

© 2016 Ali Yekkehkhany

NEAR-DATA SCHEDULING FOR DATA CENTERS WITH MULTIPLE  
LEVELS OF DATA LOCALITY

BY

ALI YEKKEHKHANY

THESIS

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Electrical and Computer Engineering  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2016

Urbana, Illinois

Adviser:

Professor Yi Lu

# ABSTRACT

Data locality is a fundamental issue for data-parallel applications. Considering MapReduce in Hadoop, the map task scheduling part requires an efficient algorithm which takes data locality into consideration; otherwise, the system may become unstable under loads inside the system's capacity region and jobs may experience longer completion times which are not of interest. The data chunk needed for any map task can be in memory, on a local disk, in a local rack, in the same cluster or even in another data center. Hence, unless there has been much work on improving the speed of data center networks, different levels of service rates still exist for a task depending on where its data chunk is saved and from which server it receives service. Most of the theoretical work on load balancing is for systems with two levels of data locality including the Pandas algorithm by Xie et al. and the JSQ-MW algorithm by Wang et al., where the former is both throughput and heavy-traffic optimal, while the latter is only throughput optimal, but heavy-traffic optimal in only a special traffic load. We show that an extension of the JSQ-MW algorithm for a system with three levels of data locality is throughput optimal, but not heavy-traffic optimal for all loads, only for a special traffic scenario. Furthermore, we show that the Pandas algorithm is not even throughput optimal for a system with three levels of data locality. We then propose a novel algorithm, Balanced-Pandas, which is both throughput and heavy-traffic optimal. To the best of our knowledge, this is the first theoretical work on load balancing for a system with more than two levels of data locality. This is more challenging than two levels of data locality as a dilemma between performance and throughput emerges.

# ACKNOWLEDGMENTS

Firstly, I would like to thank Professor Yi Lu for her invaluable support, motivation, immense knowledge, and her dedication to research. She is not only a great professor, but also a great advisor, and I feel so pleased to be her student. Working under the supervision of Professor Lu was one of the most fortunate happenings in my life. She is just the right advisor for me.

I would also like to thank my mother who has always been a great support to me. I am also grateful to my sister.

# TABLE OF CONTENTS

LIST OF FIGURES . . . . .	v
LIST OF ABBREVIATIONS . . . . .	vi
CHAPTER 1 INTRODUCTION . . . . .	1
CHAPTER 2 LITERATURE REVIEW . . . . .	7
2.1 Fluid Model Planning . . . . .	7
2.2 Generalized $c\mu$ -Rule . . . . .	8
2.3 Join the Shortest Queue-MaxWeight (JSQ-MW) . . . . .	10
2.4 Priority Algorithm for Near-Data Scheduling (Pandas) . . . . .	11
CHAPTER 3 THREE LEVELS OF DATA LOCALITY . . . . .	12
3.1 The Performance versus Throughput Dilemma . . . . .	12
3.2 Equivalent Capacity Region . . . . .	14
3.3 Join the Shortest Queue-MaxWeight (JSQ-MW) . . . . .	15
3.4 Balanced-Pandas . . . . .	19
CHAPTER 4 SIMULATION RESULTS . . . . .	29
APPENDIX A THEOREM PROOFS . . . . .	32
A.1 Proof of Theorem 1 . . . . .	32
A.2 Proof of Theorem 2 . . . . .	35
A.3 Proof of Theorem 3 . . . . .	37
A.4 Proof of Theorem 4 . . . . .	43
REFERENCES . . . . .	60

# LIST OF FIGURES

1.1	Data center architecture when the data communication cost between the storage and computing centers is affordable. . . . .	1
1.2	The state-of-the-art data center architecture. . . . .	2
3.1	A system with two racks showing how performance should be sacrificed in order to achieve throughput optimality. . . . .	13
3.2	The queueing structure needed for the Balanced-Pandas algorithm. . . . .	20
3.3	The four types of servers and their load under the ideal load decomposition. . . . .	27
4.1	The mean task completion time versus the mean arrival rate for two algorithms, the Balanced-Pandas algorithm and the JSQ-MW algorithm, under a load that both algorithms minimize the mean task completion time at high loads. . . . .	30
4.2	The mean task completion time versus the mean arrival rate for three algorithms the Balanced-Pandas, JSQ-MW, and Pandas algorithms under a general load that all four kinds of servers exist in the system. . . . .	31
4.3	The performance of the JSQ-MW algorithm versus the Balanced-Pandas algorithm at high loads. . . . .	31
A.1	The queue compositions of the four types of servers in the heavy-traffic regime with $\alpha = 1, \beta = 0.8, \gamma = 0.5$ . The workload at four types of servers maintain the ratio $\alpha : \beta : \frac{\alpha\gamma}{\beta} : \gamma = 1 : 0.8 : 0.625 : 0.5$ . . . . .	43

# LIST OF ABBREVIATIONS

Pandas: Priority Algorithm for Near-Data Scheduling

JSQ-MW: Joint the Shortest Queue-MaxWeight

FCFS: First-Come-First-Served

# CHAPTER 1

## INTRODUCTION

Today's data centers need to keep pace with the explosion of data and processing the data [1]. The emergence of large data sets by social networks such as Facebook [2], Twitter [3], LinkedIn [4], health-care industry, search engines, and scientific research has pushed researchers to change the architecture of data centers in order to adapt them with the new needs of fast processing for large data sets. The architecture of a data center is mainly determined by the storage and processor units and the way these two units are connected to each other. The structure of a data center was once depicted as in Figure 1.1. The data was stored in a large storage unit, and whenever the data was needed for a job, it was fetched by the computing unit. Hence, every time that a chunk of data is needed for a process, it has to go through the network between storage and computing units. Without the existence of large data sets in the past, this structure worked well. However, with the appearance of large data sets, this structure lost its utility as the network between the two units was not capable of responding to the real-time applications.

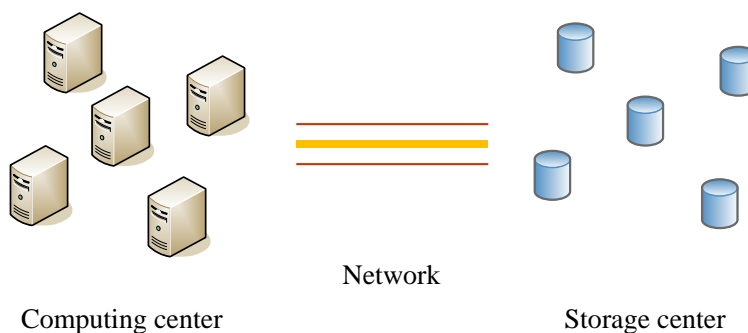


Figure 1.1: Data center architecture when the data communication cost between the storage and computing centers is affordable.

The objection to the data center architecture described above is that it is not consistent with the large data set processing applications as all the



data needed for the process should be transmitted through the network. Unless there has been a large body of research on increasing the speed of the network used in data centers, there is still a significantly large delay in data transmission compared to the service time [5–8]. Therefore, scientists changed the data center structure as the one depicted in Figure 1.2. Both the large computing and storage centers are split into smaller units, and each small computing and storage center is combined with the others that we name it a server. This way, data is moved to the computing unit and if an appropriate scheduler is used to assign tasks to servers, then very few data communication through a network is needed. Assume that there are  $M$  parallel servers in the system. The set of servers is denoted by  $\mathcal{M} = \{1, 2, 3, \dots, M\}$ .

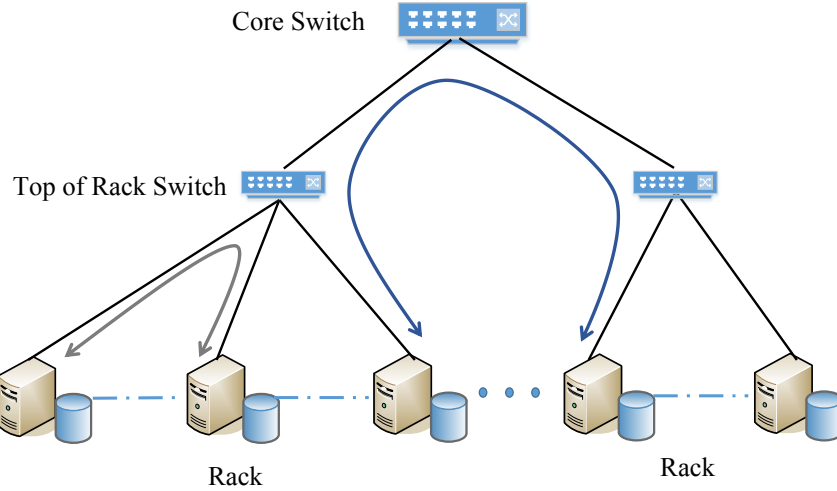


Figure 1.2: The state-of-the-art data center architecture.

A large data set is split into small chunks of ordinary sizes of 64, 128, or 256 megabytes. Each data chunk is stored on the storage of a number of servers for easier accessibility and fault resilience. If the data set is going to be processed, different servers process the data chunks stored on them, and then the results of all the servers are reduced to the final result (MapReduce). Using such an architecture for data centers, the requirement for data transmission decreases as we try to assign each server to process a task with the needed data chunk saved on the storage of the server. This concept is called *Near-Data Scheduling* as each server prefers processing a data chunk saved on itself. On the other hand, as we will see in Sections 2.3, 2.4, 3.3, and 3.4, there are cases where we still need a data chunk to be transmitted

from one server to another. In order to give the data center such a flexibility, the servers are not completely isolated from each other. Instead, there are *rack switches* on top of servers in the same rack (a *rack* consists of servers that are directly connected with each other through a switch called rack switch). Furthermore, there is a core switch which is connected to all rack switches. This structure for data centers allows the data chunks to be transmitted from a server to another server inside or outside of the rack where the data chunk is stored. The system consists of  $K$  racks, denoted with the set  $\mathcal{K} = \{1, 2, 3, \dots, K\}$ . A server  $m$  belongs to a rack denoted by  $K(m) \in \mathcal{K}$ .

From the new underlying network architecture of data centers, it is obvious that the transmission of a data chunk between two servers in different racks on average takes more time other than two servers in the same rack as the data chunk should pass through three switches rather than one switch.

A task generated by a user requires its own data chunk to be processed. The data chunk is saved on  $d$  servers for security and availability reasons. In real applications a data chunk is stored on three servers in order that if one or two of the servers fail to work or become disconnected from the network, the data chunk needed for the task will still be available in the other servers. Because of limited storage, the data chunks are usually not replicated on more than three servers. We define the type of a task by the location where its data chunk is stored, denoted by  $\bar{L} = (m_1, m_2, m_3)$  where  $m_1, m_2$ , and  $m_3$  are the servers storing the corresponding data chunk. Then, the set of all task types is  $\mathcal{L} = \{(m_1, m_2, m_3) \in \mathcal{M}^3 | m_1 < m_2 < m_3\}$ .

When a server is allocated to a task to process it, the server could have the data chunk available on its memory, local disk, or the server could not have the data. In case that the data is not available on the server, the server requires the data from another server which can be in the same rack, in a different rack, or even in another data center. Therefore, different levels of data locality exists in data centers [9]. Most of the theoretical works that have been done on load balancing (scheduling) for data centers have been done on two levels of data locality which will be illustrated in more detail later. The focus of this thesis is scheduling for three levels of data locality, but algorithms for two levels of data locality will also be discussed in the previous work found in Sections 2.3 and 2.4.

As a convention, for a task, the three servers which have the data chunk associated to it ( $\{m | m \in \bar{L}\}$ ) are called the *local servers*, and if the task

receives service from one of these local servers, we say that the task is *local* to the server and receives service *locally*. A task receives service *rack-locally* from one of the servers that does not have the data chunk stored on it, but is in the same rack as the required data is stored in one of its servers. The set of *rack-local* servers for a task of type  $\bar{L}$  is  $\bar{L}_k = \{m \notin \bar{L} | \exists n \in \bar{L} \text{ s.t. } K(m) = K(n)\}$ . Finally, a task receives service *remotely* if the server giving service not only does not have the required data, but the data is also not stored in another server in the same rack of the server. In summary,  $m \in \bar{L}$ ,  $m \in \bar{L}_k$ , and  $m \in \bar{L}_r$  denote that the server  $m$  is a local, rack-local, and remote server to the task of type  $\bar{L}$ , respectively.

We analyze the system in a discrete-time regime, where the time slots are numbered by  $t, t \geq 0$ , with the following service and task arrival processes:

**Service process:** Assume that a local service, rack-local service, and remote service follows geometric distribution with mean  $\frac{1}{\alpha}$ ,  $\frac{1}{\beta}$ , and  $\frac{1}{\gamma}$ , respectively. It is clear from the structure of the data centers that because of the fetching time, it takes on average shortest time for a task to receive service from a local server, other than a rack-local server, and longest time on a remote server. Hence,  $\alpha > \beta > \gamma$ . A server can process at most one task at a time slot, and the task processing is assumed to be non-preemptive. A task departs the system at the end of the time slot that the service is completed. Note that the completion time of a task is not only the service time, but also the waiting time of the task to be assigned to a server for service.

**Arrival process:** Task arrival occurs in the beginning of a time slot. The number of incoming tasks of type  $\bar{L}$  at the beginning of time slot  $t$  is denoted by  $A_{\bar{L}}(t)$ . The arrival process of different task types are independent of each other, with  $\mathbb{E}[A_{\bar{L}}(t)] = \lambda_{\bar{L}}$ . The arrival rate vector of all task types is denoted by  $\boldsymbol{\lambda} = (\lambda_{\bar{L}} : \bar{L} \in \mathcal{L})$ . We further assume a bounded total number of task arrival in each time slot.

The load balancing policy (consisting of routing and scheduling defined later in this section) for a data center decides which task should be assigned to an idle server for service. The main two optimality criteria for a load balancing policy (scheduler) are *throughput optimality* and *heavy-traffic optimality* defined as follows:

- **Throughput Optimality:** A load balancing algorithm is called to be throughput optimal if it stabilizes the data center for any arrival rate

vector strictly within the capacity region.

- **Heavy-Traffic Optimality:** A load balancing algorithm is said to be heavy-traffic optimal if it asymptotically minimizes the mean task completion time as the arrival rate approaches the boundary of the capacity region.

The load on data centers changes frequently, so as long as the arrival rate is within the capacity region of the data center, a throughput optimal scheduler is robust to the changes. In pick loads where the arrival rate is close to the boundary of the capacity region, a heavy-traffic optimal scheduler assigns tasks to servers efficiently, hence tasks experience the minimum mean completion time. Most of the heuristic load balancing algorithms for data centers have not been studied in theory [8, 10–14]. In this thesis, we will discuss the literature for scheduling algorithms with theoretical guarantee for their optimality for two levels of data locality and the affinity scheduling case. We claim that the extension of algorithms for two levels of data locality are either not optimal for three levels of data locality or not practical to be implemented. We will then propose a novel throughput and heavy-traffic optimal algorithm for three levels of data locality.

Note that the arriving tasks to the system do not get service immediately in case where all the servers are busy processing other tasks. Therefore, in such cases that all servers are busy, an incoming task is routed to a queue waiting for receiving service. Based on the scheduler used in the system, different queueing structures are needed. For example, only one queue is needed in order to implement *First-Come-First-Served (FCFS)* scheduler for a data center. For other algorithms fewer, the same or a greater number of queues as the number of servers may be needed. The queue structure for different algorithms will be mentioned when the algorithms are illustrated in Chapters 2 and 3.

There are two parts for any load balancing policy (scheduler), *routing* and *scheduling* policies which are described as follows:

- **Routing:** When a new task arrives at the system, the routing policy determines which queue it should be routed to in order to wait until it receives service from a server.

- **Scheduling:** When a server becomes idle, and is ready to process a task, the scheduling policy determines which task receives service from the idle server.

The capacity region realization of a data center with three levels of data locality is calculated as follows. Recall the arrival rate vector  $\boldsymbol{\lambda} = (\lambda_{\bar{L}} : \bar{L} \in \mathcal{L})$ . A decomposition of the arrival rate of a task type,  $\lambda_{\bar{L}}$  is  $(\lambda_{\bar{L},m}, m \in \mathcal{M})$ , where  $\lambda_{\bar{L},m}$  is the arrival rate of task type  $\bar{L}$  that is processed in server  $m$ . Assuming that a server can afford total local, rack-local and remote load of 1, a necessary condition that an arrival rate vector  $\boldsymbol{\lambda}$  is supportable is the following:

$$\sum_{\bar{L}:m \in \bar{L}} \frac{\lambda_{\bar{L},m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \frac{\lambda_{\bar{L},m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \frac{\lambda_{\bar{L},m}}{\gamma} < 1. \quad (1.1)$$

Then, an outer bound of the capacity region can be characterized as the set of arrival rates  $\boldsymbol{\lambda}$  such that there exists a decomposition  $(\lambda_{\bar{L},m}, m \in \mathcal{M})$  satisfying the necessary condition (1.1). The outer bound of the capacity region denoted by  $\Lambda$ , which will be shown in Section 3.3 that is the same as the capacity region itself, is formalized as below:

$$\begin{aligned} \Lambda = \{ & \boldsymbol{\lambda} = (\lambda_{\bar{L}} : \bar{L} \in \mathcal{L}) \mid \\ & \exists \lambda_{\bar{L},m} \geq 0, \forall \bar{L} \in \mathcal{L}, \forall m \in \mathcal{M}, s.t. \\ & \lambda_{\bar{L}} = \sum_{m=1}^M \lambda_{\bar{L},m}, \forall \bar{L} \in \mathcal{L}, \\ & \sum_{\bar{L}:m \in \bar{L}} \frac{\lambda_{\bar{L},m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \frac{\lambda_{\bar{L},m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \frac{\lambda_{\bar{L},m}}{\gamma} < 1, \forall m \}. \end{aligned} \quad (1.2)$$

Therefore, a linear programming optimization problem should be solved in order to find  $\Lambda$ .

# CHAPTER 2

## LITERATURE REVIEW

The near-data scheduling problem for the system illustrated in Chapter 1 is a special case of *affinity* scheduling [15–19]. In an affinity scheduling problem, instead of having a number of locality levels, a task of type  $\bar{L}$  can be processed by server  $m$  with rate  $\mu_{\bar{L},m}$  (in our system model of Chapter 1,  $\mu_{\bar{L},m}$  can only be  $\alpha, \beta$ , or  $\gamma$  according to whether server  $m$  is local, rack-local, or remote to the task of type  $\bar{L}$ , respectively, but in affinity scheduling problem  $\mu_{\bar{L},m}$  can take any non-negative value). In the following, we briefly describe the *Fluid Model Planning* and *Generalized  $c\mu$ -rule* as the affinity scheduling algorithms and discuss their shortcomings. Then we explain two algorithms for a system with two levels of data locality.

### 2.1 Fluid Model Planning

Harrison and Lopez [17, 18] proposed the fluid model planning algorithm for affinity scheduling problem. The queueing structure needed to implement this algorithm is to have separate queues for different types of tasks (each queue is associated to a task type). Then the routing and scheduling policies are as follows:

- **Routing:** An incoming task is routed to the queue associated to its type.
- **Scheduling:** The arrival rate of each task type is needed to be known to solve a linear programming optimization and find the basic activities based on which the servers are assigned to process tasks.

Fluid model planning algorithm is both throughput and heavy-traffic optimal. However, there are two main objections to this algorithm. First, distinct queues are considered for different task types. In the system model described

in Chapter 1 each task type has its data chunk stored on three servers out of a total of  $M$  servers. Therefore, there can be  $\binom{M}{3}$  different number of task types. This means that we should have in the order of  $O(M^3)$  distinct queues for tasks, where each queue is associated to a task type. However, in a data center there usually exists thousands of servers, so it is not logical to define  $O(M^3)$  number of queues as it makes the underlying network, scheduling computation, and data base infrastructure complicated. Second, the arrival rate of each task type is considered to be known to the system for the scheduling part, but in a real data center the load changes frequently and is not known as users can have random behavior. As a result, fluid model planning cannot be used for practical issues, unless it has both optimality conditions.

## 2.2 Generalized $c\mu$ -Rule

Stolyar and Mandelbaum [16,20] proposed the generalized  $c\mu$ -rule algorithm. Similar to the fluid model planning queueing structure, this algorithm also requires the existence of one queue per task type. In contrast to the fluid model planning algorithm, the arrival rate of task types are not required to be known to implement this algorithm. Instead, the generalized  $c\mu$ -rule utilizes the MaxWeight procedure for the scheduling part which makes the algorithm needless of the prior knowledge of the task types' arrival rates. Assume the cost rate incurred by the type  $\bar{L}$  tasks is  $C_{\bar{L}}(Q_{\bar{L}})$  where  $Q_{\bar{L}}$  denotes the number of tasks of type  $\bar{L}$  queued in the corresponding queue. The cost function should have fairly normal conditions for which we can mention the following (for more detail refer to [16, 20]):  $C_{\bar{L}}(\cdot)$  should be convex and continuous with  $C_{\bar{L}}(0) = 0$ . The derivative of the cost function,  $C'_{\bar{L}}(\cdot)$ , should be strictly increasing and continuous with  $C'_{\bar{L}}(0) = 0$ . The routing and scheduling policies of the generalized  $c\mu$ -rule are as below:

- **Routing:** An incoming task is routed to the queue associated to its type.
- **Scheduling:** An idle server  $m \in \mathcal{M}$  is scheduled to a task of type  $\bar{L}$  in the set below at time slot  $t$ :

$$ArgMax_{\bar{L}} \left\{ C'_{\bar{L}}(Q_{\bar{L}}(t)) \mu_{\bar{L},m} \right\}. \quad (2.1)$$

In the system model with three levels of data locality,  $\mu_{\bar{L},m}$  is  $\alpha, \beta$ , or  $\gamma$  if the task type  $\bar{L}$  is local, rack-local, or remote to the idle server  $m$ , respectively. Mandelbaum and Stolyar proved that the generalized  $c\mu$ -rule asymptotically minimizes both instantaneous and cumulative queueing costs in heavy traffic [16]. Consider the same cost function for all task types as  $C_{\bar{L}}(Q_{\bar{L}}) = Q_{\bar{L}}^{\beta+1}$ ,  $\forall \bar{L} \in \mathcal{L}$ , where  $\beta > 0$ . The function  $Q^{\beta+1}$  satisfies all the required conditions mentioned for a valid cost function. Therefore, the generalized  $c\mu$ -rule asymptotically minimizes the holding cost  $\sum_{\bar{L}} Q_{\bar{L}}^{\beta+1}$ , and as the constant  $\beta$  should strictly be greater than zero, this algorithm cannot minimize  $\sum_{\bar{L}} Q_{\bar{L}}$ . Hence, the generalized  $c\mu$ -rule is not heavy-traffic optimal. Besides, we still need many queues in the order of cubic number of servers in order to implement this algorithm which makes the underlying system complicated and is not practical.

In Sections 2.3 and 2.4, two algorithms for two levels of data locality will be discussed. In both algorithms, no prior task types' arrival rate information is needed to be known. Furthermore, only  $M$  queues, one for each server, is needed for queueing the tasks that are waiting for service. In a system with two levels of data locality there is no notion of rack structure and rack-local service, instead there is only a core switch connecting all servers to each other. A task can only get service locally with rate  $\alpha$  from one of the local servers ( $m \in \bar{L}$ ), or it gets service remotely with rate  $\gamma$  from any other servers ( $m \notin \bar{L}$ ). The capacity region of a system with two levels of data locality is given in equation (2.2) which can be driven with the same reasoning we had for a system with three levels of data locality.

$$\begin{aligned} \Lambda = \{ \boldsymbol{\lambda} = (\lambda_{\bar{L}} : \bar{L} \in \mathcal{L}) \mid \\ \exists \lambda_{\bar{L},m} \geq 0, \forall \bar{L} \in \mathcal{L}, \forall m \in \mathcal{M}, s.t. \\ \lambda_{\bar{L}} = \sum_{m=1}^M \lambda_{\bar{L},m}, \forall \bar{L} \in \mathcal{L}, \\ \sum_{\bar{L}: m \in \bar{L}} \frac{\lambda_{\bar{L},m}}{\alpha} + \sum_{\bar{L}: m \notin \bar{L}} \frac{\lambda_{\bar{L},m}}{\gamma} < 1, \forall m \}. \end{aligned} \quad (2.2)$$



First, join the shortest queue-MaxWeight (*JSQ-MW*) [21] which is a throughput optimal, but not heavy-traffic optimal algorithm will be discussed. Then the Priority Algorithm for Near-Data Scheduling (*Pandas*) algorithm proposed by Xie and Lu [22] which is both throughput and heavy-traffic optimal will be presented.

## 2.3 Join the Shortest Queue-MaxWeight (JSQ-MW)

Joining the shortest queue-MaxWeight algorithm proposed by Wang et al. [21] requires the existence of one queue per server. The length of the  $m$ -th queue at time slot  $t$  is denoted by  $Q_m(t)$ . The central scheduler that maintains all the queue lengths routes the new incoming tasks to a queue and schedules the idle servers to a task as follows:

- **Routing:** An arriving task of type  $\bar{L}$  is routed to the shortest queue of the local servers in the set  $\bar{L}$  (all ties are broken randomly throughout this paper).
- **Scheduling:** At time slot  $t$ , the idle server  $m$  is assigned to process a task from a queue in the set given in equation (2.3):

$$\arg \max_{n \in \mathcal{M}} \{ \alpha Q_n(t) I_{\{n=m\}}, \beta Q_n(t) I_{\{n \neq m\}} \}. \quad (2.3)$$

The JSQ-MW algorithm is proven to be throughput optimal for a system with two levels of data locality [21]. However, it is not heavy-traffic optimal.

As a definition, if the incoming load routed to a server exceeds the capacity of the server, the server is called to be a *beneficiary* server. On the other hand, if the incoming load routed to a server is less than what the server can process, the server is called to be a *helper*. Beneficiaries cannot process all the tasks routed to their queues, so they get help from helpers (helping servers).

Wang et al. [21] proved that the JSQ-MW algorithm can minimize the mean task completion time in a specific traffic scenario as follows. If all the incoming traffic is local to a set of servers where all of them are beneficiaries,

and the rest of servers do not receive any traffic load, so to be helpers, the JSQ-MW algorithm minimizes the mean task completion time.

## 2.4 Priority Algorithm for Near-Data Scheduling (Pandas)

To the best of our knowledge, the Pandas algorithm is the only throughput and heavy-traffic optimal algorithm for a system with two levels of data locality. Assuming the existence of one queue per server, the routing and scheduling policies are as follows:

- **Routing:** An arriving task of type  $\bar{L}$  is routed to the shortest queue of the local servers in the set  $\bar{L}$ .
- **Scheduling:** As long as there exists a task available in the queue of the idle server  $m$ , it will be assigned to a local task at the  $m$ -th queue. If there is no local task at the  $m$ -th queue, the idle server is assigned to give service to a task in the longest queue of the system (the task in the longest queue can be local or remote to the idle server).

We can improve the Pandas algorithm by adding the following two features to the algorithm. First, when an idle server  $m$  does not have any tasks queued in front of it and is assigned to process a task from the longest queue, the scheduler can assign a local task to server  $m$  queued at the longest queue if available. Second, an idle server  $m$  is assigned to serve a remote task from the longest queue in the system if  $\max_{n \in \mathcal{M}} Q_n > \frac{\alpha}{\gamma}$  in order to make sure that the remote task will experience less service time if it gets service from the idle server, other than waiting at its current queue and receiving service locally.

In Chapter 3, our theoretical analysis of a system with three levels of data locality will be discussed.

# CHAPTER 3

## THREE LEVELS OF DATA LOCALITY

In this chapter we propose our two algorithms: the JSQ-MW algorithm and the Balanced-Pandas (Weighted-Workload Routing and Priority Scheduling) algorithm [1]. Our proposed JSQ-MW algorithm is an extension of the algorithm used for two levels of data locality which is throughput optimal, but not heavy-traffic optimal. We will show that the extension of the JSQ-MW algorithm also minimizes the mean task completion time in specific workloads, but not all loads. On the other hand, the Balanced-Pandas algorithm is a novel algorithm proposed by us which is throughput optimal for three levels of data locality. It is heavy-traffic optimal in the case that  $\beta^2 > \alpha\gamma$ , which means that the rack-local service is much faster than the remote service. This condition usually holds in real systems. To the best of our knowledge the Balanced-Pandas algorithm is the only throughput and heavy-traffic optimal algorithm proposed for a system with three levels of data locality.

The difficulty of designing an algorithm for a system with three levels of data locality is illustrated in Section 3.1.

### 3.1 The Performance versus Throughput Dilemma

Under the Pandas algorithm, each server has a queue maintaining tasks local to it. The Pandas routing policy balances tasks across their local servers. An idle server processes local tasks as long as there exists one in its queue; otherwise, it processes a remote task from the longest queue of the system. In other words, the Pandas algorithm forces servers to process as many local tasks as possible, then process remote tasks if they do not have any local tasks available. Therefore, in a system with three levels of data locality, the Pandas algorithm has good performance in low and medium loads by maximizing the number of tasks served locally. However, the Pandas algorithm sacrifices

throughput optimality at high loads. The example depicted in Figure 3.1 and the explanation afterward makes the performance versus throughput dilemma clear.

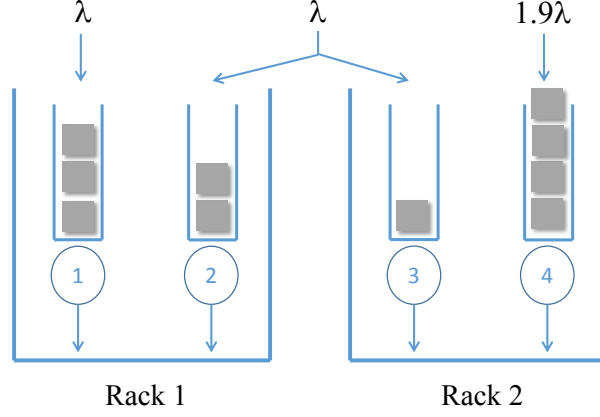


Figure 3.1: A system with two racks showing how performance should be sacrificed in order to achieve throughput optimality.

Assume each of the two racks has two servers as depicted in Figure 3.1. Three types of tasks receive service from the system as follows: one type of task with arrival rate  $\lambda$  is only local to the first server, another type with arrival rate  $\lambda$  is local to both the second and third servers, and the other type of task with arrival rate  $1.9\lambda$  is only local to the fourth server. Using the Pandas algorithm to maximize the number of tasks served locally, the second type of tasks is split between the second and third servers evenly. Assuming that local, rack-local, and remote service rates are  $\alpha = 1$ ,  $\beta = 0.9$ , and  $\gamma = 0.5$ , respectively, the Pandas algorithm can stabilize the system in the following region:

$$1.9\lambda < \alpha + \beta\left(1 - \frac{0.5\lambda}{\alpha}\right) + \gamma\left(1 - \frac{0.5\lambda}{\alpha}\right) + \gamma\left(1 - \frac{\gamma}{\alpha}\right),$$

which gives  $\lambda < 0.9355$ . However, if the second type of tasks is only routed to the second server, so the third server does not process any local tasks, but only processes rack-local tasks, then the system is stable as long as  $\lambda < 1$  which is a bigger capacity region other than the one for the Pandas algorithm. Our proposed two algorithms described in Sections 3.3 and 3.4 stabilize the system in the capacity region given in equation (1.2) without the knowledge of the tasks' arrival rates. In the next section, we draw an equivalent capacity region for a system with rack structure which will be used in the optimality

proofs of our proposed algorithms.

## 3.2 Equivalent Capacity Region

In this section we will show that the outer bound of capacity region proposed in equation (1.2) is actually the capacity region of a system with three levels of data locality. In the following lemma, we propose an equivalent capacity region with the one in equation (1.2) which will be used in our proofs.

**Lemma 1** *The following set  $\bar{\Lambda}$  is equivalent to  $\Lambda$  defined in equation (1.2):*

$$\begin{aligned} \bar{\Lambda} = \{ \boldsymbol{\lambda} = (\lambda_{\bar{L}} : \bar{L} \in \mathcal{L}) \mid \\ \exists \lambda_{\bar{L},n,m} \geq 0, \forall \bar{L} \in \mathcal{L}, \forall n \in \bar{L}, \forall m \in \mathcal{M}, s.t. \\ \lambda_{\bar{L}} = \sum_{n:n \in \bar{L}} \sum_{m=1}^M \lambda_{\bar{L},n,m}, \quad \forall \bar{L} \in \mathcal{L}, \\ \sum_{\bar{L}:m \in \bar{L}} \sum_{n:n \in \bar{L}} \frac{\lambda_{\bar{L},n,m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \sum_{n:n \in \bar{L}} \frac{\lambda_{\bar{L},n,m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \sum_{n:n \in \bar{L}} \frac{\lambda_{\bar{L},n,m}}{\gamma} < 1, \forall m \}, \end{aligned} \quad (3.1)$$

where  $\lambda_{\bar{L},n,m}$  is the arrival rate of type  $\bar{L}$  tasks that are local to server  $n$ , but are scheduled to be processed at server  $m$ .  $\lambda_{\bar{L},n,m}$  is actually the decomposition of  $\lambda_{\bar{L},m}$  and  $\lambda_{\bar{L},m} = \sum_{n \in \mathcal{M}} \lambda_{\bar{L},n,m}$ .

**proof:** In order to prove that  $\bar{\Lambda} = \Lambda$ , we show that  $\bar{\Lambda} \subset \Lambda$  and  $\Lambda \subset \bar{\Lambda}$ .

$\bar{\Lambda} \subset \Lambda$ : If  $\boldsymbol{\lambda} \in \bar{\Lambda}$ , there exists a load decomposition  $\{\lambda_{\bar{L},n,m}\}$  such that it satisfies all the conditions in (3.1). By defining  $\lambda_{\bar{L},m} \equiv \sum_{n:n \in \bar{L}} \lambda_{\bar{L},n,m}$ , it is clear that this decomposition of  $\boldsymbol{\lambda}$ , that is  $\{\lambda_{\bar{L},m}\}$  satisfies the conditions in equation (1.2), so  $\boldsymbol{\lambda} \in \Lambda$ . Hence  $\bar{\Lambda} \subset \Lambda$ .

$\Lambda \subset \bar{\Lambda}$ : If  $\boldsymbol{\lambda} \in \Lambda$ , there exists a load decomposition  $\{\lambda_{\bar{L},m}\}$  such that it satisfies all the conditions in equation (1.2). By defining  $\lambda_{\bar{L},n,m} \equiv \frac{\lambda_{\bar{L},m}}{|\bar{L}|}$ , it is clear that this decomposition of  $\boldsymbol{\lambda}$ , that is  $\{\lambda_{\bar{L},n,m}\}$  satisfies the conditions in (3.1), so  $\boldsymbol{\lambda} \in \bar{\Lambda}$ . Hence  $\Lambda \subset \bar{\Lambda}$ .

### 3.3 Join the Shortest Queue-MaxWeight (JSQ-MW)

In order to implement the JSQ-MW algorithm, the central scheduler keeps one queue per server, where the length of  $m$ -th queue associated to the  $m$ -th server at time slot  $t$  is denoted by  $Q_m(t)$ . The  $m$ -th queue only keeps tasks that are local to the  $m$ -th server. Then, the JSQ-MW routing and scheduling policies are as follows:

- **JSQ-MW Routing:** An arriving task of type  $\bar{L}$  is routed to its shortest local queue. That is, the central scheduler inserts the new task to the shortest queue in the set  $\{Q_m | m \in \bar{L}\}$ , where ties are broken randomly.
- **JSQ-MW Scheduling:** The scheduling decision  $\eta_m(t)$  of an idle server  $m$  at time slot  $t$  is chosen from the following set where the ties are broken randomly. That is, idle server  $m$  is scheduled to give service to a task queued in a queue in the following set:

$$\arg \max_{n \in \mathcal{M}} \{ \alpha Q_n(t) I_{\{n=m\}}, \beta Q_n(t) I_{\{K(n)=K(m)\}}, \gamma Q_n(t) I_{\{K(n) \neq K(m)\}} \}.$$

In order to describe the queue evolution of a system with three levels of data locality using the JSQ-MW algorithm, we define the following terminologies. Let the number of type  $\bar{L}$  tasks that are routed to the  $m$ -th queue at time slot  $t$  be denoted by  $A_{\bar{L},m}(t)$ . Then the total number of task arrivals to  $Q_m$  at time slot  $t$  is as follows:

$$A_m(t) = \sum_{\bar{L}: m \in \bar{L}} A_{\bar{L},m}(t).$$

Server  $m$  provides local, rack-local, and remote services denoted by  $S_m^l(t)$ ,  $R_m^k(t)$ , and  $R_m^r(t)$ , respectively. Service times are assumed to follow geometric distribution, so  $S_m^l(t)$ ,  $R_m^k(t)$ , and  $R_m^r(t)$  are Bernoulli random variables in each time slot as follows:  $S_m^l(t) \sim \text{Bern}(\alpha I_{\{\eta_m(t)=m\}})$ ,  $R_m^k(t) \sim \text{Bern}(\beta I_{\{K(\eta_m(t))=K(m), \eta_m(t) \neq m\}})$ , and  $R_m^r(t) \sim \text{Bern}(\gamma I_{\{K(\eta_m(t)) \neq K(m)\}})$ . Queue  $m$  can receive local, rack-local, and remote services as follows. The local service is received from server  $m$  which is  $S_m^l(t)$ , while rack-local service can be received from any other servers in the same rack of  $Q_m$  which we denote by  $S_m^k(t) = \sum_{n: K(n)=K(m), n \neq m} R_n^k(t) I_{\{\eta_n(t)=m\}}$ , and  $Q_m$  receives re-

remote services from all servers out of its rack which is denoted by  $S_m^r(t) = \sum_{n:K(n) \neq K(m)} R_n^r(t) I_{\{\eta_n(t)=m\}}$ . The total number of task departures for  $Q_m$  at time slot  $t$  is equal to the summation of local, rack-local, and remote services given to  $Q_m$  at time slot  $t$  which we denote by  $S_m(t) \equiv S_m^l(t) + S_m^k(t) + S_m^r(t)$ . Defining  $U_m(t) = \max\{0, S_m(t) - A_m(t) - Q_m(t)\}$  as the unused service allocated to the  $m$ -th queue, the queues evolve from one time slot to the next one as follows:

$$Q_m(t+1) = Q_m(t) + A_m(t) - S_m(t) + U_m(t).$$

Note that the queue length vector  $\mathbf{Q}(t) = (Q_1(t), Q_2(t), \dots, Q_M(t))$  is not a Markov chain since given the queue lengths at a time slot, the future of the queue lengths is not independent from the past. The reason is that given the queue lengths, we cannot figure out the status of each server in the system. Therefore, we define the working status of server  $m$  at time slot  $t$ ,  $f_m(t)$  as follows:

$$f_m(t) = \begin{cases} -1 & \text{if server } m \text{ is idle} \\ n & \text{if server } m \text{ processes a task from } Q_n \end{cases}.$$

If the  $m$ -th server is in idle mode, not processing any tasks, its working status is equal to  $-1$ . Otherwise, if server  $m$  processes a local task from  $Q_m$ , then  $f_m(t) = m$ . If server  $m$  processes a task from  $Q_n$ , where  $n \neq m$ , but  $K(n) = K(m)$  ( $K(n) \neq K(m)$ ), it means that it is serving a rack-local (remote) task and  $f_m(t) = n$ .

Defining the working status vector  $\mathbf{f}(t) = (f_1(t), f_2(t), \dots, f_M(t))$ , as the service times follow geometric distributions, both queue length vector  $\mathbf{Q}(t)$  and  $\mathbf{f}(t)$  together,  $\{\mathbf{Z}(t) = (\mathbf{Q}(t), \mathbf{f}(t)), t \geq 0\}$ , form an irreducible and aperiodic Markov chain. The following theorem indicates the capacity region and throughput optimality of the JSQ-MW algorithm. What we mean by the system being stabilized is that the queue lengths are bounded in steady state.

**Theorem 1** *The JSQ-MW algorithm can stabilize the system with three levels of data locality, as long as the arrival rate vector of the task types is strictly within the outer bound of the capacity region,  $\Lambda$ . This means that  $\Lambda$  is the capacity region of the system and the JSQ-MW algorithm is throughput*

*optimal.*

**proof:** An extension of the Foster-Lyapunov theorem, where the  $T$ -time slot drift of the Lyapunov function is studied, is used to prove the throughput optimality of the JSQ-MW algorithm. We choose the function  $V_1(t) = \|\mathbf{Q}(t)\|^2 = \sum_{m=1}^M Q_m^2(t)$  as the Lyapunov function. Note that this choice of the Lyapunov function satisfies the requirements of non-negativity, being equal to zero only at  $\mathbf{Q}(t) = \bar{\mathbf{0}}$ , and going to infinity as any elements of  $\mathbf{Q}(t)$  goes to infinity. In Appendix A.1, we show that as long as the arrival rate vector of task types is strictly within  $\Lambda$ , under the JSQ-MW algorithm, there exists an integer  $T > 0$  where the expected  $T$ -time slot drift of  $V_1(t)$  is negative outside of a bounded region of the state space, and is finite inside this bounded region. Therefore, the fact that  $\Lambda$  is the capacity region of the system and the throughput optimality of the JSQ-MW algorithm are followed by the extension of the Foster-Lyapunov theorem.

A corollary of Theorem 1 is that the outer bound of capacity region proposed in equation (1.2) is actually the capacity region of a system with three levels of data locality.

It was mentioned in Section 2.3 that the JSQ-MW algorithm is not heavy-traffic optimal for two levels of data locality, but minimizes the mean task completion time at high loads under a specific traffic scenario. We should also mention that it is very rare that such a traffic load occurs in real-world applications. In the simulation results in Chapter 4, it would be clear from the simulation results that the JSQ-MW algorithm is not heavy-traffic optimal for a system with three levels of data locality. In the following, we will show the traffic scenario in which the JSQ-MW algorithm can minimize the mean task completion time at high loads for a system with three levels of data locality.

Under the following three conditions, the JSQ-MW algorithm minimizes the mean task completion time at high loads:

1. All the incoming traffic concentrates on a subset of racks, so the complement set of racks does not have any incoming local tasks.
2. The racks that receive nonzero incoming local tasks cannot process their incoming local tasks without getting help from other servers in the racks with no incoming local task.



3. The servers in the racks that receive local incoming tasks either have zero incoming local tasks or are overloaded.

Here we formalize the traffic scenario in which the JSQ-MW algorithm is heavy-traffic optimal after setting some notations. The set of racks that receive nonzero local tasks is shown by  $\mathcal{O}$ , and the racks belonging to this set are called overloaded racks. The set of all servers that receive nonzero local tasks is denoted by  $\mathcal{M}_l$ . It is clear that all the servers in the set  $\mathcal{M}_l$  belong to the racks in the set  $\mathcal{O}$ . The other set of servers in the overloaded racks that receive zero local tasks is shown by  $\mathcal{M}_k = \{m \in \mathcal{M} | K(m) \in \mathcal{O}, \text{ and } m \notin \mathcal{M}_l\}$ . The set of all other servers not belonging to the overloaded racks that do not receive any local tasks is denoted by  $\mathcal{M}_r = \{m \in \mathcal{M} | K(m) \notin \mathcal{O}\}$ . Assuming a set of servers  $\mathcal{S} \subset \mathcal{M}_l$ , the set of all task types having a local server in the set  $\mathcal{S}$  is denoted by  $\mathcal{N}(\mathcal{S}) = \{\bar{L} \in \mathcal{L} | \exists m \in \mathcal{S}, \text{ s.t. } m \in \bar{L}\}$ . Likewise, for a set of racks  $\mathcal{R} \subset \mathcal{O}$ , the set of task types that have a local server in the set  $\mathcal{R}$  is denoted by  $\mathcal{N}(\mathcal{R}) = \{\bar{L} \in \mathcal{L} | \exists m, \text{ s.t. } K(m) \in \mathcal{R}, \text{ and } m \in \bar{L}\}$ . Furthermore, the set of servers belonging to a rack in the set  $\mathcal{R}$  that receive local incoming tasks (receive zero local task) is denoted by  $\mathcal{M}_l^{\mathcal{R}} = \{m \in \mathcal{M}_l | K(m) \in \mathcal{R}\}$  ( $\mathcal{M}_k^{\mathcal{R}} = \{m \in \mathcal{M}_k | K(m) \in \mathcal{R}\}$ ).

The heavy-traffic regime is characterized as follows. Any subset of servers receiving local tasks are overloaded, and any subset of racks receiving local tasks are also overloaded. However, the arrival rate vector of task types should be in the capacity region in order to be supportable, that is,  $\sum_{\bar{L} \in \mathcal{L}} \lambda_{\bar{L}} < |\mathcal{M}_l|\alpha + |\mathcal{M}_k|\beta + |\mathcal{M}_r|\gamma$ . The following three conditions characterize a heavy-traffic regime with a parameter  $\epsilon > 0$  which is the  $L^1$ -norm of the difference of the arrival rate vector and the nearest point on the boundary of the capacity region to the arrival rate vector.

$$\begin{aligned}
& \forall \mathcal{S} \subset \mathcal{M}_l, \quad \sum_{\bar{L} \in \mathcal{N}(\mathcal{S})} \lambda_{\bar{L}} > |\mathcal{S}|\alpha, \\
& \forall \mathcal{R} \subset \mathcal{O}, \quad \sum_{\bar{L} \in \mathcal{N}(\mathcal{R})} \lambda_{\bar{L}} > |\mathcal{M}_l^{\mathcal{R}}|\alpha + |\mathcal{M}_k^{\mathcal{R}}|\beta, \\
& \sum_{\bar{L} \in \mathcal{L}} \lambda_{\bar{L}} = |\mathcal{M}_l|\alpha + |\mathcal{M}_k|\beta + |\mathcal{M}_r|\gamma - \epsilon.
\end{aligned} \tag{3.2}$$

Theorem 2 clarifies the heavy-traffic optimality of the JSQ-MW algorithm.

**Theorem 2** *The JSQ-MW algorithm minimizes the mean task completion time in the steady state as long as the task arrival process  $\{A_{\bar{L}}^{\epsilon}(t), t \geq 0\}_{\bar{L} \in \mathcal{L}}$*

with arrival rate vector  $\lambda^\epsilon$  satisfies the conditions in (3.2), where servers are either overloaded or do not receive any local tasks, and the racks are also either overloaded or do not receive any local tasks.

**proof:** The complete proof of Theorem 2 appears in Appendix A.2.

As mentioned, the JSQ-MW algorithm is not heavy-traffic optimal in all loads. The reason is that if any task is local to helper servers in overloaded or under-loaded racks, their local queues grow and local tasks receive service by unnecessary delay. In the next section, we propose our novel algorithm called Balanced-Pandas, which is both throughput and heavy-traffic optimal.

### 3.4 Balanced-Pandas

In order to implement the Balanced-Pandas algorithm, the central scheduler should keep three queues per server since using this algorithm the incoming tasks are not necessarily routed to the queue of their local server, so we want to keep track of the local, rack-local and remote tasks routed to a server by considering different queues for each of them. The tasks routed to server  $m$  which are local (rack-local or remote) to it are queued in the first (second or third) queue of the server, which is denoted by  $Q_m^l$  ( $Q_m^k$  or  $Q_m^r$ ). The three queue lengths of the  $m$ -th server at time slot  $t$  are denoted by the vector notation  $\bar{Q}_m(t) = (Q_m^l(t), Q_m^k(t), Q_m^r(t))$ , and the central scheduler maintains the vector of queue lengths  $\bar{Q}(t) = (\bar{Q}_1(t), \bar{Q}_2(t), \dots, \bar{Q}_M(t))$ . As the service time of local, rack-local, and remote tasks follow geometric distributions with means  $\frac{1}{\alpha}$ ,  $\frac{1}{\beta}$ , and  $\frac{1}{\gamma}$ , respectively, the mean time needed for server  $m$  to process all the tasks queued at  $Q_m^l$ ,  $Q_m^k$ , and  $Q_m^r$  at time slot  $t$  is as follows:

$$W_m(t) = \frac{Q_m^l(t)}{\alpha} + \frac{Q_m^k(t)}{\beta} + \frac{Q_m^r(t)}{\gamma}.$$

We call  $W_m(t)$  the *workload* on the  $m$ -th server.

Figure 3.2 along with the Balanced-Pandas routing and scheduling policies presented in the following make the queueing structure and the Balanced-Pandas algorithm clear.

- **Balanced-Pandas Routing (Weighted-Workload Routing):** The routing decision for an incoming task of type  $\bar{L}$  is based on both data

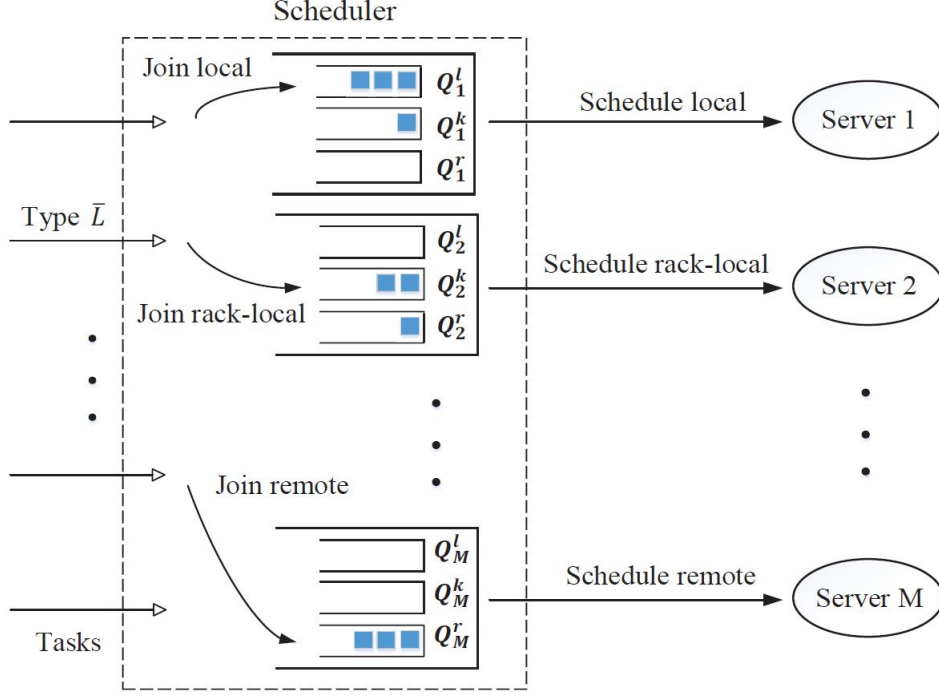


Figure 3.2: The queueing structure needed for the Balanced-Pandas algorithm.

locality and the workload on the servers. In order to decide routing for an incoming task, the workloads of servers local (rack-local or remote) to the incoming task are each divided by  $\alpha$  ( $\beta$  or  $\gamma$ ). The task is routed to the corresponding sub-queue ( $Q_{m^*}^l$ ,  $Q_{m^*}^k$ , or  $Q_{m^*}^r$ ) of the server  $m^*$  with the minimum weighted workload. That is, if the incoming task is local, rack-local, or remote to the server with the minimum weighted workload, it is routed to the first, second, or third queue, respectively. Formally speaking, an arriving task of type  $\bar{L}$  is routed to the corresponding sub-queue of a server in the set below:

$$\arg \min_{m \in \mathcal{M}} \left\{ \frac{W_m(t)}{\alpha} I_{\{m \in \bar{L}\}}, \frac{W_m(t)}{\beta} I_{\{m \in \bar{L}_k\}}, \frac{W_m(t)}{\gamma} I_{\{m \in \bar{L}_r\}} \right\}.$$

- **Balanced-Pandas Scheduling (Prioritized Scheduling):** As server  $m$  becomes idle at time  $t^-$ , the central scheduler assigns it to a local task queued at  $Q_m^l$  at time slot  $t$ , if available. However, if  $Q_m^l(t) = 0$ , server  $m$  will be assigned to a rack-local task queued at  $Q_m^k$ , if available. If both local and rack-local sub-queues of server  $m$  are empty, the server is assigned to process a remote task queued at  $Q_m^r$ . In other

words, the highest priority to process a task for an idle server is given to local, then rack-local, and finally remote tasks queued in front of it. If all sub-queues of the idle server are empty, the idle server remains idle until a new task joins any of the sub-queues.

In the following, we will first propose two optimality theorems of the Balanced-Pandas algorithm, then we will define the notations which will be used in the proof of these theorems in Appendices A.3 and A.4.

**Theorem 3** *The Balanced-Pandas algorithm stabilizes a system with three levels of data locality as long as the arrival rate is strictly inside the capacity region, which means that the Balanced-Pandas algorithm is throughput optimal.*

**proof:** We use the Foster-Lyapunov theorem to prove the throughput optimality. We use the  $l^2$ -norm of the workload vector of servers as the Lyapunov function:

$$V_3(Z(t)) = \|\mathbf{W}(t)\|^2.$$

This choice of Lyapunov function is non-negative, is equal to zero just at  $\mathbf{W}(t) = \bar{0}$ , and goes to infinity as any elements of  $\mathbf{W}(t)$  goes to infinity. We show that there exists a finite integer  $T > 0$  where the expectation of the  $T$ -time slot drift of the Lyapunov function is negative outside of a bounded region of the state space, and is positive and finite inside this bounded region. We should note that in the proof of throughput optimality, we do not use the fact of using prioritized scheduling. Therefore, to the purpose of throughput optimality, an idle server can serve any task in its three sub-queues as local, rack-local and remote tasks decrease the expected workload at the same rate. The prioritized scheduling is to minimize the mean task completion time experienced by tasks which will be of interest in heavy-traffic optimality. For the complete proof refer to Appendix A.3

**Theorem 4** *As long as  $\beta^2 > \alpha\gamma$ , the Balanced-Pandas algorithm is heavy-traffic optimal, i.e., minimizes the mean task completion time as the arrival rate vector of task types approaches the boundary of the capacity region.*

**proof:** For the complete proof look at Appendix A.4.

In the next three subsections, the queue dynamics when the Balanced-Pandas algorithm is used, overloaded servers and racks, and ideal load decomposition will be discussed which will be used in the proofs of Theorems 3 and 4.

### 3.4.1 Queue Dynamics

Recall that  $A_{\bar{L},m}(t)$  denotes the number of type  $\bar{L}$  tasks that are routed to the  $m$ -th queue, and  $\bar{L}$ ,  $\bar{L}_k$ , and  $\bar{L}_r$  denote the set of local, rack-local and remote servers to a task of type  $\bar{L}$ . Using these definitions, we can formalize the local, rack-local, and remote tasks routed to the three sub-queues of the  $m$ -th server at time slot  $t$  denoted by  $A_m^l(t)$ ,  $A_m^k(t)$ , and  $A_m^r(t)$ , respectively as follows:  $A_m^l(t) = \sum_{\bar{L}:m \in \bar{L}} A_{\bar{L},m}(t)$ ,  $A_m^k(t) = \sum_{\bar{L}:m \in \bar{L}_k} A_{\bar{L},m}(t)$ , and  $A_m^r(t) = \sum_{\bar{L}:m \in \bar{L}_r} A_{\bar{L},m}(t)$ .

Remembering the Markov chain defined in Section 3.3, the queue dynamics themselves cannot form a Markov chain alone. Therefore, we define the working status of server  $m$  at time slot  $t$  as follows:

$$f_m(t) = \begin{cases} -1 & \text{if server } m \text{ is idle} \\ 0 & \text{if server } m \text{ processes a local task from } Q_m^l \\ 1 & \text{if server } m \text{ processes a rack-local task from } Q_m^k \\ 2 & \text{if server } m \text{ processes a remote task from } Q_m^r \end{cases}.$$

When server  $m$  is done processing a task at time slot  $t - 1$ , so its working status is  $f_m(t^-) = -1$ , the scheduling decision  $\eta_m(t)$  for this server is made based on both working status of servers  $\mathbf{f}(t) = (f_1(t), f_2(t), \dots, f_M(t))$  and the queue length vector  $\bar{\mathbf{Q}}(t)$ . Note that  $\eta_m(t) = f_m(t)$  as long as server  $m$  is busy processing a task, and when the server becomes idle at the end of a time slot,  $\eta_m(t)$  is determined by the scheduling policy.

As described in Chapter 1, in a system with three levels of data locality the service distributions are as follows. If server  $m$  is working on a local (rack-local, or remote) task at time slot  $t$ , the service provided by the server at time slot  $t$  follows a Bernoulli random variable with mean  $\alpha$  ( $\beta$ , or  $\gamma$ ) which is denoted by  $S_m^l(t)$  ( $S_m^k(t)$ , or  $S_m^r(t)$ ). In other words,

the local (rack-local, or remote) service provided by the  $m$ -th server at time slot  $t$  is  $S_m^l(t) \sim \text{Bern}(\alpha I_{\{\eta_m(t)=0\}})$  ( $S_m^k(t) \sim \text{Bern}(\beta I_{\{\eta_m(t)=1\}})$ , or  $S_m^r(t) \sim \text{Bern}(\gamma I_{\{\eta_m(t)=2\}})$ ).

Defining the unused service of server  $m$  as  $U_m(t) = \max\{0, S_m^r(t) - A_m^r(t) - Q_m^r(t)\}$ , the three sub-queues of server  $m$  evolve as follows:

$$\begin{aligned} Q_m^l(t+1) &= Q_m^l(t) + A_m^l(t) - S_m^l(t), \\ Q_m^k(t+1) &= Q_m^k(t) + A_m^k(t) - S_m^k(t), \\ Q_m^r(t+1) &= Q_m^r(t) + A_m^r(t) - S_m^r(t) + U_m(t). \end{aligned} \tag{3.3}$$

The service times are all geometrically distributed, so the queue length vector  $\bar{\mathbf{Q}}(t)$  together with the working status vector of servers  $\mathbf{f}(t)$  form the irreducible and aperiodic Markov chain  $(\{\mathbf{Z}(t) = (\bar{\mathbf{Q}}(t), \mathbf{f}(t)), t \geq 0\})$ .

### 3.4.2 Overloaded Servers and Racks

Server  $m$  is overloaded if it cannot process the local tasks that are routed to it without the help of other servers, and its local load cannot be distributed with under-loaded servers by load balancing. In order to describe the overloaded servers formally, we define the following notation:

$$\psi_n = \sum_{\bar{L}: n \in \bar{L}} \sum_{m=1}^M \lambda_{\bar{L}, n, m} \quad \forall n \in \mathcal{M}.$$

$\psi_n$  is the pseudo-arrival rate of type  $\bar{L}$  tasks routed to server  $n$  under the task types' arrival rates  $\{\lambda_{\bar{L}, n, m}\}$ . Server  $n$  is overloaded under the load decomposition  $\{\lambda_{\bar{L}, n, m}\}$  if  $\psi_n > \alpha$ . For a subset of servers  $\mathcal{S} \subset \mathcal{M}$ , we define  $\mathcal{L}_{\mathcal{S}}$  as the set of task types only local to servers of the set  $\mathcal{S}$ , that is,  $\mathcal{L}_{\mathcal{S}} = \{\bar{L} \in \mathcal{L} | \bar{L} \subset \mathcal{S}\}$ . On the other hand, for the same set  $\mathcal{S}$ , we define  $\mathcal{L}_{\mathcal{S}}^*$  as the set of task types that are local to at least one server of the set  $\mathcal{S}$ , that is,  $\mathcal{L}_{\mathcal{S}}^* = \{\bar{L} \in \mathcal{L} | \bar{L} \cap \mathcal{S} \neq \emptyset\}$ . Lemma 2 shows that there exists a load decomposition under which the truly overloaded set of servers, denoted by  $\mathcal{D}$ , do not receive any task type that has at least one local task out of this set  $\mathcal{D}$ , that is tasks that are local to a server out of the set  $\mathcal{D}$  are routed to under-loaded servers which are out of set  $\mathcal{D}$ .

**Lemma 2** *There exists a load decomposition  $\{\tilde{\lambda}_{\bar{L},n,m}\}$  for any arrival rate vector  $\lambda \in \bar{\Lambda}$  that satisfies the following two conditions [22]:*

$$\sum_{\bar{L}:m \in \bar{L}} \sum_{n:n \in \bar{L}} \frac{\tilde{\lambda}_{\bar{L},n,m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \sum_{n:n \in \bar{L}} \frac{\tilde{\lambda}_{\bar{L},n,m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \sum_{n:n \in \bar{L}} \frac{\tilde{\lambda}_{\bar{L},n,m}}{\gamma} < 1, \forall m \in \mathcal{M}, \quad (3.4)$$

$$\forall n \in \mathcal{D} = \{n \in \mathcal{M} : \psi_n \geq \alpha\}, \tilde{\lambda}_{\bar{L},n,m} = 0, \forall \bar{L} \notin \mathcal{L}_{\mathcal{D}}, \forall m \in \mathcal{M}. \quad (3.5)$$

**proof:** This lemma is lent from Lemma 2 in [22]. For any arrival rate vector  $\lambda$  in the capacity region, there exists a load decomposition  $\{\lambda_{\bar{L},n,m}\}$  that satisfies (3.4). The proof iteratively refines the load decomposition  $\{\lambda_{\bar{L},n,m}\}$  as follows. In each iteration, an appropriate amount of the local load to both temporary overloaded and under-loaded servers that are routed to overloaded servers are moved to under-loaded servers. What we mean by moving an appropriate amount of local load from an overloaded server to a local under-loaded server is that we move the shared load until either both servers become overloaded or under-loaded, or there is no more shared local load between them to move. By such load movements, the load on the whole system reduces, so (3.4) still holds for the new load decomposition in each iteration. We continue the load movements from overloaded servers to under-loaded ones until there is no more load local to both kinds of servers that are routed to overloaded servers. We name the ultimate local decomposition  $\{\tilde{\lambda}_{\bar{L},n,m}\}$ . For more details refer to the proof provided in the Appendix A, Section 7.1 in [22].

A rack is overloaded if the under-loaded servers in the rack cannot process the whole extra load on the overloaded servers in the same rack, and the local load on this rack cannot be distributed on other servers in under-loaded racks by load balancing. For a subset of racks,  $\mathcal{R} \subset \mathcal{K}$ , we define  $\mathcal{L}_{\mathcal{R}}$  as the set of task types that are only local to the servers in this set of racks, that is,  $\mathcal{L}_{\mathcal{R}} = \{\bar{L} \in \mathcal{L} | \forall m \in \bar{L}, K(m) \in \mathcal{R}\}$ . Formally, rack  $k$  is overloaded under a load decomposition  $\{\lambda_{\bar{L},n,m}\}$  if the following inequality holds:

$$\sum_{m:K(m)=k, \psi_m \geq \alpha} (\psi_m - \alpha) \geq \beta \sum_{n:K(n)=k, \psi_n < \alpha} \left(1 - \frac{\psi_n}{\alpha}\right). \quad (3.6)$$

The left-hand side of the above inequality is the extra load on overloaded servers of the overloaded rack  $k$  that cannot be served locally, and the right-hand side is the maximum rack-local service that can be afforded by the

under-loaded servers in the rack. Therefore, if (3.6) holds for rack  $k$ , the rack is overloaded and needs to receive remote service in order not to overflow. The following lemma for overloaded racks is an equivalent to Lemma 2 for overloaded servers.

**Lemma 3** *Assuming  $\beta^2 > \alpha\gamma$ , there exists a load decomposition  $\{\hat{\lambda}_{\bar{L},n,m}\}$  for any arrival rate vector  $\boldsymbol{\lambda} \in \bar{\Lambda}$  that satisfies not only conditions (3.4) and (3.5), but also the following condition. Note that  $\mathcal{O}$  is the set of overloaded racks satisfying equation (3.6) under the load decomposition  $\{\hat{\lambda}_{\bar{L},n,m}\}$ .*

$$\forall n \text{ s.t. } K(n) \in \mathcal{O}, \hat{\lambda}_{\bar{L},n,m} = 0, \forall \bar{L} \notin \mathcal{L}_{\mathcal{O}}, \forall m \in \mathcal{M}. \quad (3.7)$$

Analogous to overloaded servers, an overloaded rack only receives task types that are only local to the servers in the overloaded set of racks.

**proof:** The proof is similar to proof of Lemma 2. Starting from the load decomposition  $\{\tilde{\lambda}_{\bar{L},n,m}\}$  that satisfies both conditions (3.4) and (3.5), the load local to both overloaded and under-loaded racks that are routed to overloaded racks are iteratively moved to under-loaded ones. Note that this load can be moved from beneficiary servers of overloaded racks to beneficiary servers of under-loaded racks, or from helper servers of overloaded racks to helper servers of under-loaded racks. By moving the load from under-loaded servers in overloaded racks to under-loaded servers in under-loaded racks, there is a possibility that the under-loaded servers in under-loaded racks become overloaded. By moving  $\Delta$  amount of traffic from  $\mathcal{H}_o$  to  $\mathcal{H}_u$  (where  $\mathcal{H}_u$  is about to become  $\mathcal{B}_u$ ), the added rack-local load on the under-loaded rack is  $\frac{\Delta}{\beta}$  which means that the reduced amount of remote load on the under-loaded rack is  $\gamma \frac{\Delta}{\beta}$ . On the other hand, by removing  $\Delta$  amount of traffic from  $\mathcal{H}_o$ ,  $\mathcal{H}_o$  can process additional rack-local traffic of  $\beta \frac{\Delta}{\alpha}$ . This load movement reduces traffic on the whole system in case that  $\beta \frac{\Delta}{\alpha} > \gamma \frac{\Delta}{\beta}$ , which is equivalent to  $\beta^2 > \alpha\gamma$ . In summary, the condition  $\beta^2 > \alpha\gamma$  dictates that any load local to both overloaded and under-loaded racks should be routed to the under-loaded racks regardless of the load on servers. The load movement from an overloaded rack to an under-loaded one continues until both racks become overloaded or under-loaded, or no task local to both ones is routed to the overloaded rack. After load movements, some overloaded racks may become under-loaded or some under-loaded racks may become overloaded. By such load movements, the overall load on system reduces and both conditions (3.4)



and (3.5) hold for the obtained load decomposition. We call the ultimate load decomposition  $\{\hat{\lambda}_{L,n,m}\}$  that satisfies all the conditions in Lemma 3.

### 3.4.3 Ideal Load Decomposition

Under an arrival rate vector, servers can be classified into four types: helper servers in under-loaded racks, beneficiary servers in under-loaded racks, helper servers in overloaded racks, and beneficiary servers in overloaded racks. The definition of these four types of servers is as follows:

- Helpers in under-loaded racks ( $\mathcal{H}_u$ ): The set of under-loaded servers that are in under-loaded racks form the set  $\mathcal{H}_u$ . The tasks local to this set of servers all receive service locally. The remaining capacity of this set of servers is scheduled for processing rack-local and remote tasks.
- Beneficiaries in under-loaded racks ( $\mathcal{B}_u$ ): The set of overloaded servers that are in under-loaded racks form the set  $\mathcal{B}_u$ . The tasks local to this set of servers all receive service locally or rack-locally, but not remotely. The servers in this set only process local tasks, not rack-local or remote tasks.
- Helpers in overloaded racks ( $\mathcal{H}_o$ ): The set of under-loaded servers that are in overloaded racks form the set  $\mathcal{H}_o$ . The tasks local to this set of servers all receive service locally. The remaining capacity of this set of servers is scheduled for processing only rack-local tasks, not remote tasks.
- Beneficiaries in overloaded racks ( $\mathcal{B}_o$ ): The set of overloaded servers that are in overloaded racks form the set  $\mathcal{B}_o$ . The tasks local to this set of servers receive service locally, rack-locally, or remotely. The servers in this set only process local tasks, not rack-local or remote tasks.

Figure 3.3 depicts the four types of servers and their load under ideal load decomposition.

Unless no real helper or beneficiary servers, under-loaded or overloaded racks exist in a real system, we will use this concept in the heavy-traffic optimality proof. The following lemma formalizes the definition of four types of servers.

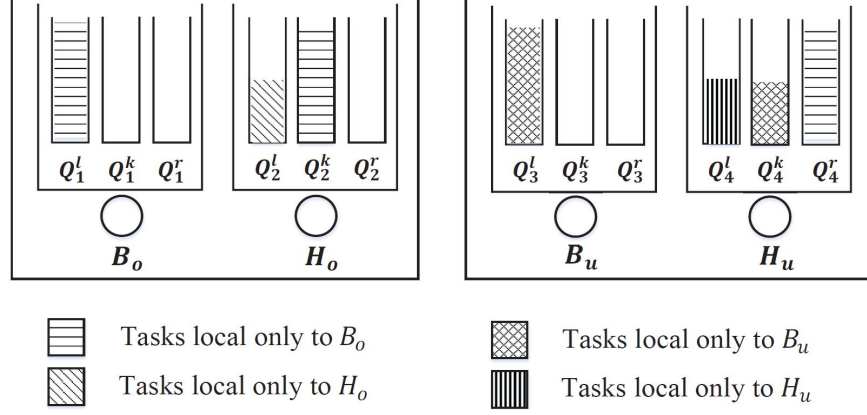


Figure 3.3: The four types of servers and their load under the ideal load decomposition.

**Lemma 4** *Assuming  $\beta^2 > \alpha\gamma$ , there exists a load decomposition  $\{\lambda_{\bar{L},n,m}^*\}$  for any arrival rate vector  $\lambda \in \Lambda$  that satisfies conditions (3.4), (3.5), and (3.7) in Lemmas 2 and 3, and under this load decomposition any server belongs to one of the four types described below. Note that  $\mathcal{O}$  and  $\mathcal{U}$  stand for the set of overloaded and under-loaded racks, respectively.*

$$\begin{aligned}
\mathcal{H}_u &= \{n : K(n) \in \mathcal{U} | \psi_n < \alpha, \text{ and } \forall \bar{L} \in \mathcal{L}, \forall m \neq n, \lambda_{\bar{L},n,m}^* = 0\}, \\
\mathcal{B}_u &= \{n : K(n) \in \mathcal{U} | \psi_n \geq \alpha, \text{ and } \forall \bar{L} \in \mathcal{L}, \forall m \neq n, \lambda_{\bar{L},n,m}^* = 0, \\
&\quad \text{and } \forall \bar{L} \in \mathcal{L}, \forall m \text{ s.t. } K(m) \neq K(n), \lambda_{\bar{L},n,m}^* = 0\}, \\
\mathcal{H}_o &= \{n : K(n) \in \mathcal{O} | \psi_n < \alpha, \text{ and } \forall \bar{L} \in \mathcal{L}, \forall m \neq n, \lambda_{\bar{L},n,m}^* = 0, \\
&\quad \text{and } \forall \bar{L} \in \mathcal{L}, \forall m \text{ s.t. } K(m) \neq K(n), \lambda_{\bar{L},n,m}^* = 0\}, \\
\mathcal{B}_o &= \{n : K(n) \in \mathcal{O} | \psi_n \geq \alpha, \text{ and } \forall \bar{L} \in \mathcal{L}, \forall m \neq n, \lambda_{\bar{L},n,m}^* = 0\}.
\end{aligned}$$

**proof:** In order to achieve the ideal load decomposition  $\{\lambda_{\bar{L},n,m}^*\}$  that satisfies the conditions in Lemma 4, we start from the load decomposition  $\{\hat{\lambda}_{\bar{L},n,m}\}$  which satisfies conditions (3.4), (3.5), and (3.7). The following four steps should be taken to achieve the ideal load decomposition:

1. If server  $n$  is an under-loaded server in an under-loaded rack, and  $\hat{\lambda}_{\bar{L},n,m} \neq 0$ , where  $m \neq n$ , we move this load to be scheduled locally at server  $n$ . This way, the local load to the under-loaded servers in under-loaded racks which were scheduled to be served rack-locally or remotely will be served locally, so the load on the whole system de-

creases (the rack-local or remote load on server  $n$  may be required to be rescheduled to other servers with removed load).

2. Under the updated load decomposition in the previous step, we offload any rack-local or remote load on any overloaded server  $n$  in an under-loaded rack. Hence, server  $n$  is only scheduled to process its local load. This way, there would be empty capacity on the servers that used to serve the overloaded load of server  $n$ . This empty capacity can be used for the previous rack-local and remote load that were being processed by server  $n$ . On the other hand, if local load to server  $n$  receive service remotely, it can be scheduled to under-loaded servers in the same rack of server  $n$  to receive service rack-locally. All these load movements reduce the load on the whole system.
3. Under the updated load decomposition in step 2, if the load local to an under-loaded server  $n$  in an overloaded rack receive rack-local or remote service, we reschedule it to be processed in its local server  $n$  (the rack-local or remote load on server  $n$  may be required to be rescheduled to other servers with removed load). Furthermore, we remove any remote load on server  $n$  to make more space for rack-local load of overloaded servers in the same rack of server  $n$  which used to receive service remotely. By these load adjustments, the overall load decreases on the whole system.
4. Under the updated load decomposition of step 3, the rack-local or remote load scheduled to overloaded servers in overloaded racks should be removed. Instead, local loads to these servers should be assigned to them. This way, the required remote service of these servers decreases more than the remote load that was removed from them. Hence, the overall load on the whole system decreases under this load movement.

# CHAPTER 4

## SIMULATION RESULTS

The performances of FCFS scheduler which is the Hadoop’s default scheduler, and Hadoop Fair Scheduler (HFS) are studied against the JSQ-MW algorithm in a system with two levels of data locality in [21]. As FCFS scheduler does not take the data locality into account, it performs worst than other algorithms like the JSQ-MW algorithm, specially at high loads. Hence, performance of FCFS is not given in the analysis. In this chapter, we compare three algorithms, the Balanced-Pandas algorithm, the JSQ-MaxWeight algorithm, and the Pandas algorithm implemented on a system with three levels of data locality through simulation. The configuration of the simulated system is as follows: we assume a continuous time system consisting of 10 racks ( $K = 10$ ), each of which consists of 50 servers, that is  $M = 500$ . The task arrival follows Poisson process, and the service time for a local, rack-local or remote task follows exponential distribution with rate  $\alpha = 1, \beta = 0.9$ , or  $\gamma = 0.5$ , respectively. The two times slowdown service for a remote task is consistent to the measurements in [8]. In our simulation environment, the three local servers to a task (the task type) is determined at the task’s arrival among a set of servers uniformly at random. The set of servers among which the local servers are chosen determines the load on the system. We investigate two traffic scenarios as follows:

1. In this traffic scenario, all the incoming task have their data chunks stored in three servers that are uniformly selected among the first five racks. This means that, the incoming load is uniformly distributed over all the 250 servers in the first five racks. If the mean arrival rate  $\lambda \equiv \sum_L \frac{\lambda_L}{M}$  is larger than or equal to 0.5, the first 250 servers in the first five racks are beneficiaries, and the first five racks are overloaded. The rest of the servers are helpers, and their five corresponding racks are under-loaded. The JSQ-MW algorithm achieves heavy-traffic optimality under this specific load. The Balanced-Pandas algorithm is also

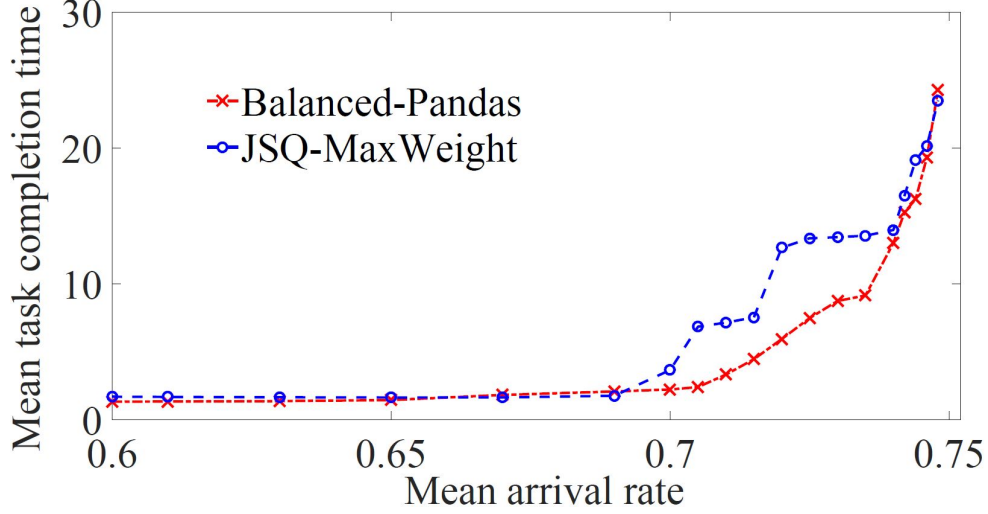


Figure 4.1: The mean task completion time versus the mean arrival rate for two algorithms, the Balanced-Pandas algorithm and the JSQ-MW algorithm, under a load that both algorithms minimize the mean task completion time at high loads.

heavy-traffic optimal in all loads. Therefore, both algorithms achieve the minimum mean task completion time in this traffic scenario. Figure 4.1 affirms the above statement.

- Under this load, 20 percent of the arriving tasks have their three local servers chosen uniformly at random from the first 10 servers of the first rack, and six percent of the incoming tasks have their three local servers chosen uniformly at random from the first 25 servers in the second rack. All the other 74 percent of the incoming tasks have their three local servers chosen uniformly at random from the rest of 465 servers in the system. This way, at high loads, the first 10 servers in the first rack and the 25 first servers in the second rack are beneficiaries, and the rest of servers are helpers. The first rack is overloaded and the rest of racks are under-loaded at high loads. Therefore, all four kinds of servers exist in the system under this traffic scenario at high loads. The mean task completion time of three algorithms is shown in Figure 4.2. As Figure 4.2 affirms, the Pandas algorithm is not throughput optimal as there exists other algorithms that can stabilize the system at higher loads. Calculating the capacity region, the system is stabilizable as long as  $\lambda < 0.9027$ . Both the Balanced-Pandas and JSQ-MW algorithms stabilize the system in this capacity region, but the Pandas algorithm

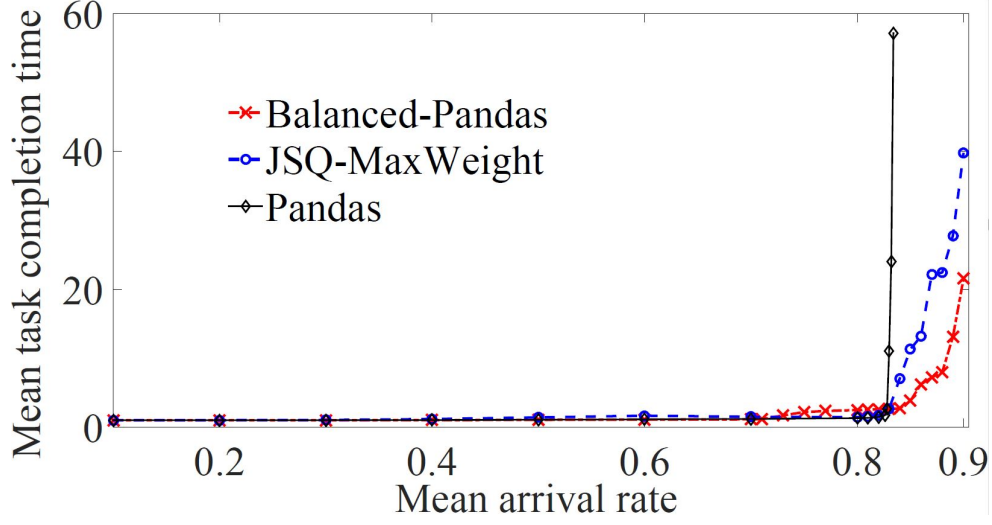


Figure 4.2: The mean task completion time versus the mean arrival rate for three algorithms the Balanced-Pandas, JSQ-MW, and Pandas algorithms under a general load that all four kinds of servers exist in the system.

makes the system unstable at load  $\lambda \approx 0.83$ , so the Pandas algorithm is not throughput optimal. Taking a more careful look at high loads, Figure 4.3 shows a significant up to fourfold outperformance of the Balanced-Pandas algorithm compared to the JSQ-MW algorithm. This fact affirms that the JSQ-MW algorithm is not a heavy-traffic optimal algorithm.

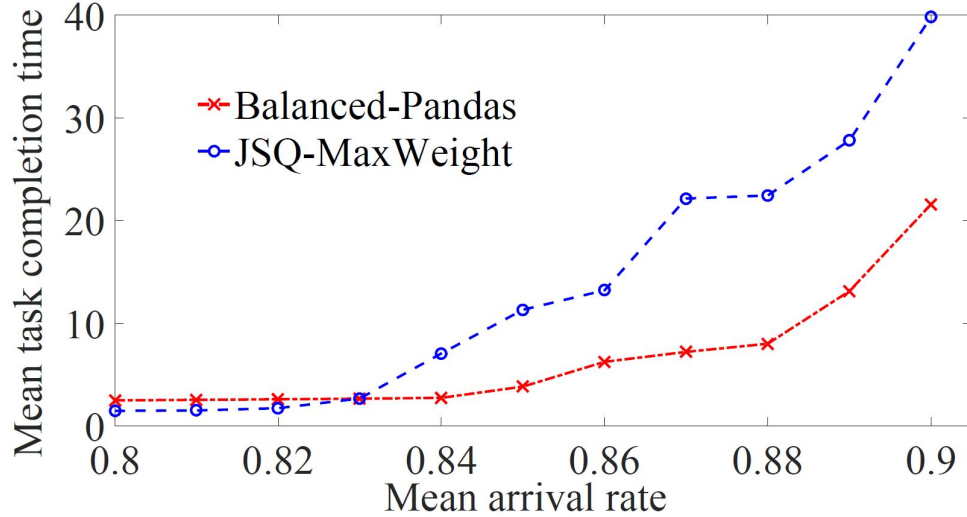


Figure 4.3: The performance of the JSQ-MW algorithm versus the Balanced-Pandas algorithm at high loads.

# APPENDIX A

## THEOREM PROOFS

### A.1 Proof of Theorem 1

We prove that the JSQ-MW algorithm stabilizes the system as long as the arrival rate vector is strictly inside the outer bound of the capacity region. This means that the outer bound  $\bar{\Lambda}$  is the capacity region and the JSQ-MW algorithm is a throughput optimal algorithm. Assume that  $\lambda \in \bar{\Lambda}$ , then there exists  $\delta > 0$  such that  $\lambda(1 + \delta) = \lambda' \in \bar{\Lambda}$ . As  $\lambda' \in \bar{\Lambda}$  there exists a load decomposition for  $\lambda'$ ,  $\{\lambda'_{\bar{L},n,m}\}$ , such that it satisfies the conditions in (3.1) specifically the following:

$$\sum_{\bar{L}:m \in \bar{L}} \sum_{n:n \in \bar{L}} \frac{\lambda'_{\bar{L},n,m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \sum_{n:n \in \bar{L}} \frac{\lambda'_{\bar{L},n,m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \sum_{n:n \in \bar{L}} \frac{\lambda'_{\bar{L},n,m}}{\gamma} < 1, \forall m.$$

By our choice of arrival rate vector, we have the following:

$$\{\lambda_{\bar{L},n,m}, \forall \bar{L} \in \mathcal{L}, \forall n, m \in \mathcal{M}\} = \left\{ \frac{\lambda'_{\bar{L},n,m}}{1 + \delta}, \forall \bar{L} \in \mathcal{L}, \forall n, m \in \mathcal{M} \right\}.$$

Hence, we conclude the following:

$$\sum_{\bar{L}:m \in \bar{L}} \sum_{n:n \in \bar{L}} \frac{\lambda_{\bar{L},n,m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \sum_{n:n \in \bar{L}} \frac{\lambda_{\bar{L},n,m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \sum_{n:n \in \bar{L}} \frac{\lambda_{\bar{L},n,m}}{\gamma} < \frac{1}{1 + \delta}, \forall m.$$

We define the pseudo arrival rate vector of servers,  $\psi = (\psi_1, \psi_2, \dots, \psi_M)$ , as follows:

$$\psi_n = \sum_{\bar{L}:n \in \bar{L}} \sum_{m=1}^M \lambda_{\bar{L},n,m}, \forall n. \quad (\text{A.1})$$

We use  $\psi$  as an intermediary to prove this theorem. In the proof, we will use the following three lemmas where the first two lemmas are analogous to Lemmas 2 and 3 in [21]. We eliminate the proofs of the first two lemmas as they mostly do not change for a system with three levels of data locality other than for a system with two levels of data locality.

**Lemma 5** *For any arrival rate vector strictly inside the capacity region,  $\lambda \in \bar{\Lambda}$ , and its corresponding pseudo arrival rate vector of servers  $\psi$  defined in (A.1), under the Joining the Shortest Queue routing policy we have the following inequality:*

$$\mathbb{E}[\langle \mathbf{Q}(t), \mathbf{A}(t) \rangle - \langle \mathbf{Q}(t), \psi \rangle | Z(t_0)] \leq 0, \quad \forall t_0, \text{ and } \forall t \geq t_0.$$

**Lemma 6** *For any arrival rate vector strictly inside the capacity region,  $\lambda \in \bar{\Lambda}$ , and its corresponding pseudo arrival rate vector of servers  $\psi$  defined in (A.1), under MaxWeight scheduling policy we have the following inequality:*

$$\begin{aligned} & \forall T > T_1, \text{ and } \forall t_0, \exists T_1 > 0 \text{ such that,} \\ & \mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( \langle \mathbf{Q}(t), \psi \rangle - \langle \mathbf{Q}(t), \mathbf{S}(t) \rangle \right) \middle| Z(t_0) \right] \leq -\theta_1 \|\mathbf{Q}(t_0)\|_1 + c_1, \end{aligned}$$

where the constants  $\theta_1 > 0$  and  $c_1$  are independent from  $Z(t_0)$ .

**Lemma 7**

$$\langle \mathbf{Q}(t), \mathbf{U}(t) \rangle < M^2, \quad \forall t.$$

**proof:** If  $U_m(t) > 0$ , then it implies that  $Q_m(t) < M$ . The reason is that queue  $m$  can receive at most  $M$  services at a time slot, so if there exists any unused services, the queue length should have been less than the whole services which is  $M$ . If  $U_m(t) = 0$ , then  $Q_m(t) \times U_m(t) = 0$ . Therefore,  $Q_m(t) \times U_m(t) < M \times U_m(t)$ . On the other hand it is clear that the summation of all the unused services in a time slot is less than or equal to the number of servers, that is  $\sum_{m \in \mathcal{M}} U_m(t) < M$ . Hence,  $\langle \mathbf{Q}(t), \mathbf{U}(t) \rangle < \sum_{m \in \mathcal{M}} M \times U_m(t) \leq M^2$  for any time slot  $t$ .

Proof of Theorem 1 mainly starts here. Choosing the Lyapunov function  $V_1(Z(t)) = \sum_{m \in \mathcal{M}} Q_m^2(t) = \|\mathbf{Q}(t)\|^2$ , it satisfies the conditions in the Foster-Lyapunov theorem to be non-negative, to be equal to zero only at  $\mathbf{Q}(t) = 0$ ,



and to go to infinity as any elements of  $\mathbf{Q}(t)$  goes to infinity. Then the expected  $T$ -time slot drift of the Lyapunov function is as follows:

$$\begin{aligned}
& \mathbb{E}[\Delta V_1(Z(t_0))] \\
&= \mathbb{E}[V_1(t_0 + T) - V_1(t_0) | Z(t_0)] \\
&= \mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( V_1(t+1) - V_1(t) \right) \middle| Z(t_0) \right] \\
&= \mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( \|\mathbf{Q}(t+1)\|^2 - \|\mathbf{Q}(t)\|^2 \right) \middle| Z(t_0) \right] \\
&= \mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( \|\mathbf{Q}(t) + (\mathbf{A}(t) - \mathbf{S}(t) + \mathbf{U}(t))\|^2 - \|\mathbf{Q}(t)\|^2 \right) \middle| Z(t_0) \right] \\
&= \mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( 2\langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle + 2\langle \mathbf{Q}(t), \mathbf{U}(t) \rangle \right. \right. \\
&\quad \left. \left. + \|\mathbf{A}(t) - \mathbf{S}(t) + \mathbf{U}(t)\|^2 \right) \middle| Z(t_0) \right].
\end{aligned}$$

We assumed that the task arrival process at a time slot is bounded with probability one and it is clear that the provided services and the unused services are also bounded. Hence,  $\|\mathbf{A}(t) - \mathbf{S}(t) + \mathbf{U}(t)\|^2$  is bounded. Also using Lemma 7, we have  $2\langle \mathbf{Q}(t), \mathbf{U}(t) \rangle + \|\mathbf{A}(t) - \mathbf{S}(t) + \mathbf{U}(t)\|^2 = c_2$ , where  $c_2 > 0$  is a constant independent of  $Z(t_0)$ . Then for any arrival rate vector  $\boldsymbol{\lambda} \in \bar{\Lambda}$  we can use the corresponding  $\boldsymbol{\psi}$  defined in (A.1) as an intermediary to write the expected Lyapunov function drift as follows:

$$\begin{aligned}
& \mathbb{E}[\Delta V_1(Z(t_0))] \\
&= \mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( 2\langle \mathbf{Q}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \right) \middle| Z(t_0) \right] + c_2 \\
&= 2\mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( \langle \mathbf{Q}(t), \mathbf{A}(t) \rangle - \langle \mathbf{Q}(t), \boldsymbol{\psi} \rangle \right) \middle| Z(t_0) \right] \\
&\quad + 2\mathbb{E} \left[ \sum_{t=t_0}^{t_0+T-1} \left( \langle \mathbf{Q}(t), \boldsymbol{\psi} \rangle - \langle \mathbf{Q}(t), \mathbf{S}(t) \rangle \right) \middle| Z(t_0) \right] + c_2 \\
&\stackrel{(a)}{\leq} -2\theta_1 \|\mathbf{Q}(t_0)\|_1 + c_2,
\end{aligned}$$

where (a) in the last inequality follows from Lemmas 5 and 6. Hence, for any  $\epsilon > 0$ , there exists  $T \geq T_0$  such that for any  $Z(t_0) \in \mathcal{P}^c$ , we have  $\mathbb{E}[V_1(Z(t_0 + T)) - V_1(Z(t_0))] \leq -\epsilon$ , where  $\mathcal{P}$  is a finite subset of state spaces and it is defined as  $\mathcal{P} = \left\{ Z = (\mathbf{Q}, \mathbf{f}) \mid \|\mathbf{Q}\|_1 \leq \frac{c_2 + \epsilon}{2\theta_1} \right\}$ . It is also obvious that the expected  $T$ -period drift of the Lyapunov function is bounded as long as  $Z(t_0) \in \mathcal{P}$ . Therefore, from the Foster-Lyapunov theorem we conclude that  $\{Z(t), t \geq 0\}$  is positive recurrent, which means that the JSQ-MaxWeight algorithm stabilizes the system under any arrival rate vector  $\boldsymbol{\lambda} \in \bar{\Lambda}$ . This means that  $\bar{\Lambda}$  and  $\Lambda$  are both the capacity region of the system.

## A.2 Proof of Theorem 2

The proof consists of three parts. First we obtain a lower bound for the expected queue length. Then, we prove the state space collapse of the queue lengths. Finally, we use a Lyapunov drift based approach presented in [23] which uses the state space collapse result to find an upper bound for the expected queue length. If the lower and upper bounds in the first and last steps match each other, the algorithm is heavy-traffic optimal in the traffic load that we considered. We summarize the proof as follows as it is similar to the proof in [21] for a system with two levels of data locality.

The lower bound on the expected sum of all queue lengths is obtained as follows. Assume we have a single server with the following arrival and service processes, respectively:

$$\sum_{\bar{L}} A_{\bar{L}}^{\epsilon}(t),$$

$$b_1(t) = \sum_{i \in \mathcal{B}_o} X_i(t) + \sum_{j \in \mathcal{H}_o} Y_j(t) + \sum_{n \in \mathcal{H}_u} V_n(t),$$

where  $\{X_i(t) \sim \text{Bern}(\alpha)\}_{i \in \mathcal{B}_o}$ ,  $\{Y_j(t) \sim \text{Bern}(\beta)\}_{j \in \mathcal{H}_o}$ , and  $\{V_n(t) \sim \text{Bern}(\gamma)\}_{n \in \mathcal{H}_u}$ .  $\{X_i\}_{i \in \mathcal{B}_o}$ ,  $\{Y_j\}_{j \in \mathcal{B}_u}$ , and  $\{V_n\}_{n \in \mathcal{H}_u}$  are independent from each other and each of them are i.i.d. processes. We define the variances of the arrival and service processes of this single server model as  $\sigma_1^{\epsilon} = \text{var}(\sum_{\bar{L}} A_{\bar{L}}^{\epsilon}(t))$  and  $\nu_1^2 = \text{var}(b_1(t))$ . It is obvious that the queue length of this single server/single queue model is stochastically smaller than the sum of the queue

lengths in the original system model with three levels of data locality. Hence, we have the following lower bound on the expected sum of queue lengths:

$$\mathbb{E}\left[\sum_{m=1}^M Q_m^\epsilon(t)\right] \geq \frac{(\sigma_1^\epsilon)^2 + \nu_1^2 + \epsilon^2}{2\epsilon} - \frac{M}{2}.$$

Therefore, if we let  $\epsilon$  go to zero to create the heavy-traffic regime, we have the following lower bound:

$$\liminf_{\epsilon \rightarrow 0^+} \epsilon \mathbb{E}\left[\sum_{m=1}^M Q_m^\epsilon(t)\right] \geq \frac{\sigma_1^2 + \nu_1^2}{2}. \quad (\text{A.2})$$

As  $\epsilon$  goes to zero, we expect the queue lengths of beneficiary servers to grow to infinity and have somehow equal lengths. We define the  $M$ -dimensional vector  $\mathbf{c}_1 \in \mathbb{R}^M$  and define the parallel and perpendicular components of  $\mathbf{Q}$  with respect to  $\mathbf{c}_1$  as follows:

$$\mathbf{c}_1(m) = \begin{cases} \frac{1}{\sqrt{M_{B_o}}} & \forall m \in \mathcal{B}_o \\ 0 & \text{else} \end{cases},$$

$$\mathbf{Q}_\parallel = \langle \mathbf{c}_1, \mathbf{Q} \rangle \mathbf{c}_1,$$

$$\mathbf{Q}_\perp = \mathbf{Q} - \mathbf{Q}_\parallel.$$

By taking  $V_2(\mathbf{Z}) = \|\mathbf{Q}_\perp\|$  as the Lyapunov function, we can show that using the JSQ-MW scheduling algorithm, the expected drift of this Lyapunov function is bounded, and becomes negative for sufficiently large  $\mathbf{Q}_\perp$ . Therefore, we have the following theorem for state space collapse (the proof for the following theorem is eliminated as it is similar to the corresponding theorem in [21]).

**Theorem 5** *There exists finite sequence of numbers  $\{C_r : r \in \mathbb{N}\}$  such that*

$$\mathbb{E}[\|\mathbf{Q}_\perp\|^r] \leq C_r, \quad \forall r \in \mathbb{N}.$$

We can then use the state space collapse result to prove the following upper bound for the mean sum of queue lengths:

$$\mathbb{E} \left[ \sum_{m=1}^M Q_m^\epsilon(t) \right] \leq \frac{(\sigma_1^\epsilon)^2 + \nu_1^2}{2\epsilon} + B^\epsilon,$$

where  $B^\epsilon = o(\frac{1}{\epsilon})$ . Hence, letting  $\epsilon$  to go to zero, we have the following upper bound in the heavy-traffic regime:

$$\limsup_{\epsilon \rightarrow 0^+} \epsilon \mathbb{E} \left[ \sum_{m=1}^M Q_m^\epsilon(t) \right] \leq \frac{\sigma^2 + \nu_1^2}{2}.$$

As the upper bound of the mean sum of queue lengths coincides with the lower bound under using the JSQ-MW algorithm, this algorithm is heavy-traffic optimal under the load we specified in Theorem 2 (but it is not heavy-traffic optimal in all traffic scenarios).

### A.3 Proof of Theorem 3

A corollary of Theorem 1 is that  $\Lambda$  is the capacity region of a system with three levels of data locality. Hence, to prove the throughput optimality of the Balanced-Pandas algorithm, it is enough to show that this scheduling algorithm can stabilize the system as long as the arrival rate vector is strictly inside the capacity region,  $\boldsymbol{\lambda} \in \Lambda$ . For any  $\boldsymbol{\lambda} \in \Lambda$ , there exists  $\delta > 0$  such that  $\boldsymbol{\lambda}' = \boldsymbol{\lambda}(1 + \delta) \in \Lambda$ . As  $\boldsymbol{\lambda}' \in \Lambda$ , there exists a load decomposition  $\{\lambda'_{\bar{L},m}\}$  such that it satisfies the following:

$$\sum_{\bar{L}:m \in \bar{L}} \frac{\lambda'_{\bar{L},m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \frac{\lambda'_{\bar{L},m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \frac{\lambda'_{\bar{L},m}}{\gamma} < 1, \quad \forall m \in \mathcal{M},$$

then by our choice of  $\boldsymbol{\lambda}'$  to be  $\boldsymbol{\lambda}(1 + \delta)$ , we have the following:

$$\sum_{\bar{L}:m \in \bar{L}} \frac{\lambda_{\bar{L},m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \frac{\lambda_{\bar{L},m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \frac{\lambda_{\bar{L},m}}{\gamma} < \frac{1}{1 + \delta}, \quad \forall m \in \mathcal{M}. \quad (\text{A.3})$$

Define the workload vector of servers,  $\mathbf{w} = (w_1, w_2, \dots, w_M)$ , under the load decomposition  $\{\lambda_{\bar{L},m}\}$  as follows:

$$w_m = \sum_{\bar{L}:m \in \bar{L}} \frac{\lambda_{\bar{L},m}}{\alpha} + \sum_{\bar{L}:m \in \bar{L}_k} \frac{\lambda_{\bar{L},m}}{\beta} + \sum_{\bar{L}:m \in \bar{L}_r} \frac{\lambda_{\bar{L},m}}{\gamma}, \quad \forall m \in \mathcal{M}. \quad (\text{A.4})$$

The workload on a server evolves as follows:

$$\begin{aligned}
W_m(t+1) &= \frac{Q_m^l(t+1)}{\alpha} + \frac{Q_m^k(t+1)}{\beta} + \frac{Q_m^r(t+1)}{\gamma} \\
&\stackrel{(a)}{=} \frac{Q_m^l(t) + A_m^l(t) - S_m^l(t)}{\alpha} + \frac{Q_m^k(t) + A_m^k(t) - S_m^k(t)}{\beta} \\
&\quad + \frac{Q_m^r(t) + A_m^r(t) - S_m^r(t) + U_m(t)}{\gamma} \\
&= W_m(t) + \left( \frac{A_m^l(t)}{\alpha} + \frac{A_m^k(t)}{\beta} + \frac{A_m^r(t)}{\gamma} \right) \\
&\quad - \left( \frac{S_m^l(t)}{\alpha} + \frac{S_m^k(t)}{\beta} + \frac{S_m^r(t)}{\gamma} \right) + \frac{U_m(t)}{\gamma},
\end{aligned}$$

where (a) follows from the queue evolution in (3.3). Define the pseudo arrival, service and unused service processes as  $\mathbf{A} = (A_1, A_2, \dots, A_M)$ ,  $\mathbf{S} = (S_1, S_2, \dots, S_M)$ , and  $\tilde{\mathbf{U}} = (\tilde{U}_1, \tilde{U}_2, \dots, \tilde{U}_M)$ , respectively, where

$$\begin{aligned}
A_m(t) &= \frac{A_m^l(t)}{\alpha} + \frac{A_m^k(t)}{\beta} + \frac{A_m^r(t)}{\gamma}, \quad \forall m \in \mathcal{M}, \\
S_m(t) &= \frac{S_m^l(t)}{\alpha} + \frac{S_m^k(t)}{\beta} + \frac{S_m^r(t)}{\gamma}, \quad \forall m \in \mathcal{M}, \\
\tilde{U}_m(t) &= \frac{U_m(t)}{\gamma}, \quad \forall m \in \mathcal{M}.
\end{aligned}$$

By the above definitions, we can write the dynamics of the queue workloads,  $\mathbf{W} = (W_1, W_2, \dots, W_M)$ , as follows:

$$\mathbf{W}(t+1) = \mathbf{W}(t) + \mathbf{A}(t) - \mathbf{S}(t) + \tilde{\mathbf{U}}(t). \quad (\text{A.5})$$

The following three lemmas will be used in the proof of Theorem 3.

**Lemma 8**

$$\langle \mathbf{W}(t), \tilde{\mathbf{U}}(t) \rangle = 0, \quad \forall t. \quad (\text{A.6})$$

**proof:** The expression simplifies as follows:

$$\langle \mathbf{W}(t), \tilde{\mathbf{U}}(t) \rangle = \sum_m \left( \frac{Q_m^l(t)}{\alpha} + \frac{Q_m^k(t)}{\beta} + \frac{Q_m^r(t)}{\gamma} \right) \frac{U_m(t)}{\gamma}.$$

Note that for any server  $m$ , if  $U_m(t) = 0$ , then  $(\frac{Q_m^l(t)}{\alpha} + \frac{Q_m^k(t)}{\beta} + \frac{Q_m^r(t)}{\gamma}) \frac{U_m(t)}{\gamma} = 0$ . Otherwise,  $U_m(t) > 0$  implies that all sub-queues of server  $m$  are empty which again results in  $U_m(t) = 0$ , then  $(\frac{Q_m^l(t)}{\alpha} + \frac{Q_m^k(t)}{\beta} + \frac{Q_m^r(t)}{\gamma}) \frac{U_m(t)}{\gamma} = 0$ . Therefore,  $\langle \mathbf{W}(t), \tilde{\mathbf{U}}(t) \rangle = 0$  for all time slots.

**Lemma 9** *For any arrival rate vector strictly inside the capacity region,  $\boldsymbol{\lambda} \in \Lambda$ , and the corresponding workload vector of servers  $\mathbf{w}$  defined in (A.4), we have the following inequality by using the Balanced-Pandas algorithm:*

$$\mathbb{E}[\langle \mathbf{W}(t), \mathbf{A}(t) \rangle - \langle \mathbf{W}(t), \mathbf{w} \rangle | Z(t)] \leq 0, \quad \forall t \geq 0. \quad (\text{A.7})$$

**proof:** We first define the minimum weighted workload for a task type,  $\bar{L} \in \mathcal{L}$  as follows:

$$W_{\bar{L}}^*(t) = \min_{m \in \mathcal{M}} \left\{ \frac{W_m(t)}{\alpha} I_{\{m \in \bar{L}\}}, \frac{W_m(t)}{\beta} I_{\{m \in \bar{L}_k\}}, \frac{W_m(t)}{\gamma} I_{\{m \in \bar{L}_r\}} \right\}. \quad (\text{A.8})$$

At the beginning of time slot  $t$ , an incoming task of type  $\bar{L}$  is routed to queue  $m^*$  with the minimum expected workload  $W_{\bar{L}}^*(t)$ . Therefore, for any task type  $\bar{L} \in \mathcal{L}$  we have the following:

$$\begin{aligned} \frac{W_m(t)}{\alpha} &\geq W_{\bar{L}}^*(t), \quad \forall m \in \bar{L}, \\ \frac{W_m(t)}{\beta} &\geq W_{\bar{L}}^*(t), \quad \forall m \in \bar{L}_k, \\ \frac{W_m(t)}{\gamma} &\geq W_{\bar{L}}^*(t), \quad \forall m \in \bar{L}_r. \end{aligned} \quad (\text{A.9})$$

In other words, task of type  $\bar{L}$  does not join a server  $m$  with weighted workload greater than  $W_{\bar{L}}^*$ . Then we have the following:

$$\begin{aligned}
& \mathbb{E}[\langle \mathbf{W}(t), \mathbf{A}(t) \rangle | Z(t)] \\
&= \mathbb{E} \left[ \sum_m W_m(t) \left( \frac{A_m^l(t)}{\alpha} + \frac{A_m^k(t)}{\beta} + \frac{A_m^r(t)}{\gamma} \right) \middle| Z(t) \right] \\
&= \mathbb{E} \left[ \sum_m W_m(t) \left( \frac{1}{\alpha} \sum_{\bar{L}: m \in \bar{L}} A_{\bar{L},m}(t) + \frac{1}{\beta} \sum_{\bar{L}: m \in \bar{L}_k} A_{\bar{L},m}(t) \right. \right. \\
&\quad \left. \left. + \frac{1}{\gamma} \sum_{\bar{L}: m \in \bar{L}_r} A_{\bar{L},m}(t) \right) \middle| Z(t) \right] \\
&\stackrel{(a)}{=} \mathbb{E} \left[ \sum_{\bar{L} \in \mathcal{L}} \left( \sum_{m: m \in \bar{L}} \frac{W_m(t)}{\alpha} A_{\bar{L},m}(t) + \sum_{m: m \in \bar{L}_k} \frac{W_m(t)}{\beta} A_{\bar{L},m}(t) \right. \right. \\
&\quad \left. \left. + \sum_{m: m \in \bar{L}_r} \frac{W_m(t)}{\gamma} A_{\bar{L},m}(t) \right) \middle| Z(t) \right] \\
&\stackrel{(b)}{=} \mathbb{E} \left[ \sum_{\bar{L} \in \mathcal{L}} W_{\bar{L}}^*(t) A_{\bar{L}}(t) \middle| Z(t) \right] \\
&= \sum_{\bar{L} \in \mathcal{L}} W_{\bar{L}}^*(t) \lambda_{\bar{L}},
\end{aligned} \tag{A.10}$$

where (a) is true by changing the order of the summations, and (b) follows from the Balanced-Pandas routing policy that task of type  $\bar{L}$  is routed to the queue with the minimum weighted workload,  $W_{\bar{L}}^*$ . On the other hand,

$$\begin{aligned}
& \mathbb{E}[\langle \mathbf{W}(t), \mathbf{w} \rangle | Z(t)] \\
&= \sum_m W_m(t) w_m \\
&= \sum_m W_m(t) \left( \sum_{\bar{L}: m \in \bar{L}} \frac{\lambda_{\bar{L},m}}{\alpha} + \sum_{\bar{L}: m \in \bar{L}_k} \frac{\lambda_{\bar{L},m}}{\beta} + \sum_{\bar{L}: m \in \bar{L}_r} \frac{\lambda_{\bar{L},m}}{\gamma} \right) \\
&\stackrel{(a)}{=} \sum_{\bar{L} \in \mathcal{L}} \left( \sum_{m: m \in \bar{L}} \frac{W_m(t)}{\alpha} \lambda_{\bar{L},m} + \sum_{m: m \in \bar{L}_k} \frac{W_m(t)}{\beta} \lambda_{\bar{L},m} + \sum_{m: m \in \bar{L}_r} \frac{W_m(t)}{\gamma} \lambda_{\bar{L},m} \right) \\
&\stackrel{(b)}{\geq} \sum_{\bar{L} \in \mathcal{L}} \sum_{m \in \mathcal{M}} W_{\bar{L}}^*(t) \lambda_{\bar{L},m} \\
&= \sum_{\bar{L} \in \mathcal{L}} W_{\bar{L}}^*(t) \lambda_{\bar{L}},
\end{aligned} \tag{A.11}$$

where (a) is true by changing the order of summations, and (b) follows from

(A.9). Lemma 9 is concluded from expressions (A.10) and (A.11).

**Lemma 10** *For any arrival rate vector strictly inside the capacity region,  $\boldsymbol{\lambda} \in \Lambda$ , and the corresponding workload vector of servers  $\mathbf{w}$  defined in (A.4), we have the following inequality by using the Balanced-Pandas algorithm:*

$$\mathbb{E}[\langle \mathbf{W}(t), \mathbf{w} \rangle - \langle \mathbf{W}(t), \mathbf{S}(t) \rangle | Z(t)] \leq -\theta_2 \|\bar{\mathbf{Q}}(t)\|_1, \quad \forall t \geq 0, \quad (\text{A.12})$$

where the constant  $\theta_2 > 0$  is independent of  $Z(t)$ .

**proof:** Using (A.3), the mean workload vector on servers defined in (A.4) can be bounded as follows:

$$w_m \leq \frac{1}{1 + \delta}, \quad \forall m \in \mathcal{M}.$$

Hence,

$$\mathbb{E}[\langle \mathbf{W}(t), \mathbf{w} \rangle | Z(t)] = \sum_m W_m(t) w_m \leq \frac{1}{1 + \delta} \sum_m W_m(t). \quad (\text{A.13})$$

We also have the following:

$$\begin{aligned} & \mathbb{E}[\langle \mathbf{W}(t), \mathbf{S}(t) \rangle | Z(t)] \\ &= \mathbb{E} \left[ \sum_m W_m(t) \left( \frac{S_m^l(t)}{\alpha} + \frac{S_m^k(t)}{\beta} + \frac{S_m^r(t)}{\gamma} \right) \middle| Z(t) \right] \\ &= \sum_m W_m(t) \mathbb{E} \left[ \frac{S_m^l(t)}{\alpha} + \frac{S_m^k(t)}{\beta} + \frac{S_m^r(t)}{\gamma} \middle| Z(t) \right] \\ &= \sum_m W_m(t) \mathbb{E} \left[ \sum_{i=0}^2 \mathbb{E} \left[ \frac{S_m^l(t)}{\alpha} + \frac{S_m^k(t)}{\beta} + \frac{S_m^r(t)}{\gamma} \middle| Z(t), \eta_m(t) = i \right] \middle| Z(t) \right] \\ &= \sum_m W_m(t) \left\{ \mathbb{E} \left[ \mathbb{E} \left[ \frac{S_m^l(t)}{\alpha} \middle| Z(t), \eta_m(t) = 0 \right] \middle| Z(t) \right] \right. \\ &\quad \left. + \mathbb{E} \left[ \mathbb{E} \left[ \frac{S_m^k(t)}{\beta} \middle| Z(t), \eta_m(t) = 1 \right] \middle| Z(t) \right] + \mathbb{E} \left[ \mathbb{E} \left[ \frac{S_m^r(t)}{\gamma} \middle| Z(t), \eta_m(t) = 2 \right] \middle| Z(t) \right] \right\} \\ &= \sum_m W_m(t). \end{aligned} \quad (\text{A.14})$$

Therefore,



$$\begin{aligned}
& \mathbb{E}[\langle \mathbf{W}(t), \mathbf{w} \rangle - \langle \mathbf{W}(t), \mathbf{S}(t) \rangle | Z(t)] \\
& \leq \frac{1}{1+\delta} \sum_m W_m(t) - \sum_m W_m(t) \\
& = -\frac{\delta}{1+\delta} \sum_m W_m(t) \\
& = -\frac{\delta}{1+\delta} \sum_m \left( \frac{Q_m^l(t)}{\alpha} + \frac{Q_m^k(t)}{\beta} + \frac{Q_m^r(t)}{\gamma} \right) \\
& \leq -\frac{\delta}{\alpha(1+\delta)} \sum_m (Q_m^l(t) + Q_m^k(t) + Q_m^r(t)) \\
& = -\theta_2 \|\bar{\mathbf{Q}}(t)\|_1.
\end{aligned}$$

Using Lemmas 8, 9, and 10, we prove Theorem 3 as follows. Assume the Lyapunov function is chosen as

$$V_3(Z(t)) = \|\mathbf{W}(t)\|^2,$$

then its expected drift is as follows:

$$\begin{aligned}
& \mathbb{E}[\Delta(Z(t))] \\
& = \mathbb{E}[V_3(t+1) - V_3(t) | Z(t)] \\
& = \mathbb{E} \left[ \|\mathbf{W}(t+1)\|^2 - \|\mathbf{W}(t)\|^2 \middle| Z(t) \right] \\
& \stackrel{(a)}{=} \mathbb{E} \left[ \|\mathbf{W}(t) + \mathbf{A}(t) - \mathbf{S}(t) + \tilde{\mathbf{U}}(t)\|^2 - \|\mathbf{W}(t)\|^2 \middle| Z(t) \right] \\
& = \mathbb{E} \left[ 2\langle \mathbf{W}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle + 2\langle \mathbf{W}(t), \tilde{\mathbf{U}}(t) \rangle + \|\mathbf{A}(t) - \mathbf{S}(t) + \tilde{\mathbf{U}}(t)\|^2 \middle| Z(t) \right] \\
& \stackrel{(b)}{=} 2\mathbb{E} \left[ \langle \mathbf{W}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle \middle| Z(t) \right] + c_3 \\
& = 2\mathbb{E} \left[ \langle \mathbf{W}(t), \mathbf{A}(t) \rangle - \langle \mathbf{W}(t), \mathbf{w} \rangle \middle| Z(t) \right] \\
& \quad + 2\mathbb{E} \left[ \langle \mathbf{W}(t), \mathbf{w} \rangle - \langle \mathbf{W}(t), \mathbf{S}(t) \rangle \middle| Z(t) \right] + c_3 \\
& \stackrel{(c)}{\leq} -2\theta_2 \|\bar{\mathbf{Q}}(t)\|_1 + c_3,
\end{aligned}$$

where (a) follows from (A.5), (b) follows from Lemma 8, and the fact that  $\mathbf{A}(t)$ ,  $\mathbf{S}(t)$ , and  $\tilde{\mathbf{U}}(t)$  are all bounded, and (c) is true by Lemmas 9 and 10. By

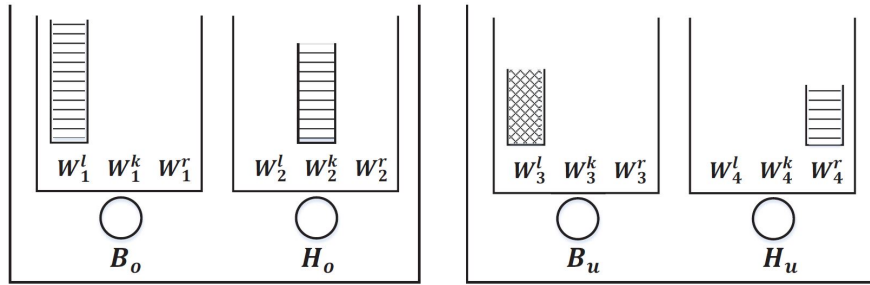
choosing any positive constant  $\epsilon > 0$ , let  $\mathcal{P} = \left\{ Z = (\mathbf{Q}, \mathbf{f}) \mid \|\mathbf{Q}\|_1 \leq \frac{c_3 + \epsilon}{2\theta_2} \right\}$ , where  $\mathcal{P}$  is a bounded subset of the state space. For any  $Z \in \mathcal{P}$ ,  $\Delta V_3(Z)$  is bounded and for any  $Z \in \mathcal{P}^c$ ,  $\Delta V_3(Z) \leq -\epsilon$ . Hence, for any  $\boldsymbol{\lambda} \in \Lambda$ , the Markov process  $\{Z(t), t \geq 0\}$  is positive recurrent and the Balanced-Pandas algorithm makes the system stable, which means that this algorithm is throughput optimal.

## A.4 Proof of Theorem 4

For simplicity, this proof is for the special case where  $\mathcal{O} \neq \emptyset$  and  $\mathcal{B}_u = \emptyset$ . For the general proof refer to [24]. The heavy-traffic optimality is driven through the following three steps:

1. Establishing the state-space collapse in the heavy traffic regime.
2. Finding a lower bound on the expected sum of the queue lengths as  $\epsilon \rightarrow 0$ .
3. Finding an upper bound on the expected sum of the queue lengths as  $\epsilon \rightarrow 0$ , which matches the lower bound found in step 2.

In heavy traffic regime, the system collapses to the one-dimensional state space vector shown in Figure A.1.



$$\alpha : \beta : \gamma = 1 : 0.8 : 0.5$$

$$w_1^l : w_2^k : w_3^l : w_4^r = \alpha : \beta : \frac{\alpha\gamma}{\beta} : \gamma = 1 : 0.8 : 0.625 : 0.5$$

Figure A.1: The queue compositions of the four types of servers in the heavy-traffic regime with  $\alpha = 1, \beta = 0.8, \gamma = 0.5$ . The workload at four types of servers maintain the ratio  $\alpha : \beta : \frac{\alpha\gamma}{\beta} : \gamma = 1 : 0.8 : 0.625 : 0.5$ .

Note that the prioritized service uniformly bounds the helper subsystem in heavy-traffic regime. This results in disappearance of local and rack-local queues of servers in the set  $\mathcal{H}_u$  and local queues of servers in the set  $\mathcal{H}_o$ . On the other hand, the weighted-workload routing policy distributes tasks that are only local to beneficiary servers in overloaded racks across  $\mathcal{B}_o$ ,  $\mathcal{H}_o$ , and  $\mathcal{H}_u$  in the ratio of  $\alpha : \beta : \gamma$  in terms of server workload. Furthermore, the tasks only local to servers in the set  $\mathcal{B}_u$  are just helped rack-locally by servers in the set  $\mathcal{H}_u$ , and the weighted-workload scheduling policy maintains the ratio  $\alpha : \beta$  in terms of workload on beneficiary and helper servers in under-loaded racks. Hence, the workload is distributed over servers in this proportion:  $W_1^l : W_2^k : W_3^l : W_4^r = \alpha : \beta : \frac{\alpha\gamma}{\beta} : \gamma$ .

Denote the local traffic on  $\mathcal{H}_u$  and  $\mathcal{H}_o$  by  $\sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} \lambda_{\bar{L}} \equiv \Phi_u \alpha$  and  $\sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_o}} \lambda_{\bar{L}} \equiv \Phi_o \alpha$ , respectively, where  $\mathcal{L}_{\mathcal{H}_u}^* = \{\bar{L} : \exists m \in \mathcal{H}_u \text{ s.t. } m \in \bar{L}\}$ , and  $\mathcal{L}_{\mathcal{H}_o} = \{\bar{L} : \forall m \in \bar{L}, m \in \mathcal{H}_o \cup \mathcal{B}_o, \text{ and } \exists n \in \mathcal{H}_o \text{ s.t. } n \in \bar{L}\}$ . Then, the heavy-traffic regime parameterized by  $\epsilon > 0$ , where  $\epsilon$  shows the distance of the arrival rate vector from the boundary of the capacity region, is defined as follows:

$$\sum_{\bar{L} \in \mathcal{L}_{\mathcal{B}_o}} \lambda_{\bar{L}} = \alpha |\mathcal{B}_o| + \beta (|\mathcal{H}_o| - \Phi_o) + \gamma (|\mathcal{H}_u| - \Phi_u) - \epsilon.$$

Consider the arrival process  $\{A_{\bar{L}}^{(\epsilon)}(t)\}_{\bar{L} \in \mathcal{L}}$  with arrival rate vector  $\boldsymbol{\lambda}^{(\epsilon)}$ . An assumption is made that the total local load for helpers is fixed, that is  $\{\lambda_{\bar{L}} : \bar{L} \in \mathcal{L}_{\mathcal{H}_u}^* \cup \mathcal{L}_{\mathcal{H}_o}\}$  is independent of  $\epsilon$ . Hence, the variance of  $\{A_{\bar{L}}^{(\epsilon)}(t)\}_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^* \cup \mathcal{L}_{\mathcal{H}_o}}$  is independent of  $\epsilon$ . On the other hand, the variance of the number of tasks that are only local to beneficiary servers in overloaded racks is denoted by  $(\sigma^{(\epsilon)})^2$  that converges to  $\sigma^2$  as  $\epsilon \downarrow 0$ , that is  $Var\left(\sum_{\bar{L} \in \mathcal{L}_{\mathcal{B}_o}} A_{\bar{L}}^{(\epsilon)}(t)\right) = (\sigma^{(\epsilon)})^2 \xrightarrow{\epsilon \rightarrow 0} \sigma^2$ . The system state under the Balanced-Pandas algorithm when the arrival rate is  $\boldsymbol{\lambda}^{(\epsilon)}$  is denoted by  $\{\mathbf{Z}^{(\epsilon)}(t) = (\bar{\mathbf{Q}}^{(\epsilon)}(t), \mathbf{f}^{(\epsilon)}(t)), t \geq 0\}$ . Then, the Markov chain  $\mathbf{Z}^{(\epsilon)}(t)$  is positive recurrent and has a steady state distribution as long as  $\boldsymbol{\lambda}^{(\epsilon)} \in \Lambda$ . The following theorem states that local and rack-local queues of  $\mathcal{H}_u$  and local queues of  $\mathcal{H}_o$  are uniformly bounded and the bound is independent of  $\epsilon$ .

**Theorem 6 (*Helper Queues*)**

$$\lim_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (Q_m^{l(\epsilon)}(t) + Q_m^{k(\epsilon)}(t)) \right] = 0,$$

$$\lim_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[ \sum_{m \in \mathcal{H}_o} Q_m^{l(\epsilon)}(t) \right] = 0.$$

The proof of Theorem 6 is given in Section A.4.1.

As the arrival rate vector approaches the boundary of the capacity region, that is  $\epsilon \downarrow 0$ , the mean sum of queue lengths approaches infinity in steady state, that is  $\mathbb{E} [\|\bar{\mathbf{Q}}\|] = \mathbb{E} [\sum_m (Q_m^l + Q_m^k + Q_m^r)] \rightarrow \infty$ . By Theorem 6, it is enough to consider the following to characterize the scaling order of  $\mathbb{E} [\bar{\mathbf{Q}}]$ :

$$\Phi = \sum_{m \in \mathcal{H}_u} Q_m^r + \sum_{m \in \mathcal{H}_o} (Q_m^k + Q_m^r) + \sum_{m \in \mathcal{B}_o} (Q_m^l + Q_m^k + Q_m^r).$$

Define  $\tilde{\mathbf{c}} \in \mathbb{R}_+^M$  as follows:

$$\tilde{c}_m = \begin{cases} \gamma, & \forall m \in \mathcal{H}_u \\ \beta, & \forall m \in \mathcal{H}_o \\ \alpha, & \forall m \in \mathcal{B}_o \end{cases}.$$

By defining  $\mathbf{c} = \frac{\tilde{\mathbf{c}}}{\|\tilde{\mathbf{c}}\|}$ , the parallel and perpendicular components of the steady-state weighted queue-length vector,  $\mathbf{W}$ , with respect to vector  $\mathbf{c}$  are as follows:

$$\mathbf{W}_{\parallel} = \langle \mathbf{c}, \mathbf{W} \rangle \mathbf{c}, \quad \mathbf{W}_{\perp} = \mathbf{W} - \mathbf{W}_{\parallel}.$$

The following theorem states that the deviation of  $\mathbf{W}$  from direction  $\mathbf{c}$  is bounded and is independent of the heavy-traffic parameter,  $\epsilon$ .

**Theorem 7 (*State Space Collapse*)**

*There exists a sequence of finite numbers  $\{C_r : r \in \mathbb{N}\}$  such that for each positive integer  $r$  we have the following:*

$$\mathbb{E} [\|\mathbf{W}_{\perp}\|^r] \leq C_r.$$

The proof of Theorem 7 is given in Section A.4.2.

Define the service process as follows:

$$b^{(\epsilon)}(t) = \sum_{i \in \mathcal{B}_o} X_i(t) + \sum_{j \in \mathcal{H}_o} Y_j(t) + \sum_{n \in \mathcal{H}_u} V_n(t),$$

where  $\{X_i(t)\}_{i \in \mathcal{B}_o}$ ,  $\{Y_j(t)\}_{j \in \mathcal{H}_o}$ , and  $\{V_n(t)\}_{n \in \mathcal{H}_u}$  are independent from each other and each process is i.i.d. and,

$$\begin{cases} X_i(t) \sim \text{Bern}(\alpha) & \forall i \in \mathcal{B}_o \\ Y_j(t) \sim \text{Bern}(\beta(1 - \rho_j^l)) & \forall j \in \mathcal{H}_o \\ V_n(t) \sim \text{Bern}(\gamma(1 - \rho_n)) & \forall n \in \mathcal{H}_u \end{cases}.$$

where  $\rho_j^l$  is the proportion of time that helper server  $j$  gives service to local tasks in steady state, and  $\rho_n$  is the proportion of time that helper server  $n$  gives service to local and rack-local tasks in steady state. Let  $\text{Var}(b^{(\epsilon)}(t)) = (\nu^{(\epsilon)})^2$  that converges to  $\nu^2$  as  $\epsilon \downarrow 0$ . Then, we have the following two theorems.

**Theorem 8 (*Lower Bound*)**

$$\mathbb{E} [\Phi^{(\epsilon)}(t)] \geq \frac{(\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2 + \epsilon^2}{2\epsilon} - \frac{M}{2}.$$

Hence,

$$\liminf_{\epsilon \downarrow 0} \epsilon \mathbb{E} [\Phi^{(\epsilon)}(t)] \geq \frac{\sigma^2 + \nu^2}{2}.$$

The proof of Theorem 8 is given in Section A.4.3.

**Theorem 9 (*Upper Bound*)**

$$\mathbb{E} [\Phi^{(\epsilon)}(t)] \leq \frac{(\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2}{2\epsilon} + B^{(\epsilon)},$$

where  $B^{(\epsilon)} = o(\frac{1}{\epsilon})$ , that is  $\lim_{\epsilon \downarrow 0} \epsilon B^{(\epsilon)} = 0$ ; hence,

$$\limsup_{\epsilon \downarrow 0} \epsilon \mathbb{E} [\Phi^{(\epsilon)}(t)] \leq \frac{\sigma^2 + \nu^2}{2}.$$

This upper bound matches with the lower bound found in Theorem 8.

The proof of Theorem 9 is given in Section A.4.4.

Note that,

$$\begin{aligned} & \mathbb{E} \left[ \sum_m (Q_m^{l(\epsilon)}(t) + Q_m^{k(\epsilon)}(t) + Q_m^{r(\epsilon)}(t)) \right] \\ &= \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (Q_m^{l(\epsilon)}(t) + Q_m^{k(\epsilon)}(t)) + \sum_{m \in \mathcal{H}_o} Q_m^{l(\epsilon)}(t) \right] + \mathbb{E} [\Phi^{(\epsilon)}(t)], \end{aligned}$$

where Theorems 8 and 9 give the coincidence of lower and upper bounds of the term  $\epsilon \mathbb{E} [\Phi^{(\epsilon)}(t)]$  as  $\epsilon \rightarrow 0$ . Using Theorem 6, the proof of heavy-traffic optimality is complete and Theorem A.4 is proved.

#### A.4.1 Proof of Theorem 6

Considering the system in steady state, define the following for any  $m \in \mathcal{H}_u$ ,

$$\begin{aligned} \hat{Q}_m(t) &= Q_m^l(t) + Q_m^k(t), \\ \hat{A}_m(t) &= A_m^l(t) + A_m^k(t), \\ \hat{S}_m(t) &= S_m^l(t) + S_m^k(t), \end{aligned}$$

where  $\hat{\mathbf{Q}}$  evolves as below:

$$\hat{\mathbf{Q}}(t+1) = \hat{\mathbf{Q}}(t) + \hat{\mathbf{A}}(t) - \hat{\mathbf{S}}(t).$$

Let  $\hat{F}_m(t) = F_m^l(t) + F_m^k(t)$ , where the ideal arrival process  $\mathbf{F}(t)$  is defined in the proof of Theorem 9. Now we can rewrite the dynamics of  $\hat{\mathbf{Q}}$  in the following way:

$$\hat{\mathbf{Q}}(t+1) = \hat{\mathbf{Q}}(t) + \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) + \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t).$$

Define the unit vector  $\mathbf{c}_h \in \mathbb{R}_+^{M_{\mathcal{H}_o}}$  as follows:

$$\mathbf{c}_h = \frac{1}{\sqrt{M_{\mathcal{H}_o}}} \underbrace{(1, 1, \dots, 1)}_{M_{\mathcal{H}_o}}.$$

The drift of the function  $\|\hat{\mathbf{Q}}_{||}\|^2 = \|\langle \mathbf{c}_h, \hat{\mathbf{Q}}_{||} \rangle\|^2$  is zero in steady state, so

$$\begin{aligned}
& 2\mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{S}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\
&= \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) \rangle^2 \right] + \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle^2 \right]
\end{aligned} \tag{A.15}$$

$$+ 2\mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) + \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle \right]. \tag{A.16}$$

The definition of the ideal arrival process yields that,

$$\langle \mathbf{c}_h, \hat{\mathbf{F}}(t) \rangle = \frac{1}{\sqrt{M_{\mathcal{H}_o}}} \sum_{m \in \mathcal{H}_u} F_m^l(t) = \frac{1}{\sqrt{M_{\mathcal{H}_o}}} \sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} A_{\bar{L}}(t).$$

Hence the sum of the ideal arrivals on  $\mathcal{H}_u$  and the queue lengths are independent.

$$\begin{aligned}
& \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{S}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\
&= \frac{1}{M_{\mathcal{H}_o}} \mathbb{E} \left[ \left( \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right) \left( \sum_{m \in \mathcal{H}_u} \hat{S}_m(t) \right) \right] - \frac{1}{M_{\mathcal{H}_o}} \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right] \left( \sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} \lambda_{\bar{L}} \right)
\end{aligned}$$

Note that  $\hat{S}_m(t) = S_m^l(t) + S_m^k(t)$  only depends on the state of the  $m$ -th queue, so

$$\begin{aligned}
& \mathbb{E} \left[ \left( \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right) \left( \sum_{m \in \mathcal{H}_u} \hat{S}_m(t) \right) \right] \\
&= \sum_{m \in \mathcal{H}_u} \mathbb{E} \left[ \hat{S}_m(t) \hat{Q}_m(t) \right] + \sum_{m \in \mathcal{H}_u} \mathbb{E} \left[ \hat{S}_m(t) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u: n \neq m} \hat{Q}_n(t) \right] \\
&= \sum_{m \in \mathcal{H}_u} \mathbb{E} \left[ \hat{S}_m(t) \hat{Q}_m(t) \right] + \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{S}_m(t) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u} \hat{Q}_n(t) \right] \\
&\quad - \sum_{m \in \mathcal{H}_u} \left( \mathbb{E} \left[ \hat{S}_m(t) \right] \mathbb{E} \left[ \hat{Q}_m(t) \right] \right).
\end{aligned} \tag{A.17}$$

The following lemma gives a lower bound on term  $\sum_{m \in \mathcal{H}_u} \mathbb{E} \left[ \hat{S}_m(t) \hat{Q}_m(t) \right]$ . For proof of the following lemma, refer to Lemma B.18 in [24].

**Lemma 11**

$$\sum_{m \in \mathcal{H}_u} \mathbb{E} \left[ \hat{S}_m(t) \hat{Q}_m(t) \right] \geq \sum_{m \in \mathcal{H}_u} \alpha \mathbb{E} \left[ \hat{Q}_m(t) \right] - C_1,$$

where  $C_1$  is a constant.

As we are studying the system in steady state,  $\mathbb{E} \left[ \hat{Q}_m(t+1) \right] = \mathbb{E} \left[ \hat{Q}_m(t) \right]$ ,  $\forall m \in \mathcal{H}_u$ , so

$$\mathbb{E} \left[ \hat{A}_m(t) - \hat{S}_m(t) \right] = \mathbb{E} \left[ \hat{Q}_m(t+1) - \hat{Q}_m(t) \right] = 0,$$

which results in  $\mathbb{E} \left[ \hat{S}_m(t) \right] = \mathbb{E} \left[ \hat{A}_m(t) \right]$ , so

$$\begin{aligned} & \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{S}_m(t) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u} \hat{Q}_n(t) \right] - \sum_{m \in \mathcal{H}_u} \left( \mathbb{E} \left[ \hat{S}_m(t) \right] \mathbb{E} \left[ \hat{Q}_m(t) \right] \right) \\ &= \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{A}_m(t) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u} \hat{Q}_n(t) \right] - \sum_{m \in \mathcal{H}_u} \left( \mathbb{E} \left[ \hat{A}_m(t) \right] \mathbb{E} \left[ \hat{Q}_m(t) \right] \right) \\ &= \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (A_m^l(t) + A_m^k(t)) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u} \hat{Q}_n(t) \right] \\ & \quad - \sum_{m \in \mathcal{H}_u} \left( \mathbb{E} \left[ A_m^l(t) + A_m^k(t) \right] \mathbb{E} \left[ \hat{Q}_m(t) \right] \right) \\ &\stackrel{(a)}{\geq} \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (A_m^l(t) + A_m^k(t)) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u} \hat{Q}_n(t) \right] - \sum_{m \in \mathcal{H}_u} \left( \alpha \rho_h^* \mathbb{E} \left[ \hat{Q}_m(t) \right] \right), \end{aligned}$$

where (a) follows from the following lemma. Refer to Lemma B.17 in [24] for the proof of the following lemma.

**Lemma 12**

$\forall m \in \mathcal{H}_o \cup \mathcal{H}_u$ ,  $\exists 0 \leq \rho_h < 1$ , where  $\rho_h$  does not depend on  $\epsilon$ , s.t.

$$\mathbb{E} \left[ \frac{A_m^l}{\alpha} \right] \leq \rho_h.$$

Then we have the following:



$$\begin{aligned}
& \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{S}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\
& \geq \frac{1}{M_{\mathcal{H}_o}} \left\{ \sum_{m \in \mathcal{H}_u} \alpha \mathbb{E} \left[ \hat{Q}_m(t) \right] - C_1 + \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (A_m^l(t) + A_m^k(t)) \right] \mathbb{E} \left[ \sum_{n \in \mathcal{H}_u} \hat{Q}_n(t) \right] \right. \\
& \quad \left. - \sum_{m \in \mathcal{H}_u} \left( \alpha \rho_h^* \mathbb{E} \left[ \hat{Q}_m(t) \right] \right) - \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right] \left( \sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} \lambda_{\bar{L}} \right) \right\} \\
& = \frac{1}{M_{\mathcal{H}_o}} \left\{ \alpha(1 - \rho_h^*) + \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (A_m^l(t) + A_m^k(t)) \right] - \sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} \lambda_{\bar{L}} \right\} \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right] \\
& \quad - \frac{C_1}{M_{\mathcal{H}_o}}.
\end{aligned}$$

The following can be driven from proof of Lemma 17 (or Lemma B.23 in [24]):

$$\begin{aligned}
& \sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} \lambda_{\bar{L}} - \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (A_m^l(t) + A_m^k(t)) \right] \\
& = \mathbb{E} \left[ A_{\mathcal{H}_u \mathcal{H}_o}^l + A_{\mathcal{H}_u \mathcal{B}_o}^l + A_{\mathcal{H}_u \mathcal{H}_o}^k + A_{\mathcal{H}_o \mathcal{H}_o}^k + A_{\mathcal{H}_u \mathcal{H}_u}^k + A_{\mathcal{H}_u \mathcal{B}_o}^r + A_{\mathcal{H}_u \mathcal{H}_o}^r + A_{\mathcal{H}_u \mathcal{H}_u}^r \right] \\
& \leq C\epsilon,
\end{aligned} \tag{A.18}$$

where  $C$  is a constant that is only a function of  $\alpha, \beta$ , and  $\gamma$ . Furthermore, the definition of the ideal arrival process yields the following:

$$\sum_{m \in \mathcal{H}_u} \hat{F}_m(t) = \sum_{\bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*} \lambda_{\bar{L}} \geq \sum_{m \in \mathcal{H}_u} \hat{A}_m(t).$$

Then we have the following:

$$\begin{aligned}
& \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{S}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\
& \geq \frac{1}{M_{\mathcal{H}_o}} [\alpha(1 - \rho_h^*) - C\epsilon] \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right] - \frac{C_1}{M_{\mathcal{H}_o}}.
\end{aligned} \tag{A.19}$$

The number of arriving tasks and services are bounded, so

$$\mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) \rangle^2 \right] \leq C_2, \quad (\text{A.20})$$

$$\mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle^2 \right] \leq C_3, \quad (\text{A.21})$$

where  $C_2 > 0$  and  $C_3 > 0$  are constants not depending on  $\epsilon$ . Then we have the following for the term in equation (A.16):

$$\begin{aligned} & \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) + \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\ &= \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{Q}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle \right] + \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\ &\stackrel{(a)}{\leq} \mathbb{E} \left[ \langle \mathbf{c}_h, \hat{\mathbf{F}}(t) - \hat{\mathbf{S}}(t) \rangle \langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle \right] \\ &\stackrel{(b)}{\leq} C_4, \end{aligned} \quad (\text{A.22})$$

where (a) is true as  $\langle \mathbf{c}_h, \hat{\mathbf{A}}(t) - \hat{\mathbf{F}}(t) \rangle \leq 0$ , and (b) is true as the number of task arrival is bounded, and  $C_4$  is a constant.

From equations (A.16), (A.19), (A.20), (A.21), and (A.22), the following is derived:

$$\frac{2}{M_{\mathcal{H}_o}} [\alpha(1 - \rho_h^*) - C\epsilon] \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right] \leq \frac{2C_1}{M_{\mathcal{H}_o}} + C_2 + C_3 + 2C_4,$$

so for any  $0 < \epsilon < \frac{\alpha(1 - \rho_h^*)}{C}$ , the following is true:

$$\mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} \hat{Q}_m(t) \right] \leq \frac{C_5}{\alpha(1 - \rho_h^*) - C\epsilon},$$

where  $C_5 = C_1 + (C_2 + C_3 + 2C_4) \frac{M_{\mathcal{H}_o}}{2}$ , so

$$\lim_{\epsilon \downarrow 0} \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (Q_m^{l(\epsilon)}(t) + Q_m^{k(\epsilon)}(t)) \right] \leq \frac{C_5}{\alpha(1 - \rho_h^*)},$$

that is equivalent to the following:

$$\lim_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[ \sum_{m \in \mathcal{H}_u} (Q_m^{l(\epsilon)}(t) + Q_m^{k(\epsilon)}(t)) \right] = 0.$$

Similarly, we can prove the following:

$$\lim_{\epsilon \downarrow 0} \epsilon \mathbb{E} \left[ \sum_{m \in \mathcal{H}_o} Q_m^{l(\epsilon)}(t) \right] = 0.$$

#### A.4.2 Proof of Theorem 7

The following lemma is given for ideal load decomposition for the case  $\mathcal{O} \neq \emptyset$ . For proof of this lemma, refer to Lemma B.19 in [24].

##### Lemma 13

$\exists \lambda_0 > 0$  not depending on  $\epsilon$ , such that:

1. *Defining*

$$\omega_m = \sum_{\bar{L}: m \in \bar{L}} \frac{\lambda_{\bar{L},m}^*}{\alpha} + \sum_{\bar{L}: m \in \bar{L}_k} \frac{\lambda_{\bar{L},m}^*}{\beta} + \sum_{\bar{L}: m \in \bar{L}_r} \frac{\lambda_{\bar{L},m}^*}{\gamma},$$

we have the following:

$$\omega_m = \begin{cases} 1 - \gamma\epsilon_0, & \forall m \in \mathcal{H}_u \\ 1 - \beta\epsilon_0, & \forall m \in \mathcal{H}_o, \\ 1 - \alpha\epsilon_0, & \forall m \in \mathcal{B}_o \end{cases}$$

where  $\epsilon_0 = \frac{\epsilon}{\|\bar{\mathbf{c}}\|^2}$ .

2. *Denote the set of task types that are only local to  $\mathcal{B}_o$  by  $\mathcal{L}_{\mathcal{B}_o}$ . Then,*

$$\begin{aligned} & \forall \bar{L} \in \mathcal{L}_{\mathcal{B}_o}, \text{ and } \forall m \in \{i \in \mathcal{M} | i \in \bar{L}, \text{ or } i \in \mathcal{H}_u, \text{ or } i \in \bar{L}_k \cap \mathcal{H}_o\}, \\ & \exists \kappa > 0, \text{ independent of } \epsilon, \text{ such that: } \lambda_{\bar{L},m}^* = \sum_{n \in \bar{L}} \lambda_{\bar{L},n,m}^* \geq \kappa. \end{aligned}$$

For the evenly loaded scenario, the following three lemmas are used. For the proof of these lemmas, refer to Lemmas B.20, B.21, and B.22 in [24].

**Lemma 14** *Using the Balanced-Pandas algorithm, we have the following:*

$$\mathbb{E} [\langle \mathbf{W}(t), \mathbf{A}(t) \rangle - \langle \mathbf{W}, \boldsymbol{\omega} \rangle | Z(t)] \leq -\lambda_{\min} \|\mathbf{W}_{\perp}(t)\|, \quad \forall t \geq 0,$$

where  $\lambda_{\min} > 0$  is a constant not depending on  $\epsilon$ .

**Lemma 15** *Using the Balanced-Pandas algorithm, we have the following:*

$$\mathbb{E}[\langle \mathbf{W}(t), \boldsymbol{\omega} \rangle - \langle \mathbf{W}(t), \mathbf{S}(t) \rangle | Z(t)] = -\frac{\epsilon}{\|\hat{\mathbf{c}}\|} \langle \mathbf{c}, \mathbf{W} \rangle, \quad \forall t \geq 0.$$

**Lemma 16**

$$\mathbb{E}[\langle \mathbf{c}, \mathbf{W}(t) \rangle \langle \mathbf{c}, \mathbf{A}(t) - \mathbf{S}(t) \rangle | Z(t)] \geq -\frac{\epsilon}{\|\hat{\mathbf{c}}\|} \langle \mathbf{c}, \mathbf{W} \rangle, \quad \forall t \geq 0.$$

In order to prove Theorem 7, consider the following Lyapunov function:

$$F(Z) = \|\mathbf{W}_\perp\|.$$

The drift of this Lyapunov function is given as below:

$$\Delta F(Z) \leq \frac{1}{2\|\mathbf{W}_\perp\|} (\Delta V(Z) - \Delta V_\parallel(Z)),$$

where  $\Delta V(Z)$  and  $\Delta V_\parallel(Z)$  are the drifts for Lyapunov functions  $V(Z) = \|\mathbf{W}\|^2$  and  $V_\parallel(Z) = \|\mathbf{W}_\parallel\|^2$ , respectively. Then, we have the following:

$$\begin{aligned} & \mathbb{E}[\Delta V(Z(t)) - \Delta V_\parallel(Z(t)) | Z(t)] \\ & \leq 2\mathbb{E}[\langle \mathbf{W}(t), \mathbf{A}(t) - \mathbf{S}(t) \rangle - \langle \mathbf{c}, \mathbf{W}(t) \rangle \langle \mathbf{c}, \mathbf{A}(t) - \mathbf{S}(t) \rangle | Z(t)] + C_1. \end{aligned}$$

Using Lemmas 14, 15, and 16, we get the following upper bound on  $\mathbb{E}[\Delta F(Z(t)) | Z(t)]$ :

$$\mathbb{E}[\Delta F(Z(t)) | Z(t)] \leq -\lambda_0 + \frac{C}{\|\mathbf{W}_\perp(t)\|},$$

where  $\lambda_0 > 0$  and  $C > 0$  are constants not depending on  $\epsilon$ . This last inequality satisfies the negative drift condition, so there exists finite series of constants  $\{C'_r\}_{r \in \mathbb{N}}$  such that  $\mathbb{E}[\|\mathbf{W}_\perp^{(\epsilon)}(t)\|^r] \leq C'_r$  for any  $\epsilon \in (0, M\alpha)$ .

#### A.4.3 Proof of Theorem 8

The lower bound on  $\mathbb{E}[\Phi^{(\epsilon)}(t)]$  can be driven by constructing a system with a single server and a single queue with arrival process  $\left\{ \sum_{\bar{L} \in \mathcal{L}_{B_o}} A_{\bar{L}}^{(\epsilon)}(t), t \geq 0 \right\}$

and service process  $\left\{b^{(\epsilon)}(t) = \sum_{i \in \mathcal{B}_o} X_i(t) + \sum_{j \in \mathcal{H}_o} Y_j(t) + \sum_{n \in \mathcal{H}_u} V_n(t), t \geq 0\right\}$ . Denote the queue length of the constructed system by  $\Psi^{(\epsilon)}(t)$ . By the definitions of  $X_i, Y_j$ , and  $V_n$ ,  $\mathbb{E}[\sum_{i \in \mathcal{B}_o} X_i(t)]$ ,  $\mathbb{E}[\sum_{j \in \mathcal{H}_o} Y_j(t)]$ , and  $\mathbb{E}[\sum_{n \in \mathcal{H}_u} V_n(t)]$  are the maximum amount of local, rack-local, and remote services that can be given to  $\sum_{\bar{L} \in \mathcal{L}_{\mathcal{B}_o}} A_{\bar{L}}^{(\epsilon)}(t)$ . Hence, it is obvious that in steady state  $\Psi^{(\epsilon)}(t)$  is stochastically smaller than or equal to  $\Phi^{(\epsilon)}(t)$ . Then the lower bound on  $\mathbb{E}[\Phi^{(\epsilon)}(t)]$  is derived by using Lemma 4 in [23].

#### A.4.4 Proof of Theorem 9

The ideal scheduling, service, and arrival processes are defined as follows:

**Ideal Scheduling Decision Process  $\eta'(t)$ :** Under ideal scheduling, a beneficiary server in an over-loaded rack is only giving service to its local tasks queued in its local sub-queue, and an idle helper server in an overloaded rack that has no local tasks in its local sub-queue is only scheduled to give service to its rack-local tasks queued at its rack-local sub-queue. In other words,

$$\begin{aligned} \forall m \in \mathcal{B}_o, \eta'_m(t) &= 0, \\ \forall m \in \mathcal{H}_o, \eta'_m(t) &= \eta_m(t) \text{ if } \eta_m(t) = 0, \text{ and } \eta'_m(t) = 1 \text{ if } f_m(t^-) = -1, Q_m^l(t) = 0 \\ \forall m \in \mathcal{H}_u, \eta'_m(t) &= \eta_m(t). \end{aligned}$$

**Ideal Service Process  $\mathbf{D}(t)$ :**

$$\forall m \in \mathcal{B}_o, D_m^l(t) = X_m^l(t), D_m^k(t) = 0, D_m^r(t) = 0,$$

where  $X_m^l(t) \sim \text{Bern}(\alpha)$ , and each process  $X_m^l(t)$  is i.i.d. and is coupled with  $S_m(t)$  as follows: If  $\eta_m(t) = 0$ ,  $X_m^l(t) = S_m^l(t)$ ; if  $\eta_m(t) = 1$ ,  $X_m^l(t) = 0$  when  $S_m^k(t) = 0$ , and  $X_m^l(t) \sim \text{Bern}(\frac{\alpha}{\beta})$  when  $S_m^k(t) = 1$ ; if  $\eta_m(t) = 2$ ,  $X_m^l(t) = 0$  when  $S_m^r(t) = 0$ , and  $X_m^l(t) \sim \text{Bern}(\frac{\alpha}{\gamma})$  when  $S_m^k(t) = 1$ . Furthermore,

$$\forall m \in \mathcal{H}_o, D_m^l(t) = S_m^l(t), D_m^k(t) = Y_m^k(t), D_m^r(t) = 0,$$

where  $Y_m^k(t) \sim \text{Bern}(\beta I_{\{\eta_m(t) \neq 0\}})$  and each process  $Y_m^k(t)$  is i.i.d. Finally,  $\forall m \in \mathcal{H}_u, \mathbf{D}_m(t) = \mathbf{S}_m(t)$ ,

$$D_m^l(t) = S_m^l(t), D_m^k(t) = S_m^k(t), D_m^r(t) = S_m^r(t).$$

**Ideal Arrival Process  $\mathbf{F}(t)$ :** Ideally, any task type that has a local server in the set  $\mathcal{H}_u$  should receive service locally. In other words,  $\forall \bar{L} \in \mathcal{L}_{\mathcal{H}_u}^*$ , task of type  $\bar{L}$  is routed to one of its local servers in the set  $\mathcal{H}_u$ . Hence, unwanted arrivals  $\sum_{m:m \notin \bar{L}, m \notin \mathcal{H}_u} A_{\bar{L},m}$  should be reassigned evenly among their local servers in  $\mathcal{H}_u$ . Similarly,  $\forall \bar{L} \in \mathcal{L}_{\mathcal{H}_o}$ , the task of type  $\bar{L}$  should ideally be assigned to its local servers in  $\mathcal{H}_o$ , that is unwanted arrivals  $\sum_{m:m \notin \bar{L}, m \notin \mathcal{H}_o} A_{\bar{L},m}$  should be reassigned evenly among their local servers in  $\mathcal{H}_o$ . On the other hand,  $\forall \bar{L} \in \mathcal{L}_{\mathcal{B}_o}$ , task of type  $\bar{L}$  should either receive service locally from a server in  $\mathcal{B}_o$ , or rack-locally from a server in  $\mathcal{H}_o$ , or remotely from a server in  $\mathcal{H}_u$ . Hence, we reassign tasks so that the above conditions hold in the ideal case. Then, the dynamics of  $\tilde{\mathbf{Q}}$  can be written as follows:

$$\tilde{\mathbf{Q}}(t+1) = \tilde{\mathbf{Q}}(t) + \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) + \tilde{\mathbf{V}}(t),$$

where  $\tilde{\mathbf{V}}(t) = \tilde{\mathbf{A}}(t) - \tilde{\mathbf{F}}(t) + \tilde{\mathbf{D}}(t) - \tilde{\mathbf{S}}(t) + \tilde{\mathbf{U}}(t)$ . Note that in steady state, we have the following:

$$\begin{aligned} & 2\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{D}}(t) - \tilde{\mathbf{F}}(t) \rangle \right] \\ &= \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle^2 \right] + \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle^2 \right] \\ & \quad + 2\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) + \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right]. \end{aligned} \tag{A.23}$$

On the other hand,

$$\begin{aligned} \Phi^{(\epsilon)}(t) &\leq \sum_{m \in \mathcal{H}_u} \gamma \frac{Q_m^r}{\gamma} + \sum_{m \in \mathcal{H}_o} \beta \left( \frac{Q_m^k}{\beta} + \frac{Q_m^r}{\gamma} \right) + \sum_{m \in \mathcal{B}_o} \alpha \left( \frac{Q_m^l}{\alpha} + \frac{Q_m^k}{\beta} + \frac{Q_m^r}{\gamma} \right) \\ &= ||\tilde{\mathbf{c}}|| \langle \mathbf{c}, \tilde{\mathbf{Q}} \rangle, \end{aligned} \tag{A.24}$$

so in order to find an upper bound on  $\mathbb{E} [\Phi^{(\epsilon)}(t)]$ , we need to find an upper bound on  $\mathbb{E} [\langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle]$ . To this aim, We start by analyzing different terms in equation (A.23). For simplicity, we omit the superscripts  $(\epsilon)$  in the following equations temporarily. The definition of ideal arrival process yields the following:

$$\langle \hat{\mathbf{c}}, \tilde{\mathbf{F}}(t) \rangle = \sum_{m \in \mathcal{B}_o} \alpha \cdot \frac{F_m^l(t)}{\alpha} + \sum_{m \in \mathcal{H}_o} \beta \cdot \frac{F_m^k(t)}{\beta} + \sum_{m \in \mathcal{H}_u} \gamma \cdot \frac{F_m^r(t)}{\gamma} = \sum_{\bar{L} \in \mathcal{L}_{\mathcal{B}_o}} A_{\bar{L}}(t).$$

Therefore,

$$\mathbb{E} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{F}}(t) \rangle \right] = \sum_{\bar{L} \in \mathcal{L}_{\mathcal{B}_o}} \lambda_{\bar{L}},$$

$$\text{Var} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{F}}(t) \rangle \right] = (\sigma^{(\epsilon)})^2.$$

The definition of ideal service process yields the following:

$$\langle \hat{\mathbf{c}}, \tilde{\mathbf{D}}(t) \rangle = \sum_{m \in \mathcal{B}_o} \alpha \cdot \frac{D_m^l(t)}{\alpha} + \sum_{m \in \mathcal{H}_o} \beta \cdot \frac{D_m^k(t)}{\beta} + \sum_{m \in \mathcal{H}_u} \gamma \cdot \frac{D_m^r(t)}{\gamma}.$$

For a server  $m$ , define  $\rho_m^l$  as the proportion of time in steady state the server spends on giving local service to the tasks queued in its local sub-queue. Then we have the following:

$$\begin{aligned} \mathbb{E} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{D}}(t) \rangle \right] &= \alpha M_{\mathcal{B}_o} + \sum_{m \in \mathcal{H}_o} \beta (1 - \rho_m^l) + \sum_{m \in \mathcal{H}_u} \gamma (1 - \rho_m^l), \\ \text{Var} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{D}}(t) \rangle \right] &= \alpha(1 - \alpha)M_{\mathcal{B}_o} + \sum_{m \in \mathcal{H}_o} \beta (1 - \rho_m^l) [1 - \beta (1 - \rho_m^l)] \\ &\quad + \sum_{m \in \mathcal{H}_u} \gamma (1 - \rho_m^l) [1 - \gamma (1 - \rho_m^l)] \\ &= (\nu^{(\epsilon)})^2. \end{aligned}$$

Then,

$$\begin{aligned} \mathbb{E} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{D}}(t) \rangle \right] - \mathbb{E} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{F}}(t) \rangle \right] &= \epsilon + \sum_{m \in \mathcal{H}_o} \beta (\rho_m^{l(\epsilon)} - \rho_m^l) + \sum_{m \in \mathcal{H}_u} \gamma (\rho_m^{l(\epsilon)} - \rho_m^l) \\ &= \epsilon + \delta, \end{aligned}$$

where  $\delta = \sum_{m \in \mathcal{H}_o} \beta (\rho_m^{l(\epsilon)} - \rho_m^l) + \sum_{m \in \mathcal{H}_u} \gamma (\rho_m^{l(\epsilon)} - \rho_m^l) \geq 0$ , and  $\delta \rightarrow 0$  as  $\epsilon \rightarrow 0$ . Hence, we have the following for the left-hand side term in equation (A.23):

$$\begin{aligned} &\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{D}}(t) - \tilde{\mathbf{F}}(t) \rangle \right] \\ &= \frac{1}{\|\hat{\mathbf{c}}\|} \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \left( \langle \hat{\mathbf{c}}, \tilde{\mathbf{D}}(t) \rangle - \langle \hat{\mathbf{c}}, \tilde{\mathbf{F}}(t) \rangle \right) \right] \\ &= \frac{\epsilon + \delta}{\|\hat{\mathbf{c}}\|} \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \right]. \end{aligned} \tag{A.25}$$

The first term on the right-hand side of equation (A.23) can be simplified

as follows:

$$\begin{aligned}
& \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle^2 \right] \\
&= \frac{1}{\|\hat{\mathbf{c}}\|^2} \left\{ \text{Var} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{D}}(t) \rangle \right] + \text{Var} \left[ \langle \hat{\mathbf{c}}, \tilde{\mathbf{F}}(t) \rangle \right] + \left( \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \right] \right)^2 \right\} \\
&= \frac{1}{\|\hat{\mathbf{c}}\|^2} \left\{ (\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2 + (\epsilon + \delta)^2 \right\}.
\end{aligned} \tag{A.26}$$

The second term on the right-hand side of equation (A.23) is upper bounded as the following lemma suggests.

**Lemma 17**

$$\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle^2 \right] \leq C\epsilon,$$

where  $C$  is a constant that does not depend on  $\epsilon$ .

In order to find an upper bound on the third term on the right-hand side of equation (A.23), we do the following. The system is in steady state, so

$$\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) + \tilde{\mathbf{V}}(t) \rangle \right] = \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t+1) - \tilde{\mathbf{Q}}(t) \rangle \right] = 0,$$

so

$$\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] = \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \right] = \frac{\epsilon}{M\alpha},$$

then

$$\begin{aligned}
& \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] \\
& \leq \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] \\
& \leq \frac{C_A}{\sqrt{M\alpha}} \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] \\
& = \frac{C_A}{M\sqrt{M\alpha^2}} \epsilon,
\end{aligned}$$

so we have the following upper bound on the third term on the right-hand side of equation (A.23):

$$\begin{aligned}
& \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) + \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] \\
&= \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] + \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] \\
& \leq \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] + \frac{C_A}{M\sqrt{M\alpha^2}} \epsilon.
\end{aligned}$$



We then simplify the term  $\langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle$  as follows:

$$\begin{aligned}
& \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \\
&= \langle \tilde{\mathbf{Q}}(t), \tilde{\mathbf{V}}(t) \rangle - \langle \tilde{\mathbf{Q}}_{\perp}(t), \tilde{\mathbf{V}}_{\perp}(t) \rangle \\
&= \langle \tilde{\mathbf{Q}}(t), \tilde{\mathbf{D}}(t) - \tilde{\mathbf{S}}(t) \rangle + \langle \tilde{\mathbf{Q}}(t), \tilde{\mathbf{A}}(t) - \tilde{\mathbf{F}}(t) \rangle + \langle \tilde{\mathbf{Q}}(t), \tilde{\mathbf{V}}(t) \rangle - \langle \tilde{\mathbf{Q}}_{\perp}(t), \tilde{\mathbf{V}}_{\perp}(t) \rangle.
\end{aligned} \tag{A.27}$$

The following two lemmas give a bound for the first two terms in equation (A.27).

**Lemma 18**

$$\mathbb{E} \left[ \langle \tilde{\mathbf{Q}}(t), \tilde{\mathbf{D}}(t) - \tilde{\mathbf{S}}(t) \rangle \right] = 0.$$

**Lemma 19**

$$\mathbb{E} \left[ \langle \tilde{\mathbf{Q}}(t), \tilde{\mathbf{A}}(t) - \tilde{\mathbf{F}}(t) \rangle \right] = o(\epsilon).$$

For the proof of Lemmas 18 and 19 refer to Lemmas B.24 and B.25 in [24].

By Lemma 8, the third term in equation (A.27) is equal to zero. In order to find an upper bound for the last term in equation (A.27), we first find an upper bound on  $\mathbb{E} \left[ \|\tilde{\mathbf{V}}(t)\|^2 \right]$ . Using lemma 17, we have the following:

$$\mathbb{E} \left[ \|\tilde{\mathbf{V}}(t)\|^2 \right] \leq R\epsilon,$$

where  $R$  is a constant not depending on  $\epsilon$ . Then we use Cauchy-Schwartz inequality and the result on state space collapse to find the following bound:

$$\mathbb{E} \left[ -\langle \tilde{\mathbf{Q}}_{\perp}(t), \tilde{\mathbf{V}}_{\perp}(t) \rangle \right] \leq \sqrt{\mathbb{E} \left[ \|\tilde{\mathbf{Q}}_{\perp}(t)\|^2 \right] \mathbb{E} \left[ \|\tilde{\mathbf{V}}_{\perp}(t)\|^2 \right]} \leq \sqrt{C'_2 R \epsilon}.$$

Hence, we have the following bound on the last term on the right-hand side of equation (A.23):

$$\mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) + \tilde{\mathbf{F}}(t) - \tilde{\mathbf{D}}(t) \rangle \langle \mathbf{c}, \tilde{\mathbf{V}}(t) \rangle \right] \leq \frac{C_A}{M\sqrt{M}\alpha^2} \epsilon + \sqrt{C'_2 R \epsilon} + o(\epsilon). \tag{A.28}$$

Using Lemma 17, equations (A.23), (A.25), (A.26), and (A.28) in equation (A.24) and bringing the superscript  $^{(\epsilon)}$  back in the equations, we have the

following:

$$\begin{aligned}
& 2 \frac{\epsilon + \delta}{\|\hat{\mathbf{c}}\|} \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \right] \\
& \leq \frac{1}{\|\hat{\mathbf{c}}\|^2} \left( (\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2 + (\epsilon + \delta)^2 \right) + C\epsilon + \frac{2C_A}{M\sqrt{M}\alpha^2} \epsilon + 2\sqrt{C'_2 R \epsilon} + 2o(\epsilon),
\end{aligned}$$

since  $\delta \geq 0$ ,

$$\begin{aligned}
& 2 \frac{\epsilon}{\|\hat{\mathbf{c}}\|} \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \right] \\
& \leq \frac{1}{\|\hat{\mathbf{c}}\|^2} \left( (\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2 + (\epsilon + \delta)^2 \right) + C\epsilon + \frac{2C_A}{M\sqrt{M}\alpha^2} \epsilon + 2\sqrt{C'_2 R \epsilon} + 2o(\epsilon),
\end{aligned}$$

so,

$$\begin{aligned}
& \|\hat{\mathbf{c}}\| \mathbb{E} \left[ \langle \mathbf{c}, \tilde{\mathbf{Q}}(t) \rangle \right] \\
& \leq \frac{(\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2 + (\epsilon + \delta)^2}{2\epsilon} + \left( \frac{C}{2} + \frac{C_A}{M\sqrt{M}\alpha^2} \right) \|\hat{\mathbf{c}}\|^2 + \|\hat{\mathbf{c}}\|^2 \sqrt{\frac{C'_2 R}{\epsilon}} + o(1).
\end{aligned}$$

Note that,

$$\begin{aligned}
& \mathbb{E} [\Phi^{(\epsilon)}(t)] \\
& = \mathbb{E} \left[ \sum_{m \in \mathcal{B}_o} (Q_m^{l(\epsilon)}(t) + Q_m^{k(\epsilon)}(t) + Q_m^{r(\epsilon)}(t)) + \sum_{m \in \mathcal{H}_o} (Q_m^{k(\epsilon)}(t) + Q_m^{r(\epsilon)}(t)) \right. \\
& \quad \left. + \sum_{m \in \mathcal{H}_u} Q_m^{r(\epsilon)}(t) \right] \\
& \leq \mathbb{E} \left[ \sum_{m \in \mathcal{B}_o} \alpha \left( \frac{Q_m^l}{\alpha} + \frac{Q_m^k}{\beta} + \frac{Q_m^r}{\gamma} \right) + \sum_{m \in \mathcal{H}_o} \beta \left( \frac{Q_m^k}{\beta} + \frac{Q_m^r}{\gamma} \right) + \sum_{m \in \mathcal{H}_u} \gamma \frac{Q_m^r}{\gamma} \right] \\
& = \mathbb{E} \left[ \|\tilde{\mathbf{c}}\| \langle \mathbf{c}, \tilde{\mathbf{Q}} \rangle \right].
\end{aligned}$$

Hence,

$$\mathbb{E} [\Phi^{(\epsilon)}(t)] \leq \frac{(\sigma^{(\epsilon)})^2 + (\nu^{(\epsilon)})^2 + (\epsilon + \delta)^2}{2\epsilon} + B^{(\epsilon)},$$

where  $B^{(\epsilon)} = \left( \frac{C}{2} + \frac{C_A}{M\sqrt{M}\alpha^2} \right) \|\hat{\mathbf{c}}\|^2 + \|\hat{\mathbf{c}}\|^2 \sqrt{\frac{C'_2 R}{\epsilon}} + o(1)$ , i.e.,  $B^{(\epsilon)} = o(\frac{1}{\epsilon})$ . This proves Theorem 9 as  $\epsilon \rightarrow 0$ .

## REFERENCES

- [1] Q. Xie, A. Yekkehkhany, and Y. Lu, “Scheduling with multi-level data locality: Throughput and heavy-traffic optimality,” in *Computer Communications, IEEE INFOCOM 2016-The 35th Annual IEEE International Conference on*. IEEE, 2016, pp. 1–9.
- [2] “Facebook.com,” <http://facebook.com>.
- [3] “Twitter.com,” <http://twitter.com>.
- [4] “Linkedin.com,” <http://linkedin.com>.
- [5] G. Ananthanarayanan, S. Agarwal, S. Kandula, A. Greenberg, I. Stoica, D. Harlan, and E. Harris, “Scarlett: Coping with skewed content popularity in mapreduce clusters,” in *Proceedings of the Sixth Conference on Computer Systems*. ACM, 2011, pp. 287–300.
- [6] G. Ananthanarayanan, A. Ghodsi, A. Wang, D. Borthakur, S. Kandula, S. Shenker, and I. Stoica, “Pacman: Coordinated memory caching for parallel jobs,” in *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*. USENIX Association, 2012.
- [7] Q. Xie, M. Pundir, Y. Lu, C. L. Abad, and R. H. Campbell, “Pandas: Robust locality-aware scheduling with stochastic delay optimality,” *IEEE/ACM Transactions on Networking*, 2016.
- [8] M. Zaharia, D. Borthakur, J. Sen Sarma, K. Elmeleegy, S. Shenker, and I. Stoica, “Delay scheduling: A simple technique for achieving locality and fairness in cluster scheduling,” in *Proceedings of the 5th European Conference on Computer Systems*. ACM, 2010, pp. 265–278.
- [9] Q. Xie and Y. Lu, “Degree-guided map-reduce task assignment with data locality constraint,” in *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium on*. IEEE, 2012, pp. 985–989.
- [10] “Apache hadoop,” June 2011.

- [11] C. He, Y. Lu, and D. Swanson, “Matchmaking: A new mapreduce scheduling technique,” in *Cloud Computing Technology and Science (CloudCom), 2011 IEEE Third International Conference on*. IEEE, 2011, pp. 40–47.
- [12] S. Ibrahim, H. Jin, L. Lu, B. He, G. Antoniu, and S. Wu, “Maestro: Replica-aware map scheduling for mapreduce,” in *Cluster, Cloud and Grid Computing (CCGrid), 2012 12th IEEE/ACM International Symposium on*. IEEE, 2012, pp. 435–442.
- [13] M. Isard, V. Prabhakaran, J. Currey, U. Wieder, K. Talwar, and A. Goldberg, “Quincy: Fair scheduling for distributed computing clusters,” in *Proceedings of the ACM SIGOPS 22nd Symposium on Operating Systems Principles*. ACM, 2009, pp. 261–276.
- [14] J. Jin, J. Luo, A. Song, F. Dong, and R. Xiong, “Bar: An efficient data locality driven task scheduling algorithm for cloud computing,” in *Cluster, Cloud and Grid Computing (CCGrid), 2011 11th IEEE/ACM International Symposium on*. IEEE, 2011, pp. 295–304.
- [15] M. S. Squillante, C. H. Xia, D. D. Yao, and L. Zhang, “Threshold-based priority policies for parallel-server systems with affinity scheduling,” in *Proceedings of the American Control Conference, 2001.*, vol. 4. IEEE, 2001, pp. 2992–2999.
- [16] A. Mandelbaum and A. L. Stolyar, “Scheduling flexible servers with convex delay costs: Heavy-traffic optimality of the generalized  $c\mu$ -rule,” *Operations Research*, vol. 52, no. 6, pp. 836–855, 2004.
- [17] J. M. Harrison and M. J. López, “Heavy traffic resource pooling in parallel-server systems,” *Queueing Systems*, vol. 33, no. 4, pp. 339–368, 1999.
- [18] J. M. Harrison, “Heavy traffic analysis of a system with parallel servers: Asymptotic optimality of discrete-review policies,” *Annals of Applied Probability*, pp. 822–848, 1998.
- [19] S. L. Bell, R. J. Williams et al., “Dynamic scheduling of a system with two parallel servers in heavy traffic with resource pooling: Asymptotic optimality of a threshold policy,” *The Annals of Applied Probability*, vol. 11, no. 3, pp. 608–649, 2001.
- [20] A. L. Stolyar, “Maxweight scheduling in a generalized switch: State space collapse and workload minimization in heavy traffic,” *Annals of Applied Probability*, pp. 1–53, 2004.

- [21] W. Wang, K. Zhu, L. Ying, J. Tan, and L. Zhang, “Maptask scheduling in mapreduce with data locality: Throughput and heavy-traffic optimality,” *IEEE/ACM Transactions on Networking*, vol. 24, no. 1, pp. 190–203, 2016.
- [22] Q. Xie and Y. Lu, “Priority algorithm for near-data scheduling: Throughput and heavy-traffic optimality,” in *Computer Communications (INFOCOM), 2015 IEEE Conference on*. IEEE, 2015, pp. 963–972.
- [23] A. Eryilmaz and R. Srikant, “Asymptotically tight steady-state queue length bounds implied by drift conditions,” *Queueing Systems*, vol. 72, no. 3-4, pp. 311–359, 2012.
- [24] Q. Xie, “Scheduling and resource allocation for clouds: Novel algorithms, state space collapse and decay of tails,” Ph.D. dissertation, University of Illinois at Urbana-Champaign, 2016.