

# Gang Scheduling Extensions for I/O Intensive Workloads

Yanyong Zhang<sup>†</sup>, Antony Yang<sup>†</sup>, Anand Sivasubramaniam<sup>‡</sup>, Jose Moreira<sup>§</sup>

<sup>†</sup> Department of Electrical & Computer Engg.  
Rutgers, The State University of New Jersey  
Piscataway NJ 08854  
{yyzhang, pheroth}@ece.rutgers.edu

<sup>‡</sup> Department of Computer Science & Engg.  
The Pennsylvania State University  
University Park PA 16802  
{anand}@cse.psu.edu

<sup>§</sup> IBM T. J. Watson Research Center  
P. O. Box 218  
Yorktown Heights NY 10598-0218  
{jmoreira}@us.ibm.com

## Abstract

Scientific applications are becoming more I/O demanding than ever. For such applications, the system with dedicated I/O nodes does not provide enough scalability. Rather, a serverless approach is a viable alternative. However, with the serverless approach, a job's execution time is decided by whether it is co-located with the file blocks it needs. Gang scheduling (GS), which is used in supercomputing centers to schedule parallel jobs, is completely not aware of the application's spatial preferences.

In this paper, we show that gang scheduling delivers poor performance towards workloads with high I/O intensities (I/O ratio higher than 50%). We propose an I/O-aware extension of gang scheduling, IOGS, which co-locates jobs with their files. While IOGS performs better for high I/O intensity workloads, its performance for workloads with lower I/O intensities is rather poor because of high system fragmentation. Further, we propose an adaptive strategy, adaptive-IOGS, which attempts to combine the advantages of both gang scheduling and GS, and we show that adaptive-IOGS is better than the other two schemes in many scenarios. Finally, we combine process migration techniques with adaptive-IOGS, and propose Migration-IOGS, which is shown to be the best among the four for a wide spectrum of workloads.

## 1 Introduction

Scheduling strategies can have a significant impact on the performance characteristics of a large parallel system [5]. Early strategies used a space-sharing approach,

wherein jobs can concurrently run on different nodes of the system, but each node is exclusively assigned to a job, until the job completes (no pre-emption is considered). The wait and response times for jobs with an exclusively space-sharing strategy can be relatively high. Gang scheduling (or coscheduling) is a technique to improve the performance by adding a time-sharing dimension to space sharing [15]. This technique virtualizes the physical machine by slicing the time axis into multiple virtual machines. Tasks of a parallel job are coscheduled to run in the same time-slices. The number of virtual machines created (equal to the number of time slices), is called the multiprogramming level (MPL) of the system. This approach opens more opportunities for the execution of parallel jobs, and is thus quite effective in reducing the wait time, at the expense of increasing the apparent job execution time. Gang scheduling can increase the system utilization as well. A considerable body of previous work has been done in investigating different variations of gang scheduling [19, 18, 20, 17]. Gang-scheduling is now often used in supercomputing centers. For example, it has been used in the prototype GangLL job scheduling system developed by IBM Research for the ASCI Blue-Pacific machine at LLNL (a large scale parallel system spanning thousands of nodes [13]).

On the application's end, over the past decade, the input data sizes that parallel applications are facing have increased dramatically, at a much faster pace than the performance improvement of the I/O system [12]. Supercomputing centers usually deploy dedicated I/O nodes to store the data and manage the data access. Data is striped across high-performance disk arrays that are attached to

the I/O nodes. The I/O bandwidth provided in such systems is limited by the number of I/O channels and the number of disks, which are normally smaller than the number of available computing units. This is not enough to serve the emerging scientific applications that are more complex than ever, and are more I/O demanding than ever. Fortunately, several research groups have proposed a serverless approach to address this problem [1, 2]. In a serverless file system, each node in the system serves as both compute node and I/O node. All the data files are distributed across the disks of the nodes in the system. With this approach, we have as many I/O nodes as the compute nodes. It is much easier to scale than the dedicated solution. The only assumption this approach has is reasonably fast interconnect between the nodes, which should not be a concern because of the rapid progress in technology.

While the serverless approach shows a lot of promise, realizing the promise is not a trivial task. It raises new challenges to the design of a job scheduler. Files are partitioned across a set of nodes in the system, and, on the same set of nodes, jobs are running as well. Suppose task  $t$  needs to access block  $b$  of file  $f$ . The associated access cost can vary considerably depending on the location of  $b$ . If  $b$  is hosted on the same node where  $t$  is running, then  $t$  only needs to go to the local disk to fetch the data. Otherwise,  $t$  needs to pay extra network cost to fetch the data from a remote disk. The disk accesses now become asymmetric. In order to avoid this extra network cost, a scheduler must co-locate tasks of parallel jobs with the file blocks they need. This will be even more urgent for I/O intensive jobs.

Although gang scheduling, which is used in many supercomputing centers, is effective for traditional workloads, most of its implementations are completely unaware of applications' spatial preferences. Instead, they focus on maximizing system usage by accommodating jobs as quickly as possible, and by assigning more CPU time to the running jobs. However, it is not yet clear which of the two factors - I/O awareness, or system utilization - is more important to the performance, at what scenarios, by how much? And, can we do even better by adaptively balancing between the two? This paper sets out to answer all the above questions by conducting numerous simulation based experiments on a wide range of workloads. This paper proposes several new scheduling heuristics, which try to schedule jobs to their preferred nodes. This paper also extensively evaluates their perfor-

mances.

In this paper, we have made the following contributions:

- We propose a new variation of gang scheduling, I/O aware gang scheduling (IOGS), which is aware of the jobs' spatial preferences, and always schedules jobs to their desired nodes.
- We quantitatively compare IOGS and gang scheduling (GS) under workloads with different I/O intensities. We find that GS is better than IOGS for workloads with low to medium I/O intensities, while IOGS is much better than gang scheduling for I/O intensive workloads.
- We propose a hybrid scheme, adaptive I/O aware gang scheduling (Adaptive-IOGS), which tries to combine the benefits of both schemes. We show that this adaptive scheme performs the best in many situations.
- Further, we propose another scheme, migration adaptive I/O aware gang scheduling (Migrate-IOGS), which combines migration technique with Adaptive-IOGS. Migration technique can move more jobs to their desired nodes, thus leading to better performance.
- Finally, we have shown that Migrate-IOGS delivers the best performance among the four across a wide range of workloads.

The rest of the paper is organized as follows: Section 2 describes the system and workload models used in this study. Section 3 presents all the proposed scheduling heuristics, and report their performance results. Section 5 presents our conclusions and possible direction for future work.

## 2 System and Workload Model

Large scale parallel applications are demanding large data inputs during their executions. Optimizing their I/O performance is becoming a critical issue. In this study, we set out to investigate the impact of application I/O intensity on the design of a job scheduler.

## 2.1 I/O Model

In the serverless file system [1, 2], each node serves as both compute node and I/O node. The data files in the system are partitioned according to some heuristics, and distributed on the disks of a subset of the nodes (clients that participate the file system). We use  $t_{IO}^{local}$  and  $t_{IO}^{remote}$  to denote the local and remote I/O costs respectively. In this paper, *I/O costs* are the costs associated with those I/O requests that are not satisfied in the file cache, and have to fetch data from disk.

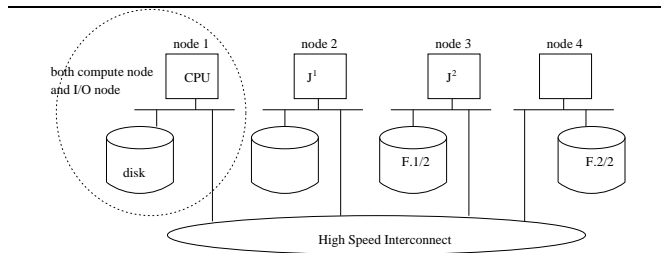


Figure 1: the I/O model example

Figure 1 illustrates such a system with 4 nodes, which are connected by high-speed interconnect. File  $F$  has two partitions,  $F.1/2$  and  $F.2/2$ , hosted by nodes 3 and 4 respectively. Parallel job  $J$  (with tasks  $J^1$  and  $J^2$ ) is running on nodes 2 and 3, and needs data from  $F$  during its execution. In this example, if the data needed by  $J^2$  belongs to partition  $F.1/2$ , then  $J^2$  spends much less time in I/O because all its requests can be satisfied by local disks, while  $J^1$  has to fetch the data from remote disks.

## 2.2 File Partition Model

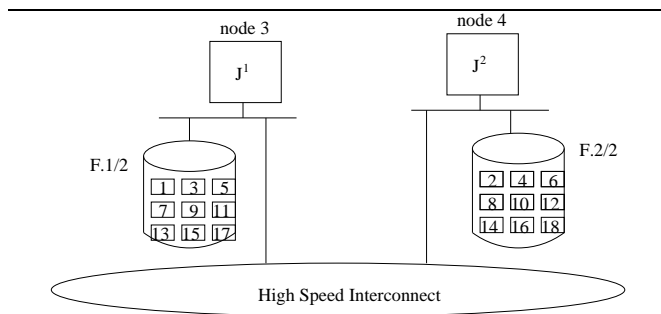


Figure 2: File partition Example

File partition plays an important role in deciding a task’s I/O portion. Even though a task is running on a node that hosts one partition of its file, it does not mean

that the task can enjoy lower local I/O costs because the data it needs may not belong to this partition. In the example shown in Figure 2,  $F$  has 18 blocks in total, and all the odd-numbered blocks belong to  $F.1/2$ , while the even-numbered blocks belong to  $F.2/2$ . Task  $J^1$  is co-located with all the odd-numbered blocks of  $F$ . Unfortunately,  $J^1$  needs the first 10 blocks of file  $F$ , which are evenly distributed between  $F.1/2$  and  $F.2/2$ . Thus, half of its I/O requests have to go to remote disks. This example shows that even if a job scheduler manages to assign tasks to the nodes where their files are hosted, their I/O requests may not be fully served by local disks. This observation suggests that we need to coordinate the file partitioning and applications’ data access pattern to better appreciate the results of good job schedulers.

As suggested by Corbett and Feitelson in [3], it is possible for the applications to pass their access patterns to the underlying file system, so that the system can partition the files accordingly, leading to a one-to-one mapping between the tasks and partitions. Thus, as long as the task is assigned to the appropriate node (hosting the corresponding partition), all its I/O accesses are local. In this study, we used the above approach. Although the hardware configurations in [3] and in this study are different, we believe that we can implement the same idea on our platform. Again, please note that our schemes can work with other partitioning heuristics as well. In that case, we need to compose the job access pattern model and quantify the ratio of the local I/Os out of the total number of I/Os made by each task.

As statistics in [14] show, file sharing across applications is rare, which makes the partitioning easier because we do not need to consider the cases where different applications have different access patterns if they share the same file. Further, the number of files accessed by each job does not affect the effectiveness of our scheduling schemes, because we can just simply partition every file a job accesses according to the corresponding access pattern.

## 2.3 Workload Model

We conduct a simulation based study of our scheduling strategies using synthetic workloads.

Before we present the workload model, we first discuss our parallel job model. As in [21], we assume that a job has a number of tasks, each task performing computation and I/O, and synchronizing with each other. The tasks of a job are co-scheduled in the time-slice so that

the synchronization between them incurs very low cost. Further, we employ a simple model assuming I/O accesses are evenly distributed throughout the execution.

We define the following parameters for each job:

1.  $t_i^a$ : arrival time of job  $i$ .
2.  $t_i^C$ : CPU time of job  $i$  (the total time spent on CPU if  $i$  runs alone).
3.  $n_i$ : number of tasks of job  $i$ .
4.  $d_i^I$ : the I/O requests interarrival time of job  $i$  (the gap between two I/O requests).
5.  $n_i^B$ : number of disk blocks per request of job  $i$ .
6.  $t_i^s$ : start time of job  $i$ .
7.  $t_i^f$ : finish time of job  $i$ .

Among all these parameters,  $t_i^a$  (job arrival time) and  $n_i$  (job size) are generated from stochastic models that fit actual workloads at the ASCI Blue-Pacific system in Lawrence Livermore National Laboratory (a 320-node RS/6000 SP) [11].

From the above parameters and our system model, we can derive the following:

1.  $r_i^{IO} = \frac{n_i^B \times t_{IO}^{local}}{n_i^B \times t_{IO}^{local} + d_i^I}$ : I/O intensity of job  $i$ .
2.  $n_i^I = \lceil \frac{t_i^C}{d_i^I} \rceil \times n_i^B$ : number of disk blocks accessed by job  $i$ .
3.  $t_i^e = t_i^C + \sum_{j=0}^{j < n_i^I} (l_{i,j} * t_{IO}^{local} + (1 - l_{i,j}) * t_{IO}^{remote})$ :  $i$ 's execution time on a dedicated setting, where  $l_{i,j}$  is 1 when the  $j$ th I/O of job  $i$  is local, and 0 otherwise. Next we need to elaborate on deciding the value of  $l_{i,j}$ , i.e., how to decide if an I/O is a local or remote. In our job model, we do not differentiate tasks of the same job. We assume they are identical. Thus, even though only one task of a job is not running on the file node, and the other tasks are all on the file nodes, the job, as a whole, is considered to have remote I/O accesses because the tasks from the same parallel job need to synchronize with each other and progress at the same pace.
4.  $t_i^r = t_i^f - t_i^s$ : response time for job  $i$ .
5.  $t_i^w = t_i^s - t_i^a$ : wait time for job  $i$ .

6.  $s_i = \frac{\max(t_i^r, T)}{\max(t_i^e, T)}$ : the slowdown for job  $i$ , where  $T$  is the time-slice for gang-scheduling. To reduce the statistical impact of very short jobs, it is common practice [7, 8] to adopt a minimum execution time. We adopt a minimum of one time slice. That is the reason for the  $\max(\cdot, T)$  terms in the definition of slowdown. Another thing to notice is that we cannot calculate slowdown using  $t_i^e$  to compare different scheduling schemes because different schemes result in different  $t_i^e$ . Instead, we use a normalized job execution time  $t_i^{e'} = t_i^C + \sum_{j=0}^{j < n_i^I} t_{IO}^{local}$  whose value is invariant from scheme to scheme.

Each synthetic workload contains 10000 parallel jobs whose parameters are described above. The characteristics of a workload can be described by two factors: load and I/O intensity. The load is decided by the job interarrival time ( $\lambda^{-1}$ ) and average job execution time ( $t^e$ ). By fitting the actual workloads at the ASCI Blue-Pacific system in Lawrence Livermore National Laboratory (a 320-node RS/6000 SP) [11], we have 9 different values for  $\lambda^{-1}$  and  $t^e$ , resulting in 81 workloads. The I/O intensity of a workload is decided by the average I/O intensity of its jobs. For each workload, we choose its I/O intensity ( $r^{IO}$ ) from the following seven values 0.1, 0.15, 0.2, 0.25, 0.33, 0.4, 0.5. Thus, we have  $81 \times 7 = 567$  workloads in total.

## 3 Scheduling Strategies

### 3.1 Gang-Scheduling

Gang scheduling [4, 6, 9, 10, 16] is a timed-shared parallel job scheduling technique. This technique virtualizes the physical machine by slicing the time axis into multiple virtual machines. Tasks of a parallel job are coscheduled to run in the same time-slices (same as virtual machines). The number of virtual machines is called the multiprogramming level (MPL) of the system. This multiprogramming level in general depends on how many jobs can be executed concurrently, but is typically limited by system resources. This approach opens more opportunities for the execution of parallel jobs, and is thus quite effective in reducing the wait time. Gang-scheduling has been used in the prototype GangLL job scheduling system developed by IBM Research for the ASCI Blue-Pacific machine at LLNL (a large scale parallel system spanning thousands of nodes [13]).

As an example, gang-scheduling an 8-processor system with a multiprogramming level of four is shown in 3. The figure shows the Ousterhout matrix that defines the tasks executing on each processor and each time-slice.  $J_i^j$  represents the  $j$ -th task of job  $J_i$ . The matrix is cyclic in that time-slice 3 is followed by time-slice 0. One cycle through all the rows of the matrix defines a *scheduling cycle*. Each row of the matrix defines an 8-processor virtual machine, which runs at 1/4th of the speed of the physical machine. We use these four virtual machines to run two 8-way parallel jobs ( $J_1$  and  $J_2$ ) and several smaller jobs ( $J_3, J_4, J_5, J_6$ ). All tasks of a parallel job are always coscheduled to run concurrently. This approach gives each job the impression that it is still running on a dedicated, albeit slower, machine. This type of scheduling is commonly called *gang-scheduling* [4]. Note that some jobs can appear in multiple rows (such as jobs  $J_4$  and  $J_5$ ).

	$P_0$	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	$P_6$	$P_7$
time-slice 0	$J_1^0$	$J_1^1$	$J_1^2$	$J_1^3$	$J_1^4$	$J_1^5$	$J_1^6$	$J_1^7$
time-slice 1	$J_2^0$	$J_2^1$	$J_2^2$	$J_2^3$	$J_2^4$	$J_2^5$	$J_2^6$	$J_2^7$
time-slice 2	$J_3^0$	$J_3^1$	$J_3^2$	$J_3^3$	$J_4^0$	$J_4^1$	$J_5^0$	$J_5^1$
time-slice 3	$J_6^0$	$J_6^1$	$J_6^2$	$J_6^3$	$J_4^0$	$J_4^1$	$J_5^0$	$J_5^1$

Figure 3: The scheduling matrix defines spatial and time allocation.

Every job arrival or departure constitutes a *scheduling event* in the system. For each scheduling event, a new scheduling matrix is computed for the system. Even though we analyze various scheduling strategies in this paper, they all follow an overall organization for computing that matrix, which can be divided into the following steps:

1. **CleanMatrix:** The first phase of a scheduler removes every instance of a job in the Ousterhout matrix that is not at its assigned home row. Removing duplicates across rows effectively opens the opportunity of selecting other waiting jobs for execution.
2. **CompactMatrix:** This phase moves jobs from less populated rows to more populated rows. It further increases the availability of free slots within a single row to maximize the chances of scheduling a large job.
3. **Schedule:** This phase attempts to schedule new jobs. We traverse the queue of waiting jobs as dic-

tated by the given priority policy until no further jobs can be fitted into the scheduling matrix.

4. **FillMatrix:** This phase tries to fill existing holes in the matrix by replicating jobs from their home rows into a set of replicated rows. This operation is essentially the opposite of **CleanMatrix**.

In the rest of this paper, we describe each scheduling strategy based on their implementations of these scheduling steps.

### 3.2 Plain Gang Scheduling (GS)

The simulation model is based on the implementation of the GangLL scheduler [13] on the Blue Pacific machine at Lawrence Livermore National Labs. For GS, we implement the four scheduling steps of Section 3.1 as follows:

1. **CleanMatrix:** The implementation of CleanMatrix is best illustrated with the following algorithm:

```

for i = first row to last row
  for all jobs in row i
    if row i is not home of job, remove it

```

It eliminates all occurrences of a job in the scheduling matrix other than the one in its home row. A job is assigned a home row when it is first scheduled into the system, and its home row may change when the matrix is re-calculated.

2. **CompactMatrix:** We implement the CompactMatrix step in GS according to the following algorithm:

```

do{
  for i = least to most populated row
    for j = most to least populated row
      for all jobs in row i
        if they can be moved to row j
          move and break
}while matrix changes

```

We traverse the scheduling matrix from the least populated row to the most populated row. We attempt to find new homes for the jobs in each row. The goal is to pack as many jobs in as few rows as possible.

3. **Schedule:** The Schedule phase for GS traverses the waiting queue in FCFS order. For each job, it looks for the row with the least number of free slots in the scheduling matrix that has enough free columns to hold the job. This corresponds to a best fit algorithm. This algorithm opens more opportunities for

the larger jobs to be accommodated into the matrix. The row to which the job is assigned becomes its home row. Within that row, we look for the columns that are free across multiple rows, increasing the likelihood of this job being replicated to other rows, which helps with the matrix utilization. We stop when the next job in the queue cannot be scheduled right away.

4. **FillMatrix:** We use the following algorithm in executing the FillMatrix phase.

```
do {
  for each job in starting time order
    for all rows in matrix,
      if job can be replicated in same columns
        do it and break
  } while matrix changes
```

### 3.3 IO-Aware Gang Scheduling (IOGS)

GS tries to reduce job wait times by accommodating them as soon as possible, and to keep the matrix efficiently utilized. This goal is attempted by its allocation heuristics, including searching for the target row with BestFit algorithm, and then searching for the target columns for the possibility of future replication. These optimizations would translate into performance improvement if the jobs do not have any co-location preferences. Nonetheless, this is not true for the motivating applications in this study. The job execution times are significantly affected by their locations. Suppose job  $i$  needs file  $f$  during its execution, and partitions of  $f$  reside on nodes  $n_{i_1}, n_{i_2}, \dots, n_{i_{n_i}}$ . These nodes are called *file nodes* of  $i$ .  $i$ 's execution time will be shortened if it is allocated to its file nodes. If an allocation scheme is not aware of this, and thus causes longer execution times, it may hurt the overall performance. To address this issue, we propose a new scheme, called I/O aware gang scheduling (IOGS). IOGS schedules every job to its file nodes, with the cost of longer wait times in the queue.

At every scheduling event, IOGS also re-calculates the matrix using the following four steps:

1. **CleanMatrix:** Same as in GS.
2. **CompactMatrix:** Same as in GS.
3. **Schedule:** The Schedule phase for IOGS traverses the waiting queue in FCFS order. For each job, it looks for the row with the least number of free slots where the job's file nodes are available. Similar to the Schedule phase for GS, this also corresponds to

a best fit algorithm. However, here, the criteria of "fit" is having all file nodes available (to the interests of applications), instead of having enough free columns as in GS (to the interests of matrix usage). Within that row, we choose the file nodes to schedule the job. Please note that the number of file nodes equals the corresponding job size (Section 2.2). We stop when the next job in the queue cannot find a row that has its files nodes available.

4. **FillMatrix:** Same as in GS.

## 3.4 Comparing IOGS and GS

### 3.4.1 Experimental Setup

It is always a challenging issue to design a good scheduler because numerous parameters are involved, and these parameters can have an infinite design space. In sections 3.4.2-3.6, without loss of generality, we run the experiments using the following values for the parameters unless explicitly stated otherwise.

- $t_{IO}^{local} = 0.03$  second,  $t_{IO}^{remote} = 0.09$  second.
- workload. As mentioned in Section 2.3, we have  $9 \times 9 \times 7$  workloads corresponding to 9 different values of  $\lambda^{-1}$ , 9 values of  $t^e$ , and 7 values of  $r^{IO}$ . Due to space limit, we only show results for 27 workloads by having only one  $\lambda^{-1}$  and three  $r^{IO}$ . Since the goal of this study is to investigate the impact of workload I/O intensity, we organize these 27 workloads into three classes according to their I/O intensity: high I/O intensity ( $r^{IO} = 0.5$ ), medium I/O intensity ( $r^{IO} = 0.33$ ) and low I/O intensity ( $r^{IO} = 0.1$ ). Within each class, we have nine workloads with increasing load (increasing  $t^e$ ). For all 27 workloads, we have  $\lambda^{-1} = 583$  seconds.
- system parameters. Time quantum ( $T$ ) is 200 seconds. Maximum multiprogramming level (MPL) is 5.

### 3.4.2 Results

Looking at the heuristics of the two schemes, we know that GS tries to maximize the matrix usage, while IOGS provides local I/O accesses. Before we report our results, let us first take a closer look at what happens to a job so that we can understand the pros and cons for both schemes better for the motivating workloads.

Once a job starts executing, how soon it could finish is determined by its execution time (job duration on a dedicated setting) as well as how much CPU time it receives. A job’s execution consists of two components: computation and I/O (Section 2.3), and the latter is affected by whether the data can be accessed locally or not. IOGS yields short execution times by scheduling jobs to their files nodes. On the other hand, GS allocation scheme opens more opportunities for the jobs to be replicated to multiple rows, thus receiving more CPU time. Moreover, in IOGS, a job is scheduled into the matrix only if all its file nodes are free during at least one time quantum, leading to longer job wait time in the queue.

Intuitively, the relative importance of these two factors should be related to the I/O intensity of the workload.

The results of the two scheduling strategies for the three classes of workloads are shown in Figures 4(a-f).

For high I/O intensity workload  $r^{IO}=0.5$  (Figures 4(a),(d)), IOGS is clearly better than GS for all nine different loads for both performance metrics (slow down and wait times). With respect to slowdown, IOGS is always 100% better. Under most of the loads (i.e.,  $t^{ef} \geq 2000$  seconds), the improvement is over 200%. With respect to wait time, the improvement of IOGS over GS is not as obvious as under low loads because the wait time is low for both schemes, and we observe significant improvement under very high loads (i.e.,  $t^{ef} \geq 2600$  seconds). When the workload is I/O intensive (more than 50% of I/O ratio), the drawbacks of longer execution times due to remote I/O accesses far outweigh the benefits of more CPU time.

For the other two classes of workloads ( $r^{IO}=0.33,0.1$ ), we observe the opposite trend. GS is always better. The reason is, for medium to low I/O intensity workloads, because of the less importance of I/O, the benefits of shorter execution times are overshadowed by the drawbacks of less CPU resources received.

From this set of experiments, we learn that neither IOGS nor GS is the best for all the workloads. For I/O intensive workloads, IOGS is better, while GS is a much better choice for low to medium I/O workloads.

### 3.5 Adaptive I/O Aware Gang Scheduling (Adaptive-IOGS)

From the previous section, we learn that IOGS delivers better performance for I/O intensive workloads that motivated this work, as it yields much shorter execution times by scheduling jobs to their file nodes. However,

the downside of doing so is that it causes fragmentation in the scheduling matrix. Figure 5 illustrates this prob-

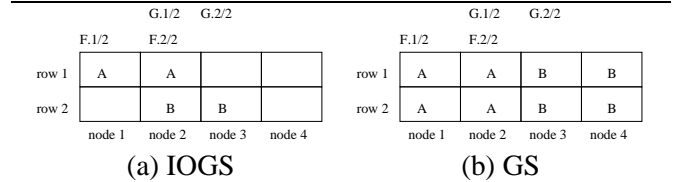


Figure 5: An example: job  $A$  needs file  $F$ , and job  $B$  needs file  $G$ .  $F$  is distributed on nodes 1 and 2, and  $G$  on 2 and 3.

lem. In this example, node 2 hosts partitions of both files ( $F$  and  $G$ ), which means both jobs will be assigned to node 2. As a result,  $A$  and  $B$  will be allocated to different rows (Figure 5(a)), and neither can be replicated to the second row because they share a common node. The matrix utilization in this case is only 50%. Moreover, this allocation heuristic leaves two small holes in both rows. If job  $C$ , with 3 tasks, arrives, it has to be kept waiting. On the other hand, GS can utilize the matrix 100%, as shown in Figure 5(b). Further, IOGS makes a job wait if it cannot find a row where all its file nodes are free, even though there are rows that have enough free slots to hold the job, which makes it more difficult to accommodate a job, thus leading to another type of matrix fragmentation: a job cannot start execution even though there are free slots because of strict allocation heuristics. To make matters even worse, because this job is the first in the waiting queue, all the jobs that arrive later than it have to wait as well, until its file nodes are freed after the jobs that occupy them finish.

Next, we propose a new strategy that alleviates the fragmentation caused by IOGS, while keeping its benefit for I/O intensive workloads. We call it Adaptive-IOGS. Adaptive-IOGS tries to first allocate a job as IOGS does. If no row has all the files nodes free, then we allocate this job as GS does, looking for rows/columns so that the matrix usage is maximized, instead of making it wait. If IOGS and GS are the two extremes, then this approach stands in the middle, using an adaptive allocation heuristic.

A formal description of its four scheduling steps is as follows:

1. **CleanMatrix:** Same as in IOGS and GS.
2. **CompactMatrix:** Same as in IOGS and GS.
3. **Schedule:** The Schedule phase for Adaptive-IOGS

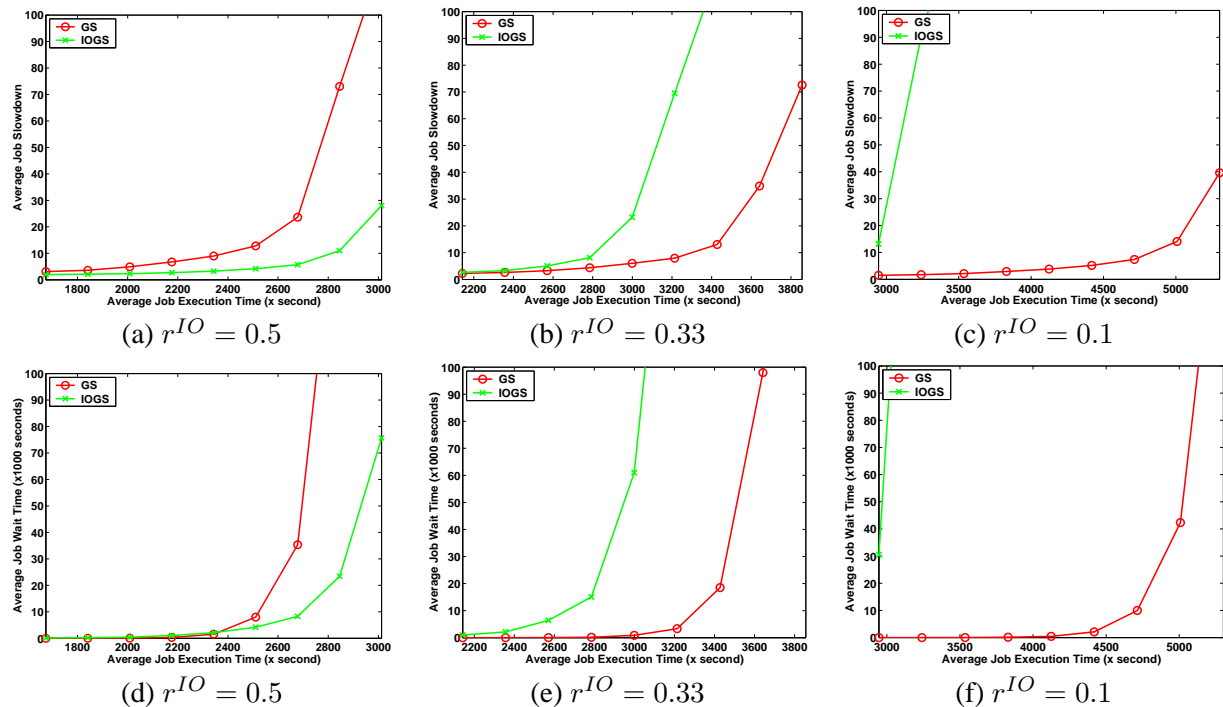


Figure 4: Comparison of GS and IOGS.

traverses the waiting queue in FCFS order. For each job, it first looks for the row with the least number of free slots where the job’s file nodes are available. If such a row is found, then we schedule the job to its file nodes within that row. Otherwise, we go back to GS: we look for the row with the least number of free slots that has enough free slots to hold the job. Within that row, we look for the columns/nodes that are empty across the most number of rows to open opportunities of replicating the job to other rows.

- 4. FillMatrix:** Adaptive-IOGS puts every running job into one of the two queues: local job (jobs that are assigned to their file nodes) or remote job (jobs that are not assigned to their file nodes). The FillMatrix phase for Adaptive-IOGS examines local jobs first to see whether they can be replicated to other rows. After no local job can be replicated, it examines all the remote jobs. Within each queue, the jobs are examined in the FCFS order.

We conduct experiments to compare Adaptive-IOGS with both IOGS and GS, using the same configurations described in Section 3.4.1. We compare these three strategies under workloads with three different I/O intensities. The results are shown in Figures 6(a-f).

We observe that Adaptive-IOGS is the best among the

all for (i) medium I/O workloads (for all nine different loads), and (ii) I/O intensive workloads for low to medium loads (i.e.,  $t^{el} \leq 2500$  seconds). Further, it is very close to the best under low I/O intensity workloads.

For I/O intensive workloads with  $r^{IO} \geq 0.5$ , under low to medium loads, Adaptive-IOGS enjoys both benefits of local I/O accesses and high matrix utilization. For instance, under load  $t^{el} = 1800$  seconds, Adaptive-IOGS manages to assign 9101 local jobs out of 10000, while GS only has 20 local jobs. (IOGS has 10000 local jobs out of 10000.) Please note that local I/O accesses are the key for I/O intensive workloads, and we find that Adaptive-IOGS can do almost as well as IOGS, with the additional benefit of less fragmentation. On the other hand, when the loads become high, very few jobs can be assigned to their files nodes when they are first scheduled. Adaptive-IOGS thus approaches GS, and its performance starts to degrade dramatically. Quantitatively, for the highest load  $t^{el} = 3000$  seconds, Adaptive-IOGS only has 831 local jobs, which is comparable to GS with 5 local jobs.

For medium I/O workloads  $r^{IO}=0.33$ , the key is efficient matrix usage, which is achieved by Adaptive-IOGS by scheduling the jobs of which the files nodes are not free to the rows/columns that minimize matrix fragmentation. In addition, Adaptive-IOGS can reduce the execu-



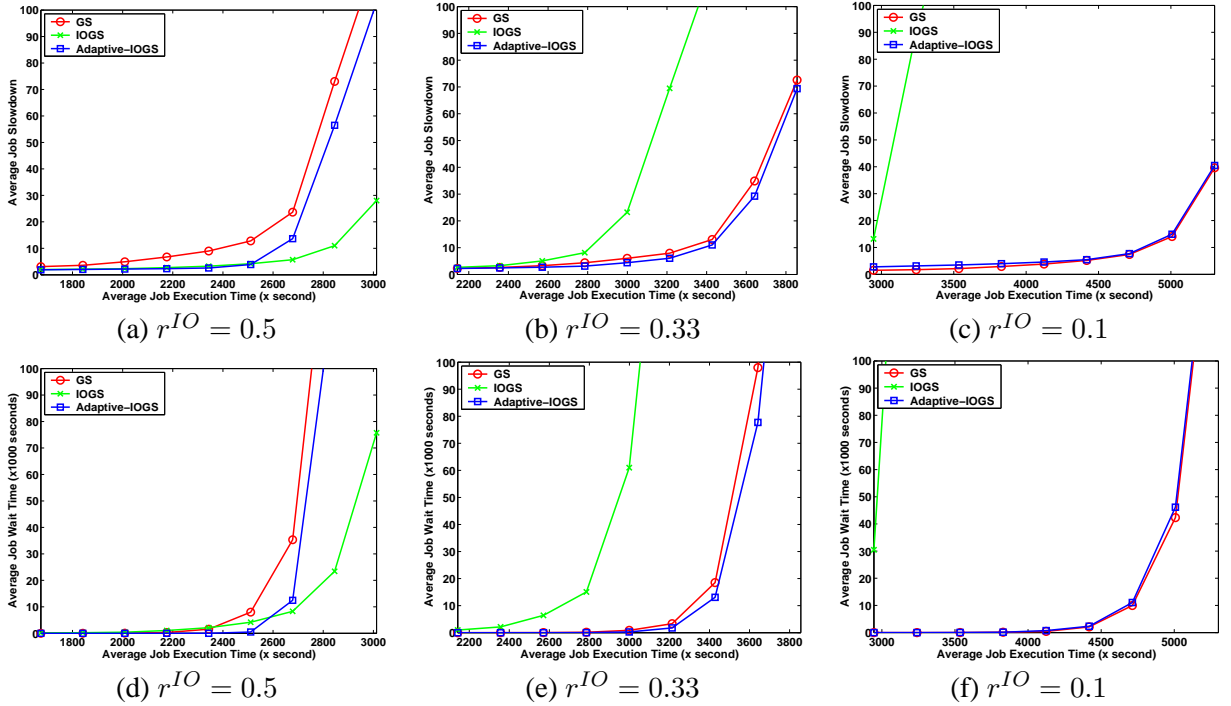


Figure 6: Comparison of GS, IOGS and Adaptive-IOGS.

tion times of the jobs that are assigned to the local nodes. As a result, it has the best performance throughout the entire load range for medium I/O workloads.

For workloads with low I/O intensity  $r^{IO}=0.1$ , I/O cost is a small fraction of job execution time. Co-locating jobs with their files is thus not important, and it will even hurt the performance because of the high system fragmentation. This explains why IOGS performs poorly. On the other hand, Adaptive-IOGS adopts GS Schedule heuristics often enough that it is very close to GS. Figures 6(c), (f) also show that Adaptive-IOGS is closer to GS under high loads because it is more unlikely to schedule a job to its file nodes at the first time when the load is high.

Looking at all three workloads, Adaptive-IOGS is either the best, or very close to the best, except under high loads for I/O intensive workloads.

### 3.6 Migration Adaptive I/O Aware Gang Schedule (Migrate-IOGS)

From the earlier results, we conclude that the challenge here is *how to achieve local I/O requests and low matrix fragmentation at the same time*. These two goals seem contradicting with each other. IOGS and GS works on one of the two goals respectively. Adaptive-IOGS tries to balance the two factors by adaptively change the goal

based on the context. Although Adaptive-IOGS outperforms both GS and IOGS in most of the scenarios, it is not as good as IOGS under high loads for I/O intensive workloads. However, high loads of I/O intensive workloads are the most interesting scenarios to this study. We need to further investigate on how to improve its performance.

Adaptive-IOGS is approaching GS under high loads. In these scenarios, the matrix is already full, so few jobs can find their file nodes available when they are first scheduled. Thus, most of the jobs are scheduled into the system using GS Schedule heuristic. These jobs stay on the remote nodes once they are allocated there, paying higher remote I/O costs throughout the execution, even though their file nodes might have been freed immediately after they are scheduled.

	F.1/2	F.2/2		
row 1	A	A	C	C
row 2	D	B	D	
	node 1	node 2	node 3	node 4

Figure 7: An example illustrating Adaptive-IOGS: job  $D$  needs file  $F$ , which is distributed on nodes 1 and 2.

In the example shown by Figure 7, When  $D$  was first scheduled, its file nodes (nodes 1 and 2) were occupied by jobs  $A$  in row 1 and  $B$  in row 2, so it was allocated to row 2, on nodes 1 and 3. After a short period of time, job  $A$  finishes execution and leaves the system,  $D$ 's file nodes are now free in row 1, but  $D$  still stays on nodes 1 and 3, and pays remote I/O costs, which can also be considered as one type of fragmentation. In order to address this type of fragmentation, we propose to use migration technique. In the above example, if we can migrate job  $D$  from row 2 to row 1, to its file nodes,  $D$ 's remaining execution time will be significantly shortened.

The next scheduling strategy we propose uses migration technique to enhance Adaptive-IOGS, and we call the new scheme Migration Adaptive I/O Aware Gang Scheduling (Migrate-IOGS). Every time when a job leaves the system, we will examine each remote job to see if its file nodes are freed by then. If so, we move the job to its file nodes. We put this migration phase after we clean the matrix, and before we compact the scheduling matrix.

The formal description of the scheme is as follows:

1. **CleanMatrix:** Same as in Adaptive-IOGS, IOGS and GS.
2. **MigrateMatrix:** MigrateMatrix is a new phase added to Migrate-IOGS, where remote jobs can be migrated to their file nodes if those nodes are freed due to this new scheduling event.

As explained in [18, 20], there is a cost associated with migration. In systems like IBM RS/6000 SP, migration involves a checkpoint/restart operation. Let  $C$  denote the total migration cost, including checkpointing and restarting a process. After a task is migrated to another node, it can resume its execution only after  $C$ . Further, even though only one task of a job is migrated, we make the other tasks also wait for  $C$  to model the synchronization between the tasks. Although all tasks can perform a checkpoint in parallel, resulting in a  $C$  that is independent of job size, there is a limit to the capacity and bandwidth that the file system can accept. Therefore we introduce a parameter  $Q$  that controls the maximum number of tasks that can be migrated in any time-slice.

In this phase, for each remote job, Migrate-IOGS looks for the row that has all its file nodes available. If such a row is found, then Migrate-IOGS

moves the job there so that its remaining execution will be significantly shortened. The job needs to pay the migration overhead only once (at the beginning of the next time quantum), and the value of  $C$  is negligible compared to the savings on the job's execution time, so we do not expect to observe any negative impact of migration. On the other hand, if there are more than one such rows (in which the file nodes are free), we have two options: (1) randomly choose one row, or (2) we can choose the row with the least number of free slots so that we can accommodate larger waiting jobs. In our study, we simply chose the first option due to the time limit. We will conduct further investigations for the final version.

We implement the MigrateMatrix step in Migrate-IOGS according to the following algorithm:

```
do{
  for each running job that is remote
    for all rows in matrix,
      if the job can become local in that row
        move it and break
}while matrix changes
```

3. **CompactMatrix:** Same as in Adaptive-IOGS, IOGS and GS.
4. **Schedule:** Same as in Adaptive-IOGS.
5. **FillMatrix:** Same as in Adaptive-IOGS.

Migrate-IOGS is intended to shorten the job execution times and reduce matrix fragmentation at the same time. However, two factors may limit its effectiveness in achieving this goal: the migration cost ( $C$ ), and the maximum number of tasks that can be migrated in any time-slice ( $Q$ ). We conduct experiments to study the impact of these two factors, and the results are shown in Figures 8(a-f). Within each class (we have three classes), we choose the highest load (the one with the highest execution time).

Keeping the time quantum length as 200 seconds, we vary the migration cost from 0, 10, 20, and 30 seconds, where  $C = 0$  represents the ideal situation where migration is free. For all three workloads, we observe that the impact of migration cost is negligible. The reason is that each job being migrated needs to pay the overhead only once (at the beginning of the next time quantum), and this overhead is negligible compared to the savings on the job's execution time. We adopt  $C = 20$  seconds in the following experiments.

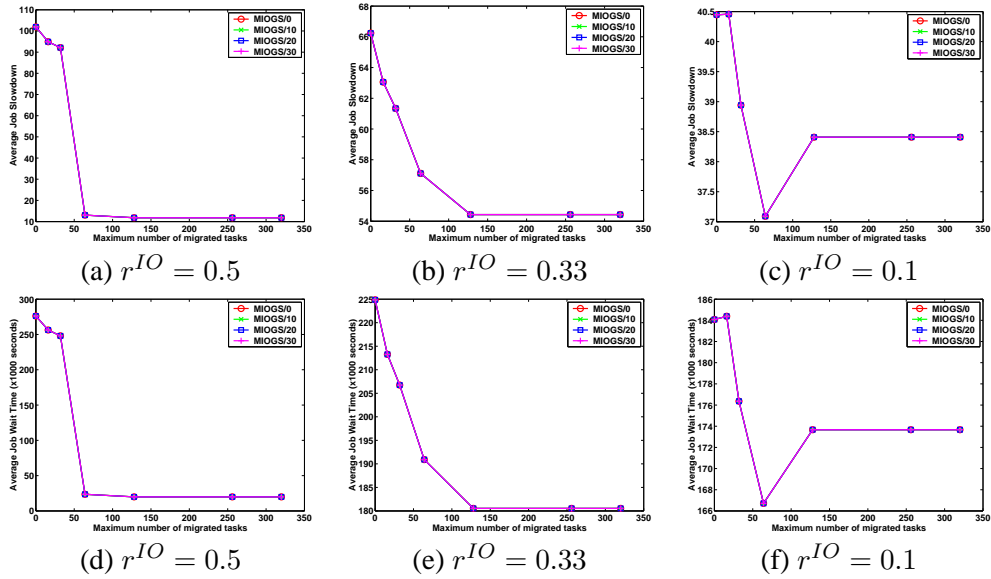


Figure 8: The impact of migration cost and maximum number of migrated tasks. For all three workloads, we have  $t^e = 3102$  seconds.

With respect to the limit of migrated tasks  $Q$ , we vary  $Q = 0, 16, 32, 64, 128, 256$ , and  $320$ , where  $Q = 0$  means it goes back to Adaptive-IOGS, and  $Q = 320$  means there is no limit because we only have 320 nodes in the system. We find that the value of  $Q$  has a big impact on the performance. For I/O intensive workloads,  $Q=320$  improves the performance at least 90% compared to  $Q=0$ . However, most of the benefits is accomplished at  $Q=64$ , which is a reasonably small number compared to the total number of nodes in the system. We observe similar trends in medium I/O workloads. For low I/O workloads, migrating more than 64 tasks makes the performance worse. In summary, if we migrate 64 tasks, we can achieve most of the benefits that migration technique can have. At the same time, migrating upto 64 nodes during any time-slice will not saturate the limit of the system capacity. As a result, we adopt  $Q = 64$  in the following experiments.

From these results, we conclude that Migrate-IOGS is a practical approach.

### 3.7 Comparing Migrate-IOGS, Adaptive-IOGS, IOGS, and GS

Finally, we can compare all four schemes under the same configurations. As in Sections 3.4.2 and 3.5, we use three workloads:  $r^{IO}=0.5$ ,  $r^{IO}=0.33$ , and  $r^{IO}=0.1$ . In the experiments, we choose migration cost as 10% of the time quantum ( $10\% \times 200 = 20$  seconds), and we assume no

more than 64 tasks can be migrated during a quantum.

For medium to high I/O workloads (figures 9(a), (b), (d), and (e)), Migrate-IOGS outperforms the other three scheduling strategies. With respect to I/O intensive workloads, under high loads (i.e.,  $t^{el} \geq 2500$  seconds), Migrate-IOGS improves all three performance metrics by 100% compared to the second best scheme, IOGS. With the help of migration, Migrate-IOGS has 6524 local jobs at load  $t^{el} = 3000$  seconds, which is a significant improvement compared to its counter part without migration (Adaptive-IOGS) which only has 831 local jobs. With respect to medium I/O workloads, it improves the performance by at least 60% for all three metrics compared to the second best scheme - Adaptive-IOGS.

For low I/O workloads (figures 9 (c) and (f)), for the same reason mentioned in section 3.5, Migrate-IOGS is not as good as GS. However, the difference is very small, especially at high loads.

## 4 Impact of workloads

In the earlier sections, we assume that all the jobs in a workload are homogeneous. For example, any job  $i$  in the I/O intensive workload ( $r^{IO}=0.5$ ), has  $r_i^{IO}=0.5$ . In this section, we consider a mixed workload including jobs with different I/O intensities. Namely, every job in the mixed workload is equally likely to have I/O intensity  $r^{IO}$  of 0.1, 0.33, or 0.5.

Figures 10(a),(b) show the results comparing these

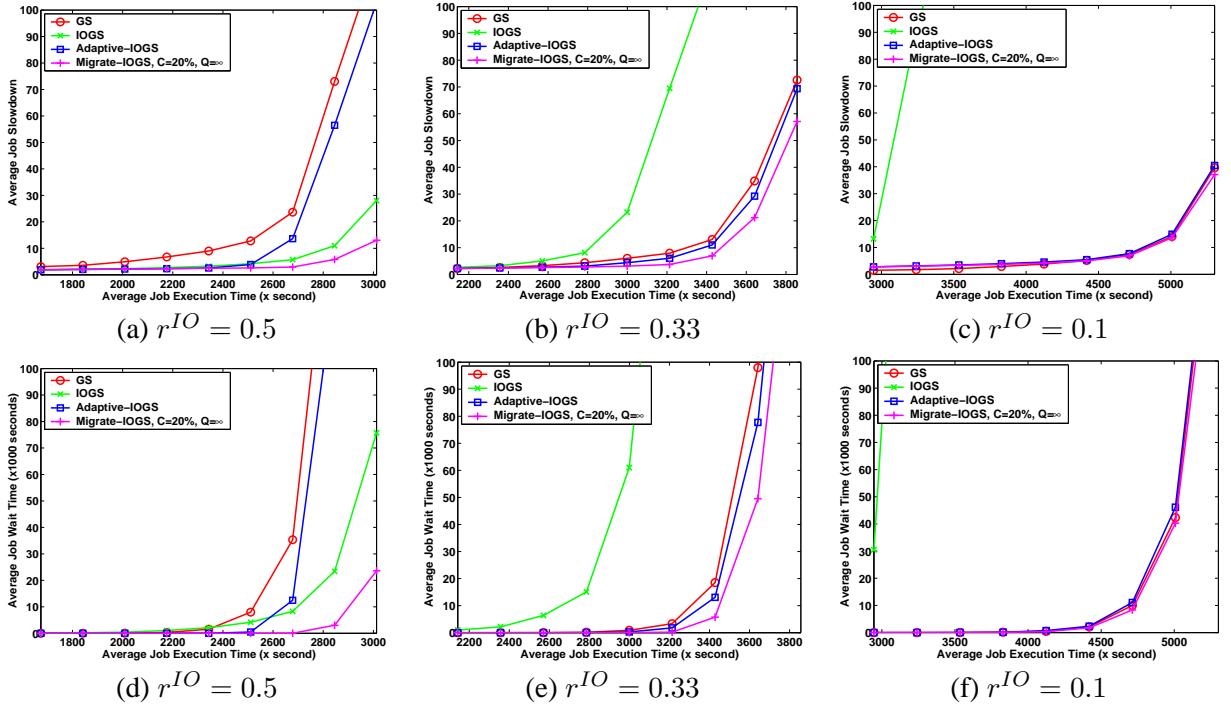


Figure 9: Comparison of GS, IOGS, Adaptive-IOGS and Migrate-IOGS. For Migrate-IOGS, we assume the migration cost 10% of the time quantum, and not more than 64 tasks can be migrated in the same quantum.

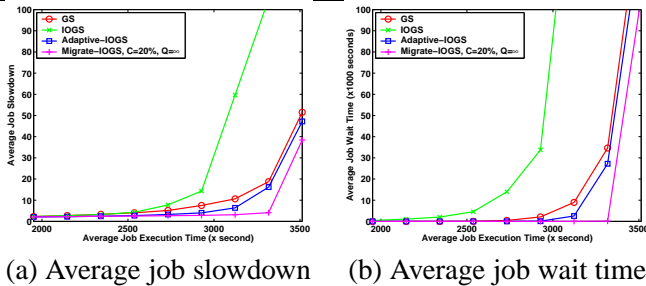


Figure 10: Comparison of GS, IOGS, Adaptive-IOGS and Migrate-IOGS under a mixed workload.

four schemes under such a mixed workload. We observe that Migrate-IOGS outperforms the other three for both metrics.

## 5 Concluding Remarks

This paper studies the impact of job I/O intensity on the design of a job execution scheduler. Nowadays, scientific applications are becoming more complex and more I/O demanding than ever. For such applications, the system with dedicated I/O nodes does not provide enough scalability. Rather, a serverless approach [1, 2] is a viable alternative. However, with the serverless approach, a job’s execution

time is decided by whether it is co-located with the file blocks it needs.

Gang scheduling (GS), which is widely used in supercomputing centers to schedule parallel jobs, is completely not aware of the application’s spatial preferences. We find that gang scheduling does not do a good job if the applications are I/O intensive.

We propose a new scheme, called IOGS (I/O aware gang scheduling), which assigns jobs to their file nodes, leading to a much shortened execution. Our results show that this scheme can improve the performance of I/O intensive workloads significantly, while it does not fare well for workloads with lower I/O intensity.

Further, we propose Adaptive-IOGS (Adaptive I/O aware gang scheduling), which takes the middle ground between gang scheduling and IOGS. Based on the matrix utilization, it adaptively assigns jobs either to their files nodes or just random nodes. We show that Adaptive-IOGS is the best for most of the workloads, except for those with high I/O intensity, and imposing very high load on the system.

Finally, we integrate migration technique to Adaptive-IOGS, by migrating jobs to their file nodes during their execution if those nodes are freed by others, leading to Migrate-IOGS (Migration adaptive I/O aware gang

scheduling). Our results clearly show that Migrate-IOGS outperforms the other three significantly for a wide spectrum of workloads.

## References

- [1] T. Anderson, M. Dahlin, J. Neeffe, D. Roselli, D. Patterson, and R. Wang. Serverless Network File Systems. *ACM Transactions on Computer System*, 14(1):41–79, 1996.
- [2] W. J. Bolosky, J. R. Douceur, D. Ely, and M. Theimer. Feasibility of a serverless distributed file system deployed on an existing set of desktop PCs. In *Proceedings of the ACM SIGMETRICS 2000 Conference on Measurement and Modeling of Computer Systems*, pages 34–43, 2000.
- [3] P. F. Corbett and D. G. Feitelson. The Vesta parallel file system. *ACM Transactions on Computer System*, 14(3):225–264, 1996.
- [4] D. G. Feitelson. A Survey of Scheduling in Multiprogrammed Parallel Systems. Technical Report RC 19790 (87657), IBM T. J. Watson Research Center, October 1994.
- [5] D. G. Feitelson. A Survey of Scheduling in Multiprogrammed Parallel Systems. Technical Report Research Report RC 19790(87657), IBM T. J. Watson Research Center, October 1994.
- [6] D. G. Feitelson and M. A. Jette. Improved Utilization and Responsiveness with Gang Scheduling. In *IPPS'97 Workshop on Job Scheduling Strategies for Parallel Processing*, volume 1291 of *Lecture Notes in Computer Science*, pages 238–261. Springer-Verlag, April 1997.
- [7] D. G. Feitelson, L. Rudolph, U. Schwiegelshohn, K. C. Sevcik, and P. Wong. Theory and Practice in Parallel Job Scheduling. In *IPPS'97 Workshop on Job Scheduling Strategies for Parallel Processing*, volume 1291 of *Lecture Notes in Computer Science*, pages 1–34. Springer-Verlag, April 1997.
- [8] D. G. Feitelson and A. M. Weil. Utilization and predictability in scheduling the IBM SP2 with backfilling. In *12th International Parallel Processing Symposium*, pages 542–546, April 1998.
- [9] H. Franke, P. Pattnaik, and L. Rudolph. Gang Scheduling for Highly Efficient Multiprocessors. In *Sixth Symposium on the Frontiers of Massively Parallel Computation, Annapolis, Maryland*, 1996.
- [10] N. Islam, A. L. Prodromidis, M. S. Squillante, L. L. Fong, and A. S. Gopal. Extensible Resource Management for Cluster Computing. In *Proceedings of the 17th International Conference on Distributed Computing Systems*, pages 561–568, 1997.
- [11] J. Jann, P. Pattnaik, H. Franke, F. Wang, J. Skovira, and J. Riordan. Modeling of Workload in MPPs. In *Proceedings of the 3rd Annual Workshop on Job Scheduling Strategies for Parallel Processing*, pages 95–116, April 1997. In Conjunction with IPPS'97, Geneva, Switzerland.
- [12] X. Ma, J. Jiao, M. Campbell, and M. Winslett. Flexible and efficient parallel i/o for large-scale multi-component simulations. In *Proceedings of The 4th Workshop on Parallel and Distributed Scientific and Engineering Computing with Applications, in conjunction with the 2003 International Parallel and Distributed Processing Symposium*, 2003.
- [13] J. E. Moreira, H. Franke, W. Chan, L. L. Fong, M. A. Jette, and A. Yoo. A Gang-Scheduling System for ASCI Blue-Pacific. In *High-Performance Computing and Networking, 7th International Conference*, volume 1593 of *Lecture Notes in Computer Science*, pages 831–840. Springer-Verlag, April 1999.
- [14] N. Nieuwejaar and D. Kotz. Low-level interfaces for high-level parallel I/O. In *Proceedings of the IPPS '95 Workshop on Input/Output in Parallel and Distributed Systems*, pages 47–62, April 1995.
- [15] J. K. Ousterhout. Scheduling Techniques for Concurrent Systems. In *Third International Conference on Distributed Computing Systems*, pages 22–30, 1982.
- [16] K. Suzaki and D. Walsh. Implementation of the Combination of Time Sharing and Space Sharing on AP/Linux. In *IPPS'98 Workshop on Job Scheduling Strategies for Parallel Processing*, March 1998.
- [17] Y. Wiseman and D. G. Feitelson. Paired Gang Scheduling. *IEEE Transactions on Parallel and Distributed Systems*.
- [18] Y. Zhang, H. Franke, J. Moreira, and A. Sivasubramaniam. The Impact of Migration on Parallel Job Scheduling for Distributed Systems. In *Proceedings of 6th International Euro-Par Conference Lecture Notes in Computer Science 1900*, pages 245–251, Aug/Sep 2000.
- [19] Y. Zhang, H. Franke, J. Moreira, and A. Sivasubramaniam. Improving Parallel Job Scheduling by Combining Gang Scheduling and Backfilling Techniques. In *Proceedings of the International Parallel and Distributed Processing Symposium*, pages 133–142, May 2000.
- [20] Y. Zhang, H. Franke, J. Moreira, and A. Sivasubramaniam. An Integrated Approach to Parallel Scheduling Using Gang-Scheduling, Backfilling and Migration. *IEEE Transactions on Parallel and Distributed Systems*, 14(3):236–247, March 2003.

- [21] Y. Zhang, A. Sivasubramaniam, J. Moreira, and H. Franke. A Simulation-based Study of Scheduling Mechanisms for a Dynamic Cluster Environment. In *Proceedings of the ACM 2000 International Conference on Supercomputing*, pages 100–109, May 2000.