

Performance Analysis of Work-Conserving Schedulers in Minimizing the Total Flow Time with Phase Precedence

Yousi Zheng*, Prasun Sinha†, Ness B. Shroff*†

*Department of Electrical and Computer Engineering

†Department of Computer Science and Engineering

The Ohio State University, Columbus, OH 43210

Abstract—We consider the problem of minimizing the total flow time of multiple jobs in a pool of multiple homogenous machines, where the jobs arrive over time and have to be served with phase precedence. This is a common occurrence in job scheduling for the increasingly popular data center oriented systems, where jobs need to be processed through a MAP and Reduce procedure before leaving the system. For this problem, one can construct an arrival pattern such that no scheduler can achieve a constant competitive ratio. However, what we find is that by using a slightly weaker metric of performance, which we call the efficiency ratio. We say that a scheduler achieves an efficiency ratio of γ when the flow time incurred by that scheduler divided by the minimum flow time achieved over all possible schedulers is less than or equal to γ almost surely, when the time slots or job arrivals go to infinity. We show a surprising property that any work-conserving scheduler for the flow time problem with phase precedence has a constant efficiency ratio in both preemptive and non-preemptive scenarios under some weak assumptions. We provide numerical results to support our analysis.

I. INTRODUCTION

In many widely used systems, the workload of arriving jobs can be divided into several phases. Some phases have to be finished before other phases start to process. One famous example of such system is MapReduce framework, which is designed to process massive amounts of data in a cluster of machines [1]. It is widely used for applications such as search indexing, distributed searching, web statistics generation, and data mining. MapReduce has two elemental phases: Map and Reduce. A critical consideration for the design of the scheduler is the *precedence* between the Map and Reduce tasks. For each job, the Map tasks need to be finished before starting any of its Reduce tasks [1], [2].

To minimize the total flow time, it is well known that the Shortest-Remaining-Processing-Time (SRPT) scheduler is optimal in one machine case [3] or in multiple machines with work preserving malleable tasks [4]. The related scenarios are studied: scheduling malleable tasks [5], scheduling chain-structured tasks [6], and scheduling with release time constraint [7]. However, they are focus on the complexity analysis, and they cannot directly applied in our multiple phases scenario. For the MapReduce framework, some scheduling solutions has also been proposed [2], [8]–[11], but analytical

bounds on performance have been derived only in some of these works [2], [10], [11]. However, rather than focusing directly on the flow-time, for deriving performance bounds, [10], [11] have considered a slightly different problem of minimizing total completion time and [2] has assumed speed-up of the machines.

In this paper, we directly analyze the performance of total delay (flow-time) in the system. In an attempt to minimize this, we introduce a new metric to analyze the performance of schedulers called *efficiency ratio*. Based on this new metric, we analyze and design several schedulers, which can guarantee the performance.

The contributions of this paper are as follows:

- To directly analyze the total delay in the system, we propose a new metric to measure the performance of schedulers, which we call the *efficiency ratio*. (Section II)
- For bounded workload of Phase 2, we then show a surprising property that for the flow time problem any work-conserving scheduler has a constant efficiency ratio in the preemptive scenario (precise definitions provided in Section II). (Section III)
- For bounded workload of Phase 2, we then show the same property in the non-preemptive scenarios (precise definitions provided in Section II). (Section IV)
- We relax the assumption of Phase 2 from bounded workload to light-tailed distributed workload, and get the same properties in both preemptive and non-preemptive scenarios. (Section V)
- For the typical work-conserving schedulers, we show that the convergence of efficiency ratios through simulations. (Section VI)

II. SYSTEM MODEL AND EFFICIENCY RATIO

A. System Model

Consider a pool of N machines. A machine could represent a processor, a core in a multi-core processor or a virtual machine. Assume that there are n jobs arriving into the system, and each machine can only process one job at a time. We assume the scheduler periodically collects the information on the state of jobs running on the machines, which is used to

make scheduling decisions. Time is slotted and each machine can run one unit of workload in each time slot. We assume that the distribution of job arrivals in each time slot is i.i.d., and the arrival rate is λ . Each job i contains several tasks, which can be classified into two phases: Phase 1 and Phase 2. Its tasks of Phase 1 need to be finished before starting any of its tasks of Phase 2. For each job i , Phase 1 brings M_i units of workload and Phase 2 brings R_i units of workload. Each task of Phase 1 has 1 unit of workload², however, each task of Phase 2 can have multiple units of workload. Assume that $\{M_i\}$ are i.i.d. with expectation \overline{M} , and $\{R_i\}$ are i.i.d. with expectation \overline{R} . We assume that the traffic intensity $\rho < 1$, i.e., $\lambda < \frac{N}{\overline{M} + \overline{R}}$. Assume the moment generating function of workload of arriving jobs in a time slot has finite value in some neighborhood of 0. In time slot t for job i , $m_{i,t}$ and $r_{i,t}$ machines are scheduled for the tasks of Phase 1 and 2, respectively. We assume that job i contains K_i tasks of Phase 2, and the workload of the Phase 2 task k of job i is $R_i^{(k)}$. Thus, for any job i , $\sum_{k=1}^{K_i} R_i^{(k)} = R_i$. In time slot t for job i , $r_{i,t}^{(k)}$ machines are scheduled for the task k of Phase 2. As each task of Phase 2 may consist of multiple units of workload, it can be processed in either preemptive or non-preemptive fashion based on the type of scheduler.

Definition 1. A scheduler is called *preemptive* if the tasks of Phase 2 belonging to the same job can run in parallel on multiple machines, can be interrupted by any other task, and can be rescheduled to different machines in different time slots.

A scheduler is called *non-preemptive* if each task of Phase 2 can only be scheduled on one machine and, once started, it must keep running without any interruption.

The flow time F_i of job i is equal to $f_i^{(r)} - a_i + 1$. The objective of the scheduler is to determine the assignment of jobs in each time slot, such that the cost of delaying the jobs or equivalently the flow time is minimized.

For the preemptive scenario, the problem definition is as follows:

$$\begin{aligned} \min_{m_{i,t}, r_{i,t}} \quad & \sum_{i=1}^n (f_i^{(r)} - a_i + 1) \\ \text{s.t.} \quad & \sum_{i=1}^n (m_{i,t} + r_{i,t}) \leq N, \quad r_{i,t} \geq 0, \quad m_{i,t} \geq 0, \quad \forall t, \\ & \sum_{t=a_i}^{f_i^{(m)}} m_{i,t} = M_i, \quad \sum_{t=f_i^{(m)}+1}^{f_i^{(r)}} r_{i,t} = R_i, \quad \forall i \in \{1, \dots, n\}. \end{aligned} \quad (1)$$

In the non-preemptive scenario, the tasks of Phase 2 cannot be interrupted by other jobs. Once a task of Phase 2 begins execution on a machine, it has to keep running on that machine without interruption until all its workload is finished. Also, the optimization problem in this scenario is similar to Eq. (1), with

²For example, in the MapReduce framework, the Map tasks are independent and have small workload [9], such assumption is valid.

additional constraints representing the non-preemptive nature, as shown below:

$$\begin{aligned} \min_{m_{i,t}, r_{i,t}^{(k)}} \quad & \sum_{i=1}^n (f_i^{(r)} - a_i + 1) \\ \text{s.t.} \quad & \sum_{i=1}^n (m_{i,t} + \sum_{k=1}^{K_i} r_{i,t}^{(k)}) \leq N, \quad \forall t, \\ & \sum_{t=a_i}^{f_i^{(m)}} m_{i,t} = M_i, \quad m_{i,t} \geq 0, \quad \forall i \in \{1, \dots, n\}, \\ & \sum_{t=f_i^{(m)}+1}^{f_i^{(r)}} r_{i,t}^{(k)} = R_i^{(k)}, \quad \forall i \in \{1, \dots, n\}, \quad \forall k \in \{1, \dots, K_i\}, \\ & r_{i,t}^{(k)} = 0 \text{ or } 1, \quad r_{i,t}^{(k)} = 1 \text{ if } 0 < \sum_{s=0}^{t-1} r_{i,s}^{(k)} < R_i^{(k)}. \end{aligned}$$

Theorem 1. The scheduling problem (both preemptive and non-preemptive) is NP-complete in the strong sense.

Proof: The proof of NP-completeness follows a fairly standard reduction from 3-Partition and is described in technical report [12]. ■

B. Efficiency Ratio

The competitive ratio is often used as a measure of performance in a wide variety of scheduling problems. For our problem, the scheduling algorithm S has a competitive ratio of c , if for any total time T , any number of arrivals n in the time T , any arrival time a_i of each job i , any workload M_i and R_i of the Phase 1 and 2 tasks with respect to each arrival job i , the total flow time $F^S(T, n, \{a_i, M_i, R_i; i = 1 \dots n\})$ of scheduling algorithm S satisfies the following:

$$\frac{F^S(T, n, \{a_i, M_i, R_i; i = 1 \dots n\})}{F^*(T, n, \{a_i, M_i, R_i; i = 1 \dots n\})} \leq c,$$

where

$$\begin{aligned} & F^*(T, n, \{a_i, M_i, R_i; i = 1 \dots n\}) \\ & = \min_S F^S(T, n, \{a_i, M_i, R_i; i = 1 \dots n\}), \end{aligned}$$

is the minimum flow time of an optimal off-line algorithm.

For our problem, we can construct an arrival pattern such that no scheduler can achieve a constant competitive ratio c . The quick example is given in technical report [12]. We now introduce a slightly weaker notion of performance, called the **efficiency ratio**.

Definition 2. We say that the scheduling algorithm S has an efficiency ratio γ , if the total flow time $F^S(T, n, \{a_i, M_i, R_i; i = 1 \dots n\})$ of scheduling algorithm S satisfies the following:

$$\lim_{T \rightarrow \infty} \frac{F^S(T, n, \{a_i, M_i, R_i; i = 1 \dots n\})}{F^*(T, n, \{a_i, M_i, R_i; i = 1 \dots n\})} \leq \gamma, \text{ with probability } 1.$$

Later, we will show that for the quick example, a constant efficiency ratio γ still can exist (e.g., the non-preemptive

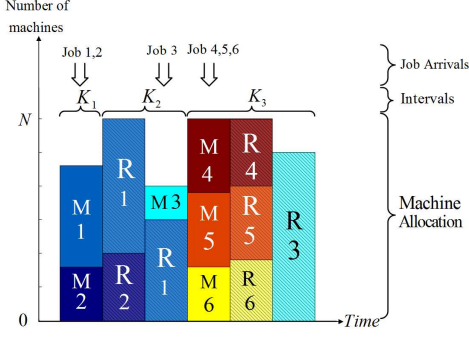


Fig. 1. A example schedule of a work-conserving scheduler

scenario with light-tailed distributed Reduce workload in Section IV).

III. WORK-CONSERVING SCHEDULER IN THE PREEMPTIVE SCENARIO

In this section, we analyze the performance of work-conserving scheduler in the preemptive scenario. In this section, we study the case in which the workload of Phase 2 of each job is bounded by a constant, i.e., there exists a constant R_{max} , s.t., $R_i \leq R_{max}$, $\forall i$.

Consider the total scheduled number of machines over all the time slots. If all the N machines are scheduled in a time slot, we call this time slot a “complete” time slot; otherwise, we call this time slot a “incomplete” time slot. We define the j^{th} “interval” to be the interval between the $(j-1)^{st}$ incomplete time slot and the j^{th} incomplete time slot. (We define the first interval as the interval from the first time slot to the first incomplete time slot.) Thus, the last time slot of each interval is the only incomplete time slot in the interval. Let K_j be the length of the j^{th} interval, as shown in Fig. 1.

Firstly, we show that any moment of the length K_j of interval is bounded by some constants, as shown in Lemma 1.

Lemma 1. *If the scheduler is work-conserving, then for any given number H , there exists a constant B_H , such that $E[K_j^H] < B_H$, $\forall j$.*

Proof: Since K_j is a positive random variable, then $E[K_j^H] \leq 1$ if $H \leq 0$. Thus, we only focus on the case that $H > 0$.

In the j^{th} interval, consider the t^{th} time slot in this interval. Assume that the total workload of the arrivals in the t^{th} time slot of this interval is $W_{j,t}$. Also assume that the unavailable workload of Phase 2 at the end of the t^{th} time slot in this interval is $R_{j,t}$, i.e., the workload of Phase 2 whose corresponding tasks of Phase 1 are finished in the t^{th} time slot in this interval is $R_{j,t}$. The remaining workload of the Reduce function left by the previous interval is $R_{j,0}$. Then, the distribution of K_j can be shown as below.

$$P(K_j = k) = P(W_{j,1} + R_{j,0} - R_{j,1} \geq N, \dots, \sum_{t=1}^{k-1} W_{j,t} + R_{j,0} - R_{j,k-1} \geq (k-1)N, \sum_{t=1}^k W_{j,t} + R_{j,0} - R_{j,k} < kN) \quad (2)$$

For each job i , $R_i \leq R_{max}$. And in the last slot of the previous interval, the maximum number of finished jobs of Phase 1 is $N-1$. Thus, $R_{j,0} \leq (N-1)R_{max}$. Now, let's consider the case of $k > 1$.

$$\begin{aligned} P(K_j = k) &\leq P\left(\sum_{t=1}^{k-1} W_{j,t} + R_{j,0} - R_{j,k-1} \geq (k-1)N\right) \\ &\leq P\left(\sum_{t=1}^{k-1} W_{j,t} + (N-1)R_{max} \geq (k-1)N\right) \\ &= P\left(\frac{\sum_{t=1}^{k-1} W_{j,t}}{k-1} \geq N - \frac{(N-1)R_{max}}{k-1}\right). \end{aligned} \quad (3)$$

By the assumptions, the total (both Phase 1 and 2) arriving workload $W_{j,t}$ in each time slot is i.i.d., and the distribution does not depend on the interval index j or the index of time slot t in this interval. Then, for any interval j , we know that

$$P(K_j = k) \leq P\left(\frac{\sum_{s=1}^{k-1} W_s}{k-1} \geq N - \frac{(N-1)R_{max}}{k-1}\right),$$

where $\{W_s\}$ is a sequence of i.i.d. random variables with the same distribution of $W_{j,t}$.

Since the traffic intensity $\rho < 1$, $E[W_s] = \lambda(\bar{M} + \bar{R}) < N$.

For any $\epsilon > 0$, there exists a k_0 , such that for any $k > k_0$, $\frac{(N-1)R_{max}}{k-1} \leq \epsilon$. Then for any $k > k_0$, we have

$$P(K_j = k) \leq P\left(\frac{\sum_{s=1}^{k-1} W_s}{k-1} \geq N - \epsilon\right),$$

We take $\epsilon < N - E[W_s]$ and corresponding k_0 . For any $k > k_0 = 1 + \lceil \frac{(N-1)R_{max}}{\epsilon} \rceil$, by Cramer-Chernoff Theorem [13], we can achieve that

$$P\left(\frac{\sum_{s=1}^{k-1} W_s}{k-1} \geq N - \epsilon\right) \leq e^{-kl(N-\epsilon)},$$

where the rate function $l(a)$ is shown as below:

$$l(a) = \sup_{\theta \geq 0} (\theta a - \log(E[e^{\theta W_s}])). \quad (4)$$

Since the moment generating function of workload in a time slot has finite value in some neighborhood around 0, for any

$0 < \epsilon < N - \lambda(\overline{M} + \overline{R})$, $l(N - \epsilon) > 0$. Thus,

$$\begin{aligned} E[K_j^H] &= \sum_{k=0}^{\infty} k^H P(K_j = k) \\ &= \sum_{k=0}^{k_0} k^H P(K_j = k) + \sum_{k=k_0+1}^{\infty} k^H P(K_j = k) \\ &\leq k_0^H + \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} \\ &= \left(1 + \left\lceil \frac{(N-1)R_{max}}{\epsilon} \right\rceil\right)^H + \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)}. \end{aligned}$$

Since $l(N - \epsilon) > 0$, $\sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)}$ is bounded.

Thus, given H , $E[K_j^H]$ is bounded by a constant B_H , for any j , where B_H is given as below.

$$B_H \triangleq \min_{\epsilon \in (0, N - \lambda(\overline{M} + \overline{R}))} \left\{ \left(1 + \left\lceil \frac{(N-1)R_{max}}{\epsilon} \right\rceil\right)^H + \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} \right\}. \quad (5)$$

Now, we directly obtain the expression of first and second moment of K_j .

Corollary 1. $E[K_j]$ is bounded by a constant B_1 , $E[K_j^2]$ is bounded by a constant B_2 , for any j . The expressions of B_1 and B_2 are shown as below, where the rate function $l(a)$ is defined in Eq. (4).

$$\begin{aligned} B_1 &= \min_{\epsilon \in (0, N - \lambda(\overline{M} + \overline{R}))} \left\{ 1 + \frac{\left\lceil \frac{(N-1)R_{max}}{\epsilon} \right\rceil}{\frac{2e^{l(N-\epsilon)} - 1}{e^{l(N-\epsilon)}(e^{l(N-\epsilon)} - 1)^2}} \right\}, \\ B_2 &= \min_{\epsilon \in (0, N - \lambda(\overline{M} + \overline{R}))} \left\{ \left(1 + \left\lceil \frac{(N-1)R_{max}}{\epsilon} \right\rceil\right)^2 + \frac{4e^{2l(N-\epsilon)} - 3e^{l(N-\epsilon)} + 1}{e^{l(N-\epsilon)}(e^{l(N-\epsilon)} - 1)^3} \right\}. \end{aligned}$$

Proof: By substituting $H = 2$ and $H = 1$ in Eq. (5), we directly achieve the corollary. ■

We add some dummy workload of Phase 2 in the beginning time slot of each interval, such that the remaining workload of the previous interval is equal to $(N - 1)R_{max}$. The added workload are only scheduled when on real job can be scheduled. Then, in the new artificial system, the total flow time of the real jobs doesn't change. However, the total number of intervals may decrease, because the added workload may merge some adjacent intervals.

Corollary 2. In the artificial system, the length of new intervals $\{\widetilde{K}_j, j = 1, 2, \dots\}$ also satisfies the same results of Lemma 1 and Corollary 1.

Proof: The proof is similar to the proofs of Lemma 1 and Corollary 1. ■

Remark 1. In the artificial system constructed in Corollary 2, the length of each interval has no effect on the length of other intervals, and it is only dependent on the job arrivals and the workload of each job in this interval. Based on our assumptions, the job arrivals in each time slot are i.i.d.. Also, the workload of jobs are also i.i.d.. Thus, the random variables $\{\widetilde{K}_j, j = 1, 2, \dots\}$ are i.i.d..

Theorem 2. For any work-conserving scheduling policy, the efficiency ratio is not greater than $\frac{B_2 + B_1}{\max\{2, \frac{1-p_0}{\lambda}, \frac{\rho}{\lambda}\} \max\{1, \frac{1}{N(1-\rho)}\}}$ with probability 1, when n goes to infinity. B_1 and B_2 can be achieved by Corollary 1, and p_0 is the probability that no job arrives in a time slot.

Proof: Consider the artificial system constructed in Corollary 2. Since we have more dummy workload of Phase 2 in the artificial system, the total flow time \widetilde{F} of the artificial system is not less than the total flow time F of the work-conserving scheduling policy.

Note that, in each interval, all the job arrivals in this interval must finish Phase 1 in this interval, and all their workload of Phase 2 will finish before the end of the next interval. In other words, for all the arrivals in the j^{th} interval, Phase 1 is finished in K_j time slots, and Phase 2 is finished in $K_j + K_{j+1}$ time slots, as shown in Fig. 1. This is also true for the artificial system with dummy workload of Phase 2, i.e., all the arrivals in the j^{th} interval of the artificial system can be finished in $\widetilde{K}_j + \widetilde{K}_{j+1}$ time slots.

Let A_j be the number of arrivals in the j^{th} interval. The total flow time F_j of the work-conserving scheduler for all the arrivals in the j^{th} interval is not greater than $A_j(K_j + K_{j+1})$. Similarly, the total flow time \widetilde{F}_j for all the arrivals in the j^{th} interval of the artificial system is not greater than $\widetilde{A}_j(\widetilde{K}_j + \widetilde{K}_{j+1})$, where \widetilde{A}_j is the number of arrivals in the j^{th} interval of the artificial system.

Let F be the total flow time of any work-conserving scheduler, and F^* be the minimized total flow time. Assume that there are a total of m intervals for the n arrivals. Correspondingly, there are a total of \widetilde{m} intervals in the artificial system. Also, based on the construction of the artificial system, $m \geq \widetilde{m}$ and $F = \widetilde{F} = \sum_{j=1}^{\widetilde{m}} \widetilde{F}_j$. $T \triangleq \sum_{j=1}^m K_j$ is equal to the finishing time of all the n jobs. For any scheduler, the total flow time is not less than \hat{T} , which is the number of time slots that scheduled machines is not equal to 0. Thus, for the optimal scheduler, its the total flow time F_j^* should be not less than \hat{T} .

For any scheduling policy, each job needs at least 1 time slot to schedule its tasks of Phase 1 and 1 time slot to schedule its tasks of Phase 2. Thus, the lower bound of the each job's flow time is 2. Thus, for the optimal scheduling method, its total flow time F_j^* should not be less than $2A_j$.

Since $E(M_i + R_i) > 0$, $\rho < 1$, then $\lim_{n \rightarrow \infty} m = \infty$ and

$$\lim_{n \rightarrow \infty} \tilde{m} = \infty.$$

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{F}{F^*} &= \lim_{n \rightarrow \infty} \frac{\tilde{F}}{F^*} \leq \lim_{\tilde{m} \rightarrow \infty} \frac{1}{\hat{T}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j (\tilde{K}_j + \tilde{K}_{j+1}) \\ &\leq \frac{\lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j \tilde{K}_j + \lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j \tilde{K}_{j+1}}{\lim_{\tilde{m} \rightarrow \infty} \hat{T}/m}. \end{aligned}$$

Since $\{\tilde{K}_j, j = 1, 2, \dots\}$ are i.i.d. distributed random variables, then $\{\tilde{A}_j \tilde{K}_j, j = 1, 2, \dots\}$ are also i.i.d.. Based on Strong Law of Large Number (SLLN), we can achieve that

$$\lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j \tilde{K}_j \stackrel{\text{w.p.1}}{=} E[\tilde{A}_j \tilde{K}_j]. \quad (6)$$

For $\{\tilde{A}_j \tilde{K}_{j+1}, j = 1, 2, \dots\}$, they are identically distributed, but not independent. However, all the odd term are independent, and all the even term are independent. Based on SLLN, we can achieve that

$$\lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j \tilde{K}_{j+1} \stackrel{\text{w.p.1}}{=} E[\tilde{A}_j \tilde{K}_{j+1}]. \quad (7)$$

Based on Eqs. 6 and 7, we can achieve that

$$\begin{aligned} \lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j \tilde{K}_j + \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j \tilde{K}_{j+1} \\ &\stackrel{\text{w.p.1}}{=} E[\tilde{A}_j \tilde{K}_j] + E[\tilde{A}_j \tilde{K}_{j+1}] \\ &= E[E[\tilde{A}_j \tilde{K}_j | \tilde{K}_j]] + E[\tilde{A}_j] E[\tilde{K}_{j+1}] \\ &= E[\tilde{K}_j E[\tilde{A}_j | \tilde{K}_j]] + E[E[\tilde{A}_j | \tilde{K}_j]] E[\tilde{K}_{j+1}] \\ &= \lambda(E[\tilde{K}_j^2] + E[\tilde{K}_j]^2). \end{aligned}$$

For these n jobs, there are m incomplete time slots before the time T . For each incomplete time slot, there are at most $N-1$ machines are assigned. Then, $N \sum_{j=1}^m (K_j - 1) + (N-1)m \geq \sum_{s \in \{\text{Arrivals before } T\}} W_s$. Thus,

$$\left(N - \frac{\sum_{s \in \{\text{Arrivals before } T\}} W_s}{T} \right) \frac{\sum_{j=1}^m K_j}{m} \geq 1. \quad (8)$$

Since the workload W_s of each job is i.i.d., by SLLN, we have $(N - \lambda(\bar{M} + \bar{R})) \lim_{m \rightarrow \infty} \frac{\sum_{j=1}^m K_j}{m} \geq 1$ w.p.1, i.e.,

$$\lim_{m \rightarrow \infty} \frac{\sum_{j=1}^m K_j}{m} \geq \frac{1}{N(1-\rho)} \text{ w.p.1. Thus,}$$

$$\lim_{m \rightarrow \infty} \frac{T}{m} = \lim_{m \rightarrow \infty} \frac{\sum_{j=1}^m K_j}{m} \geq \max \left\{ 1, \frac{1}{N(1-\rho)} \right\} \text{ w.p.1.}$$

For any work-conserving scheduler, \hat{T} is not less than the

total number of time slots with at least 1 arrival, we have

$$\lim_{m \rightarrow \infty} \frac{\hat{T}}{m} \geq (1-p_0) \lim_{m \rightarrow \infty} \frac{T}{m} \text{ w. p. 1.}$$

At the same time, for each time slot that the number of scheduled machines is greater than 0, there are at most N machines are scheduled. Then, we can get $N\hat{T} \geq \sum_{i=1}^n (M_i + R_i)$. Then,

$$\lim_{m \rightarrow \infty} \frac{\hat{T}}{m} \geq \frac{\lambda(\bar{M} + \bar{R})}{N} \lim_{m \rightarrow \infty} \frac{T}{m} \text{ w.p.1.} = \rho \lim_{m \rightarrow \infty} \frac{T}{m}. \text{ Thus,}$$

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{F}{F^*} &\leq \frac{\lambda(E[\tilde{K}_j^2] + E[\tilde{K}_j]^2)}{\max\{\rho, 1-p_0\} \max\left\{1, \frac{1}{N(1-\rho)}\right\}} \text{ w.p.1.} \\ &\leq \frac{\lambda(B_2 + B_1^2)}{\max\{\rho, 1-p_0\} \max\left\{1, \frac{1}{N(1-\rho)}\right\}} \text{ w.p.1.} \end{aligned} \quad (9)$$

Similarly,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{F}{F^*} &= \lim_{n \rightarrow \infty} \left(\sum_{j=1}^m F_j / \sum_{j=1}^m F_j^* \right) \\ &\leq \frac{\lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j (\tilde{K}_j + \tilde{K}_{j+1})}{\lim_{m \rightarrow \infty} \sum_{j=1}^m 2A_j} \leq \frac{\lim_{\tilde{m} \rightarrow \infty} \frac{1}{\tilde{m}} \sum_{j=1}^{\tilde{m}} \tilde{A}_j (\tilde{K}_j + \tilde{K}_{j+1})}{\lim_{m \rightarrow \infty} \frac{1}{m} \sum_{j=1}^m 2A_j} \end{aligned}$$

With Strong Law of Large Number (SLLN), we can get

$$\begin{aligned} \lim_{m \rightarrow \infty} \frac{\sum_{j=1}^m A_j}{m} &= \lim_{m \rightarrow \infty} \left(\sum_{j=1}^m A_j / \sum_{j=1}^m K_j \right) \lim_{m \rightarrow \infty} \left(\frac{1}{m} \sum_{j=1}^m K_j \right) \\ &= \lambda \lim_{m \rightarrow \infty} \left(\frac{1}{m} \sum_{j=1}^m K_j \right) \text{ w. p. 1.} \geq \lambda \max \left\{ 1, \frac{1}{N(1-\rho)} \right\} \end{aligned}$$

Thus,

$$\lim_{n \rightarrow \infty} \frac{F}{F^*} \leq \frac{\lambda(B_2 + B_1^2)}{2\lambda \max\left\{1, \frac{1}{N(1-\rho)}\right\}} \text{ w.p.1.} \quad (10)$$

Based on the Eqs. 9 and 10, we can achieve that the efficiency ratio of any work-conserving scheduling policy is $\frac{B_2 + B_1^2}{\max\left\{2, \frac{1-p_0}{\lambda}, \frac{\rho}{\lambda}\right\} \max\left\{1, \frac{1}{N(1-\rho)}\right\}}$. ■

IV. WORK-CONSERVING SCHEDULER IN THE NON-PREEMPTIVE SCENARIO

The definitions of $B_1, B_2, p_0, F, \tilde{F}, F^*, F_j, \tilde{F}_j, K_j, \tilde{K}_j, A_j, \tilde{A}_j, W_{j,t}$, and $R_{j,t}$ in this section are same as Section III.

Lemma 2. *Lemma 1, Corollary 1, Corollary 2 in Section III are also valid in the non-preemptive scenario.*

Proof: Different from preemptive scenario in Section III, the constitution of $R_{j,t}$ are different, because it contains two parts. First, it contains the unavailable workload of Phase 2 which are just released by Phase 1 which are finished in the last time slot of each interval. Second, it also contains the workload of Phase 2 which cannot be processed in this time

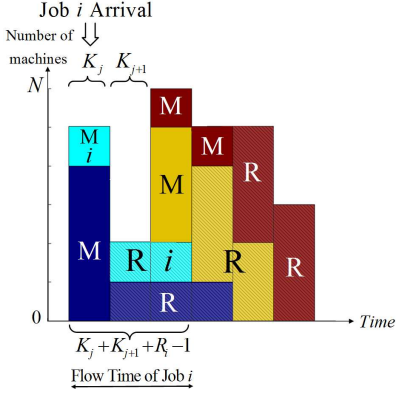


Fig. 2. The arrivals in the j^{th} interval finish Phase 1 functions in K_j time slots and Phase 2 in $K_j + K_{j+1} + R_{max} - 1$ time slots

slot, even if there are idle machines, because the tasks of Phase 2 is non-preemptive.

However, we can also achieve that the remaining workload of Phase 2 from previous interval satisfies $R_{j,0} \leq (N - 1)R_{max}$, because the total number of remaining unfinished tasks of Phase 2 and new available tasks of Phase 2 in this time slot is less than N . All the following steps are same. ■

Theorem 3. *In the non-preemptive scenario, any work-conserving scheduling policy has a efficiency ratio*

$$\frac{B_2 + B_1^2 + B_1(\bar{R} - 1)}{\max\{2, \frac{1-p_0}{\lambda}, \frac{\rho}{\lambda}\} \max\{1, \frac{1}{N(1-\rho)}\}}.$$

Proof: In the non-preemptive scenario, in each interval, all the workload of Phase 1 of job arrivals in this interval is also finished in this interval, and all their tasks of Phase 2 will be started before the end of the next interval. In other words, for the job i in the j^{th} interval, the tasks of Phase 1 are finished in K_j time slots, and the tasks of Phase 2 are finished in $K_j + K_{j+1} + R_i - 1$ time slots, as shown in Fig. 2. And this is also true for the artificial system with dummy workload of Phase 2, i.e., all the arrivals in the j^{th} interval of the artificial system can be finished in $\widetilde{K}_j + \widetilde{K}_{j+1} + R_i - 1$. So, F_i in the j^{th} interval is less than $A_j(K_j + K_{j+1} + R_i - 1)$. Similarly, \widetilde{F}_i in the j^{th} interval of the artificial system is less than $\widetilde{A}_j(\widetilde{K}_j + \widetilde{K}_{j+1} + R_i - 1)$. Compared to Section III, the total flow time we can use the same proving method in Theorem 2 to achieve that

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{F}{F^*} &\leq \lim_{m \rightarrow \infty} \frac{\sum_{j=1}^m \widetilde{A}_j(\widetilde{K}_j + \widetilde{K}_{j+1}) + \sum_{i=1}^n (R_i - 1)}{\sum_{j=1}^m (K_j - S_j)} \\ &= \frac{\lim_{m \rightarrow \infty} \frac{1}{m} \sum_{j=1}^m \widetilde{A}_j(\widetilde{K}_j + \widetilde{K}_{j+1}) + \lim_{m \rightarrow \infty} \frac{n}{m} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (R_i - 1)}{\lim_{m \rightarrow \infty} \frac{1}{m} \sum_{j=1}^m (K_j - S_j)} \\ &\leq \frac{\lambda[B_2 + B_1^2 + B_1(\bar{R} - 1)]}{\max\{\rho, 1 - p_0\} \max\{1, \frac{1}{N(1-\rho)}\}} \text{ w.p.1,} \end{aligned} \quad (11)$$

and

$$\lim_{n \rightarrow \infty} \frac{F}{F^*} \leq \frac{\lambda[B_2 + B_1^2 + B_1(\bar{R} - 1)]}{2\lambda \max\{1, \frac{1}{N(1-\rho)}\}} \text{ w.p.1.} \quad (12)$$

Thus, based on the Eqs. 11 and 12, we can achieve that any work-conserving scheduler has the efficiency ratio

$$\frac{B_2 + B_1^2 + B_1(\bar{R} - 1)}{\max\{2, \frac{1-p_0}{\lambda}, \frac{\rho}{\lambda}\} \max\{1, \frac{1}{N(1-\rho)}\}}.$$

■

V. LIGHT-TAILED DISTRIBUTED WORKLOAD OF PHASE 2

In this part, we assume the workload R of Phase 2 are light-tailed distributed, i.e., $\exists r_0$, such that

$$P(R \geq r) \leq \alpha \exp(-\beta r), \quad \forall r \geq r_0,$$

where $\alpha, \beta > 0$ are two constants. We use the same definitions of $W_{j,t}$, $R_{j,t}$ and K_j as Section III.

Lemma 3. *If the scheduling algorithm is work-conserving, then for any given number H , there exists a constant B_H , such that $E[K_j^H] < B_H, \forall j$.*

Proof: Similarly to Eqs. 2 and 3, we can achieve

$$P(K_j = k) \leq P\left(\sum_{t=1}^{k-1} W_{j,t} + R_{j,0} \geq (k-1)N\right).$$

Let R' be the maximal remaining workload of jobs in the previous interval, then

$$\begin{aligned} P(K_j = k) &\leq P\left(\sum_{t=1}^{k-1} W_{j,t} + (N-1)R' \geq (k-1)N\right) \\ &= \sum_{r=0}^{\infty} \left[P\left(\sum_{t=1}^{k-1} W_{j,t} + (N-1)R' \geq (k-1)N \mid R' = r\right) \right. \\ &\quad \left. P(R' = r) \right] \\ &= \sum_{r=0}^{\infty} \left[P\left(\frac{\sum_{t=1}^{k-1} W_{j,t}}{k-1} \geq N - \frac{(N-1)r}{k-1}\right) P(R' = r) \right]. \end{aligned} \quad (13)$$

Thus, we can achieve

$$\begin{aligned} E[K_j^H] &= \sum_{k=0}^{\infty} k^H P(K_j = k) \\ &= \sum_{k=0}^{\infty} \left\{ k^H \sum_{r=0}^{\infty} \left[P\left(\frac{\sum_{t=1}^{k-1} W_{j,t}}{k-1} \geq N - \frac{(N-1)r}{k-1}\right) P(R' = r) \right] \right\} \\ &= \sum_{r=0}^{\infty} \left\{ P(R' = r) \sum_{k=0}^{\infty} \left[k^H P\left(\frac{\sum_{t=1}^{k-1} W_{j,t}}{k-1} \geq N - \frac{(N-1)r}{k-1}\right) \right] \right\}. \end{aligned}$$

The next steps are similar to the proof of Lemma 1, we skip the details and directly show the following result. For

any $0 < \epsilon < N - \lambda(\overline{M} + \overline{R})$, we know that $l(N - \epsilon) > 0$. Thus, we can achieve

$$\begin{aligned}
& E[K_j^H] \\
& \leq \sum_{r=0}^{\infty} \left\{ P(R' = r) \left[\left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H + \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} \right] \right\} \\
& = \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} + \sum_{r=0}^{\infty} \left\{ P(R' = r) \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H \right\} \\
& \leq \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} + \sum_{r=0}^{\infty} \left\{ P(R' \geq r) \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H \right\} \\
& \leq \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} + \sum_{r=0}^{\infty} \left\{ P(R \geq r) \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H \right\} \\
& \leq \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} + \sum_{r=0}^{r_0-1} \left\{ \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H \right\} \\
& \quad + \sum_{r=r_0}^{\infty} \left\{ \alpha e^{-\beta r} \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H \right\}.
\end{aligned}$$

Thus, given H , $E[K_j^H]$ is bounded by a constant B_H , for any j , where B_H is given as below.

$$B_H \triangleq \min_{\epsilon \in (0, N - \lambda(\overline{M} + \overline{R}))} \left\{ \begin{aligned} & \sum_{k=2}^{\infty} k^H e^{-kl(N-\epsilon)} + \\ & \sum_{r=0}^{r_0-1} \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H + \\ & \sum_{r=r_0}^{\infty} \alpha e^{-\beta r} \left(1 + \lceil \frac{(N-1)r}{\epsilon} \rceil\right)^H \end{aligned} \right\}. \quad (14)$$

Corollary 3. $E[K_j]$ is bounded by a constant B'_1 , $E[K_j^2]$ is bounded by a constant B'_2 , for any j . The expressions of B'_1 and B'_2 are shown as below, where the rate function $l(a)$ is defined in Eq. (4).

$$B'_1 = \min_{\epsilon \in (0, N - \lambda(\overline{M} + \overline{R}))} \left\{ \begin{aligned} & 2r_0 + \frac{(N-1)r_0(r_0-1)}{2\epsilon} + \\ & \frac{2\alpha}{e^\beta - 1} + \frac{\alpha(N-1)}{4\epsilon} \operatorname{csch}^2\left(\frac{\beta}{2}\right) \\ & + 2\alpha + \frac{2e^{l(N-\epsilon)} - 1}{e^{l(N-\epsilon)}(e^{l(N-\epsilon)} - 1)^2} \end{aligned} \right\},$$

$$B'_2 = \min_{\epsilon \in (0, N - \lambda(\overline{M} + \overline{R}))} \left\{ \begin{aligned} & 4r_0 + \frac{2(N-1)r_0(r_0-1)}{\epsilon} + 4\alpha \\ & + \frac{(N-1)^2(r_0-1)r_0(2r_0-1)}{6\epsilon^2} \\ & + \frac{\alpha(N-1)^2}{8\epsilon^2} \sinh(\beta) \operatorname{csch}^4\left(\frac{\beta}{2}\right) \\ & + \frac{4\alpha}{e^\beta - 1} + \frac{\alpha(N-1)}{\epsilon} \operatorname{csch}^2\left(\frac{\beta}{2}\right) \\ & + \frac{4e^{2l(N-\epsilon)} - 3e^{l(N-\epsilon)} + 1}{e^{l(N-\epsilon)}(e^{l(N-\epsilon)} - 1)^3} \end{aligned} \right\}.$$

Proof: By substituting $H = 2$ and $H = 1$ in Eq. (14), we directly achieve the result by similar steps of Corollary 1. ■

Remark 2. Following the same proof steps, we can get the same results as Theorem 2 and Theorem 3. The difference is that we need use B'_1 and B'_2 instead of B_1 and B_2 . Also, following the same steps, we can directly extend the results to more than 2 phases.

VI. SIMULATION RESULTS

A. Simulation Setting

We evaluate the convergence of efficiency ratio for both preemptive and non-preemptive scenarios. We consider a data center with $N = 100$ machines, and choose Poisson process with arrival rate $\lambda = 2$ jobs per time slot as the job arrival process. We choose uniform distribution and exponential distribution as examples of bounded workload and light-tailed distributed workload, respectively. For short, we use $Exp(\mu)$ to represent an exponential distribution with mean μ , and use $U[a, b]$ to represent a uniform distribution on $\{a, a+1, \dots, b-1, b\}$. We choose the total time slots to be $T = 800$, and the number of tasks in each job is up to 10.

We choose 4 typical schedulers to evaluate the convergence of efficiency ratio:

The **SRPT** scheduler: Jobs with smaller unfinished workload are always scheduled first. Since the exact SRPT scheduler is not always feasible in the problem with phase precedence, the SRPT scheduler here only schedules the available part of the workload.

The **FIFO** scheduler: It is the default scheduler in many systems, like Hadoop. All the jobs are scheduled in their order of arrival.

The **Fair** scheduler: It is also a widely used scheduler. The assignment of machines are scheduled to all the waiting jobs in a fair manner. However, if some jobs need fewer machines than others in each time slot, then the remaining machines are scheduled to the other jobs, to avoid resource wastage and to keep the scheduler work-conserving.

The **LRPT** scheduler: Jobs with larger unfinished workload are always scheduled first. Roughly speaking, the performance of this scheduler represents in a sense how poorly even some work-conserving schedulers can perform.

B. Convergence of Efficiency Ratio

In the simulations, the efficiency ratio of a scheduler is obtained by the total flow time of the scheduler over the lower bound of the total flow time in T time slots. Thus, the real efficiency ratio should be smaller than the efficiency ratio given in the simulations. However, the proportion of all schedulers would remain the same.

First, we evaluate the exponentially distributed workload. We choose the workload distribution of Phase 1 for each job as $Exp(5)$ and the workload distribution of Phase 2 for each job as $Exp(40)$. The convergence of efficiency ratios of different schedulers are shown in Fig. 3. For different workload, we choose workload distribution of Phase 1 as $Exp(30)$ and the

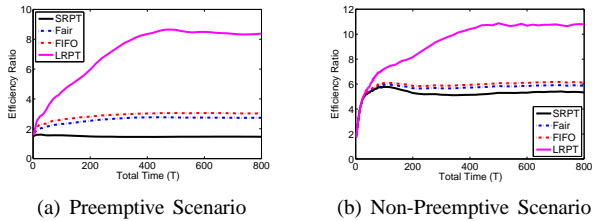


Fig. 3. Convergence of Efficiency Ratios (Exponential Distribution, Large Reduce)

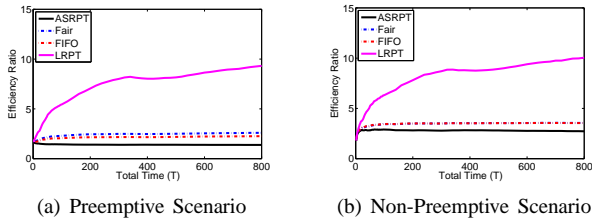


Fig. 4. Convergence of Efficiency Ratios (Exponential Distribution, Small Reduce)

workload distribution of Phase 2 as $Exp(15)$. The efficiency ratios of schedulers are shown in Fig. 4.

Then, we evaluate the uniformly distributed workload. We choose the workload distribution of Phase 1 for each job as $U[1, 9]$ and the workload distribution of Phase 2 for each job as $U[10, 70]$. The convergence of efficiency ratios of different schedulers are shown in Fig. 5. To evaluate for a smaller workload of Phase 2, we choose workload distribution of Phase 1 as $U[10, 50]$ and the workload distribution of Reduce as $U[10, 20]$. The convergence of efficiency ratios of different schedulers are shown in Fig. 6.

From Figs. 3-6, we can see that these typical work-conserving schedulers have constant efficiency ratios.

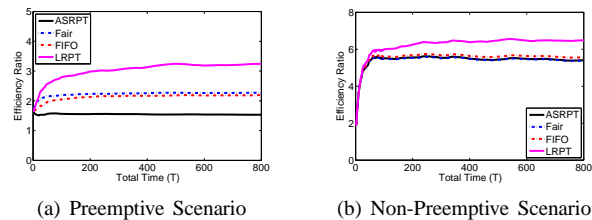


Fig. 5. Convergence of Efficiency Ratios (Uniform Distribution, Large Reduce)

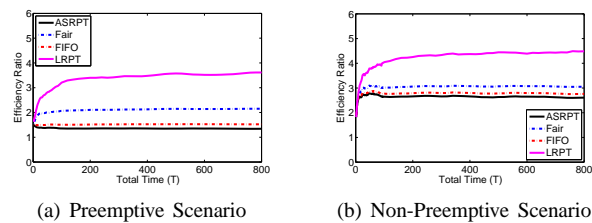


Fig. 6. Convergence of Efficiency Ratios (Uniform Distribution, Small Reduce)

VII. CONCLUSION

In this paper, we study the problem of minimizing the total flow time of a sequence of jobs with phase precedence, where the jobs arrive over time and need to be processed through Phase 1 and Phase 2 before leaving the system. Since no on-line algorithm can achieve a constant competitive ratio in some scenarios, we define weaker metric of performance called the efficiency ratio and propose a corresponding technique to analyze on-line schedulers. Under assumption of bounded workload of Phase 2, we then show a surprising property that for the flow time problem any work-conserving scheduler has a constant efficiency ratio in both preemptive and non-preemptive scenarios. Moreover, we relax the assumption to light-tailed distributed workload of Phase 2. The simulation results show that the efficiency ratios of typical work-conserving schedulers converge when time goes to infinity, for both preemptive and non-preemptive scenarios, and for both bounded and light-tailed distributed workload of Phase 2.

REFERENCES

- [1] J. Dean and S. Ghemawat, "Mapreduce: Simplified data processing on large clusters," in *Proceedings of Sixth Symposium on Operating System Design and Implementation, OSDI*, pp. 137–150, December 2004.
- [2] B. Moseley, A. Dasgupta, R. Kumar, and T. Sarlos, "On scheduling in map-reduce and flow-shops," in *Proceedings of the 23rd ACM symposium on Parallelism in algorithms and architectures, SPAA*, pp. 289–298, June 2011.
- [3] K. R. Baker and D. Trietsch, *Principles of Sequencing and Scheduling*. Hoboken, NJ, USA: John Wiley & Sons, 2009.
- [4] edited by Joseph Y-T. Leung, *Handbook of Scheduling: Algorithms, Models, and Performance Analysis*. Boca Raton, FL, USA: Chapman and Hall/CRC, 2004.
- [5] M. Caramia and M. Drozdowski, "Scheduling malleable tasks for mean flow time criterion," tech. rep., Institute of Computing Science, Poznan University of Technology, 2005.
- [6] J. Du, J. Y-T. Leung, and G. H. Young, "Scheduling chain-structured tasks to minimizing makespan and mean flow time," *Information and Computation*, vol. 92, pp. 219–236, June 1991.
- [7] J. Du, J. Y-T. Leung, and G. H. Young, "Minimizing mean flow time with release time constraint," *Theoretical Computer Science*, vol. 75, pp. 347–355, October 1990.
- [8] M. Zaharia, D. Borthakur, J. S. Sarma, K. Elmeleegy, S. Shenker, and I. Stoica, "Job scheduling for multi-user mapreduce clusters," tech. rep., University of California, Berkley, April 2009.
- [9] J. Tan, X. Meng, and L. Zhang, "Performance analysis of coupling scheduler for mapreduce/hadoop," in *Proceedings of IEEE Infocom*, March 2012.
- [10] H. Chang, M. Kodialam, R. R. Kompella, T. V. Lakshman, M. Lee, and S. Mukherjee, "Scheduling in mapreduce-like systems for fast completion time," in *Proceedings of IEEE Infocom*, pp. 3074–3082, March 2011.
- [11] F. Chen, M. Kodialam, and T. Lakshman, "Joint scheduling of processing and shuffle phases in mapreduce systems," in *Proceedings of IEEE Infocom*, March 2012.
- [12] Y. Zheng, P. Sinha, and N. Shroff, "Performance Analysis of Work-Conserving Scheduler in Minimizing Flowtime Problem with Two-Stages Precedence," tech. rep., Ohio State University, June 2012. <http://www2.ece.ohio-state.edu/~zhengy/TR3.pdf>.
- [13] A. Shwartz and A. Weiss, *Large Deviations for Performance Analysis: Queues, Communication and Computing*. New York, NY, USA: Chapman and Hall, 1995.

ACKNOWLEDGEMENTS

We gratefully acknowledge Dr. Zizhan Zheng for his helpful suggestions.