

DATA REPLICATION FOR COPING WITH UNCERTAINTY IN SCHEDULING

by

Manmohan Chaubey

A thesis submitted to the faculty of
The University of North Carolina at Charlotte
in partial fulfillment of the requirements
for the degree of Master of Science in
Computer Science

Charlotte

2014

Approved by:

Dr. Erik Saule

Dr. Yu Wang

Dr. Aidong Lu

©2014
Manmohan Chaubey
ALL RIGHTS RESERVED

ABSTRACT

MANMOHAN CHAUBEY. Data Replication for coping with Uncertainty in Scheduling. (Under the direction of DR. ERIK SAULE)

Scheduling theory is a common tool to analyze the performance of parallel and distributed computing systems, such as their load balance. How to distribute the input data to be able to execute a set of tasks in a minimum amount of time can be modeled as a scheduling problem. Often these models assume that the computation time required for each task is known accurately. However in many practical case, only approximate values are available at the time of scheduling.

This thesis research investigates how replicating the data required by the tasks can help coping with the inaccuracies of the processing times. In particular, it investigates the problem of scheduling independent tasks to optimize the makespan on a parallel system where the processing times of tasks are only known up to a multiplicative factor. The problem is decomposed in two phases: a first offline phase where the data of the tasks are placed and a second online phase where the tasks are actually scheduled.

For this problem, this thesis investigates three different strategies, each allowing a different degree of replication of jobs: a) No Replication b) Replication everywhere and c) Replication in groups, and proposes approximation algorithms and theoretical lower bound on achievable approximation ratios. This allows us to study the tradeoff between the number of replication and the guarantee on the makespan. Replication improves performance but incurs a cost in terms of memory consumption. The objec-

tive is then to develop scheduling algorithm with good competitive ratio to minimize both the makespan of the schedule and the memory consumption of the machines.

ACKNOWLEDGMENTS

First and foremost, I offer my sincerest gratitude to my advisor, Dr. Eric Saule, for his guidance, patience, and constant support. It would not have been possible to pursue my research interests without his help and encouragement. While I was impressed with his knowledge from the start my admiration for him has only increased with time. I have learned many things while pursuing this research under his guidance. It has been a pleasure and an honor working with him.

My special thanks to supervisory committee members, Dr. Yu Wang and Dr. Aidong Lu for reviewing the thesis and valuable feedback.

Finally, I offer heartfelt thanks to my parents for their affection and support, and always believing in me and encouraging me to be the best I can be.

TABLE OF CONTENTS

LIST OF FIGURES	viii
LIST OF TABLES	ix
CHAPTER 1: INTRODUCTION	1
1.1 Motivation	1
1.2 Scheduling Preliminaries	2
1.3 Research Contribution	4
1.4 Thesis Outline	5
CHAPTER 2: PROBLEM DEFINITION	6
CHAPTER 3: RELATED WORK	9
3.1 Classical Scheduling Problem	9
3.2 Uncertainty and Robustness	10
3.3 Data Placement and Replication	12
CHAPTER 4: REPLICATED DATA PLACEMENT STRATEGIES	14
4.1 Strategy 1: No Replication	14
4.2 Strategy 2: Replicate Data Everywhere	19
4.3 Strategy 3: Replication in Groups	21
4.4 Summary	26
CHAPTER 5: MEMORY AWARE REPLICATION UNDER UNCERTAINTY	29
5.1 Preliminaries	29
5.2 The $SABO_{\Delta}$ Algorithm	30

	vii
5.3 The ABO_{Δ} Algorithm	32
5.4 Summary	36
CHAPTER 6: CONCLUSION AND FUTURE WORK	39
REFERENCES	41

LIST OF FIGURES

FIGURE 1: Instance constructed by the adversary in the proof of theorem 5 with $\lambda = 3$ and $m = 6$. In the online solution, the adversary increases the processing time of a task of the most loaded machine by a factor of α . If that information was available beforehand, an optimal offline algorithm could have distributed these longer tasks to other processors.	16
FIGURE 2: An example of replication in groups with $m = 6$, $k = 2$. In phase 1, the data of the tasks are assigned to one of the groups. Phase 2 schedules each task assigned to a machine within its group.	22
FIGURE 3: Ratio-Replication graph with $m = 210$ and $\alpha \in \{1.1, 1.5, 2\}$.	28
FIGURE 4: An example of two phases of the schedule generated by the $SABO_{\Delta}$. The uncolored parts represent tasks scheduled according π_2 . The colored parts represents tasks scheduled according π_1	30
FIGURE 5: An example of the schedule generated by the ABO_{Δ} algorithm. The uncolored parts represent the memory intensive tasks scheduled according π_2 . The colored parts represent the processing time intensive tasks and scheduled using LS after replicated	33
FIGURE 6: Memory-Makespan graph for $SABO_{\Delta}$ and ABO_{Δ} . The bold lines represent impossibilities in tradeoff between guarantees.	38

LIST OF TABLES

TABLE 1: Summary of the contribution of this chapter. Three proposed algorithms have guaranteed performance. One lower bound on approximability has been established.	26
TABLE 2: Summary of the results of the algorithm $SABO_{\Delta}$ and the algorithm ABO_{Δ} .	37

CHAPTER 1: INTRODUCTION

This chapter provides the foundations for the principle objective of this dissertation, which is to investigate the effect of task replication on scheduling under uncertainty of processing times. The subject matter of this thesis falls in the intersection of several areas of current research interest. These includes: (1) Scheduling under uncertainty, (2) Data placement and Replication strategies to improve performance of a schedule, and (3) Bi-objective optimization for simultaneously optimizing makespan as well as memory consumption

1.1 Motivation

In real world scheduling problems often the parameters such as processing time of a task is not known exactly in advance. The goal of an scheduling algorithm is to generate robust schedule against uncertainty. Dealing with uncertainty is difficult as in real world problems a task can be processed on particular computing systems, otherwise a task could be moved as system sees it fit without incurring extra cost and the problem would be vastly alleviated. But in practice a task has to run on particular machine especially in ‘out of core’ computing applications which involve very large data sets. For example, solving systems of linear equations and computing eigenvalues – where matrices involved are very large. When the data sets are too large to fit in the main memory of a computer, it must be stored on any external memory source such

as disks. Disk storage is significantly cheaper than main memory storage. However, accessing data from a disk is relatively slower than accessing the main memory. So, a scheduling algorithm in ‘out of core’ execution places data to main memory of different systems such that data access from disk or any external storage is reduced. That means a task is pinned to a particular computing unit. So handling uncertainty in out of core execution is having added overhead. Hence, developing a scheduling strategy which can guarantee performance under uncertainty of processing times of the tasks with restriction that a task can be scheduled to particular set of machines motivates this research.

Scheduling tasks on distributed memory is particularly prevalent in Hadoop. Hadoop-MapReduce constitute a powerful Computation Model for processing large data sets on distributed clusters [21]. Hadoop stores large amount of data across multiple machines and processes them using MapReduce. Uncertainty in Hadoop system is related to a node failure or tasks failure. To cope with these uncertainties Hadoop uses data replication across multiple nodes. One of the main goal of a Hadoop system is to maintain node locality which means running data on the node that contains it [32]. Therefore, a data intensive scheduling incorporating data location and choosing popular data sets to replicate would be beneficial [13]; and serves to provide motivation for this research.

1.2 Scheduling Preliminaries

Parallel and distributed computing systems are often modeled using tasks that are processed simultaneously on different machines. Studying the balance of the load of

the various component of the system is often key in understanding the performance one obtains in practice. A system typically schedules the set of tasks with the goal of optimizing the load balance (or makespan) of the system or some other metric. A key information these system use to plan the execution is the time tasks will take to be processed. However, this information is typically not precisely known in practice: because the user can only make a wild guess on the runtime of her task [18], because prediction is hard in the general case [29], or because underlying models of a particular algorithm can only predict runtime within a given range [8]. Whichever the reason is, not knowing accurately the processing time can significantly impact the performance obtained from the machine.

For instance in out-of-core sparse linear algebra, executing a task where the data are not locally available would have a prohibitive overhead [34, 33].

One approach for dealing with the uncertainty of processing time is to build a robust schedule [2, 10, 6], that is, building a schedule that can naturally cope with variations in the processing times. These techniques often use sensitivity analysis to determine the robustness of the schedule. However, a better approach would be to be able to dynamically change the schedule.

The thesis pursues the idea of replicating the input data of the tasks onto multiple machines. This way, when the actual processing times of the tasks are too different from their estimations, the system will have some room to adapt at runtime. This is certainly feasible in practice as many system have more memory than the computation use. For instance, most Hadoop system replicates the data for the purpose of tolerating hardware faults [27]. And it has been shown that launching the same task

multiple times can help cope with hardware differences [26] but increases resource usage. The cost of replicating the data might be amortized in many applications where the application will iterated over the data multiple times (*e.g.*, in an iterative solver [33, 34]). This research answer the question “can data replication help cope with the uncertainty of processing time?” And the answer is that it can.

1.3 Research Contribution

This thesis proposes strategies and presents algorithms to cope with uncertainty in processing times of the tasks. The research provides three replication strategies and studies the tradeoff between the number of replication and the guarantee on the makespan. The strategy *No Replication* investigates what can be done if the tasks can only be deployed on a single machine, we provide a guaranteed algorithm and provide a lower bound on the best guarantee that one can achieve in this case. The strategy *Replicate data everywhere* takes the reverse case and investigates what can be achieved if the data are replicated everywhere, leaving the maximum flexibility at runtime. We investigate one algorithm in this case and analyze its performance guarantee. The strategy *Replication in groups* investigates grouping processors together and replicating data in these groups as an intermediate between the previous two strategies and provide a guaranteed algorithm in that case.

To alleviate the cost of replication in terms of memory consumption the thesis presents two memory-aware algorithms to optimize the makespan as well as the memory consumption. The proposed algorithms divides the tasks into two sets: memory intensive tasks and processing time intensive tasks and schedule differently to mini-

mize both the objectives.

1.4 Thesis Outline

The remaining of this thesis is organized as follows: we describe system model and notations in Chapter 2. Related works are presented in Chapter 3. Chapter 4 investigates the effect of replication on processing time uncertainty through three strategies which offer different degree of replication of the tasks. The chapter summarizes the various results derived for each strategy and studies the tradeoff between performance guarantee and data replication. Chapter 5 investigates bi-objective problem of minimizing the makespan as well as the memory usage and proposes two memory-aware algorithms which simultaneously optimizes both the objectives. Chapter 6 concludes the thesis with remarks and raises few challenging questions which could be future research topics.

CHAPTER 2: PROBLEM DEFINITION

Let J be a set of n jobs which need to be scheduled onto a set M of m machines. We will use interchangeably the terms machines and processors. Also we will use interchangeably the terms jobs and tasks. Each task j occupies s_j space in memory. We are considering the problem where the scheduler does not know the processing time p_j of task j exactly before the task completes. But the scheduler has access to some estimation of the processing time \tilde{p}_j of task j before making any scheduling decisions. We assume that the actual processing time p_j of a task j is within a multiplicative factor α of the estimated processing time \tilde{p}_j . α is a quantity known to the scheduler. In other words the scheduler knows that:

$$\frac{\tilde{p}_j}{\alpha} \leq p_j \leq \alpha \tilde{p}_j \tag{1}$$

Assuming that the processing time of the tasks is known to be in an interval is reasonable in many application scenarios. One could derive bounds experimentally using machine learning techniques: for instance [31] used Support Vector Machines to predict the time it will take to run graph traversal algorithms. Models of runtime of algorithms can also be derived analytically: in [8] the authors provide bounds for the performance of various sparse linear algebra operations using only the size of the matrices and vectors involved.

The scheduling for the problem is performed in two phases. Phase 1 chooses where

data are replicated using the estimated processing time \tilde{p}_j , for each of the task j . The phase takes \tilde{p}_j , m and α as inputs and outputs sets of machines, $M_j \subseteq M$ where each task j can be scheduled. This phase is purely offline and corresponds to the operations performed to prepare the execution of the application.

Phase 2 takes the output of phase 1 as its input and maps each task j to a machine within the set of machines M_j . For each machine i , let $E_i \subseteq J$ be the set of tasks assigned to machine i . This phase chooses the actual schedule following an online semi-clairvoyant process. Only the approximate processing time is known when a task is placed, but the scheduler can wait for a machine to become idle, to place the next one. Therefore, can dynamically schedule the tasks and the actual processing time of the tasks are known once they complete.

The parallel system scheduling can be modeled into different objective functions with different parameters to optimize. A makespan minimization problem has objective to minimize completion time of last task of the system. Memory is another parameter for objective function. A memory aware scheduling aims at minimizing total memory consumption $\sum_j s_j$ or memory consumption of most occupied machine $\max_i \sum_{j \in E_i} s_j$. Replication improves processing of tasks but increases memory consumption in the system. So, objective function attached with replication can be where to replicate tasks and which tasks to replicate so that performance can improve without violating any memory constraint or with bi-objective to minimize memory also along with improving processing time of the tasks.

In Chapter 4, the problem is to optimize the makespan, $C_{max} = \max_i \sum_{j \in E_i} p_j$ which is the completion time of the last task of the system. C_{max}^* denotes the optimal

makespan of an instance of the problem (knowing the actual processing times). The memory objective is constrained by allowing different degree of replication by choosing where (on which set of machines) a task to be replicated. An offline algorithm is said to be a ρ -approximation algorithm (or to have an approximation ratio of ρ) if it guarantees for all the instances that $C_{max} \leq \rho C_{max}^*$. When the problem is online, we are talking about competitive ratios.

In Chapter 5, we tackle the bi-objective problem of simultaneously minimizing makespan C_{max} as well as memory usage, $M_{max} = \max_i \sum_{j \in E_i} s_j$ which denotes the maximum memory usage of a machine. As a task occupies fixed amount of memory but its processing time is uncertain, both objectives are asymmetrical. M_{max}^* denotes optimal maximum memory consumption of a machine. An algorithm generates a schedule which is ρ^C -approximated on makespan and ρ^M -approximated on memory.

There are two ways to deal with multi objective optimization [23] [7]:

1. Epsilon-constraint method: This approach optimizes the primary objective setting the other objective within some constraint . We use this approach in chapter 4 to optimize makespan setting the memory objective by allowing different degree of replication of the tasks.
2. Zenith approximation: This approach optimizes both the objective at the same time. We use this approach to optimize both makespan and memory usage in chapter 5.

CHAPTER 3: RELATED WORK

This chapter provides the literature review on related research areas such as uncertainty in scheduling, data placement and replication. For better understanding the core concept of this thesis research proofs of some classical scheduling algorithms is presented along with a brief introduction in the context they appear while literature review.

3.1 Classical Scheduling Problem

When $\alpha = 1$, the problem is exactly the classical independent tasks scheduling problem on identical machines, which is known to be NP-Hard [9]. We use Graham's List Scheduling (LS) [11] and Largest Processing Time (LPT) algorithms [12] to derive approximation ratios in different scenarios. The LS algorithm takes tasks one at a time and assigns them to the processor having the least load at that time. LS is a 2-approximation algorithm and is widely used in online scheduling problems. LPT sorts the tasks in a non-increasing order of processing time and assigns them one at a time in this order to the processor with the smallest current load. The LPT algorithm has a worst case approximation ratio of $\frac{4}{3} - \frac{1}{3m}$ in the offline setting. One can even obtain an arbitrarily good approximation algorithm for this problem by increasing its complexity with a dual approximation algorithm [14].

We begin with by recalling the formal proofs of the guarantees of LS and LPT

algorithms:

Property 1. [11] List Scheduling has an approximation ratio of $2 - \frac{1}{m}$.

Proof: Let l be the last task in the system which is processed on machine r and it starts on r at time t . Clearly, the makespan C_{max} of the schedule is $t + p_l$. As in LS a new task is scheduled on the least loaded machine at that time, for each machine i , we have $t \leq \sum_{j \in E_i} p_j$. Adding this for all the machines including r , we get $mt \leq \sum_{i \in M - \{r\}} \sum_{j \in E_i} p_j + \sum_{j \in E_r} p_j - p_l \Rightarrow t \leq (\sum_j p_j - p_l)/m$. Hence, $C_{max} \leq \frac{\sum_j p_j + (m-1)p_l}{m}$.

The optimal makespan of a schedule C_{max}^* must be greater or equal to the average load over all the m machines, $C_{max}^* \geq \frac{\sum_j p_j}{m}$. Also, C_{max}^* cannot be smaller than any task in the system, hence $C_{max}^* \geq \text{Max} p_j \geq p_l$. Therefore, $C_{max} \leq C_{max}^* + C_{max}^*(m - 1)/m$. Hence, $C_{max}/C_{max}^* \leq 2 - 1/m$. \square

Property 2. [12] The LPT algorithm has an approximation ratio of $\frac{4}{3} - \frac{1}{3m}$.

Proof: LPT always generates an optimal schedule if no machine has more than 2 tasks. So to derive an approximation ratio we can assume that there are at least 3 tasks in a machine. As LPT assigns tasks to machines in non-increasing order of their processing times, the last task l is the smallest task in the machine. Since there are at least 3 tasks in a machine $C_{max}^* \geq 3p_l$. Also, $C_{max}^* \geq \frac{\sum_j p_j}{m}$ and $C_{max} \leq \frac{\sum_j p_j + (m-1)p_l}{m}$ as shown in previous proof. Therefore, $C_{max} \leq C_{max}^* + C_{max}^*(m - 1)/3m$. Hence, $C_{max}/C_{max}^* \leq 4/3 - 1/3m$. \square

3.2 Uncertainty and Robustness

Based on various models for describing the uncertain input parameter, various methodologies can be used including reactive, stochastic, fuzzy and robust approach [16].

We are using the bounded uncertainty model which assumes that an input parameter have value between a lower and upper bound. Wierman and Nuyens [28] introduce SMART, a classification to understand size-based policies and draw analytic co-relation between response time and estimated job size in single server problem. Robust approaches to deal with uncertainty are widely used on MapReduce systems [15] [25], in Hadoop [30] [27], on databases [17] and on web servers [3]. The HSFS and FLEX schedulers provide robustness in scheduling against uncertain job size [30, 19]. Cannon and Jeannot [2] analyzed the correlation between various metrics used to measure robustness and provided scheduling heuristics that optimizes both makespan and robustness for scheduling task graph on heterogeneous system.

Most of the work on robust scheduling use scenarios to structure the variability of uncertain parameters. Daniels and Kouvelis [5] used them to optimize the flow-time using a single machine. Davenport, Gefflot, and Bek analyzed slack based technique (adding extra idle time) to cope with uncertainty [6]. Gatto and Widmayer derives bounds on competitive ratio of Graham’s online algorithm in scenario where processing times of jobs either increase or decrease arbitrarily due to perturbations [10]. These works considered augmenting or decreasing of job processing times as different problem scenario that need to be optimized. We approach the problem using worst case analysis where some tasks may increase and some may decrease within the same schedule.

3.3 Data Placement and Replication

Data placement and replication methodologies are highly used in distributed systems including peer-to-peer and Grid systems to achieve effective data management and improve performance [4][1][20]. Tse [24] used selective replication of documents to increase the available bandwidth to serve files using web servers and study the problem through bi-criteria optimization techniques to maximize the quality of service and minimize the memory occupation. Our approach for bi-criteria optimization of makespan and memory consumption is based on the SBO_{Δ} Algorithm proposed by Saule, Dutot and Mounie [22]. The algorithm uses ρ_1 and ρ_2 approximated independent schedules on makespan and memory consumption respectively, and it computes a $((1+\Delta)\rho_1, (1+\frac{1}{\Delta})\rho_2)$ - approximated schedule with Δ as a parameter of the algorithm.

Property 3. [22] The SBO_{Δ} Algorithm generates a $(1+\Delta)\rho_1$ -approximated schedule on makespan.

Proof: The algorithm schedules a task j according to a π_2 schedule generated by the ρ_2 -approximated algorithm on memory if it satisfies this condition: $\frac{p_j}{C_{max}^{\pi_1}} \leq \Delta \frac{s_j}{M_{max}^{\pi_2}}$. Where $C_{max}^{\pi_1}$ is the makespan obtained using a π_1 schedule generated by the ρ_1 -approximated algorithm on makespan, and $M_{max}^{\pi_1}$ is the memory consumption of the most occupied machine obtained using π_2 . If this condition is not satisfied the task is scheduled according π_1 . Let k be the machine reaching makespan C_{max} of the schedule generated by the SBO_{Δ} algorithm. Let S_1 be the set of tasks scheduled according π_1 and S_2 be the set of tasks scheduled according π_2 schedule. C_{max} can be decomposed as the sum of the processing times of the tasks in set S_1 and S_2 scheduled

on machine k .

$$C_{max} = \sum_{j \in S_1 \cap E_k} p_j + \sum_{j \in S_2 \cap E_k} p_j$$

Since $C_{max}^{\pi_1} \geq \sum_{j \in S_1 \cap E_k} p_j$ and $\sum_{j \in S_2 \cap E_k} \Delta C_{max}^{\pi_1} \frac{s_j}{M_{max}^{\pi_2}} \geq \sum_{j \in S_2 \cap E_k} p_j$ by definition of S_2 , we have

$$C_{max} \leq C_{max}^{\pi_1} + \sum_{j \in S_2 \cap E_k} \Delta C_{max}^{\pi_1} \frac{s_j}{M_{max}^{\pi_2}}$$

Since $\sum_{j \in S_2 \cap E_k} \frac{s_j}{M_{max}^{\pi_2}} \leq 1$, we have

$$C_{max} \leq (1 + \Delta) C_{max}^{\pi_1}$$

Since $C_{max}^{\pi_1} \leq \rho_1 C_{max}^*$, the algorithm has an approximation ratio of $(1 + \Delta)\rho_1$ on the makespan. \square

Property 4. [22] The SBO_Δ Algorithm generates a $(1 + \frac{1}{\Delta})\rho_2$ -approximated schedule on memory.

Proof: Let k be the machine with most memory consumption. Similar to the previous proof, M_{max} can be written as the sum of memory usage of the tasks in sets S_1 and S_2 scheduled on machine k , $\sum_{j \in S_1 \cap E_k} s_j + \sum_{j \in S_2 \cap E_k} p_j$. Since, $\sum_{j \in S_2 \cap E_k} s_j \leq M_{max}^{\pi_1}$ and by definition of S_1 , $\sum_{j \in S_1 \cap E_k} s_j \leq \frac{M_{max}^{\pi_1}}{\Delta C_{max}^{\pi_2}} \sum_{j \in S_2 \cap E_k} p_j \leq \frac{M_{max}^{\pi_1}}{\Delta}$, we have

$$M_{max} \leq (1 + \frac{1}{\Delta}) M_{max}^{\pi_1}$$

Since $M_{max}^{\pi_1} \leq \rho_2 M_{max}^*$, the algorithm has an approximation ratio of $(1 + \frac{1}{\Delta})\rho_2$ on memory. \square

CHAPTER 4: REPLICATED DATA PLACEMENT STRATEGIES

This chapter provides three strategies, each offering different degree of replication to study the tradeoff between the number of replication and the guarantee on the makespan. The strategy *No Replication* restricts that each task can be scheduled to only one machine and allows no replication of the tasks. The strategy *Replicate data everywhere* replicates data everywhere and studies what can be achieved in doing so. The strategy *Replication in groups* replicates data in group of processors and act an intermediate strategy between the previous two strategies.

4.1 Strategy 1: No Replication

This section considers the situation where the data of each task is restricted to be on only one machine, *i.e.*, $\forall j, |M_j| = 1$. We have a set J of n jobs, and a set M of m machines. Let $f : J \mapsto M$ be a function that assigns each job to exactly one machine. The restriction that the data of each task is deployed on a single machine puts all the decision in phase 1: each task can only be scheduled on one machine in phase 2.

Theorem 5. When $|M_j| = 1$, there is no online algorithm having competitive ratio better than $\frac{\alpha^2 m}{\alpha^2 + m - 1}$.

Proof: We use the adversary technique to prove the lower bound of this theorem. An adversary discloses the input instance piece by piece. It analyzes the choices made by the algorithm to change the part of the instance that has not been disclosed yet. That

way it can build an instance that maximizes the competitive ratio of the algorithm.

Let us consider an instance with λm tasks of equal estimated processing time $\forall j, \tilde{p}_j = 1$. After phase 1, let i be the most loaded processor which has B tasks. Obviously, $B \geq \lambda$. In phase 2 the adversary increases the processing time of the tasks on processor i by a factor of α and changes the processing time of the other tasks by a factor of $\frac{1}{\alpha}$. So, $C_{max} = \alpha B$. C_{max}^* will be no worse than any feasible solution. In particular, the solution that distributes equally the jobs of size α and the jobs of size $\frac{1}{\alpha}$. Therefore $C_{max}^* \leq \frac{1}{\alpha} \lceil \frac{\lambda m - B}{m} \rceil + \alpha \lceil \frac{B}{m} \rceil$. Figure 1 depicts the online solution and the offline optimal. We have,

$$\frac{C_{max}}{C_{max}^*} \geq \frac{\alpha^2 B}{\lceil \frac{\lambda m - B}{m} \rceil + \alpha^2 \lceil \frac{B}{m} \rceil}$$

Since $\frac{\lambda m - B}{m} + 1 \geq \lceil \frac{\lambda m - B}{m} \rceil$ and $\frac{B}{m} + 1 \geq \lceil \frac{B}{m} \rceil$, we have

$$\frac{C_{max}}{C_{max}^*} \geq \frac{\alpha^2 B}{\frac{\lambda m - B}{m} + 1 + \alpha^2 \frac{B}{m} + \alpha^2}$$

From above expression it is clear that the smaller the value of B , the more the value of the expression decreases. So, any algorithm should minimize B to achieve better performance. For a schedule to be feasible the condition $B \geq \lambda$ must be satisfied. For $B = \lambda$ the value of $\frac{C_{max}}{C_{max}^*}$ is minimum and is equal to $\frac{\alpha^2 m \lambda}{\lambda(\alpha^2 + m - 1) + m(\alpha^2 + 1)}$. The adversary can maximize the ratio of the algorithm by arbitrarily increasing λ . When λ tends to ∞ , we have

$$\frac{C_{max}}{C_{max}^*} \geq \frac{\alpha^2 m}{\alpha^2 + m - 1}$$

□

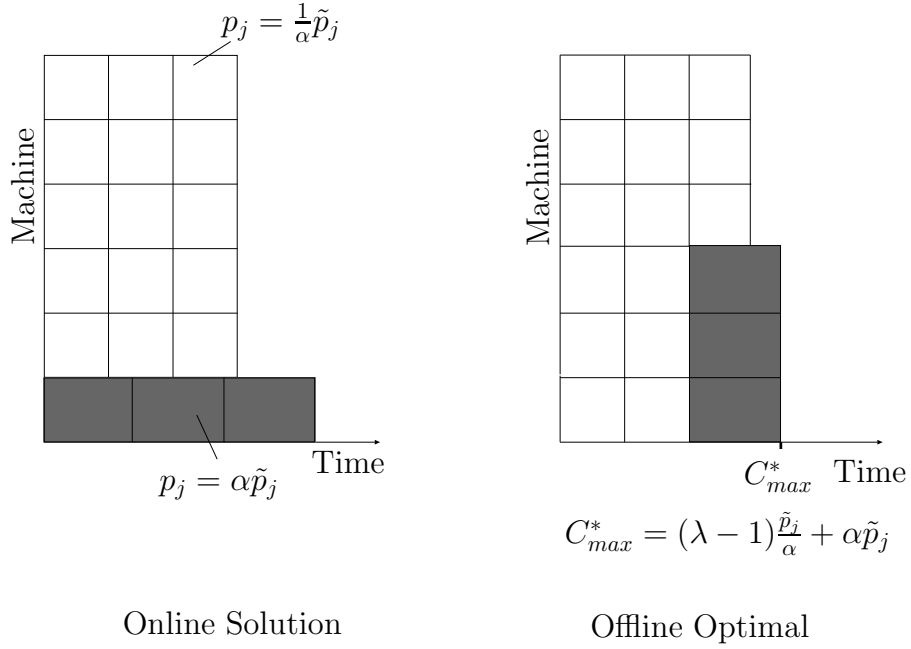


Figure 1: Instance constructed by the adversary in the proof of theorem 5 with $\lambda = 3$ and $m = 6$. In the online solution, the adversary increases the processing time of a task of the most loaded machine by a factor of α . If that information was available beforehand, an optimal offline algorithm could have distributed these longer tasks to other processors.

Corollary 5.1. When m goes to ∞ there is no online algorithm having competitive ratio better than α^2 .

4.1.1 Algorithm

We present the algorithm *LPT-No Choice*. In phase 1, the algorithm distribute the data of the tasks to the processor using their estimated processing times according to Graham's LPT algorithm [12]: The tasks are sorted in non-increasing order of their processing time and are greedily scheduled on the processor that minimizes the sum of \tilde{p}_j of the tasks allocated on that processor. Since there is no replication, there is no decision to take in phase 2.

The performance of the algorithm depends mostly on how much the actual processing times of the tasks differ from their estimation. It also depends on the existence of a better arrangement if the actual processing times were known. The following theorem states the theoretical guarantee of the algorithm.

Theorem 6. The *LPT-No Choice* has a competitive ratio of $\frac{2\alpha^2 m}{2\alpha^2 + m - 1}$.

Proof: The algorithm assigns the jobs to processors based on their estimated processing times using LPT in Phase 1. So, the planned makespan considering the estimated processing times of tasks, \tilde{C}_{max} have the following relation with the total estimated processing time, \tilde{p}_j and estimated processing time of the task l that reaches \tilde{C}_{max} .

$$\tilde{C}_{max} \leq \frac{\sum \tilde{p}_j + (m-1)\tilde{p}_l}{m} \quad (2)$$

The actual makespan of a schedule, C_{max} , obtained using the actual processing times of all the jobs, must be smaller than $C_{max} \leq \alpha \tilde{C}_{max}$ (thanks to Equation 1).

We have following inequality:

$$C_{max} \leq \alpha \tilde{C}_{max} \leq \alpha \left(\frac{\sum \tilde{p}_j + (m-1)\tilde{p}_l}{m} \right) \quad (3)$$

The worst case situation is when the task of the processor where the sum of estimated processing time is \tilde{C}_{max} sees the actual processing time of its task being α times larger than their estimate; meanwhile the processing time of the task on the rest of the processors is $\frac{1}{\alpha}$ times their estimation. The argument behind this statement is that greater the value of ratio $\frac{C_{max}}{\sum p_j}$, the worse the algorithm approximation ratio will be. So the total actual processing time is given by the following equation.

$$\sum p_j = \frac{\sum \tilde{p}_j - C_{max}}{\alpha} + \alpha \tilde{C}_{max} \quad (4)$$

Also the actual optimal makespan have following constraint

$$C_{max}^* \geq \frac{\sum p_j}{m}$$

Substituting for $\sum p_j$, we have

$$\begin{aligned} mC_{max}^* &\geq \frac{\sum \tilde{p}_j - C_{max}}{\alpha} + \alpha \tilde{C}_{max} \\ mC_{max}^* &\geq \frac{\sum \tilde{p}_j - \left(\frac{\sum \tilde{p}_j + (m-1)\tilde{p}_l}{m} \right)}{\alpha} + C_{max} \\ mC_{max}^* &\geq \frac{m-1}{\alpha m} \left(\sum \tilde{p}_j - \tilde{p}_l \right) + C_{max} \end{aligned}$$

By the property of LPT, $\sum \tilde{p}_j - \tilde{p}_l \geq m(\tilde{C}_{max} - \tilde{p}_l)$, we have,

$$mC_{max}^* \geq \frac{m-1}{\alpha} (C_{max} - \tilde{p}_l) + C_{max}$$

All instances where there is only one task per processor is always optimal. There-

fore, we can restrict our analysis without loss of generality to instances with at least two jobs per processor. (Notice that in the original proof of Graham's LPT [12], an argument is made that all instances with two tasks per machine are optimal. However, the argument does not port in our case where only estimated processing times are known.) For at least two jobs on the processor that reaches \tilde{C}_{max} , the (estimated) processing time of last job is smaller than half the estimated makespan, $\tilde{p}_l \leq \frac{\tilde{C}_{max}}{2}$.

Substituting this expression in the above equation, we have

$$mC_{max}^* \geq \frac{m-1}{\alpha} \left(\tilde{C}_{max} - \frac{\tilde{C}_{max}}{2} \right) + C_{max}$$

Using equation 3,

$$\begin{aligned} mC_{max}^* &\geq \frac{m-1}{2\alpha} \frac{C_{max}}{\alpha} + C_{max} \\ mC_{max}^* &\geq \left(\frac{m-1}{2\alpha^2} + 1 \right) C_{max} \\ \frac{C_{max}}{C_{max}^*} &\leq \frac{2\alpha^2 m}{2\alpha^2 + m - 1} \end{aligned}$$

□

4.2 Strategy 2: Replicate Data Everywhere

With this strategy, we put no restriction on phase 2. The tasks are replicated everywhere i.e. $\forall j, |M_j| = |M|$. We introduce the *LPT-No Restriction* which replicates the data of all the tasks on each machine in the first phase. In the second phase we simply use the Longest Processing Time algorithm (LPT) in an online fashion using the estimated processing time of the task. That is to say, the tasks are sorted in non-increasing order of their estimated processing time. Then the task are greedily allocated on the first processor that becomes available. Note that this is done in

phase 2, the processor become available with when the actual processing time of the task scheduled onto it elapse.

Lemma 7. Let l be the task that reaches C_{max} in the solution constructed by *LPT-No Restriction*. If there are at least two tasks on the machine that executes l in *LPT-No Restriction*, then $C_{max}^* \geq \frac{2}{\alpha^2} p_l$.

Proof: Since there are at least two tasks on the machine that executes l in *LPT-No Restriction*, there are at least $m + 1$ tasks i such that $\tilde{p}_j \geq \tilde{p}_l$. Therefore in any solution at least one machine gets two tasks c and d , such that $\tilde{p}_c \geq \tilde{p}_l$ and $\tilde{p}_d \geq \tilde{p}_l$. C_{max}^* must be greater than sum of the processing time of these two tasks.

$$C_{max}^* \geq p_c + p_d$$

As the actual processing time of a task must be greater than $\frac{1}{\alpha}$ times of its estimated value, we have $p_c \geq \frac{1}{\alpha} \tilde{p}_c$ and $p_d \geq \frac{1}{\alpha} \tilde{p}_d$. Using this

$$C_{max}^* \geq \frac{1}{\alpha} \tilde{p}_c + \frac{1}{\alpha} \tilde{p}_d \geq \frac{2}{\alpha} \tilde{p}_l$$

Since, $\tilde{p}_l \geq \frac{1}{\alpha} p_l$, we have

$$C_{max}^* \geq \frac{2}{\alpha^2} p_l$$

□

Theorem 8. *LPT-No Restriction* has a competitive ratio of $\frac{C_{max}}{C_{max}^*} \leq 1 + (\frac{m-1}{m}) \frac{\alpha^2}{2}$

Proof: The optimal makespan, C_{max}^* must be at least equal to the average load on the m machines. We have

$$C_{max}^* \geq \frac{\sum p_j}{m} \tag{5}$$

By the property of LPT (actually, it is a property of List Scheduling which LPT is a refinement of) the load on each machine i is greater than the load on the machine which reach C_{max} before the last task l is scheduled. So for each machine i , $C_{max} \leq \sum_{j \in E_i} p_j + p_l$ holds true. Summing for all the machines we have

$$\begin{aligned} mC_{max} &\leq \sum p_j + (m-1)p_l \\ C_{max} &\leq \frac{\sum p_j}{m} + \frac{(m-1)}{m}p_l \end{aligned} \tag{6}$$

Using 5 and 6, we have

$$\frac{C_{max}}{C_{max}^*} \leq 1 + \frac{m-1}{m} \left(\frac{p_l}{C_{max}^*} \right)$$

Using Lemma 7, we have

$$\frac{C_{max}}{C_{max}^*} \leq 1 + \left(\frac{m-1}{m} \right) \frac{\alpha^2}{2}$$

□

Graham's List Scheduling algorithm always has a competitive ratio of $2 - \frac{1}{m}$. For $\alpha^2 < 2$, the *LPT-No Restriction* algorithm has better approximation than List Scheduling. For $\alpha^2 > 2$ List Scheduling has better guarantee than the one expressed in Theorem 8. Since *LPT-No Restriction* is a variant of List Scheduling, the algorithm has a competitive ratio of $\min(1 + \frac{m-1}{2m}\alpha^2, 2 - \frac{1}{m})$.

4.3 Strategy 3: Replication in Groups

This strategy partitions the processors into k groups $G1, G2 \dots Gk$. The size of each group is equal and have $\frac{m}{k}$ processors within each group. For the sake of simplicity,

we assume that we will only use values of k such that k divides m . In the first phase, the data of each task is replicated on all the processors of one of the k groups, i.e. $\forall j, |M_j| = \frac{m}{k}$. In the second phase the tasks are scheduled within the group they are assigned to in first phase. Figure 2 shows the construction of two phases.

We propose the *LS-Group* algorithm which is based on Graham's List Scheduling algorithm. In phase 1, we use List Scheduling to distribute the tasks to the k groups of processors. In phase 2 each task is scheduled to a particular processor within the group it was allocated in phase 1 using the online List Scheduling algorithm.

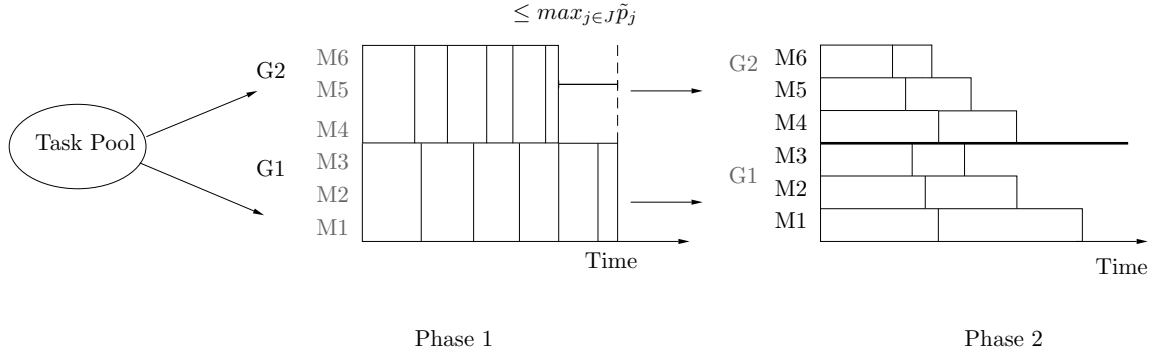


Figure 2: An example of replication in groups with $m = 6$, $k = 2$. In phase 1, the data of the tasks are assigned to one of the groups. Phase 2 schedules each task assigned to a machine within its group.

Theorem 9. With k groups, the competitive ratio of *LS-Group* is $\frac{k\alpha^2}{\alpha^2+k-1}(1 + \frac{k-1}{m}) + \frac{m-k}{m}$

Proof: We assume without loss of generality that C_{max} comes from group $G1$. C_{max}^* must be greater than the average of the loads on the machines.

$$C_{max}^* \geq \frac{\sum_{j \in J} p_j}{m}$$

$\sum_{j \in J} p_j$ can be written as sum of load on $G1$ and load on rest of groups.

$$C_{max}^* \geq \frac{\sum_{j \in G1} p_j + \sum_{l=2}^k \sum_{j \in Gl} p_j}{m} \quad (7)$$

As in phase 1 tasks are allocated to different groups using List Scheduling with the estimated processing times of the tasks, the (estimated) load difference between any two groups cannot be greater than the estimated value of largest task $\max_{j \in J} \tilde{p}_j$. So, for any group $Gl \neq G1$, We have

$$\forall l \in \{2, 3, \dots, k\}, \left| \sum_{j \in G1} \tilde{p}_j - \sum_{j \in Gl} \tilde{p}_j \right| \leq \max_{j \in J} \tilde{p}_j$$

Adding for all values of l leads to

$$\left| (k-1) \sum_{j \in G1} \tilde{p}_j - \sum_{l=2}^k \sum_{j \in Gl} \tilde{p}_j \right| \leq (k-1) \max_{j \in J} \tilde{p}_j$$

Case 1: If $(k-1) \sum_{j \in G1} \tilde{p}_j > \sum_{l=2}^k \sum_{j \in Gl} \tilde{p}_j$.

$$\sum_{l=2}^k \sum_{j \in Gl} \tilde{p}_j \geq (k-1) \left(\sum_{j \in G1} \tilde{p}_j - \max_{j \in J} \tilde{p}_j \right)$$

As the actual processing time of the tasks can vary within a factor α and $\frac{1}{\alpha}$ of their estimated processing time, the following inequality holds

$$\begin{aligned} \alpha \sum_{l=2}^k \sum_{j \in Gl} p_j &\geq (k-1) \left(\frac{1}{\alpha} \sum_{j \in G1} p_j - \alpha \max_{j \in J} p_j \right) \\ \sum_{l=2}^k \sum_{j \in Gl} p_j &\geq (k-1) \left(\frac{1}{\alpha^2} \sum_{j \in G1} p_j - \max_{j \in J} p_j \right) \end{aligned} \quad (8)$$

Phase 2 applies List Scheduling in the online mode. We assumed that C_{max} comes

from $G1$. Using the guarantees of List Scheduling we can write,

$$C_{max} \leq \frac{\sum_{j \in G1} p_j}{m/k} + \frac{m/k - 1}{m/k} p_{max} \quad (9)$$

where p_{max} is actual processing time of longest task in $G1$.

From Equation 8 and 7, we derive

$$\begin{aligned} C_{max}^* &\geq \frac{\sum_{j \in G1} p_j + (k-1) \left(\frac{1}{\alpha^2} \sum_{j \in G1} p_j - \max_{j \in J} p_j \right)}{m} \\ \alpha^2(mC_{max}^* + (k-1)\max_{j \in J} p_j) &\geq (\alpha^2 + k - 1) \sum_{j \in G1} p_j \\ \frac{\alpha^2}{\alpha^2 + k - 1} (mC_{max}^* + (k-1)\max_{j \in J} p_j) &\geq \sum_{j \in G1} p_j \end{aligned} \quad (10)$$

Using 9 and 10, We have

$$\begin{aligned} C_{max} &\leq \frac{k\alpha^2}{\alpha^2 + k - 1} \left(C_{max}^* + \frac{k-1}{m} \max_{j \in J} p_j \right) \\ &\quad + \frac{m/k - 1}{m/k} p_{max} \end{aligned}$$

As $C_{max}^* \geq \max_{j \in J} p_j \geq p_{max}$, we have

$$\begin{aligned} C_{max} &\leq \frac{k\alpha^2}{\alpha^2 + k - 1} \left(C_{max}^* + \frac{k-1}{m} C_{max}^* \right) \\ &\quad + \frac{m-k}{m} C_{max}^* \end{aligned}$$

So, in Case 1 the algorithm has a competitive ratio of,

$$\frac{C_{max}}{C_{max}^*} \leq \frac{k\alpha^2}{\alpha^2 + k - 1} \left(1 + \frac{k-1}{m} \right) + \frac{m-k}{m}$$

Case 2: If $(k-1) \sum_{j \in G1} \tilde{p}_j \leq \sum_{l=2}^k \sum_{j \in Gl} \tilde{p}_j$.

Since the processing times of the tasks can vary within a factor α and $\frac{1}{\alpha}$ of their estimated values, the expression for case 2 can be written as

$$\sum_{l=2}^k \sum_{j \in Gl} p_j \geq \frac{1}{\alpha^2} (k-1) \sum_{j \in G1} p_j$$

Putting this value in Equation 7, we have

$$C_{max}^* \geq \frac{\alpha^2 + k - 1}{m\alpha^2} \sum_{j \in G1} p_j \quad (11)$$

Using Equations 9 and 11, and as $C_{max}^* \geq p_{max}$, we have

$$C_{max} \leq \frac{k\alpha^2}{\alpha^2 + k - 1} C_{max}^* + \frac{m - k}{m} C_{max}^*$$

So, in case 2 the algorithm has a competitive ratio of $\frac{k\alpha^2}{\alpha^2 + k - 1} + \frac{m - k}{m}$.

Clearly, the algorithm has a worst competitive ratio in case 1. So, the algorithm has a competitive approximation ratio of $\frac{C_{max}}{C_{max}^*} \leq \frac{k\alpha^2}{\alpha^2 + k - 1} \left(1 + \frac{k-1}{m}\right) + \frac{m-k}{m}$. \square

LS-Group uses List Scheduling in both its phases. A LPT-based algorithm may have better guarantee. But without performing any replication, *i.e.* when $k = m$, the *LS-Group* algorithm has a competitive ratio almost equal to *LPT-No choice*'s when the number of machines m is large and the value of α is within practical range. This indicates an LPT-based algorithm for strategy 3 would likely not have a much more interesting guarantee.

4.4 Summary

Table 1 summarizes the results of this chapter in term of approximation theory. Based on adversary technique, Theorem 5 states that there is no algorithm which can give performance better than $\frac{\alpha^2 m}{\alpha^2 + m - 1}$ for the model where no replication is allowed. *LPT-No Choice* is a $\frac{2\alpha^2 m}{2\alpha^2 + m - 1}$ -approximation that uses that strategy. For the second strategy that replicates the data of all tasks everywhere ($|M_j| = |M|$), *LPT-No Restriction* achieves a competitive ratio of $1 + (\frac{m-1}{m})\frac{\alpha^2}{2}$. The third strategy uses replication within k groups of size m/k (i.e., $|M_j| = m/k$). Using this strategy, the *LS-Group* algorithm has a competitive ratio of $\frac{k\alpha^2}{\alpha^2 + k - 1} \left(1 + \frac{k-1}{m}\right) + \frac{m-k}{m}$.

Replication	Approximation ratio
$ M_j = 1$	$\frac{C_{max}}{C_{max}^*} \leq \frac{2\alpha^2 m}{2\alpha^2 + m - 1}$ (Th. 6) No approximation better than $\frac{\alpha^2 m}{\alpha^2 + m - 1}$ (Th. 5)
$ M_j = m$	$\frac{C_{max}}{C_{max}^*} \leq 1 + (\frac{m-1}{m})\frac{\alpha^2}{2}$ (Th. 8) $\frac{C_{max}}{C_{max}^*} \leq 2 - \frac{1}{m}$ [11]
$ M_j = \frac{m}{k}$	$\frac{C_{max}}{C_{max}^*} \leq \frac{k\alpha^2}{\alpha^2 + k - 1} \left(1 + \frac{k-1}{m}\right) + \frac{m-k}{m}$ (Th. 9)

Table 1: Summary of the contribution of this chapter. Three proposed algorithms have guaranteed performance. One lower bound on approximability has been established.

Of course, there is an inherent tradeoff between replicating data and obtaining good values for the makespan. To better understand the tradeoff we show in Figure 3 how the expressions of the guarantees (or impossibility) translate to actual values in a approximation ratio / replication space. We picked 3 values of α while keeping the number of machines fixed $m = 210$.

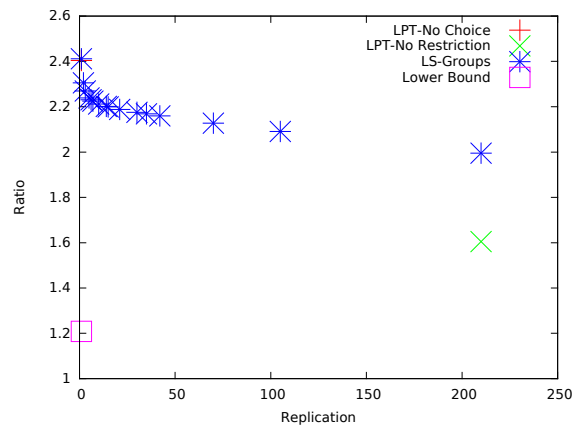
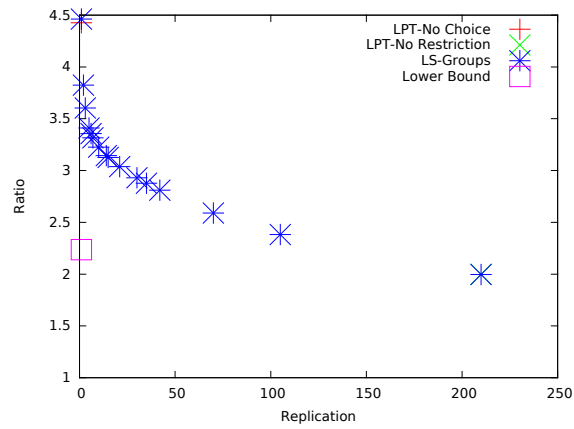
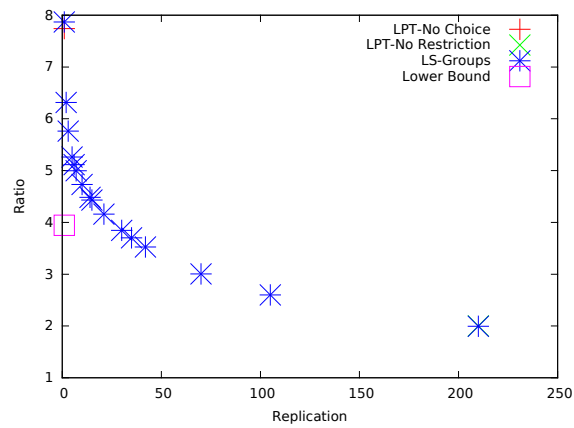
When $\alpha = 1.1$, even with multiple groups *LS-Group* provides little improvement

over *LPT-No Choice*. However there is a significant gap between the guarantee of *LPT-No Choice* and the lower bound on possible approximation. When α is small, there is a significant improvement in using *LPT-No Restriction* over using simply *LS-Group* with only 1 group.

When α increases to 1.5, there is no more differences in the guarantees of *LS-Group* with 1 group and *LPT-No Restriction*. Also *LS-Group* provides many intermediate solution between deploying the data on a single machine and deploying them everywhere.

When $\alpha = 2$, the range of the approximation ratios increase and the value of the lower bound increases. Now *LS-Group* is able to get a better approximation using less than 50 replications than is possible by deploying data on a single machine. Also, the approximation ratio quickly improves from more than 7.5 with the data being replicated on 1 machine to a ratio of less than 6 with only replicating the data on 3 machines.

Overall, when α is large, only few replication improve the performance significantly.

(a) $m = 210, \alpha = 1.1$ (b) $m = 210, \alpha = 1.5$ (c) $m = 210, \alpha = 2$ Figure 3: Ratio-Replication graph with $m = 210$ and $\alpha \in \{1.1, 1.5, 2\}$.

CHAPTER 5: MEMORY AWARE REPLICATION UNDER UNCERTAINTY

Replication improves performance but incurs a cost in terms of memory consumption. Replication allows to obtain a better load balancing by reducing the effect of uncertainties in processing times of tasks. But each replica occupies memory, and increases the memory consumption. So, replicating all the tasks is not possible in real scenarios. This justifies the need for an efficient replication strategy which allows an algorithm to choose which tasks are to be replicated and where. In this chapter we investigate the bi-objective problem of minimizing the makespan as well as the memory consumption. A memory-aware replication strategy improves execution times with little increase in memory consumption.

5.1 Preliminaries

The problem is to schedule a set J of n tasks on m machines such that both makespan C_{max} as well as memory usage M_{max} is optimized. Let π_1 be the schedule which minimizes makespan and π_2 be the memory-aware schedule. $\tilde{C}_{max}^{\pi_1}$ and \tilde{C}_{max}^* are makespan and optimal makespan when all the tasks are scheduled according to π_1 . Similarly, $M_{max}^{\pi_2}$ is memory consumption of the most occupied machine and M_{max}^* is its optimal value. The strategy is to divide tasks into two sets S_1 and S_2 such that set S_1 contains the processing time intensive tasks and set S_2 contains the memory intensive tasks, and schedule them differently and in such a way that it optimizes

both the objectives.

We propose two algorithms $SABO_{\Delta}$ (stands for static asymmetric bi-objective) and ABO_{Δ} (stands for asymmetric bi-objective), which are based on SBO_{Δ} algorithm. SBO_{Δ} [22] is bi-objective algorithm for minimizing makespan and memory usage for independent tasks by combining results of two symmetric schedules each dedicated to a single objective.

5.2 The $SABO_{\Delta}$ Algorithm

We propose $SABO_{\Delta}$ Algorithm which is static in nature and restrict each task to be scheduled to only one machine. Similar to SBO_{Δ} this algorithm assigns tasks to all the machines in phase 1 such that it minimizes both the objectives. As each task is restricted to only one machine, there is no task replication. Based on similar condition as in SBO_{Δ} , a processing-time intensive task is assigned to π_1 schedule and a memory intensive task is assigned to π_2

In phase 2, the algorithm loads the tasks to the machines they were assigned in phase 1.

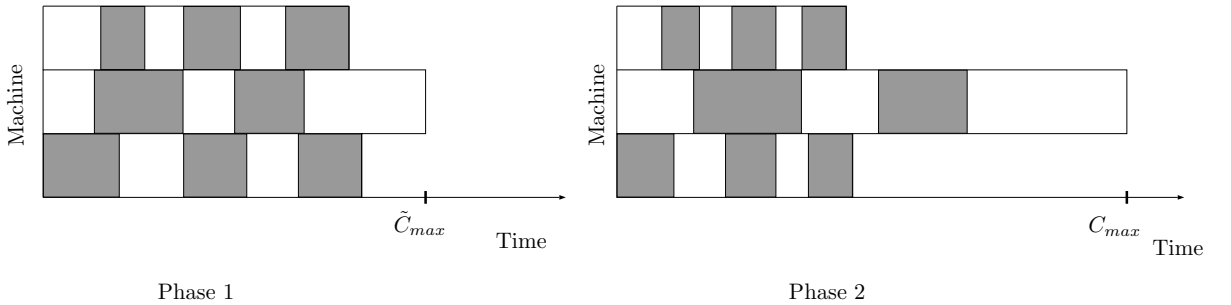


Figure 4: An example of two phases of the schedule generated by the $SABO_{\Delta}$. The uncolored parts represent tasks scheduled according π_2 . The colored parts represents tasks scheduled according π_1

Algorithm 1 $SABO_{\Delta}$

Input: m machines

Set J of n tasks

Let π_1 be a ρ_1 -approximated schedule on makespan \tilde{C}_{max}

Let π_2 be a ρ_2 - approximated schedule on memory M_{max}

Phase 1: [Uses SBO_{Δ}]

for all $j \in J$ **do**

if $\frac{\tilde{p}_j}{\tilde{C}_{max}^{\pi_1}} \leq \Delta \frac{s_j}{M_{max}^{\pi_2}}$ **then**

 Assign j to a machine according to π_2 schedule

 Add j to S_2

else

 Assign j to a machine according to π_1 schedule

 Add j to S_1

end if

end for

End of Phase 1

Phase 2:

Schedule tasks to machines to which they were assigned during phase 1

End of Phase 2

Theorem 10. The $SABO_{\Delta}$ Algorithm generates a $(1+\Delta)\alpha^2\rho_1$ - approximated schedule on makespan.

Proof: Let k be the machine reaching the makespan C_{max} of the schedule. C_{max} can be written as the sum of processing times of tasks in set S_1 and S_2 scheduled on machine k .

$$C_{max} = \sum_{j \in S_1 \cap E_k} p_j + \sum_{j \in S_2 \cap E_k} p_j$$

$$\text{Since, } \sum_{j \in S_2 \cap E_k} p_j \leq \alpha \sum_{j \in S_2 \cap E_k} \tilde{p}_j$$

$$C_{max} \leq \sum_{j \in S_1 \cap E_k} p_j + \alpha \sum_{j \in S_2 \cap E_k} \tilde{p}_j$$

Let $C_{max}^{\pi_1}$ denotes the makespan obtained after phase 2 when tasks are loaded and

actual processing time of a task is known to scheduler. Since $C_{max}^{\pi_1} \geq \sum_{j \in S_1 \cap E_k} p_j$ and

$\sum_{j \in S_2 \cap E_k} \Delta \tilde{C}_{max}^{\pi_1} \frac{s_j}{M_{max}^{\pi_2}} \geq \sum_{j \in S_2 \cap E_k} \tilde{p}_j$ by definition of S_2 , we have

$$C_{max} \leq C_{max}^{\pi_1} + \alpha \sum_{j \in S_2 \cap E_k} \Delta \tilde{C}_{max}^{\pi_1} \frac{s_j}{M_{max}^{\pi_2}}$$

Since, $C_{max}^{\pi_1} \leq \alpha \tilde{C}_{max}^{\pi_1}$ and $\sum_{j \in S_2 \cap E_k} \frac{s_j}{M_{max}^{\pi_2}} \leq 1$, we have

$$C_{max} \leq (1 + \Delta) \alpha \tilde{C}_{max}^{\pi_1}$$

Since $\tilde{C}_{max}^{\pi_1} \leq \rho_1 \tilde{C}_{max}^* \leq \alpha \rho_1 C_{max}^*$ the algorithm has an approximation ratio of $(1 + \Delta) \alpha^2 \rho_1$ on makespan. \square

Theorem 11. The $SABO_{\Delta}$ Algorithm generates $(1 + \frac{1}{\Delta}) \rho_2$ - approximated schedule on memory

Proof: The proof is identical to SBO_{Δ} algorithm and is presented in chapter 3. \square

5.3 The ABO_{Δ} Algorithm

We propose a two phase algorithm. In phase 1 the algorithm assigns tasks to all the machines such that it minimizes both the makespan as well as memory consumption. The tasks having more memory value in comparison to its processing time are scheduled using memory intensive schedule which aim at minimizing memory. Similarly tasks which incur more processing time cost compared to memory cost are assigned to machines according to the makespan intensive schedule. These tasks are replicated to all machines in order to provide better load balancing and hence minimized makespan. The algorithm in its phase 1 assigns all the memory intensive tasks to machines first, then chooses tasks having more processing time values compared to memory they consume.

In phase 2, the algorithm loads the memory intensive tasks to the machines they were assigned in phase 1 respecting the tasks assignment during phase 1. The algorithm schedule the time intensive tasks (replicated tasks) using Graham's List Scheduling after all the memory intensive tasks are scheduled. Figure 5 shows a schedule instance using the algorithm.

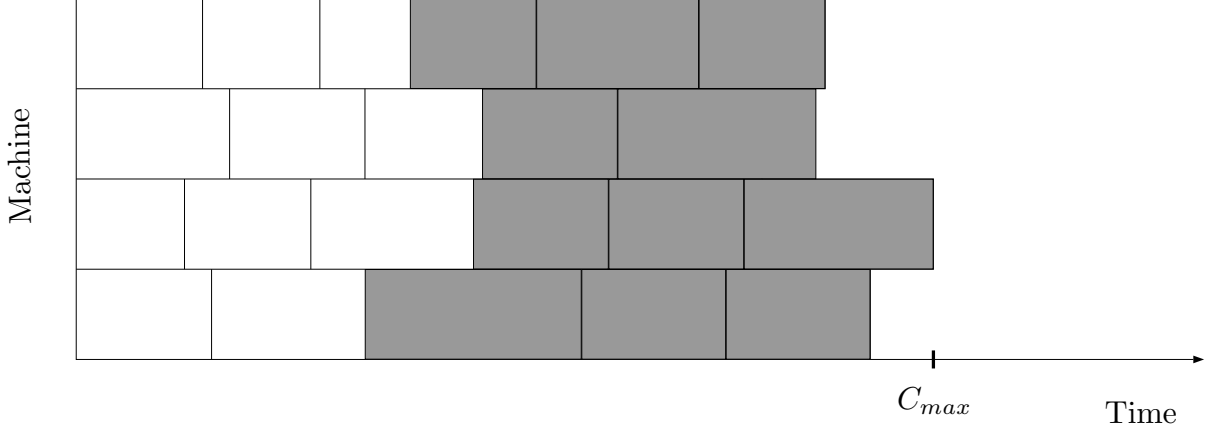


Figure 5: An example of the schedule generated by the ABO_{Δ} algorithm. The uncolored parts represent the memory intensive tasks scheduled according π_2 . The colored parts represent the processing time intensive tasks and scheduled using LS after replicated

Theorem 12. The ABO_{Δ} Algorithm generates a $(2 - \frac{1}{m} + \Delta\alpha^2\rho_1)$ - approximated schedule on makespan.

Proof: Let k be the machine reaching the makespan C_{max} of the schedule. C_{max} can be written as the sum of the processing times of tasks in sets S_1 and S_2 scheduled on machine k .

$$C_{max} = \sum_{j \in S_1 \cap E_k} p_j + \sum_{j \in S_2 \cap E_k} p_j$$

Algorithm 2 ABO_{Δ}

Input: m machines

Set J of n tasks

Let π_1 be a ρ_1 -approximated schedule on makespan \tilde{C}_{max}

Let π_2 be a ρ_2 -approximated schedule on memory M_{max}

Phase 1:

for all $j \in J$ **do**

if $\frac{\tilde{p}_j}{\tilde{C}_{max}^{\pi_1}} \leq \Delta \frac{s_j}{M_{max}^{\pi_2}}$ **then**

 Assign j to a machine according to π_2 schedule

 Add task j to set S_2

end if

end for

for all $j \in J$ **do**

if $\frac{\tilde{p}_j}{\tilde{C}_{max}^{\pi_1}} \geq \Delta \frac{s_j}{M_{max}^{\pi_2}}$ **then**

 Add j to set S_1

 Replicate j everywhere

end if

end for

End of Phase 1

Phase 2:

 Schedule tasks from set S_2 respecting job assignment during phase 1

 Schedule all replicated tasks from set S_1 using Graham's LS Algorithm

End of Phase 2

Since, $\sum_{j \in S_2 \cap E_k} p_j \leq \alpha \sum_{j \in S_2 \cap E_k} \tilde{p}_j$

$$C_{max} \leq \sum_{j \in S_1 \cap E_k} p_j + \alpha \sum_{j \in S_2 \cap E_k} \tilde{p}_j$$

Let C_{max}^R denotes makespan obtained by scheduling the replicated tasks using LS.

Since $C_{max}^R \geq \sum_{j \in S_1 \cap E_k} p_j$ and $\sum_{j \in S_2 \cap E_k} \Delta \tilde{C}_{max}^{\pi_1} \frac{s_j}{M_{max}^{\pi_2}} \geq \sum_{j \in S_2 \cap E_k} \tilde{p}_j$ by definition of S_2 , we have

$$C_{max} \leq C_{max}^R + \alpha \sum_{j \in S_2 \cap E_k} \Delta \tilde{C}_{max}^{\pi_1} \frac{s_j}{M_{max}^{\pi_2}}$$

Using the property of LS, the approximation ratio of the schedule incorporating only replicated tasks is $2 - \frac{1}{m}$. So, $C_{max}^R \leq (2 - \frac{1}{m})C_{max}^*$. Also, $\sum_{j \in S_2 \cap E_k} \frac{s_j}{M_{max}^{\pi_2}} \leq 1$.

$$C_{max} \leq (2 - \frac{1}{m})C_{max}^* + \alpha \Delta \tilde{C}_{max}^{\pi_1}$$

Also, $\tilde{C}_{max}^{\pi_1} \leq \rho_1 \tilde{C}_{max}^*$. Since \tilde{C}_{max}^* is the optimal makespan obtained after phase 1 considering estimated processing times of the tasks, we have, $\tilde{C}_{max}^* \leq \alpha C_{max}^*$. So, $\tilde{C}_{max}^{\pi_1} \leq \alpha \rho_1 C_{max}^*$. Using this, we have

$$C_{max} \leq (2 - \frac{1}{m})C_{max}^* + \alpha^2 \Delta \rho_1 C_{max}^*$$

Hence, we proved that the algorithm generates a $(2 - \frac{1}{m} + \Delta \alpha^2 \rho_1)$ - approximated schedule on makespan. \square

Theorem 13. The ABO_{Δ} Algorithm generates a $(1 + \frac{m}{\Delta})\rho_2$ - approximated schedule on memory.

Proof: When a task is replicated all its replica occupies space in memory and increase memory consumption. For m replicas the total memory consumption is m times of

the replicated tasks. Similar to proof of previous theorem, the highest maximum memory occupied by any machine k can be written as

$$M_{max} = \sum_{j \in S_1 \cap E_k} s_j + \sum_{j \in S_2 \cap E_k} s_j$$

As each task in set S_1 is replicated over all the machines, $\sum_{j \in S_1 \cap E_k} s_j = \sum_{j \in S_1} s_j$.

$$M_{max} = \sum_{j \in S_1} s_j + \sum_{j \in S_2 \cap E_k} s_j$$

$\sum_{j \in S_2 \cap E_k} s_j$ at most be equal to $M_{max}^{\pi_2}$ and $\sum_{j \in S_1} s_j$ is bounded by $\sum_{j \in J} M_{max}^{\pi_2} \frac{\tilde{p}_j}{\Delta \tilde{C}_{max}^{\pi_1}}$ as per condition for π_1 scheduling, using this we have

$$M_{max} \leq \sum_{j \in J} M_{max}^{\pi_2} \frac{\tilde{p}_j}{\Delta \tilde{C}_{max}^{\pi_1}} + M_{max}^{\pi_2}$$

Since $\sum_{j \in J} \tilde{p}_j \leq m \tilde{C}_{max}^{\pi_1}$, we have

$$M_{max} \leq \frac{m}{\Delta} M_{max}^{\pi_2} + M_{max}^{\pi_2}$$

Also, $M_{max}^{\pi_2} \leq \rho_2 M_{max}^*$. Hence, The Algorithm generate $(1 + \frac{m}{\Delta})\rho_2$ - approximated schedule on memory. \square

5.4 Summary

Table 2 summarizes the results for $SABO_\Delta$ and ABO_Δ algorithms. $SABO_\Delta$ is similar to SBO_Δ algorithm in its first phase and has a approximation ratio of $[(1 + \Delta)\alpha^2\rho_1, (1 + \frac{1}{\Delta})\rho_2]$ on makespan and memory. ABO_Δ is a $[(2 - \frac{1}{m} + \Delta\alpha^2\rho_1), (1 + \frac{m}{\Delta})\rho_2]$ - approximated algorithm on makespan and memory and replicate processing time intensive tasks to improve makespan.

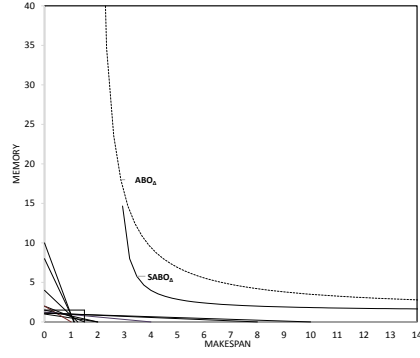
Algorithm	Approx. on makespan	Approx. on memory
$SABO_{\Delta}$	$(1 + \Delta)\alpha^2\rho_1$ (Th. 10)	$(1 + \frac{1}{\Delta})\rho_2$ (Th. 11)
ABO_{Δ}	$(2 - \frac{1}{m} + \Delta\alpha^2\rho_1)$ (Th. 12)	$(1 + \frac{m}{\Delta})\rho_2$ (Th. 13)

Table 2: Summary of the results of the algorithm $SABO_{\Delta}$ and the algorithm ABO_{Δ} .

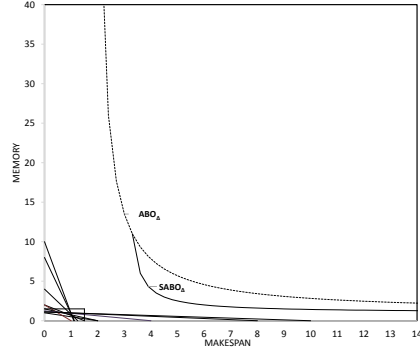
To better understand the tradeoff between memory consumption and makespan Figure 6 shows memory-makespan guarantees for the two algorithms. The bold lines shows impossibilities in the tradeoff between makespan and memory and means that no algorithm can guarantee better tradeoff than this. [22] discusses about these impossibilities in context of SBO_{Δ} algorithm.

The graph shows that for higher values for α the algorithm ABO_{Δ} have better tradeoff between memory-makespan than that of $SABO_{\Delta}$. For $\alpha\rho_1 \geq 2$, ABO_{Δ} always have better guarantee on makespan than $SABO_{\Delta}$. So, a schedule more centric to optimize makespan should follow ABO_{Δ} algorithm. And a memory centric schedule should follow $SABO_{\Delta}$ as the algorithm always has better guarantee on memory.

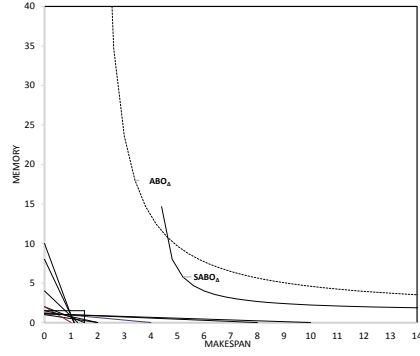
As a system designer, one always want to pick the algorithm and parameters that is best tradeoff between makespan and memory consumption. Depending on the guarantee, one should either pick ABO_{Δ} or $SABO_{\Delta}$ for scheduling tasks. For instance, if you want to guarantee a makespan less than 3 as in the case depicted in Figure 6b, you should use ABO_{Δ} . However if you want a better guarantee on memory, you should always use $SABO_{\Delta}$ for task scheduling.



(a) $m = 5, \alpha^2 = 2, \rho_1 = \rho_2 = 4/3$



(b) $m = 5, \alpha^2 = 3, \rho_1 = \rho_2 = 1$



(c) $m = 5, \alpha^2 = 3, \rho_1 = \rho_2 = 4/3$

Figure 6: Memory-Makespan graph for $SABO_{\Delta}$ and ABO_{Δ} . The bold lines represent impossibilities in tradeoff between guarantees.

CHAPTER 6: CONCLUSION AND FUTURE WORK

This thesis studies the effect on uncertainty in the processing time of tasks on scheduling for parallel and distributed machines. In particular, it investigates how allowing tasks to execute on different machines can help dealing with not knowing the processing time of tasks accurately. The thesis proposes three replication strategies, provides approximation algorithm in each case and a lower bound on the best achievable approximation in one of the case. Further to limit memory consumption the thesis presents two memory-aware bi-objective algorithms, one of which chooses only critical tasks to replicate and limits memory consumption.

The various strategies allow to trade the number of replication for a better guarantee. The results of these strategies show that a better guarantee can be achieved with fewer replication than that can be achieved by putting the data of a task on only one machine and even a small amount of replications can improve the guarantee significantly. These observations concludes that deploying the data on multiple machines can be an effective way of dealing with processing time uncertainties.

The bi-objective algorithms proposed in this thesis, schedule the memory intensive tasks and the processing time intensive tasks differently and optimizes both the objectives. One of the algorithms, chooses processing time intensive tasks to replicate and achieves better guarantee for higher values of α .

There are some open problems which can be explored further. Better lower bounds

might help understanding the problem better: clearly when α is low, the problem is no different than the offline problem, and when it is large, the problem converges to the non-clairvoyant online problem. Having a clearer idea of where the boundary is will certainly prove useful in understanding how much can be gained using data replication. Also, while replicating data using groups of processor proved effective, more general replication policies can certainly lead to better guarantees.

REFERENCES

- [1] ABAWAJY, J. Placement of file replicas in data grid environments. In *Computational Science - ICCS 2004*, M. Bubak, G. van Albada, P. Sloot, and J. Dongarra, Eds., vol. 3038 of *Lecture Notes in Computer Science*. 2004, pp. 66–73.
- [2] CANON, L.-C., AND JEANNOT, E. Evaluation and optimization of the robustness of dag schedules in heterogeneous environments. *IEEE Transactions on Parallel and Distributed Systems* 21, 4 (2010), 532–546.
- [3] CARDELLINI, V., I, R., COLAJANNI, M., AND YU, P. S. Dynamic load balancing on web-server systems. *IEEE Internet Computing* 3 (1999), 28–39.
- [4] CIRNE, W., BRASILEIRO, F., PARANHOS, D., GÓES, L. F. W., AND VOORSLUYS, W. On the efficacy, efficiency and emergent behavior of task replication in large distributed systems. *Parallel Computing* 33, 3 (2007), 213 – 234.
- [5] DANIELS, R. L., AND KOUVELIS, P. Robust Scheduling to Hedge Against Processing Time Uncertainty in Single-Stage Production. *Management Science* 41, 2 (Feb. 1995), 363–376.
- [6] DAVENPORT, A. J., GEFFLOT, C., AND BECK, J. C. Slack-based techniques for robust schedules. *Proceedings of the Sixth European Conference on Planning (ECP)* (2001).
- [7] DUTOT, P.-F., RZADCA, K., SAULE, E., AND TRYSTRAM, D. *Multi-objective scheduling*. Introduction to scheduling. Chapman and Hall/CRC Press, Nov. 2009, ch. 9. ISBN: 978-1420072730.
- [8] ERLEBACHER, G., SAULE, E., FLYER, N., AND BOLLIG, E. Acceleration of derivative calculations with application to radial basis function - finite-differences on the Intel MIC architecture. In *Proc. of International Conference on Supercomputing (ICS)* (2014).
- [9] GAREY, M. R., AND JOHNSON, D. S. *Computers and Intractability*. Freeman, San Francisco, 1979.
- [10] GATTO, M., AND WIDMAYER, P. On the robustness of graham’s algorithm for online scheduling. In *Proc of WADS* (2007).
- [11] GRAHAM, R. L. Bounds for certain multiprocessing anomalies. *Bell System Technical Journal* 45 (1966), 1563–1581.
- [12] GRAHAM, R. L. Bounds on multiprocessing timing anomalies. *SIAM Journal On Applied Mathematics* 17, 2 (1969), 416–429.

- [13] GUO, Z., FOX, G., AND ZHOU, M. Investigation of data locality in mapreduce. In *Proceedings of the 2012 12th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (Ccgriid 2012)* (Washington, DC, USA, 2012), CCGRID '12, IEEE Computer Society, pp. 419–426.
- [14] HOCHBAUM, D. S., AND SHMOYS, D. B. Using dual approximation algorithms for scheduling problems: Practical and theoretical results. *Journal of ACM* 34 (1987), 144–162.
- [15] KAVULYA, S., TAN, J., GANDHI, R., AND NARASIMHAN, P. An analysis of traces from a production MapReduce cluster. In *Proceedings of the 2010 10th IEEE/ACM International Conference on Cluster, Cloud and Grid Computing* (2010), CCGRID '10, pp. 94–103.
- [16] LI, Z., AND IERAPETRITOU, M. G. Process scheduling under uncertainty: Review and challenges. *Computers & Chemical Engineering* 32, 4-5 (2008), 715–727.
- [17] LIPTON, R., AND NAUGHTON, J. Query size estimation by adaptive sampling. *Journal of Computer and System Sciences* 51, 1 (1995), 18 – 25.
- [18] LUONG, D., DEOGUN, J., AND GODDARD, S. Feedback scheduling of real-time divisible loads in clusters. *SIGBED Rev.* 5, 2 (July 2008), 2:1–2:4.
- [19] PASTORELLI, M., BARBUZZI, A., CARRA, D., DELL'AMICO, M., AND MICHIARDI, P. Hfsp: Size-based scheduling for hadoop. In *Big Data, 2013 IEEE International Conference on* (Oct 2013), pp. 51–59.
- [20] RAHMAN, R., BARKER, K., AND ALHAJJ, R. Study of different replica placement and maintenance strategies in data grid. In *Cluster Computing and the Grid CCGRID* (2007), pp. 171–178.
- [21] RAO, B. T., AND REDDY, L. S. S. Survey on improved scheduling in hadoop mapreduce in cloud environments. *CoRR abs/1207.0780* (2012).
- [22] SAULE, E., DUTOT, P.-F., AND MOUNIE, G. Scheduling with storage constraints. *Parallel and Distributed Processing Symposium, International 0* (2008), 1–8.
- [23] T'KINDT, V., AND BILLAUT, J. *Multicriteria Scheduling*. Springer, 2007.
- [24] TSE, S. S. H. Online bounds on balancing two independent criteria with replication and reallocation. *IEEE Trans. Computers* 61, 11 (2012), 1601–1610.
- [25] VERMA, A., CHERKASOVA, L., AND CAMPBELL, R. H. ARIA: Automatic resource inference and allocation for MapReduce environments. In *Proceedings of the 8th ACM International Conference on Autonomic Computing* (New York, NY, USA, 2011), ICAC '11, ACM, pp. 235–244.

- [26] WANG, D., JOSHI, G., AND WORNELL, G. W. Efficient task replication for fast response times in parallel computation. In *proc. of SIGMETRICS* (2014).
- [27] WHITE, T. *Hadoop: The Definitive Guide*, 1st ed. O'Reilly Media, Inc., 2009.
- [28] WIERMAN, A., AND NUYENS, M. Scheduling despite inexact job-size information. In *Proc. of SIGMETRICS* (2008), pp. 25–36.
- [29] WILHELM, R., ENGBLOM, J., ERMEDAHL, A., HOLSTI, N., THESING, S., WHALLEY, D., BERNAT, G., FERDINAND, C., HECKMANN, R., MITRA, T., MUELLER, F., PUAUT, I., PUSCHNER, P., STASCHULAT, J., AND STENSTRÖM, P. The worst-case execution-time problem - overview of methods and survey of tools. *ACM Trans. Embed. Comput. Syst.* 7, 3 (May 2008), 36:1–36:53.
- [30] WOLF, J., RAJAN, D., HILDRUM, K., KHANDEKAR, R., KUMAR, V., PAREKH, S., WU, K.-L., AND BALMIN, A. FLEX: A slot allocation scheduling optimizer for MapReduce workloads. In *Proceedings of the ACM/IFIP/USENIX 11th International Conference on Middleware* (2010), Middleware '10, pp. 1–20.
- [31] YOU, Y., BADER, D. A., AND DEHNAVI, M. M. Designing a heuristic cross-architecture combination for breadth-first search. In *Proc. of the 43rd International Conference on Parallel Processing* (2014).
- [32] ZAHARIA, M., BORTHAKUR, D., SEN SARMA, J., ELMELEEGY, K., SHENKER, S., AND STOICA, I. Job scheduling for multi-user mapreduce clusters. Tech. Rep. UCB/EECS-2009-55, EECS Department, University of California, Berkeley, Apr 2009.
- [33] ZHOU, Z., SAULE, E., AKTULGA, H. M., YANG, C., NG, E. G., MARIS, P., VARY, J. P., AND ÇATALYÜREK, Ü. V. An out-of-core dataflow middleware to reduce the cost of large scale iterative solvers. In *2012 International Conference on Parallel Processing (ICPP) Workshops, Fifth International Workshop on Parallel Programming Models and Systems Software for High-End Computing (P2S2)* (Sept. 2012).
- [34] ZHOU, Z., SAULE, E., AKTULGA, H. M., YANG, C., NG, E. G., MARIS, P., VARY, J. P., AND ÇATALYÜREK, Ü. V. An out-of-core eigensolver on SSD-equipped clusters. In *Proc. of IEEE Cluster* (Sept. 2012).