



DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Job Scheduling for Adaptive Applications in Future HPC Systems

Nishanth Nagendra





DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Job Scheduling for Adaptive Applications in Future HPC Systems

Job Scheduling für Adaptive Anwendungen auf Zukünftigen HPC Systemen

Author:	Nishanth Nagendra
Supervisor:	Prof. Dr. Michael Gerndt
Advisor:	Isaías A. Comprés Ureña
Submission Date:	Jul 15, 2016



I confirm that this master's thesis in informatics is my own work and I have documented all sources and material used.

Munich, Jul 15, 2016

Nishanth Nagendra

Acknowledgments

I would like to sincerely thank *Isaías A. Comprés Ureña* for his constant support and guidance throughout the thesis through our long and fruitful discussions. He has been extremely patient and generous in providing his time for discussions on topics that were challenging for me but necessary to understand for working in this research direction. The constructive discussions along with his advice have been really valuable to me as a student.

I am extremely grateful to *Prof. Dr. Michael Gerndt* for providing me with this wonderful opportunity to be involved in such an interesting project and to work very closely with his research group. He has been very supportive throughout the length of the time i have been involved as a student working in the Chair for Architecture of Parallel and Distributed Systems. Despite his busy schedule, I have had the opportunity for numerous discussions with him that have helped me to gain critical advice for my thesis and other related tasks.

Abstract

Invasive Computing is a novel paradigm for the design and resource-aware programming of future parallel computing systems. It enables the programmer to write resource aware programs and the goal is to optimize the program for the available resources. Traditionally, parallel applications implemented using MPI are submitted with a fixed number of MPI processes to execute on a HPC (High Performance Computing) system. This results in a fixed allocation of resources for the job. Modern techniques in scientific computing such as AMR (Adaptive Mesh Refinement) result in applications exhibiting complex behaviors where their resource requirements change during execution. Invasive MPI is an ongoing research effort to provide MPI extensions for the development of Invasive MPI applications that will result in jobs which are resource-aware for the HPC systems and can utilize such AMR techniques. Unfortunately, using only static allocations result in these applications being forced to execute using their maximum resource requirements that may lead to inefficient resource utilization. In order to support such kind of parallel applications at HPC centers, there is an urgent need to investigate and implement extensions to existing resource management systems or develop a new system. This thesis has extended the previous work during which a negotiation protocol was developed for the integration of invasive resource management into existing batch systems. In this work, we have explored the idea of separating the concerns of batch and runtime scheduling into two different software layers / components in contrast to the existing batch systems where both are merged together. Specifically, this thesis has investigated and implemented a job scheduling algorithm by separation of concerns in accordance with the new protocol developed earlier for supporting such an invasive resource management. An early prototype that can simulate the negotiation between batch and runtime scheduler using their respective scheduling algorithms for a HPC workload comprising different job types has been accomplished.

Contents

Acknowledgments	iii
Abstract	iv
1 Introduction	1
1.1 Invasive Computing	3
1.2 Dynamic Resource Management	3
1.3 Document Structure	5
2 Related Work	7
3 Invasive Computing	10
3.1 Traditional Resource Management	12
3.1.1 Classification	12
3.1.2 Job Scheduling	13
3.1.3 SLURM	16
3.2 Resource Aware Programming	19
3.2.1 Job Classification	19
3.2.2 Invasive Programming Models	20
3.3 Invasive Resource Management	23
3.3.1 Invasive MPI	23
3.3.2 Resource Management Extensions	25
4 Architecture	28
4.1 Dynamic Resource Management	31
4.1.1 Invasive Batch Scheduler	32
4.1.2 Invasive Runtime Scheduler	33
4.1.3 iMPI Process Manager	34
4.2 Negotiation Protocol	35
4.2.1 Protocol Sequence Diagrams	35
4.3 Invasive Jobs	38

5	Design	39
5.1	Entities	39
5.2	Scheduling Algorithms	40
5.2.1	Batch Scheduling	41
5.2.2	Runtime Scheduling	44
5.3	Negotiation	50
6	Implementation	53
6.1	Plugin	53
6.2	Data Structures	54
6.3	Important APIs	59
6.4	Control Flow Diagrams	62
7	Evaluation	65
7.1	Method of Evaluation	67
7.2	Setup	68
7.3	Experiments and Results	69
8	Conclusion and Future Work	70
8.1	Future Work	70
	Bibliography	71

1 Introduction

Over the last two decades, the landscape of Computer Architecture has changed radically from sequential to parallel . Due to the limiting factors of technology we have moved from single core processors to multicore processors having a network interconnecting them. Traditionally, the approach of designing algorithms has been sequential, but designing algorithms in parallel is gaining more importance now to better utilize the computing power available at our disposal. Another important trend that has changed the face of computing is an enormous increase in the capabilities of the networks that connect computers with regards to speed, reliability etc. These trends make it feasible to develop applications that use physically distributed resources as if they were part of the same computer. A typical application of this sort may utilize processors on multiple remote computers, access a selection of remote databases, perform rendering on one or more graphics computers, and provide real-time output and control on a workstation. Computing on networked computers ("Distributed Computing") is not just a subfield of parallel computing as the basic task of developing programs that can run on many computers at once is a parallel computing problem. In this respect, the previously distinct worlds of parallel and distributed computing are converging.

As technology advances, we have new applications that demand larger computing capabilities which push the limits of technology giving rise to newer advancements. The performance of a computer depends directly on the time required to perform a basic operation and the number of these basic operations that can be performed concurrently. A metric used to quantify the performance of a computer is FLOPS (floating point operations per second). The time to perform a basic operation is ultimately limited by the "clock cycle" of the processor, that is, the time required to perform the most primitive operation. The term *High Performance Computing (HPC)* refers to the practice of aggregating computing power (multiple nodes with processing units interconnected by a network in a certain topology) or the use of parallel processing for running advanced application programs efficiently, reliably and quickly. The term applies especially to systems that function above a *teraflop* or 10^{12} floating-point operations per second. The term HPC is occasionally used as a synonym for Supercomputer that works at more than a *petaflop* or 10^{15} floating-point operations per second. Future systems will reach *exaflop* or 10^{18} floating-point operations per second. The most common users

of HPC systems are scientific researchers, engineers, government agencies including the military, and academic institutions. In general, HPC systems can refer to Clusters, Supercomputers, Grid Computing etc. and they are usually used for running complex applications.

A *Batch System* is used to manage the resources in a HPC System. It is a middleware that comprises of two major components namely the *Resource Manager* and *Scheduler*. The role of a Resource Manager is to act like a glue for a parallel computer to execute parallel jobs. It should make a parallel computer as easy to use as a Personal Computer (PC). A programming model such as *Message Passing Interface (MPI)* for programming on distributed memory systems would typically be used to manage communications within a parallel program by using the MPI library functions. A Resource Manager allocates resources within a HPC system, launches and otherwise manages Jobs. Some of the examples of widely used open source as well as commercial resource managers are *SLURM*, *TORQUE*, *OMEGA*, *IBM Platform LSF* etc. Together with a scheduler it is termed as a Batch System. The role of a job scheduler is to manage queue(s) of work when there is more work than resources. It supports complex scheduling algorithms which are optimized for network topology, energy efficiency, fair share scheduling, advanced reservations, preemption, gang scheduling (time-slicing jobs) etc. It also supports resource limits (by queue, user, group, etc.). Many batch systems provide both resource management and job scheduling within a single product (e.g. LSF) while others use distinct products(e.g. Torque Resource Manager and Moab Job Scheduler). Some other examples of Job Scheduling Systems are *LoadLeveler*, *OAR*, *Maui*, *SLURM* etc.

Existing Batch Systems usually support only static allocation of resources to an application before they start which means the resources once allocated are fixed for the lifetime of the application. The complexity of applications have been growing, However, especially when we consider advanced techniques in Scientific Computing like *Adaptive Mesh Refinement (AMR)* where applications exhibit complex behavior by changing their resource requirements during execution. The Batch Systems of today are not equipped to deal with such kind of complex applications in an intelligent manner apart from giving them the maximum number of resources before it starts that will result in a sheer wastage of resources leading to a poor resource utilization. In order to support such adaptive applications at HPC centers there is an urgent need to investigate and implement extensions to existing resource management systems or develop an entirely new system. These supporting infrastructures must be able to handle the new kind of applications and the legacy ones intelligently keeping in mind that they should now be able to achieve much higher system utilization, throughput, energy efficiency etc.

compared to their predecessors due to the elasticity of the applications.

1.1 Invasive Computing

Invasive Computing is a novel paradigm for the design and resource-aware programming of future parallel computing systems. It enables the programmer to write efficient resource aware programs. This approach can be used to allocate, execute on and free resources during execution of the program. The result is an adaptive application which can expand and shrink in the number of its resources at runtime. HPC infrastructures like clusters, supercomputers execute a vast variety of jobs, majority of which are parallel applications. These centers use intelligent resource management systems that should not only perform tasks of job management, resource management and scheduling but also satisfy important metrics like higher system utilization, job throughput and responsiveness. Traditionally, MPI applications are executed with a fixed number of MPI processes but with Invasive MPI applications they can evolve dynamically at runtime in the number of their MPI processes. This in turn supports advanced techniques like AMR where the working set size of applications change at runtime. Such kind of adaptive programming paradigms need to be complemented with intelligent resource management systems that can achieve much higher system utilization, energy efficiency, throughput etc. compared to their predecessors due to elasticity of the applications.

Under the collaborative research project funded by the **German Research Foundation (DFG)** in the **Transregional Collaborative Research Centre 89 (TRR89)**, research efforts are being made to investigate this Invasive Computing approach at different levels of abstraction right from the hardware up to the programming model and its applications. **Invasive MPI** is an effort towards invasive programming with MPI where the application programmer has MPI extensions available for specifying at certain safe points in the program, the possibility of changing the resource set of the application during runtime or adapt to the available resources during execution.

1.2 Dynamic Resource Management

Two of the most widely used resource managers on HPC systems are **SLURM** and **TORQUE**. The two major components in general of any sophisticated resource manager are the batch scheduler and the process manager. The Process Manager is responsible for launching the jobs on the allocated resources and managing them throughout their lifetime. Examples of process manager are *Hydra*, *SLURM Daemon (slurmd)* etc. The

process managers interact with the processes of a parallel application via the **Process Management Interface (PMI)**. In order to support Invasive Resource Management, The following components will be implemented: *iBSched* (Batch Scheduler for Invasic Jobs) built as an extension into an existing batch system and *iRTSched* (Invasive Distributed Run Time Scheduler) similar to a controller daemon which will sit between the batch scheduler and the process manager. **SLURM** is the choice of an existing batch system on which this prototype will be implemented for demonstrating Invasive Computing and the Figure 3.4 shows a high level illustration of the architecture for such an Invasive Resource Management.

In addition to a job queue for legacy static jobs, we now have an additional job

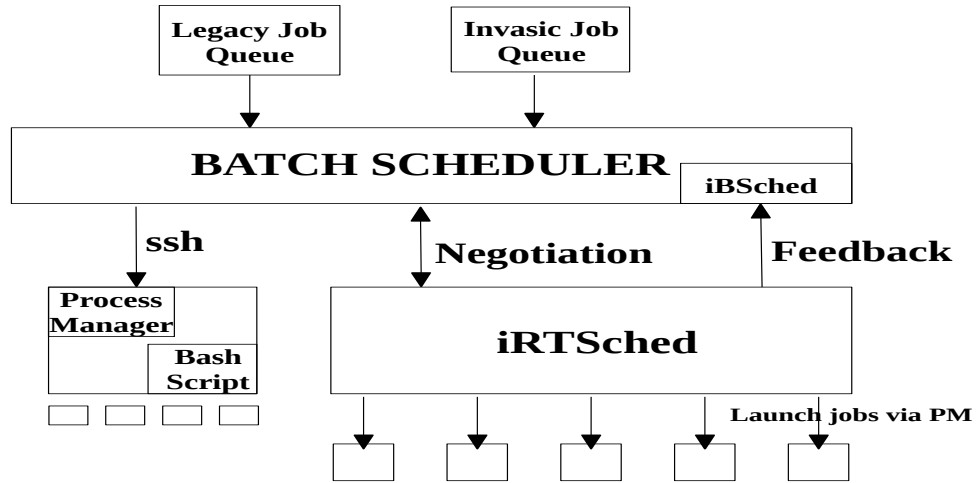


Figure 1.1: Invasive Resource Management Architecture

queue for invasic jobs. The objective of such a multi-level approach is to avoid modifying the existing system which will be a substantially large effort and rather have an independent component that caters specifically to invasic jobs. *iRTSched* is responsible for managing the resources present in the partition used specifically for running invasic jobs. With this approach, the existing legacy jobs can be served via the existing batch scheduler and the invasic jobs can be served by via *iBSched*. There could be also be a possibility that due to load balancing and all other partitions being busy, the legacy static jobs could be sent via *iBSched* as well. *iRTSched* talks to *iBSched* via a negotiation protocol to receive jobs dispatched from *iBSched* which it then launches by performing some run time scheduling like pinning of jobs, expand/shrink etc.

iBSched is responsible for scheduling invasic jobs and forwards the same via a negotiation protocol to iRTSched for execution. The decisions will be made on the basis of resource offers sent by iRTSched which creates them based on the state of the partition. Upon receiving a resource offer, iBSched will select jobs from the queue that can be mapped to this offer and decide whether to accept or reject it based on some criteria. A resource offer can represent real or virtual resources because the iRTSched can also present a virtual view of resources in the hope of getting a mapping of jobs to offer that is more suitable to satisfy its metrics such as resource utilization, throughput, energy efficiency etc. It can either accept or reject the mapping received from iBSched. Similarly, The iBSched makes its decisions to optimize for certain metrics such as reduced job waiting times, deadlines, priorities etc. This highlights the mismatching policies/metrics for which both iRTSched and iBSched make their decisions on and therefore both will be involved in a negotiation to reach a common agreement.

1.3 Document Structure

This is the end of the first section which gave an introduction to this Master Thesis and the kind of problem it deals with. The rest of this report is organized as follows:

- **Related Work:** This section will briefly mention some of the earlier research efforts that have been made in the direction of batch job scheduling, runtime scheduling specifically to support adaptive applications and resource-aware programming paradigms to implement adaptive applications.
- **Invasive Computing:** This section will first introduce the concept of invasive computing in brief and explain the motivation behind this concept. This is followed by an elaborate description of the traditional resource management approach in order to contrast it with the following section on invasive resource management that is necessary to support invasive computing.
- **Architecture:** This section will present an abstraction of the complete system at a high level showing all the components and how they will interact with each other like a skeleton. It deals with what is being done and where is it being done but not how. This "how" is tackled in the following section of design.
- **Design:** This section will present the details on how we are building the system whose architecture was illustrated in the previous section. It deals with the internal details of the individual modules / components, flow charts and illustrations. It describes what it can do and what it cannot.

- **Evaluation:** This section will cover the evaluation of the work presented in this thesis. It will describe the approach used for evaluating the system in order to test its functionality, correctness and performance.
- **Conclusion and Future Work:** This last section concludes the report on this thesis with a highlight of what was successfully achieved along with the possible scope of what can be done as a part of future research work.

2 Related Work

The problem of job scheduling for adaptive applications as introduced in the previous chapter is really an intersection of many different problems. In this section, we will look at some of the relevant work that has happened in the past with respect to batch job scheduling and also briefly look at a subset of the relevant work that has happened with resource-aware programming models and runtime scheduling.

Substantial amount of work has been done in exploring efficient resource management and scheduling for HPC systems in the past focusing on rigid and moldable jobs. There has been growing interest in the recent years to explore dynamic resource management and scheduling techniques due to the increasing complexity of applications, ex: Scientific applications utilizing Adaptive Mesh Refinement(AMR) techniques change their working set size as the mesh is refined/coarsened. Dynamically allocating resources to such applications at runtime or other applications which are potentially scalable has the potential to improve resource utilization, energy efficiency, fault tolerance and a host of other metrics[15][32]. Wolf.et.al[28] demonstrated such a dynamic resource management approach for supporting evolving jobs<ref to section no.> by extending the Torque / Maui Batch System to allow dynamic allocations and a dynamic fairness strategy implemented in the Maui Scheduler to efficiently service static and dynamic allocations. Their results showed reduced turnaround and waiting times for applications, while increasing system utilization and throughput. Wolf.et.al[29] also proposed an extension to this work to further support malleable<ref to section no.> type of jobs with the help of a communication protocol between the Batch System and the Charm++ runtime which enables the malleability of applications. In combination with their earlier work, The Scheduler now supported a mix of jobs with different types such as rigid, malleable and evolving.

In contrast to the research works mentioned above that focused on the Charm++ programming model, there have been several efforts directed towards the MPI programming model[5][25][47][27][26] some of which explored on how to first provide such an adaptive behavior on MPI applications and some on how to support the same with the help of libraries, resource management systems[18] etc. In the current collaborative research project earlier efforts were directed towards supporting such adaptive pro-

programming paradigms in shared memory programming model such as openMP[14] and work was also done in MPI[46]. The ongoing efforts in this project is to now continue forward in the direction of MPI but to approach the problem vertically at different levels of abstraction such as the programming model, middleware and runtime tuning of applications.

This thesis will mainly focus on the middleware part of this solution stack to support adaptive applications in future HPC systems. Middlewares provide job scheduling for HPC systems and job scheduling has always remained a very active area of research both for theoretical purposes and for practical systems. The most widely used and proven technique running in most of the HPC systems today is backfilling. Backfilling was first developed by lifka[23] for the Argonne National Laboratory's IBM SP system to address the need for a new scheduling system on supercomputers and was named as **EASY**(Extensible Argonne Scheduling System). Feitelson.et.al[12] proposed an extension to backfilling by providing a reservation for every job that could not be started and allowing lower priority jobs from behind in the queue to start if they would not delay these reservations. This was called as *Conservative Backfilling* whereas the one proposed by lifka was called *Aggressive Backfilling* as it provided a reservation to the job only at the front of the queue. Feitelson.et.al[11] came up with a new approach to backfilling algorithm where a set of jobs from the job queue are looked at once for making a scheduling decision using dynamic programming instead of traversing the queue one job at a time. The claim was that this resulted in better packing of jobs resulting in higher utilization, reduced mean response time and mean slowdown of all jobs. This scheduler was named as *LOS(Lookahead Optimizing Scheduler)*.

In addition to the standard batch job scheduling for rigid jobs, efforts have also been made in the direction of supporting moldable jobs. Cirne.et.al[6] proposed **SA(Supercomputer AppLes - Application level scheduler)** that would select on behalf of the user, the appropriate size for a given job from its available sizes. This decision was made based on the state of the resources, job characteristics etc to minimize turnaround time. Sadayappan.et.al[31][35][36] proposed several works for moldable scheduling of parallel jobs by using different policies like fair-share, overbooking, and finally considering job's efficiency also in determining the best partition size for the job. They developed an iterative algorithm in order to make the appropriate choice of values for overbooking and efficiency which is inturn based on the scalability characteristics of the job mix.

Traditionally, Batch Systems have had both batch and runtime scheduling merged together into a single component. In this work, we explore the idea of separating

the two concerns into different components for providing a dynamic and flexible scheduling functionality. Such a multi-level approach has usually been observed in the grid and cloud based infrastructures[34][22][38]. These infrastructures are distributed geographically and employ a local resource manager and scheduler at each of these locations and all of these locations are coordinating by a single centralized broker / service provider. The broker usually engages in some sort of a negotiation[49][20][51][30] with the different locations / resource providers / sites for availability of resources in order to make the best possible decision for serving its users.

3 Invasive Computing

Invasive Computing[17] is a novel paradigm for designing and programming future parallel systems. Decreasing feature sizes are motivating a redesign of multi million transistor system-on-chip architectures. This can lead to a dramatic increase in the rates of temporary and permanent faults as well as feature variations. SoCs with 1000 or more processors on a single chip in the year 2020 are foreseen, hence static and central management concepts to control the execution of all resources are no longer appropriate. Invasive Computing allows a program to be resource-aware by which it can explore and dynamically spread its computations to neighbour processors in a phase called as *invade*, then to execute portions of code of high parallelism degree in parallel based on the available invisable region in a given multi-processor architecture in a phase called as *infect*. Later, once the degree of parallelism should be lower again or if it terminates, the program may enter a *retreat* phase where it can deallocate resources and resume execution again, for example, on a single processor. With the help of such resource awareness, the program has the ability to self-organise itself and be immune to faults, feature variations, be highly scalable, show performance gain and record a higher resource utilization metric.

This concept would require not just new programming concepts, languages, compilers and operating systems but a radical change in the architectural design of MPSoCs (*Multi-Processor Systems-On-a-Chip*) so as to efficiently support invasion, infection and retreat operations. Some of the main motivations behind the idea of invasive computing are enumerated below:

- **Programmability** How to map algorithms and programs to 1000 processors or more and how to benefit from the massive parallelism available by tolerating manufacturing defects, feature variations etc.
- **Adaptivity** Modern applications have unpredictable resource requirements most of which may not be known at compile-time. In addition to this, when different applications are running on a single chip, resource distribution will have to happen dynamically keeping up a high resource utilization and performance. These factors show the need for some sort of hardware / software reconfigurability of the MPSoC.

- **Scalability** How to efficiently run algorithms and programs on different number of resources?
- **Physical Constraints** Heat dissipation will be a major bottleneck. Intelligent methods and architectural support to run algorithms at different speeds to exploit parallelism under power reduction is needed.
- **Reliability and Fault-Tolerance** Applications must be immune to temporal or permanent faults that may be caused due to manufacturing defects, feature variations, degradation etc. This especially has a higher likelihood to happen in the case of future MPSoCs.

The paradigm of invasive computing offers a new perspective for programming large scale HPC systems. Current resource management systems manage resources via static partitioning among parallel jobs. This is a very rigid approach considering that an application will then be limited to a fixed amount of parallelism it can utilize. This, however, will not be beneficial especially in the case of future exascale systems where if one needs to derive maximum performance then the maximum number of resources will have to be allocated. The application can benefit from invasive programming during certain phases of its runtime by running at maximum parallelism and in the remaining time it can run at a lower parallelism.

Another motivation is for a specific classes of applications like multi-grid and adaptive grid. Multi-grid applications work on multiple grid levels ranging from fine to coarse grids. On fine grids, more resources could yield better performance and efficiency, whereas on coarse grids fewer resources would be sufficient. In the case of adaptive grid applications, the grid is dynamically refined according to the current solution and the application may go through different levels of parallelism in different phases.

Scaling the systems to exaflop level would consume significantly more power that would very likely cross a gigawatt. Reducing the power requirement by a factor of atleast 100 is a challenge for future hardware and software technologies. Invasive computing concept with invasive programming models combined with intelligent resource management and flexible scheduling mechanisms can possibly help in addressing this challenge.

Coping with run-time errors would be another major challenge. Due to design and power constraints, the clock frequency is unlikely to change and feature sizes would continue to decrease as per moore's law for the next few years. By 2020, it is envisaged that exascale systems can possibly have approximately one billion processing elements.

An immediate consequence is that the frequency of errors will increase while timely identification and correction of errors would be much more difficult. Fault tolerance would be one of the most important challenges in this regard.

Exploiting massive parallelism for current and emerging scientific applications would also be another major challenge.

3.1 Traditional Resource Management

The role of a resource manager is to act like a *glue* for a parallel computer to execute parallel jobs. MPI would typically be used to manage communications within the parallel program. A resource manager allocates resources within a cluster, launches and otherwise manages the jobs. The combination of *Scheduler+Resource Manager* makes it possible to run parallel jobs and is together termed as a Batch System. These systems are classified[38] depending on whether they work towards serving job requests to satisfy the available resources or they also plan for the future.

3.1.1 Classification

The process of computing a schedule may be done by a queueing or planning based scheduler. A *Schedule* is computed for the job requests that are present in the job queue. Every request contains information such as the number of requested resources and a duration for how long the resources are requested for. There can also be some reservation requests present. These request for resources at a specified time for a given duration. Once the scheduler accepts such a request, it is a reservation and those exact resources are then blocked for that specified time and are unavailable for any scheduling purposes.

Queueing systems try to utilize the currently available resources in order to satisfy the job requests. Future resource planning for all the pending requests is not done. Hence, the pending requests have no proposed start times. Planning systems in contrast, plan for the present and the future. Planned start times are assigned to all requests and a complete schedule about the future resource usage is computed.

Queueing Systems These systems have several queues with different limits on the number of requested resources and the runtime limit for the job. Jobs within a queue are ordered according to a scheduling policy, e. g. FCFS (first come, first serve). Queues might be activated only for specific times (e. g. prime time, non prime time, or weekend). The task of a queueing system is to assign free resources to waiting requests.

The highest prioritized request is always the queue head. If it is possible to start more than one queue head, further criteria like queue priority or best fit (e. g. leaving less resources idle) are used to select a request. There might also exist a high priority queue whose jobs are preferred at any time. If not enough resources are available to start any of the queue heads, the system waits until enough resources become available. These idle resources may be utilized with less prioritized requests by backfilling mechanisms.

In queuing systems, no information about future job starts are available. Consequently guarantees can not be given and resources can not be reserved in advance. Resource reservation will have to be done manually by the administrative staff. Job requests also come with run time limits. If a job runs for more than the run time limit it specified or the partition limit, then the system will usually kill such a job.

Planning Systems Planning systems schedule for the present and future. They assign start times to all requests and a full schedule is generated. Runtime estimates for jobs are mandatory for this planning. With this knowledge advanced reservations are easily made possible. The re-planning process is the key element of a planning system. Each time a new request is submitted or a running request ends before it was estimated to end, a new schedule has to be computed and this function is invoked. At the beginning of a re-plan, all pending requests are sorted according to a scheduling policy in addition to clearing their previously planned start times. All the pending requests are then re-inserted at the earliest possible start time in the schedule. After this step each request is assigned a planned start and end time.

As planning systems work with a full schedule and assign start times to all requests, resource usage is guaranteed and advanced reservations are possible. If any reservation request is accepted then it is stored in an extra list for accepted reservations. During the re-planning process this list is processed before the list of standard job requests which can float around in the schedule. One drawback in a planning system is the cost of scheduling.

3.1.2 Job Scheduling

Typical resource management systems store job requests in list-like structures. A scheduling policy consists of two parts: inserting a new request in the data structure at its submission and taking requests out during the scheduling. Different sorting criteria are used for inserting new requests and some examples are (either in increasing or decreasing order)[45]:

- by arrival time: FCFS (first come first serve).

- by duration: Both increasing and decreasing orders are used. Sorting by increasing order leads to SJF (shortest job first). Accordingly LJF (longest job first) sort by decreasing run time. This requires job runtime estimates. SJF and LJF are both not fair, as very long (SJF) and short (LJF) jobs potentially wait forever.
- by area: The jobs area is the product of the width (requested resources) and length (estimated duration).
- by given job weights: Jobs may come with weights which are used for sorting. Job weights consist of user or system given weights or a combination of both. For example: all jobs receive default weights of one and only very important jobs receive higher weights, i. e. they are scheduled prior to other jobs.
- by many others: e. g. smith ratio, number of requested resources, current slowdown, ...

In the scheduling process, Jobs are taken out of the ordered data structure for either a direct start in queuing systems or for placing the job in a full schedule (planning system):

- front: The first job in the data structure is always processed. Most scheduling policies use this approach as only with this a sorting policy makes sense. FCFS, SJF, and LJF use this approach.
- first fit: The first job that fits into the available resources.
- best fit: All jobs are tested to see whether they can be scheduled. According to a quality criterion the best suited job is chosen. Commonly the job which leaves the least resources idle in order to increase the utilization is chosen. If more than one job is best suited an additional rule is required, e. g. always take the first, the longest/shortest job, or the job with the most weight.

If fairness in common sense has to be met, i. e. the starting order equals the arrival order, only the combination of sorting by increasing arrival time and always processing the front of the job structure can be used. All other combinations do not generate fair schedules. However, such a fair scheduler is not very efficient, as jobs usually have to wait until enough free resources are available. Therefore, basic scheduling policies are extended by backfilling, a method to avoid excessive idleness of resources.

Backfilling The default algorithms used by job schedulers for parallel supercomputers select jobs for execution in FCFS order, and run each job to completion, in batch mode. The problem with this simplistic approach is that it causes significant fragmentation,

as jobs with arbitrary sizes/arrivals do not pack perfectly. Specifically, if the first queued job requires many processors, it may have to wait a long time until enough are freed. During this time, processors stand idle as they accumulate, despite the fact there may very well be enough of them to accommodate the requirements of other, smaller, waiting jobs.

To solve the problem, most schedulers therefore employ the following algorithm. Whenever the system status changes (job arrivals or terminations), the scheduler scans the queue of waiting jobs in order of arrival (FCFS) and starts the traversed jobs if enough processors are available. Upon reaching the first queued job that cannot be started immediately, the scheduler makes a reservation on its behalf for the earliest future-time at which enough free processors would accumulate to allow it to run. This time is also called the *shadow time*. The scheduler then continues to scan the queue for smaller jobs (require fewer processors) that have been waiting less, but can be started immediately without interfering with the reservation. In other words, a job is started out of FCFS order only if it terminates before the shadow time and therefore does not delay the first queued job, or if it uses extra processes that would not be needed by the first queued job. The action of selecting smaller jobs for execution before their time provided they do not violate the reservation constraint is called backfilling.

This approach was initially developed for the IBM SP1 supercomputer installed

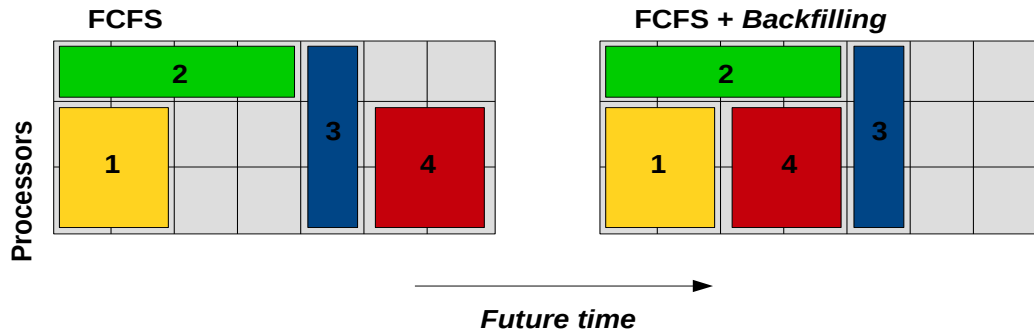


Figure 3.1: FCFS with and without Backfilling[45]

at the Argonne National Laboratory as part of EASY (Extensible Argonne Scheduling System), which was the first backfilling scheduler[lifka]. In terms of performance, backfilling has shown to be a close second to more sophisticated algorithms that involve preemption (time slicing), migration, and dynamic partitioning. The down side of backfilling is that it requires the scheduler to know in advance how long each job (*user runtime estimates*) will run. This is needed for two reasons:

- to compute the shadow time for the longest-waiting job (e.g. in the example given in 3.1, we need to know the runtimes of job 1 and job 2 to determine when their processors will be freed in favor of job 3), and
- to know if smaller jobs positioned beyond the head of the wait-queue are short enough to be backfilled (we need to make sure backfilling job 4 will not delay job 3, namely, that job 4 will terminate before the shadow time of job 3).

Jobs that exceed their estimates are killed, so as not to violate subsequent commitments (the reservation). The combination of simplicity, effectiveness, and FCFS semantics has made EASY a very attractive and a very popular job scheduling strategy. Nowadays, virtually all major commercial and open-source production schedulers support EASY backfilling. Figure 3.2 briefly mentions some of the various tunable knobs of backfilling algorithms.

3.1.3 SLURM

The prime focus of this work will be on **SLURM(Simple Linux Utility For Resource Management)** which will be the choice of batch system upon which the support for Invasive Computing will be demonstrated. SLURM is a sophisticated open source batch system written in C whose development started in the year 2002 at Lawrence Livermore National Laboratory as a simple resource manager for Linux Clusters. A few years ago it spawned into an independent firm under the name SchedMD. SLURM has since its inception also evolved into a very capable job scheduler through the use of optional plugins. It is used on many of the world's largest supercomputers and is used by a large fraction of the world's TOP500 Supercomputer list. It supports many UNIX flavors like AIX, Linux, Solaris and is also fault tolerant, highly scalable, and portable.

SLURM has a centralized manager called *slurmctld*(controller daemon) that is the main nerve center of SLURM. SLURM operates in a style similar to the Master-Slave paradigm where the Master is the *slurmctld*. It takes centralized decisions to monitor resources and work. In the event of a failure, there may also be a backup controller. Each of the nodes in the cluster has a daemon running on it called as *slurmd* and these are the slaves. These daemons are started on every node and they are responsible for monitoring them. This can resemble a remote shell: it waits for work from the controller, executes that work, returns status and waits for more work. The *slurmd* daemons provide fault-tolerant hierarchical communications and also are responsible for spawning an additional daemon called *slurmstepd*. The step daemon as it is called is responsible for the node local part of the job step that are the subset of processes running on the local node. A job step in SLURM refers to an application started with

Parameter	Description
<i>Number of reservations</i>	Default is 1 . This is called " Aggressive Backfilling " where only the first queued job receives a reservation. This can cause delays in execution of the other waiting jobs. The alternative is " Conservative Backfilling " where other waiting jobs are also allocated reservations. The number of reservations to allow is configurable by the administrator.
<i>Looseness of reservations</i>	This refers to a " Selective Reservation " strategy depending on the extent different jobs have been delayed by previous backfilling decisions. This is similar to " Flexibe Backfilling " strategy where the backfilling is allowed to violate the reservation upto a certain slack.
<i>Order of Queued Jobs</i>	Usually FCFS order is used. An alternative is where the jobs are prioritized in some way and the scheduler picks the jobs according to this order including for backfilling. The factors on which priority could be calculated are: job characteristics, user, group, user priority, fairness, weight, duration etc.
<i>Partitioning of Reservations</i>	Machine is divided into several disjoint partitions with the freedom to move around idle processors dynamically based on current needs. Each partition is associated with its own job class, runtime limit for jobs, wait-queue and reservation. Backfilling candidate is chosen in a round-robin fashion, each time from a different partition respecting the reservations.
<i>Adaptiveness of Backfilling</i>	Simulates the execution of recently submitted jobs under various scheduling disciplines and switches the algorithm to the ones which gives the highest performance.
<i>Lookahead into the queue</i>	Default behavior is to consider the queued jobs one at a time that may lead to loss of resources. Alternative is to look at a window jobs at a time and pick a mapping that gives the maximum utilization but respects all the reservations.
<i>Speculative Backfilling</i>	The Scheduler is allowed to backfill jobs even if it interferes with an existing reservation as it speculates that the job will finish earlier than estimated time.
<i>Minimize the electric power demand</i>	The current direction of research in job scheduling is targeted towards power efficient exascale HPC systems. Scheduling decisions will have to be made considering power constraints of the machine, power stability, energy efficiency.

Figure 3.2: Backfilling Variations[45]

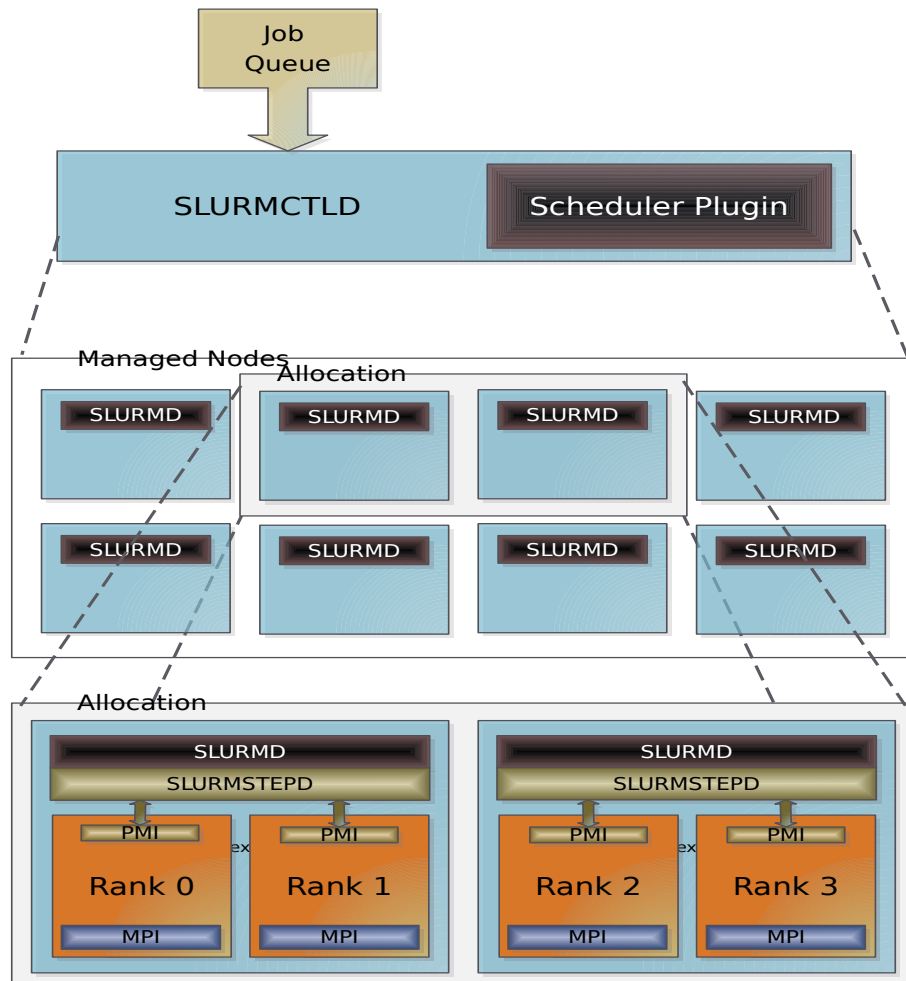


Figure 3.3: SLURM Architecture

the help of *srun* and its allocated resources. *srun* could be used independently to launch jobs or one can specify the same within a batch script while using *sbatch*. *srun* is one of the tools SLURM provides that allows the user to launch interactive jobs on the cluster, *sbatch* to launch batch jobs and several others relating to accounting, job status, cancellation operation etc.

The Figure 3.3 shows the high level architecture of SLURM with the interaction between the several of its key components. It also shows the interaction between an MPI application through the PMI (Process Management Interface), *slurmd* daemon of a node and the *slurmstepd*. **Plugins** are dynamically linked objects loaded at runtime based upon configuration file and/or user options. Approximately 80 plugins of different varieties are currently available. Some of them are listed below:

- **Accounting storage:** MySQL, PostgreSQL, textfile.
- **Network Topology:** 3D-Torus, tree.
- **MPI:** OpenMPI, MPICH1, MVAPICH, MPICH2, etc.

PLugins are typically loaded when the daemon or command starts and persist indefinitely. They provide a level of indirection to a configurable underlying function.

3.2 Resource Aware Programming

3.2.1 Job Classification

The throughput of HPC Systems depends not only on efficient job scheduling but also on the type of jobs forming the workload. As defined by Feitelson, and Rudolph[9], Jobs can be classified into four categories based on their flexibility:

- **Rigid Job:** Requires a fixed number of resources throughout its execution.
- **Moldable Job:** The resource requirement of the job can be molded or modified by the batch system before starting the job (e.g. to effectively fit alongside other rigid jobs). Once started its resource set cannot be changed anymore.
- **Evolving Job:** These kind of jobs request for resource expansion or shrinkage during their execution. Applications that use Multi-Scale Analysis or Adaptive Mesh Refinement (AMR) exhibit this kind of behavior typically due to unexpected increases in computations or having reached hardware limits (e.g. memory) on a node.
- **Malleable Job:** The expansion and shrinkage of resources are initiated by the batch system in contrast to the evolving jobs. The application adapts itself to the changing resource set.

The first two types fall into the category of what is called as the static allocation since the allocation of rigid and moldable jobs must be finalized before the job starts.

Whereas, the last two types fall under the category of dynamic allocation since this property of expanding or shrinking evolving and malleable jobs (together termed adaptive jobs) happens at runtime. Adaptive Jobs hold a strong potential to obtain high system performance. Batch systems can substantially improve the system utilization, throughput and response times with efficient shrink or expand strategies for running jobs that are adaptive. Similarly, applications also profit when expanded with additional resources as this can increase application speedup and improve load balance across the job's resource set.

3.2.2 Invasive Programming Models

In this section, we will briefly look into the details of the earlier invasive extensions done to OpenMP and MPI done as a part of this ongoing research project. This will give us an insight into the earlier approach taken towards realizing such resource-aware programming models. These invasive extensions provide us with a new parallel programming model that allows us to implement resource aware programs. Depending on the semantics of these new extensions, the resulting application can either be evolving (application dictates the changes to its resource set) or a malleable job (resource manager dictates the change in resources to which the application must adapt). Resource awareness could mean that either the program can allocate or free resources according to the amount of available parallelism / the dynamic size of the data or it could mean that it can adapt to the available resources for execution.

Parallel applications that are resource aware will invade or retreat from resources depending on their availability and on the load imbalances encountered during their runtime. To support this, some form of dynamic process management of the parallel application is necessary. And, in order to realize this in practice, the most basic requirement would be the need for a library that will serve as an application programming interface for programmers to implement such invasive applications that are capable of adapting to a changing set of resources. This requirement needs to be complemented by the extension of the resource management systems which would need to allocate or deallocate resources and coordinate with the library to allow for such adaptive operations of an invasive parallel application.

iOMP[14]

The OpenMP parallel programming model for shared memory systems was extended to support the programming of resource aware applications and is named as Invasive OpenMP or iOMP. Parallelization using OpenMP is done by inserting compiler directives into the application's source code to define parallel regions that are executed in

parallel by a team of threads whereas the sequential region would be executed by a single master thread. There are different ways to control the number of threads in a parallel region and the most common approach is through the environment variable, or through OpenMP library call or as an additional clause in its directives. iOMP has been implemented as a library in C++ using an object oriented approach and provides two important methods / operations available in its class *Claim*:

- **Invade**: This operation allocates additional resources / PE's (processing elements). A constraint parameter passed as an argument to this operation specifies the details such as which resources and how many of them [range] are additionally required from the resource manager.
- **Retreat**: This operations deallocates resources / PE's. A constraint parameter passed as an argument to this operation specifies the details such as which resources and how many of them must be freed to the resource manager.

A *Claim* (not the C++ Class) in iOMP refers to all the resources / PE's allocated to the application. This means that an iOMP program will always have a single claim. Initially, the claim size is 1 but it will increase and decrease during the runtime of the application. The constraint parameter mentioned before also allows the programmer to specify several other constraints such as memory, pinning strategy, architecture specific optimizations etc. Below is a small snippet of code from that shows an example of iOMP program.

```
int main()
{
    Claim claim;
    int sum = 0;
    /* Acquire resources according to the given constraints */
    claim.invade(PEQuantity(1, 3));

    /* Executing a parallel for loop on the given resources */
    #pragma omp for reduce reduction(+:sum)
    for (int i=0; i < 100000; i++)
        sum += i;

    /* Free resources and delete pinning */
    claim.retreat();
}
```

As another important part of the iOMP implementation, a resource manager has also been implemented. This has a global view of the resources in the shared memory system and acts like a server to every other running application that are its clients. Every client-server communication happens over a message queue. The resource manager handles the redistribution of the resources over time to all running applications based on their invade or retreat operations.

iMPI[46]

Similar to iOMP, previous research effort in this project was also directed towards extending parallel programming models for distributed memory systems. iMPI which stands for Invasive MPI is an extension to the MPI library that can support resource aware programming. The Single-chip Cloud Computer (SCC) from Intel Labs was an experimental CPU that integrates 48 cores and is basically a distributed memory system on the chip. This hardware platform along with its interesting memory features was used in order to evaluate this invasive programming model.

Message Passing has for long remained the dominant programming model for distributed memory systems. MPI stands for Message Passing Interface. It is a standardized and portable message passing system designed to function on a wide variety of parallel computers. It implements a message passing type of parallel programming model where the application consists of a set of processes with separate address spaces. The processes exchange messages by explicit send / receive operations. The following are the invasive extensions to MPI:

- *MPI_Comm_invalidate*: The main purpose of this operation is to reserve resources and for this it looks into what resources are currently available and invades them.
- *MPI_Comm_infect*: This operation is used by the application to specify the total number of cores to infect and which ones are preferred. This number can be less than or equal to the total number of cores that were reserved by the invade operation.
- *MPI_Comm_retreat*: This operation does the reverse of the invade+infect sequence. Instead of reserving and claiming resources, it returns them so that other invasive applications can claim them.

The above extensions were based on the MPICH2 library. A new process manager called **Invasive Process Manager (IPM)** was also developed as a part of the iMPI implementation. It was responsible for launching of the MPI jobs, as well as spawn operations (invade + infect) with low latency and functionality for resource awareness.

3.3 Invasive Resource Management

In this section, we will look at the latest extensions done to the MPI library for programming on HPC systems like clusters, supercomputers etc and also the extensions required for resource managers managing these HPC systems. This is in contrast to the earlier version of the iMPI which was targeted towards the Intel SCC platform. The MPI and Resource Manager extensions are not accomplished as a part of this thesis but are essential to be described here in order to get the right context for adaptive applications and their scheduling.

3.3.1 Invasive MPI

Traditionally MPI applications are static in nature which means that they are executed with a fixed number of processes. Dynamic process support is available through MPI spawn and its related operations. MPI spawn operation creates child processes in a separate child process group and thereafter both the parent and child process groups are connected by an intercommunicator. Many drawbacks of this dynamic process functionality have motivated the need for such an extended programming model. Some of these drawbacks are: spawn is a collective operation across the parent and child process group, intercommunicators complicate development effort as more spawn operations create more disjoint process groups, destruction of entire process groups is the only option limiting the granularity and location of releasing resources, spawn usually operates within the same unmodified resource allocation.

To overcome these drawbacks, Invasive MPI is being developed as an extended version of the MPI library that provides new API calls in order to allow the programmer to create an invasive MPI application. These extensions are necessary to make the application resource aware and to adapt according to a change in the resource set by performing data / load redistributions. Following are the proposed extensions being implemented in MPICH:

MPI Initialization in Adaptive Mode This allows the application to be initialized in adaptive mode. It is an extension of the standard MPI_Init operation and is now called as MPI_Init_adapt. The difference now is that a new parameter called local status is being passed. Upon the return of this Init function, local status will hold a value of *new* if the process doing this MPI initialization was created using the mpiexec command

or it will hold a value of *joining* if the process was created by the resource manager as a part of the expansion of an already running invasive MPI application. The *joining* processes will then begin an adaptation window after completing the initialization.

```
int
MPI_Init_adapt( int *argc,
               char ***argv,
               int *local_status,
               );
```

Probing Adaptation Data The resource manager decides when and how the adaptation of a running application will be initiated. This operation will allow the application to probe the resource manager for adaptation instructions. This operation is called `MPI_Probe_adapt` and instructs the application on whether there is an adaptation pending.

```
int
MPI_Probe_adapt( int *current_operation,
               int *local_status,
               int *nfailed,
               int *failed_ranks,
               MPI_Info *info
               );
```

A value of false returned from the current operation parameter will simply tell the application to continue doing progress normally. A true or fault value indicates that there is an adaptation to be done. In the case of a fault, the application will receive information of the failed MPI ranks, since failed processes may no longer be reachable. This operation will also provide the application with additional information on whether it is a joining process if the process was created by the resource manager to represent newly allocated resources as a part of an expansion operation. Joining processes can skip calling the probe operation. If the information returned is staying then it means it should remain in the process group after the adaptation, otherwise it is retreating.

Beginning an Adaptation Window This operation marks the start of an adaptation window. It provides two communicators as output: one intercommunicator that is equivalent to what is provided by standard spawn operations, and one intracommunicator that gives an early view of how the `MPI_COMM_WORLD` communicator will look like after the adaptation is committed. It is up to the application to make calls to this operation in a safe location. Each process is required to read its future rank and the future size of the process group from the helper `new_comm_world` commu-

nicator to perform an adaptation consistently. This new size and local rank of the process will persist after the `MPI_COMM_ADAPT_COMMIT` operation. Processes that are retreating during the adaptation window will not have access to the future `MPI_COMM_WORLD`, since a retreating process will be removed from the process group, their `new_comm_world` will be set to `MPI_COMM_NULL`. These processes will need to be reached over the provided intercomm from the children, or their current `MPI_COMM_WORLD` from the parents, during the adaptation window.

```
int
MPI_Comm_adapt_begin(
    MPI_Comm *intercomm,
    MPI_Comm *new_comm_world,
);
```

Committing an Adaptation Window This operation commits the adaptation. This operation affects `MPI_COMM_WORLD`: any process that has retired is eliminated from it, and any new joining process is inserted into it. After the commit, the `MPI_COMM_WORLD` communicator will match exactly the `new_comm_world` intracommunicator provided by the previously mentioned `MPI_COMM_ADAPT_BEGIN` operation. This operation also notifies the resource manager that the current adaptation is complete.

```
int
MPI_Comm_adapt_commit();
```

3.3.2 Resource Management Extensions

Existing batch systems usually support only static allocation of resources to applications before they start. We need to integrate invasive resource management into these existing batch systems in order to change the allocated resources dynamically at runtime. This will allow for an elastic execution of MPI applications. Such efforts have already been initiated in the Flux project[10]. Existing systems like SLURM allow a job to have extra resources by expanding its allocation. But, this does not fully satisfy the use case here as we need to either grow or shrink the application. Another important factor is the support needed from a programming model that would allow applications to be adaptive to such allocations. The extensions needed on the MPI library side have already been mentioned in the previous section. In order to achieve the extensions on the resource manager side the following SLURM components have been extended:

SLURM

The resource manager needs to closely coordinate with the invasive MPI library to

support invasive applications. It needs to fork processes on new resources when an application is expanding or destroy them in case it is shrinking. Both of which needs to be done in coordination with MPI. New processes could be created in the existing resource allocation of the running application, possibly allowing for oversubscription of CPU cores, but that would be of little benefit to most HPC application's performance and scalability. The following are the extensions specific to SLURM. Figure 3.4 shows an example of job adaptation using the below extensions.

Slurmctld Extensions A new operation that initiates an adaptation through the *srun* command: *srun_realloc_message* is introduced. This message is sent by slurmctld to the *srun* responsible for this job. The *srun_realloc_message* provides *srun* the following information: the list of new nodes allocated to the application and the number of processes to create on them, the list of nodes from where processes need to be destroyed and how many processes to destroy in them, the full list of nodes that compose the new allocation, and other necessary details. Currently adaptations are based on full nodes but future developments will support fine grain scheduling. When a transformation is triggered on a job, its status changes from *RUNNING* to *ADAPTING*. Each application notifies the resource manager when its adaptation is completed and its job record will get updated from the status *ADAPTING* to *RUNNING*; this state change marks the application as eligible for adaptations again and its released resources available for other jobs.

Slurmd Extensions In order to support the *MPI_Probe_adapt* operation, the PMI plugin for PMI2 interface towards MPI applications has been extended. This plugin is loaded by slurmd daemon after it starts up. The extensions are: Notifying the local daemon that joining processes are waiting in the *MPI_Comm_adapt_begin* operation, Notifying the local daemon that both the joining and preexisting processes have completed their adaptation and exited *MPI_Comm_adapt_commit*. The first extension to the PMI is used by the leader process of the joining group. It notifies its local Slurmd daemon, which then notifies the *srun* instance of the job step. The *Srun* instance then proceeds to notify each of the Slurmd daemons running in the preexisting nodes. These daemons then proceed to update their local *MPI_Probe_adapt* metadata. The second extension to the PMI is used by the leader process of the new adapted process group, that was created as a result of a successful completion of *MPI_Comm_adapt_commit*. The local daemon sends a notification message to *Srun*, which then forwards it to Slurmctld. The controller handles this message by updating the status of the job from *JOB_ADAPTING* to *JOB_RUNNING*.

Srun Extensions Most of the operations that are initiated by either the controller or any

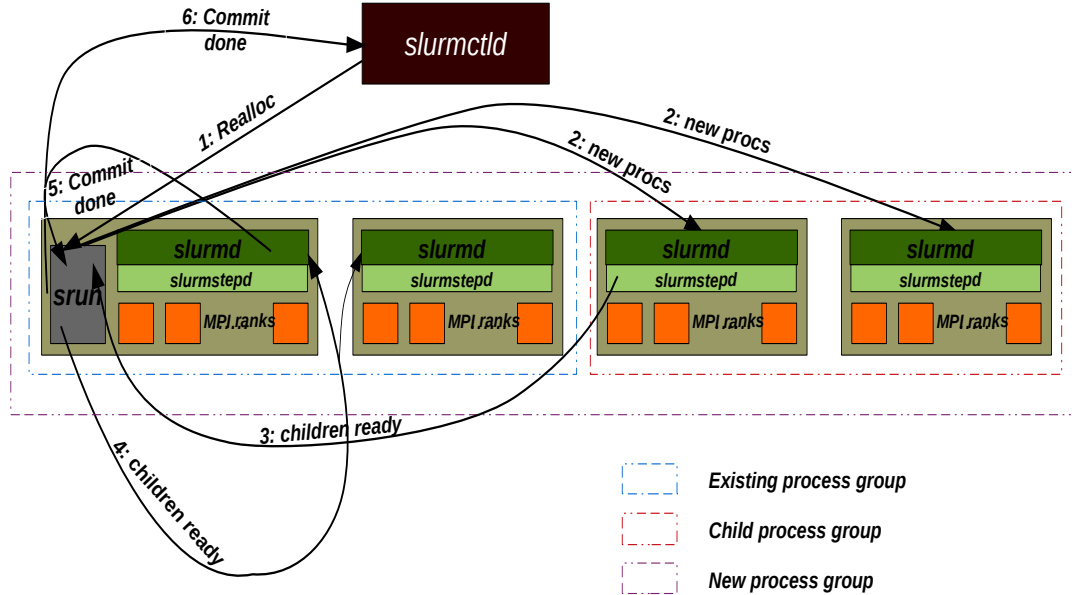


Figure 3.4: Job Adaptation Flow

application process (via the PMI and the SLURMD daemons) is handled partially by `srun`. It will handle the reallocation message received from the controller, notification that joining processes are ready and waiting in the `MPI_Comm_adapt_begin` operation, notification that the adaptation was completed through a successful `MPI_Comm_adapt_commit`. In addition to these handlers, `srun` has also been extended to manage the IO redirection of joining processes. In the original implementation, these were setup only at launch time; it can now manage redirections dynamically as processes are created and destroyed.

4 Architecture

This section illustrates and describes the high level architecture of the software implemented along with a few protocol sequence diagrams. It will help to understand at a high level about the components which form a part of this system to support invasive resource management and how will they interact with each other in order to integrate such an invasive resource management into existing batch systems. The following page shows the software architecture for how invasive resource management can be supported with SLURM and how exactly the new components fit into the existing software hierarchy.

- The top layer is that of the core resource management component which has access to job queues. In this architecture, it will now have access to not only the queue for the legacy (static) jobs but also invasive job queue.
- In a traditional setup the top layer will perform the task of job scheduling as well. This means that it will select a job(s) from the queue of jobs based on the current state of resources and many other factors to dispatch it to the traditional process manager below in the hierarchy. The process manager then takes the responsibility of launching these jobs on the allocated resources in the partition and managing them for their full lifetime. In case of parallel jobs, it will manage the job in a parallel environment along with facilitating the communication amongst the parallel tasks/processes with the help of a PMI (Process Manager Interface). The process manager may also spawn slave daemons on each of the nodes which are a part of the resource allocation for a single job to manage them more effectively.
- As discussed in the previous chapter, an independent invasive resource management component by the name "iRTSched" will be implemented which needs to communicate with the batch scheduler and influence the scheduling decisions taken by it. The iRTSched sits between the top layer and the process manager.
- A new job scheduler specifically for invasive jobs needs to be integrated into the existing batch system. In case of SLURM which has a modular design with several optional plugins, a new plugin by name "iBSched" will be implemented for SLURM to handle job scheduling specifically for invasive jobs.

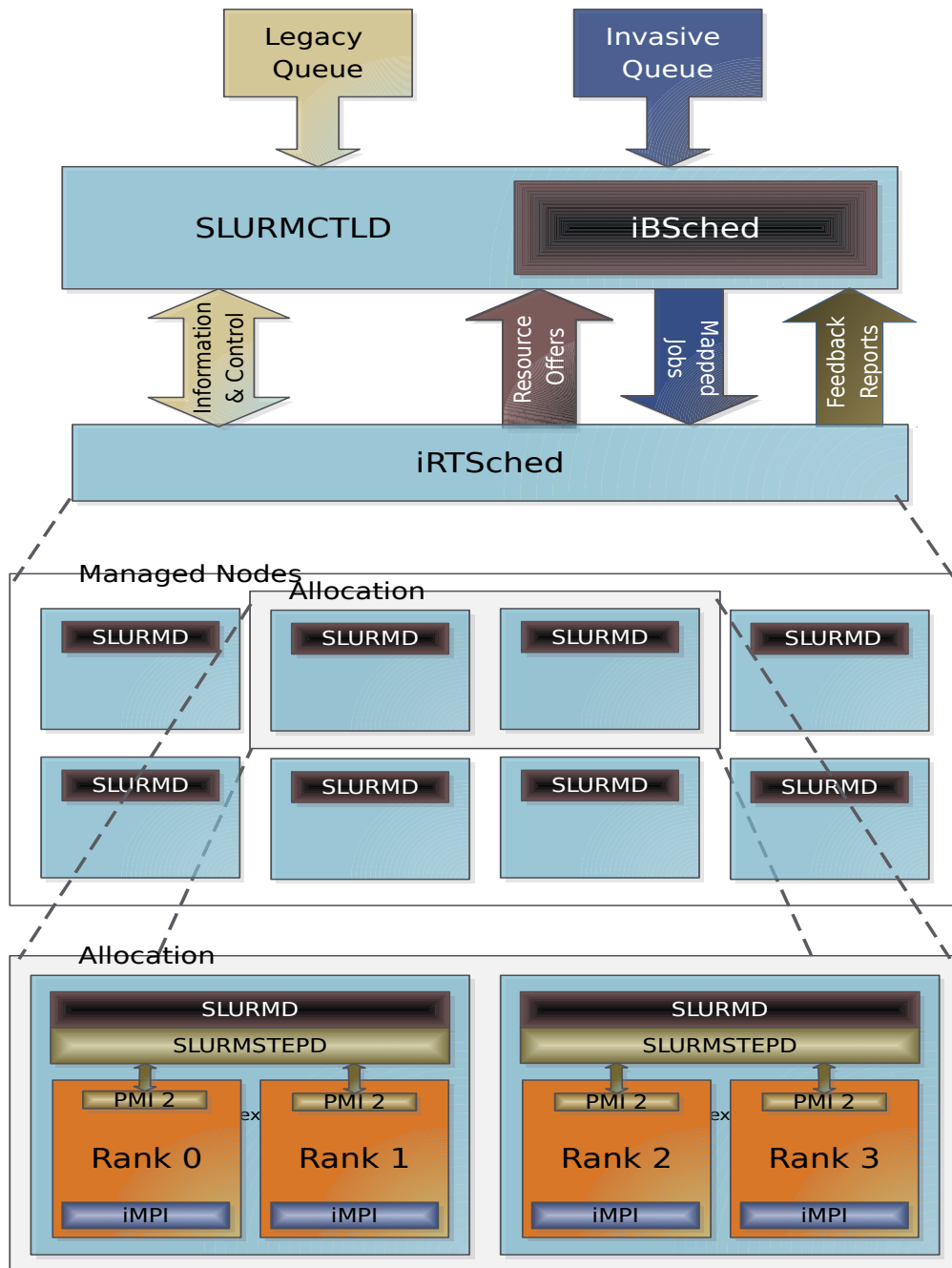


Figure 4.1: Invasive Resource Management Architecture

- Communication between iRTSched and iBSched will involve the negotiation protocol as explained in the previous chapter but will also include periodic / event-driven feedbacks being sent by iRTSched to iBSched. These will contain some useful information about the current state of the running jobs, their energy consumption, other job characteristics etc. This communication will also additionally support a means to service urgent jobs immediately.

Communication Phases

- **Protocol Initialization:** This phase basically establishes the initial environment between the communicating parties (iBSched and iRTSched) for proper communication later on. This happens only once at the start when both batch and runtime schedulers connect with each other. Successful initialization of this phase prepares both the parties to start negotiating based on the negotiation protocol described in the following points. During this protocol initialization various parameters such as protocol version, maximum attempts for negotiation, timer intervals and several others could be exchanged to set up the internal data structures and configuration tables for both the communicating parties. This protocol is a bi-directional communication.
- **Protocol Finalization:** This phase signals the end of the communication between iRTSched and iBSched using negotiation protocol. It leads to a safe termination of this communication followed by the release of any internal data structures allocated earlier along with configuration parameters. This results in consistent behaviour of both the communicating parties which can then proceed to safely terminate and exit. This protocol is a bi-directional communication.
- **Negotiation:** This is the most important phase in this whole approach to support invasive resource management. It is the phase during which both iRTSched and iBSched are negotiating with each other till they reach an agreement. If they do not then they continue till a certain limit is reached after which both of them just agree in their final attempt closing the current negotiation. After this a new transaction of negotiation will begin at some later point of time.
- **Feedback:** This concerns the periodic / event-driven feedback sent by the iRTSched to the iBSched containing useful information as mentioned earlier. iRTSched will also send a performance model of every completed job in the feedback that can be stored in some database as a part of the history of executions for this job. This will help the iBSched in the future if the same job is submitted when there will be additional performance specific information available about this job that can

<EVENT> <=> <PACKET>	DIRECTION OF COMMUNICATION
<REQUEST_RESOURCE_OFFER>	iBSched → iRTSched
<RESOURCE_OFFER>	iRTSched → iBSched
<RESPONSE_RESOURCE_OFFER>	iBSched → iRTSched
<NEGOTIATION_START>	iBSched → iRTSched
<RESPONSE_NEGOTIATION_START>	iRTSched → iBSched
<NEGOTIATION_END>	iBSched → iRTSched iRTSched → iBSched
<RESPONSE_NEGOTIATION_END>	iBSched → iRTSched iRTSched → iBSched
<STATUS_REPORT>	iRTSched → iBSched
<URGENT_JOB>	iBSched → iRTSched
<RESPONSE_URGENT_JOB>	iRTSched → iBSched

Figure 4.2: Message Types

be used by the batch scheduling algorithm to make better decisions for this job. This protocol is a uni-directional communication.

- **Urgent Jobs:** This protocol concerns the support for urgent jobs. At any given point of time a cluster or supercomputing center may want to support very high priority jobs immediately without any further delay. By introducing support for invasive computing, it makes it all the more feasible to help run these urgent jobs immediately by either shrinking the resources of running jobs or suspending / killing them.

4.1 Dynamic Resource Management

Separation of Concerns: In this thesis, we explore the idea of separating the concerns of batch and runtime scheduling into two different software layers or components in contrast to the existing systems where both are merged together. A negotiation protocol will be implemented as a means of communication between the two. The

motivation behind the negotiation protocol is the conflicting set of objectives between batch scheduler (user perspective: faster response time, fairness for jobs etc.) and runtime scheduler (system perspective: maximize utilization, throughput, energy efficiency etc.). The role of the batch scheduler is to forward jobs to a runtime scheduler which is managing the resources in the partition(s) and supports runtime scheduling of jobs. Runtime scheduler will also make intelligent expand or shrink decisions by observing scalability behavior of the running applications which are adaptive and use it to predict future resource requirements. The proposed protocol will also allow us to integrate such adaptive resource management systems into existing batch systems and thereby allowing for an easy migration from legacy systems to invasive resource management systems. Negotiation also helps us to realize a dynamic and flexible scheduling strategy by balancing the conflicting objectives at the two layers.

We will briefly look at the important components and their respective roles in this architecture in order to support adaptive applications on HPC systems by dynamically managing the resource allocation of running jobs.

4.1.1 Invasive Batch Scheduler

This component will be an extension to the batch scheduler in the existing batch systems. The scheduling decisions are communicated via the negotiation protocol to iRTSched in response to its resource offers. Batch scheduler is making its decisions to optimize for certain metrics such as reduced job waiting times, fairness, deadlines, priorities etc. Due to the separation of batch and runtime scheduling, it now looks analogous to a long term scheduler (batch scheduler) that admits jobs into the HPC system and a short term scheduler (runtime scheduler) which will manage the running jobs. Batch scheduler runs less frequently and there may be quite a lot of time gap until the next job may be admitted into the system. This depends upon available resources, any events like job termination, completion where we can expect to receive a resource offer from the runtime scheduler.

This batch scheduler is termed as invasive batch scheduler because it supports the negotiation interface to talk to an invasive runtime scheduler. It will now consider a mix of job types which are: *malleable*, *moldable*, *rigid*, and *evolving* and perform the following tasks as per the existing functionalities provided by SLURM:

- It is based on an event-driven batch scheduling where scheduling decisions would be made only when an offer is received from iRTSched. Scheduling is performed here at the granularity level of nodes.

- Dispatch jobs according to a scheduling algorithm by considering the resource constraints of the job that can include the number[range] of nodes (exclusive / shared), memory size, duration, quality of service parameters etc.
- With every negotiation attempt, the batch scheduler will try to relax the node count constraints of those jobs which could not be mapped to the available resource offer. As the negotiation attempts increase the degree to which the constraints are relaxed also increases. By doing this the batch scheduler is trying to bargain as best as it can in order to get an offer to serve as many jobs as possible.
- It will process the feedbacks received from iRTSched and update the status of running jobs in its queue. These feedbacks can be useful in understanding the scalability characteristics of running jobs that can possibly help in fixing the node count constraints for batch jobs that are still waiting.
- It can also process jobs which have an urgent priority. These jobs will receive special service and are going to be dispatched immediately to iRTSched. It will also consider jobs that may have some quality of service requirements such as deadlines and start time.
- Based on the offers received from the iRTSched, the batch scheduler may try to bias its scheduling decisions towards certain jobs (io/cpu bound).
- If a job that is in the queue has already been submitted in the past and there is some history stored in the system about it, The batch scheduler will then use this historical information from the database to update the details of the currently submitted job. This can lead to better scheduling decisions specific to this job.

4.1.2 Invasive Runtime Scheduler

This is an independent component which can talk to the existing batch systems via the negotiation protocol to receive forwarded batch jobs for execution. The runtime scheduler manages the running jobs via space sharing, time sharing or both. The runtime scheduler here is invasive because of the ability to support elastic execution for jobs that are adaptive. It can now redistribute resources to running jobs by either expanding or shrinking them based on their performance and scalability. Following points describe the role of iRTSched:

- It manages all the resources in the partition and is responsible for runtime scheduling at the granularity of cores / sockets, resource management and process management (infrastructure to launch, adapt and monitor parallel jobs).

- It will use an intelligent expand / shrink algorithm to dynamically adjust the resources of running jobs which are adaptive in nature.
- The runtime scheduler can bias its decisions towards any or all of the following criterias: maximizing throughput, energy efficiency, optimize for topology, resource utilization etc.
- Similar to iBSched, iRTSched will increase the degree of transformation of the running jobs to better serve the batch scheduler. As the negotiation attempts increase, the scheduler will try its best to fit every batch job forwarded by performing an aggressive transformation of the running jobs. Aggressive here means that the scheduler will try to either expand or shrink jobs to their maximum or minimum respectively of its constraints.
- In addition to the expand / shrink strategy, The runtime scheduler can support scheduling algorithms like backfilling, gang scheduling, preemption etc.
- iRTSched will build a performance model of the running application or will refine an existing model if it was forwarded by the batch scheduler in the job details. It will also make its scheduling decisions using this performance model by estimating the scalability, performance and efficiency of the application at different job sizes.
- It will send the generated models and other relevant information about a completed job like runtime, energy consumption, IO consumption in a feedback report to the iBSched.
- The runtime scheduler can send resource offers to batch scheduler for the following events: job termination (due to a signal), job completion (normal job termination), request for a resource offer, periodic runtime transformations of running jobs resulting in free resources.

4.1.3 iMPI Process Manager

The iMPI process manager in this architecture refers to multiple components as they all are necessary to successfully provide an elastic execution framework for adaptive applications. These components are *srun*, *slurmd* and *slurmstepd*. The previous chapter already described as to how the extended MPI library called as iMPI allows the programmer to create invasive MPI applications and how the extensions on the resource manager side in coordination with iMPI provide a complete infrastructure necessary for the elastic execution of adaptive applications.

4.2 Negotiation Protocol

This protocol forms the core of the interaction between the iBSched and iRTSched. It allows for iRTSched to make one or a set of resource offers to iBSched which then needs to select jobs from its job queue to be mapped to these resource offers and send back the mapping to the iRTSched. The iRTSched will then decide whether to accept or reject this mapping based on whether it satisfies its local metrics. If it accepts, it will launch them based on some runtime scheduling and if it rejects then it informs this to iBSched in addition to sending it a new resource offer that will also contain possible future start times for pending jobs. The iBSched can also reject the resource offer in which case it will forward the previous job mapping (if any) again with relaxed resource constraints for jobs that could not be mapped. On accepting an offer, the iBSched will again send back a mapping to iRTSched. This exchange of messages continue until both reach an agreement. If the number of such exchanges reach a threshold then iBSched will just accept whatever offer it receives and iRTSched will also accept the final mapping it received and try its best to satisfy all the jobs forwarded. This will close the transaction.

4.2.1 Protocol Sequence Diagrams

- Figure 4.3 illustrates a scenario where both iBSched and iRTSched are negotiating with each other. The scenario is continued in Figure 4.4. Figure 4.5 illustrates another scenario where negotiations may stop when job queue becomes empty and iRTSched will then wait for a request from iBSched for a resource offer that will happen when new jobs arrive.
- iBSched makes scheduling decisions at a coarser level of granularity which is nodes whereas iRTSched does at the granularity of cores and sockets. Both will negotiate with each other till they reach an agreement.
- The granularity difference is due to the fact that it is the iRTSched that is managing the resources in the partition where the jobs are running. It can try to do time sharing (more than one job sharing a node in time), fine granular space sharing (more than one job sharing the same node but mapped to different cores) etc. But what is known to iBSched is just a simplistic view of the number of nodes available to be used by its batch jobs.
- It is an event based scheduling which means iBSched makes a scheduling decision only when it is triggered by receiving a resource offer from iRTSched. It is only when the invasive job queue becomes empty that iBSched will have to explicitly send a request message to iRTSched for a resource offer otherwise at all other times scheduling is event based.

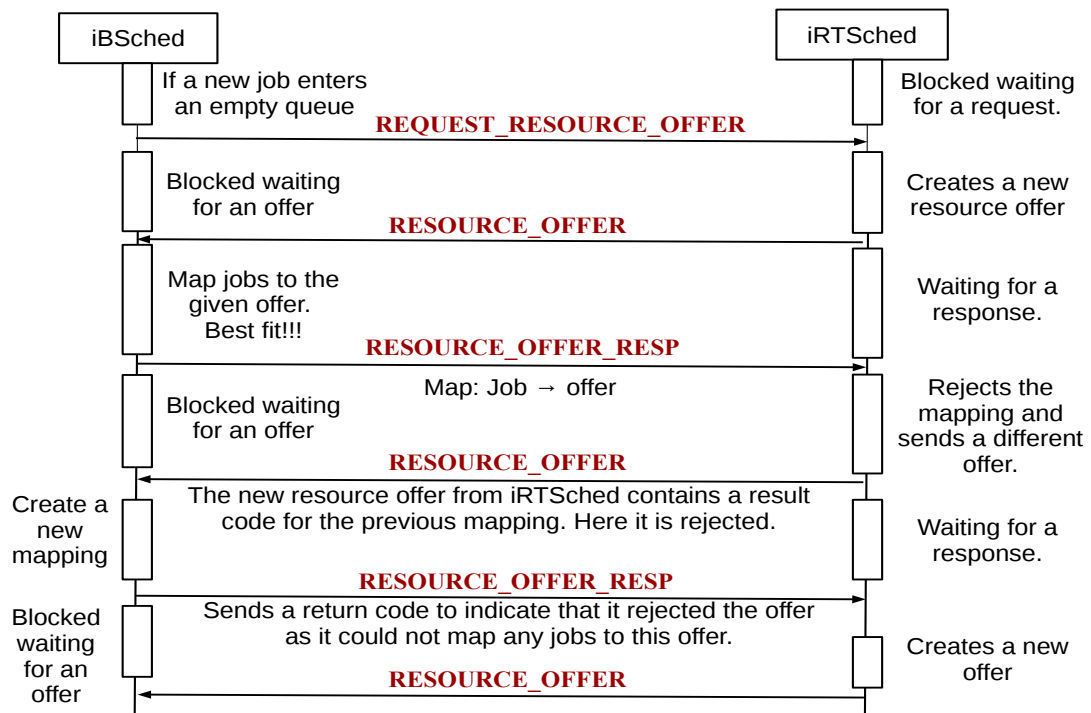


Figure 4.3: Scenario 1

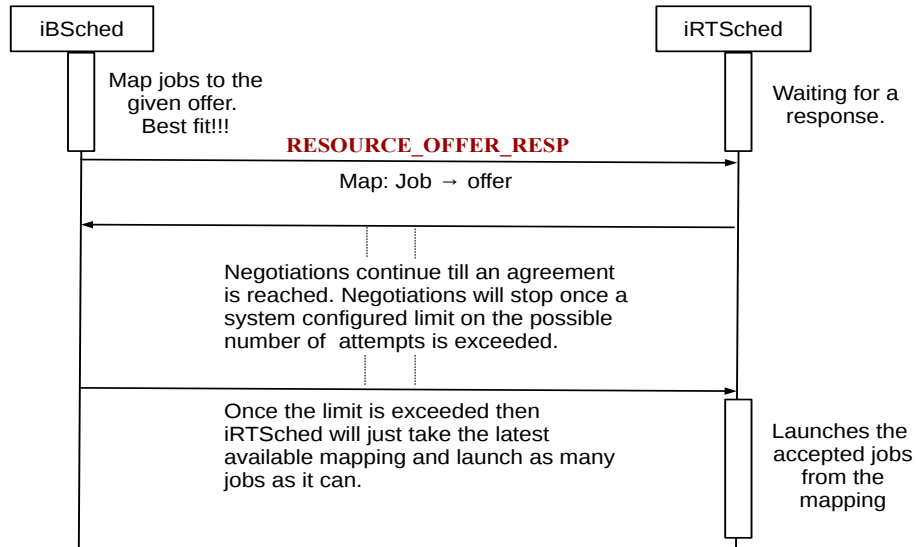


Figure 4.4: Scenario 1 contd.

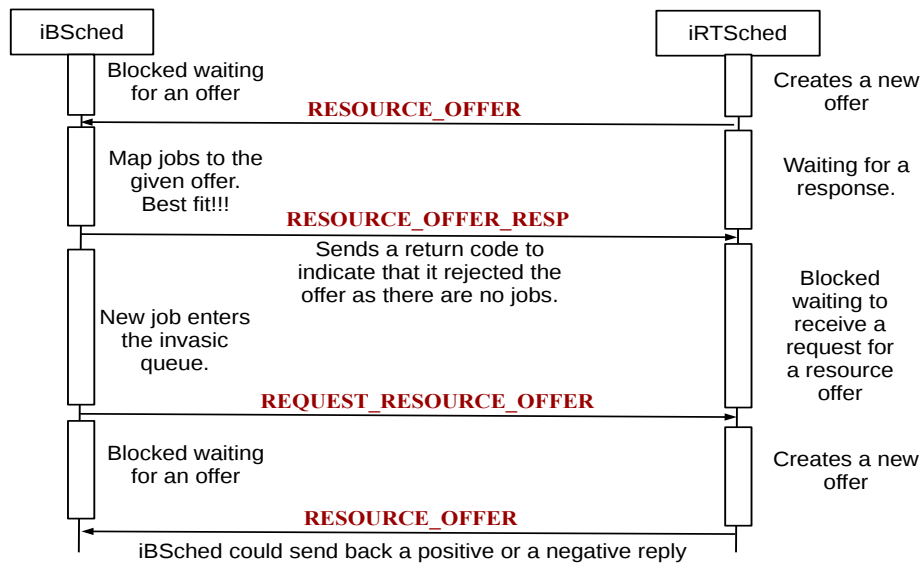


Figure 4.5: Scenario 2

4.3 Invasive Jobs

Invasive jobs refer to those jobs whose resource set can be changed at runtime. As mentioned before in the job classification, *malleable* and *evolving jobs* fall under this category and we will refer to these as invasive jobs in this thesis. These kind of jobs will usually come with a set of constraints when submitted to a HPC system. An example of how such constraints may look like is given below. Some of these constraints will already be supported by current batch systems for job submission.

Node Count: <min nodes><max nodes>

Runtime Estimate: <duration for min nodes>

Memory Per Node: <min MB>

Hint: <io bound | cpu bound | communication bound | na>

Urgent: <yes | no>

Deadline: <timestamp | na>

Start Time: <timestamp | na>

Budget: <core hours>

Performance Model: <if available>

na → not available

As mentioned above, these constraints are specific to a job provided by the user at its submission time. However, these can change or be updated after some feedback is received by the runtime scheduler once this job has finished running. Also, during the runtime of the job, the node count constraints may not hold good because the application is invasive in nature and may be much more scalable than the estimate provided by the user. These feedbacks can then be used to update the available job details in order to make better scheduling decisions in the future if the same job was to arrive again.

5 Design

In this section, we will describe the internal details of the components in the architecture specifically on how the scheduling algorithms work in order to realize the negotiation, the meaning of job mappings, resource offers and feedback reports.

5.1 Entities

A design entity is an element (component) of a design that is structurally and functionally distinct from other elements and that is separately named and referenced. They result from the decomposition of the software system requirements. The objective is to divide the system into separate components that can be considered, implemented, changed, and tested with minimal effect on other entities. Examples of design entities are: *protocol, application layer, state machine, data model, object, task, sub-systems, modules* etc.

Job Mappings

A job mapping is basically a job and its set of allocated resources. A forward job record is a job mapping and other related job details important for its launch. A list of such forward job records are sent by the batch scheduler in response to receiving a resource offer from the runtime scheduler.

$$\begin{aligned} Map &:= < Map_1 > < Map_2 > \dots < Map_n > \\ Map_i &:= Job\ Record := \{M_i, details_i\} \\ M_i &:= Job\ Mapping := \{jobid, nodecount\} \end{aligned}$$

In the above description $details_i$ refers to the details of any i^{th} job such as *constraints, priority, runtime estimate, max nodes, min nodes, job name, start time, deadline*. In these details, the constraints are specific to the purpose of negotiation and this can be related to node count that the job is requesting for, memory size, exclusive / shared resource access etc. Currently, we support only the constraints for node count. The node count constraint will again have a min and max both of which will fall within the window of max and min nodes of the job. It is using these constraints that the batch scheduler will negotiate with the runtime scheduler and with every negotiation attempt it will

relax the constraints to bargain better.

Resource Offers

A reservation is basically a job being guaranteed a set of resources for a certain period of time. A resource offer is a list of such reservations and the set of nodes which are free in the partition after the runtime scheduler has guaranteed these reservations. The runtime scheduler will send such a resource offer to the batch scheduler if it rejected its mapping or it is sending a new offer that it generated.

$$\begin{aligned} Offer &:= \langle Resv_1 \rangle \langle Resv_2 \rangle \dots \langle Resv_n \rangle, IdleNodes \\ Resv_i &:= Job\ Reservation := \{jobid, nodecount, start\ time, end\ time\} \\ IdleNodes &:= Free\ Nodes := \{nodecount\} \end{aligned}$$

The $Resv_i$ entry in the resource offer corresponding to a particular job id must match the job id of the Map_i entry in the forwarded jobs. This means that the sequence in which jobs were forwarded to the runtime scheduler will be the same sequence of jobs for which reservations will be sent to the batch scheduler.

Feedback Reports

Periodic or event-driven feedback reports on the latest status of the running jobs are sent by the runtime scheduler to the batch scheduler. It is event-driven for events such as job termination, job completion, acceptance of a mapping received from batch scheduler.

$$\begin{aligned} Feedback &:= \langle Report_1 \rangle \langle Report_2 \rangle \dots \langle Report_n \rangle \\ Report_i &:= Status\ Report := \{jobid, nodecount, runtime, end\ time, state, P_i\} \\ P_i &:= Performance\ Data := \{model, energy, io, memory, hint\} \end{aligned}$$

The feedback reports can also contain performance specific data such as energy consumption, performance model of completed jobs that can be used by the batch scheduler for making better decisions on same jobs submitted again in the future. For running jobs, the reports can contain similar information but most important ones would be the current assigned node count, expected end time, memory usage etc.

5.2 Scheduling Algorithms

Following pages present the pseudo code for both the batch and runtime scheduling algorithms.

5.2.1 Batch Scheduling

Algorithm 1 presents a high level pseudo code of the batch scheduling algorithm. Following points describe the algorithm:

- It takes as input a resource offer from iRTSched and decides whether to accept or reject it. Irrespective of the decision, the batch scheduler has to send a list of jobs to the runtime scheduler.
- In this list some or all of the jobs will be mapped to the resources offered and some of them may not have any mapping but will just be forwarded along to respect the order of jobs in the queue. This allows for *fairness* and *avoids starvation* when other small jobs may potentially keep overtaking these waiting jobs in getting mapped.
- We will refer to jobs which do not get mapped as a *reserved job* for our convenience whereas those which have resources allocated or have been mapped are called as *mapped jobs*. The mapped jobs and reserved jobs together form a list of jobs to be forwarded.
- At a high level, the algorithm can be viewed as a sequence of 3 important steps that will be described in the following points along with references to where they are located in the algorithm.
- *Line 1-33*: basically represents the *step 1* of the algorithm. *Line 1*: A separate job queue is constructed by scanning the main job queue for jobs which have been submitted specifically to the invasive partition. These jobs are then sorted according to their priorities.
- *Line 8-34*: For every job in the invasive queue: If it is found in the reservation list from the resource offer and if it can start immediately (has a node count allocated) then we will just append this to a new list (Map) after creating a new forward job record using the details of this job.
- If it starts in the future (has a zero node count) then we will relax its node count constraints by a step size calculated as below and try to map it to the freenodes provided in the offer.

$$step\ size = \frac{(Job.MaxNodes - Job.MinNodes)}{MAX_NEGOTIATION_ATTEMPTS}$$

- If *step size* in the above calculation comes up as 0 then we will set it as 1.

Algorithm 1: Batch Scheduling Algorithm

```

Input: Resource Offer
Output: Map : Jobs  $\rightarrow$  Offer
1 /* Invasive job queue consists of jobs submitted specifically to the
   invasive partition */
2 build invasive job queue
3 if empty queue then
4   | return empty Map
5 end if
6 sort queue by job priorities
7 reservations  $\leftarrow$  Offer.Reservations
8 freenodes  $\leftarrow$  Offer.ResidualNodes
9 while not at end of the invasive job queue do
10  | read next queue entry
11  | if entry.jobId in reservations then
12  |   | resv  $\leftarrow$  reservations(entry.jobId)
13  |   | if resv.node_count EQ 0 then
14  |   |   | entry.start_time  $\leftarrow$  resv.start_time
15  |   |   | entry.node_count  $\leftarrow$  0
16  |   |   | adjustNodeCount(entry.jobId)
17  |   | else
18  |   |   | entry.start_time  $\leftarrow$  0
19  |   |   | entry.node_count  $\leftarrow$  resv.node_count
20  |   |   | create forward job record using entry
21  |   |   | append record to Map
22  |   |   | mapped  $\leftarrow$  true
23  |   | end if
24  |   | end if
25  |   | if mapped NEQ true then
26  |   |   | try mapping to the freenodes
27  |   |   | if successfully mapped then
28  |   |   |   | update the count of freenodes
29  |   |   | end if
30  |   | end if
31  |   | create forward job record using entry
32  |   | append record to Map
33 end while
34 try_to_best_fit(Map, invasive job queue, Offer)
35 return Map

```

- $job.node_count.min = job.node_count.min - step\ size$. If the computed value is negative or less than the job's minimum number of nodes then we just set $job.node_count.min = job.min_nodes$.
- If the job is not found in the reservation list then this must be a new job that was submitted to the batch scheduler in between the negotiations and was not considered during the previous negotiation attempt. And, this job has been placed near the front of the queue owing to high priority or because it is a small job etc. We will directly try to map it to the freenodes in the offer.
- If a job could not be mapped, then we will just append it to the Map by creating a new forward job record. If it was mapped, then we will set the current node count of that job to the count of the allocated nodes. The node count constraints remain the same.
- **Line 34:** This routine basically represents the *step* 2 and 3 of this algorithm. The result of these two steps is an approximation of a best fit mapping of the invasive jobs ready to be forwarded to iRTSched.
- **Algorithm 2** This algorithm shows how the best fit mapping is computed. **Line 8:** We will try to map reserved jobs from the newly constructed job list (Map) by reducing the node counts of jobs that have been currently mapped within their constraints. This operation considers mapped jobs in the reverse order of their priority so that the operation can start from the lowest priority mapped job till the highest one.
- If we were successful in mapping any new reserved jobs by the above operations, then we must rescan the invasive job queue to update the new job list by jobs which had been previously skipped in *step* 1 due to a limit on the *reservation depth*. Reservation depth means that we have a limit on the number of reserved jobs that can be forwarded to the runtime scheduler. This is the *step* 3 of the algorithm.
- **Algorithm 3** This algorithm shows how the analysis of mapped jobs are done in order to reduce their allocated node counts to satisfy a reserved job.

$$step\ size = \frac{(Job.node_cnt - Job.node_count.min)}{MAX_NEGOTIATION_ATTEMPTS}$$

Batch scheduler will provide *fairness* to jobs while scheduling by considering them in the order of decreasing priorities. It also avoids any *job starvation* by forwarding jobs that could not get mapped to the available resource offer. This will result in the runtime

scheduler giving these jobs future start times as per its backfill algorithm and only then process other job requests after it in the list of forwarded jobs. Batch scheduler will also adjust the expected end times of those jobs that it will transform by expand / shrink. This adjustment will take into consideration the performance model of the running application.

Decision Logic The batch scheduler currently has a very simple decision making routine. It tries to maximize the sum of the priorities(objective) of jobs that have been mapped to some resources in the offer. Below are the key points:

- Batch scheduler will keep sending a reject if the job at the front of the queue could not be mapped or atleast some fraction of the total number of jobs in the queue have not been mapped(ex: one-third of the total jobs). But, if the number of negotiations have crossed half the limit, then it will ignore these constraints.
- If the batch scheduler had sent an accept decision on the previous negotiation attempt, then it will send an accept again in the current attempt if the objective value has remained the same or improved. It will also stop relaxing the node count constraints of jobs for further negotiations in case the number of jobs which have been mapped are more than the previous attempt.
- If the batch scheduler had sent a reject on the previous negotiation attempt, then it will send a reject again in the current attempt if the objective has reduced. If the objective has remained the same or improved, then it will accept if it does not violate the first point.

5.2.2 Runtime Scheduling

Algorithm 4 presents a high level pseudo code of the algorithm to generate a resource offer. Following points describe the algorithm:

- If batch scheduler accepted the previously sent offer, then we will repeat the system transformation of the previous negotiation attempt to check if runtime scheduler will accept this mapping. If the mapping was accepted, then the runtime scheduler will commit the mapped jobs to the running list.
- If this is the first attempt and there are no forwarded jobs then it means that the runtime scheduler is generating a new offer. In this case it will perform a partial transformation and will return back from the function to send the offer.

Algorithm 2: Best Fit Algorithm**Input:** Map, Invasive Job Queue, Offer**Output:** Updated Map : Jobs \rightarrow Offer

```
1 avail_bitmap  $\leftarrow$  Offer.ResidualNodes
2 while not at end of the Map do
3   read Map entry
4   if (entry.start_time EQ 0) AND (entry.bitmap NEQ NULL) then
5     | continue
6   end if
7   /* Try to map this job by shrinking other mapped jobs */
8   try_sched(entry, Map, avail_bitmap)
9   if successfully scheduled then
10    | update avail_bitmap
11  end if
12 end while
13 if avail_bitmap not empty then
14   | Rescan the invasive job queue to fill the residual nodes with some new jobs
15 end if
```

Algorithm 3: Try Schedule Algorithm**Input:** Entry, Map, avail_bitmap**Output:** Updated entry: Mapped to some nodes

```
1 /* The mapped jobs in the Map would be analyzed in the reverse order
   which is increasing in the priority */
2 Analyze in increasing order of priority if the mapped jobs in the Map can be
  shrunk to find enough nodes for entry
3 if sufficient nodes available then
4   | shrink the mapped jobs as per the analysis
5 end if
6 entry.bitmap  $\leftarrow$  bitmap(available nodes)
```

Algorithm 4: Algorithm for generating a resource offer

Input: Attempts, Jobs2Map, Error Code, Offer<empty>
Output: Offer<Reservations, Residual nodes>

```

1 if Error Code is SUCCESS then
2   /* Batch Scheduler accepted the previously sent offer */
3   if (Jobs2Map NEQ NULL) AND (Attempts GT 1) then
4     initialize runtime state
5   end if
6   /* Repeat the transformation of the running jobs which was done
   for the previous negotiation attempt */
7   schedule((Attempts GT 1) ? (Attempts - 1) : Attempts, Jobs2Map, Error
Code, Offer)
8   /* Offer being generated for the first time if attempts is equal
   to 1 */
9   if Attempts EQ 1 then
10    return SUCCESS
11  end if
12  if schedule was successful AND Jobs2Map NEQ NULL AND Attempts GT 1
then
13    commit the mapped jobs to the running list
14    return SUCCESS
15  end if
16 end if
17 /* Batch Scheduler rejected the previously sent offer */
18 if Error Code NEQ SUCCESS then
19   initialize runtime state
20   schedule(Attempts, Jobs2Map, Error Code, Offer)
21 end if
22 return error code

```

- If the mapping was rejected after repeating the previous transformation, then the runtime scheduler will reset the runtime state of the system and perform a new transformation for generating an offer to satisfy the forwarded job requests.
- If the response from batch scheduler was reject, then the runtime scheduler will go through the same step as described in the previous point.

Algorithm 5 presents a high level pseudo code of the runtime scheduling algorithm. The runtime scheduling algorithm in this thesis is based on the *DBES* algorithm[29]. It has been adapted for negotiation based scheduling in this work. Following points describe the algorithm:

- *Line 1:* Run the backfill algorithm to schedule the forwarded job requests. Some jobs can immediately start whereas for others it will give reservations. Both of these are written into the resource offer.
- *Line 2:* If there were no forwarded job requests then this was a new offer being generated.
- *Line 5:* Update job dependencies according to the new system state. This means that some running malleable jobs may be maintaining invalid job dependencies. If a forwarded job can immediately start and a running malleable job has a dependency from the same job id, then this dependency should be cleared now as it is no longer valid.
- *Line 6-9:* For every reserved job, Analyze running malleable jobs according to the approach mentioned in the pseudo code for shrink operations to generate sufficient resources for immediate start of reserved jobs. If analysis was successful, then shrink the jobs.
- *Line 14:* Reschedule forwarded job requests using backfill algorithm and compute reservations(without job start). In this step we will not consider jobs which could get immediate start time in the previous steps. Update the resource offer accordingly.
- *Line 15:* Update job dependencies according to the new system state similar to point 3.
- *Line 16-20:* For every reserved job, if it has a dependency on a running malleable job then we will expand that running job to the maximum possible node count and update the dependency on that running malleable with the job id of this forwarded job. Expansion of the running job will also consider shrinking other running jobs for its needs. For the purpose of shrinking, it will follow the same approach as shown in *Line 7* except that it will follow only the second and third steps.
- *Line 25:* Once all the previous steps of the algorithm have been completed, runtime scheduler will finally decide whether to accept or reject the mapping based on some decision making logic. If the number of negotiation attempts have

reached the limit then it will just accept. If it rejects then it will send back the constructed offer to iBSched.

Runtime scheduler can generate a resource offer for the following events: *job termination, job completion, request for a resource offer, periodic runtime transformations of the running jobs resulting in free resources.*

A resource offer is sent to the batch scheduler on every alternate periodic time step. This is just a very simple *fairness* policy to equally favor the running jobs and batch jobs. Between any two transactions, the runtime scheduler will perform a transformation of running jobs based on their performance and also sometimes based on their dependencies (due to the DBES algorithm). During the other periodic time steps, The runtime scheduler will initiate a transaction with batch scheduler by creating an offer and sending it. It will do this by performing a partial transformation which means that it will shrink those running jobs which are shrinkable based on their performance but not use the free resources to expand other jobs and instead use it to send an offer. These partial transformations will only be done on every alternate periodic time step in order to favor running jobs and batch jobs equally. Sending an offer at every periodic time step may degrade the performance of running applications by driving them towards their minimum node counts whereas sending an offer too infrequently may effect throughput and quality of service to batch scheduler.

Runtime Transformation Due to the usage of DBES algorithm, a running malleable job always keeps track of the job id of a batch job that did not start till now but had its dependency on this running job. During the periodic runtime transformations (not partial), if there are any running malleable jobs which have such dependencies then they will be given complete preference for expansion at the cost of shrinking all other jobs without any dependencies. In case there are no such running jobs with dependencies then transformation will be done solely on the basis of performance and scalability.

Algorithm 5: Runtime Scheduling Algorithm

Input: Attempts, Jobs2Map, Error Code, Offer
Output: Updated Offer

```

1 schedule requests and create reservations(without job start)
2 if Jobs2Map EQ NULL then
3   | return /* new resouce offer generated */
4 end if
5 update job dependencies according to the new system state
6 for each reserved job do
7   | Prioritize malleable jobs in the order: (1) malleable job expanded for this
8   | reserved job, (2) malleable job expanded for no specific reserved job, (3)
9   | malleable job expanded for other reserved jobs
10  | Analyze if expanded malleable jobs can be shrunk in the above order to
11  | make enough nodes available to start the reserved job
12  | if enough nodes found then then
13  |   | Shrink the selected malleable jobs
14  |   | Insert the job as an entry in Offer.Reservations
15  | end if
16 end for
17 reschedule requests and create reservations(without job start)
18 update job dependencies according to the new system state
19 for each reserved job do
20   | if job depends on a malleable job then
21   |   | Expand the malleable job with the available nodes
22   | end if
23 end for
24 if Error Code is SUCCESS then
25   | /* Runtime scheduler will decide now whether to accept / reject
26   | the mapping */
27   | if Attempts LT MAX_LIMIT then
28   |   | decision  $\leftarrow$  decision_logic(Offer, Jobs2Map, Attempts, count(idle nodes))
29   | end if
30 end if
31 if Attempts EQ MAX_LIMIT then
32   | decision  $\leftarrow$  accept
33 end if

```



```
30 if Error Code is SUCCESS AND decision is reject then  
31   |   return  
32 end if  
33 equipartition the available idle nodes among other remaining running malleable  
   jobs  
34 return from the function with newly constructed resource offer the job mapping  
   received
```

Decision Logic: Runtime scheduler like batch scheduler uses a very simple decision making routine. This is due to the fact that it does not manage any resources now but only simulates it for the early prototype developed in this thesis. Below are the rules it follows to accept or reject a mapping:

- If there are any idle nodes left after the complete runtime scheduling algorithm has run, then the received mapping will be rejected.
- If only one job could be immediately started out of the many that have been forwarded, then it will reject the mapping else accept. This constraint will be ignore, however, if the number of negotiations have crossed half the limit.
- If it rejected on any negotiation attempt, then on subsequent attempts it will keep rejecting if the number of jobs immediately starting does not improve.
- If the batch scheduler had rejected the previous offer but accepted the current offer, then runtime scheduler will again employ the same rules for making a decision as enumerated in the last three points.

5.3 Negotiation

The figure 7.1 illustrates one possible scenario of the negotiation between batch and runtime scheduler. Following points describe it in detail:

- The alphabets **A,B,C,D,E** represent runtime state (running jobs and idle nodes) of the partition. **TRF** means transformation of the running jobs through expand and shrink operations. This happens when a list of forwarded job requests are received from batch scheduler and the runtime scheduler runs its algorithm (based on DBES) for expanding or shrinking running jobs in order to fit as many batch jobs as possible.

- **INIT** means initialize state. It will restore the transformed jobs back to their original state which they were at the beginning of the transaction (set of negotiations). This is done in order to perform a new transformation from the original state during every negotiation attempt otherwise the transformation will reach a saturation very soon.
- Saturation means that we can no longer perform any significant expand or shrink operations on the running jobs as they will already be either be at their min or max node counts due to a result of the previous transformations. This would result in making no progress.
- The green boxes labelled "**Algo 1**" represent the batch scheduling algorithm which will run every time the batch scheduler receives a resource offer. On every such attempt, it will relax the node count constraints of those jobs which could not be mapped due to lack of sufficient resources.
- The box labelled "**Update**" will update the details of jobs in its queue once it receives a feedback. This is important as a subsequent negotiation attempt must not result in the batch scheduling algorithm dispatching a job that is already running. The feedback is the only way for the batch scheduler to recognize if the runtime scheduler accepted its mapping and started any of the jobs from the mapping it had sent.
- Once the negotiation has completed, the runtime scheduler will then accept the mapping and commit the jobs. The box labelled "**COMM**" represents the step where the commit happens. The forwarded jobs are added to the list of running jobs and would be started very soon. The runnings jobs will now take up their transformed state going forward as shown by the box **E**.

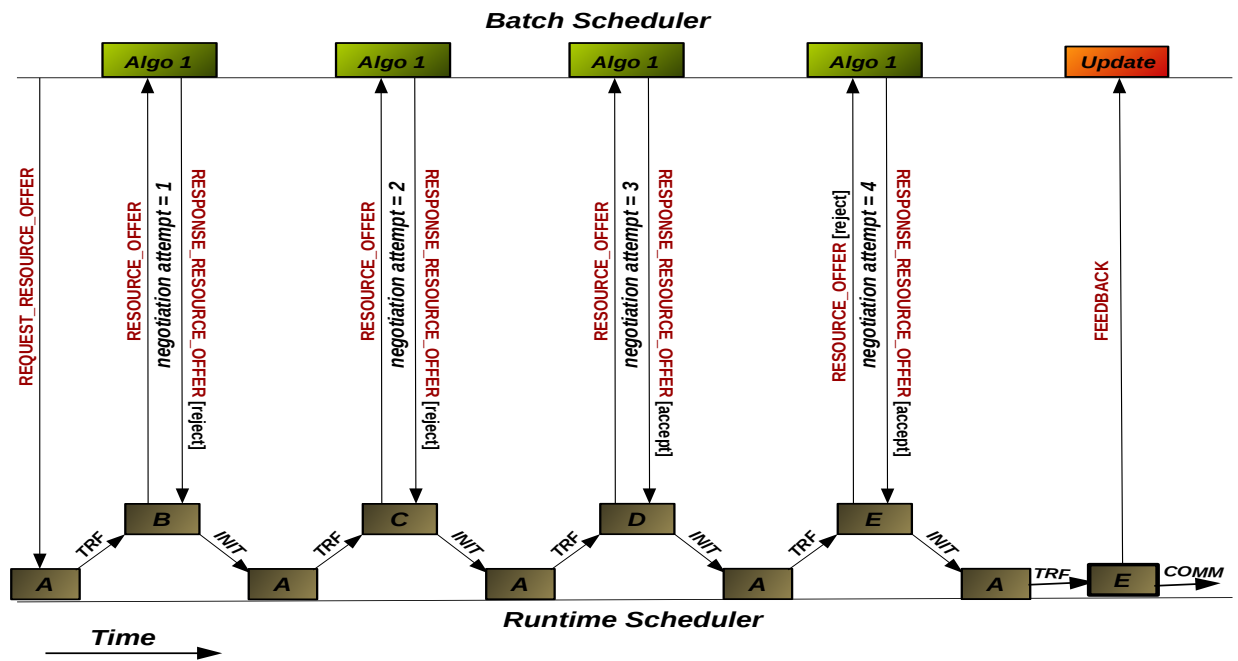


Figure 5.1: Negotiation protocol

6 Implementation

6.1 Plugin

iBSched is implemented as an extension to SLURM in the form of a scheduler plugin. A Slurm plugin is a dynamically linked code object which is loaded explicitly at run time by the Slurm libraries. A plugin provides a customized implementation of a well-defined API connected to tasks such as authentication, interconnect fabric, and task scheduling. There are several other existing scheduler plugins of SLURM such as: the default(FCFS) plugin, Backfill, wiki and wiki2(Interfaces to external schedulers) etc. A related plugin that is of importance to scheduling is the priority plugin because the job queue is first constructed and then sorted based on priorities before the scheduling algorithm runs through it. Slurm provides a basic priority plugin and a multifactor priority plugin. The basic version provides a standard FIFO based job priority whereas the multifactor version provides priorities to job based on several factors such as age, fair-share, job size(also considers the job duration to calculate a weight factor for job size relative to time), partition, quality of service etc. The multifactor priority plugin has been refactored for this work to be better suitable for our implementation. Specifically, the multifactor priority plugin now assigns priority to jobs based on age and size(relative to time). As its default behavior it periodically re-calculates the priorities of the jobs.

For the current scope of the work, iRTSched has been implemented as an independent binary that talks to iBSched via the negotiation protocol. The purpose of this independent binary is to fake an actual runtime scheduler by simulating the runtime states of the system in time by employing the runtime scheduling algorithm and also negotiating with iBSched. In reality, a full-fledged runtime scheduler will be managing the resources in the invasic partition, will launch the jobs and perform dynamic resource management. In order to evaluate the negotiation-based separation of batch and runtime scheduling system, a simulation of the same would be very useful and sufficient for our objectives.

6.2 Data Structures

This section will enumerate some of the important data structures that are a part of the implementation.

Job Mappings

Below is enum that represents the possible characteristics of a job. This hint can be supplied at job submission time or can be inferred during runtime and profiling the application performance. By default it will be considered as NONE.

```
enum job_hint {  
    COMP_BOUND,  
    MEMORY_BOUND,  
    IO_BOUND,  
    NONE  
};
```

Below is an enum that represents the type of a job. This needs to be specified at submission time by the user. By default it will be considered as rigid.

```
enum job_type {  
    RIGID,  
    MOLDABLE,  
    EVOLVING,  
    MALLEABLE  
};
```

Below structure represents the node count constraint that is used by the batch scheduler during the negotiations with runtime scheduler. At the submission time max will be set to the job's max nodes and min will be set to a value of max nodes - step_size where step_size is calculated from the formula given in x.x. This step will be repeated when a new transaction starts.

```
typedef struct qos_rsrc {  
    uint32_t min;  
    uint32_t max;  
} qos_rsrc_t;
```

Below structure represents the details of a job necessary for its launch. This is packed within the following structure and it gets forwarded to the runtime scheduler.

```
struct forward_job_details {  
    qos_rsrc_t node_count; /* Job resource requirements */  
};
```

```

uint16_t cpus_per_task; /* number of processors required for
                        * each task */
uint8_t hint; /* Job characteristic: IO / Compute / Memory bound
              */
uint8_t job_type; /* Type of job: Static or Dynamic */
/* Ideally this structure should further contain all the necessary
   details required for launching this job. For the purpose of
   this thesis, we are only doing a simulation without running
   jobs, hence other details in this structure are not important
   */
};

```

The below structure makes up a single entry in the Map that is constructed by batch scheduler. A list of such entries make up the mapping of jobs to offer which is sent to the runtime scheduler. In addition to the details of the job as mentioned for the previous structure, this structure include details such as the bitmap of nodes allocated, priority, job id, min / max nodes etc.

```

struct forward_job_record {
    struct forward_job_details *details; /* job details */
    uint32_t job_id; /* job ID */
    uint32_t magic; /* magic cookie for data integrity */
    char *name; /* name of the job */
    bitstr_t *node_bitmap; /* bitmap of nodes allocated to job */
    uint32_t node_cnt; /* Current node count assigned */
    uint32_t min_nodes;
    uint32_t max_nodes;
    uint32_t priority; /* relative priority of the job,
                      * zero == held (don't initiate) */
    time_t start_time; /* Expected or Actual start time */
    uint32_t time_limit; /* time_limit minutes or INFINITE,
                       * NO_VAL implies partition max_time
                       */
    uint32_t time_min; /* minimum time_limit minutes or
                      * INFINITE,
                      * zero implies same as time_limit */
};

```

Below structure represents the message used by batch scheduler to request for a resource offer. It contains a value field used for the purpose of protocol communication followed by the most important field which is the Map.

```
typedef struct request_resource_offer_msg {
    uint16_t value; /* info */
    List jobs2map; /* This is the list of jobs waiting to be
        dispatched. And communicates the current requirements to rt
        scheduler which
            can then try to suitably construct a resource
            offer to satisfy the requirements as best as
            possible */
} request_resource_offer_msg_t;
```

Below structure represents the response to a resource offer sent by the batch scheduler. It resembles very closely to the previous message but has in addition the error code and the error msg fields. These fields which identify the error are important for batch scheduler to convey its response to the runtime scheduler for the offer it sent.

```
typedef struct resource_offer_resp_msg {
    uint16_t value; /* info */
    List map_jobs2offer; /* Jobs mapped to the given offer. This
        depends on whether the batch scheduler accepted / rejected /
        countered
            the resource offer it received */
    uint32_t error_code; /* error code on failure */
    char * error_msg; /* error message on failure */
} resource_offer_resp_msg_t;
```

Resource Offers

Below structure represents the reservation for a job which can begin immediately or in the future. This depends on the response from the runtime scheduler.

```
typedef struct job_resv {
    time_t end_time; /* end time of reservation */
    uint8_t full_nodes; /* when reservation uses full nodes or not */
    uint32_t job_id; /* job ID */
    bitstr_t *node_bitmap; /* bitmap of reserved nodes */
    uint32_t node_cnt; /* count of nodes required */
    time_t start_time; /* start time of reservation */
} job_resv_t;
```

Below structure represents the message format for resource offer message sent by runtime scheduler. It contains a list of reservations for each of the jobs that were sent by batch scheduler in its mapping. It contains the set of residual or free nodes available.

```
typedef struct resource_offer_msg {
```

```

uint16_t value; /* info */
uint8_t negotiation; /* if negotiation is ongoing then this value
    is 1 else it becomes 0 to indicate ischeduler that previous
    negotiat
        ion is over */
List resource_offer; /* List of node space records, currently it
    has only one entry for residual nodes */
List resv_jobs; /* List of jobs with actual reservations. Can also
    include jobs with virtual reservations that are those with
        future start / service times */
uint32_t error_code; /* error code on failure */
char * error_msg; /* error message on failure */
} resource_offer_msg_t;

```

Feedback Reports

Below structures represent the format of a feedback report and a list of reports respectively. For the scope of this current work, the feedback does not send back any performance model of the running / completed application or any kind of job characteristic and energy consumption. It sends across basic details about the status of running / completed jobs.

```

typedef struct job_status {
    uint32_t job_id;
    time_t run_time;
    time_t end_time;
    uint32_t node_cnt;
    bitstr_t *node_bitmap;
    uint32_t job_state;
} job_status_t;

typedef struct status_report_msg {
    List status_reports;
} status_report_msg_t;

```

Running Job Record

The below structure is very important as the runtime scheduler uses this as the job record for each of the running jobs. Since the scope of this work is simulation, We do not include many other details of a job that would be necessary for its launch, logging etc. Dependency of a forwarded job onto a running malleable job is saved in the variables `depend_job_id` and `depend_job_prio` of a running job record. This dependency information is used by the PDBES algorithm.


```
struct run_job_record {
    uint32_t job_id;
    bitstr_t *node_bitmap;
    time_t start_time;
    uint32_t priority;
    uint32_t time_limit;
    uint8_t hint;
    uint8_t job_type;
    time_t end_time;
    time_t orig_end_time;
    uint32_t job_state;
    uint32_t orig_node_cnt;
    uint32_t node_cnt;
    uint32_t last_node_cnts[2];
    bitstr_t *next_node_bitmap;
    uint32_t next_node_cnt;
    uint32_t depend_job_id;
    uint32_t depend_job_prio;
    uint32_t save_depend_job_id;
    uint32_t save_depend_job_prio;
    uint32_t min_nodes;
    uint32_t max_nodes;
    uint8_t adapt; /* 0 - No change, 1 - Expand, 2 - Shrink */
    uint8_t transformed;
    time_t exp_end_time;
};
```

Node Space Map

Below is an existing structure used by SLURM for its backfilling algorithm. This is used to form a chain of records chronologically ordered in time that represent the set of available nodes in those timeslots. The backfill algorithm uses this structure to prepare the view of resources in time and resembles a two dimensional space of nodes and time. This is used in order to compute when and where a job can start by looking at the reservations that higher priority jobs already have and these must be respected by lower priority jobs in case they can start now.

```
typedef struct node_space_map {
    time_t begin_time;
    time_t end_time;
    bitstr_t *avail_bitmap;
```

```
    int next; /* next record, by time, zero termination */  
} node_space_map_t;
```

6.3 Important APIs

```
extern List schedule_invasic_jobs(resource_offer_msg_t *, List, uint16_t  
    *, uint32_t *);
```

This routine is responsible for creating a map of jobs to the available offer. It will fill the offer that has been passed as a an argument to it with the list of invasive jobs that have been selected according to its algorithm. Jobs in this list would be mapped to some number of nodes that satisfy the job constraints. Those jobs which could not be mapped will have their constraints relaxed and just forwarded along with other jobs. Also, The batch scheduler has a limit on the number of jobs that can be forwarded but have not been mapped to any resources as it did not satisfy its constraints. This is the reservation depth similar to what is used in backfill algorithms. Once the limit has been reached the scheduling algorithm will then just scan the rest of the invasive job queue to look for jobs that can fit the remaining available nodes in the offer.

```
extern uint32_t adjustQoS_node_count(struct job_record *);
```

This routine relaxes the node count constraint of a job. The original constraints are min and max nodes supplied at the job submission time. Negotiation begins by setting a node count constraint(node_count.min, node_count.max) within this window of min and max. With more negotiation attempts in a single transaction, the constraints would keep getting relaxed and node_count.min would approach the original min value.

```
extern void resetQoS_node_counts();
```

This routine is called before the start of a new set of negotiations. It will reset the node count constraint for every job in the queue to their original values.

```
extern int _decision_logic(List, int, uint32_t);
```

This routine is where the batch scheduler makes a final decision on whether the received resource offer from runtime scheduler is accepted / rejected. It does so by checking through the mapping that has been constructed after the scheduling algorithm

processed the offer. The mapping is supplied as the first argument to this function.

```
extern void *irm_agent(void *);
```

This is the routine that is associated with the main runtime scheduler thread. It gets called once the thread has been created. It basically runs an infinite control loop for the scheduler to initiate various operations such as processing the response for an offer from batch scheduler, performing a transformation of the runtime system state, waiting for a request for resource offer in a separate spawned thread, sending back an offer to batch scheduler, terminating the scheduler etc.

```
extern int process_resource_offer(resource_offer_msg_t *, uint16_t *, int  
    *, List *, List);
```

This routine initiates the processing of a resource offer, calls the appropriate scheduling algorithm.

```
extern int _request_resource_offer(slurm_fd_t); /* Wrapper for  
    request_resource_offer to construct the list of job requirements to be  
    sent */
```

This routine sends a request for resource offer to the runtime scheduler. It does so by constructing a list of job requests to be forwarded.

```
extern int process_rsrc_offer_resp(resource_offer_resp_msg_t *, int,  
    resource_offer_msg_t **);
```

This routine is responsible for processing the mapping received from the batch scheduler and invoking the appropriate runtime scheduling algorithm. If the response from batch scheduler was accept for its previous offer, then it will recreate the transformation of the previous attempt else it will create a new transformation. It will run the algorithm to satisfy as many jobs forwarded as possible by using the PDBES algorithm and then determine whether it will accept / reject this mapping. If it does not accept then it will send a new offer created as a result of this new transformation.

```
extern int create_resource_offer(int, List, uint16_t,  
    resource_offer_msg_t **);
```

This routine is called from within the "process_rsrc_offer_resp" routine to initiate the scheduling algorithm to run through the mapping. It is also called in the "irm_agent" routine in situations when a request for a resource offer is received. Or, a new offer needs to be generated by the runtime scheduler since some running applications have shrunk to create a gap in the resources that can be utilized to serve some new batch jobs.

```
static int
_schedule(int attempts, List job_queue,
          uint16_t *recv_err_code, resource_offer_msg_t **
          rsrc_offer_msg);
```

This is the runtime scheduling algorithm whose steps have been described in the form of a pseudo code x.x.x. It will use the PDBES algorithm to scan through the mapping and provide start times for each of them.

```
extern bitstr_t * _try_transf(bool);
```

This routine performs a runtime transformation of the system state by performing expand / shrink of the running malleable jobs considering their performance and scalability. In case the result of this transformation are some free resources, then those will be sent up to batch scheduler as a new resource offer. The boolean argument passed here represents the negotiation flag. If it is true then it means we have received a request for an offer and must perform a partial transformation by doing only the shrink operations to create a gap in the resources. If it is false then it means that we must do the full transformation as no offer is being sent out.

```
extern int _commit_state(bool);
```

This routine will commit the forwarded jobs(with immediate start time) to the running jobs list. It will also change the node bitmaps of running jobs to their new node bitmaps if they have been subjected to some sort of transformation during the phase of runtime scheduling algorithm. In alignment with the new node counts for many of the running malleable jobs, their expected end times will also be updated based on whatever performance model is known until this point.

```
extern bool _get_delta(uint32_t *, time_t, time_t);
```

This routine will compute the next time slice for which the runtime scheduler will sleep. This could be a fixed time interval for periodic wake ups of the runtime scheduler or it

can wake up earlier if there was some job that is expected to complete within this time interval.

```
extern void _initialize_state(void);
```

This routine resets the state of the system back to the point before the start of the negotiation. During the negotiation process, a lot of the details in the running job records would be temporarily updated to account for negotiations and the associated expand / shrink operations. With every negotiation attempt within a transaction, the state will have to be reset before the runtime scheduler can run its algorithm again.

```
extern int _decision_logic(resource_offer_msg_t *, List, int, bool);
```

The runtime scheduler decides whether to accept / reject the mapping received from batch scheduler by calling this routine. This routine checks through the mapping to see how many jobs can be started immediately and based on how well it is suitable for its metrics a decision will be taken.

6.4 Control Flow Diagrams

Following pages give the state machines diagrams of both the batch and the runtime scheduler. These diagrams illustrate their implementation and the general flow of their working highlighting the important details. Many of the error handling or subordinate threads related to feedback, urgent job handling, signal handling for termination etc. have not been shown here due to space constraints. The diagrams in both the pages represent the main control thread for both the batch and runtime scheduler. This control thread is responsible to spawn other subordinate threads, drive the control loop responsible for doing / initiating all the work such as receiving processing and sending protocol messages.

iBSched is a SLURM plugin, hence it is dynamically loaded by the SLURM controller(slurmctld) due to which a plugin initialization happens for iBSched during which a main plugin thread is created. This plugin thread is then responsible to create the main control thread for iBSched. The plugin thread is responsible for starting and shutting down this plugin(in case slurmctld sends a shutdown signal) and other threads it spawned. The main control thread is responsible for shutting down all the slave threads it had spawned. iRTSched also has multithreaded design exactly similar to iBSched except that it is an independent binary and is launched independently but not by any SLURM application.

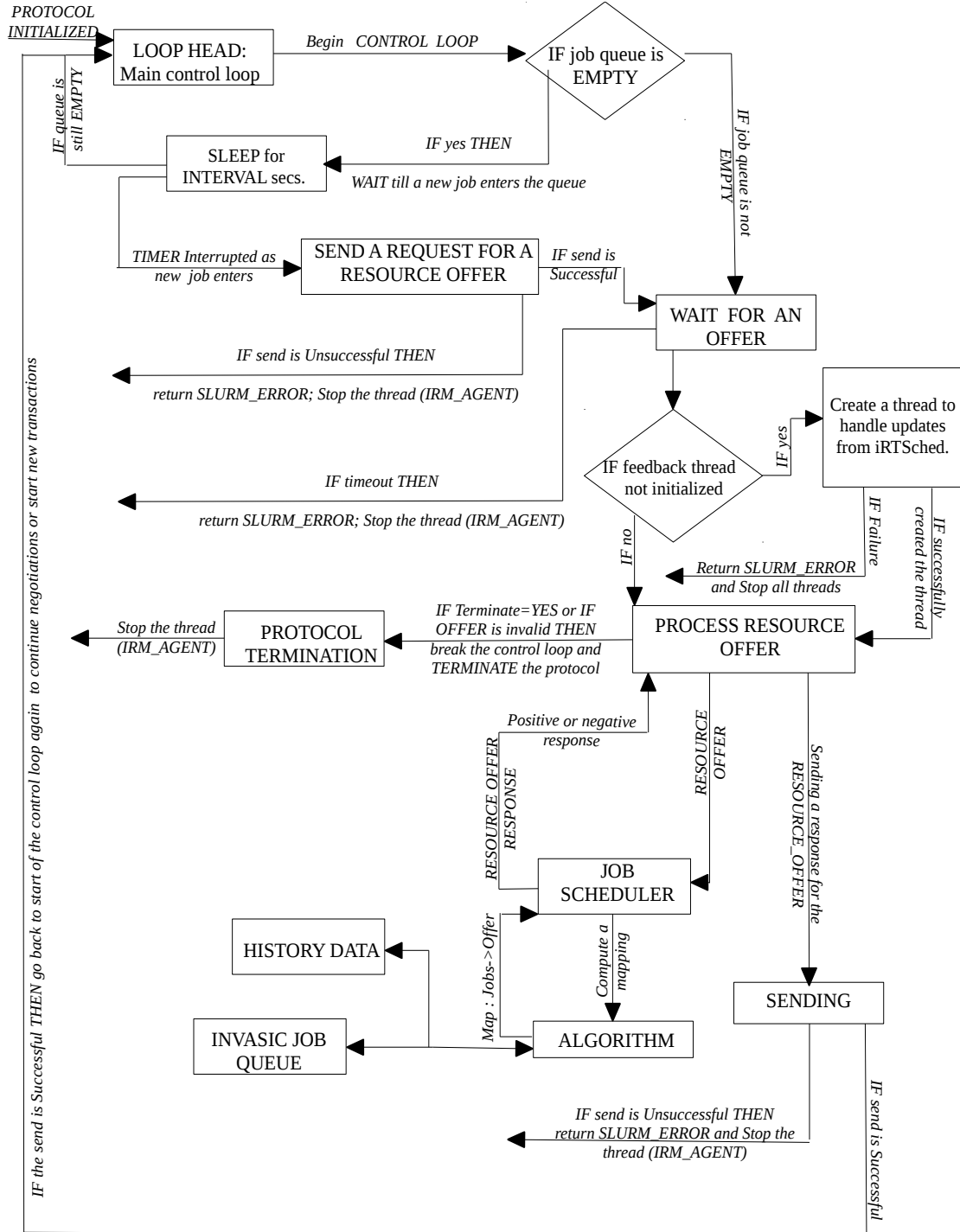


Figure 6.1: iBSched

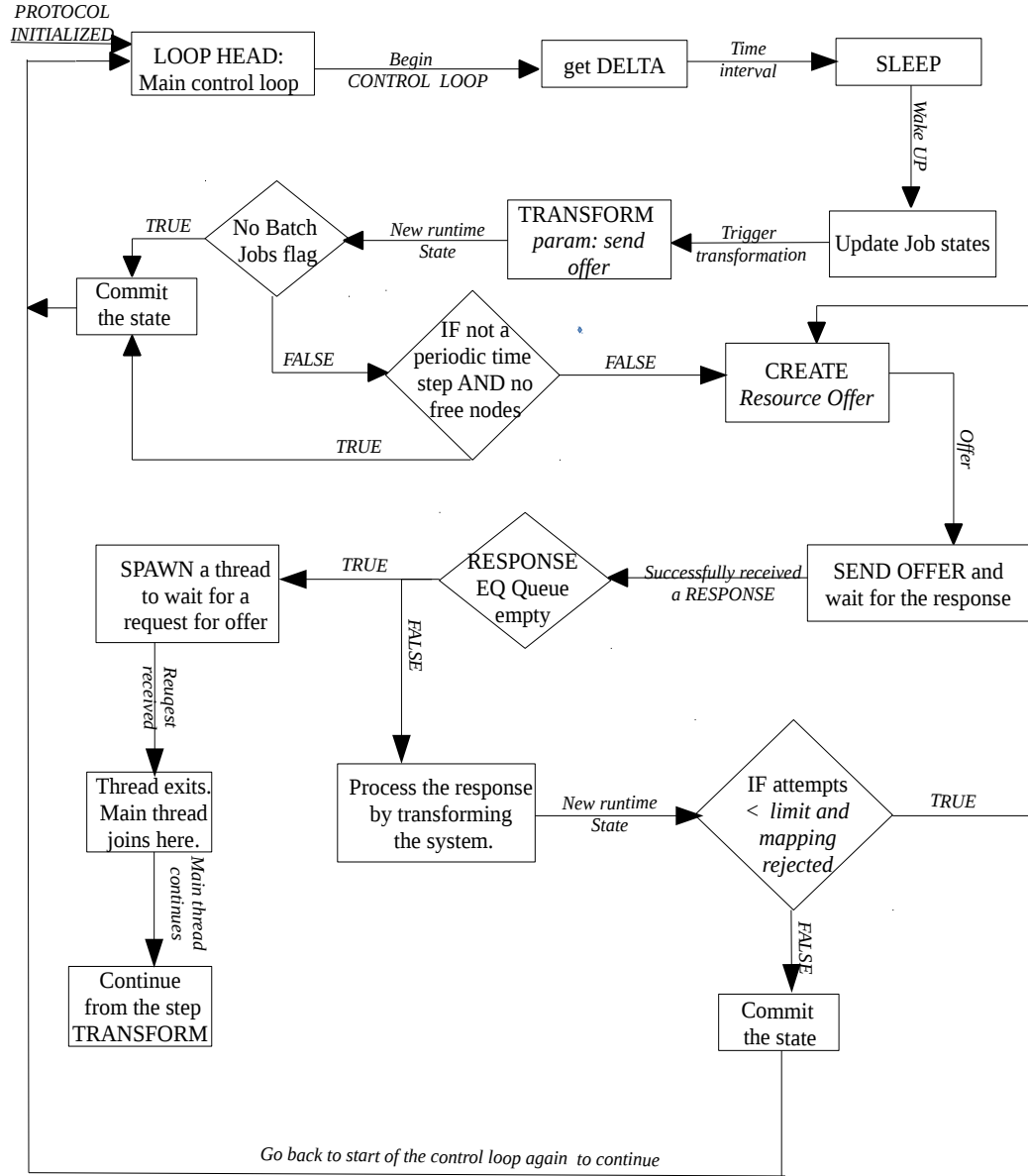


Figure 6.2: iRTSched

7 Evaluation

<May be a page illustrating a scenario of negotiation with 2d node space map diagrams will be good.>

The evaluation of job schedulers depend on two important factors: use of appropriate metrics, and the use of appropriate workloads on which the scheduler can operate. Before we consider these two factors in details, the following definitions will be assumed going forward.

- t_i^a is the arrival or submission time of job i.
- t_i^s is the start time of job i.
- t_i^e is the end time of job i.
- w_i is the width (number of requested / used resources) of job i.

From these parameters are computed:

- $l_i = t_i^e - t_i^s$ is the runtime(length) or duration of job i.
- $t_i^w = t_i^s - t_i^a$ is the waiting time of job i.
- $t_i^r = t_i^w + l_i$ is the response time of job i.
- $s_i = \frac{t_i^r}{l_i} = 1 + \frac{t_i^w}{l_i}$ is the slowdown of job i.
- $a_i = w_i \cdot l_i$ is the area of job i.

The performance metrics on which one can evaluate a parallel job scheduler falls into two categories usually with some overlap tolerated between the two:

User-Centric metrics[38]

These type of metrics refer to the actual job performance. The performance of the scheduler directly affects the waiting time of jobs. Below mentioned are some of the commonly used user-centric metrics. We will assume N as the number of jobs for which the metric is computed.

- **Average Waiting Time:** It refers to the average over the waiting time for all jobs.

$$AWT = \frac{1}{N} \cdot \sum_{i=1}^N t_i^w$$

- **Average Response Time:** It refers to the time when job results are available. It is computed from the waiting time plus the execution time (time interval from job submission to job completion).

$$ART = \frac{1}{N} \cdot \sum_{i=1}^N t_i^r$$

- **Average Response Time Weighted by Job Width:**

$$ART_{wW} = \frac{\sum_{i=1}^N w_i \cdot t_i^r}{\sum_{i=1}^N w_i}$$

- **Average Slowdown Weighted by Job Area:**

$$SLD_{wA} = \frac{\sum_{i=1}^N a_i \cdot s_i}{\sum_{i=1}^N a_i}$$

The last 2 metrics would not be suitable to use for adaptive jobs or malleable / evolving jobs since their width changes dynamically at runtime.

System-Centric metrics[38]

These type of metrics focus on the effective usage of the resources in the system. Below mentioned are some of the commonly used metrics:

- **Utilization** Percentage of all resources were actually used on average over a specific time frame
- **Throughput** Number of jobs that were processed during a specific time frame
- **Makespan** The time of the completion of the last job in the workload
- **Loss of Capacity** Percentage of resources that were idle, although workload for processing was available

7.1 Method of Evaluation

In order to select a scheduling algorithm for a system, it needs to be evaluated based on certain criterias, many of which have been defined earlier. The performance of such an algorithm depends a lot on the workload which it is processing, hence a workload benchmark maybe useful in evaluating a scheduler or even comparing it with other schedulers. There are various methods used to perform such an evaluation based on a workload[2]:

- **Deterministic Modeling** This evaluation takes a pre-determined workload and defines the performance of each algorithm for that workload. For example: 5 jobs arriving at time 0 in a particular order. One can then use a criteria such as average waiting time to find which algorithm performs the best. This method is simple and fast in giving exact results if the inputs are the same. But, the results only apply to these cases. The main usage of deterministic modeling can be to describe algorithms or serve as examples and over a set of inputs may indicate a trend or behavior that can help for further analysis.
- **Queueing Models** The drawback of deterministic modeling is that there is no static set of processes, this varies from day to day. In the method of queueing models, a model of the system is constructed by describing the arrival of jobs having some inter-arrival time, job sizes and job runtime estimates as probability distributions respectively. With all these details, one can then mathematically derive serveral metric values such as average waiting time, utilization, average queue length etc. This method of evaluation is useful in comparing algorithms but has its own limitations. It is only an approximation of the real system and only certain class of algorithms and distributions can be handled as of now in addition to a number of independent assumptions made that may not be accurate.
- **Simulations** A better and a more accurate approach to evaluating scheduling algorithms is a simulation. In such a method a model of the complete system is implemented using programming and the simulation is driven in time by moving the clock(variable representing time) forward. With time, the state of the system is modified by the simulator to reflect the activities taking place. As the simulation executes, statistics that indicate algorithm performance are gathered and printed. The data to drive the simulation can be generated in several ways: probability distributions for generating job arrivals, job sizes, runtime estimates etc or trace data generated by monitoring real systems and recording the actual sequence of events. Simulations can often be expensive requiring huge hours of computer time, effort to implement the simulator, a more detailed simulator will give a more accurate result but increases the effort.

- **Implementation** Simulation is of limited accuracy and hence the last approach requires the entire algorithm to be implemented and tested in a real system. This requires much more effort but also brings more uncertainty into the performance of the scheduler as the environment is constantly changing in with real jobs on real physical resources.

7.2 Setup

The scheduling algorithm implemented in this thesis will be tested by the approach of simulation. This means that the jobs are not actually launched and iRTSched instead will simulate this by changing the state of the system and modifying the necessary metadata to run the simulation in time. This also means that we do not need actual physical resources or a real cluster to do this testing and just need to emulate a cluster using SLURM to serve our purpose. For the purpose of doing this simulation, below are the steps taken:

- Configure the SLURM using `--enable-front-end` and install it. This allows us to emulate a cluster by this front end node which in practical systems will actually dispatch the jobs it receives to all the resources in a cluster managed by a different middleware.
- Run the SLURM by configuring the scheduler plugin as iBSched and the priority plugin as multifactor.
- Run iRTSched as an independent binary(it is not a SLURM application).
- iBSched and iRTSched will connect with each other and complete a handshake to prepare for future negotiations.
- All jobs submitted to SLURM with partition as `invasic` will be ignored by the main batch scheduler but picked up by the plugin iBSched. This is just a simple hack to redirect the stream of incoming jobs towards iRTSched and not towards the configured front end node.
- This has been done to ease the effort for the purpose of simulation. In reality `invasic` jobs must be directed towards `invasic` partition and due to some load balancing by the batch system, legacy(rigid) jobs can also be assigned to `invasic` partition when other partitions are busy.
- Simulation continues until we shutdown iRTSched and subsequently SLURM. Feedback data from iRTSched will be used to update job states in SLURM. Output

data corresponding to job start and end times will be written in a separate file by iRTSched for post-simulation analysis.

The experiments have been executed on a laptop with a 64-bit CPU(8 cores) each having a max. speed of approx. 3GHz. The 64-bit version of Ubuntu 14.04 LTS is running on the laptop and the version of SLURM used for this thesis is 15.08.

7.3 Experiments and Results

<TBA> <Mention about the assumption of jobs having linear speedup curves for time adjustments>

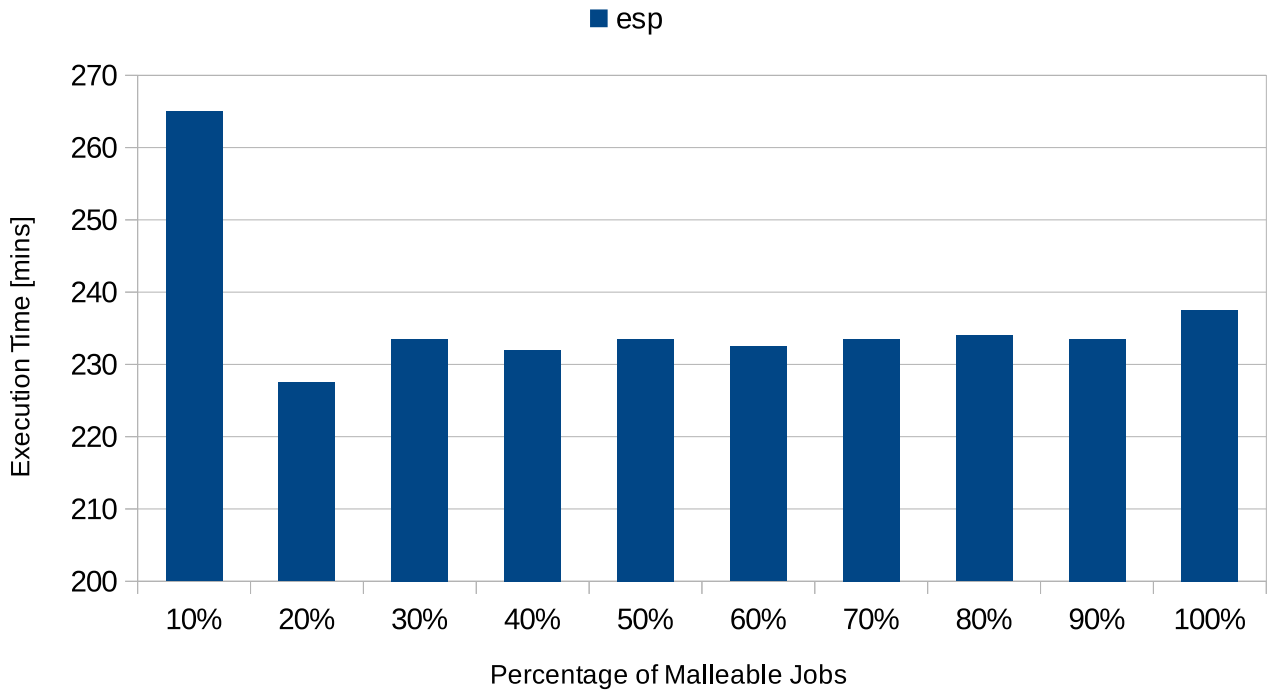


Figure 7.1: Time for completion of the modified ESP workload with varying amounts of rigid and malleable jobs

8 Conclusion and Future Work

<TBA>

8.1 Future Work

<Dummy citations of papers to test if bibliography section is working fine> [28] [29] [nikolas] [9] [46] [14] [5] [25] [47] [18] [pavan] [52] [15] [23] [joseph] [16] [4] [53] [24] [43] [54] [42] [11] [12] [dinesh] [33] [45] [38] [37] [achim1] [39] [achim2] [striet2] [8] [tiachao] [49] [alain] [51] [30] [viktor1] [javier] [34] [22] [ariel] [20] [cirne] [43] [michal] [36] [35] [song] [31] [19] [32] [41] [calvin] [40] [ribbens] [13] [27] [26] [gonzalo1] [maria] [srinidhi] [21] [zhiling] [48] [deshmeh2010adept] [6] [50] [daniel1] [7] [44] [3] [striet2] [10] [1] [17]

<Refer to the works for Prof. Bungartz and Emily on topics relating to AMR applications>

Bibliography

- [1] L. O. A. T. Wong, W. T. C. Kramer, T. L. Kaltz, and D. H. Bailey. "ESP: A System Utilization Benchmark." In: *ACM/IEEE Conference on Supercomputing* (Nov. 2000).
- [2] P. B. G. Abraham Silberschatz and G. Gagne. "Operating System Concepts Essentials." In: *John Wiley & Sons Inc.* (2011).
- [3] J. Buisson, O. Sonmez, H. Mohamed, W. Lammers, and D. Epema. "Scheduling Malleable Applications in MultiCluster Systems." In: *International Conference on Cluster Computing* (Sept. 2007).
- [4] Y. Cao, H. Sun, W.-J. Hsu, and D. Qian. "Malleable-Lab: A Tool for Evaluating Adaptive Online Schedulers on Malleable Jobs." In: *Euromicro Conference on Parallel, Distributed and Network-Based Processing(PDP)* (Feb. 2010).
- [5] M. C. Cera, Y. Georgiou, O. Richard, N. Maillard, and P. Navaux. "Supporting malleability in parallel architectures with dynamic cpusets mapping and dynamic MPI." In: *International Conference on Distributed Computing and Networking(ICDCN)* (Jan. 2010).
- [6] W. Cirne. "Using Moldability to Improve the Performance of Supercomputer Jobs." In: *Journal of Parallel and Distributed Computing* (Oct. 2002).
- [7] G. Da Costa, T. Fahringer, J.-A. Rico-Gallego, I. Grasso, A. Hristov, H. D. Karatza, A. Lastovetsky, F. Marozzo, D. Petcu, G. L. Stavrinides, et al. "Exascale Machines Require New Programming Paradigms and Runtimes." In: *Supercomputing Frontiers and Innovations* (Feb. 2015).
- [8] G. Deshmeh. "ADEPT Runtime/Scalability Predictor in Support of Adaptive Scheduling." PhD thesis. University of Windsor, Sept. 2013.
- [9] D.G.Feitelson and L.Rudolph. "Towards convergence in job schedulers for parallel supercomputers." In: *Workshop on Job Scheduling Strategies for Parallel Processing* (Apr. 1996).
- [10] J. G. Dong H. Ahn, M. Grondona, D. Lipari, B. Springmeyer, and M. Schulz. "Flux: A Next-Generation Resource Management Framework for Large HPC Centers." In: *International Conference on Parallel Processing Workshops* (Sept. 2014).

- [11] D. G. Feitelson and E. Shmueli. "Backfilling with Lookahead to Optimize the Packing of Parallel Jobs." In: *Journal of Parallel and Distributed Computing* (May 2005).
- [12] D. G. Feitelson and A. M. Weil. "Utilization and Predictability in Scheduling the IBM SP2 with Backfilling." In: *Parallel Processing Symposium and Symposium on Parallel and Distributed Processing(IPPS/SPDP)* (Apr. 1998).
- [13] C. George and S. S. Vadhiyar. "AdFT: An Adaptive Framework for Fault Tolerance on Large Scale Systems using Application Malleability." In: *International Conference on Computational Science (ICCS)* (May 2012).
- [14] M. Gerndt, A. Hollmann, M. Meyer, M. Schrieber, and J. Weidendorfer. "Invasive Computing with iOMP." In: *Forum on Specification and Design Languages(FDL)* (Feb. 2012).
- [15] A. Gupta, B. Acun, O. Sarood, and L. V. Kale. "Towards Realizing the Potential of Malleable Jobs." In: *High Performance Computing(HiPC)* (Dec. 2014).
- [16] J. Hungershofer. "On the Combined Scheduling of Malleable and Rigid Jobs." In: *International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)* (Oct. 2004).
- [17] J. H. Jürgen Teich, A. Herkersdorf, D. Schmitt-Landsiedel, W. Schröder-Preikschat, and G. Snelting. "Invasive Computing: An Overview." In: *Book Title: Multiprocessor System-on-Chip*, pp 241-268 (Nov. 2010).
- [18] C. Klein and C. Perez. "An RMS for Non-predictably Evolving Applications." In: *Cluster Computing(CLUSTER)* (Sept. 2001).
- [19] D. Klusacek and H. Rudova. "Alea 2 - Job Scheduling Simulator." In: *International ICST Conference on Simulation Tools and Techniques (SIMUTools)* (Mar. 2010).
- [20] R. Kübert. "Service Level Agreements for Job Submission and Scheduling in High Performance Computing." PhD thesis. Universität Stuttgart, Aug. 2014.
- [21] D. Kumar, Z.-y. Shae, and H. Jamjoom. "Scheduling Batch and Heterogeneous Jobs with Runtime Elasticity in a Parallel Processing Environment." In: *International Parallel and Distributed Processing Symposium Workshops and PhD Forum* (May 2012).
- [22] K. Kurowski, A. Oleksiak, and W. Piatek. "Hierarchical Scheduling Strategies for Parallel Tasks and Advance Reservations in Grids." In: *Journal of Scheduling* (Aug. 2013).
- [23] D. A. Lifka. "The ANL/IBM SP scheduling system." In: *Job Scheduling Strategies for Parallel Processing* (Apr. 1995).

- [24] A. Lucero. "Simulation of Batch Scheduling using Real Production Ready Software Tools." In: *Iberian Grid Infrastructure Conference* (June 2011).
- [25] K. E. Maghraoui, T. J. Desell, B. K. Szymanski, and C. A. Varela. "Dynamic Malleability in Iterative MPI Applications." In: *Concurrency and Computation: Practice and Experience* (Jan. 2008).
- [26] G. Martin, M.-C. Marinescu, D. E. Singh, and J. Carretero. "Enhancing the Performance of Malleable MPI Applications by Using Performance-Aware Dynamic Reconfiguration." In: *Journal of Parallel Computing* (Apr. 2015).
- [27] G. Martin, M.-C. Marinescu, D. E. Singh, and J. Carretero. "FLEX-MPI: An MPI Extension for Supporting Dynamic Load Balancing on Heterogeneous Non-dedicated Systems." In: *Euro-Par 2013 Parallel Processing* (Aug. 2013).
- [28] S. Prabhakaran, M. Iqbal, S. Rinke, C. Windisch, and F. Wolf. "A Batch System with Fair Scheduling for Evolving Applications." In: *International Conference on Parallel Processing(ICPP)* (Sept. 2014).
- [29] S. Prabhakaran, M. Neumann, S. Rinke, F. Wolf, A. Gupta, and L. V. Kale. "A Batch Sytem with Efficient Adaptive Scheduling for Malleable and Evolving Applications." In: *International Parallel and Distributed Processing Symposium(IPDPS)* (May 2015).
- [30] S. Rizos and V. Yarmolenko. "Job Scheduling on the Grid: Towards SLA based Scheduling." In: *High Performance Computing Workshop* (Apr. 2006).
- [31] G. Sabin, M. Lang, and P. Sadayappan. "Moldable Parallel Job Scheduling Using Job Efficiency: An Iterative Approach." In: *Job Scheduling Strategies for Parallel Processing (JSSPP)* (June 2006).
- [32] O. Sarood, A. Langer, A. Gupta, and L. Kale. "Maximizing Throughput of Over-provisioned HPC Data Centers Under a Strict Power Budget." In: *International Conference for High Performance Computing, Networking, Storage and Analysis (SC)* (Nov. 2014).
- [33] SchedMD. *SLURM Workload Manager*. www.slurm.schedmd.com.
- [34] K. M. Sim. "Grid Resource Negotiation: Survey and New Directions." In: *IEEE Transactions on Systems Man and Cyberbetics* (May 2010).
- [35] S. Srinivasan, V. Subramani, R. Kettimuthu, P. Holenarsipur, and P. Sadayappan. "Effective Selection of Partition Sizes for Moldable Scheduling of Parallel Jobs." In: *High Performance Computing(HiPC)* (Dec. 2002).
- [36] S. Srinivasan, S. Krishnamoorthy, and P. Sadayappan. "A Robust Scheduling Strategy for Moldable Scheduling of Parallel Jobs." In: *IEEE International Conference on Cluster Computing* (Dec. 2003).

- [37] A. Streit. "A Self-Tuning Job Scheduler Family with Dynamic Policy Switching." In: *Job Scheduling Strategies for Parallel Processing (JSSPP)* (July 2002).
- [38] A. Streit. "Self Tuning Job Scheduling Strategies for the Resource Management of HPC Systems and Computational Grids." PhD thesis. Universität Paderborn, Oct. 2003.
- [39] A. Streit. "The Self-Tuning dynP Job Scheduler." In: *Proceedings of the International Parallel and Distributed Symposium (IPDPS)* (Apr. 2001).
- [40] R. Sudarshan. "ReSHAPE: A Framework for Dynamic Resizing of Parallel Applications." PhD thesis. Virginia Polytechnic Institute and State University, Sept. 2009.
- [41] R. Sudarshan and C. J. Ribbens. "Design and performance of a scheduling framework for resizable parallel applications." In: *Journal of Parallel Computing* (Jan. 2010).
- [42] W. Tang, N. Desai, D. Buettner, and Z. Wan. "Analyzing and Adjusting User Runtime Estimates to Improve Job Scheduling on the Blue Gene/P." In: *International Symposium on Parallel and Distributed Processing Symposium (IPDPS)* (Apr. 2010).
- [43] W. Tang, Z. Lan, and N. Desai. "Job Scheduling with Adjusted Runtime Estimates on Production Supercomputers." In: *Journal of Parallel and Distributed Computing* (Mar. 2013).
- [44] W. Tang, D. Ren, Z. Lan, and N. Desai. "Adaptive Metric-Aware Job Scheduling for Production Supercomputers." In: *International Conference on Parallel Processing Workshops* (Sept. 2012).
- [45] D. Tsafir. "Modeling, Evaluating and Improving the Performance of Supercomputer Scheduling." PhD thesis. Hebrew University, Sept. 2006.
- [46] I. A. C. Urena, M. Riepen, M. Konow, and M. Gerndt. "Invasive MPI on Intel's single-chip cloud computer." In: *Architecture of Computing Systems (ARCS)* (Feb. 2012).
- [47] G. Utrera, J. Corbalan, and J. Albarta. "Implementing Malleability on MPI Jobs." In: *International Conference on Parallel Architecture and Compilation Techniques (PACT)* (Oct. 2004).
- [48] G. Utrera, S. Tabik, J. Corbalan, and J. Labarta. "A Job Scheduling Approach for Multi Core Clusters Based on Virtual Malleability." In: *Euro-Par 2012 Parallel Processing* (Aug. 2012).
- [49] O. Waeldrich, D. Battre, F. Brazier, K. Clark, M. Oey, A. Papaspyrou, P. Wieder, and W. Ziegler. "WS-Agreement Negotiation Version 1.0." In: *Open Grid Forum* (Jan. 2011).

- [50] J. B. Weissman, L. R. Abburi, and D. England. "Integrated scheduling: the best of both worlds." In: *Journal of Parallel and Distributed Computing* (June 2003).
- [51] V. Yarmolenko, R. Sakellariou, D. Ouelhadj, and J. M. Garibaldi. "SLA Based Job Scheduling: A Case Study on Policies for Negotiation with Resources." In: *Proceedings of e-Science All Hands Meeting (AHM2005)* (Sept. 2005).
- [52] A. B. Yoo, M. A. Jette, and M. Gondona. "SLURM: Simple Linux Utility for Resource Management." In: *Job Scheduling Strategies for Parallel Processing* (June 2003).
- [53] Z. Zhou, Z. Lan, W. Tang, and N. Desai. "Reducing Energy Costs for IBM Blue Gene/P via Power-Aware Job Scheduling." In: *Job Scheduling Strategies for Parallel Processing* (May 2013).
- [54] Z. Zhou, X. Yang, D. Zhao, P. Rich, W. Tang, J. Wang, and Z. Wan. "I/O Aware Batch Scheduling for Petascale Computing Systems." In: *International Conference on Cluster Computing* (Sept. 2015).