by the time-loop. The overall procedure

may be thought of as a Gauss-Seidel iteration.

~~~~~% , ~~~~~7@◆◆◆◆

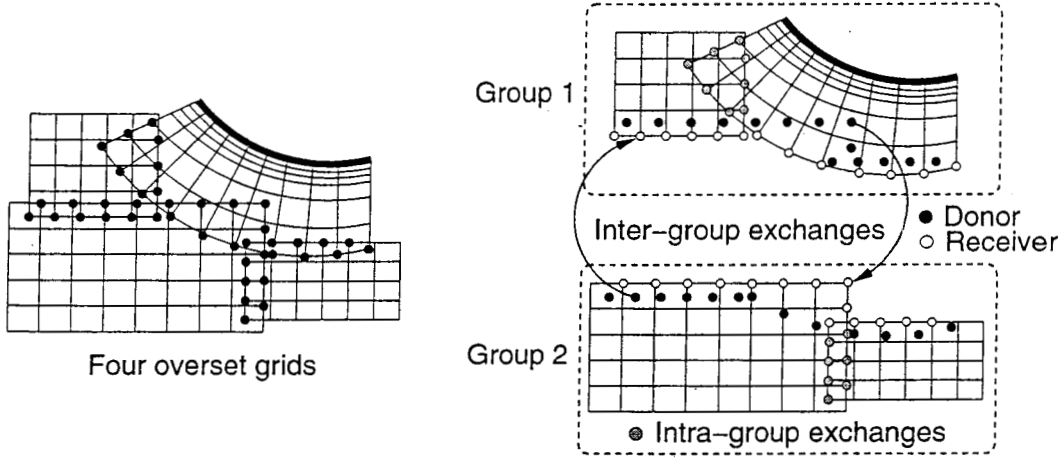


Fig. 2. Overset grid intra-group and inter-group communication.

A message passing programming model based on the MPI library was implemented using the single program multiple data (SPMD) paradigm. To facilitate parallel execution, a grouping strategy is required to assign each grid to an MPI process. The total number of groups,  $G$ , is equal to the total number of MPI processes,  $P$ . Since a grid can only belong in one group, the total number of grids,  $Z$ , must be at least equal to  $P$ . If  $Z$  is larger than  $P$ , a group will consist of more than one grid. Two techniques for clustering grids into groups are discussed later in Section 3.3.

The logic in the MPI programming model differs slightly from that of the sequential case (where  $G = P = 1$ ). Here the grid-loop is subdivided into two procedures, a loop over groups (“group-loop”) and a loop over the grids within each group. Since each MPI process is assigned to only one group, the group-loop is performed in parallel, with each group performing its own sequential grid-loop. The inter-grid boundary updates among the grids within each group (these are also called intra-group updates) are performed as in the serial case. Chimera updates are also necessary for overlapping grids that are in different groups, and are known as inter-group exchanges (see Fig. 2). The inter-group donor points from grids in group  $G_i$  to grids in group  $G_j$  are stored in a send buffer and exchanged between the corresponding processes via MPI calls. These inter-group exchanges are transmitted at the beginning of every time step based on the interpolatory updates from the previous time step.

### 3 Performance Enhancement Techniques

We have developed and utilized various performance enhancement techniques to improve the parallel efficiency of the OVERFLOW-D application. Specifically, the role of asynchronous communication, grid splitting, and grid grouping strategies are presented and discussed in this section. Superior parallel performance of such large-scale realistic applications on state-of-the-art commercial supercomputers is critical to advance our scientific understanding

and problem solving capability.

### 3.1 Asynchronous Communication

Almost all of the communication that is required in the OVERFLOW-D application concerns the exchange of inter-grid boundary data, and is contained in the subroutine, `qbc_exchange`. The message passing can be synchronous or asynchronous, but the choice significantly affects the MPI programming model. The synchronous communication is performed with blocking MPI send/receive calls, while the asynchronous communication uses non-blocking calls.

With synchronous communication, the total number of send/receive calls is  $P \times (P - 1)$ , counting even the messages of zero length. A send call is blocked until the receiving processor is ready to accept the message, i.e., until the matching receive call is posted. The increase in execution time caused by this communication pattern is analogous to the introduction of an implicit serialization into the code. The initial parallel version of OVERFLOW-D was implemented with synchronous message passing and tested with a relatively small dataset on 16 processors [21]. As a result, the communication time was quite insignificant and therefore acceptable. However, performance analysis using larger datasets and more processors (presented in Section 4) indicate a serious communication bottleneck for the exchange of boundary data via the

```

/* Send data from group ND to group NR */
do ND = 1, G
  if (myrank .eq. ND) then
    do NR = 1, G
      MLEN_SEND = ISND (NR)      /* Set length of send array */
      if (myrank .ne. NR) then
        call MPI_SEND (QBCSND, MLEN_SEND, ...)
      else
        do I = 1, MLEN_SEND
          QBCRCV (I) = QBCSND (I)      /* Memory copy */
        end do
      end if
    end do
  else
    /* Receive data from group ND */
    MLEN_RECV = IRCV (ND)      /* Set length of receive array */
    call MPI_RECV (QBCRCV, MLEN_RECV, ...)
  end if
end do

```

Fig. 3. Outline of the synchronous communication model in the original OVERFLOW-D code.

synchronous approach.

In order to be better able to compare the original synchronous and our new asynchronous communication strategies, we present in Fig. 3 an outline of the synchronous model. The group boundary data arrays are specified by QBCSND and QBCRCV with total lengths of MLEN\_SEND and MLEN\_RECV, respectively. The values of MLEN\_SEND and MLEN\_RECV are determined by arrays ISND and IRCV, respectively, for each group. As mentioned in Section 2, all boundary related arrays (QBCSND, QBCRCV, ISND, and IRCV) are determined by the DCF procedure prior to the start of a new time step and updated by Chimera interpolations.

Our first performance enhancement technique is to use asynchronous communication for inter-grid boundary data exchange within the `qbc_exchange` subroutine. The asynchronous strategy is an attempt to relax the communication schedule in order to hide latency. Asynchronous communication consists of non-blocking MPI send/receive calls. Unlike the corresponding blocking calls of the synchronous method, these invocations place no constraints on each other in terms of completion. Non-blocking receives complete immediately,

```

/* Post receives in group NR from group ND */
do ND = 1, G
  MLEN_RECV = IRCV (ND)           * Set length of receive array */
  if (MLEN_RECV .ne. 0) then
    if (myrank .ne. ND) then
      call MPI_Irecv (QBCRCV, MLEN_RECV, ...)
    end if
  end if
end do
/* Send data from group ND to group NR */
do NR = 1, G
  MLEN_SEND = ISND (NR)           * Set length of send array */
  if (myrank .ne. NR) then
    if (MLEN_SEND .ne. 0) then
      call MPI_Isend (QBCSND, MLEN_SEND, ...)
    end if
  else
    do I = 1, MLEN_SEND
      QBCRCV (I) = QBCSND (I)      /* Memory copy */
    end do
  end if
end do
/* Check that all receives have completed */
call MPI_WAITALL

```

Fig. 4. Outline of our asynchronous communication model in OVERFLOW-D.

even if no messages are available, and hence allow maximal concurrency; they are posted by receiving processors prior to the pertinent sends from the sending processors. Furthermore, messages of zero length are not sent to decrease the communication overhead. We have implemented this asynchronous message passing model in the current version of OVERFLOW-D.

In general, however, control flow and debugging can become a serious problem if, for instance, the order of messages needs to be preserved. Fortunately, in the overset grid application, the Chimera boundary updates take place at the completion of each time step, and the computations are independent of the order in which messages are sent or received. Being able to exploit this fact allows us to easily use asynchronous communication within OVERFLOW-D. Figure 4 gives an outline of the asynchronous approach that we have implemented. The same naming convention discussed with respect to the synchronous case is also adopted here. Comparisons with Fig. 3 shows that the outer loop containing the MPI\_SEND/MPI\_RECV pair in the synchronous case is now broken into two separate loops. The first posts non-blocking MPI\_RECV calls while the second posts non-blocking MPI\_SEND. Note that for the asynchronous strategy, the MPI barrier function (MPI\_WAITALL) must be invoked to ensure completion of the operations and to release the message buffers.

### 3.2 Grid Splitting

Load balancing is critically important for efficient parallel computing. The objective is to distribute equal computational workloads among the processors while minimizing the inter-processor communication cost. On a given platform, the primary procedure that affects the load balancing of an overset grid application is the grid grouping strategy. To facilitate parallel execution, each grid must be assigned to an MPI process. Since the total number of grids,  $Z$ , is at least equal to the number of processes,  $P$ , a proper clustering of the grids into  $G$  groups is required ( $G = P$ ).

Unfortunately, the sizes of the  $Z$  blocks in an overset grid system may vary substantially, thereby complicating the grouping procedure and significantly affecting the overall load balance. For instance, each near-body block is a three-dimensional curvilinear structured grid generated about the geometric components of an aerodynamics configuration. The dimensions of each block are primarily selected to introduce proper refinement into the grid spacing in an effort to maintain certain features of the physical solution, but have no bearing on the type of computations used, serial or parallel. Consequently, there may be orders of magnitude differences in near-body block sizes for the initial grid system. Recall that these near-body grids overlap the Cartesian wake (off-body) grid system to cover the entire computational domain. The indices of each grid  $z_i$ ,  $i = 1, 2, \dots, Z$ , varies from (1,1,1) to a maximum of  $(I_i, J_i, K_i)$ , for a total of  $I_i \times J_i \times K_i$  grid points.

A smart mechanism is therefore needed to limit the size of the individual blocks. One option is to add some control during the grid generation pro-



cess, but this would further complicate an already complex task. The strategy would also require information about the number of groups ( $G$ ) which may vary from one simulation run to the next depending on the chosen number of processors ( $P$ ), since  $G$  must be equal to  $P$ . The second approach, which we have implemented as part of this work, is to split the largest blocks into sub-blocks of desired sizes prior to grouping them. This performance enhancement technique is independent of the grid generation procedure and is automatically implemented at runtime, prior to the start of the time-loop.

The original version of OVERFLOW-D has the ability to perform some automatic grid splitting without any user input, but it was only to ensure that there were enough blocks  $Z$  to form  $G$  groups with  $G = P$ . However, for large test cases such as those used in this paper, further control is required. In particular, we must maintain exactly the same blocks in the grid system when executing on different numbers of processors to examine code scalability. In our latest version of OVERFLOW-D, the user specifies two input parameters, *maxnb* and *maxgrd*, for splitting the largest blocks based on some knowledge of the initial grid system and the maximum target value of  $P$ . All near- and off-body grids larger in size than *maxnb* and *maxgrd*, respectively, are then split into overlapping sub-blocks of smaller but equal sizes. For near-body grids, the split is done in one dimension only, depending on the values of  $I$ ,  $J$ , and  $K$ . The type of imposed boundary condition (periodic, reflecting, etc.) also plays a role in determining the splitting direction. For the uniform Cartesian off-body grids, splitting can be performed in multiple dimensions if necessary.

Conceptually, having smaller block sizes simplifies the load balancing procedure and leads to a more computationally balanced workload; however, a limiting drawback is an increase in the ratio of surface-to-volume grid points. A large value of this ratio indicates that the amount of overlap boundary data to be transferred via point-to-point communication between pairs of processors has increased disproportionately relative to the computational workload. Furthermore, since it is necessary to maintain at least a single (one-cell) and sometimes even a double (two-cell) fringe overlap between adjacent blocks, the total number of grid points increases during the splitting process, resulting in a larger computational and communication load per processor.

### 3.3 Grid Grouping

As mentioned in Section 3.2, the grid grouping strategy has a substantial effect on the quality of load balancing for an overset grid application like OVERFLOW-D. The grouping is a function of the following parameters: execution time per grid point, total number of grid points per block, number of blocks  $Z$ , volume of the total boundary data to be exchanged per processor, rate of communication, and total number of processors  $P$ . In principle, grouping depends only on the characteristics of the grids and their connectivity; it does not take into account the topology of the physical processors. The assignment of groups to processors is somewhat random, and is handled by the



```

/* LTF_MFT_ACC heuristic */
1:  sort tasks  $z_i$ ,  $i = 1, 2, \dots, Z$  in descending order by size (LTF)
2:  for each processor  $p_j$ ,  $j = 1, 2, \dots, P$ 
    set  $T(p_j) = 0$ 
3:  for each sorted task  $z_i$ ,  $i = 1, 2, \dots, Z$ 
3.1:  assign  $z_i$  to  $p_j$  with minimum  $T(p_j)$  (MFT)
      compute  $T(p_j) = T(p_j) + X_i$ 
3.2:  for each  $z_r \in R(z_i)$  assigned to  $p_k \neq p_j$  (ACC)
      set  $T(p_j) = T(p_j) + C_{ir}$ 
3.3:  for each  $z_d \in D(z_i)$  assigned to  $p_m \neq p_j$  (ACC)
      set  $T(p_m) = T(p_m) + C_{di}$ 
    end for

```

Fig. 5. Outline of the LTF\_MFT\_ACC task assignment heuristic from EVAH.

within the context of distributed grid computing across multiple resources [4]. In this work, we have modified EVAH to cluster overset grids into groups while taking into account their relative overlaps.

Among several heuristics that are available within EVAH, we have used the one called *largest task first with minimum finish time and available communication costs* (LTF\_MFT\_ACC). In the context of the current work, a task is synonymous with a block in the overset grid system. The size of a task is defined as the computation time for the corresponding block. An outline of the LTF\_MFT\_ACC procedure is presented in Fig. 5. It is constructed from the basic *largest task first* (LTF) heuristic that sorts the tasks in descending order by size. LTF is then enhanced by the systematic integration of the status of the processors in terms of their *minimum finish time* (MFT). Because of the overhead involved due to data exchanges between neighboring zones and their impact on overall execution time, the assignment heuristic is further augmented by including *available communication costs* (ACC). A procedure has been developed to interface the DCF subroutine of OVERFLOW-D with EVAH heuristics.

It is easiest to explain the LTF\_MFT\_ACC grouping strategy by using a simple example and then stepping through the procedure in Fig. 5. Figure 6(a) shows a graph representation of the overset grid system in Fig. 2 that is being partitioned across two processors,  $p_0$  and  $p_1$ . The computational time for block  $z_i$  is denoted as  $X_i$  and shown for all four blocks in Fig. 6(a). Similarly, the communication overhead from  $z_d$  (donor) to  $z_r$  (receiver) is denoted as  $C_{dr}$  and shown for all inter-grid data exchanges along the graph edges. In step 1, the four blocks are sorted in descending order by computational time; hence the order is:  $z_3, z_2, z_0, z_1$ . In step 2, the total execution times of the two processors are initialized:  $T(p_0) = T(p_1) = 0$ . Step 3 has to be executed four times since we have four grids that must be grouped. Block  $z_3$  is assigned to  $p_0$  and  $T(p_0) = 75$  in step 3.1. Since no other blocks have yet been assigned,

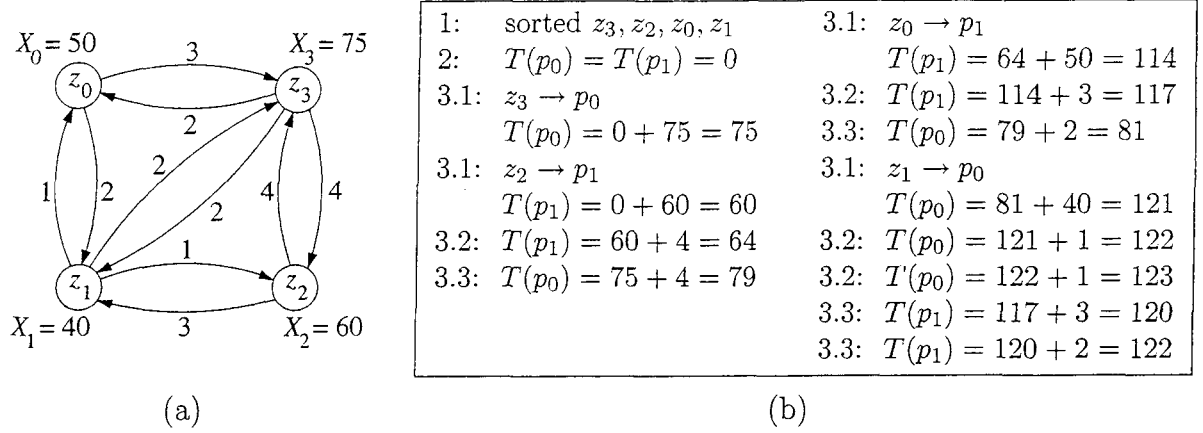


Fig. 6. (a) Graph representation of overset grid system in Fig. 2; (b) stepping through the LTF\_MFT\_ACC procedure in Fig. 5.

steps 3.2 and 3.3 are not executed.

Now  $z_2$  must be mapped to a processor that has the smallest total execution time; thus,  $z_2$  goes to  $p_1$  and  $T(p_1) = 60$  in step 3.1. In step 3.2, we need to look at all grids assigned to processors other than  $p_1$  that are also receivers of inter-grid data from  $z_2$ . This set of grids is denoted as  $R(z_2)$  in Fig. 5. The graph in Fig. 6(a) shows that  $z_1$  and  $z_3$  are in  $R(z_2)$ ; however,  $z_1$  is still unassigned. Since  $z_3$  belongs to  $p_0$ , the communication overhead  $C_{23}$  is added to  $T(p_1)$ ; hence,  $T(p_1) = 64$ . Similarly, in step 3.3, the set  $D(z_2)$  consists of grids that are donors of inter-grid data to  $z_2$ . Because  $z_1$  is unassigned and  $z_3$  is mapped to  $p_0$ ,  $T(p_0)$  is updated with  $C_{32}$ ; thus,  $T(p_0) = 79$ . The remainder of the grouping procedure is shown in Fig. 6(b).

#### 4 Parallel Performance Results

The CFD problem used for the experiments in this paper is a Navier-Stokes simulation of vortex dynamics in the complex wake flow region for hovering rotors. Figure 7 shows sectional views of the test application grid system. The Cartesian off-body wake grids surround the curvilinear near-body grids with uniform resolution, but become gradually coarser upon approaching the outer boundary of the computational domain. Specifically, the spacing of the off-body grid nearest the rotor blade is  $\Delta s$ , that for the next surrounding level is  $2\Delta s$ , and so on for every successive level. Figure 8 shows a cut plane through the computed vortex wake system including vortex sheets as well as a number of individual tip vortices. A complete description of the underlying physics and an extensive analysis of the numerical simulations pertinent to this test problem can be found in [20]. We have used the following three cases to evaluate our performance enhancement techniques discussed in Section 3:

- Case 1:  $Z = 454$ ,  $\sim 63$ M grid points,  $maxnb = 250$ K,  $maxgrd = 300$ K.
- Case 2:  $Z = 857$ ,  $\sim 69$ M grid points,  $maxnb = 100$ K,  $maxgrd = 100$ K.
- Case 3:  $Z = 1679$ ,  $\sim 78$ M grid points,  $maxnb = 60$ K,  $maxgrd = 70$ K.

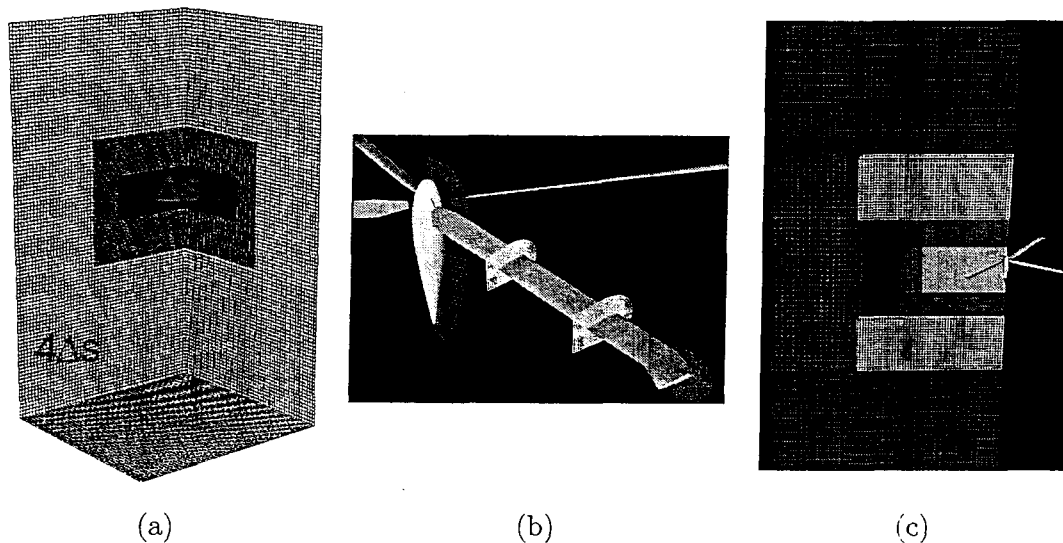


Fig. 7. Sectional views of the test application grid system: (a) off-body Cartesian wake grids, (b) near-body curvilinear grids, and (c) cut plane through the off-body wake grids surrounding the hub and rotors.

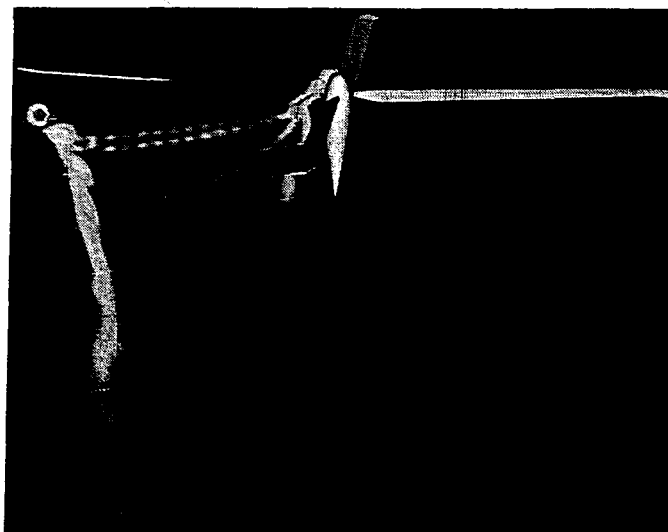


Fig. 8. Computed vorticity magnitude contours on a cutting plane located  $45^\circ$  behind the rotor blade.

All experiments were run on the 512-processor SGI Origin2000 shared-memory system at NASA Ames Research Center. Each Origin2000 node is a symmetric multiprocessor (SMP) containing two 400 MHz MIPS R12000 processors and 512 MB of local memory. Due to the memory requirements of the test cases, runs could not be conducted on less than 32 processors. Our timing results are averaged over 100 iterations and reported in seconds.

#### 4.1 Asynchronous Communication Results

Table 1 shows a comparison of various timings for Case 1 using synchronous (blocking send/receive) and asynchronous (non-blocking) communication. The grouping algorithm is the basic bin-packing strategy that is available with the original version of OVERFLOW-D. The execution time  $T_{exec}$  is the average time required to solve every time step of the application, and includes the computation, communication, Chimera interpolation, and processor idle times. The average computation ( $T_{comp}^{avg}$ ) and communication ( $T_{comm}^{avg}$ ) times over  $P$  processors are also shown. Finally, the maximum computation ( $T_{comp}^{max}$ ) and communication ( $T_{comm}^{max}$ ) times are reported and used to measure the quality of load balancing for each run.

The computation load balance factor ( $LB_{comp}$ ) is the ratio of  $T_{comp}^{max}$  to  $T_{comp}^{avg}$ , while the communication load balance factor ( $LB_{comm}$ ) is the ratio of  $T_{comm}^{max}$  to  $T_{comm}^{avg}$ . For a given  $P$ ,  $T_{comp}^{avg}$  is close to  $T_{comp}^{max}$  when the per-processor computation times are comparable; this in turn implies that  $LB_{comp}$  is almost unity. The degree to which the  $LB_{comp}$  is larger than unity is a measure of the computational load imbalance that adversely impacts the overall performance. The same argument applies to the communication times; therefore, the closer  $LB_{comp}$  and  $LB_{comm}$  are to unity, the higher is the quality of load balancing.

Notice that  $T_{comp}^{avg}$  is essentially the same since the computational workload for both runs is identical (same grid system and grid-to-processor assignment). However,  $T_{exec}$  for asynchronous communication is consistently lower, and shows bigger improvements as the number of processors increases. In fact, for  $P > 256$ , the non-blocking communication strategy reduces  $T_{exec}$  by more than a factor of two. The reason for this reduction can be found in the communication times. A comparison shows that  $T_{comm}^{max}$  and  $T_{comm}^{avg}$  for the synchronous runs are an order of magnitude larger than the corresponding times for the asynchronous communication. This is reflected in  $T_{exec}$  where communication usually accounts for less than 6% for the asynchronous case, but is more than 50% for many of the synchronous runs. The communication problem is exacerbated with increasing  $P$  since the number of messages exchanged is  $O(P^2)$  (although the individual message sizes decrease). Similar comparisons for Cases 2 and 3 would show even larger benefits when using the asynchronous strategy.

Scalability for the asynchronous case, with  $P \leq 256$ , is significantly better than its synchronous counterpart. For  $P \geq 320$ , scalability suffers for both cases, not only due to the relatively larger communication overhead, but also because of workload imbalance. The latter can be observed from the increasing value of  $LB_{comp}$ . However,  $LB_{comm}$  shows that communication is well-balanced across all processors for most runs, particularly for blocking communication (which is due to the very nature of synchronous communication). Overall results indicate the general superiority of the non-blocking asynchronous approach over synchronous communication for this application.

Table 1

Runtimes (in seconds) and load imbalance factors with synchronous and asynchronous communication, and bin-packing grouping strategy for Case 1

| $P$ | Synchronous |                  |                  |                  |                  |             |             |
|-----|-------------|------------------|------------------|------------------|------------------|-------------|-------------|
|     | $T_{exec}$  | $T_{comp}^{max}$ | $T_{comp}^{avg}$ | $T_{comm}^{max}$ | $T_{comm}^{avg}$ | $LB_{comp}$ | $LB_{comm}$ |
| 32  | 37.7        | 31.9             | 24.1             | 4.4              | 4.3              | 1.32        | 1.02        |
| 64  | 22.1        | 17.0             | 12.5             | 4.3              | 4.2              | 1.36        | 1.02        |
| 128 | 14.0        | 8.8              | 6.3              | 4.3              | 4.3              | 1.40        | 1.00        |
| 256 | 13.0        | 6.0              | 3.2              | 6.4              | 6.4              | 1.87        | 1.00        |
| 320 | 14.8        | 5.3              | 2.7              | 9.2              | 8.0              | 1.96        | 1.15        |
| 384 | 16.6        | 5.2              | 2.0              | 9.9              | 9.9              | 2.60        | 1.00        |
| 448 | 18.3        | 6.5              | 1.8              | 11.5             | 11.4             | 5.50        | 1.01        |

| $P$ | Asynchronous |                  |                  |                  |                  |             |             |
|-----|--------------|------------------|------------------|------------------|------------------|-------------|-------------|
|     | $T_{exec}$   | $T_{comp}^{max}$ | $T_{comp}^{avg}$ | $T_{comm}^{max}$ | $T_{comm}^{avg}$ | $LB_{comp}$ | $LB_{comm}$ |
| 32  | 34.6         | 32.7             | 24.5             | 0.70             | 0.61             | 1.33        | 1.15        |
| 64  | 18.0         | 16.8             | 12.4             | 0.41             | 0.35             | 1.35        | 1.17        |
| 128 | 9.8          | 8.8              | 6.3              | 0.36             | 0.31             | 1.40        | 1.16        |
| 256 | 7.0          | 5.9              | 3.2              | 0.37             | 0.30             | 1.84        | 1.23        |
| 320 | 6.9          | 5.1              | 2.6              | 0.98             | 0.68             | 1.96        | 1.44        |
| 384 | 6.8          | 6.0              | 2.0              | 0.48             | 0.36             | 3.00        | 1.33        |
| 448 | 7.0          | 6.3              | 1.8              | 0.49             | 0.43             | 3.50        | 1.14        |

#### 4.2 Grid Splitting Results

The impact of grid splitting on load balancing quality is investigated for all the three cases. Case 1 models the situation where the number of splits per block is low, or equivalently, the sizes of the newly-created sub-blocks are quite large. In Case 2, the sizes of the sub-blocks are somewhat smaller, while in Case 3, they are even more so. All runs use asynchronous communication and the bin-packing strategy to cluster grids into groups. Timings for Case 1 are shown in Table 1, while those for the other two cases are presented in Table 2.

The overall quality of computational workload balancing for the three cases can be observed by comparing  $LB_{comp}$  from Tables 1 and 2. Obviously, the factor increases with the number of processors as load balancing becomes more challenging with a fixed problem size. As expected, Case 3 exhibits the best quality for any given value of  $P$  because it has the largest number of grids

Table 2

Runtimes (in seconds) and load imbalance factors with asynchronous communication and bin-packing grouping strategy for Cases 2 and 3

| Case 2 |            |                  |                  |                  |                  |             |             |
|--------|------------|------------------|------------------|------------------|------------------|-------------|-------------|
| $P$    | $T_{exec}$ | $T_{comp}^{max}$ | $T_{comp}^{avg}$ | $T_{comm}^{max}$ | $T_{comm}^{avg}$ | $LB_{comp}$ | $LB_{comm}$ |
| 32     | 32.0       | 29.3             | 23.0             | 1.30             | 0.90             | 1.27        | 1.44        |
| 64     | 13.5       | 11.9             | 10.8             | 0.67             | 0.55             | 1.10        | 1.22        |
| 128    | 7.6        | 6.2              | 5.4              | 0.60             | 0.52             | 1.15        | 1.15        |
| 256    | 5.5        | 3.7              | 2.8              | 0.88             | 0.50             | 1.32        | 1.76        |
| 320    | 4.7        | 2.9              | 2.2              | 0.57             | 0.46             | 1.32        | 1.24        |
| 384    | 4.7        | 2.9              | 1.9              | 0.99             | 0.56             | 1.53        | 1.77        |
| 448    | 4.5        | 3.0              | 1.7              | 0.85             | 0.46             | 1.76        | 1.85        |

| Case 3 |            |                  |                  |                  |                  |             |             |
|--------|------------|------------------|------------------|------------------|------------------|-------------|-------------|
| $P$    | $T_{exec}$ | $T_{comp}^{max}$ | $T_{comp}^{avg}$ | $T_{comm}^{max}$ | $T_{comm}^{avg}$ | $LB_{comp}$ | $LB_{comm}$ |
| 32     | 39.9       | 34.5             | 27.8             | 3.10             | 2.20             | 1.24        | 1.41        |
| 64     | 13.8       | 11.4             | 10.4             | 0.95             | 0.75             | 1.10        | 1.27        |
| 128    | 7.9        | 6.4              | 5.2              | 0.85             | 0.70             | 1.23        | 1.21        |
| 256    | 4.5        | 3.1              | 2.6              | 0.95             | 0.68             | 1.19        | 1.40        |
| 320    | 4.3        | 2.8              | 2.1              | 0.90             | 0.61             | 1.33        | 1.48        |
| 384    | 4.0        | 2.4              | 1.8              | 0.65             | 0.57             | 1.33        | 1.14        |
| 448    | 3.8        | 2.3              | 1.6              | 0.71             | 0.60             | 1.44        | 1.18        |

which are all generally smaller, i.e., it has the finest granularity. In fact, the significant improvement in load balancing for  $P \geq 256$  causes  $T_{exec}$  to be reduced by almost 40% over that for Case 1. It should be noted here that though we are evaluating the level of workload imbalance from the runtimes, the grid splitting and grouping strategies are based on the number of weighted grid points. Computed from that perspective, the load balance quality is somewhat better but follows the same trend.

Let us now look at the communication times. Notice that both  $T_{comm}^{max}$  and  $T_{comm}^{avg}$  generally increase with increasing number of blocks (Case 1 through Case 3). (The communication time also depends on the topology and connectivity of the grid system.) This is because even though grid splitting has a positive impact on the computational load balance, it adversely affects the communication time. Basically, the surface area increases with the number of blocks, thereby increasing the volume of the boundary exchange data. For example, the ratio of surface-to-volume grid points for the three cases are



11%, 14%, and 18%, respectively. Communication therefore also accounts for a larger percentage of the total execution time. The ratio of  $T_{comm}^{avg}$  to  $T_{exec}$  is 2–10% for Case 1, 3–12% for Case 2, and 5–16% for Case 3.

Conceptually, the splitting of grids into smaller blocks should also improve the communication load balance  $LB_{comm}$ . Clearly, Case 3 does not have the best overall communication load balance, and Case 2 has poorer quality than Case 1. These results indicate that the optimal choice of the splitting parameters  $maxnb$  and  $maxgrd$  depends on the number of processors used. However, in our experiments, we wanted a fixed grid system independent of the processor count. Note that because of the complexity of OVERFLOW-D, grid splitting has been implemented in only one coordinate direction at this time for the near-body grids. Even if it were available in multiple directions, most grids would not benefit due to boundary condition and viscous direction splitting restrictions.

Finally, parallel scalability also improves with more blocks. This can be observed by comparing the  $T_{exec}$  times in Tables 1 and 2. In fact, we obtain superlinear speedup between 32 and 64 processors for Cases 2 and 3. This occurs partly because as  $P$  increases, a larger fraction of the problem fits in cache. Also,  $T_{exec}$  decreases consistently for Case 3 all the way to the maximum number of processors used. Overall, our grid splitting investigations show that a larger number of smaller blocks improves computational load balance and parallel scalability; however, there is a tradeoff since a large number of splits adversely affects efficiency due to an increase in the surface-to-volume ratio of grid points.

### 4.3 Grid Grouping Results

We compare our EVAH-based heuristic grid grouping strategy with naive bin-packing only for Case 2, and investigate their role on load balance quality. All timing results are presented in Table 3 (the performance data for bin-packing is reproduced from Table 2 for easier comparison).

The results in Section 4.2 showed that the quality of communication load balancing  $LB_{comm}$  is a big drawback of the bin-packing strategy. Table 3 demonstrates that the EVAH heuristic technique improves this factor considerably. In fact, except for  $P = 384$ ,  $LB_{comm}$  using EVAH is at most 1.17. However,  $T_{comm}^{avg}$  is always larger for EVAH, and is 5–15% of  $T_{exec}$  (compared to 3–12% for bin-packing). This is expected since the overall goal of the heuristic grouping strategy is to reduce  $T_{exec}$  by balancing inter-processor communication, while bin-packing tries to merely minimize total communication.

Results in Table 3 also show that except when  $P = 384$ , EVAH returns the better values for  $LB_{comp}$ . However, it should be noted here that though we are evaluating  $LB_{comp}$  from the computational run times, both grid grouping strategies are based on the number of weighted grid points. The three left plots in Fig. 9 show the distribution of weighted grid points across 64, 128, and 256 processors for bin-packing and EVAH. The “predicted” values of  $LB_{comp}$  are

Table 3

Runtimes (in seconds) and load imbalance factors with bin-packing and EVAH heuristic grouping strategies, and asynchronous communication for Case 2

| $P$ | Bin-packing |                  |                  |                  |                  |             |             |
|-----|-------------|------------------|------------------|------------------|------------------|-------------|-------------|
|     | $T_{exec}$  | $T_{comp}^{max}$ | $T_{comp}^{avg}$ | $T_{comm}^{max}$ | $T_{comm}^{avg}$ | $LB_{comp}$ | $LB_{comm}$ |
| 32  | 32.0        | 29.3             | 23.0             | 1.30             | 0.90             | 1.27        | 1.44        |
| 64  | 13.5        | 11.9             | 10.8             | 0.67             | 0.55             | 1.10        | 1.22        |
| 128 | 7.6         | 6.2              | 5.4              | 0.60             | 0.52             | 1.15        | 1.15        |
| 256 | 5.5         | 3.7              | 2.8              | 0.88             | 0.50             | 1.32        | 1.76        |
| 320 | 4.7         | 2.9              | 2.2              | 0.57             | 0.46             | 1.32        | 1.24        |
| 384 | 4.7         | 2.9              | 1.9              | 0.99             | 0.56             | 1.53        | 1.77        |
| 448 | 4.5         | 3.0              | 1.7              | 0.85             | 0.46             | 1.76        | 1.85        |

| $P$ | EVAH heuristic |                  |                  |                  |                  |             |             |
|-----|----------------|------------------|------------------|------------------|------------------|-------------|-------------|
|     | $T_{exec}$     | $T_{comp}^{max}$ | $T_{comp}^{avg}$ | $T_{comm}^{max}$ | $T_{comm}^{avg}$ | $LB_{comp}$ | $LB_{comm}$ |
| 32  | 26.2           | 23.3             | 21.8             | 1.50             | 1.28             | 1.07        | 1.17        |
| 64  | 13.0           | 11.3             | 10.8             | 0.76             | 0.66             | 1.05        | 1.15        |
| 128 | 7.3            | 5.8              | 5.4              | 0.99             | 0.89             | 1.07        | 1.11        |
| 256 | 4.8            | 3.2              | 2.7              | 0.77             | 0.67             | 1.19        | 1.15        |
| 320 | 4.4            | 3.0              | 2.3              | 0.76             | 0.65             | 1.30        | 1.17        |
| 384 | 4.7            | 3.1              | 1.9              | 0.97             | 0.55             | 1.63        | 1.76        |
| 448 | 4.3            | 3.0              | 1.7              | 0.73             | 0.64             | 1.76        | 1.14        |

also reported and are typically somewhat better than those computed from actual run times (see Table 3), but demonstrate the same overall trend.

The three plots on the right in Fig. 9 present a more detailed report of the execution, computation, and communication times per processor for  $P = 64$ , 128, and 256. Due to the synchronization of MPI processes,  $T_{exec}$  is independent of the processor ID and shown at the top of the plot area for each case. For the sake of clarity,  $T_{comm}$  is shown at a different scale indicated by the right vertical axis. Observe that the EVAH  $T_{comp}$  and  $T_{comm}$  curves are consistently much smoother than those for bin-packing, indicating a much more uniform distribution across processors.

Performance scalability for both strategies when using more than 256 processors is low due to our fixed problem size. For example, when  $P = 448$ , each group contains, on average, only two grids (since  $Z = 857$ ). With such a low number of blocks per group, the effectiveness of any strategy is diminished; moreover, the communication overhead relative to computation may increase

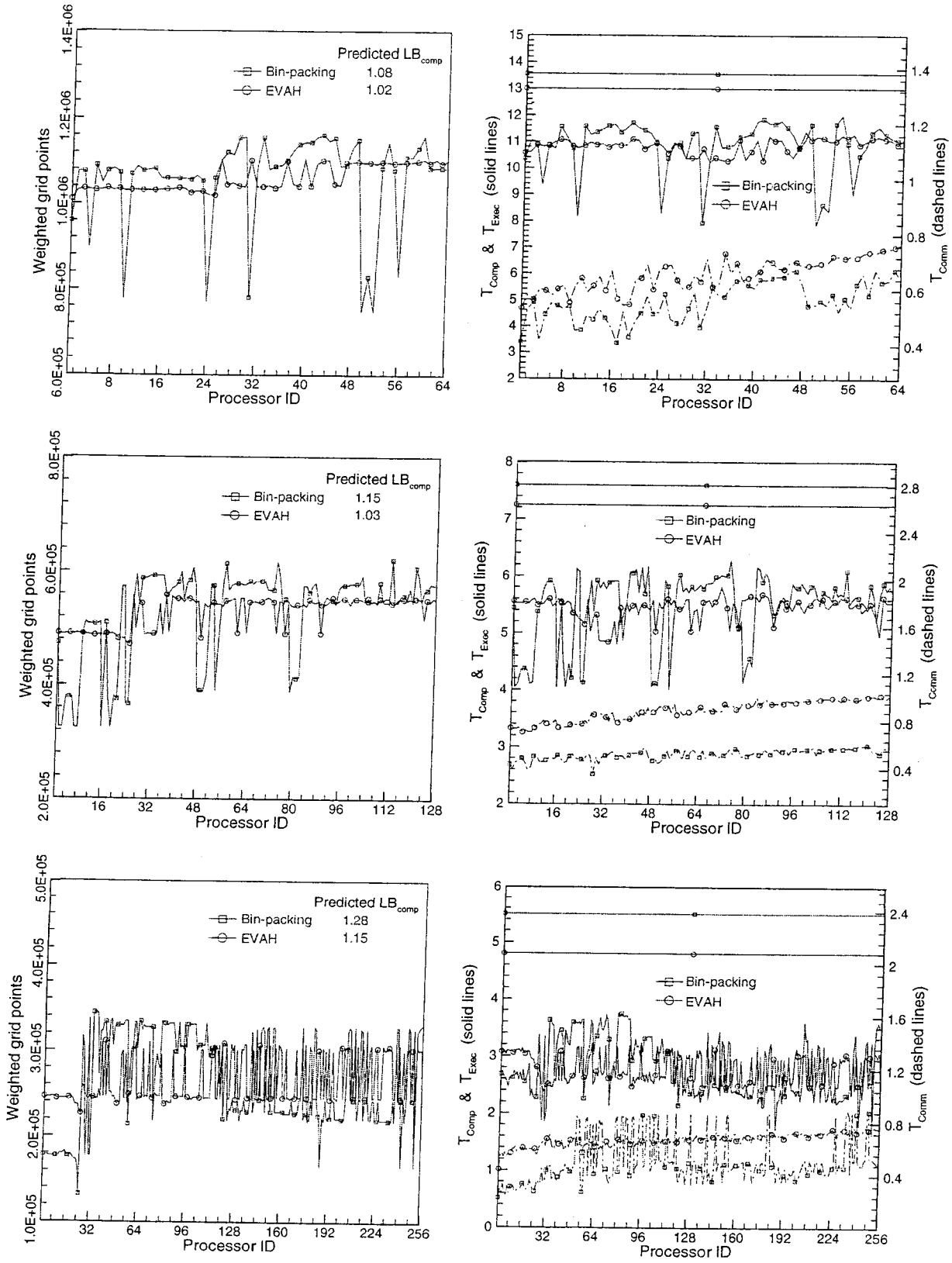


Fig. 9. Distribution of weighted grid points (left) and execution, computation, and communication times (right) per processor for  $P = 64, 128, \text{ and } 256$  (top to bottom).

substantially. For instance, with low processor counts, the communication-to-computation ratio is less than 6%, but grows to more than 37% with higher counts.

## 5 Summary and Conclusions

The overset grid method is a powerful technique for high-fidelity CFD simulations about complex aerospace configurations. In this paper, we presented and analyzed three parallel performance enhancement techniques for efficient computations of such large-scale realistic applications on state-of-the-art supercomputers. Specifically, the role of asynchronous communication, grid splitting, and grid grouping strategies were discussed. The asynchronous exchange relaxed the communication schedule in order to hide latency. Grid splitting was used to improve computational load balance while retaining the same grid system on different numbers of processors. Finally, a heuristic grid clustering technique balanced inter-processor communication with the goal of reducing the overall execution time.

All experiments were performed with the OVERFLOW-D Navier-Stokes code on a 512-processor Origin2000 system at NASA Ames Research Center. The CFD problem was the simulation of vortex dynamics in the complex flow region for hovering rotors. The grid systems for our three test cases consisted between 454 and 1679 overset grids, and varied in size from 63 million to 78 million grid points. The asynchronous communication strategy reduced execution time by more than a factor of two by significantly reducing the communication overhead. Grid splitting improved the workload balance by increasing the number of grids; however, the relative communication cost was adversely affected due to a larger surface-to-volume ratio of grid points. The heuristic grid grouping strategy compared extremely favorably with the original bin-packing technique. It improved the communication balance considerably while reducing the execution time. Overall results indicated that all three performance enhancement techniques are very effective in improving the quality of load balance and reducing execution time for overset grid applications.

Further improvements in the scalability of the overset grid methodology could be sought by using a more sophisticated parallel programming paradigm especially when the number of blocks  $Z$  is comparable to the number of processors  $P$ , or even when  $P > Z$ . One potential strategy that can be exploited on SMP clusters is to use a hybrid MPI+OpenMP multilevel programming style [5]. This approach is currently under investigation.

## Acknowledgements

The work of the first author was supported by NASA Ames Research Center under Contract Number DTTS59-99-D-00437/A61812D with AMTI/CSC. The authors would like to thank Prof. N. Lopez-Benitez at Texas Tech Uni-

versity for his help with EVAH, and Drs. M. Potsdam and R. Strawn of Army/NASA Rotorcraft Division for providing the test cases.

## References

- [1] American Institute of Aeronautics and Astronautics conferences, Available from URL: <http://www.aiaa.org/>.
- [2] P.G. Buning, D.C. Jespersen, T.H. Pulliam, W.M. Chan, J.P. Slotnick, S.E. Krist, and K.J. Renze, *Overflow User's Manual, Version 1.8g*, NASA Langley Research Center, Hampton, VA, 1999.
- [3] M.J. Djomehri, R. Biswas, M. Potsdam, and R.C. Strawn, An analysis of performance enhancement techniques for overset grid applications, in: *Proceedings 17th International Parallel and Distributed Processing Symposium* (Nice, France, 2003).
- [4] M.J. Djomehri, R. Biswas, R.F. Van der Wijngaart, and M. Yarrow, Parallel and distributed computational fluid dynamics: Experimental results and challenges, in: *Proceedings 7th International Conference on High Performance Computing* (Bangalore, India, 2000) 183–193.
- [5] M.J. Djomehri and H. Jin, Hybrid MPI+OpenMP programming of an overset CFD solver and performance investigations, *Technical Report NAS-02-002*, NASA Ames Research Center, Moffett Field, CA, 2002.
- [6] J. Häuser, M. Spel, J. Muylaert, and R. Williams, ParNSS: An efficient parallel Navier-Stokes solver for complex geometries, in: *Proceedings 25th AIAA Fluid Dynamics Conference* (Colorado Springs, CO, 1994) Paper 94-2263.
- [7] International Parallel and Distributed Processing Symposium series, Available from URL: <http://www.ipdps.org/>.
- [8] N. Lopez-Benitez, M.J. Djomehri, and R. Biswas, Task assignment heuristics for distributed CFD applications, in: *Proceedings 30th International Conference on Parallel Processing Workshops* (Valencia, Spain, 2001) 128–133.
- [9] R. Meakin, A new method for establishing inter-grid communication among systems of overset grids, in: *Proceedings 10th AIAA Computational Fluid Dynamics Conference* (Honolulu, HI, 1991) 662–676.
- [10] R. Meakin, On adaptive refinement and overset structured grids, in: *Proceedings 13th AIAA Computational Fluid Dynamics Conference* (Snowmass, CO, 1997) Paper 97-1858.
- [11] R. Meakin and A.M. Wissink, Unsteady aerodynamic simulation of static and moving bodies using scalable computers, in: *Proceedings 14th AIAA Computational Fluid Dynamics Conference* (Norfolk, VA, 1999) Paper 99-3302.
- [12] C. de Nicola, R. Tognaccini, and P. Visingardi, Multiblock structured algorithms in parallel CFD, in: *Proceedings Parallel Computational Fluid Dynamics Conference* (Pasadena, CA, 1995) 1–8.

- [13] R. Pankajakshan and W.R. Briley, Parallel solution of viscous incompressible flow on multi-block structured grids using MPI, in: *Proceedings Parallel Computational Fluid Dynamics Conference* (Pasadena, CA, 1995) 601–608.
- [14] Parallel Computational Fluid Dynamics Conference series, Available from URL: <http://www.parcfd.org/>.
- [15] Y. Saad, *Iterative Methods for Sparse Linear Systems*, PWS Publishing Company, Boston, MA, 1996.
- [16] K. Schloegel, G. Karypis, and V. Kumar. Graph partitioning for high performance scientific simulations, in: *CRPC Parallel Computing Handbook* Morgan Kautmann, San Francisco, CA, 2000.
- [17] SCxy Conference series, Available from URL: <http://www.sc-conference.org>
- [18] J.L. Steger, F.C. Dougherty, and J.A. Benek, A Chimera grid scheme, *Advances in Grid Generation*, ASME FED-5 (1983).
- [19] R.C. Strawn and J.U. Ahmad, Computational modeling of hovering rotors and wakes, in: *Proceedings 38th AIAA Aerospace Sciences Meeting* (Reno, NV, 2000) Paper 2000-0110.
- [20] R.C. Strawn and M.J. Djomehri, Computational modeling of hovering rotor and wake aerodynamics, *Journal of Aircraft* 39 (2002) 786–793.
- [21] A.M. Wissink and R. Meakin, Computational fluid dynamics with adaptive overset grids on parallel and distributed computer platforms, in: *Proceedings International Conference on Parallel and Distributed Processing Techniques and Applications* (Las Vegas, NV, 1998) 1628–1634.