



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: M.Sc. Isaias Alberto Compres Urena
Submission Date: May 15, 2015



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

Munich, May 15, 2015

Nishanth Nagendra

Acknowledgments

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 executed with a fixed number of MPI processes before submitting to a HPC(High Performance Computing) system. This results in a fixed allocation of resources for the job. Newer techniques in scientific computing such as AMR(Adaptive Mesh Refinement) result in applications exhibiting complex behavior where their resource requirements change during execution. Invasive MPI which is a part of an ongoing research effort to provide MPI extensions for the development of Invasive MPI applications will result in evolving jobs for the HPC systems during runtime that utilize such AMR techniques. Unfortunately, using only static allocations result in the evolving applications being forced to execute using their maximum resource requirements that may lead to an inefficient resource utilisation. In order to support such parallel evolving applications at HPC centers there is an urgent need to investigate and implement extensions to existing resource management systems or develop an entirely new one. This thesis will extend the work done over the last few months during which an early prototype was implemented by developing a protocol for the integration of invasive resource management into existing standard batch systems. Specifically, This thesis will now investigate and implement a job scheduling algorithm in accordance with the new protocol developed earlier for supporting such an invasive resource management.

Contents

Acknowledgments	iv
Abstract	v
1 Introduction	1
1.1 Invasive Computing	3
1.2 Dynamic Resource Management	4
1.3 Master Thesis	7
1.3.1 Tentative Plan	7
1.4 Motivation	8
1.5 Document Structure	8
2 Related Work	10
2.1 Batch Scheduling	10
2.2 Runtime Scheduling	10
2.3 Adaptive Programming Paradigms	10
3 Invasive Computing	11
3.1 Job Classification	11
3.2 Resource Aware Programming	11
3.3 Traditional Resource Management	11
3.4 Support for Invasive Computing	12
3.4.1 Invasive MPI	12
3.4.2 Resource Management Extensions	12
4 Dynamic Resource Management Architecture	15
4.1 System Design	15
4.1.1 Batch Scheduler	19
4.1.2 Distributed Run Time Scheduler	19
4.1.3 MPI Process Manager	19
4.2 Negotiation Protocol	19
4.2.1 Protocol Sequence Diagrams	19

Contents

4.2.2	State Machine Diagrams	19
4.3	Invasive Jobs	26
5	iScheduler Design	27
5.1	Job Description	27
5.2	Resource Offers	27
5.3	Feedback Reports	27
5.4	Negotiation Protocol	27
5.5	Job Scheduling Algorithm	27
5.5.1	Problem Formulation	27
5.5.2	Pseudo Code	27
6	Implementation	28
6.1	Plugins for SLURM	28
6.2	Data Structures	28
6.3	State Machine Diagrams	28
6.3.1	iScheduler	28
6.3.2	iHypervisor	28
6.4	Important APIs	28
7	Evaluation	29
7.1	Method of Evaluation	29
7.1.1	Emulation of Workload	29
7.1.2	Real Invasive Applications	29
7.2	Setup	29
7.3	Experiments and Results	29
7.4	Performance and Graphs	29
8	Conclusion and Future Work	30
8.1	Future Work	30
Glossary		31
Acronyms		32
List of Figures		33
List of Tables		34
Bibliography		35

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 multi core 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 newer problems or 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. 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 Jobs/Applications 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 the Job the maximum number of resources before it starts that will result in sheer wastage of resources leading to a poor resource utilisation. In order to support such parallel 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 adaptive applications and the legacy static jobs 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

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, 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/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.

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 results 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 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 to allow for elasticity which means that the application can evolve at runtime.

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: *iScheduler* (Batch Scheduler for Invasive Jobs) built as an extension into an existing batch system and *iDRScheduler* (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 1.1 shows a high level illustration of the architecture for such an Invasive Resource Management.

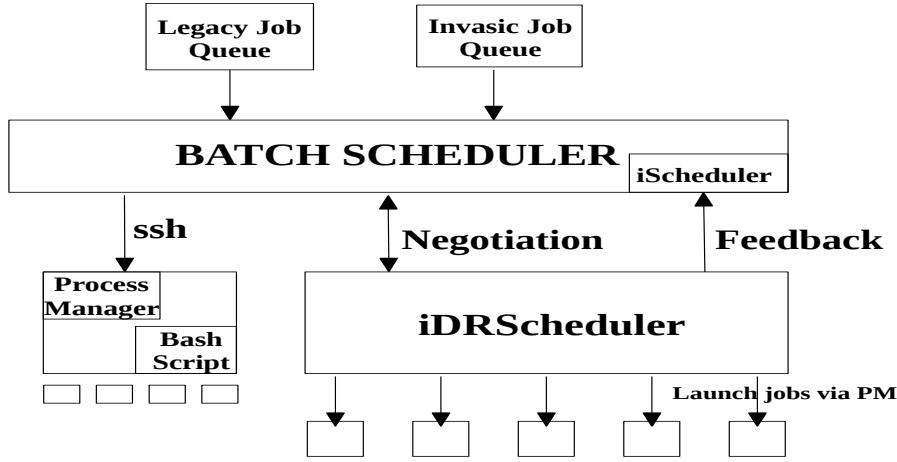


Figure 1.1: Invasive Resource Management Architecture

The above figure illustrates the proposed invasive resource management architecture. In addition to a job queue for legacy static jobs, we now have an additional job queue for invasive jobs. The existing batch scheduler needs to be extended in order to schedule these new type of Invasive Jobs and a new component called *Invasive Scheduler (iScheduler)* is responsible for this. In contrast to modifying an existing system to support Invasive Resource Management, a new component called *Invasive Distributed Run Time Scheduler (iDRScheduler)* is proposed which sits between the Batch Scheduler and Process Manager. 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 Invasive Jobs. It is responsible for managing the resources present in the Invasive partition used specifically for running Invasive Jobs. With this approach, the existing Legacy Jobs can be served via the existing batch scheduler and the new Invasive Jobs can be served by via iScheduler and iDRScheduler. iDRScheduler talks to the iScheduler via a protocol called the Negotiation protocol to receive Invasive Jobs dispatched from iScheduler which it then launches for execution by performing some run time scheduling like pinning of jobs, expand/shrink etc.

The new components proposed in the architecture help achieve the objective of supporting dynamic resource management for invasive applications. iScheduler is responsible for scheduling Invasive Jobs. The scheduling decisions are communicated via negotiation protocol to iDRScheduler and these decisions are basically job(s) selected via a scheduling algorithm to be submitted for execution. The decisions will be made on the basis of Resource Offers sent by iDRScheduler which creates these resource offers based on the state of the partition. Upon receiving a resource offer, iScheduler will either accept it by

selecting jobs from the queue that can be mapped to this offer or reject it. A Resource Offer can represent real or virtual resources because the iDRScheduler 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 local metrics such as resource utilization, power stability, energy efficiency etc. It can either accept or reject the mapping received from iScheduler. Similar to iDRScheduler, the iScheduler makes its decisions to optimize for certain local metrics such as high job throughput, reduced job waiting times, deadlines, priorities etc. This highlights the mismatching policies/metrics for which both the iDRScheduler and iScheduler make their decisions on and hence both will be involved in some kind of a negotiation via the protocol to reach a common agreement. iDRScheduler is an independent entity introduced with the purpose of inter-operating with existing batch systems rather than replacing them with an entirely new one. It may be possible that in the future this component will not be a separate entity but will be built into the batch system itself.

Shrink/Expand Strategies: There are several strategies one may employ to tackle adaptive applications during runtime and take decisions that can lead to higher system utilisation, energy efficiency etc.:

- Let us consider that at any given point of time there are some Invasive Jobs running in the system. If the parallel runtime environment is able to anticipate that in the near future, there may be a large window of free resources available because some applications may shrink according to a prediction of their scalability behavior with the help of run time profiling and collected performance data then iDRScheduler can provide a resource offer to iScheduler. This offer can specify a virtual list of nodes more than what is available in order to get a mapping of jobs. These jobs can then be shrunk and fit into the existing space of resources with the knowledge that they may be able to expand later.
- Another scenario is where there is an anticipation of a smaller window of resources in the future because some of the currently running applications may expand. In such a case iDRScheduler will provide a resource offer to iScheduler with a virtual list of nodes smaller than what is currently available in order to get a mapping of jobs. These jobs can then be expanded and fit into the existing space of resources with the knowledge that they may have to shrink later.
- The third variation could be the case where the runtime anticipates the state of resources to remain the same as the current state for the near future since none of the running applications may expand or shrink. In such a scenario iDRScheduler will send a resource offer that is a list of nodes which is exactly the

same as the available resources to iScheduler in order to get a mapping of jobs. These jobs can then be fit into the space of resources available as it is without expanding/shrinking them.

1.3 Master Thesis

The focus of this Master Thesis is to extend the early prototype developed as a part of the guided research in the last few months. It will now give concrete meanings to the negotiation protocol, defining the format of the invasic job records, its constraints, definition of resource offers, feedback reports and most important of all to investigate and develop an efficient job scheduling algorithm at the iScheduler level. Following this closely from the Guided Research would be to continue having an automated testing in place that will help in simulating a workload of jobs, testing the job scheduling algorithm for its correctness, evaluating and analysing the performance of such a prototype for various metrics. Given below is a tentative plan of the activities proposed monthwise starting from November 15 till May 15 to be carried out for the duration of 6 months in this Master Thesis.

1.3.1 Tentative Plan

- Literature survey for current/recent related work on batch job scheduling from research groups focusing on problems in the areas of batch scheduling, resource management, middleware etc. In addition to this, defining resource offers, invasic job records and job constraints that can come from a pre-computed performance model of an application during its runtime or could be user defined. Once these tasks have been completed, it will follow up with extending the early prototype by the implementation of these ideas and testing them for correctness using the automated testing feature.
- Integration of the fake iHypervisor with the real iHypervisor by making it as its plugin. Test this integration for all the earlier development using the automated testing feature for correctness. Review, analyse, fix any errors and repeat the process till the integration is stable. Follow this up by understanding the SLURM's existing scheduling algorithms including the recent Machine Learning approach for a possible reuse or enhancement. Define the format of feedback reports sent by iHypervisor and how they can be processed on the iScheduler side including its storage mechanism (possibly using SLURM's existing database functionalities). Implement these ideas in the prototype and test it thoroughly.

- Investigate and experiment with basic scheduling algorithms first to later proceed with complex ones. Design and implement a sophisticated scheduling algorithm for mapping jobs from Invasive job queue to the resource offers sent by iHypervisor. This algorithm must perform its decision making by also using the collected history/feedback data about many statistical measures provided by iHypervisor periodically.
- Implementation of the Job Scheduling Algorithm continues for realizing all the necessary requirements as proposed earlier till it has been successfully implemented and later followed by a lot of functional testing.
- Test the full prototype along with the other component iHypervisor (Including its Run Time scheduling and Invasive Resource Management) for simple to complex workloads (emulating a queue of invasive jobs possibly including legacy static jobs). Testing also needs to be done with real Invasive applications developed by the *Chair Of Scientific Computing*. Find issues/errors, review, analyse, correct and repeat the process till it is stable and then collect all the results necessary.
- Coming up with the draft version of the Master Thesis Report followed by reviews, corrections. This process is repeated till the final version is decided. Also prepare the slides for the Master thesis to present them at a later stage.

The above timeline highlights a tentative plan for the activities to be taken up during the Master Thesis and the 6 items above correspond to these 6 months of the Thesis in a chronological order. The steps may overlap or shift depending on the progress but the same overall structure will be followed for the Thesis. It will also include in parallel small amounts of documentation in the report as and when necessary during the 6 month period and not necessarily everything at the end.

1.4 Motivation

1.5 Document Structure

This is 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:

- **Modeling:** This section will explain in brief on how the DNDP will be mathematically formulated using a bilevel linear program with all its constraints, decision variables and objective function. It will explain the approach to approximate

the non-linear objective function of the inner problem of DNDP which is TAP and an introduction to metaheuristics that will be used to solve such kind of a combinatorial optimization problem.

- ***Genetic Algorithm:*** This section will dive into the details of the genetic algorithm that includes all the important aspects which fall under GA that needs to be tackled in order to come up with a correct implementation yielding good performance. It will explain in detail some of the many choices one has for implementation during every step of GA.
- ***Implementation:*** This section will illustrate with the help of a flow chart the high level view of the GA implementation followed by some pseudo codes to demonstrate in a simple language the implementation details of the inner workings of GA.
- ***Experiments and Visualization:*** This section will present all the results of the experiments conducted on different types(sizes) of datasets using the implemented GA in the form of tables. It will also mention in brief some of the observations and inferences that have been drawn by looking at these results.
- ***Conclusion:*** This section concludes the report on this project 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. This is followed by a list of some useful references that played an important role in the understanding of many of the concepts towards the realization of this project.

2 Related Work

2.1 Batch Scheduling

2.2 Runtime Scheduling

2.3 Adaptive Programming Paradigms

3 Invasive Computing

3.1 Job Classification

3.2 Resource Aware Programming

3.3 Traditional Resource Management

The role of a resource manager is to acts like a *glue* for a parallel computer to execute parallel jobs. It should make a parallel computer as easy to use as almost a PC. MPI would typically be used to manage communications within the parallel program. A resource manager allocates resources within a cluster, 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 Batch System serves as a middleware for managing supercomputing resources. The combination of *Scheduler+Resource Manager* makes it possible to run parallel jobs.

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.

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 with about 500,000 lines of C code whose development started in the year 2002 at Lawrence Livermore National Laboratory as a simple resource manager for Linux Clusters and a few years ago 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.

Plugins are dynamically linked objects loaded at run time based upon configuration file and/or user options. 3.1 shows where these plugins fit inside SLURM. 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.

SLURM Kernel				
Authentication Plugin	MPI Plugin	Checkpoint Plugin	Topology Plugin	Accounting storage Plugin
Munge	mvapich	BLCR	Tree	MySQL

Figure 3.1: SLURM with optional Plugins

3.4 Support for Invasive Computing

3.4.1 Invasive MPI

3.4.2 Resource Management Extensions

Existing batch systems usually only support static allocation of resources to applications before job start. Hence we need some kind of an Invasive Resource Management to be integrated into these existing batch systems so that we can support malleable jobs

allowing us to change the allocated resources dynamically at runtime. In order to achieve this we will follow the below approach:

- **Invasive Resource Manager (alias iHypervisor):** An independent component which will talk to the current batch systems via a communication mechanism(protocol) to obtain invasive job(s) submitted specifically to the invasic partition that can support invasive computing. The iHypervisor will then take these jobs and perform some kind of runtime scheduling for pinning these jobs to the resources in the partition and makes these decisions in order to optimize certain local metrics such as resource utilization, power stability, energy efficiency etc. The scheduling here is done at the granularity of cores and sockets. iHypervisor is the one that has the complete information of the resources in the invasic partition and also manages them. This component in an independent entity with the purpose of inter-operating with existing batch systems rather than replacing them with an entirely new one. It may be possible that in the future this component will not be a separate entitiy but will be built into the batch system itself.
- **Invasive Job Scheduler (alias iScheduler):** This component will be a new extension built into the existing batch systems for performing job scheduling. The scheduling decisions are communicated via the protocol used to speak to iHypervisor and these decisions are basically job(s) selected via a scheduling algorithm to be submitted to the iHypervisor for execution. The scheduling decisions will be made on the basis of available resources in the partition and it is the iHypervisor that communicates this to iScheduler in the form of resource offers (Real/Virtual). It can be a virtual resource offer because the iHypervisor can hide the real resources and present a rather fake view of them to iScheduler in the hope of getting a mapping of jobs to offer that is more suitable to satisfy its local metrics. Similar to iHypervisor, the iScheduler makes its decisions to optimize for certain local metrics such as high job throughput, reduced job waiting times, deadlines, priorities etc. This highlights the mismatching policies/metrics for which both the iHypervisor and iScheduler make their decisions on and hence both will be involved in some kind of a negotiation via the protocol to reach a common agreement.
- **Negotiation Protocol:** This protocol forms the core of the interaction between the iScheduler and iHypervisor. It allows for iHypervisor to make one or a set of resource offers to iScheduler which then needs to select jobs from its job queue to be mapped to these resource offers and finally sent back to the iHypervisor. The iHypervisor will then decide whether to accept/reject this mapping to satisfy its local metrics. If it accepts it will launch them based on some run time scheduling

and if it rejects then it informs this to iScheduler in addition to sending it a new resource offer. The iScheduler can also reject the resource offers in which case it will be sent a new one. On accepting an offer, the iScheduler then repeats its tasks to send back a mapping to iHypervisor and this interaction continues until both reach a common agreement. If the number of such attempts reach a threshold then iScheduler will just accept whatever offer it receives and iHypervisor will also accept whatever Map:Jobs→Offers it receives closing this transaction of negotiating. After this a new transaction will start.

- **SLURM:** SLURM is our choice of an existing batch system upon which this new implementation will be demonstrated as a proof of concept to support the new paradigm of resource-aware programming in the domain of invasive computing. In the near future, This can motivate further such supporting infrastructures with other batch systems using such Invasive Resource Management and Scheduling components.

This project implements a testing prototype for demonstrating how such an approach to support invasive computing with the above entities may work. It will involve implementing the communication infrastructure using SLURM API due to which iHypervisor and iScheduler will interact with each other using protocol messages filled with dummy values. It will also involve implementing the iScheduler as a new plugin(multithreaded) for SLURM and iHypervisor as a fake multithreaded Invasive Resource Manager daemon with bare minimum functionalities that in the near future will be an enhanced version of the daemon slurmd(built on top of it) found in SLURM. For the purposes of testing this prototype, the different scenarios that can be observed most of the time with the kind of negotiation protocol described earlier would be verified.

4 Dynamic Resource Management Architecture

4.1 System Design

This section illustrates and describes a high level design of the software implemented with the help of protocol sequence diagrams and state machine diagrams. It will help to understand at a high level as to how the system has been designed to support this new approach of invasive computing and how will many of its components in the software hierarchy interact with each other with new protocols or extensions of existing protocols to integrate such an invasive resource management into current batch systems.

The following page shows the software architecture of how Invasive Resource Management can be supported with a traditional resource manager and how exactly the new software components will fit in the existing software hierarchy. The 4.1 relates closely to how SLURM is organized since the intention of this work would be to demonstrate the support for Invasive Computing with the help of SLURM as a resource manager.

- 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(jobs submitted to invasive partition that supports invasive computing).
- 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) server. 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 "iHypervisor" will be implemented which needs to communicate with a new scheduling component iScheduler and influence the scheduling decisions taken by it. The iHypervisor 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. This is due to the reason that the scheduler for invasive jobs will work in a different manner based on the approach described earlier in comparison to the legacy job scheduler for static jobs. In case of SLURM which has a modular design with several optional plugins, a new plugin by name "iScheduler" will be implemented for SLURM to handle job scheduling specifically for invasive jobs.
- Communication between iHypervisor and iScheduler will involve the negotiation protocol as explained in the previous chapter but will also include periodic feedbacks being sent by iHypervisor to iScheduler having some useful statistical measures about current state of resources, resource utilization, job throughput etc. that may help influence the decision making of iScheduler. 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 (iScheduler and iHypervisor) for proper communication later on. 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 communication between iHypervisor and iScheduler 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 computing as discussed in the previous chapter. It is the phase during

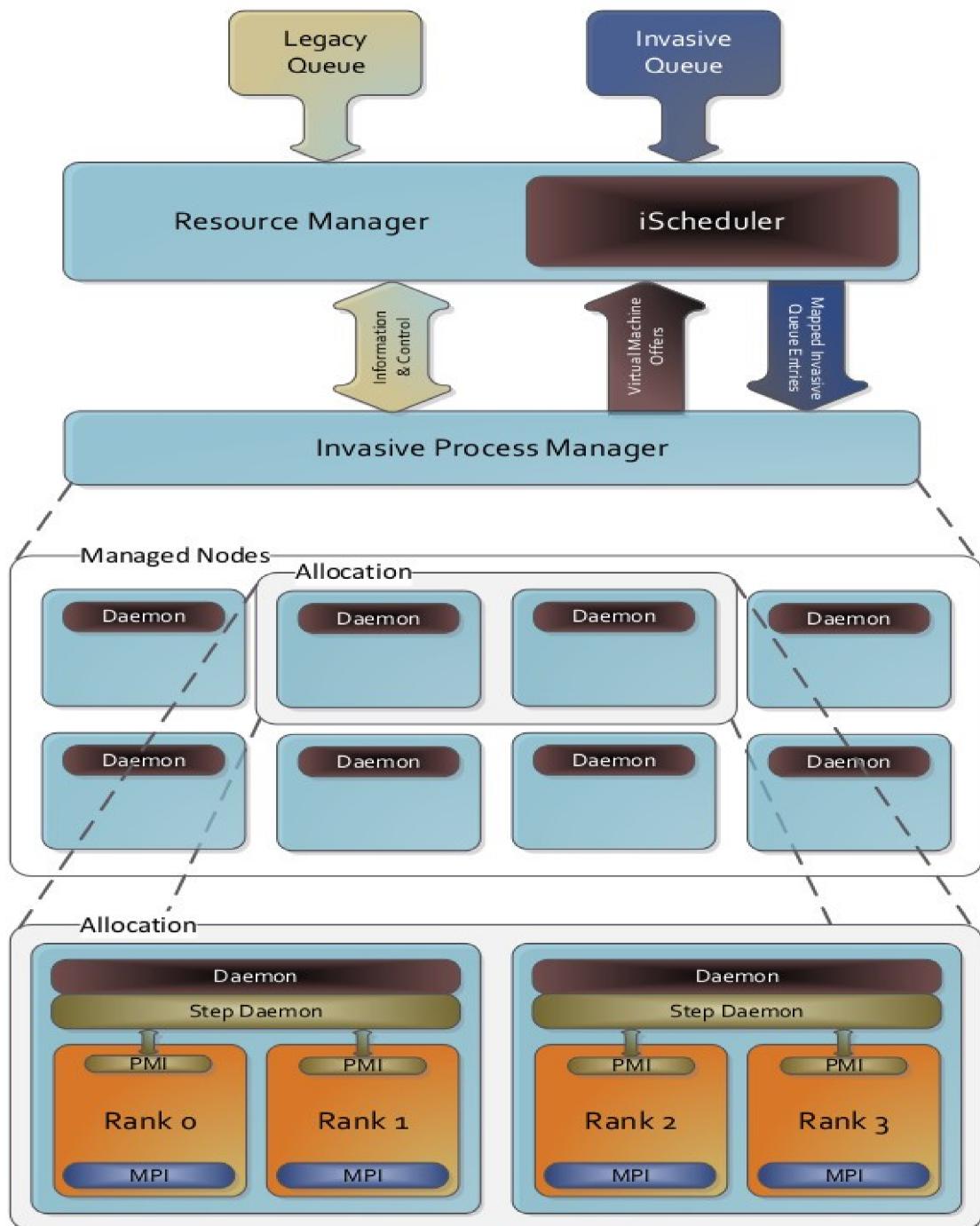


Figure 4.1: Invasive Resource Management Architecture

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

Figure 4.2: Message Types

which both iHypervisor and iScheduler are negotiating with each other till they reach an agreement. If they do not then they continue till a certain limit to the number of negotiating attempts are reached after which both of them just agree in their final attempt closing the current negotiation. After this a new transaction of negotiation begins.

- **Feedback:** This concerns the periodic feedback sent by the iHypervisor to the iScheduler containing useful information such as the job states, latest snapshot of the resources in the invasic partition and many other statistical measures not limited to system utilization, job throughput, waiting times of jobs to help and influence the iScheduler in its decision making for scheduling jobs during its future transactions of negotiation. 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 other jobs or suspending/Killing them.

4.1.1 Batch Scheduler

4.1.2 Distributed Run Time Scheduler

4.1.3 MPI Process Manager

4.2 Negotiation Protocol

4.2.1 Protocol Sequence Diagrams

This section focuses on iScheduler and a thread iRM_AGENT that it spawns which is the one responsible for all the communication with the iHypervisor including spawning other agent threads for handling feedbacks and urgent jobs.

- Above diagram and the ones in the following pages illustrate state machine diagrams for few of the communication phases described earlier starting first with a general diagram of how the multithreaded component iRM agent inside iScheduler starts up and shuts down.

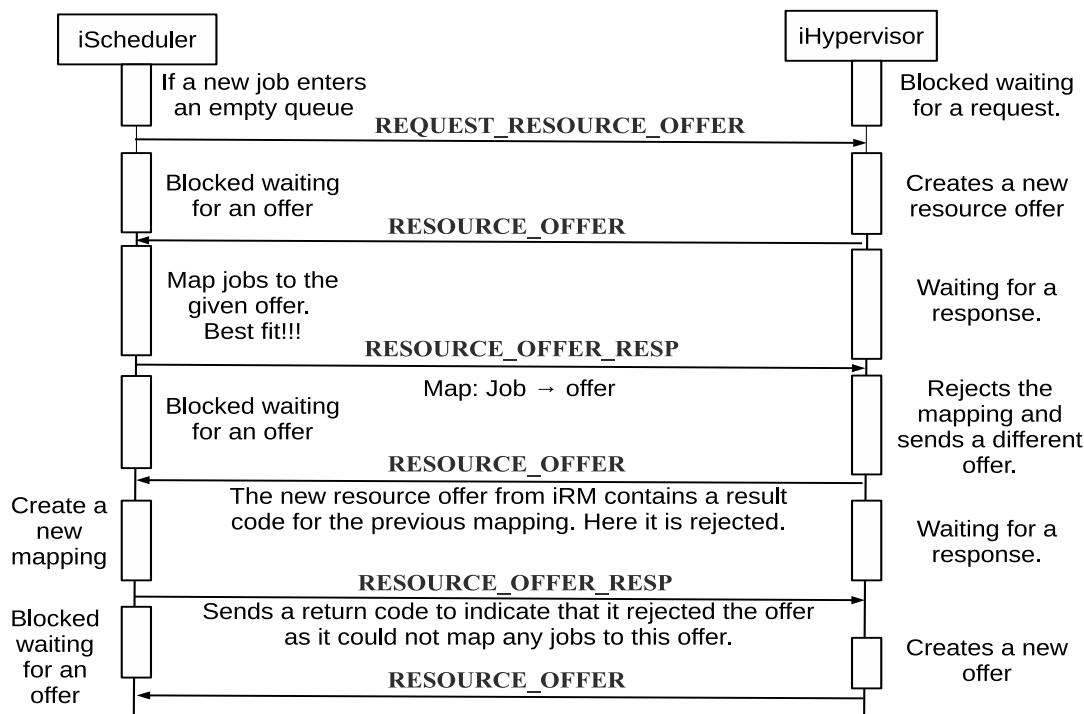


Figure 4.3: Scenario 1

- Above diagram illustrates a scenario where both iScheduler and iHypervisor are negotiating with each other. The scenario is continued in the next page. 4.5 illustrates another scenario where negotiations may stop when job queue becomes empty and iHypervisor then will wait for a request from iScheduler for a resource offer that will happen when new jobs arrive.
 - iScheduler makes scheduling decisions at a coarser level of granularity which is nodes whereas iHypervisor does at the granularity of cores and sockets. Both will negotiate with each other till they reach an agreement.
 - It is an event based scheduling which means iScheduler makes a scheduling decision only when it is triggered by receiving a resource offer from iHypervisor. It is only at the start when there are no jobs in the queue and during the operations when the queue may become empty that the iScheduler will have to explicitly send a request message to iHypervisor for a resource offer otherwise at all other times scheduling is event based.

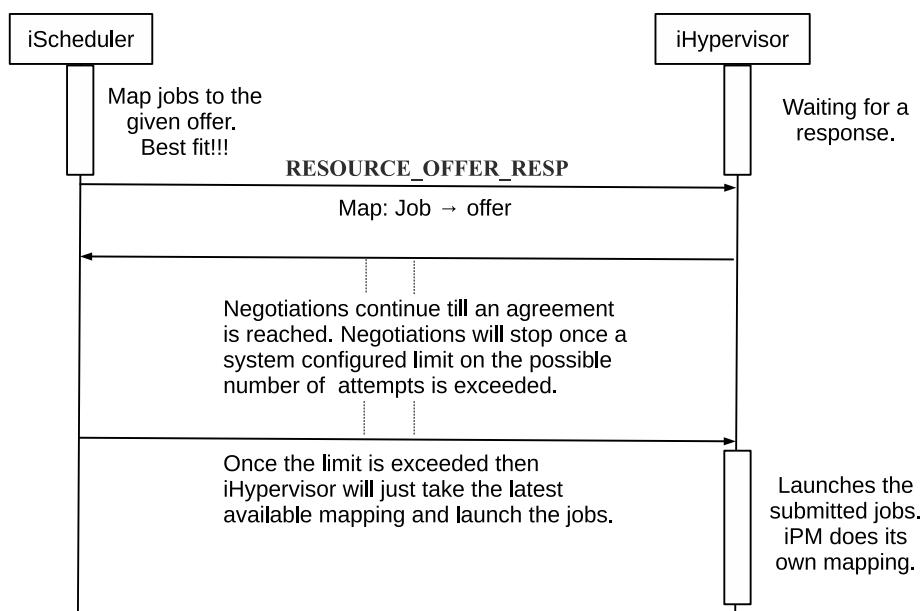


Figure 4.4: Scenario 1 contd.

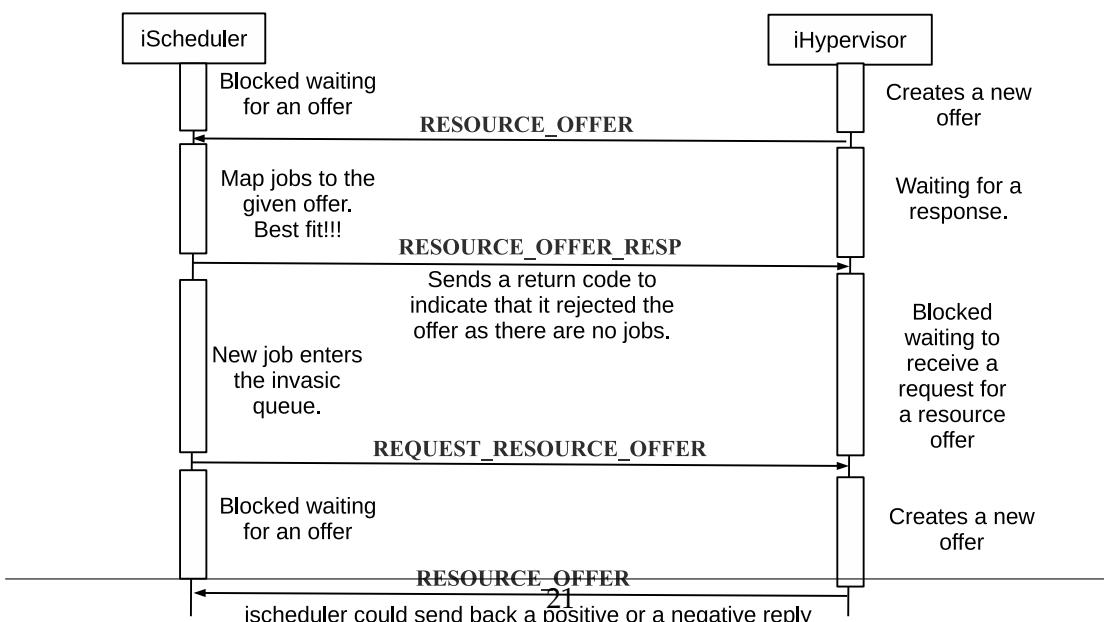


Figure 4.5: Scenario 2

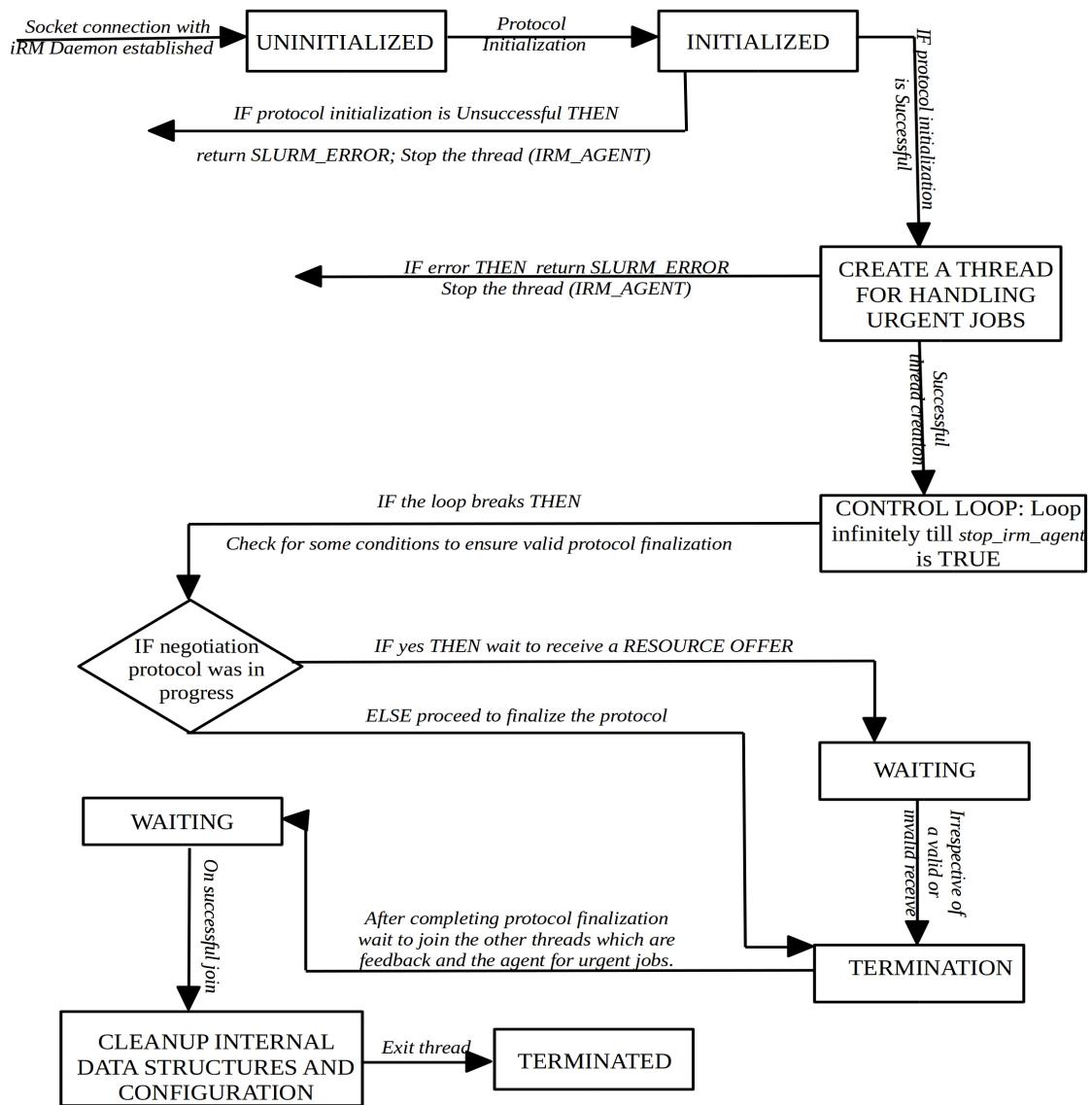


Figure 4.6: iRM Agent

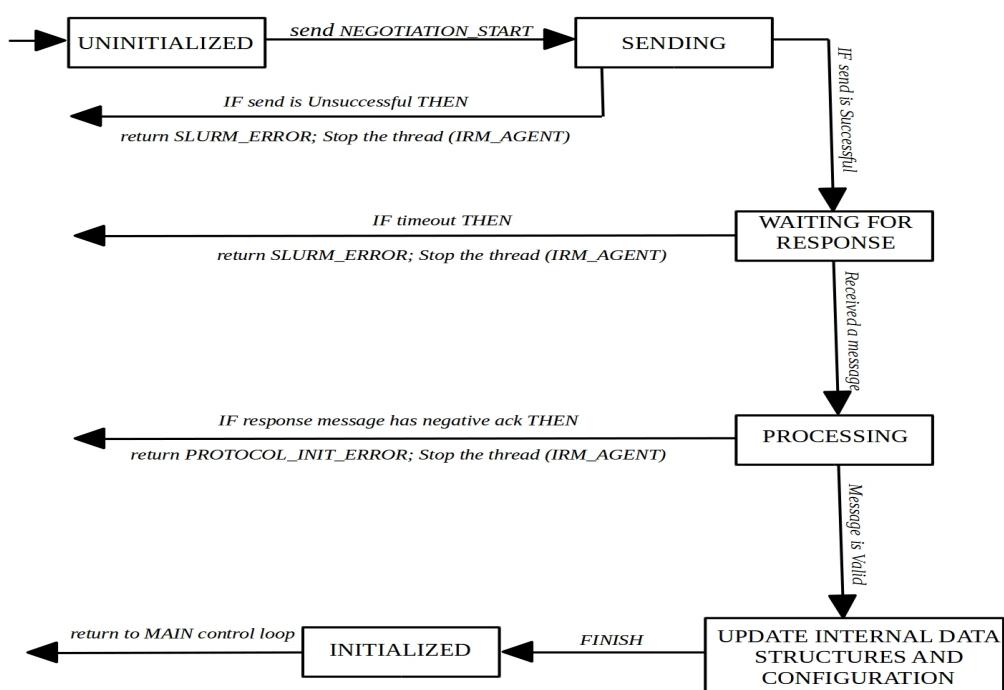


Figure 4.7: Protocol Initialization

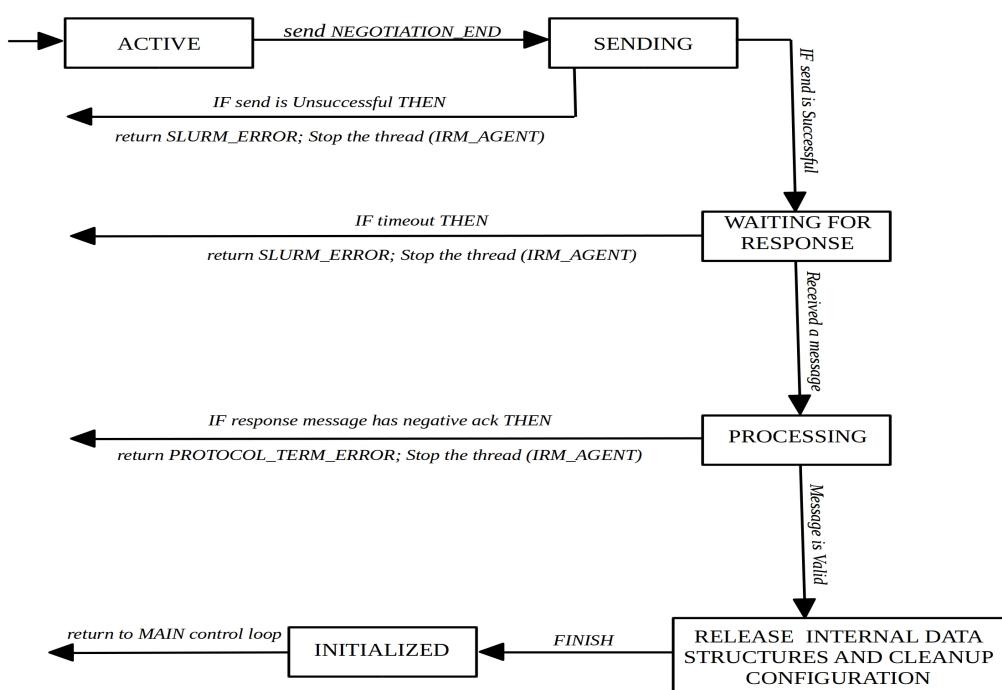


Figure 4.8: Protocol Termination

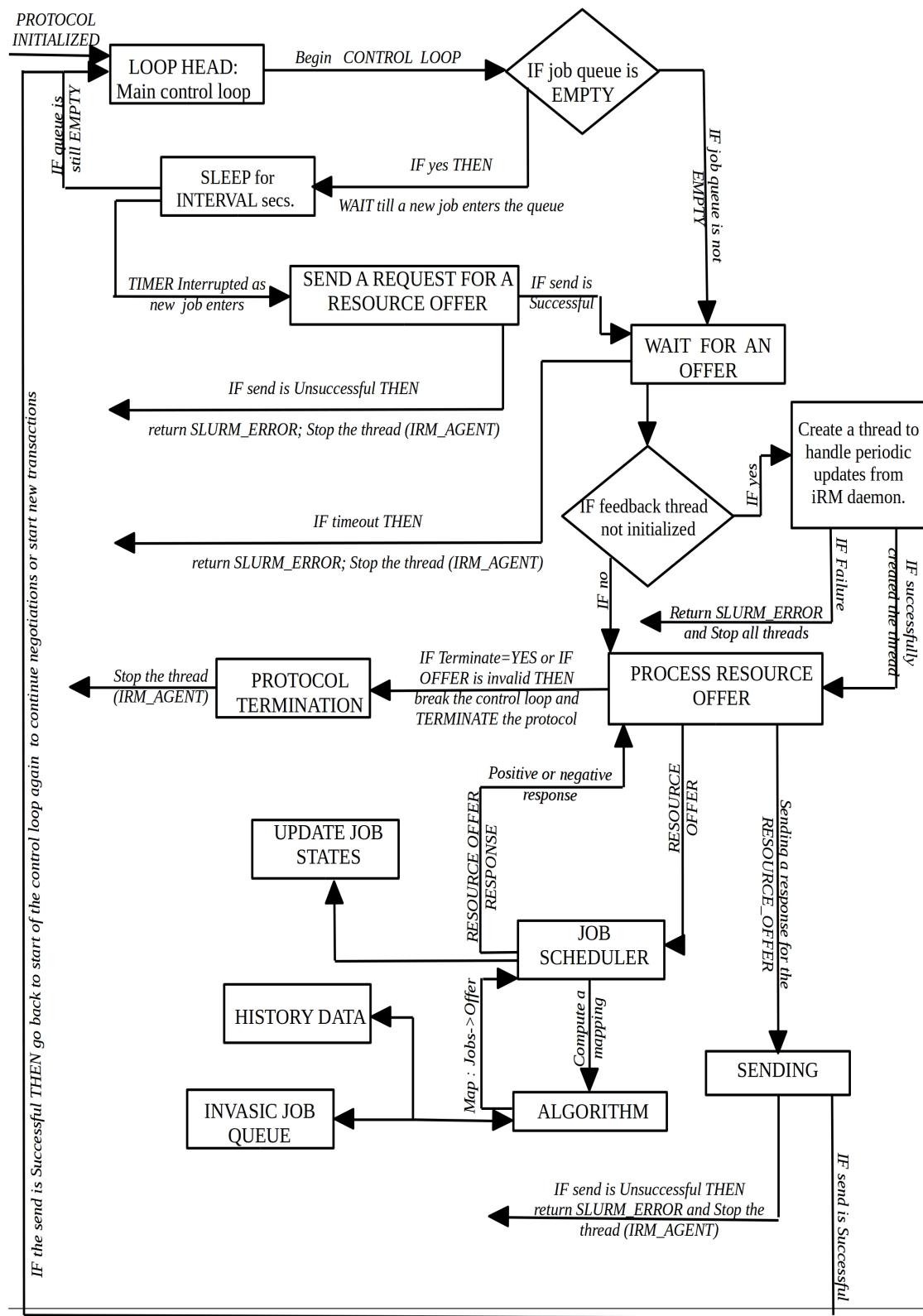


Figure 4.9: Negotiation

4.3 Invasive Jobs

5 iScheduler Design

5.1 Job Description

5.2 Resource Offers

5.3 Feedback Reports

5.4 Negotiation Protocol

5.5 Job Scheduling Algorithm

5.5.1 Problem Formulation

5.5.2 Pseudo Code

6 Implementation

6.1 Plugins for SLURM

6.2 Data Structures

6.3 State Machine Diagrams

6.3.1 iScheduler

6.3.2 iHypervisor

6.4 Important APIs

7 Evaluation

7.1 Method of Evaluation

7.1.1 Emulation of Workload

7.1.2 Real Invasive Applications

7.2 Setup

7.3 Experiments and Results

7.4 Performance and Graphs

8 Conclusion and Future Work

8.1 Future Work

[Pra+14] [Pra+15] [Ioa+11] [DL96] [Ure+12] [Ger+12] [Cer+10] [Mag+08] [UCA04] [KP01] [Bal+10] [YJG03] [Gup+14] [Lif95] [Sko+96] [Hun04] [Cao+10] [Zho+13] [Luc11] [TLD13] [Zho+15] [Tan+10] [Tan+05] [FW98] [FW12] [Sch]

Glossary

computer is a machine that....

Acronyms

TUM Technische Universität München.

List of Figures

1.1	Invasive Resource Management Architecture	5
3.1	SLURM with optional Plugins	12
4.1	Invasive Resource Management Architecture	17
4.2	Message Types	18
4.3	Scenario 1	20
4.4	Scenario 1 contd.	21
4.5	Scenario 2	21
4.6	iRM Agent	22
4.7	Protocol Initialization	23
4.8	Protocol Termination	24
4.9	Negotiation	25

List of Tables

Bibliography

- [Bal+10] P. Balaji, D. Buntinas, D. Goodwell, W. Gropp, J. Krishna, E. Lusk, and R. Thakur. “PMI: A Scalable Parallel Process-Management Interface for Extreme-Scale Systems.” In: *Recent Advances in the Message Passing Interface* (Sept. 2010).
- [Cao+10] 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).
- [Cer+10] 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).
- [DL96] 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).
- [FW12] D. G. Feitelson and A. M. Weil. “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).
- [FW98] 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).
- [Ger+12] 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).
- [Gup+14] A. Gupta, B. Acun, O. Sarood, and L. V. Kale. “Towards Realizing the Potential of Malleable Jobs.” In: *High Performance Computing(HiPC)* (Dec. 2014).
- [Hun04] 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).

Bibliography

- [Ioa+11] N. Ioannou, M. Kauschke, M. Gries, and M. Cintra. "Phase-Based Application-Driven Hierarchical Power Management on the Single-chip cloud Compupter." In: *Parallel Architectures and Compilation Techniques(PACT)* (Oct. 2011).
- [KP01] C. Klein and C. Perez. "An RMS for Non-predictably Evolving Applications." In: *Cluster Computing(CLUSTER)* (Sept. 2001).
- [Lif95] D. A. Lifka. "The ANL/IBM SP scheduling system." In: *Job Scheduling Strategies for Parallel Processing* (Apr. 1995).
- [Luc11] A. Lucero. "Simulation of Batch Scheduling using Real Production Ready Software Tools." In: *Iberian Grid Infrastructure Conference* (June 2011).
- [Mag+08] 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).
- [Pra+14] 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).
- [Pra+15] S. Prabhakaran, M. Neumann, S. Rinke, F. Wolf, A. Gupta, and L. V. Kale. "A Batch System with Efficient Adaptive Scheduling for Malleable and Evolving Applications." In: *International Parallel and Distributed Processing Symposium(IPDPS)* (May 2015).
- [Sch] SchedMD. *SLURM Workload Manager*. www.slurm.schedmd.com.
- [Sko+96] J. Skovira, W. Chan, H. Zhou, and D. Lifka. "The EASY - LoadLeveler API Project." In: *Job Scheduling Strategies for Parallel Processing* (Apr. 1996).
- [Tan+05] W. Tang, N. Desai, D. Buettner, and Z. Wan. "Backfilling with Lookahead to Optimize the Packing of Parallel Jobs." In: *Journal of Parallel and Distributed Computing* (May 2005).
- [Tan+10] 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).
- [TLD13] 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).
- [UCA04] G. Utrera, J. Corbalan, and J. Albarta. "Implementing Malleability on MPI Jobs." In: *International Conference on Parallel Architecture and Compilation Techniques(PACT)* (Oct. 2004).

Bibliography

- [Ure+12] 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).
- [YJG03] 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).
- [Zho+13] 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).
- [Zho+15] 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).