

Article

Analysis of Random Scheduling Policy of Multi-Server Parallel System with Redundancy: Quality of Exponent

Jixiang Zhang

¹ School of Information Science and Engineering, Southeast University; 230179077@seu.edu.cn

* Correspondence: 230179077@seu.edu.cn;

Abstract: A multi-server parallel system dispatches the incoming job, which contains k_n tasks into n servers. A job is considered to be computed if all the tasks associated with the job are processed. One job's tasks can be encoded into at least k_n "replicas" such that the job is considered to be served if any k_n replicas finishing computation. In this paper, we analyze the random scheduling policy of a multi-server computing system under discrete time model in terms of Quality of Exponent (QoE), which is defined as the probability exponent that a typical job can be computed within a given number of time slots. We let k_n/n be a constant. Assuming that any task of any job can be randomly dispatched by a "scheduler" to any server, and computing each task takes exactly one time slot. We divide the calculation of probability exponent into two parts, exponent of numerator and exponent of denominator. For the denominator, we give the almost exact exponent using Lagrange multiplier method, while for the numerator, an upper bound of the numerator's exponent is provided. In addition, we also express the exponent in terms of information theoretical quantities and reconsider both of exponents in the context of large deviation theory.

Keywords: probability exponent; multi-server parallel system; discrete time model; arrangement of multiple sets; large deviation theory

1. Introduction

Parallel computing has become the mainstream for most computing platforms in order to support today's increasing demand for handling large-scale data and computational workloads. Splitting the information processing job into multiple tasks and allocating them to multiple computing servers, computing in parallel reduces the system delay due to the advantage of parallel processing of job-associated tasks. Furthermore, the inherent scalability, which means that the servers can easily join into or leave from the system as the demand for computing varies, leads parallel computing to the prevailing standard for many modern computing architectures.

There are mainly two kinds of classical settings when considering parallel computing. In the Fork/Join form, there is a fork point and a join point. When K independent tasks arrive at the fork point, they are sent immediately to K parallel servers. The k -th task is sent to the respective server, where it is served in first come first served (FCFS) fashion. That is, there may be some other task dispatched to the same server and waiting for service in its queue buffer when a task is currently in service. Completed tasks are saved in synchronization queues and waiting for the finishing of their siblings. After all the tasks reach the synchronization queue, they are resynchronized in join point and then, a Fork/Join computation is considered to be completed. Another setting is called Split-Merge, which is different from the F/J queueing system in that all of the tasks of the former job need to complete service before the next job can be processed. This feature is very important in many information processing tasks, especially those computation tasks including multiple iterations, where the former outputs are the inputs of current computing. For instance, training machine learning model using gradient descent method needs to iteratively compute the gradient decreases, where the

current obtained gradient vector is the input of the next iteration. F/J systems have been analyzed in the context of queueing network models which are always complex due to the existence of multiple parallel queues correspond to the servers, which are in general correlated with each other. SM system is slightly easier to be dealt with in that the next job cannot be served before current job departs the system, so that there is no interference between different tasks from different jobs or competition for service resource which always exists in real computing systems. In [1] the author gives a detailed overview of Fork/Join system in the context of network queueing model.

Using codes, for example MDS code, to facilitate various distributed computation problems or distributed storage problems have been extensively studied in the last few years, such as [2–12]. The advantage of introducing redundancy in multi-server systems are as follows. For each incoming job, multiple servers are used to perform computing works in parallel. If the original tasks are directly processed without coding, when any task is computed successfully and service output reaches the synchronization queue in join point, the locations of this output in original sequence, need to be recorded. Besides, if a task execution fails, then the task must be handled in some later service until it has been completed successful. A task may be repeated multiple times and processed by several servers simultaneously. In that case, the task is completed successful if the task execution is successful on at least one of the servers. These repetitions greatly increase the workload of parallel computing system and consequently impair system's processing efficiency. On the other hand, due to the property of certain codes, such as MDS codes or rateless codes, a predetermined number of "replicas" are sufficient to reconstruct the original computation result. Thus, what matters is the number of replicas that are served successfully and this threshold is independent of specific tasks, which provides great facilitation during entire computations. The encoding and decoding overhead are ignored in this paper due to excluding these factors do not affect our results. However, when considering practical design of parallel computation system, they are both important factors.

Varieties of assessment metrics are used to evaluate the performance of such parallel computing systems, including the mean of response time, moments of server queue lengths in the context of QNM, the distribution of response time and their variations. For computing delays, most of the prior works on parallel service systems with redundancy have been made in the case where, exactly a single job is served every time under the continuous independent run-times model, which assumes that the time each replica required to complete their service are i.i.d. and, independent of the correspond job. In this model, lots of works have focused on characterizing the optimal mean response time of a typical job [13,14]. In addition, in [15], the authors introduce a more realistic model where service times of replicas associated with the same job can be correlated. In this model, the optimal scheduling policy is much harder to analyze. A scheduling policy identifies the number of tasks that are dispatched to each server and their service order. Parallel service system analysis based on a discrete time model was considered in [16], where the service times of replicas are i.i.d. and geometrically distributed, and where replicas can be created and deleted after each time slot, so that the number of servers working on a job differs from slot to slot. In this setting, the results show that the delay is minimized when all servers are used all of the time, and the numbers of servers handling each job are all equal or differ no more than 1. Random analysis of multi-server parallel systems in asymptotic regime is also conducted in recent works [17,18], which focus on the convergence of mean response time of a job to its maximum task's delay.

In this paper, we consider a random scheduling of a multi-server computing system in a simplified way where we are not using the stochastic job arrivals model and assume that the arrival is predictable. This can be achieved by using a large storage cache to save the stochastic arrival jobs. Suppose that a "scheduler" dispatches the tasks from different jobs to n servers in a random manner. We make a simple random analysis for a multi-server system in an asymptotic perspective under the discrete processing time model. Unlike the common performance metrics used before, in this paper we intend to characterize the exponent of the probability that a typical job is completed within a given number of time slots. We define the obtained exponent as system's Quality of Exponent (QoE).

The rest of this paper is organized as follows. In Section II, we describe the multi-server computing system model and give the problem formulation. In Section III.A, first we give the definition of QoE metric. Then, Lagrange multiplier method is used in III.B and almost accurate exponent of denominator is derived explicitly. For the exponent of numerator, in Section III.C, we propose an upper bound in virtue of solving an arrangement of multiple sets problem. In the end of this section, we express the exponent in terms of information theoretical quantities. Both of exponents of denominator and numerator are reconsidered from a large deviation point of view in Section IV. Finally, in Section IV, we conclude this paper and discuss some problems in future works.

2. Multi-Server Computing System Model and Problem Formulation

Assume that the arrived jobs are saved in large caches. Each job J_l contains k_n identical tasks $\{TK(l, j), j \in [k_n]\}$. The parallel system processes these tasks using n servers $s_i, i \in [n]$ and the system computes M jobs at a time. A job J_l is considered to be completed if all the tasks associated with it are computed completely. Meanwhile, coding techniques can be used, where “replicas”, other than original tasks, are computed. By introducing extra redundancy, a job is considered to be completed as long as at least k_n replicas associated with it are processed.

In this paper, we assume that there is an auxiliary “scheduler”, who randomly dispatches each task to a server independently according to the following predetermined probability distributions

$$\{p_l, l \in [M]\}, \{q_t, t \in [n]\}, \quad (1)$$

satisfying

$$\sum_{l=1}^M p_l = 1 \quad \text{and} \quad \sum_{t=1}^n q_t = 1, \quad (2)$$

where

$$[M] = \{1, \dots, M\}, [n] = \{1, \dots, n\} \quad (3)$$

Assume that the computing system was empty before computation starts. At the beginning of each time slot, the scheduler chooses a task coming from job J_l according to the probability distribution $\{p_l, l \in [M]\}$, and with probability q_t , this task is allocated to the server s_t . Probability distributions $\{p_l, l \in [M]\}$ show the possible difference between M jobs, such as different response priorities. For simplicity, we consider a homogeneous job model in this paper, i.e., $p_l = 1/M, l \in [M]$. Multiple job classification cases have been considered in a few of earlier works [19] and [20].

Dispatching a task to a potential server in a deterministic manner is inefficient in that a task may be enqueued at a server, whereas other servers are idle. Instead, we assume a random dispatching method and the optimal scheduling policy in the sense of QoE can be obtained by maximizing the QoE metric of system over all the feasible schemes. Generally, the distribution $\{q_t, t \in [n]\}$ may depend on the current state of each server and need to be adjusted in real time, such as its workload, the duration of service time. It is very likely that the longer the server works, the worse the performance gets. For example, we may regard probability q_t as a nonincreasing function $g(L_{que})$ of queue length L_{que} in the queue buffer of server s_t .

For each server s_t , dividing its processing timeline into multiple identical time slots, in which we are interested is the QoE of system for a typical job, e.g., the exponent of the probability that a job is processed spending time no longer than T time slots using n parallel servers. In the following, we give its formal definition.

Definition 1. *The Quality of Exponent (QoE) of a parallel system in discrete time model is the server-number n normalized probability exponent that a typical job is completed within T time slots,*

$$QoE = \frac{1}{n} \log \Pr \{a \text{ typical job is completed within } T \text{ time slots}\}. \quad (4)$$

3. QoE Metric in Random Scheduling

3.1. Probability of completion of job J_l within T time slots

In this setting, for any time slot of a server s_t , either exactly a task $TK(l, j)$ of job J_l is processed or this time slot is occupied by a task coming from another job. With probability q_t/M , a task of job J_l is delivered to a time slot of server s_t . Then, the probability that job J_l is finished within T time slots equals the probability that there are at least k_n replicas processed in no more than T time slots. We have

$$\begin{aligned} & \Pr\{\text{job } J_l \text{ is completed within } T \text{ time slots}\} \\ &= \frac{\sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \prod_{i=1}^n (q_i/M)^{d_i}}{\sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i/M)^{d_i}} \end{aligned} \quad (5)$$

Since in this paper we consider the homogeneous job model, on average there are same number of tasks coming from M different jobs in every server. Then, if job J_l is completed in no more than T time slots, in each server, less than $\lceil T/M \rceil$ time slots are occupied by tasks from job J_l . We sum up the probabilities of all the feasible dispatches in numerator, which satisfy the condition that at each server, at most $\lceil T/M \rceil$ time slots are selected for tasks coming from J_l . In denominator, all the dispatching policies using no more than nT time slot are added up.

3.2. Almost Exact QoE Metric Using Lagrange Multiplier Method and KKT Conditions

In this part, we determine the almost exact probability exponent using Lagrange multiplier method and KKT conditions.

First, we have a brief discussion on the order of magnitude of T before providing the detailed derivation. Since scheduler chooses a job from the M bulk arrivals in a random manner and in this paper we consider the homogeneous job model, it is reasonable that typically, a task from job J_l will be selected every M time slots. Suppose what the first task scheduler chooses is from job J_l . Then, before the k th selection, on average, $nT - (k-1)M$ slots remain idle. At least one time slot is available for the last task of job J_l to be served, so we have

$$nT - (\tilde{k} - 1)M \geq 1. \quad (6)$$

T has the lower bound

$$T \geq \left\lceil \frac{1 + (\tilde{k} - 1)M}{n} \right\rceil. \quad (7)$$

On the other hand, in the worst case, all the tasks from all jobs are delivered to a single server. In this case, we have

$$T \leq \tilde{k}M. \quad (8)$$

Practically, in order to maintain computing system stable, the number of bulk arrival jobs M is always bounded. Here, we let $M = O(1)$, and k_n/n be a constant. Hence, the value of T satisfies

$$\left\lceil \frac{1 + (\tilde{k} - 1)M}{n} \right\rceil \leq T \leq \tilde{k}M, \quad (9)$$

The above estimation tells us $O(1) \leq T \leq O(n)$, and with large probability T reaches the lower bound. This result is useful in the calculation of the probability exponent.

Theorem 2. *The exponent of denominator in (5) satisfies*

$$\begin{aligned} & \log \left[\sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i / M)^{d_i} \right] \\ & \approx E_{de}(d_1^*, \dots, d_n^*, \log M), \end{aligned} \quad (10)$$

where function E_{de} is defined as

$$\begin{aligned} E_{de}(d_1, \dots, d_n, \lambda) = & \frac{1}{2} \log(2\pi\tilde{k}) + \tilde{k} \log \tilde{k} - \sum_{i=1}^n \left[\frac{1}{2} \log(2\pi d_i) + d_i \log d_i \right] \\ & + \sum_{i=1}^n d_i \log(q_i / M) + \lambda (\sum_{i=1}^n d_i - \tilde{k}), \end{aligned} \quad (11)$$

and $d_1^*, d_2^*, \dots, d_n^*$ are asymptotically determined by

$$\frac{d_i^*}{\tilde{k}} = q_i, \quad i \in [n]. \quad (12)$$

We shall give the detailed derivation in Appendix A. In (10) we show that the exponent of denominator approximately equals to E_{de} , because when using Lagrange multiplier method we relax the condition where d_i / \tilde{k} can be any real number.

For the exponent of numerator of probability expression (5), we have following Theorem.

Theorem 3. *The exponent of (5) satisfies*

$$\begin{aligned} & \log \left[\sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \prod_{i=1}^n (q_i / M)^{d_i} \right] \\ & \approx E_{ne}(d_1^*, \dots, d_n^*, \lambda^*, \mu_1^*, \dots, \mu_n^*) + \log O(n^2), \end{aligned} \quad (13)$$

where function $E_{ne}(d_1, \dots, d_n, \lambda, \mu_1, \dots, \mu_n)$ is defined as

$$\begin{aligned} E_{ne}(d_1, \dots, d_n, \lambda, \mu_1, \dots, \mu_n) = & \frac{1}{2} \log(2\pi\tilde{k}) + \tilde{k} \log \tilde{k} \\ & - \sum_{i=1}^n \left[\frac{1}{2} \log(2\pi d_i) + d_i \log d_i \right] + \sum_{i=1}^n d_i \log(q_i / M) \\ & + \lambda (\sum_{i=1}^n d_i - \tilde{k}) + \sum_{i=1}^n \mu_i (d_i - \lceil T/M \rceil), \end{aligned} \quad (14)$$

and the parameters $d_1, \dots, d_n, \lambda, \mu_1, \dots, \mu_n$ are determined by

$$\begin{cases} \frac{\partial E_{ne}}{\partial d_i} = \frac{1}{2 \ln 2} \left(\frac{1}{\tilde{k}} - \frac{1}{d_i} \right) + \log \frac{\tilde{k}}{d_i} + \log \frac{q_i}{M} + \lambda + \mu_i = 0 & i \in [n] \\ \mu_i (d_i - \lceil T/M \rceil) = 0 & i \in [n] \\ \frac{\partial E_{ne}}{\partial \lambda} = \sum_{i=1}^n d_i - \tilde{k} = 0 \end{cases} \quad (15)$$

The calculation of this exponent can be found in Appendix B. Provided these results, we immediately have

$$QoE = \frac{1}{n} \log \Pr\{\text{job } J_l \text{ is completed within } T \text{ time slots}\} \approx \frac{1}{n} (E_{ne} - E_{de}). \quad (16)$$

3.3. An Upper Bound of Numerator's Exponent

Explicit optimal exponent of numerator cannot be obtained easily by solving (15). Instead, in this subsection we establish an upper bound.

We know that the numerator is less than

$$O(n^2) \max_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} N_2 \cdot \exp_2 \left\{ \sum_{i=1}^n d_i \log(q_i/M) \right\}, \quad (17)$$

where

$$N_2 = \left| \left\{ (d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\} \right\} \right|. \quad (18)$$

This quantity can be regarded as the number of nonnegative integer solutions (y_1, y_2, \dots, y_n) of following equation

$$y_1 + y_2 + \dots + y_n = \tilde{k}, \quad (19)$$

under conditions where

$$0 \leq y_i \leq \lceil T/M \rceil, \quad i \in [n]. \quad (20)$$

In addition, since the servers are in general heterogeneous, we should consider the partition's permutation, other than combination. Furthermore, in this problem we do not assign the size of each partitioned set, so that the number N_2 is independent of specific assignment of $y_i, i \in [n]$. Thus, the numerator is upper bounded by

$$O(n^2) N_2 \max_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \exp_2 \left\{ \sum_{i=1}^n d_i \log(q_i/M) \right\}. \quad (21)$$

Obtaining the exact value of N_2 and write it as an exponential is not easy, therefore we postpone the derivation to Appendix C. Here, we first give a simple upper bound of the maximization problem in (21) and then describe our upper bound.

$$\begin{aligned} & \max_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \exp_2 \left\{ \sum_{i=1}^n d_i \log(q_i/M) \right\} \\ & \leq \exp_2 \left\{ \sum_{j=1}^{\lceil \tilde{k}M/T \rceil} (\lceil T/M \rceil) \log(q_{i_j}/M) \right\}. \end{aligned} \quad (22)$$

Because $d_i \leq \lceil T/M \rceil$ for all $i \in [n]$, in order to get the above upper bound, we apply a simple greedy strategy. Rearrange $\{q_i, i \in [n]\}$ in a decreasing order as

$$q_{i_1} \geq q_{i_2} \geq \dots \geq q_{i_n}, \quad (23)$$

at most $\lceil \tilde{k}M/T \rceil$ terms in summation if all the corresponding d_{i_j} 's achieve $\lceil T/M \rceil$.

Theorem 4. N_2 defined in (18) satisfies

$$N_2 \leq n! \cdot \sum_{t=0}^n (-1)^t \binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + n - 1)}{n - 1} \quad (24)$$

which is upper bounded by

$$\begin{aligned} & \frac{1}{2} \frac{n}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \left(1 - \frac{n - t^*}{t^* + 1} \exp_2 \left\{ -\frac{n (\lceil T/M \rceil + 1)}{\tilde{k} - (t^* + 1) (\lceil T/M \rceil + 1) + n} \right\} \right) \\ & \cdot \exp_2 \left\{ \frac{1}{2} \log \frac{(n+1) [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}{(2\pi) t^* (n - t^*) [\tilde{k} - t^* (\lceil T/M \rceil + 1)]} + (n+1) \log(n+1) - \frac{n}{\ln 2} \right. \\ & \left. + [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] H(t^*, n - t^*, [\tilde{k} - t^* (\lceil T/M \rceil + 1)]) + O(1/n) \right\}. \quad (25) \end{aligned}$$

where $H(t^*, n - t^*, [\tilde{k} - t^* (\lceil T/M \rceil + 1)])$ denote the entropy function

$$H \left(\frac{t^*}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n}, \frac{n - t^*}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n}, \frac{\tilde{k} - t^* (\lceil T/M \rceil + 1)}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \right) \quad (26)$$

and t^* can be obtained by solving equation

$$\log \frac{t^* + 1}{n - t^*} + \sum_{l=0}^{\lceil T/M \rceil} \log \left(1 + \frac{n}{\tilde{k} - t^* (\lceil T/M \rceil + 1) - l} \right) = 0. \quad (27)$$

Combined Theorem 4 and simple upper bound (22), we have our upper bound of numerator's exponent.

Corollary 5. For the numerator of (5), we have

$$\begin{aligned} & \log \left[\sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \prod_{i=1}^n (q_i/M)^{d_i} \right] \\ & \leq \log A + B + C + \log O(n^2) \end{aligned} \quad (28)$$

the parameter A , B and C represent the follows

$$A = \frac{1}{2} \frac{n}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \left(1 - \frac{n - t^*}{t^* + 1} \exp_2 \left\{ -\frac{n (\lceil T/M \rceil + 1)}{\tilde{k} - (t^* + 1) (\lceil T/M \rceil + 1) + n} \right\} \right) \quad (29)$$

$$\begin{aligned} B = & \frac{1}{2} \log \frac{(n+1) [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}{(2\pi) t^* (n - t^*) [\tilde{k} - t^* (\lceil T/M \rceil + 1)]} + (n+1) \log(n+1) - \frac{n}{\ln 2} \\ & + [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] H(t^*, n - t^*, [\tilde{k} - t^* (\lceil T/M \rceil + 1)]) + O(1/n), \quad (30) \end{aligned}$$

where the entries in entropy function should be normalized, and

$$C = \sum_{j=1}^{\lceil \tilde{k}M/T \rceil} (\lceil T/M \rceil) \log(q_j/M). \quad (31)$$

Remember that we let k_n/n be a constant. when server number n tends to infinity, $(1/n) \log A$ vanishes if $n \rightarrow \infty$.

Similarly,

$$\lim_{n \rightarrow \infty} \frac{1}{n} B = \log \frac{n+1}{e} + \left[1 + \frac{\tilde{k} - t^* (\lceil T/M \rceil + 1)}{n} \right] H(t^*, n - t^*, \tilde{k} - t^* (\lceil T/M \rceil + 1)). \quad (32)$$

and

$$\lim_{n \rightarrow \infty} \frac{1}{n} C \leq \frac{\tilde{k}}{n} \max_{t \in [n]} \log(q_t/M). \quad (33)$$

From our upper bound, we obtain

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} E_{de} \leq \log \frac{n+1}{e} + \left[1 + \frac{\tilde{k} - t^*(\lceil T/M \rceil + 1)}{n} \right] H(t^*, n - t^*, \tilde{k} - t^*(\lceil T/M \rceil + 1)) \\ + \frac{\tilde{k}}{n} \max_{t \in [n]} \log(q_t/M). \quad (34) \end{aligned}$$

3.4. QoE Metric Using Information Theoretical Quantities

In the last part of this section, we express exponents of both numerator and denominator using information theoretical quantities. We describe the results in the following corollary.

Corollary 6. *Let the exponents of denominator and numerator be E_{de} and E_{ne} , respectively. Then we have the their upper bound*

$$\begin{aligned} \max_{(d_1, d_2, \dots, d_n)} \left[\frac{n}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) - \tilde{k} D \left(\left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \parallel (q_1, \dots, q_n) \right) \right. \\ \left. - \frac{n}{2} \log \frac{2\pi\tilde{k}}{n} - \tilde{k} \log M + \frac{1}{2} \log(2\pi\tilde{k}) + \frac{1}{(12\tilde{k}) \ln 2} + \log O(n^2) \right], \quad (35) \end{aligned}$$

and lower bound

$$\begin{aligned} \max_{(d_1, d_2, \dots, d_n)} \left[\frac{n}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) - \tilde{k} D \left(\left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \parallel (q_1, \dots, q_n) \right) \right. \\ \left. - \frac{n}{2} \log \frac{2\pi\tilde{k}}{n} - \tilde{k} \log M + \frac{1}{2} \log(2\pi\tilde{k}) - \frac{1}{(12) \ln 2} \sum_{i=1}^n \frac{1}{d_i} \right]. \quad (36) \end{aligned}$$

according to different maximization problems. For exponent E_{de} , d_1, d_2, \dots, d_n can take any integer from 0 to \tilde{k} , while for the exponent E_{ne} , each d_i cannot be larger than $\lceil T/M \rceil$.

Proof. We first handle the denominator. As shown above,

$$\begin{aligned} \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i/M)^{d_i} \\ \leq O(n^2) \max_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \binom{\tilde{k}}{d_1 \ d_2 \ \dots \ d_n} \exp_2 \left\{ \sum_{i=1}^n d_i \log \frac{q_i}{M} \right\} \quad (37) \end{aligned}$$

$$= O(n^2) \max_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \exp_2 \left\{ \log \tilde{k}! - \sum_{i=1}^n d_i! + \sum_{i=1}^n d_i \log \frac{q_i}{M} \right\} \quad (38)$$

Using Sterling approximation, we have

$$\begin{aligned} \log \tilde{k}! - \sum_{i=1}^n \log d_i! \\ \leq \frac{1}{2} \log(2\pi\tilde{k}) + \tilde{k} \log \tilde{k} - \sum_{i=1}^n \left[\frac{1}{2} \log(2\pi d_i) + d_i \log d_i \right] + \frac{1}{(12\tilde{k}) \ln 2} \quad (39) \end{aligned}$$

$$= \frac{1}{2} \log(2\pi\tilde{k}) - \frac{1}{2} \sum_{i=1}^n \log(2\pi d_i) - \sum_{i=1}^n d_i \log \frac{d_i}{\tilde{k}} + \frac{1}{(12\tilde{k}) \ln 2} \quad (40)$$

$$= \frac{1}{2} \log(2\pi\tilde{k}) - \frac{1}{2} \sum_{i=1}^n \log(2\pi d_i) + \tilde{k} H \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) + \frac{1}{(12\tilde{k}) \ln 2} \quad (41)$$

where the second term in (41) is tackled below.

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^n \log(2\pi d_i) \\ &= \frac{1}{2} \sum_{i=1}^n \log \frac{d_i}{\tilde{k}} + \frac{1}{2} \sum_{i=1}^n \log(2\pi\tilde{k}) \end{aligned} \quad (42)$$

$$= \frac{n}{2} \sum_{i=1}^n \frac{1}{n} \log \frac{d_i}{\tilde{k}} + \frac{n}{2} \log(2\pi\tilde{k}) \quad (43)$$

$$= \frac{n}{2} \sum_{i=1}^n \frac{1}{n} \log \frac{d_i/\tilde{k}}{1/n} + \frac{n}{2} \sum_{i=1}^n \frac{1}{n} \log \frac{1}{n} + \frac{n}{2} \log(2\pi\tilde{k}) \quad (44)$$

$$= -\frac{n}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) + \frac{n}{2} \log \frac{2\pi\tilde{k}}{n} \quad (45)$$

Substituting (45) into (41), we obtain

$$\begin{aligned} & \log \tilde{k}! - \sum_{i=1}^n \log d_i! \\ & \leq \frac{1}{2} \log(2\pi\tilde{k}) + \frac{n}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) - \frac{n}{2} \log \frac{2\pi\tilde{k}}{n} + \tilde{k} H \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) + \frac{1}{(12\tilde{k}) \ln 2}, \end{aligned} \quad (46)$$

where U_n denotes the uniform distribution $(1/n, 1/n, \dots, 1/n)$.

Next, for the last term in (38),

$$\begin{aligned} & \sum_{i=1}^n d_i \log \frac{q_i}{M} \\ &= \sum_{i=1}^n d_i \log q_i - \tilde{k} \log M \end{aligned} \quad (47)$$

$$= \tilde{k} \sum_{i=1}^n \frac{d_i}{\tilde{k}} \log q_i - \tilde{k} \log M \quad (48)$$

$$= \tilde{k} \sum_{i=1}^n \frac{d_i}{\tilde{k}} \log \frac{q_i}{(d_i/\tilde{k})} + \tilde{k} \sum_{i=1}^n \frac{d_i}{\tilde{k}} \log \frac{d_i}{\tilde{k}} - \tilde{k} \log M \quad (49)$$

$$= -\tilde{k} D \left(\left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \parallel (q_1, \dots, q_n) \right) - \tilde{k} H \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) - \tilde{k} \log M. \quad (50)$$

Combining two part of exponent, the denominator's exponent is upper bounded by

$$\begin{aligned} & \max_{(d_1, d_2, \dots, d_n)} \frac{n}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) - \tilde{k} D \left(\left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \parallel (q_1, \dots, q_n) \right) \\ & \quad - \frac{n}{2} \log \frac{2\pi\tilde{k}}{n} - \tilde{k} \log M + \frac{1}{2} \log(2\pi\tilde{k}) + \frac{1}{(12\tilde{k}) \ln 2} + \log O(n^2), \end{aligned} \quad (51)$$

where d_1, d_2, \dots, d_n can take any integer from 0 to \tilde{k} .

Similarly, we can also give a lower bound,

$$\begin{aligned} & \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i/M)^{d_i} \\ & \geq \max_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \binom{\tilde{k}}{d_1 \ d_2 \ \dots \ d_n} \exp_2 \left\{ \sum_{i=1}^n d_i \log \frac{q_i}{M} \right\}. \end{aligned} \quad (52)$$

using the Sterling approximation, we obtain that the exponent of (52) is lower bounded by

$$\max_{(d_1, d_2, \dots, d_n)} \left\{ \frac{n}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) - \tilde{k} D \left(\left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \parallel (q_1, \dots, q_n) \right) - \frac{n}{2} \log \frac{2\pi\tilde{k}}{n} - \tilde{k} \log M + \frac{1}{2} \log(2\pi\tilde{k}) - \frac{1}{(12) \ln 2} \sum_{i=1}^n \frac{1}{d_i} \right\}. \quad (53)$$

The exponent of numerator differs from that of denominator in that the optimization problem in numerator has extra restrictions that each d_i is not larger than $\lceil T/M \rceil$. The proof is completed. \square

If normalized by server number n and let n approach to infinity, the upper bound (51) of denominator and its lower bound (53) asymptotically meet. Thus, we obtain

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_{de} = \max_{d_1, d_2, \dots, d_n} \left\{ \frac{1}{2} D \left(U_n \parallel \left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \right) - \frac{\tilde{k}}{n} D \left(\left(\frac{d_1}{\tilde{k}}, \dots, \frac{d_n}{\tilde{k}} \right) \parallel (q_1, \dots, q_n) \right) - \frac{1}{2} \log \frac{2\pi\tilde{k}}{n} - \frac{\tilde{k}}{n} \log M \right\}. \quad (54)$$

Recall that in Theorem 2, the maxima of (54) are reached by setting $d_i/\tilde{k} = q_i$, so that we have following result,

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_{de} = \frac{1}{2} D(U_n \parallel (q_1, \dots, q_n)) - \frac{1}{2} \log \frac{2\pi\tilde{k}}{n} - \frac{\tilde{k}}{n} \log M. \quad (55)$$

Eventually, combing (34) and (55) we derive that

$$\lim_{n \rightarrow \infty} QoE \approx \lim_{n \rightarrow \infty} \frac{1}{n} (E_{ne} - E_{de}) \quad (56)$$

$$\leq \log \frac{n+1}{e} - \frac{1}{2} D(U_n \parallel (q_1, \dots, q_n)) + D, \quad (57)$$

where D is a constant not depending on n ,

$$D = \left[1 + \frac{\tilde{k} - t^* (\lceil T/M \rceil + 1)}{n} \right] H(t^*, n - t^*, [\tilde{k} - t^* (\lceil T/M \rceil + 1)]) + \frac{\tilde{k}}{n} \max_{t \in [n]} \log q_t + \frac{1}{2} \log \frac{2\pi\tilde{k}}{n}. \quad (58)$$

4. A Large Deviation Theory Analysis

In this section, we reconsider the probabilities of denominator and numerator of (5) in a large deviation theory perspective. From Theorem 2 we know that for the denominator

$$\log \left[\sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i/M)^{d_i} \right] \approx E_{de}(d_1^*, \dots, d_n^*, \log M),$$

where (d_1, d_2, \dots, d_n) satisfy $d_i^*/\tilde{k} = q_i, i \in [n]$.

Denote $P_d = (d_1/\tilde{k}, d_2/\tilde{k}, \dots, d_n/\tilde{k})$ as the dispatch type of the allocating scheme adopted by multi-server system. In fact, Theorem 2 tells that the probability of denominator is mainly determined by those dispatch types approach probability distribution $\{q_t, t \in [n]\}$, which is a direct result of Weak law of large Numbers.

While for the probability of numerator, there are other restrictions for the selections of dispatch policies where the number of tasks in each server cannot exceed $\lceil T/M \rceil$. This is a large deviation and

the probability exponent can be determined by following Sanov Theorem.

Theorem 7. Assume that X_1, X_2, \dots, X_n are i.i.d. random variables that follow a probability distribution $Q(x)$. \mathcal{P} denotes the collection of all the probability distributions and \mathcal{P}_n is the collection of all the types. If $E \subseteq \mathcal{P}$, then we have

$$Q^n(E) = Q^n(E \cap \mathcal{P}_n) \leq (n+1)^{|\mathcal{X}|} 2^{-nD(P^* \| Q)}, \quad (59)$$

where

$$P^* = \arg \min_{P \in E} D(P \| Q), \quad (60)$$

is the distribution which is closest to Q in the probability distribution collection E . Furthermore, if the interior of E is not empty, then

$$\frac{1}{n} \log Q^n(E) \rightarrow -D(P^* \| Q), \quad (61)$$

Provided this Theorem, for the probability exponent of the numerator of (5), we have following Corollary.

Corollary 8. Denote E as the collection of probability distributions defined below,

$$E = \left\{ P_d(P_{d,1}, P_{d,2}, \dots, P_{d,n}) : \tilde{k} \cdot \max_{i \in [n]} P_{d,i} \leq \lceil T/M \rceil \right\}, \quad (62)$$

Then the exponents of (52) and (53) for the numerator are the same when $n \rightarrow \infty$ and are achieved when the dispatch type P_d^* satisfies

$$P_d^* = \arg \min_{P_d \in E} -D(P_d \| (q_1, q_2, \dots, q_n)) \quad (63)$$

This corollary is a direct application of Sanov Theorem. Recall that in Theorem 3, we can now claim that the exponent E_{ne} define in (14) reaches its maxima when set $d_i/\tilde{k}, i \in [n]$ according to (64).

5. Conclusion and Discussion

In this paper, we considered the random scheduling problem in parallel computing with redundancy. Based on a discrete time model, we characterized the exponent of the probability that a typical job is completed within time slots. In this work, we gave an upper bound for the exponent of numerator in probability expression to avoid using combinatorial expressions in computation. By the way, our upper bound should be improved if we choose

$$\max_t [f(t) - f(t+1) + f(t+2) - f(t+3)], \quad (64)$$

in (A39). Information quantities representations of both exponents and a large deviation theory analysis are also provided in this paper.

There are still many open problems concerning parallel computing system. Queueing Network Model based analysis is quite involved, however using discrete time model may avoid some hardship. Analyzing scheduling policies under certain discrete time assumptions using reinforcement learning is prevailing in recent works, such as minimizing age of information, which is an information timeliness metric, in various settings. A possible further work is considering the case where tasks coming from the same job are correlated, either in processing time or in the servers. Dispatching a job's task to a large number of servers shall inevitably increase communication load of the computing system. Considering communication-delay tradeoff is also meaningful for practical design of computing systems.

Funding: Please add: "This research received no external funding" or "This research was funded by NAME OF FUNDER grant number XXX." and and "The APC was funded by XXX". Check carefully that the details given

are accurate and use the standard spelling of funding agency names at <https://search.crossref.org/funding>, any errors may affect your future funding.

Acknowledgments: In this section you can acknowledge any support given which is not covered by the author contribution or funding sections. This may include administrative and technical support, or donations in kind (e.g., materials used for experiments).

Conflicts of Interest: The authors declare no conflict of interest.

Appendix A

In this appendix, we give the detailed calculation of exponent E_{de} of denominator in (5).

Define

$$N_1 = \left| \left\{ (d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N} \right\} \right|. \quad (\text{A1})$$

Then,

$$\begin{aligned} & \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i/M)^{d_i} \\ & \leq O(n^2) \max_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} N_1 \cdot \exp_2 \left\{ \sum_{i=1}^n d_i \log(q_i/M) \right\}. \end{aligned} \quad (\text{A2})$$

Because the external summation has at most $O(n^2)$ terms. It is known that

$$N_1 = \binom{\tilde{k}}{d_1 \ d_2 \ \dots \ d_n} = \exp_2 \left\{ \log \tilde{k}! - \sum_{i=1}^n \log d_i! \right\}. \quad (\text{A3})$$

To continue, we apply the strengthen version of Sterling formula,

$$\sqrt{2\pi n} \left(\frac{n}{e}\right)^n \leq n! \leq \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}}, \quad (\text{A4})$$

which shows that

$$\frac{1}{2} \log(2\pi n) + n \log n - \frac{n}{\ln 2} \leq \log n! \leq \frac{1}{2} \log(2\pi n) + n \log n - \frac{n}{\ln 2} + \frac{1}{12n(\ln 2)}. \quad (\text{A5})$$

when n tends to be large, we have

$$\log n! \approx \frac{1}{2} \log(2\pi n) + n \log n - \frac{n}{\ln 2}. \quad (\text{A6})$$

Two expressions differ at most $O(1/n)$. Thus, we have

$$\begin{aligned} & \log \tilde{k}! - \sum_{i=1}^n \log d_i! \\ & \approx \frac{1}{2} \log(2\pi \tilde{k}) + \tilde{k} \log \tilde{k} - \sum_{i=1}^n \left[\frac{1}{2} \log(2\pi d_i) + d_i \log d_i \right]. \end{aligned} \quad (\text{A7})$$

Define Lagrange function

$$\begin{aligned} E_{de}(d_1, \dots, d_n, \lambda) = & \frac{1}{2} \log(2\pi \tilde{k}) + \tilde{k} \log \tilde{k} - \sum_{i=1}^n \left[\frac{1}{2} \log(2\pi d_i) + d_i \log d_i \right] \\ & + \sum_{i=1}^n d_i \log(q_i/M) + \lambda \left(\sum_{i=1}^n d_i - \tilde{k} \right), \end{aligned} \quad (\text{A8})$$

and take the partial derivatives of E_{de} with respect to d_i and λ , let

$$\begin{cases} \frac{\partial E_{de}}{\partial d_i} = \frac{1}{2\ln 2} \left(\frac{1}{\tilde{k}} - \frac{1}{d_i} \right) + \log \frac{\tilde{k}}{d_i} + \log \frac{q_i}{M} + \lambda = 0 \\ \frac{\partial E_{de}}{\partial \lambda} = \sum_{i=1}^n d_i - \tilde{k} = 0 \end{cases} \quad (\text{A9})$$

we obtain the optimal d_i 's satisfying

$$\log \frac{\tilde{k}}{d_i} = \frac{1}{(2 \ln 2) \tilde{k}} \left(\frac{\tilde{k}}{d_i} - 1 \right) + \log \frac{M}{2^\lambda q_i}. \quad (\text{A10})$$

Set

$$\frac{\tilde{k}}{d_i} = x, \quad (\text{A11})$$

then, (A10) is rewritten as

$$\log x = \frac{1}{(2 \ln 2) \tilde{k}} (x - 1) + \log \frac{M}{2^\lambda q_i}. \quad (\text{A12})$$

In the above expression, LHS is logarithmic function of x , while RHS is a line passing through fix point $(1, \log(M/2^\lambda q_i))$, whose slope $1/(2 \ln 2) \tilde{k}$ tends to be 0. So, the two curves intersect approximately at $(M/2^\lambda q_i, \log(M/2^\lambda q_i))$, which shows that the optimal d_i 's asymptotically satisfy

$$\frac{d_i^*}{\tilde{k}} = q_i, \quad i \in [n]. \quad (\text{A13})$$

Thus, the denominator is upper bounded by

$$\exp_2 \left\{ E_{de}(d_1^*, \dots, d_n^*, \log M) + \log O(n^2) \right\}, \quad (\text{A14})$$

At the same time,

$$\begin{aligned} & \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \prod_{i=1}^n (q_i/M)^{d_i} \\ & \geq \max_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \mathbb{N}} \exp_2 \{ E_{de} \} \end{aligned} \quad (\text{A15})$$

$$= \exp_2 \{ E_{de}(d_1^*, \dots, d_n^*, \log M) \}. \quad (\text{A16})$$

When normalized by n and let n be sufficiently large, both upper and lower bounds are asymptotically tight.

Appendix B

In the follows, we deal with the numerator of (5).

$$\begin{aligned} & \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \prod_{i=1}^n (q_i/M)^{d_i} \\ & = \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \exp_2 \left\{ \sum_{i=1}^n d_i \log(q_i/M) \right\}. \end{aligned} \quad (\text{A17})$$

Define

$$N_2 = \left| \left\{ (d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\} \right\} \right|. \quad (\text{A18})$$

As discussed before, $T \leq O(n)$, then there are no more than $O(n^2)$ terms in external summation, (A17) is less than

$$O(n^2) \max_{(d_1, d_2, \dots, d_n): \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} N_2 \cdot \exp_2 \left\{ \sum_{i=1}^n d_i \log(q_i/M) \right\}. \quad (\text{A19})$$

Define Lagrange function

$$E_{ne}(d_1, \dots, d_n, \lambda, \mu_1, \dots, \mu_n) = \frac{1}{2} \log(2\pi\tilde{k}) + \tilde{k} \log \tilde{k} - \sum_{i=1}^n \left[\frac{1}{2} \log(2\pi d_i) + d_i \log d_i \right] + \sum_{i=1}^n d_i \log(q_i/M) + \lambda(\sum_{i=1}^n d_i - \tilde{k}) + \sum_{i=1}^n \mu_i(d_i - \lceil T/M \rceil). \quad (\text{A20})$$

Due to the KKT condition, the optimal solution (d_1^*, \dots, d_n^*) and other parameters satisfying

$$\begin{cases} \frac{\partial E_{ne}}{\partial d_i} = \frac{1}{2\ln 2} \left(\frac{1}{\tilde{k}} - \frac{1}{d_i^*} \right) + \log \frac{\tilde{k}}{d_i^*} + \log \frac{q_i}{M} + \lambda^* + \mu_i^* = 0 & i \in [n] \\ \mu_i^*(d_i^* - \lceil T/M \rceil) = 0 & i \in [n] \\ \frac{\partial E_{ne}}{\partial \lambda} = \sum_{i=1}^n d_i^* - \tilde{k} = 0 \end{cases}. \quad (\text{A21})$$

In fact, we have the following result,

$$\begin{aligned} & \exp_2 \{E_{ne}(d_1^*, \dots, d_n^*, \lambda, \mu_1^*, \dots, \mu_n^*)\} \\ & \leq \sum_{\tilde{k}=k_n}^{nT} \sum_{(d_1, d_2, \dots, d_n) : \sum_{i=1}^n d_i = \tilde{k}, d_1, d_2, \dots, d_n \in \{0, 1, \dots, \lceil T/M \rceil\}} \prod_{i=1}^n (q_i/M)^{d_i} \end{aligned} \quad (\text{A22})$$

$$\leq \exp_2 \left\{ E_{ne}(d_1^*, \dots, d_n^*, \lambda, \mu_1^*, \dots, \mu_n^*) + \log O(n^2) \right\}. \quad (\text{A23})$$

Appendix C

We give the cumbersome derivation of computing N_2 in this appendix.

Denote \tilde{N}_2 as the combination counterpart of N_2 . We have

$$N_2 \leq n! \cdot \tilde{N}_2. \quad (\text{A24})$$

If some partitioned set is empty, then N_2 is strictly smaller than $n! \cdot \tilde{N}_2$. Because in these cases, less than n partitioned set need be permuted. It is glad that \tilde{N}_2 can be determined exactly using Inclusion-Exclusion Principle.

Let S be the set of all nonnegative integer solutions of equation (19) without any other restriction. It is easy to obtain

$$|S| = \binom{\tilde{k} + n - 1}{n - 1}. \quad (\text{A25})$$

Denote P_i is the property that $d_i \geq \lceil T/M \rceil + 1$, let

$$A_i = \{(d_1, \dots, d_n) \in S : (d_1, \dots, d_n) \text{ satisfies } P_i\}, \quad i \in [n]. \quad (\text{A26})$$

then we have

$$N_2 = |\bar{A}_1 \cap \bar{A}_2 \cap \dots \cap \bar{A}_n|, \quad (\text{A27})$$

According to inclusion-exclusion principle,

$$\begin{aligned} |\bar{A}_1 \cap \bar{A}_2 \cap \dots \cap \bar{A}_n| &= |S| - \sum |A_i| + \sum |A_i \cap A_j| \\ &\quad - \sum |A_i \cap A_j \cap A_k| + \dots + (-1)^n |A_1 \cap A_2 \cap \dots \cap A_n|, \end{aligned} \quad (\text{A28})$$

where $|A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_m}|$ is the number of nonnegative integer solutions that are coincident with properties $P_{i_1}, P_{i_2}, \dots, P_{i_m}$. Substituting y_1, y_2, \dots, y_n by z_1, z_2, \dots, z_n ,

$$z_i = \begin{cases} y_i - (\lceil T/M \rceil + 1), & i \in \{i_1, i_2, \dots, i_m\} \\ y_i, & \text{otherwise} \end{cases}, \quad (\text{A29})$$

the equation turns

$$z_1 + z_2 + \dots + z_n = \tilde{k} - m(\lceil T/M \rceil + 1). \quad (\text{A30})$$

The number of nonnegative integer solutions of this new equation is exactly that of original equation with given restrictions. This quantity is

$$\binom{\tilde{k} - m(\lceil T/M \rceil + 1) + n - 1}{n - 1}. \quad (\text{A31})$$

Provided above results, we can compute \tilde{N}_2 now.

$$\begin{aligned} \tilde{N}_2 = |\bar{A}_1 \cap \bar{A}_2 \cap \dots \cap \bar{A}_n| &= |S| - \sum |A_i| + \sum |A_i \cap A_j| \\ &\quad - \sum |A_i \cap A_j \cap A_k| + \dots + (-1)^n |A_1 \cap A_2 \cap \dots \cap A_n|, \end{aligned} \quad (\text{A32})$$

which equals

$$\begin{aligned} &\binom{\tilde{k} + n - 1}{n - 1} - \binom{n}{1} \binom{\tilde{k} - (\lceil T/M \rceil + 1) + n - 1}{n - 1} + \binom{n}{2} \binom{\tilde{k} - 2(\lceil T/M \rceil + 1) + n - 1}{n - 1} \\ &\quad - \binom{n}{3} \binom{\tilde{k} - 3(\lceil T/M \rceil + 1) + n - 1}{n - 1} + \dots + (-1)^n \binom{n}{n} \binom{\tilde{k} - n(\lceil T/M \rceil + 1) + n - 1}{n - 1} \end{aligned} \quad (\text{A33})$$

$$= \sum_{t=0}^n (-1)^t \binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1}{n - 1}. \quad (\text{A34})$$

From (A34) we have an upper bound of N_2 .

$$N_2 \leq n! \cdot \sum_{t=0}^n (-1)^t \binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1}{n - 1}. \quad (\text{A35})$$

It is relative cumbersome to obtain an exponential bound of (A35). We give the details below.

$$\begin{aligned} N_2 &\leq n! \sum_{t=0}^n (-1)^t \binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1}{n - 1} \\ &\leq \frac{(n+1)!}{2} \max_t \left[\binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1}{n - 1} - \binom{n}{t+1} \binom{\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n - 1}{n - 1} \right]. \end{aligned} \quad (\text{A36})$$

Define

$$f(t) = \binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1}{n - 1}, \quad (\text{A37})$$

be a function defined on $t \in [n] \cup \{0\}$. With this definition, N_2 does not greater than

$$\frac{(n+1)!}{2} \max_t [f(t) - f(t+1)]. \quad (\text{A38})$$

Furthermore,

$$\begin{aligned} f(t) &= \binom{n}{t} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1}{n - 1} \\ &= \frac{n!}{t!(n-t)!} \frac{[\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1]!}{(n-1)! [\tilde{k} - t(\lceil T/M \rceil + 1)]!} \end{aligned} \quad (\text{A39})$$

$$= n \frac{[\tilde{k} - t(\lceil T/M \rceil + 1) + n - 1]!}{t!(n-t)! [\tilde{k} - t(\lceil T/M \rceil + 1)]!} \quad (\text{A40})$$

$$= \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{[\tilde{k} - t(\lceil T/M \rceil + 1) + n]!}{t!(n-t)! [\tilde{k} - t(\lceil T/M \rceil + 1)]!} \quad (\text{A41})$$

$$= \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \binom{\tilde{k} - t(\lceil T/M \rceil + 1) + n}{t \quad n-t \quad \tilde{k} - t(\lceil T/M \rceil + 1)} \quad (\text{A42})$$

So far, we write $f(t)$ as a form of trinomial coefficient. Continue,

$$f(t) = \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{[\tilde{k} - t(\lceil T/M \rceil + 1) + n]!}{t!(n-t)! [\tilde{k} - t(\lceil T/M \rceil + 1)]!} \quad (\text{A43})$$

$$= \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{[\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n]! \prod_{l=0}^{\lceil T/M \rceil} [\tilde{k} - t(\lceil T/M \rceil + 1) + n - l]}{\frac{(t+1)!}{t+1} (n-t)! [\tilde{k} - (t+1)(\lceil T/M \rceil + 1)]! \prod_{l=0}^{\lceil T/M \rceil} [\tilde{k} - t(\lceil T/M \rceil + 1) - l]} \quad (\text{A44})$$

$$= \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{t+1}{n-t} \left(\prod_{l=0}^{\lceil T/M \rceil} \frac{\tilde{k} - t(\lceil T/M \rceil + 1) + n - l}{\tilde{k} - t(\lceil T/M \rceil + 1) - l} \right) \cdot \frac{[\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n]!}{(t+1)!(n-t-1)! [\tilde{k} - (t+1)(\lceil T/M \rceil + 1)]!} \quad (\text{A45})$$

$$= \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{t+1}{n-t} \left(\prod_{l=0}^{\lceil T/M \rceil} \frac{\tilde{k} - t(\lceil T/M \rceil + 1) + n - l}{\tilde{k} - t(\lceil T/M \rceil + 1) - l} \right) \cdot \binom{\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n}{t+1 \quad n-t-1 \quad \tilde{k} - (t+1)(\lceil T/M \rceil + 1)} \quad (\text{A46})$$

$$= \frac{n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{t+1}{n-t} \left(\prod_{l=0}^{\lceil T/M \rceil} \frac{\tilde{k} - t(\lceil T/M \rceil + 1) + n - l}{\tilde{k} - t(\lceil T/M \rceil + 1) - l} \right) \cdot \frac{\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n}{n} f(t+1) \quad (\text{A47})$$

$$= \frac{\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{t+1}{n-t} \left(\prod_{l=0}^{\lceil T/M \rceil} \frac{\tilde{k} - t(\lceil T/M \rceil + 1) + n - l}{\tilde{k} - t(\lceil T/M \rceil + 1) - l} \right) f(t+1). \quad (\text{A48})$$

That is, we have obtained the following relation

$$f(t) = \frac{\tilde{k} - (t+1)(\lceil T/M \rceil + 1) + n}{\tilde{k} - t(\lceil T/M \rceil + 1) + n} \frac{t+1}{n-t} \left(\prod_{l=0}^{\lceil T/M \rceil} \frac{\tilde{k} - t(\lceil T/M \rceil + 1) + n - l}{\tilde{k} - t(\lceil T/M \rceil + 1) - l} \right) f(t+1). \quad (\text{A49})$$

Thus,

$$\begin{aligned} & f(t) - f(t+1) \\ &= f(t) [1 - f(t+1)/f(t)] \end{aligned} \quad (\text{A50})$$

$$= f(t) \left\{ 1 - \frac{\bar{k} - t(\lceil T/M \rceil + 1) + n}{\bar{k} - (t+1)(\lceil T/M \rceil + 1) + n} \frac{n-t}{t+1} \prod_{l=0}^{\lceil T/M \rceil} \left(1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n - l} \right) \right\}. \quad (\text{A51})$$

When n tends to be large enough,

$$\frac{\bar{k} - t(\lceil T/M \rceil + 1) + n}{\bar{k} - (t+1)(\lceil T/M \rceil + 1) + n} = 1 + \frac{\lceil T/M \rceil + 1}{\bar{k} - (t+1)(\lceil T/M \rceil + 1) + n} \geq 1, \quad (\text{A52})$$

and

$$1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n - l} \leq 1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n}, \quad (\text{A53})$$

for all $l \in \{0, 1, \dots, \lceil T/M \rceil\}$. we continue the derivation as follows.

$$\begin{aligned} & \frac{(n+1)!}{2} \max_t [f(t) - f(t+1)] \\ & \leq \frac{(n+1)!}{2} \max_t f(t) \left[1 - \frac{n-t}{t+1} \left(1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1} \right]. \end{aligned} \quad (\text{A54})$$

In order to determine above maxima, set

$$g(t) = f(t) \left[1 - \frac{n-t}{t+1} \left(1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1} \right]. \quad (\text{A55})$$

We estimate

$$\frac{g(t)}{g(t+1)} = \frac{f(t)}{f(t+1)} \frac{1 - \frac{n-t}{t+1} \left(1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1}}{1 - \frac{n-t-1}{t+2} \left(1 - \frac{n}{\bar{k} - (t+1)(\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1}}. \quad (\text{A56})$$

When $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty} \frac{1 - \frac{n-t}{t+1} \left(1 - \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1}}{1 - \frac{n-t-1}{t+2} \left(1 - \frac{n}{\bar{k} - (t+1)(\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1}} = 1, \quad (\text{A57})$$

Thus, in asymptotic regime,

$$\begin{aligned} \frac{g(t)}{g(t+1)} & \approx \frac{f(t)}{f(t+1)} \\ & = \frac{\bar{k} - (t+1)(\lceil T/M \rceil + 1) + n}{\bar{k} - t(\lceil T/M \rceil + 1) + n} \frac{t+1}{n-t} \prod_{l=0}^{\lceil T/M \rceil} \frac{\bar{k} - t(\lceil T/M \rceil + 1) + n - l}{\bar{k} - t(\lceil T/M \rceil + 1) - l} \end{aligned} \quad (\text{A58})$$

$$\approx \frac{t+1}{n-t} \prod_{l=0}^{\lceil T/M \rceil} \left(1 + \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) - l} \right). \quad (\text{A59})$$

which is an increasing function of t . We can see as t increases, $g(t)/g(t+1)$ is less than 1 and when t exceeds some point t^* , the ratio becomes larger than 1. This fact shows that function first goes up and then it goes down as t increases. Let

$$\frac{g(t)}{g(t+1)} \approx \frac{t+1}{n-t} \prod_{l=0}^{\lceil T/M \rceil} \left(1 + \frac{n}{\bar{k} - t(\lceil T/M \rceil + 1) - l} \right) = 1. \quad (\text{A60})$$

we can determine the point t^* by equation

$$\log \frac{t^* + 1}{n - t^*} + \sum_{l=0}^{\lceil T/M \rceil} \log \left(1 + \frac{n}{\tilde{k} - t (\lceil T/M \rceil + 1) - l} \right) = 0. \quad (\text{A61})$$

Finally, we derive

$$N_2 \leq n! \sum_{t=0}^n (-1)^t \binom{n}{t} \binom{\tilde{k} - t (\lceil T/M \rceil + 1) + n - 1}{n-1} \quad (\text{A62})$$

$$\leq \frac{(n+1)!}{2} \max_t [f(t) - f(t+1)] \quad (\text{A63})$$

$$= \frac{(n+1)!}{2} [f(t^*) - f(t^* + 1)] \quad (\text{A64})$$

$$\leq \frac{(n+1)!}{2} f(t^*) \left\{ 1 - \frac{n - t^*}{t^* + 1} \left(1 - \frac{n}{\tilde{k} - (t^* + 1) (\lceil T/M \rceil + 1) + n} \right)^{\lceil T/M \rceil + 1} \right\} \quad (\text{A65})$$

$$\leq \frac{(n+1)!}{2} f(t^*) \left\{ 1 - \frac{n - t^*}{t^* + 1} \exp_2 \left\{ - \frac{n (\lceil T/M \rceil + 1)}{\tilde{k} - (t^* + 1) (\lceil T/M \rceil + 1) + n} \right\} \right\}. \quad (\text{A66})$$

where

$$\begin{aligned} \frac{(n+1)!}{2} f(t^*) &= \frac{(n+1)!}{2} \binom{n}{t^*} \binom{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n - 1}{n-1} \\ &= \frac{(n+1)!}{2} \frac{n}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \frac{[\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]!}{t^*! (n - t^*)! [\tilde{k} - t^* (\lceil T/M \rceil + 1)]!} \end{aligned} \quad (\text{A67})$$

which equals

$$\frac{1}{2} \frac{n}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \exp_2 \left\{ \sum_{b \in \{n+1, \tilde{k} - t^* (\lceil T/M \rceil + 1) + n\}} \log b! - \sum_{d \in \{t^*, n - t^*, \tilde{k} - t^* (\lceil T/M \rceil + 1)\}} \log d! \right\}. \quad (\text{A68})$$

Since

$$\frac{1}{2} \log(2\pi n) + n \log n - \frac{n}{\ln 2} \leq n! \leq \frac{1}{2} \log(2\pi n) + n \log n - \frac{n}{\ln 2} + \frac{1}{12n \ln 2}. \quad (\text{A69})$$

we have

$$\sum_{b \in \{n+1, \tilde{k} - t^* (\lceil T/M \rceil + 1) + n\}} \log b! - \sum_{d \in \{t^*, n - t^*, \tilde{k} - t^* (\lceil T/M \rceil + 1)\}} \log d!, \quad (\text{A70})$$

is upper bounded by

$$\begin{aligned} &\frac{1}{2} \log \left\{ (2\pi)^2 (n+1) [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] \right\} + (n+1) \log(n+1) \\ &\quad + [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] \log [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] - \frac{n}{\ln 2} + O(1/n) \\ &\quad - \frac{1}{2} \log \left\{ (2\pi)^3 t^* (n - t^*) [\tilde{k} - t^* (\lceil T/M \rceil + 1)] \right\} - t^* \log t^* - (n - t^*) \log(n - t^*) \\ &\quad - [\tilde{k} - t^* (\lceil T/M \rceil + 1)] \log [\tilde{k} - t^* (\lceil T/M \rceil + 1)]. \end{aligned} \quad (\text{A71})$$

$$\begin{aligned}
&= \frac{1}{2} \log \frac{(n+1) [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}{(2\pi)^{t^*} (n - t^*) [\tilde{k} - t^* (\lceil T/M \rceil + 1)]} + (n+1) \log(n+1) - \frac{n}{\ln 2} + O(1/n) \\
&\quad - t^* \log \frac{t^*}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} - (n - t^*) \log \frac{n - t^*}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \\
&\quad - [\tilde{k} - t^* (\lceil T/M \rceil + 1)] \log \frac{[\tilde{k} - t^* (\lceil T/M \rceil + 1)]}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n}. \quad (\text{A72})
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \log \frac{(n+1) [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}{(2\pi)^{t^*} (n - t^*) [\tilde{k} - t^* (\lceil T/M \rceil + 1)]} + (n+1) \log(n+1) - \frac{n}{\ln 2} \\
&\quad + [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] H(t^*, n - t^*, \tilde{k} - t^* (\lceil T/M \rceil + 1)) + O(1/n). \quad (\text{A73})
\end{aligned}$$

The entropy term denotes the following expression,

$$H \left(\frac{t^*}{[\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}, \frac{n - t^*}{[\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}, \frac{\tilde{k} - t^* (\lceil T/M \rceil + 1)}{[\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]} \right). \quad (\text{A74})$$

Substituting (A73) into (A68) eventually we obtain an upper bound of N_2 ,

$$\begin{aligned}
N_2 \leq & \frac{1}{2} \frac{n}{\tilde{k} - t^* (\lceil T/M \rceil + 1) + n} \left(1 - \frac{n - t^*}{t^* + 1} \exp_2 \left\{ - \frac{n (\lceil T/M \rceil + 1)}{\tilde{k} - (t^* + 1) (\lceil T/M \rceil + 1) + n} \right\} \right) \\
& \cdot \exp_2 \left\{ \frac{1}{2} \log \frac{(n+1) [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n]}{(2\pi)^{t^*} (n - t^*) [\tilde{k} - t^* (\lceil T/M \rceil + 1)]} + (n+1) \log(n+1) - \frac{n}{\ln 2} \right. \\
& \left. + [\tilde{k} - t^* (\lceil T/M \rceil + 1) + n] H(t^*, n - t^*, \tilde{k} - t^* (\lceil T/M \rceil + 1)) + O(1/n) \right\}. \quad (\text{A75})
\end{aligned}$$

References

1. Thomasian, A. Analysis of Fork/Join and Related Queueing Systems. *ACM Computing Surveys*, 47.
2. Dean, J.; Ghemawat, S. MapReduce: Simplified Data Processing on Large Clusters. *Communications of the ACM* **2008**, 51, 107–113.
3. Lee, K., e.a. The MDS Queue: Analysing the Latency Performance of Erasure Codes. *IEEE Transactions on Information Theory* **2017**, 63, 2822–2842.
4. Lee, K., e.a. Speeding Up Distributed Machine Learning Using Codes. *IEEE Transactions on Information Theory* **2018**, 64, 1514–1529.
5. Lee, K.C.S.; Ramchandran, K. High-dimensional coded matrix multiplication. 2017 IEEE International Symposium on Information Theory (ISIT), 2017, pp. 2418–2422.
6. Yu, Q.M.M.A.; Avestimehr, A. Polynomial Codes: an Optimal Design for High-Dimensional Coded Matrix Multiplication. *Advances in Neural Information Processing Systems* 30 (Nips 2017), 2017.
7. Baharav, T., e.a. Straggler-proofing massive-scale distributed matrix multiplication with d-dimensional product codes. *IEEE International Symposium on Information Theory (ISIT)*, 2018, pp. 1993–1997.
8. Fahim, M., e.a. On the Optimal Recovery Threshold of Coded Matrix Multiplication. 55th Annual Allerton Conference on Communication, Control, and Computing (Allerton), 2017, pp. 1264–1270.
9. Park, H.; Moon, J. Irregular Product Coded Computation for High-Dimensional Matrix Multiplication. *IEEE International Symposium on Information Theory (ISIT)*, 2019, pp. 1782–1786.
10. Park, H., e.a. Hierarchical Coding for Distributed Computing. *IEEE International Symposium on Information Theory (ISIT)*, 2018, pp. 1630–1634.
11. Li, S.M.M.A.; Avestimehr, A. Coding for Distributed Fog Computing. *IEEE Communications Magazine* **2017**, 55, 34–40.

12. Gupta, S.; Lalitha, V. Locality-Aware Hybrid Coded MapReduce for Server-Rack Architecture. *IEEE Information Theory Workshop (ITW)*, 2017, pp. 459–463.
13. Joshi, G., E.S.; Wornell, G. Efficient Redundancy Techniques for Latency Reduction in Cloud Systems. *ACM Transactions on Modeling and Performance Evaluation of Computing Systems* **2017**, *2*.
14. Sun, Y., K.C.E.; Shroff, N.B. On Delay-Optimal Scheduling in Queueing Systems with Replications. arXiv:1603:07322.
15. Gardner, K., e.a. A Better Model for Job Redundancy: Decoupling Server Slowdown and Job Size. *IEEE-ACM Transactions on Networking* **2017**, *25*, 3353–3367.
16. Borst, S., e.a. Task allocation in a multi-server system. *Journal of Scheduling* **2003**, *6*, 423–436.
17. Weina W., e.a. Delay Asymptotics and Bounds for Multi-Task Parallel Jobs. *Queueing Systems* **2019**.
18. M, Z. Delay-Optimal Policies in Partial Fork/Join Systems with Redundancy and Random Slowdowns. arXiv:1910:09602v1.
19. Baynat, B.; Dallery, Y. A Decomposition Approximation Method for Closed Queueing Networks with Fork/Join Subnetworks. In *Proceedings of the IFIP WG10.3 International Conference on Decentralized and Distributed Systems*, 1993, pp. 199–210.
20. Thomasian, A.; Bay, P. Analytic Queueing Network Models for Parallel Processing of Task Systems. *IEEE Transactions on Computers* **1986**, *35*, 1045–1054.