## **Provably Efficient Online Non-clairvoyant Adaptive Scheduling**

Yuxiong He <sup>1,2</sup> Wen Jing Hsu <sup>1,2</sup> Charles E. Leiserson <sup>1,3</sup>

- <sup>1</sup> Singapore-MIT Alliance Program
- <sup>2</sup> Nanyang Technological University
- <sup>3</sup> Massachusetts Institute of Technology

#### **Abstract**

Scheduling competing jobs on multiprocessors has always been an important issue for parallel and distributed systems. The challenge is to ensure global, system-wide efficiency while offering a level of fairness to user jobs. Various degrees of successes have been achieved over years. However, few existing schemes address both efficiency and fairness over a wide range of work loads. Moreover, in order to obtain analytical results, most of them require prior information about jobs, which may be difficult to obtain in real applications.

This paper presents a novel adaptive scheduling algorithm GRAD that ensures fair allocation under all levels of workload, and it offers provable efficiency without requiring prior information of job's parallelism. Moreover, it provides effective control over the scheduling overhead and ensures efficient utilization of processors. Specifically, we show that GRAD is O(1)-competitive against an optimal offline scheduling algorithm with respect to both mean response time and makespan for batched jobs and non-batched jobs respectively.

To the best of our knowledge, GRAD is the first nonclairvoyant scheduling algorithm that offers such guarantees. We also believe that our new approach of resource requestallotment protocol deserves further exploration.

The simulation results show that, for non-batched jobs, the makespan produced by GRAD is no more than 1.39 times of the optimal on average. For batched jobs, the mean response time produced by GRAD is no more than 2.37 times of the optimal on average.

## 1. Introduction

Parallel computers are an expensive resource that often must be shared among a large community of users. One major issue of parallel job scheduling is how to efficiently share multiple processors among a number of competing jobs, while ensuring each job a required quality of services (see e.g. [15, 8, 7, 11, 9, 18, 24, 19, 25, 22, 31, 28, 23, 13, 14, 5]). Efficiency and fairness are two important performance measures,

where efficiency is often quantified in terms of makespan and mean response time. This paper presents a scheduling algorithm — GRAD, which offers provable efficiency in terms of makespan and mean response time by allotting each job a fair share of processor resources. Our algorithm is *non-clairvoyant* [11, 9, 18, 13], i.e. it assumes nothing about the release time, the execution time, and the parallelism profile of jobs. We model the execution of each job as a dynamically unfolding directed acyclic graph [8, 7].

A parallel job can be classified as adaptive or non-adaptive. An *adaptively parallel job* [27] may change its parallelism, and it allows the number of the allotted processors to vary during its execution. A job is *nonadaptive* if it runs on a fixed number of processors over its lifetime. With adaptivity, new jobs can enter the system by simply recruiting processors from the already executing jobs. Moreover, in order to improve the system utilization, schedulers can shift processors from jobs that do not require many processors to the jobs in need. However, since the parallelism of adaptively parallel jobs can change during the execution and the future parallelism is usually unknown, how a scheduler decides the processor allotments for jobs is a challenging problem. We developed GRAD to schedule such adaptively parallel jobs.

Scheduling parallel jobs on multiprocessors can be implemented in two levels [15]: a kernel-level job scheduler which allots processors to jobs, and a user-level thread scheduler which maps the threads of a given job to the allotted processors. The processor reallocation occurs periodically between scheduling quanta. The thread scheduler provides parallelism feedback to the job scheduler. The feedback is an estimation of the number of processors that its job can effectively use during the next quantum. The job scheduler follows some processor allocation policy to determine the allotment to the job. It may implement a policy that is either space-sharing, where jobs occupy disjoint processor resources, or time-sharing, where different jobs may share the same processor resources at different points in time. Once a job is allotted its processors, the allotment does not change within the quantum.

Previous theoretical work [11, 9, 18] of adaptive scheduling uses "dynamic equi-partitioning" (DEQ) [28, 22] as job scheduler. DEQ allots processors to jobs through space shar-

ing. However, it assumes that, at any time step, the number of jobs is always less than or equal to the total number of processors. Therefore, the scheduling schemes based on DEQ only work for this special case. Since multiprocessors are usually shared among many user applications, the case of having more numerous jobs than processors is very common. We will design a job scheduler that can handle both cases.

This paper presents a novel job scheduler RAD, which unifies the space-sharing job scheduling algorithm DEQ with the time-sharing round robin (RR) algorithm. When the total number of jobs is smaller than or equal to the total number of processors, it uses DEQ job scheduler, which allots each job with an equal number of processors unless the job requests for less. When the total number of jobs is greater than the total number of processors, RAD applies time-sharing round robin algorithm, which assigns each job with a single processor for an equal slice of scheduling time.

Based on the "equalized allotment" scheme for processor allocation, and by using the utilization in the past quantum as feedback, we show that GRAD is provably efficient. The performance is measured in terms of both makespan and mean response time. GRAD achieves O(1)-competitiveness with respect to makespan for job sets with arbitrary release times, and O(1)-competitiveness with respect to mean response time for batched job sets where all jobs are released simultaneously. Unlike many previous results, which either assume clairvoyance [24, 19, 25] or use instantaneous parallelism [11, 9], GRAD removes these restrictive assumptions. Moreover, because the quantum length can be adjusted to amortize the cost of context-switching during processor reallocation, GRAD provides effective control over the scheduling overhead and ensures efficient utilization of processors.

Our simulation results also suggest that GRAD performs well in practice. For job sets with arbitrary release time, their makespan scheduled by GRAD is no more than 1.39 times of the optimal on average (geometric mean). For batched job sets, their mean response time scheduled by GRAD is no more than 2.37 times of the optimal on average.

The remainder of this paper is organized as follows. Section 2 describes the job model, scheduling model, and objective functions. Section 3 describes the GRAD algorithm. Section 4 analyzes the competitiveness of GRAD with respect to makespan. Section 5 shows the competitiveness of GRAD with respect to mean response time for batched jobs, while its detailed analysis is presented in Appendix A. Section 6 presents the empirical results. Section 7 discusses the related work, and Section 8 gives some concluding remarks.

# 2. Scheduling and Analytical Model

Our scheduling input consists of a collection of independent jobs  $\mathcal{J} = \left\{J_1, J_2, \dots, J_{|\mathcal{J}|}\right\}$  to be scheduled on a collection of P identical processors. Time is broken into a sequence of equal-sized **scheduling quanta**  $1, 2, \dots$ , each of length L, where each quantum q includes the interval  $[L \cdot q, L \cdot q + 1, \dots, L(q+1) - 1]$  of time steps. The quantum length L is

a system configuration parameter chosen to be long enough to amortize scheduling overheads. In this section, we formalize the job model, define the scheduling model, and present the optimization criteria of makespan and mean response time.

We model the execution of a multithreaded job  $J_i$  as a dynamically unfolding directed acyclic graph ( $\pmb{DAG}$ , for short). Each vertex of the DAG represents a unit-time instruction. The  $\pmb{work}$   $T_1(J_i)$  of the job  $J_i$  corresponds to the total number of vertices in the dag. Each edge represents a dependency between the two vertices. The  $\pmb{span}$   $T_\infty(J_i)$  corresponds to the number of nodes on the longest chain of the precedence dependencies. The  $\pmb{release}$  time  $r(J_i)$  of the job  $J_i$  is the time at which  $J_i$  becomes first available for processing. Each job is handled by a dedicated thread scheduler, which operates in an online manner, oblivious to the future characteristics of the dynamically unfolding DAG.

The job scheduler and the thread schedulers interact as follows. The job scheduler may reallocate processors between scheduling quanta. Between quantum q-1 and quantum q, the thread scheduler of a given job  $J_i$  determines the job's **desire**  $d(J_i,q)$ , which is the number of processors  $J_i$  wants for quantum q. Based on the desire of all running jobs, the job scheduler follows its processor-allocation policy to determine the **allotment**  $a(J_i,q)$  of the job with the constraint that  $a(J_i,q) \leq d(J_i,q)$ . Once a job is allotted its processors, the allotment does not change during the quantum.

Our scheduler uses makespan and mean response time as the performance measurement.

**Definition 1** The *makespan* of a given job set  $\mathcal{J}$  is the time taken to complete all the jobs in  $\mathcal{J}$ , i.e.  $T(\mathcal{J}) = \max_{J_i \in \mathcal{J}} T(J_i)$ , where  $T(J_i)$  denotes the completion time of job  $J_i$ .

**Definition 2** The *response time* of a job  $J_i$  is  $T(J_i) - r(J_i)$ , which is the duration between its release time  $r(J_i)$  and the completion time  $T(J_i)$ . The *total response time* of a job set  $\mathcal{J}$  is given by  $R(\mathcal{J}) = \sum_{J_i \in \mathcal{J}} (T(J_i) - r(J_i))$  and the *mean response time* is  $\overline{R}(\mathcal{J}) = R(\mathcal{J})/|\mathcal{J}|$ .

The competitive analysis of an online scheduling algorithm is to compare the algorithm against an optimal clairvoyant algorithm. Let  $T^*(\mathcal{J})$  denote the makespan of an arbitrary jobset  $\mathcal{J}$  scheduled by an optimal scheduler, and  $T(\mathcal{J})$  denote the makespan produced by an algorithm A for the job set  $\mathcal{J}$ . A deterministic algorithm A is said to be c-competitive if there exists a constant b such that  $T(\mathcal{J}) \leq c \cdot T^*(\mathcal{J}) + b$  holds for the schedule of any job set. We will show that our algorithm is c-competitive in terms of the makespan, where c is a small constant. Similarly, for the mean response time, we will show that our algorithm is also constant-competitive for any batched jobs.

# 3. Algorithms

This section presents the job scheduler – RAD, and overviews the thread scheduler – A-GREEDY [1].

## **RAD Job Scheduler**

The job scheduler RAD unifies the space-sharing job scheduling algorithm DEQ [28, 22] with the time-sharing RR algorithm. When the number of jobs is greater than the number of processors, GRAD schedules the jobs in a batched, roundrobin fashion, which allocates one processor to each job with an equal share of time. When the number of jobs is not more than the number of processors, GRAD uses DEQ as the job scheduler. DEQ gives each job an equal share of spatial allotments unless the job requests for less.

When a batch of jobs are scheduled in the round-robin fashion, RAD maintains a queue of jobs. At the beginning of each quantum, if there are more than P jobs, it pops P jobs from the top of the queue, and allots one processor to each of them during the quantum. At the end of the quantum, RAD pushes the P jobs back to the bottom of the queue if they are uncompleted. The new jobs can be put into the queue once they are released

DEQ attempts to give each job a fair share of processors. If a job requires less than its fair share, however, DEQ distributes the extra processors to the other jobs. More precisely, upon receiving the desires  $\{d(J_i,q)\}$  from the thread schedulers of all jobs  $J_i \in \mathcal{J}$ , DEQ executes the following processor-allocation algorithm:

- 1. Set  $n = |\mathcal{J}|$ . If n = 0, return.
- 2. If the desire of every job  $J_i \in \mathcal{J}$  satisfies  $d(J_i, q) \ge P/n$ , assign each job  $a(J_i, q) = P/n$  processors.
- 3. Otherwise, let  $\mathcal{J}'=\{J_i\in\mathcal{J}:d(J_i,q)< P/n\}$ . Assign  $a(J_i,q)=d(J_i,q)$  processors to each  $J_i\in\mathcal{J}'$ . Update  $\mathcal{J}=\mathcal{J}-\mathcal{J}'$ , and  $P=P-\sum_{J_i\in\mathcal{J}'}d(J_i,q)$ . Go to Step 1.

Note that, at any quantum where the number of jobs is equal to the number of processors, DEQ and RR give exactly the same processor allotment, and allocate each of P jobs with one processor.

## **Adaptive Greedy Thread Scheduler**

A-GREEDY [1] is an adaptive greedy thread scheduler with parallelism feedback. Between quanta, it estimates its job's desire, and requests processors from the job scheduler. During the quantum, it schedules the ready threads of the job onto the allotted processors greedily [16, 8]. If there are more than  $a\left(J_{i},q\right)$  ready threads, A-GREEDY schedules any  $a\left(J_{i},q\right)$  of them. Otherwise, it schedules all of them.

A-GREEDY's desire-estimation algorithm is parameterized in terms of a *utilization parameter*  $\delta > 0$  and a *responsiveness parameter*  $\rho > 1$ , both of which can be adjusted for different levels of guarantees for waste and completion time.

Before each quantum, A-GREEDY for a job  $J_i \in \mathcal{J}$  provides parallelism feedback to the job scheduler based on the  $J_i$ 's history of utilization in the previous quantum. A-GREEDY classifies quanta as "satisfied" versus "deprived" and "efficient" versus "inefficient." A quantum q is **satisfied** if  $a(J_i,q) = d(J_i,q)$ , in which case  $J_i$ 's allotment is equal to

its desire. Otherwise, the quantum is *deprived*. <sup>1</sup> The quantum q is *efficient* if A-GREEDY utilizes no less than a  $\delta$  fraction of the total allotted processor cycles during the quantum, where  $\delta$  is the utilization parameter. Otherwise, the quantum is *inefficient*. Under the four-way classification, however, A-GREEDY only uses three: inefficient, efficient-and-satisfied, and efficient-and-deprived.

Using this three-way classification and the job's desire for the previous quantum, A-GREEDY computes the desire for the next quantum as follows:

- If quantum q-1 was inefficient, decrease the desire, setting  $d(J_i,q)=d(J_i,q-1)/\rho$ , where  $\rho$  is the responsiveness parameter.
- If quantum q-1 was efficient-and-satisfied, increase the desire, setting  $d(J_i,q)=\rho d(J_i,q-1)$ .
- If quantum q-1 was efficient-and-deprived, keep desire unchanged, setting  $d(J_i,q)=d(J_i,q-1)$ .

## 4. Makespan

This section shows that GRAD is c-competitive with respect to makespan, where c denotes a constant. The exact value of c is related to the choice of A-GREEDY's utilization and responsiveness parameter, as will be explained shortly.

We first review the lower bounds of makespan. Given a job set  $\mathcal{J}$  and P processors, lower bounds on the makespan of any job scheduler can be obtained based on release time, work, and span. Recall that, for a job  $J_i \in \mathcal{J}$ , the quantities  $r(J_i)$ ,  $T_1(J_i)$ , and  $T_\infty(J_i)$  represent the release time, work, and span of  $J_i$ , respectively. Let  $T^*(\mathcal{J})$  denote the makespan produced by an optimal scheduler on a job set  $\mathcal{J}$  on P processors. Let  $T_1(\mathcal{J}) = \sum_{J_i \in \mathcal{J}} T_1(J_i)$  denote the total work of the job set. The following two inequalities give two lower bounds on the makespan [9]:

$$T^*(\mathcal{J}) \geq \max_{J_i \in \mathcal{J}} \{ r(J_i) + T_{\infty}(J_i) \} , \qquad (1)$$

$$T^*(\mathcal{J}) \geq T_1(\mathcal{J})/P. \tag{2}$$

To facilitate the analysis, we state a lemma from [1] that bounds the satisfied steps and the waste of a single job scheduled by A-GREEDY. Recall that, the parameter  $\rho>1$  denotes A-GREEDY's responsiveness parameter,  $\delta>0$  its utilization parameter, and L the quantum length.

**Lemma 1** [1] For a job  $J_i$  with work  $T_1(J_i)$  and span  $T_\infty(J_i)$  on a machine with P processors, A-GREEDY produces at most  $2T_\infty(J_i)/(1-\delta) + L\log_\rho P + L$  satisfied steps, and it wastes at most  $(1+\rho-\delta)T_1(J_i)/\delta$  processor cycles in the course of the computation.

The following theorem analyzes the makespan of a job set  $\mathcal{J}$  scheduled by GRAD.

 $<sup>^1</sup>$ We can extend the classification of "satisfied" versus "deprived" from quanta to time steps. A job  $J_i$  is *satisfied* (or *deprived*) at step  $t \in [L \cdot q, L \cdot q+1, \ldots, L(q+1)-1]$  if  $J_i$  is satisfied (resp. deprived) at the quantum q.

**Theorem 2** Let  $\rho$  denote A-GREEDY's responsiveness parameter,  $\delta$  its utilization parameter, and L the quantum length. Then, GRAD completes a job set  $\mathcal{J}$  on P processors in

$$T(\mathcal{J}) \le \frac{\rho+1}{\delta} \frac{T_1(\mathcal{J})}{P} + \frac{2}{1-\delta} \max_{J_i \in \mathcal{J}} \left\{ T_{\infty}(J_i) + r(J_i) \right\} + L \log_{\rho} P + 2L$$
(3)

time steps.

*Proof.* Suppose job  $J_k$  is the last job completed among the jobs in  $\mathcal{J}$ . Let  $S(J_k)$  denote the set of satisfied steps for  $J_k$ , and  $D(J_k)$  denote its set of deprived steps. The job  $J_k$  is scheduled to start its execution at the beginning of the quantum q where  $Lq < r(J_k) \leq L(q+1)$ , which is the quantum immediately after  $J_k$ 's release. Therefore, we have  $T(\mathcal{J}) \leq r(J_k) + L + |S(J_k)| + |D(J_k)|$ . We now bound  $|S(J_k)|$  and  $|D(J_k)|$  respectively.

From Lemma 1, we know that the number of satisfied steps attributed to  $J_k$  is at most  $|S(J_k)| \leq 2T_{\infty}(J_k)/(1-\delta) + L\log_{\rho}P + L$ .

We now bound the total number of deprived steps  $D(J_k)$  of job  $J_k$ . For each step  $t \in D(J_k)$ , GRAD applies either DEQ or RR as job scheduler. RR always allots all processors to jobs. By definition, DEQ must have allotted all processors to jobs whenever  $J_k$  is deprived. Thus, the total allotment on such a step t is always equal to the total number of processors P. Moreover, the total allotment of  $\mathcal J$  over  $J_k$ 's deprived steps  $D(J_k)$  is  $a(\mathcal J, D(J_k)) = \sum_{t \in D(J_k)} \sum_{J_i \in \mathcal J} a(J_i, t) = P |D(J_k)|$ . Since any allotted processor is either working productively or wasted, the total allotment for any job  $J_i$  is bounded by the sum of its total work  $T_1(J_i)$  and total waste  $w(J_i)$ . By Lemma 1, the waste for the job  $J_i$  is at most  $(\rho - \delta + 1)/\delta$  times of its work. Thus, the total number of allotted processor cycles for job  $J_i$  is at most  $T_1(J_i) + w(J_i) \leq (\rho + 1)T_1(J_i)/\delta$ . The total number of allotted processor cycles for all jobs is at most  $\sum_{J_i \in \mathcal J} (\rho + 1)T_1(J_i)/\delta = ((\rho + 1)/\delta)T_1(\mathcal J)$ . Given  $a(\mathcal J, D(J_k)) \leq ((\rho + 1)/\delta)T_1(\mathcal J)$  and  $a(\mathcal J, D(J_k)) = P|D(J_k)|$ , we have  $|D(J_k)| \leq \frac{\rho+1}{\delta} \frac{T_1(\mathcal J)}{P}$ .

Therefore, we can get

$$T(\mathcal{J}) < r(J_k) + L + |D(J_k)| + |S(J_k)|$$

$$\leq \frac{\rho + 1}{\delta} \frac{T_1(\mathcal{J})}{P} + \frac{2}{1 - \delta} \max_{J_i \in \mathcal{J}} (T_{\infty}(J_i) + r(J_i))$$

$$+ L \log_{\rho} P + 2L.$$

Since both  $T_1(\mathcal{J})/P$  and  $\max_{J_i \in \mathcal{J}} \{T_\infty(J_i) + r(J_i)\}$  are lower bounds of  $T^*(\mathcal{J})$ , we obtain the following corollary.

**Corollary 3** GRAD completes a job set  $\mathcal{J}$  in

$$T(\mathcal{J}) \le \left(\frac{\rho+1}{\delta} + \frac{2}{1-\delta}\right) T^*(\mathcal{J}) + L \log_{\rho} P + 2L$$

time steps, where  $T^*(\mathcal{J})$  denotes the makespan of  $\mathcal{J}$  produced by an optimal clairvoyant scheduler.

Since both the quantum length L and the processor number P are independent variables with respect to any job set  $\mathcal{J}$ , Corollary 3 shows that GRAD is O(1)-competitive with respect to makespan.

## 5. Mean Response Time

Mean response time is an important measure for multiuser environments where we desire as many users as possible to get fast response from the system. In this section, we first introduce the lower bounds. Then, we show that GRAD is O(1)-competitive for batched jobs with respect to the mean response time.

#### **Lower Bounds and Preliminaries**

We first introduce some definitions.

**Definition 3** Given a finite list  $\mathcal{A} = \langle \alpha_i \rangle$  of  $n = |\mathcal{A}|$  integers, define  $f: \{1,2,\ldots,n\} \to \{1,2,\ldots,n\}$  to be a permutation satisfying  $\alpha_{f(1)} \leq \alpha_{f(2)} \leq \cdots \leq \alpha_{f(n)}$ . The **squashed sum** of  $\mathcal{A}$  is defined as

$$\operatorname{sq-sum}(\mathcal{A}) = \sum_{i=1}^{n} (n-i+1)\alpha_{f(i)}.$$

The *squashed work area* of a job set  $\mathcal{J}$  on a set of P processors is

$$swa(\mathcal{J}) = \frac{1}{P} sq\text{-sum}(\langle T_1(J_i) \rangle),$$

where  $T_1(J_i)$  is the work of job  $J_i \in \mathcal{J}$ . The **aggregate span** of  $\mathcal{J}$  is

$$T_{\infty}(\mathcal{J}) = \sum_{J_i \in \mathcal{J}} T_{\infty}(J_i) ,$$

where  $T_{\infty}(J_i)$  is the span of job  $J_i \in \mathcal{J}$ .

The research in [29, 30, 11] establishes two lower bounds for the mean response time:

$$\overline{\mathbb{R}^*}(\mathcal{J}) \geq T_{\infty}(\mathcal{J})/|\mathcal{J}|,$$
 (4)

$$\overline{\mathbf{R}^*}(\mathcal{J}) \geq \operatorname{swa}(\mathcal{J})/|\mathcal{J}|,$$
 (5)

where  $\overline{R}^*(\mathcal{J})$  denotes the mean response time of  $\mathcal{J}$  scheduled by an optimal clairvoyant scheduler. Both the aggregate span  $T_\infty(\mathcal{J})$  and the squashed work area  $\mathrm{swa}\,(\mathcal{J})$  are lower bounds of the total response time  $\mathrm{R}^*(\mathcal{J})$  under an optimal clairvoyant scheduler.

#### **Analysis**

The proof is divided into two parts. In the first part where  $|\mathcal{J}| \leq P$ , GRAD always uses DEQ as job scheduler. In this case, we apply the result in [18], and show that GRAD is O(1)-competitive. In the second part where  $|\mathcal{J}| > P$ , GRAD uses both RR and DEQ. Since we consider batched jobs, the number of incomplete jobs decreases monotonically. When the number of incomplete jobs drops to P, GRAD switches its job scheduler from RR to DEQ. Therefore, we prove the second case based on the properties of round robin scheduling and the results of the first case. The following theorem shows

the total response time bound for the batched job sets scheduled by GRAD. Please refer to Appendix A for the complete proof.

**Theorem 4** Let  $\rho$  be A-GREEDY's responsiveness parameter,  $\delta$  its utilization parameter, and L the quantum length. The total response time  $R(\mathcal{J})$  of a job set  $\mathcal{J}$  produced by GRAD is at most

$$R(\mathcal{J}) = \left(2 - \frac{2}{|\mathcal{J}| + 1}\right) \left(\frac{\rho + 1}{\delta} \operatorname{swa}(\mathcal{J}) + \frac{2}{1 - \delta} T_{\infty}(\mathcal{J})\right) + O(|\mathcal{J}| L \log_{\rho} P),$$
(6)

where  $swa(\mathcal{J})$  denotes the squashed work area of  $\mathcal{J}$ , and  $T_{\infty}(\mathcal{J})$  denotes the aggregate span of  $\mathcal{J}$ .

Since both swa  $(\mathcal{J})/|\mathcal{J}|$  and  $T_{\infty}(\mathcal{J})/|\mathcal{J}|$  are lower bounds on  $\overline{\mathbb{R}}(\mathcal{J})$ , we obtain the following corollary. It shows that GRAD is O(1)-competitive with respect to mean response time for batched jobs.

**Corollary 5** *The mean response time*  $\overline{\mathbb{R}}(\mathcal{J})$  *of a batched job set*  $\mathcal{J}$  *produced by* GRAD *satisfies* 

$$\overline{\mathbf{R}}(\mathcal{J}) = \left(2 - \frac{2}{|\mathcal{J}| + 1}\right) \left(\frac{\rho + 1}{\delta} + \frac{2}{1 - \delta}\right) \overline{\mathbf{R}^*}(\mathcal{J}) + \ O(L \log_{\rho} P) \ ,$$

where  $\overline{\mathbb{R}^*}(\mathcal{J})$  denotes the mean response time of  $\mathcal{J}$  scheduled by an optimal clairvoyant scheduler.

## 6. Experimental Results

To evaluate the performance of GRAD, we conducted four sets of experiments, which are summarized below.

- The **makespan experiments** compares the makespan produced by GRAD against the theoretical lower bound for over 10000 runs of job sets.
- The mean response time experiments investigate how GRAD performs with respect to mean response time for over 8000 batched job sets.
- The **load experiments** investigate how the system load affects the performance of GRAD.
- The proactive RAD experiments compare the performance of RAD against its variation proactive RAD.
   The proactive RAD always allots all processors to jobs even if the overall desire is less than the total number of processors.

## 6.1. Simulation Setup

To study GRAD, we build a Java-based discrete-time simulator using DESMO-J [12]. Our benchmark application is the Fork-Join jobs, where each job alternates between a *serial phase* and a *parallel phase*. Fork-Join jobs arise naturally in jobs that exhibit "data parallelism", and apply the same computation to a number of different data points. The repeated fork-join cycle in the job reflects the iterative nature of these computations.

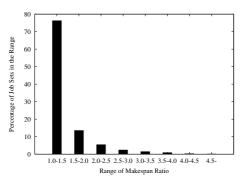
GRAD requires some parameters as input. We set the responsiveness parameter to be  $\rho = 2.0$ , and the utilization parameter  $\delta = 0.8$  unless otherwise specified. GRAD is designed for moderate-scale and large-scale multiprocessors, and we set the number of processors to be P = 128. The quantum length L represents the time between successive reallocations of processors by the job scheduler, and is selected to amortize the overheads due to the communication between the job scheduler and the thread scheduler, and the reallocation of processors. In conventional computer systems, a scheduling quantum is typically between 10 and 20 milliseconds. The execution time of a task is decided by the granularity of the job. If a task takes approximately 0.5 to 5 microseconds, then the quantum length L should be set to values between  $10^3$  and  $10^5$  time steps. Our theoretical bounds indicate that as long as  $T_{\infty} \gg L \log P$ , the length of L should have little effect on our results. In our experiments, we set L = 1000.

## 6.2. Makespan Experiments

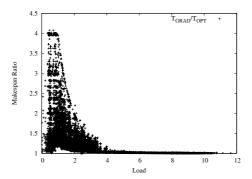
The competitive ratio of makespan derived in Section 4, though asymptotically strong, has a relatively large constant multiplier. The makespan experiments were designed to evaluate the constants that would occur in practice and compare GRAD to an optimal scheduler. The experiments are conducted on more than 10,000 runs of job sets using many combinations of jobs and different loads.

Figure 1 shows how GRAD performs compared to an optimal scheduler. The makespan of a job set  $\mathcal{J}$  has two lower bounds  $\max_{J_i \in \mathcal{J}}(r(J_i) + T_\infty(J_i))$  and  $T_1(\mathcal{J})/P$ . The makespan produced by an optimal scheduler is lower-bounded by the larger of these two values. The makespan ratio in Figure 1 is defined as the makespan of a job set scheduled by GRAD divided by the theoretical lower bounds. Its X-axis represents the range of the makespan ratio, while the histogram shows the percentage of the job sets whose makespan ratio falls into the range. Among more than 10,000 runs, 76.19% of them use less than 1.5 times of the theoretical lower bound, 89.70% use less than 2.0 times, and none uses more than 4.5 times. The average makepsan ratio is 1.39, which suggests that, in practice, GRAD has a small competitive ratio with respect to the makespan.

We now interpret the relation between the theoretical bounds and experimental results as follows. When  $\rho=2$  and  $\delta=0.8$ , from Theorem 2, GRAD is 13.75-competitive in the worst case. However, we anticipate that GRAD's makespan ratio would be small in practical settings, especially when the jobs have total work much great than the span and with the machine moderately- or highly- loaded. In this case, the term on  $T_1(\mathcal{J})/P$  in Inequality (3) of Theorem 2 is much larger than the term  $\max_{J_i \in \mathcal{J}} \{T_\infty(i) + r(i)\}$ , i.e. the term  $T_1(\mathcal{J})/P$  generally dominates the makespan bound. The proof of Theorem 2 calculates the coefficient of  $T_1(\mathcal{J})/P$  as the ratio of the total allotment (total work plus total waste) versus the total work. When the job scheduler is RAD, which is not a true adversary, our simulation results indicate that the ratio of the waste versus the total work is only about 1/10 of the total



**Figure 1:** Comparing the makespan of GRAD with the theoretical lower bound for job sets with arbitrary job release time.



**Figure 3:** Comparing GRAD against the theoretical lower bound for makespan with varying load.

work. Thus, the coefficient of  $T_1(\mathcal{J})/P$  in Inequality (3) is about 1.1. It explains why the makespan produced by GRAD is less than 2 times of the lower bound on average as shown in Figure 1.

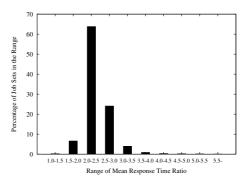
#### 6.3. Mean Response Time Experiments

This set of experiments is designed to evaluate the mean response time of the batch job sets scheduled by GRAD. Figure 2 shows the distribution of the mean response time normalized w.r.t. the larger of the two lower bounds – the squashed work bound swa  $(\mathcal{J})/|\mathcal{J}|$  and the aggregated critical path bound  $T_{\infty}(\mathcal{J})/|\mathcal{J}|$ . The histogram in Figure 2 shows that, among more than 8,000 runs, 94.65% of them use less than 3 times of the theoretical lower bound, and none of them uses more than 5.5 times. The average mean response time ratio is 2.37. Similar to the discussion in Section 6.2, we can relate the theoretical bounds for mean response time to the experimental results.

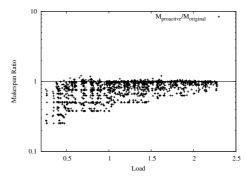
## **6.4. Load Experiments**

This set of experiments is designed to investigate how the load affects the performance of GRAD. The *load* of a job set  $\mathcal{J}$  on a machine with P processors indicates how heavily the jobs compete for processors on the machine, which is calculated as follows

$$load = \frac{T_1(\mathcal{J})}{P \cdot \left(\max_{J_i \in \mathcal{J}} r(J_i) - \min_{J_i \in \mathcal{J}} r(J_i) + T_{\infty}(\mathcal{J}) / |\mathcal{J}|\right)}$$



**Figure 2:** Comparing the mean response time of GRAD with the theoretical lower bound for batched job sets.



**Figure 4:** Comparing the proactive RAD against the original for makespan with varying load.

For a batched job set, the load is just the average parallelism of the set divided by the total number of processors.

Figure 3 shows how GRAD performs against the theoretical lower bound with respect to makespan by varying system load. The makespan ratio in this figure is defined as the makespan of a job set scheduled by GRAD divided by the larger of the two lower bounds. Each data point represents the makespan ratio of a job set. The testing results suggest that the makespan ratio becomes smaller when the load gets heavier. Specifically, the makespan generated by GRAD is very close to the lower bound when the load is greater than 4; it never exceeds 1.5 times of the makespan produced when the system load is greater than 3. However, when the load is less than 2, the makespan ratio spreads in the range from 1 to 4. To improve the performance of GRAD under light load, we will explore such a variation of the job scheduler RAD (called proactive RAD) in the next section.

#### **6.5. Proactive RAD Experiments**

Proactive RAD always allocates all processors to jobs even if the total requests are less than the total number of processors. At a quantum q, when the total requests  $d(\mathcal{J},q) = \sum_{J_i \in J} d(J_i,q)$  are greater than or equal to the total number P of processors, the proactive RAD works exactly the same as the original one. However, if  $d(\mathcal{J},q) < P$ , the proactive RAD evenly allots the remaining  $P - d(\mathcal{J},q)$  processors to all the jobs.

Figure 4 shows the makespan ratio of proactive RAD against its original algorithm by varying system load. Each data point in the figure represents a job set's makespan ratio, defined as the makespan produced by the proactive RAD divided by that of the original. We can see that the makespan ratio is less than 1 for most of the runs, indicating that the proactive RAD out-performs the original one in most of these job sets. Moreover, the difference between them becomes more pronounced under light load, and diminishes with the increase of the system load. The reason is that the proactive RAD generally allocates more processors to jobs, especially when the load is light. The increased allotment allows faster execution of jobs which shortens the makespan of the job set. Figure 4 gives evidences that the proactive RAD improves the performance of our scheduling algorithm under light load.

#### 7. Related Work

Adaptive parallel job scheduling has been studied both empirically [22, 31, 28, 21] and theoretically [17, 10, 23, 13, 14, 5]. McCann, Vaswani, and Zahorjan [22] introduce the notion of dynamic equipartitioning (DEQ), which gives each job a fair allotment of processors based on the job's request, while allowing processors that cannot be used by a job to be reallocated to the other jobs. Brecht, Deng, and Gu [9] prove that DEQ with instantaneous parallelism as feedback is 2-competitive with respect to the makespan. Later, Deng and Dymond [11] prove that DEQ with instantaneous parallelism is also 4-competitive for batched jobs with respect to the mean response time.

Even though using instantaneous parallelism as feedback is intuitive, it can either cause gross misallocation of processor resources [26] or introduce significant scheduling overhead. For example, the parallelism of a job may change substantially during a scheduling quantum, alternating between parallel and serial phases. Depending on which phase is currently active, the sampling of instantaneous parallelism may lead the task scheduler to request either too many or too few processors. Consequently, the job may either waste processor cycles or take too long to complete. On the other hand, if the quantum length is set to be small enough to capture frequent changes in instantaneous parallelism, the proportion of time spent reallotting processors among the jobs increases, resulting in a high scheduling overhead.

Our previous work in [18] presents a two-level adaptive scheduler AGDEQ, which uses DEQ as the job scheduler, and A-GREEDY as the thread scheduler. Instead of using instantaneous parallelism, AGDEQ uses the job's utilization in the past as feedback. AGDEQ is O(1)-competitive for makespan, and in a batched setting, O(1)-competitive for mean response time. However, as with other prior work [9, 11] that uses DEQ as the job scheduler, AGDEQ can only be applied to the case where the total number of jobs in the job set is less than or equal to the number of processors.

## 8. Conclusions

We have presented a non-clairvoyant adaptive scheduling algorithm GRAD that ensures provable efficiency, fairness and minimal overhead.

The request-allotment protocol in GRAD can be applied to other thread schedulers. For example, A-GREEDY is suitable for scheduling jobs in more centralized setting such as data parallel applications. However, for applications with more distributed settings, we can use A-STEAL [2, 3] as thread scheduler, which is a distributed adaptive thread scheduler using randomized work stealing. Analogously, one can develop a two-level scheduler by applying the request-allotment protocol in GRAD, and application-specific thread schedulers. Such a two-level scheduler may provide both system-wide performance guarantees such as minimal makespan and mean response time, and optimization of individual applications.

With respect to mean response time, our schedulers are O(1)-competitiveness only for batched parallel jobs. Some researchers [20, 4, 6] have studied the online non-clairvoyant scheduling of serial jobs with arbitrary release time. It remains an interesting open problem to develop non-clairvoyant schedulers that minimizes mean response time for parallel jobs with arbitrary release times

### References

- [1] Kunal Agrawal, Yuxiong He, Wen Jing Hsu, and Charles E. Leiserson. Adaptive task scheduling with parallelism feedback. In *PPoPP*, pages 100 – 109, New York City, NY, USA, 2006.
- [2] Kunal Agrawal, Yuxiong He, and Charles E. Leiserson. An empirical evaluation of work stealing with parallelism feedback. In *ICDCS*, pages 19 29, Lisboa, Portugal, 2006.
- [3] Kunal Agrawal, Yuxiong He, and Charles E. Leiserson. Work stealing with parallelism feedback. Unpublished manuscripts, 2006.
- [4] Nir Avrahami and Yossi Azar. Minimizing total flow time and total completion time with immediate dispatching. In SPAA, pages 11–18, New York, NY, USA, 2003.
- [5] Nikhil Bansal, Kedar Dhamdhere, Jochen Konemann, and Amitabh Sinha. Non-clairvoyant scheduling for minimizing mean slowdown. *Algorithmica*, 40(4):305–318, 2004.
- [6] Luca Becchetti and Stefano Leonardi. Nonclairvoyant scheduling to minimize the total flow time on single and parallel machines. *Journal of the ACM*, 51(4):517–539, 2004.
- [7] Guy E. Blelloch and John Greiner. A provable time and space efficient implementation of NESL. In *ICFP*, pages 213–225, Philadelphia, Pennsylvania, 1996.
- [8] Robert D. Blumofe and Charles E. Leiserson. Space-efficient scheduling of multithreaded computations. SIAM Journal on Computing, 27(1):202–229, 1998.
- [9] T. Brecht, Xiaotie Deng, and Nian Gu. Competitive dynamic multiprocessor allocation for parallel applications. In *Parallel* and *Distributed Processing*, pages 448 – 455, San Antonio, TX, 1995.
- [10] Xiaotie Deng and Patrick Dymond. On multiprocessor system scheduling. In SPAA, pages 82–88, Padua, Italy, 1996.
- [11] Xiaotie Deng, Nian Gu, Tim Brecht, and KaiCheng Lu. Preemptive scheduling of parallel jobs on multiprocessors. In SODA, pages 159–167, Philadelphia, PA, USA, 1996.

- [12] DESMO-J: A framework for discrete-event modelling and simulation. http://asi-www.informatik.uni-hamburg.de/desmoj/.
- [13] Jeff Edmonds. Scheduling in the dark. In STOC, pages 179– 188, Atlanta, Georgia, United States, 1999.
- [14] Jeff Edmonds, Donald D. Chinn, Timothy Brecht, and Xiaotie Deng. Non-clairvoyant multiprocessor scheduling of jobs with changing execution characteristics. *Journal of Scheduling*, 6(3):231–250, 2003.
- [15] Dror G. Feitelson. Job scheduling in multiprogrammed parallel systems (extended version). Technical report, IBM Research Report RC 19790 (87657) 2nd Revision, 1997.
- [16] R. L. Graham. Bounds on multiprocessing anomalies. SIAM Journal on Applied Mathematics, pages 17(2):416–429, 1969.
- [17] Nian Gu. Competitive analysis of dynamic processor allocation strategies. Master's thesis, York University, 1995.
- [18] Yuxiong He, Wen Jing Hsu, and Charles E. Leiserson. Provably efficient two-level adaptive scheduling. In *JSSPP*, Saint-Malo, France, 2006.
- [19] Klaus Jansen and Hu Zhang. Scheduling malleable tasks with precedence constraints. In SPAA, pages 86–95, New York, NY, USA, 2005.
- [20] Bala Kalyanasundaram and Kirk R. Pruhs. Minimizing flow time nonclairvoyantly. *Journal of the ACM*, 50(4):551–567, 2003
- [21] Scott T. Leutenegger and Mary K. Vernon. The performance of multiprogrammed multiprocessor scheduling policies. In SIG-METRICS, pages 226–236, Boulder, Colorado, United States, 1990.
- [22] Cathy McCann, Raj Vaswani, and John Zahorjan. A dynamic processor allocation policy for multiprogrammed sharedmemory multiprocessors. ACM Transactions on Computer Systems, 11(2):146–178, 1993.
- [23] Rajeev Motwani, Steven Phillips, and Eric Torng. Non-clairvoyant scheduling. In SODA, pages 422–431, Austin, Texas, United States, 1993.
- [24] Gregory Mounie, Christophe Rapine, and Dennis Trystram. Efficient approximation algorithms for scheduling malleable tasks. In *SPAA*, pages 23–32, New York, NY, USA, 1999.
- [25] Uwe Schwiegelshohn, Walter Ludwig, Joel L. Wolf, John Turek, and Philip S. Yu. Smart smart bounds for weighted response time scheduling. SIAM Journal of Computing, 28(1):237–253, 1998.
- [26] Siddhartha Sen. Dynamic processor allocation for adaptively parallel jobs. Master's thesis, Massachusetts Institute of technology, 2004.
- [27] B. Song. Scheduling adaptively parallel jobs. Master's thesis, Massachusetts Institute of Technology, 1998.
- [28] Andrew Tucker and Anoop Gupta. Process control and scheduling issues for multiprogrammed shared-memory multiprocessors. In SOSP, pages 159–166, New York, NY, USA, 1989.
- [29] John Turek, Walter Ludwig, Joel L. Wolf, Lisa Fleischer, Prasoon Tiwari, Jason Glasgow, Uwe Schwiegelshohn, and Philip S. Yu. Scheduling parallelizable tasks to minimize average response time. In SPAA, pages 200–209, Cape May, New Jersey, United States, 1994.
- [30] John Turek, Uwe Schwiegelshohn, Joel L. Wolf, and Philip S. Yu. Scheduling parallel tasks to minimize average response time. In SODA, pages 112–121, Philadelphia, PA, USA, 1994.
- [31] K. K. Yue and D. J. Lilja. Implementing a dynamic processor allocation policy for multiprogrammed parallel applications in the Solaris<sup>TM</sup> operating system. *Concurrency and Computation-Practice and Experience*, 13(6):449–464, 2001.

## Appendix A. Proof of Theorem 4

The proof is divided into two cases — when  $|\mathcal{J}| \leq P$  and when  $|\mathcal{J}| > P$ .

#### Case 1: when $|\mathcal{J}| \leq P$

For the first case where  $|\mathcal{J}| \leq P$ , GRAD always use DEQ as job scheduler. In our previous work [18], we show that AGDEQ (the combination of DEQ and A-GREEDY) is O(1)-competitive with respect to mean response time for batched jobs when  $|\mathcal{J}| \leq P$ . The following lemma from [18] bounds the mean response time of a batched job set with  $|\mathcal{J}| \leq P$ .

**Lemma 7** [18] A job set  $\mathcal{J}$  is scheduled by GRAD on P processors where  $|\mathcal{J}| \leq P$ . The total response time  $R(\mathcal{J})$  of the schedule is at most

$$R(\mathcal{J}) \leq c \cdot \left(\frac{\rho+1}{\delta} \operatorname{swa}(\mathcal{J}) + \frac{2}{1-\delta} T_{\infty}(\mathcal{J}) + |\mathcal{J}| L(\log_{\rho} P + 1)\right)$$
where  $c = 2 - 2/(|\mathcal{J}| + 1)$ .

#### Case 2: when $|\mathcal{J}| > P$

We now derive the mean response time of GRAD for batched jobs for the second case where  $|\mathcal{J}| > P$ . Since all jobs in the job set  $\mathcal{J}$  arrive at time step 0, the number of uncompleted jobs decreases monotonically. When the number of uncompleted jobs drops down to P or below, GRAD switches its job scheduler from RR to DEQ. We divide the analysis into three parts. In **Part** (a), we prove two technical lemmas (Lemmas 8 and 9) which show the properties of round robin as the job scheduler. In **Part** (b), we analyze the completion time of the jobs which are scheduled by RR during their entire execution. In **Part** (c), we combine results and give response time of GRAD in general.

A batched job set  $\mathcal{J}$  can be divided into two subsets - RR set and **DEQ** set. The RR set, denoted as  $\mathcal{J}_{RR}$ , includes all the jobs in  $\mathcal{J}$ which are entirely scheduled by RR for their execution. The DEQ set, denoted as  $\mathcal{J}_{DEQ}$ , includes all the jobs in  $\mathcal{J}$  which are scheduled by RR at the beginning, and by DEQ eventually. There exists a unique quantum q called the *final RR quantum* such that q is the last quantum scheduled by RR, and from quanta q + 1 onwards are all scheduled by DEQ. According to RAD, there must be greater than P uncompleted jobs at the beginning of q, and less than or equal to Puncompleted jobs immediately after the execution of q. Let  $\sigma$  denote the total number of uncompleted jobs immediately after the execution of the final RR quantum. We know that  $\sigma = |\mathcal{J}_{DEQ}|$ , and  $\sigma \leq P$ . Let  $\pi$  denote a permutation that lists the jobs according to the ascending order of their completion time, i.e.  $T(J_{\pi(1)}) \leq T(J_{\pi(2)}) \leq$  $\dots \le T(J_{\pi(|\mathcal{J}|)})$ . We have  $\mathcal{J}_{RR} = \{J_{\pi(i)} \mid 1 \le i \le |\mathcal{J}| - \sigma\}$  and  $\mathcal{J}_{DEQ} = \left\{ J_{\pi(i)} \mid i > |\mathcal{J}| - \sigma \right\}$ , i.e.  $\mathcal{J}_{DEQ}$  includes the  $\sigma$  jobs that are completed last, and  $\mathcal{J}_{RR}$  includes the other  $|\mathcal{J}| - \sigma$  jobs.

We define two notations – **t-suffix** and **t-prefix**, and use them to simplify the notations. For any time step t, **t-suffix** denoted as  $\overrightarrow{t}$  represents the set of time steps from t to the completion of  $\mathcal J$  by  $\overrightarrow{t}=\{t,t+1,\ldots,T(\mathcal J)\}$ , while **t-prefix** denoted as  $\overleftarrow{t}$  represents set of time steps from 1 to t by  $\overleftarrow{t}=\{1,2,\ldots,t\}$ . We shall be interested in the suffixes of jobs. Define the t-suffix of a job  $J_i\in \mathcal J$  to be the job  $J_i\left(\overrightarrow{t}\right)$ , which is the portion of job  $J_i$  that remains after t-1 number of time steps have been executed. The t-suffix of the job set  $\mathcal J$  is

$$\mathcal{J}\left(\overrightarrow{t}\right) = \left\{J_i\left(\overrightarrow{t}\right) : J_i \in \mathcal{J} \text{ and } J_i\left(\overrightarrow{t}\right) \neq \emptyset\right\}.$$

Thus, we have  $\mathcal{J} = \mathcal{J}\left(\overrightarrow{1}\right)$ , and the number of uncompleted jobs at time step t is the number  $\left|\mathcal{J}\left(\overrightarrow{t}\right)\right|$  of nonempty jobs in  $\mathcal{J}\left(\overrightarrow{t}\right)$ .

Similarly, we can define the t-prefix of a job  $J_i$  as  $J_i\left(\overleftarrow{t}\right)$ , and the t-prefix of a job set  $\mathcal{J}$  as  $\mathcal{J}$  ( $\overleftarrow{t}$ ).

#### Case 2 - Part (a)

The following two technical lemmas present the properties of round robin as a job scheduler. The first lemma shows that jobs make almost the same progress on the execution of their work when they are scheduled by RR. The second lemma relates the work of jobs to their completion time.

**Lemma 8** A batched job set  $\mathcal{J}$  is scheduled by GRAD on a machine with P processors where  $|\mathcal{J}| > P$ . At any time step t scheduled by RR, for any two uncompleted jobs  $J_i$  and  $J_j$ , we have  $\left|T_1\left(J_i\left(\overleftarrow{t}\right)\right) - T_1\left(J_j\left(\overleftarrow{t}\right)\right)\right| \leq L$ , where L is the length of the scheduling quantum.

Proof. Since RR gives an equal share of processors to all uncompleted jobs, for any two jobs that arrive at the same time, their allotments differ by at most L at any time. When a job's allotment is 1, its allotted processor is always making useful work. Then the work done for any two uncompleted jobs differs by at most L at any time before their completion.

**Lemma 9** A batched job set  $\mathcal{J}$  is scheduled by GRAD on a machine with P processors where  $|\mathcal{J}| > P$ . The following two statements are true:

- $$\begin{split} 1. & \text{ If } J_i \in \mathcal{J}_{RR}, \ J_j \in \mathcal{J}_{RR}, \ \text{and } T_1\left(J_i\right) < T_1\left(J_j\right), \ \text{then } T(J_i) \leq \\ & T(J_j). \\ 2. & \text{ If } J_i \in \mathcal{J}_{RR}, \ \text{and } J_j \in \mathcal{J}_{DEQ}, \ \text{then } T_1\left(J_i\right) \leq T_1\left(J_j\right). \end{split}$$

*Proof.* We now prove the first statement. Let  $t = T(J_i)$ . At time step t, job  $J_i$  completes work  $T_1(J_i)$ . From Lemma 8, we know that  $T_1\left(J_j\left(\overleftarrow{t}\right)\right) \geq T_1\left(J_i\left(\overleftarrow{t}\right)\right) - L = T_1\left(J_i\right) - L$ . Since job  $J_j$  completes after job  $J_i$ , job  $J_j$  takes at least one more scheduling quantum than  $J_i$  to complete its execution. Thus the work done for  $J_j$  during the period from t to  $T(J_j)$  is at least L. Therefore, we have 
$$\begin{split} T_1\left(J_j\right) &= T_1\left(J_i\left(\overleftarrow{T(J_j)}\right)\right) \geq T_1\left(J_i\left(\overleftarrow{t}\right)\right) + L \geq T_1\left(J_i\right). \\ \text{For any two jobs } J_i \in \mathcal{J}_{RR} \text{, and } J_j \in \mathcal{J}_{DEQ} \text{, we have } T(J_i) < T_1 \in \mathcal{J}_{DEQ} \text{.} \end{split}$$

 $T(J_i)$ . By using a similar analysis, we can prove the second statement.

Lemma 9 relates the work of jobs to their completion time. Its second statement tells us that only the  $\sigma$  jobs with largest work are scheduled by DEQ eventually, and the other  $|\mathcal{J}| - \sigma$  jobs are scheduled by RR for their overall execution. Moreover, according to its first statement, under the schedule of RR, the jobs with less work are completed more quickly than those with more work. Consider the jobs according to their work such that  $T_1(J_1) \leq T_1(J_2) \leq \cdots \leq$  $T_1(J_{|\mathcal{J}|})$ . From Lemma 9, we have  $\mathcal{J}_{RR} = \{J_i | 1 \le i \le |\mathcal{J}| - \sigma\}$ and  $\mathcal{J}_{DEQ} = \{J_i | i > |\mathcal{J}| - \sigma\}.$ 

## Case 2 - Part (b)

The following lemma bounds the completion time of the jobs in  $\mathcal{J}_{RR}$ where  $T_1(J_i)$  denotes the work of a job  $J_i$ .

**Lemma 10** GRAD schedules a batched job set  $\mathcal{J}$  on a machine with P processors where  $|\mathcal{J}| > P$ . Consider the jobs according to their work such that  $T_1(J_1) \leq T_1(J_2) \leq \cdots \leq T_1(J_{|\mathcal{J}|})$ . For  $1 \leq i \leq |\mathcal{J}| - \sigma$ , the completion time  $T(J_i)$  of a job  $J_i$  is  $T(J_i) \le \left( \left( \left| \mathcal{J} \right| - i + 1 \right) T_1(J_i) + \sum_{1 \le j \le i} T_1(J_j) \right) / P + L.$ 

Since we consider the jobs according to their work, from Proof. Lemma 9, we have  $J_i \in \mathcal{J}_{RR}$  where  $1 \leq i \leq |\mathcal{J}| - \sigma$ . Such a job  $J_i$  completes its overall execution under the schedule of RR as job scheduler.

We first evaluate  $T_1(\mathcal{J}(t))$ , which is the work done for  $\mathcal{J}$  up to a time step t. Suppose that the job  $J_i$  terminates at the end of a quantum q where  $T(J_i) = q(L+1) - 1$ . Let t = qL - 1 be the end of the quantum q-1, which is L steps before the completion of  $J_i$ . The work done for  $J_i$  in interval  $\overleftarrow{t}$  is  $T_1\left(J_i\left(\overleftarrow{t}\right)\right) = T_1\left(J_i\right) - L$ . According to Lemma 8, no job completes more than  $T_1(J_i(t))+L$ amount of work in interval  $\overleftarrow{t}$ . Therefore, for any job  $\overrightarrow{J_j}$  with  $\overleftarrow{j} > i$ , we have

$$T_1\left(J_j\left(\overleftarrow{\tau}\right)\right) \leq T_1\left(J_i\left(\overleftarrow{\tau}\right)\right) + L$$
  
=  $T_1\left(J_i\right)$ . (7)

For each job  $J_j$  where j < i, by definition, we always have

$$T_1\left(J_j\left(\overleftarrow{t}\right)\right) \le T_1\left(J_j\right)$$
 (8)

Thus, at time step t, from Inequalities (7) and (8), the total work done for the job set  $\mathcal{J}$  is

$$T_{1}\left(\mathcal{J}\left(\overleftarrow{t}\right)\right)$$

$$= \sum_{1 \leq j < i} T_{1}\left(J_{j}\left(\overleftarrow{t}\right)\right) + T_{1}\left(J_{i}\left(\overleftarrow{t}\right)\right)$$

$$+ \sum_{i < j \leq |\mathcal{J}|} T_{1}\left(J_{j}\left(\overleftarrow{t}\right)\right)$$

$$\leq \left(|\mathcal{J}| - i + 1\right)T_{1}\left(J_{i}\right) + \sum_{1 \leq i \leq i} T_{1}\left(J_{j}\right). \tag{9}$$

Since RR always allots all processors to jobs, and all allotted processors are making useful work, RR executes P ready threads at any time step. Thus, the total work done for job set  $\mathcal J$  increases by P at each time step. From Inequality (9), we have

$$\begin{array}{lcl} t & = & T_1 \left( \mathcal{J} \left( \overleftarrow{t} \right) \right) / P \\ \\ & \leq & \left( \left( |\mathcal{J}| - i + 1 \right) T_1 \left( J_i \right) + \sum_{1 < j < i} T_1 \left( J_j \right) \right) / P \; . \end{array}$$

Since  $T(J_i) = t + L$ , we complete the proof.

#### Case 2 - Part (c)

The following lemma bounds the total response time of job sets scheduled by GRAD when  $|\mathcal{J}| > P$ , where swa  $(\mathcal{J})$  denotes squashed work area, and  $T_{\infty}(\mathcal{J})$  denotes the aggregate span.

**Lemma 11** Suppose that a job set  $\mathcal{J}$  is scheduled by GRAD on a machine with P processors where  $|\mathcal{J}| > P$ . The response time  $R(\mathcal{J})$ of  $\mathcal{J}$  is bounded by

$$R(\mathcal{J}) = \left(2 - \frac{2}{|\mathcal{J}| + 1}\right) \left(\frac{\rho + 1}{\delta} \operatorname{swa}(\mathcal{J}) + \frac{2}{1 - \delta} T_{\infty}(\mathcal{J})\right) + |\mathcal{J}| L + O(LP \log_{\rho} P).$$
 (10)

The jobs in  $\mathcal{J}$  can be divided into RR set  $\mathcal{J}_{RR}$  and DEQ set  $\mathcal{J}_{DEQ}$ . Let  $n = |\mathcal{J}|$  denote the number of jobs in  $\mathcal{J}$ . Recall that  $\sigma$  denotes the number of jobs in  $\mathcal{J}_{DEQ}$ , i.e.  $\sigma \leq P$ . Consider the jobs in the ascending order of their completion time such that  $T(J_1) \leq T(J_2) \leq \cdots \leq T(J_n)$ . From Lemma 9, we have  $\mathcal{J}_{RR} =$  $\{J_i|1\leq i\leq n-\sigma\}$  and  $\mathcal{J}_{DEQ}=\{J_i|i>n-\sigma\}$ . We will calculate the total response time of the jobs in  $\mathcal{J}_{RR}$  and  $\mathcal{J}_{DEQ}$  respectively.

**Step 1:** To calculate  $R(\mathcal{J}_{RR})$ , we apply Lemma 10. For any job  $J_i \in \mathcal{J}_{RR}$ , its completion time is  $T(J_i) \leq (1/P)((n-i+1)T_1(J_i) + (1/P)((n-i+1)T_1(J_i))$ 

 $\sum_{1\leq j< i}T_1\left(J_j\right)$ ) + L according to Lemma 10. Thus, the total response time of the jobs in  $\mathcal{J}_{RR}$  is

$$R(\mathcal{J}_{RR}) \le \frac{1}{P} \sum_{1 \le i \le n-\sigma} (2n - \sigma - 2i + 1) T_1(J_i) + Ln. \quad (11)$$

**Step 2:** We now calculate  $R(\mathcal{J}_{DEQ})$ . The  $\sigma$  jobs in  $\mathcal{J}_{DEQ}$  are scheduled by RR until the time step  $t = T(J_{n-\sigma})$  at which the job  $J_{n-\sigma}$  completes, and scheduled by DEQ afterwards. The total response time of  $\mathcal{J}_{DEQ}$  is

$$R\left(\mathcal{J}_{DEQ}\right) = R\left(\mathcal{J}_{DEQ}\left(\overrightarrow{t+1}\right)\right) + \sigma \cdot t. \tag{12}$$

From Lemma 10, we know that the completion time of the job  $J_{n-\sigma}$ 

$$t \le \left( (\sigma + 1)T_1 \left( J_{n-\sigma} \right) + \sum_{1 \le i < n-\sigma} T_1 \left( J_i \right) \right) / P + L . \quad (13)$$

To get R  $(\mathcal{J}_{DEQ})$ , we only need to calculate R  $(\mathcal{J}_{DEQ}(\overline{t+1}))$ .

Since the job set  $\mathcal{J}_{DEQ}$  is scheduled by DEQ as the job scheduler from time step t onwards, we can apply the total response time bound in Lemma 7 to calculate R  $(\mathcal{J}_{DEQ}(\overrightarrow{t+1}))$ . During the interval  $\overleftarrow{t}$ , job  $J_{n-\sigma}$  completes  $T_1(J_{n-\sigma})$  amount of work. From Lemma 8, we know that each job  $J_i$  with  $i>n-\sigma$  has completed at least  $T_1(J_{n-\sigma}) - L$  amount of work. Thus, such a job  $J_i$  has remaining work  $T_1\left(J_i\left(\overrightarrow{t+1}\right)\right) \leq T_1\left(J_i\right) - T_1\left(J_{n-\sigma}\right) + L$ . The squashed work of  $\mathcal{J}_{DEQ}\left(\overrightarrow{t+1}\right)$  is

$$\operatorname{swa}\left(\mathcal{J}_{DEQ}\left(\overline{t+1}\right)\right)$$

$$= \frac{1}{P}\operatorname{sq-sum}\left(\left\langle T_{1}\left(J_{i}\left(\overline{t+1}\right)\right) \mid n-\sigma+1 \leq i \leq n\right\rangle\right)$$

$$\leq \frac{1}{P}\operatorname{sq-sum}\left(\left\langle T_{1}\left(J_{i}\right) - T_{1}\left(J_{n-\sigma}\right) + L \mid n-\sigma+1 \leq i \leq n\right\rangle\right)$$

$$= \frac{1}{P}\sum_{n-\sigma+1 \leq i \leq n}(n-i+1)\left(T_{1}\left(J_{i}\right) - T_{1}\left(J_{n-\sigma}\right) + L\right)$$

$$\leq \frac{1}{P}\sum_{n-\sigma+1 \leq i \leq n}(n-i+1)T_{1}\left(J_{i}\right)$$

$$-\frac{(1+\sigma)\sigma}{2P}T_{1}\left(J_{n-\sigma}\right) + PL. \tag{14}$$

Let the constant c = 2 - 2/(1 + P) < 2. According to Lemma 7, we

$$R(\mathcal{J}_{DEQ}(\overrightarrow{t+1})) \le c \cdot \frac{\rho+1}{\delta} swa(\mathcal{J}_{DEQ}(\overrightarrow{t+1})) + E_1,$$
 (15)

where  $E_1 = c \cdot \frac{2}{1-\delta} T_{\infty}(\mathcal{J}) + cPL(\log_{\rho} P + 1)$ . We will now calculate the response time of  $\mathcal{J}_{DEQ}$ . Since we know c = 2 - 2/(1 + P) > 1, the responsiveness parameter  $\rho > 1$ , and the utilization parameter  $\delta \leq 1$ , we have  $c(\rho + 1)/\delta > 2$ . Given Equation (12), and Inequalities (13), (14) and (15), the response time of  $\mathcal{J}_{DEQ}$  is

$$R(\mathcal{J}_{DEQ})$$

$$= R(\mathcal{J}_{DEQ}(\overline{t+1})) + \sigma \cdot t$$

$$\leq c \cdot \frac{\rho+1}{\delta} \operatorname{swa} \left(\mathcal{J}_{DEQ}(\overline{t+1})\right) + E_1 + \sigma \cdot t$$

$$\leq c \cdot \frac{\rho+1}{\delta P} \sum_{n-\sigma+1 \leq i \leq n} (n-i+1)T_1(J_i)$$

$$+ \frac{\sigma}{P} \sum_{1 \leq i \leq n-\sigma} T_1(J_i) + E_2, \qquad (16)$$

where  $E_2 = E_1 + \left(c \cdot \frac{\rho+1}{\delta} + 1\right) PL$ .

Step 3: Given  $R(\mathcal{J}_{RR})$  in Inequality (11),  $R(\mathcal{J}_{DEQ})$  in Inequality (16), and  $c(\rho+1)/\delta > 2$ , the response time of  $\mathcal{J}$  is the sum of them as follows:

$$\begin{split} & \mathrm{R}(\mathcal{J}) \\ & = & \mathrm{R}(\mathcal{J}_{RR}) + \mathrm{R}(\mathcal{J}_{DEQ}) \\ & < & \frac{1}{P} \sum_{1 \leq i \leq n - \sigma} (2n - \sigma - 2i + 1) T_1\left(J_i\right) + Ln \\ & + c \cdot \frac{\rho + 1}{\delta P} \sum_{n - \sigma + 1 \leq i \leq n} (n - i + 1) T_1\left(J_i\right) \\ & + \frac{\sigma}{P} \sum_{1 \leq i < n - \sigma} T_1\left(J_i\right) + E_2 \\ & = & \frac{1}{P} \sum_{1 \leq i \leq n - \sigma} (2n - 2i + 1) T_1\left(J_i\right) + E_2 + Ln \\ & + c \cdot \frac{\rho + 1}{\delta P} \sum_{n - \sigma + 1 \leq i \leq n} (n - i + 1) T_1\left(J_i\right) \\ & + \frac{\sigma}{P} \sum_{1 \leq i < n - \sigma} T_1\left(J_i\right) - \frac{\sigma}{P} \sum_{1 \leq i \leq n - \sigma} T_1\left(J_i\right) \\ & \leq & c \cdot \frac{\rho + 1}{\delta P} \sum_{J_i \in \mathcal{J}} (n - i + 1) T_1\left(J_i\right) \\ & + c \cdot \frac{2}{1 - \delta} T_{\infty}(\mathcal{J}) + Ln + E_2 \\ & = & \left(2 - \frac{2}{n + 1}\right) \left(\frac{\rho + 1}{\delta} \operatorname{swa}(\mathcal{J}) + \frac{2}{1 - \delta} T_{\infty}(\mathcal{J})\right) \\ & + Ln + O(PL \log_{\rho} P) \; . \end{split}$$

Lemmas 7 and 11 bound the total response time of a batched job set  $\mathcal{J}$  when  $|\mathcal{J}| \leq P$  and  $|\mathcal{J}| > P$  respectively. Combining them, we have completed the proof of Theorem 4.