

# An Efficient and Fair Multi-Resource Allocation Mechanism for Heterogeneous Servers

Jalal Khamse-Ashari\*, Ioannis Lambadaris\*, George Kesidis†, Bhuvan Urgaonkar† and Yiqiang Zhao‡

\*Dept. of Systems and Computer Engineering, Carleton University, Ottawa, Canada

†School of EECS, Pennsylvania State University, State College, PA, USA

‡School of Math and Statistics, Carleton University, Ottawa, Canada

Emails: \*{jalalkhamseashari,ioannis}@sce.carleton.ca, †{gik2, buu1}@psu.edu ‡zhao@math.carleton.ca



**Abstract**—Efficient and fair allocation of multiple types of resources is a crucial objective in a cloud/distributed computing cluster. Users may have diverse resource needs. Furthermore, diversity in server properties/capabilities may mean that only a subset of servers may be usable by a given user. In platforms with such heterogeneity, we identify important limitations in existing multi-resource fair allocation mechanisms, notably Dominant Resource Fairness (DRF) and its follow-up work. To overcome such limitations, we propose a new *server-based approach*; each server allocates resources by maximizing a per-server *utility function*. We propose a specific class of utility functions which, when appropriately parameterized, adjusts the trade-off between efficiency and fairness, and captures a variety of fairness measures (such as our recently proposed Per-Server Dominant Share Fairness). We establish conditions for the proposed mechanism to satisfy certain properties that are generally deemed desirable, e.g., envy-freeness, sharing incentive, bottleneck fairness, and Pareto optimality. To implement our resource allocation mechanism, we develop an iterative algorithm which is shown to be globally convergent. Finally, we show how the proposed mechanism could be implemented in a distributed fashion. We carry out extensive trace-driven simulations to show the enhanced performance of our proposed mechanism over the existing ones.

## 1 INTRODUCTION

Cloud computing has become increasingly popular as it provides a cost-effective alternative to proprietary high performance computing systems. As the workloads to data-centers housing cloud computing platforms are intensively growing, developing an efficient and fair allocation mechanism which guarantees quality-of-service for different workloads has become increasingly important. Efficient and fair resource allocation in such a shared computing system is particularly challenging because of (a) the presence of multiple types of resources, (b) diversity in the workloads' needs for these resources, (c) heterogeneity in the resource capacities of servers, and (d) placement constraints on which servers may be used by a workload. In the following four paragraphs we briefly elaborate on each of these complexities.

The *multi-resource needs* of cloud workloads imply that conventional single-resource oriented notions of fairness

are inadequate [1]. Dominant Resource Fairness (DRF) is the first allocation mechanism which describes a notion of fairness for allocating multiple types of resources for a single server system. Using DRF users receive a *fair share* of their *dominant resource* [1]. Of all the resources requested by the user (for every unit of work called a *task*), its dominant resource is the one with the highest demand when demands are expressed as fractions of the overall resource capacities. DRF is shown to achieve several properties that are commonly considered desirable from a multi-resource *fair* allocation mechanism.

*Heterogeneity of workloads' resource demands* is another complexity which results in a trade-off between efficiency and fairness. Specifically, heterogeneity of users' demands may preclude some resources from being fully utilized. Hence, the DRF allocation may result in a poor resource utilization even when there is only one server [2], [3], [4]. To address this issue, [2] proposed to allocate resources by applying the so-called  $\alpha$ -proportional fairness (instead of max-min fairness [5]) on dominant shares. The proposed mechanism, when appropriately parameterized, adjusts the trade-off between efficiency and fairness. However, it is applicable only to a *single server/resource-pool*.

In the case of *multiple heterogeneous servers*, there are several studies investigating/extending DRF allocation when there is *no placement constraint* [6], [7], [8]. In all of these works, fairness is defined in terms of a *global metric*, a scalar parameter defined in terms of different resources across all servers. E.g., [7] presents an extension to DRF where the dominant resource for each user is identified as if all resources were concatenated at one server, and subsequently the resources are allocated by applying max-min fairness on the dominant shares. Since such a global metric may not perfectly capture the impact of server heterogeneity, such approaches may lead to an inefficient resource utilization (see Section 2.2 for further discussions and Section 2.3 for an illustrative example). Moreover, such mechanisms may not be readily implementable in a distributed fashion [9], as each

server needs information on the available resources over all servers. Such information may not be available at each server, especially in a cloud computing environment where the resource capacities (and even activity of servers) might be churning.

There are limited works in the literature investigating multi-resource fair allocation in the presence of *user placement constraints* [10], [11]. In this case, it is yet unclear how to globally identify the dominant resource as well as the dominant share for different users, as each one may have access only to a subset of servers. Work in [11] presents an extension to DRF identifying the user share by ignoring placement constraints and applying a similar approach as in an unconstrained setting. We show that this approach may not achieve fairness in the specific case that one of the resources serves as a bottleneck (see Section 2.2).

In [12] we proposed a multi-resource *fair* allocation mechanism, called PS-DSF, which is applicable to heterogeneous servers in the presence of placement constraints. The intuition behind PS-DSF is to capture the impact of server heterogeneity by measuring the total allocated resources to each user explicitly from the perspective of each server. Specifically, PS-DSF identifies a virtual dominant share (VDS) for each user *with respect to each server* (as opposed to a single system-wide dominant share in DRF). The VDS for user  $n$  with respect to server  $i$  is defined as the ratio of  $x_n$  - the total number of tasks allocated to user  $n$  - over the number of tasks executable by user  $n$  when monopolizing server  $i$ . Then the resources at each server are allocated by applying max-min fairness on VDS (see Section 2.3 for a detailed discussion). This approach is amenable to a distributed implementation. It results in an enhanced performance over the existing mechanisms, and satisfies certain properties essential for fair allocation of resources [12].

## 1.1 Contributions

In this paper, we build upon and generalize our proposed PS-DSF allocation mechanism [12] to *capture the trade-off between efficiency and fairness*. We concisely summarize our contributions.

- We propose a new *server-based formulation* (which includes PS-DSF as a special case) to allocate resources while capturing server heterogeneity. The new formulation can be viewed as a *concave game* among different servers, where each server allocates resources by maximizing a *per-server* utility function (Section 3.1).
- We study a specific class of utility functions which results in an extension of  $\alpha$ -proportional fairness on VDS. We show how the resulting allocation, which we call  $\alpha$ PF-VDS, captures the trade-off between efficiency and fairness by adjusting the parameter  $\alpha$ . We show that  $\alpha$ PF-VDS satisfies bottleneck fairness, envy-freeness and sharing incentive properties (as defined in Section 2.1) for  $\alpha \geq 1$ , and Pareto optimality for  $\alpha = 1$  (Section 3.2 and 3.3).

- We develop a (centralized) convergent algorithm to implement our proposed mechanism. Towards this, we introduce an equivalent formulation for which we derive an iterative solution (Section 4 and 5.1).
- We propose a simple heuristic to develop a distributed implementation for our resource allocation mechanism (Section 5.2).
- We carry-out extensive simulations, driven by real-world traces, to show the enhanced performance of our proposed mechanism (Section 6).

## 1.2 Related Work

### Resource allocation with a game-theoretic approach.

There are several works in the literature which study the resource allocation problem in a cloud computing environment with a game-theoretic approach [13], [14], [15], [16], [17], [18]. Among these, [13], [14], [15] are limited to a *single-resource* setting, while [16], [17], [18] consider a *multi-resource* environment. In these studies the multi-resource allocation problem is formulated as a *game*, where players are different servers. Since they choose DRF as the underlying notion of fairness, they will have the same limitations as DRF for heterogeneous servers (see Section 2.2 for a discussion of such limitations). Moreover, some of them need to solve an extensive form game with a huge strategy set space, e.g., [16], which may not be implementable in a distributed fashion.

### Scheduling in the presence of placement constraints.

There are some recent works investigating max-min fair allocation/scheduling for one type of resource while respecting placement constraints [19], [20], [21], [22], [23]. These single-resource schedulers could be useful in a multi-resource setting when one of the resources is dominantly requested by all users. Otherwise, they might result in a poor resource utilization [1], [19].

## 2 MODEL AND BACKGROUND WORK

Consider a set  $\mathcal{K}$  of  $K$  heterogeneous servers/resource-pools<sup>1</sup> each containing  $M$  types of resources. We denote by  $c_{i,r} \geq 0$ , the capacity (i.e., amount) of resource  $r$  ( $1, 2, \dots, M$ ) on server  $i$ . We make the reasonable assumption that all resources on each server are arbitrarily divisible among the users running on it. Let  $\mathcal{N}$  denote the set of  $N$  active users. Let  $\phi_n > 0$  denote the weight associated with user  $n$ . The weights reflect the priority of users with respect to each other. Let  $\mathbf{d}_n = [d_{n,r}]$  denote the per task *demand vector* for user  $n \in \mathcal{N}$ , i.e., the amount of each resource required for executing one task for user  $n$ . Let  $x_{n,i} \in \mathbb{R}^+$  denote the number of tasks that are allocated to user  $n$  from server  $i$ . Assuming linearly proportionate resource-needs<sup>2</sup>,  $x_{n,i} \mathbf{d}_n = [x_{n,i} d_{n,r}]$  gives

1. A resource-pool may consist of a group of *homogeneous servers*.

2. The assumption of linearly proportionate resource needs is admittedly an idealization. We are following convention set by DRF and used by follow up works.

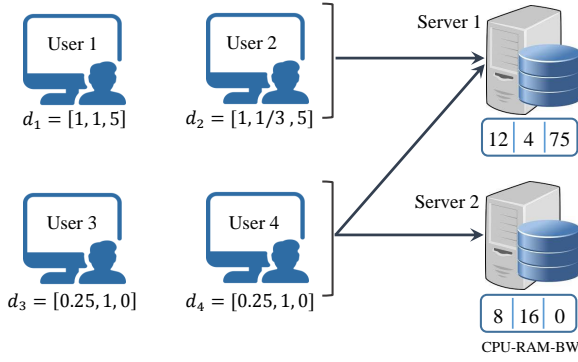


Fig. 1: A heterogeneous multi-resource system with two servers and four equally weighted users.

the amounts of different resources demanded by user  $n$  from server  $i$ .

Due to heterogeneity of users and servers, each user may be restricted to get service *only from a subset of servers*. For example, users may not run tasks on servers which lack some required resources. Furthermore, each user may have some special hardware/software requirements (e.g., public IP address, a particular kernel version, GPU, etc.) which further restrict the set of servers that the user's tasks may run on. Let  $\mathcal{N}_i \neq \emptyset$  denote the set of *eligible users* for server  $i$ . The placement constraints imply that  $x_{n,i} = 0$ ,  $n \notin \mathcal{N}_i$ ,  $\forall i$ .

For instance, consider the example in Fig. 1, where three types of resources, CPU, RAM, and network bandwidth are available over two servers in the amounts of  $\mathbf{c}_1 = [12 \text{ cores}, 4\text{GB}, 75\text{Mb/s}]$  and  $\mathbf{c}_2 = [8 \text{ cores}, 16\text{GB}, 0\text{Mb/s}]$ , where no communication bandwidth is available over the second server; four users with their corresponding demand vectors are also shown in the figure. In this example, the first two users require network bandwidth for execution of their tasks, so they are not *eligible* to run tasks on the second server. However, the last two users may run tasks on both servers.

## 2.1 Dominant resource fairness

Multi-resource fair allocation was originally studied in [1] under the assumption that all resources are aggregated at one resource-pool. Specifically, let  $c_r$  denote the total capacity of resource  $r$ . Let  $\mathbf{a}_n = [a_{n,r}]$  denote the amounts of different resources allocated to user  $n$  under some allocation mechanism. The utilization of user  $n$  of its allocated resources,  $U_n(\mathbf{a}_n)$ , is defined as the number of tasks,  $x_n$ , which could be executed using  $\mathbf{a}_n$ , that is:

$$U_n(\mathbf{a}_n) \triangleq x_n = \min_r \frac{a_{n,r}}{d_{n,r}}. \quad (1)$$

In [1] the following properties are deemed desirable for a multi-resource allocation mechanism.

- *Sharing incentive*: Each user is able to run more tasks compared to a *uniform allocation* where each user  $n$  is allocated a  $\phi_n / \sum_m \phi_m$  fraction of each resource.

- *Envy freeness*: A user should not prefer the allocation vector of another user when adjusted according to their weights, i.e., it should hold that  $U_n(\mathbf{a}_n) \geq U_n(\frac{\phi_n}{\phi_m} \mathbf{a}_m)$  for all  $n, m$ .
- *Bottleneck fairness*: If there is one resource which is *dominantly requested by every user*, then the allocation satisfies *max-min fairness* for that resource.
- *Pareto optimality*: It should not be possible to increase the number of tasks  $x_n$  for any user  $n$ , without decreasing  $x_m$  for some other user(s).
- *Strategy proofness*: Users should not be able to increase their utilization by erroneously declaring their resource demands.

The reader is referred to [1] or [24] for further details. Sharing incentive provides some sort of performance isolation, as it guarantees a minimum utilization for each user irrespective of the demands of the other users. Envy freeness embodies the notion of fairness. Bottleneck fairness describes a necessary condition which applies to a specific case that one resource is dominantly requested by every user, so that a *single-resource* notion of fairness is applicable. These three properties are essential to achieve fairness. So, we refer to them as *essential fairness-related* properties. Pareto optimality is a benchmark for maximizing system utilization. Finally, strategy proofness prevents users from gaming the allocation mechanism. In our view these properties are applicable mainly for private settings. In public settings, users pay explicit costs for their usage or allocations and the provider's goal is to maximize its profits subject to allocation guarantees for users. Even for private clouds, strategy proofness would only be necessary in settings where users act selfishly. In many private settings, users are cooperative and here strategy proofness is not needed. In view of this, we will not consider strategy proofness.

DRF is the first multi-resource allocation mechanism satisfying all the above properties. Specifically, for every user  $n$ , the *Dominant Resource* (DR) is defined as [1]:

$$\rho(n) := \arg \max_r d_{n,r} / c_r, \quad (2)$$

that is, the resource whose greatest portion is required for execution of one task for user  $n$ . The fraction of the DR that is allocated to user  $n$  is defined as its *dominant share*:

$$s_n := \frac{a_{n,\rho(n)}}{c_{\rho(n)}}. \quad (3)$$

Without loss of generality, we may restrict ourselves to non-wasteful allocations, i.e.,  $\mathbf{a}_n = x_n \mathbf{d}_n$ ,  $\forall n$ . Hence, an allocation  $\{x_n\}$  is feasible when:

$$\sum_n x_n d_{n,r} \leq c_r, \quad \forall r. \quad (4)$$

**Definition 1.** An allocation  $\{x_n\}$  satisfies DRF, if it is feasible and the weighted dominant share for each user,  $s_n / \phi_n$  cannot be increased while maintaining feasibility without decreasing  $s_m$  for some user  $m$  with  $s_m / \phi_m \leq s_n / \phi_n$  [1].

DRF is a restatement of *max-min fairness* in terms of *dominant shares*. What make it appealing are the desirable properties which are satisfied under this allocation mechanism.

## 2.2 Existing challenges with heterogeneous servers and placement constraints

In case of heterogeneous servers (whether there are any placement constraints or not), a natural approach to extend DRF is to identify a system-wide dominant resource for each user, *as if* all resources were concatenated within a *single virtual server*. Specifically, let  $c_r := \sum_i c_{i,r}$  denote the total capacity of resource  $r$  within such a virtual server. Then, one may identify the dominant resource for each user  $n$  according to (2). Furthermore, the *global dominant share* for user  $n$  is given by:

$$s_n = x_n \max_r \frac{d_{n,r}}{c_r}, \quad (5)$$

where  $x_n$  is the total number of tasks that are allocated to user  $n$  from different servers, that is  $x_n := \sum_i x_{n,i}$ . As in Definition 1, one may find an allocation  $\{x_{n,i}\}$  which satisfies max-min fairness in terms of the global dominant shares [7]. Such an allocation mechanism, referred to as DRFH, is shown to achieve *Pareto optimality* and *envy freeness*. However, it *fails* to provide *sharing incentive* [7]. We believe that the definition of bottleneck fairness employed by DRFH (with respect to a single virtual server that aggregates all resources) is also controversial. Specifically, if all users have the same dominant resource (with respect to the above mentioned virtual server), then DRFH satisfies max-min fairness with respect to such a resource [7]. In case of heterogeneous servers with placement constraints, however, one may consider other conditions under which a resource serves as a *bottleneck*.

**Definition 2.** A resource  $\rho$  is said to be a bottleneck if for every server  $i$ :

$$\frac{d_{n,\rho}}{c_{i,\rho}} \geq \frac{d_{n,r}}{c_{i,r}}, \quad \forall r, n \in \mathcal{N}_i. \quad (6)$$

If there exists a bottleneck resource, then the allocation should satisfy max-min fairness with respect to that resource.

Unfortunately, DRFH does not satisfy bottleneck fairness in the sense of Definition 2. To appreciate this shortcoming of the DRFH mechanism, consider the example in Fig. 1, where the second resource (RAM) is dominantly requested by eligible users at each server. According to Definition 2, RAM is identified as the bottleneck resource in this example. To allocate the RAM resources in a fair manner, each user should be allocated  $x_1 = x_{1,1} = 2$ ,  $x_2 = x_{2,1} = 6$ ,  $x_3 = x_{3,2} = 8$  and  $x_4 = x_{4,2} = 8$  tasks, respectively (This allocation results from our proposed PS-DSF allocation mechanism [12]). On the other hand, the DRFH mechanism would instead identify network bandwidth as the dominant resource for the first two users and RAM as the dominant resource for the last two

users. To achieve max-min fairness in terms of dominant shares, the DRFH mechanism allocates  $x_1 = x_2 = 3$  and  $x_3 = x_4 = 8$  tasks to each user. Under such an allocation, the RAM resources are not allocated in a fair manner to the first two users.

Yet another extension of DRF, which applies to heterogeneous servers in the presence of placement constraints, is TSF [11]. As in [11], we let  $\gamma_{n,i}$  denote the number of tasks that user  $n$  may execute when monopolizing server  $i$  (i.e., when  $n$  is the only user). Let  $\gamma_n := \sum_i \gamma_{n,i}$  be defined as the number of tasks executable for user  $n$  when monopolizing all servers as if there were no placement constraints. An allocation is said to satisfy Task Share Fairness (TSF), when  $x_n/\gamma_n$  satisfies max-min fairness [11]. When there is only one server, then  $x_n/\gamma_n$  results in the dominant share for each user  $n$ . In such case, TSF reduces to DRF. In case of heterogeneous servers with placement constraints, TSF is shown to satisfy Pareto optimality, envy freeness and sharing incentive properties [11]. However, we show by example that this mechanism may not satisfy bottleneck fairness (neither in the sense of Definition 2, nor in the conventional sense based on considering a single virtual server introduced above [12]).

For instance, consider again the example in Fig. 1, where the second resource is identified as a bottleneck according to Definition 2. The number of tasks that each user may run in the whole cluster is  $\gamma_1 = 4$ ,  $\gamma_2 = 12$ , and  $\gamma_3 = \gamma_4 = 4 + 16 = 20$  tasks, respectively. Hence, each user is allocated  $x_1 = x_{1,1} = 5/3$ ,  $x_2 = x_{2,1} = 5$ ,  $x_3 = x_{3,1} + x_{3,2} = 8 + 1/3 = 25/3$  and  $x_4 = x_{4,1} + x_{4,2} = 8 + 1/3 = 25/3$  tasks, according to the TSF mechanism, which differs from the fair allocation in this example.

## 2.3 Per-server dominant share fairness (PS-DSF)

In this subsection, we describe PS-DSF which we introduced in [12]. PS-DSF is an extension to DRF which is applicable for heterogeneous servers in the presence of placement constraints. The core idea of this mechanism is to introduce a *virtual dominant share* for every user, with respect to each server. Towards this, we first identify the dominant resource for every user  $n$  with respect to each server  $i$ ,

$$\rho(n, i) := \arg \max_r \frac{d_{n,r}}{c_{i,r}}. \quad (7)$$

Let  $\gamma_{n,i}$  denote the number of tasks which could be executed by user  $n \in \mathcal{N}_i$  when monopolizes server  $i$ ,

$$\gamma_{n,i} := \min_r \frac{c_{i,r}}{d_{n,r}} = \frac{c_{i,\rho(n,i)}}{d_{n,\rho(n,i)}}, \quad n \in \mathcal{N}_i. \quad (8)$$

It is assumed that  $\gamma_{n,i} > 0$  for all  $n \in \mathcal{N}_i$ . We set  $\gamma_{n,i} = 0$  if  $n \notin \mathcal{N}_i$ .

**Definition 3.** The *Virtual Dominant Share (VDS)* for user  $n$  with respect to server  $i$ ,  $s_{n,i}$ , is defined as:

$$s_{n,i} := \frac{x_n}{\gamma_{n,i}} = \frac{x_n d_{n,\rho(n,i)}}{c_{i,\rho(n,i)}}, \quad (9)$$

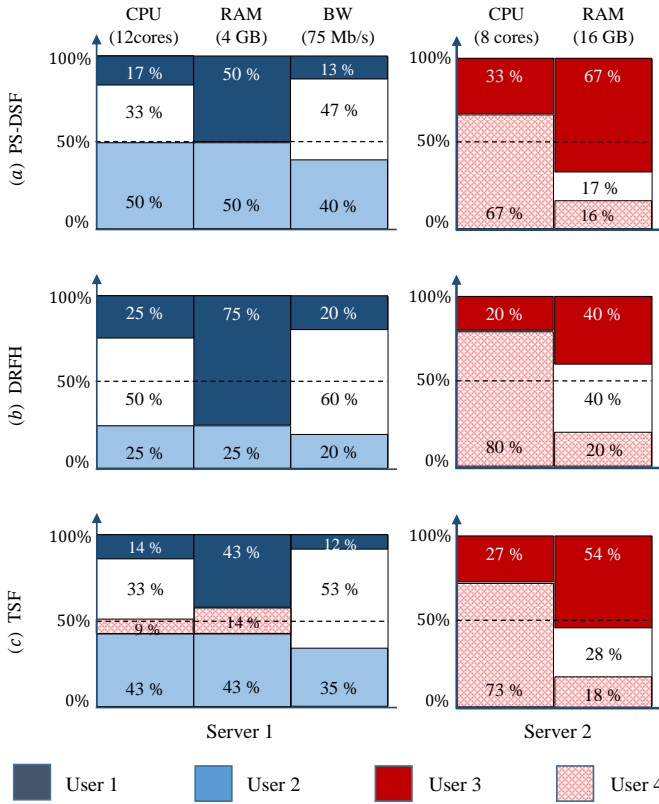


Fig. 2: Comparing the PS-DSF allocation with the DRFH and TSF allocations. The PS-DSF allocation mechanism is more efficient in utilizing different resources.

where  $x_n = \sum_j x_{n,j}$  is the total number of tasks that are allocated to user  $n$  (whether or not these tasks are actually allocated using server  $i$ ).

We have the following conditions on an allocation,  $\mathbf{x} := \{x_{n,i} \in \mathbb{R}^+ \mid n \in \mathcal{N}, i \in \mathcal{K}\}$ , to be feasible:

$$\sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r} \leq c_{i,r}, \quad \forall i, r. \quad (10)$$

$$x_{n,i} = 0, \quad n \notin \mathcal{N}_i, \quad \forall i. \quad (11)$$

**Definition 4.** An allocation  $\mathbf{x}$  satisfies PS-DSF, if it is feasible and the allocated tasks to each user,  $x_n$  cannot be increased (while maintaining feasibility) without decreasing  $x_{m,i}$  for some user  $m$  and server  $i$  with  $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$ .

Intuitively,  $s_{n,i}$  gives the normalized share of the dominant resource for user  $n$  with respect to server  $i$  which should be allocated to it as if  $x_n$  tasks were allocated resources solely from server  $i$  (see the right hand side of (9)). The reader may note that  $s_{n,i}$  could be possibly greater than 1, as some tasks might be allocated to user  $n$  from other servers. According to PS-DSF, the available resources at each server  $i$  are allocated by applying (weighted) max-min fairness on  $\{s_{n,i}\}$ . It can be seen that PS-DSF reduces to DRF when there is only one server.

To gain more intuition, consider again the example in Fig. 1, but this time let  $d_4 = [1, 0.5, 0]$ . In this case, each user may run  $\gamma_{1,1} = 4$ ,  $\gamma_{2,1} = 12$ ,  $\gamma_{3,1} = 4$ ,  $\gamma_{4,1} = 8$  tasks when monopolizing server 1. The third and the fourth users each may run  $\gamma_{3,2} = 16$  and  $\gamma_{4,2} = 8$  tasks when monopolizing server 2. In order to satisfy PS-DSF, each user should be allocated  $x_1 = x_{1,1} = 2$ ,  $x_2 = x_{2,1} = 6$ ,  $x_3 = x_{3,2} = 32/3$  and  $x_4 = x_{4,2} = 16/3$  tasks, respectively. Therefore, the VDS (c.f. Definition 3) for each user with respect to the first server is  $s_{1,1} = s_{2,1} = 0.5$ ,  $s_{3,1} = 8/3$  and  $s_{4,1} = 2/3$ . Also, the VDS for user 3 and 4 with respect to the second server is  $s_{3,2} = s_{4,2} = 2/3$ . The reader can verify that for each server  $i$  the allocated tasks to any user may not be increased without decreasing the allocated tasks to another user with a less or equal VDS. The resulting PS-DSF allocation is shown in Fig. 2. The DRFH and TSF allocations for this example are also illustrated in Fig. 2. It can be seen that the PS-DSF allocation mechanism is more efficient in utilizing different resources compared to the DRFH and TSF mechanisms.

Table 1: Properties of different allocation mechanisms in case of heterogeneous servers with placement constraints: sharing incentive (SI), envy freeness (EF), Pareto optimality (PO), and bottleneck fairness (BF).

Property	DRFH	TSF	PS-DSF
SI		✓	✓
EF	✓	✓	✓
PO	✓	✓	
BF			✓

The reader may note that PS-DSF does not satisfy Pareto optimality in general. It is worth noting that Pareto optimality may not also be satisfied in other works, e.g., [17], [18], which aim at developing a distributed implementation for DRFH. PS-DSF not only is amenable to distributed implementation (as we show in [12]), but also may lead to more efficient utilization of resources compared to the DRFH and TSF mechanisms [12] (as also can be observed in Fig. 2, or in the trace-driven simulations in Section 6). The intuitive reason for this is that each of the DRFH and TSF allocation mechanisms allocates resources based on a *global metric*. Since a global metric throws away information about the actual distribution of resources across servers, approaches based on it may not perfectly capture the impact of server heterogeneity, and therefore may lead to an inefficient resource utilization in heterogeneous settings.

In summary, PS-DSF has been shown to satisfy the essential fairness-related properties, i.e., envy-freeness, sharing incentive and bottleneck fairness, has been observed to offer highly efficient utilization of resources, and is amenable to distributed implementation [12].

### 3 A SERVER-BASED APPROACH FOR MULTI-RESOURCE ALLOCATION

As already discussed, in most of the existing multi-resource allocation mechanisms, fairness is defined in terms of a global metric, a scalar parameter defined for each user in terms of different resources across all servers. Such mechanisms may not succeed in satisfying all the essential fairness-related properties (c.f. Section 2.2), may not readily be implementable in a distributed fashion, and may lead to inefficient resource utilization. In this section, we propose a new formulation for multi-resource allocation problem which is based on a *per-server metric* (as opposed to a global metric) for different users, so that server heterogeneity is captured. The proposed allocation mechanism is built upon our proposed PS-DSF allocation mechanism [12], which was briefly described in the previous section. It generalizes PS-DSF in order to address the trade-off between efficiency and fairness. Furthermore, it inherits all the properties that are satisfied by PS-DSF.

#### 3.1 Problem formulation

As defined in Section 2.3, the VDS is a per-server metric which gives a measure of the allocated resources to each user from the perspective of each server. According to PS-DSF, the available resources at each server are allocated by applying *max-min fairness* on VDS. In order to address the trade-off between efficiency and fairness, we may choose to allocate resources at each server by applying the so-called  *$\alpha$ -proportional fairness* [25] on VDS. To this end, we propose a general formulation, where each server  $i$  strives to maximize a “*per-server utility function*”. Specifically, each server  $i$  tries to find an allocation  $\mathbf{x}_i := [x_{n,i}]$  which solves the following problem<sup>3</sup>.

**Problem 1.** For every server  $i$ :

$$\max_{\mathbf{x}_i} U_i(\mathbf{x}_i, \mathbf{x}_{-i}) := \sum_{n \in \mathcal{N}_i} \phi_n g_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right) \quad (12)$$

$$\text{Subject to: } \sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r} \leq c_{i,r}, \quad \forall r, \quad (13)$$

$$x_{n,i} \geq 0, \quad \forall n \in \mathcal{N}_i, \quad (14)$$

$$x_{n,i} = 0, \quad \forall n \notin \mathcal{N}_i, \quad (15)$$

where  $x_n = \sum_j x_{n,j}$ , and  $g_i(\cdot)$ , as can be also seen in [25], is a scalar function, which is twice-differentiable, strictly concave, and increasing.

**Definition 5.** An allocation  $\mathbf{x}$  is said to be feasible if it satisfies the feasibility conditions in (13)-(15) for all servers.

In Section 3.2, we will present specific choices for  $g_i(\cdot)$ , which capture the *trade-off between efficiency and fairness*,

3. The utility of each server, as defined by (12), depends on its-own allocation/action,  $\mathbf{x}_i$ , as well as actions taken by other servers,  $\mathbf{x}_{-i} := \{\mathbf{x}_j \mid j \neq i\}$ . This is the standard notation used in the context of game theory.

and span a variety of allocations, including the so-called *proportional fair allocation*, and the PS-DSF allocation.

Solving Problem 1 concurrently over different servers is a game, where each server strives to maximize its-own utility. In fact, Problem 1 describes a *concave game* whose players are different servers. It is well-known that a Nash Equilibrium<sup>4</sup> (NE) always exists for such a concave game [26]. Further discussions on the structure of the solution set (Nash equilibriums), and some conditions governing uniqueness of the solution, will be described in Section 4.

#### 3.2 $\alpha$ -proportional fairness on virtual dominant shares

At the optimal solution(s) to Problem 1, *not all* capacity constraints may be active. In fact, there exist trade-offs between efficiency and fairness, which depend on the specific choice of  $g_i(\cdot)$ . To capture the trade-off between efficiency and fairness, one may choose  $g_i(\cdot)$  from the class of  $\alpha$ -fair utility functions [25]. Specifically, we choose  $g_i(z)$  such that  $g'_i(z) = z^{-\alpha}$ , for some fixed parameter  $\alpha$ . For this class of utility functions, the optimal solution to Problem 1 satisfies an extension of  $\alpha$ -proportional fairness in terms of virtual dominant shares, which we call “ *$\alpha$ -Proportional Fairness on VDS*”, or in short  $\alpha$ PF-VDS.

**Definition 6.** A feasible allocation,  $\mathbf{x}$ , satisfies  $\alpha$ PF-VDS, if for every feasible allocation  $\mathbf{y}$ , and for every server  $i$ :

$$\sum_{n \in \mathcal{N}_i} \frac{(y_{n,i} - x_{n,i})/\gamma_{n,i}}{\tilde{s}_{n,i}^\alpha} \leq 0, \quad (16)$$

where  $\tilde{s}_{n,i} := s_{n,i}/\phi_n = x_n/\gamma_{n,i}\phi_n$  is the weighted VDS for user  $n$  with respect to server  $i$ .

**Theorem 1.** Let  $g_i(z)$  be from the class of  $\alpha$ -fair utility functions with  $g'_i(z) = z^{-\alpha}$ ,  $\alpha > 0$ . A feasible allocation,  $\mathbf{x}$ , is a solution to Problem 1 if and only if it satisfies  $\alpha$ PF-VDS.

The proof is given in the appendix. The following theorem, again proven in the appendix, describes how  $\alpha$ PF-VDS is related to other notions of fairness.

**Theorem 2.** The  $\alpha$ PF-VDS allocation is weighted proportionally fair<sup>5</sup> for  $\alpha = 1$ , and approaches a PS-DSF allocation as  $\alpha \rightarrow \infty$ .

Consider again the example in Fig. 1, but let  $d_4 = [1, 0.5, 0]$ . In this example,  $\alpha$ PF-VDS results in the same allocation at server 1, for every  $\alpha > 0$ . In other words, the  $\alpha$ PF-VDS allocation coincides with the PS-DSF allocation at server 1, for every  $\alpha > 0$ . The reason is that RAM is dominantly requested by both of users 1 and 2 which

4. A feasible allocation  $(\mathbf{x}_i^*, \mathbf{x}_{-i}^*)$  is a Nash Equilibrium if no unilateral deviation in action by any single server/player is profitable for that server. That is,  $\forall i, \forall$  (feasible)  $\mathbf{x}_i : U_i(\mathbf{x}_i^*, \mathbf{x}_{-i}^*) \geq U_i(\mathbf{x}_i, \mathbf{x}_{-i}^*)$

5. An allocation,  $\mathbf{x}$ , is weighted proportionally fair if it is feasible, and if for any other feasible allocation,  $\mathbf{y}$ , the weighted summation of proportional changes is not positive, i.e.,  $\sum_n \phi_n \frac{y_n - x_n}{x_n} \leq 0$  (see [27]).

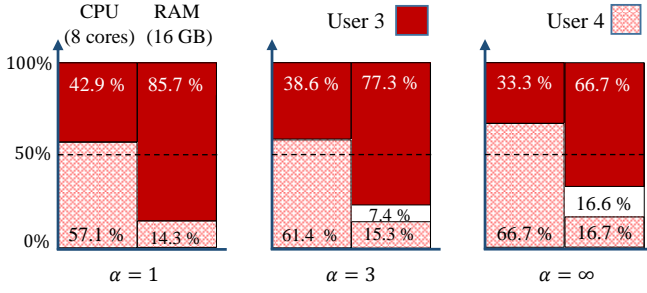


Fig. 3: An illustration of how the  $\alpha$ PF-VDS allocation mechanism can be parameterized (via  $\alpha$ ) to capture the tradeoff between efficiency and fairness.

are allocated resources using server 1 (see Corollary 1). The  $\alpha$ PF-VDS allocation for server 2 is depicted in Fig. 3, for  $\alpha = 1$  (proportional fair allocation),  $\alpha = 3$ , and  $\alpha = \infty$  (PS-DSF allocation). It can be observed that the proportional fair allocation is more efficient in utilizing different resources. However, eligible users for this server tend to get the same portion of their respective dominant resources as  $\alpha \rightarrow \infty$ . This shows how  $\alpha$ -PF-VDS captures the tradeoff between efficiency and fairness.

### 3.3 The properties of the $\alpha$ PF-VDS allocation mechanism

In this section, we investigate different properties which are satisfied by the  $\alpha$ PF-VDS mechanism. In case of heterogeneous servers with placement constraints, we need to extend the notion of sharing incentive property. The notion of bottleneck fairness has been extended by Definition 2. Other properties, Pareto optimality and envy freeness follow the same definitions as described in Section 2. To generalize the sharing incentive property, consider a *uniform allocation*, where a fraction  $\phi_n / \sum_m \phi_m$  of the available resources over each server (whether this server is eligible or not) is allocated to each user  $n$ . An allocation is said to satisfy *sharing incentive*, when each user is able to run more tasks compared to such a uniform allocation.

**Theorem 3.** *The  $\alpha$ PF-VDS allocation mechanism satisfies envy-freeness and sharing incentive properties for every  $\alpha \geq 1$ . It also satisfies Pareto optimality for  $\alpha = 1$ .*

In case that all resources are integrated at one server, [2] has proposed an allocation mechanism which applies  $\alpha$ -proportional Fairness on Dominant Shares (FDS). It can be observed that  $\alpha$ PF-VDS reduces to  $\alpha$ -proportional FDS when there exists only one server. In [2] it is shown that sharing incentive and envy freeness properties are not necessarily satisfied under  $\alpha$ -proportional FDS when  $\alpha < 1$ . Hence,  $\alpha$ PF-VDS may also violate sharing incentive and/or envy freeness properties for  $\alpha < 1$ .

The last property that we consider in this section is bottleneck fairness. According to Definition 2, a resource is considered as a *bottleneck in the whole system* if it is dominantly requested by eligible users at each server.

For the resulting allocation of Problem 1, we may show bottleneck fairness in a more general sense. Specifically, assume that there exists one resource  $\rho(i)$  at each server  $i$  for which the inequality in (6) is satisfied when substituting  $\rho$  with  $\rho(i)$ . We refer to  $\rho(i)$  as the *bottleneck resource at server  $i$* .

**Theorem 4.** *Assume there exists a bottleneck resource at each server. An allocation,  $\mathbf{x}$ , is a solution to Problem 1 if and only if it satisfies PS-DSF.*

The proof appears in the appendix. Under the conditions in Theorem 4, different notions of fairness (including different variants of  $\alpha$ PF-VDS for  $g_i(z) = z^{-\alpha}$ ,  $\alpha > 0$ ) coincide with PS-DSF. Hence, the resulting allocation of Problem 1 satisfies max-min fairness with respect to the bottleneck resource at each server [12]. The following corollary follows directly from the proof of Theorem 4.

**Corollary 1.** *If there exists a bottleneck resource at server  $i$ , then the resulting allocation of Problem 1 satisfies max-min fairness with respect to the bottleneck resource at this server.*

### 3.4 Extensions

**Non-divisible servers:** Problem 1 is formulated based on the assumption that the available resources over each server are arbitrarily divisible. For a data-center comprising of a plurality of servers, it is sometimes of practical interest to assume that servers may not be divided to finer partitions [19], so that each server may only be time-shared by different users. In this case, let  $x_{n,i}$  denote the average number of tasks that are executed by server  $i$  for user  $n$  per unit of time. Accordingly,  $x_{n,i} / \gamma_{n,i}$  gives the percentage of time unit that server  $i$  is allocated to user  $n$ . Hence, we have the following condition on an allocation,  $\mathbf{x}$ , to be feasible [12]:

$$\sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}} \leq 1, \quad \forall i. \quad (17)$$

**Theorem 5.** *Let  $g_i(z)$  be a continuously differentiable, strictly concave and increasing function. An allocation,  $\mathbf{x}$ , satisfies PS-DSF if and only if it is a solution to the following problem.*

**Problem 2.** *For every server  $i$ :*

$$\max_{\mathbf{x}_i} \sum_{n \in \mathcal{N}_i} \phi_n g_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right) \quad (18)$$

$$\text{Subject to: } \sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}} \leq 1, \quad (19)$$

$$x_{n,i} \geq 0, \quad \forall n \in \mathcal{N}_i, \quad (20)$$

$$x_{n,i} = 0, \quad \forall n \notin \mathcal{N}_i. \quad (21)$$

Theorem 5 implies that there exists an *ideal* allocation for non-divisible servers, which results from Problem 2 for any arbitrary  $g_i(z)$ , provided that  $g_i(z)$  is continuously differentiable, strictly concave, and increasing. In fact, different notions of fairness (including different variants of  $\alpha$ PF-VDS, PS-DSF ( $\alpha \rightarrow \infty$ ) and proportional

fairness ( $\alpha = 1$ )) coincide in this case. The major advantage of the PS-DSF allocation mechanism (or  $\alpha$ PF-VDS in general) is giving a simple per-server criterion to find such an ideal allocation [12].

**Servers with heterogeneous objectives:** In Section 3.2, we chose  $g_i(z)$  from the class of  $\alpha$ -fair utility functions with  $g'_i(z) = z^{-\alpha}$ , for some  $\alpha$  that is fixed for all servers. In general, different objectives may be followed for resource allocation at different servers. For instance, in a cloud computing environment, some servers may be owned by a group of users. For such proprietary servers, there might be more emphasis on *fair* allocation of the resources among proprietors (eligible users). On the other hand, there might be some public servers for which there might be more emphasis on *efficient* allocation of the resources.

To address this issue, one may choose  $g_i(z)$  from the class of  $\alpha$ -fair utility functions, but (possibly) with different  $\alpha_i$  for different servers. In this case, we may extend all the results shown in Theorem 1-3. In particular, we may extend the definition of  $\alpha$ PF-VDS, so that the resulting allocation of Problem 1 satisfies the inequality in (16), where we consider different  $\alpha_i$  for different servers. Also, it can be observed that the resulting allocation of Problem 1 satisfies PS-DSF with respect to each server  $i$  as  $\alpha_i \rightarrow \infty$  (see proof of Theorem 2 in the appendix). Furthermore, careful inspection of the proof of Theorem 3 indicates that the resulting allocation of Problem 1 in such a heterogeneous setting satisfies sharing incentive and envy freeness properties, provided that  $\alpha_i \geq 1$ ,  $\forall i$ .

Finally, we introduce a broader class of utility functions which includes the class of  $\alpha$ -fair utility functions. Specifically, consider  $g_i(z)$  such that:

$$g'_i(z) = \left(\frac{A_i}{z}\right)^{\alpha_i} + \frac{B_i}{z}, \quad (22)$$

where  $A_i, B_i \geq 0$ , and  $\alpha_i > 0$ . For this class of utility functions, we can show uniqueness of the solution to Problem 1 (in terms of  $\{x_n\}$ ), when  $B_i > 0$  (see Section 4.4). Again, the trade-off between efficiency and fairness can be adjusted by  $\alpha_i$ . Furthermore, we may extend Theorem 3 to show envy freeness and sharing incentive properties for this class of utility functions.

**Corollary 2.** *The resulting allocation of Problem 1 satisfies sharing incentive and envy-freeness properties, provided that  $g_i(z)$  is from the class of functions in (22), and  $\alpha_i \geq 1$ .*

**Corollary 3.** *Let  $g_i(z)$  be from the class of utility functions specified by (22). The resulting allocation of Problem 1 satisfies PS-DSF with respect to server  $i$  as  $\alpha_i \rightarrow \infty$ , provided that  $\tilde{\alpha}_{n,i}/A_i \leq 1$ ,  $\forall n$ .*

The proofs appear in the appendix.

## 4 TOWARDS A SOLUTION TO PROBLEM 1: AN EQUIVALENT FORMULATION

In this section, we reformulate Problem 1 as a nonlinear complementary problem. This equivalent formulation forms the basis to develop an iterative algorithm to solve this problem in the next section.

### 4.1 Formulation as a non-linear complementary problem

For Problem 1 describing a concave game, it is well known that  $\mathbf{x}$  is a solution (Nash equilibrium) if and only if there exists a set of multipliers,  $\lambda$  and  $\nu$ , such that KKT conditions are satisfied<sup>6</sup> [28]:

$$0 \leq \lambda_{i,r} \perp (c_{i,r} - \sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r}) \geq 0, \quad \forall r, i, \quad (23)$$

$$0 \leq \nu_{n,i} \perp x_{n,i} \geq 0, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (24)$$

$$\frac{\partial \mathcal{L}_i(\mathbf{x}, \lambda, \nu)}{\partial x_{n,i}} = 0, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (25)$$

where,  $\mathcal{L}_i(\mathbf{x}, \lambda, \nu)$  is the Lagrangian function for the local problem at server  $i$ ,

$$\begin{aligned} \mathcal{L}_i(\mathbf{x}, \lambda, \nu) := & \sum_{n \in \mathcal{N}_i} \left[ \phi_n g_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right) + \nu_{n,i} x_{n,i} \right] \\ & + \sum_r \lambda_{i,r} (c_{i,r} - \sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r}). \end{aligned} \quad (26)$$

Hence, the first order optimality condition in (25) implies that:

$$\frac{1}{\gamma_{n,i}} g'_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right) - \sum_r \lambda_{i,r} d_{n,r} + \nu_{n,i} = 0, \quad n \in \mathcal{N}_i, \quad \forall i. \quad (27)$$

We may solve the system of KKT conditions in (23)-(25) for  $\nu_{n,i}$ , and reach the following simplified set of conditions:

$$0 \leq \lambda_{i,r} \perp f_{i,r}(\mathbf{x}_i) \geq 0, \quad \forall r, i, \quad (28)$$

$$0 \leq x_{n,i} \perp f_{n,i}(\mathbf{x}, \lambda_i) \geq 0, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (29)$$

where,

$$f_{i,r}(\mathbf{x}_i) := c_{i,r} - \sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r}, \quad (30)$$

$$f_{n,i}(\mathbf{x}, \lambda_i) := \sum_r \lambda_{i,r} d_{n,r} - \frac{1}{\gamma_{n,i}} g'_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right). \quad (31)$$

The problem of finding  $(\mathbf{x}, \lambda)$ , such that the complementary conditions in (28)-(29) are satisfied, is known as a Non-linear Complementary Problem (NCP). A brief introduction to this family of problems is given in the next subsection.

<sup>6</sup> Conditions of the form  $0 \leq x \perp y \geq 0$  are referred to as *complementary conditions*, where  $x \perp y$  means that  $xy = 0$ .



## 4.2 Background on nonlinear complementary problems

Let  $\mathcal{F}(\mathbf{z})$  be a continuously differentiable function (mapping),  $\mathcal{F} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ . The nonlinear complementary problem,  $\text{NCP}(\mathcal{F})$ , is to find  $\mathbf{z} \in \mathbb{R}^m$  such that:

$$0 \leq \mathbf{z}, \mathcal{F}(\mathbf{z}) \geq 0, \mathbf{z}^T \mathcal{F}(\mathbf{z}) = 0, \quad (32)$$

where the inequalities are taken componentwise. This problem can be best described by introducing a *complementary function*. Specifically,  $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$  is said to be a *complementary function* when [29]:

$$\psi(a, b) = 0 \Leftrightarrow a \geq 0, b \geq 0, \text{ and } ab = 0. \quad (33)$$

Based on a complementary function,  $\text{NCP}(\mathcal{F})$  can be reformulated as:

$$\psi(z_l, F_l(\mathbf{z})) = 0, \forall l. \quad (34)$$

Several complementary functions have been proposed in the literature [30], but the most prominent one is the Fischer-Burmeister function [29]:

$$\psi_{FB}(a, b) := \sqrt{a^2 + b^2} - a - b. \quad (35)$$

Here we use the following complementary function [30]:

$$\psi(a, b) = \frac{1}{2} \psi_{FB}^2(a, b), \quad (36)$$

which is continuously differentiable, and  $\psi(a, b) \geq 0$ ,  $\forall (a, b) \in \mathbb{R}^2$ .

## 4.3 A constrained merit function

In order to solve the NCP described by (28)-(29), we confine the feasible region such that (28) is always satisfied.

**Lemma 1.**  $(\mathbf{x}, \lambda)$  is a solution to the NCP described by (28)-(29), if and only if it is a solution to the following problem.

**Problem 3.**

$$\min \Psi(\mathbf{x}, \lambda) := \sum_{i \in \mathcal{K}} \sum_{n \in \mathcal{N}_i} \psi(x_{n,i}, f_{n,i}(\mathbf{x}, \lambda_i)), \quad (37)$$

$$\text{s.t.} \quad \sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r} \leq c_{i,r}, \forall r, i, \quad (38)$$

$$\lambda_{i,r} \geq 0, r \in \mathcal{R}_i(\mathbf{x}), \forall i, \quad (39)$$

$$\lambda_{i,r} = 0, r \notin \mathcal{R}_i(\mathbf{x}), \forall i. \quad (40)$$

where,  $\mathcal{R}_i(\mathbf{x}) := \{r \mid \sum_{n \in \mathcal{N}_i} x_{n,i} d_{n,r} = c_{i,r}\}$  denotes the set of saturated resources at server  $i$  under the allocation  $\mathbf{x}$ .

*Proof:* The constraints in (38)-(40) are established if and only if  $(\mathbf{x}, \lambda)$  satisfies (28). According to (33), the complementary conditions in (29) are satisfied if and only if  $\psi(x_{n,i}, f_{n,i}(\mathbf{x}, \lambda_i)) = 0$ ,  $n \in \mathcal{N}_i$ ,  $\forall i$ . On the other hand (36) implies that  $\psi(\cdot, \cdot)$  has a lower bound of zero. Hence, a feasible point  $(\mathbf{x}, \lambda)$  is a solution to Problem 3 if and only if (29) is satisfied for all  $n \in \mathcal{N}_i$ ,  $\forall i$ .  $\square$

The following theorem, proven in the appendix, suggests how to find global minima for Problem 3.

**Theorem 6.** A feasible point,  $(\mathbf{x}, \lambda)$ , satisfying the conditions in (38)-(40), is a solution to Problem 3 if and only if it is a stationary point of  $\Psi(\mathbf{x}, \lambda)$  in (37), i.e.,  $\nabla \Psi(\mathbf{x}, \lambda) = 0$ .

## 4.4 Structure of the solution

In this section, we investigate the structure of the solution set for Problem 1, or the equivalent NCP described by Problem 3, when  $g_i(\cdot)$  is specified by (22) and  $\alpha_i \geq 1$ ,  $\forall i$ . In Problem 1, the feasible region for  $\mathbf{x}$ , described by (13)-(15), is a bounded region. As we showed in Section 3.3, the solution to this problem satisfies sharing incentive property. That is,

$$x_n \geq \frac{\psi_n}{\sum_m \phi_m} \sum_i \gamma_{n,i} =: x_n^{\min}, \forall n. \quad (41)$$

For the equivalent NCP described by (28)-(29), it can also be shown that the solution  $(\mathbf{x}, \lambda)$  is contained within a bounded region. To observe boundedness of Lagrange multipliers for each server  $i$ , consider some user  $n$  for which  $x_{n,i} d_{n,r} > 0$ , for some resource  $r$ . From the complementary condition in (29), it follows that  $f_{n,i}(\mathbf{x}, \lambda_i) = 0$ , and therefore,

$$\sum_r \lambda_{i,r} d_{n,r} = \frac{1}{\gamma_{n,i}} g'_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right). \quad (42)$$

Since  $g_i(\cdot)$  is a concave function, the lower bound in (41) implies that:

$$\sum_r \lambda_{i,r} d_{n,r} \leq \frac{1}{\gamma_{n,i}} g'_i \left( \frac{x_n^{\min}}{\phi_n \gamma_{n,i}} \right) < \infty, \quad (43)$$

which clearly shows boundedness of Lagrange multipliers for server  $i$ .

**Lemma 2.** The solution set for the NCP described by (28)-(29), or equivalently Problem 3, is connected.

*Proof:* For Problem 3 we know that: (a) if  $(\mathbf{x}, \lambda)$  is a solution, then  $\nabla \Psi(\mathbf{x}, \lambda) = 0$  (c.f. Theorem 6), (b) the function  $\Psi(\mathbf{x}, \lambda)$ ,  $\Psi : \mathbb{R}^{(N+M)K} \mapsto \mathbb{R}$ , is continuously differentiable, and (c) the solutions are contained within a bounded region. Given these conditions, the proofs for Lemma 3.1, and Corollary 3.5 of [31] can be applied to show connectedness of the solution set for this problem.  $\square$

The following result, borrowed from [31], is a direct consequence of Lemma 2.

**Corollary 4.** Problem 3 has a unique solution if and only if it has a locally unique solution [31].

In the same way, it can be observed that Problem 3 has a unique solution in terms of the total allocated tasks to each user, provided that it has a locally unique solution in terms of  $\{x_n\}$ . In the following we draw conditions under which the optimal solution to Problem 3 will be locally unique in terms of  $\{x_n\}$ .

**Theorem 7.** Let  $(\mathbf{x}, \lambda)$  denote an optimal solution to Problem 3, when  $g_i(\cdot)$  is from the class of functions specified in (22), and  $\alpha_i \geq 1$ ,  $A_i \geq 0$  and  $B_i > 0$ . If  $\lambda_{i,r} > 0$  only for one resource at each server  $i$ , then  $(\mathbf{x}, \lambda)$  will be locally unique in terms of the total allocated tasks to each user,  $\{x_n\}$ .

## 5 ITERATIVE SOLUTION AND DISTRIBUTED IMPLEMENTATION

In this section, we first develop an iterative (centralized) algorithm that is globally convergent to an optimal solution (Nash equilibrium) to Problem 1. Next, we will propose a simple heuristic to solve this problem in a distributed fashion.

### 5.1 Centralized solution

As discussed in Section 4,  $\mathbf{x}$  is a solution to Problem 1, if and only if there exists a set of multipliers such that  $(\mathbf{x}, \lambda)$  is a solution to Problem 3. According to Theorem 6, in order to find a solution to Problem 3, we may employ an iterative descent algorithm which converges to a stationary point where  $\nabla\Psi(\mathbf{x}, \lambda) = 0$ . In the following, we propose an iterative algorithm inspired by projected-gradient method.

Initially, we begin with a feasible point,  $(\mathbf{x}^1, \lambda^1)$ , which satisfies (38)-(40). Then, in each iteration  $h \geq 1$ , we update  $(\mathbf{x}^h, \lambda^h)$  such that  $\Psi(\mathbf{x}^{h+1}, \lambda^{h+1})$  is decreased compared to  $\Psi(\mathbf{x}^h, \lambda^h)$ , while  $(\mathbf{x}^{h+1}, \lambda^{h+1})$  remains feasible. Specifically, let  $\mathcal{R}_i^h$  denote the set of saturated resources at server  $i$  under the allocation  $\mathbf{x}^h$ ,

$$\mathcal{R}_i^h := \{r \mid \sum_{n \in \mathcal{N}_i} x_{n,i}^h d_{n,r} = c_{i,r}\}. \quad (44)$$

The opposite of gradient,  $-\nabla\Psi(\mathbf{x}^h, \lambda^h)$ , is a descent direction. However, by moving  $\mathbf{x}^h$  in the direction of  $-\nabla_{\mathbf{x}}\Psi(\mathbf{x}^h, \lambda^h)$ , the capacity constraint would be violated for resource  $r \in \mathcal{R}_i^h$  if:

$$\beta_{i,r}^h := -(\nabla_{\mathbf{x}_i}\Psi(\mathbf{x}^h, \lambda^h))^T \mathbf{d}_r > 0, \quad (45)$$

where  $\mathbf{d}_r := [d_{n,r}]_{N \times 1}$ , and  $\nabla_{\mathbf{x}_i}\Psi$  is the gradient of  $\Psi$  with respect to  $\mathbf{x}_i$ . So, the moving direction should be chosen in a way that the capacity constraints are not violated for any resource  $r \in \mathcal{R}_i^h$ . Furthermore, to reach a feasible point satisfying (40), we need to preserve equality for the capacity constraints corresponding to resources  $r \in \mathcal{R}_i^h$  with  $\lambda_{i,r}^h > 0$ . Hence, we choose the moving direction for updating  $\mathbf{x}_i^h$ ,  $(\mathbf{v}_x^h)_i$ , as the projection of  $-\nabla_{\mathbf{x}_i}\Psi(\mathbf{x}^h, \lambda^h)$  onto:

$$\Omega_i := \{\mathbf{v} \in \mathbb{R}^N \mid \mathbf{v}^T \mathbf{d}_r \leq 0, r \in \mathcal{R}_i^h : \lambda_{i,r}^h = 0, \text{ and } \mathbf{v}^T \mathbf{d}_r = 0, r \in \mathcal{R}_i^h : \lambda_{i,r}^h > 0\}. \quad (46)$$

To update  $\lambda_{i,r}^h$ , the following is a descent direction<sup>7</sup>:

$$(\tilde{\mathbf{v}}_\lambda^h)_{i,r} := -\frac{\partial\Psi}{\partial\lambda_{i,r}} + \beta_{i,r}^h \mathbf{U}\left(-\beta_{i,r}^h \frac{\partial\Psi}{\partial\lambda_{i,r}}\right). \quad (47)$$

However, to maintain feasibility, we need to choose the moving direction  $(\mathbf{v}_\lambda^h)_{i,r}$  such that:

$$(\mathbf{v}_\lambda^h)_{i,r} := [(\tilde{\mathbf{v}}_\lambda^h)_{i,r}]_{\lambda_{i,r}^h}^+, \quad (48)$$

7. The function  $\mathbf{U}(\cdot)$  is  $\mathbf{U}(z) = 1$  if  $z \geq 0$ , and  $\mathbf{U}(z) = 0$  otherwise.

### Algorithm I: PS-MFA Algorithm

Initially, begin with a feasible point  $(\mathbf{x}^1, \lambda^1)$  satisfying (38)-(40). Then, in each iteration  $h \geq 1$ , take the following steps.

- 1) Choose the moving direction for  $\mathbf{x}_i$ , as the projection of  $-\nabla_{\mathbf{x}_i}\Psi(\mathbf{x}^h, \lambda^h)$  onto  $\Omega_i$ ,

$$(\mathbf{v}_x^h)_i = \text{Proj}_{\Omega_i}(-\nabla_{\mathbf{x}_i}\Psi(\mathbf{x}^h, \lambda^h)), \quad (54)$$

where  $\Omega_i$  is given by (46).

- 2) Choose the moving direction for  $\lambda_{i,r}$  according to (48).
- 3) Choose the step size,  $\eta^h$ , sufficiently small, so that the conditions in (51)-(53) are satisfied.
- 4) Let  $\mathbf{x}^{h+1} = \mathbf{x}^h + \eta^h \mathbf{v}_x^h$ , and  $\lambda^{h+1} = \lambda^h + \eta^h \mathbf{v}_\lambda^h$ .
- 5) Stop when  $\|\mathbf{v}^h\| \leq \epsilon$ , where  $\mathbf{v}^h = [\mathbf{v}_x^h; \mathbf{v}_\lambda^h]$ .

where,  $[v]_z^+ = v$  if  $z > 0$ , and otherwise  $[v]_z^+ = \max\{v, 0\}$ . Finally, we update

$$\mathbf{x}^{h+1} = \mathbf{x}^h + \eta^h \mathbf{v}_x^h, \quad (49)$$

$$\lambda^{h+1} = \lambda^h + \eta^h \mathbf{v}_\lambda^h, \quad (50)$$

where, the step size  $\eta^h$  is chosen such that:

$$\Psi(\mathbf{x}^{h+1}, \lambda^{h+1}) - \Psi(\mathbf{x}^h, \lambda^h) < 0, \quad (51)$$

$$\sum_{n \in \mathcal{N}_i} x_{n,i}^{h+1} d_{n,r} \leq c_{i,r}, \quad \forall r, i, \quad (52)$$

$$\lambda_{i,r}^{h+1} \geq 0, \quad \forall r, i. \quad (53)$$

The algorithm terminates when  $\|\mathbf{v}^h\| < \epsilon$ , for some arbitrary  $\epsilon > 0$ . The above described algorithm, referred to as *Per-Server Multi-resource Fair Allocation* (PS-MFA) algorithm has been summarized in Algorithm I.

**Lemma 3.** *The moving direction,  $\mathbf{v}^h = [\mathbf{v}_x^h; \mathbf{v}_\lambda^h]$ , is a strictly descent direction, unless  $\|\mathbf{v}^h\| = 0$ .*

The proof appears in the appendix. To further analyze the convergence point of the PS-MFA algorithm, we need to make one of the following assumptions.

**Assumption 1.** *The PS-MFA algorithm terminates at a point where  $\lambda_{i,r} > 0$  for every  $r \in \mathcal{R}_i$ ,  $\forall i$ .*

**Assumption 2.** *For every server  $i$  and resource  $r$  there exists at least one user  $n \notin \mathcal{N}_i$  with  $d_{n,r} > 0$ .*

Assumption 1 requires that the algorithm terminates at a *non-degenerate point*, where  $\lambda_{i,r} > 0$  for every  $r \in \mathcal{R}_i$ . Assumption 2 does not restrict the solution, but it gently restricts the model. In particular, it holds when all users demand all types of resources, and there exists one user  $n \notin \mathcal{N}_i$  for each server  $i$ .

**Lemma 4.** *Under either of Assumptions 1 or 2,  $\mathbf{v}^h = 0$  only when  $\nabla_{\mathbf{x}}\Psi = 0$  (the proof is given in the appendix).*

**Theorem 8.** *The PS-MFA algorithm terminates at a stationary point of  $\Psi$  under either of Assumptions 1 or 2.*

*Proof:* The facts that  $\Psi$  is lower bounded and  $\mathbf{v}^h$  is a descent direction imply that the algorithm

converges/terminates. When the algorithm terminates,  $\|\mathbf{v}^h\| \rightarrow 0$ , and therefore  $\nabla_{\mathbf{x}}\Psi = 0$  (see Lemma 4). This in turn implies that  $\psi(x_{n,i}, f_{n,i}(\mathbf{x}, \lambda_i)) = 0$ ,  $\forall n, i$  (see the proof of Theorem 6), and therefore  $\nabla\Psi = 0$ .  $\square$

According to Theorem 6 and 8, the sequence of allocations  $\{\mathbf{x}^h\}$  generated by the PS-MFA algorithm globally converges to an optimal solution to Problem 1 under either of Assumption 1 or 2.

## 5.2 Distributed implementation

The PS-MFA algorithm, as described by Algorithm I, may run in parallel on different servers, where each server iteratively updates its own allocation parameters,  $(\mathbf{x}_i, \lambda_i)$ . However, to find the gradient vector at each server  $i$ , one needs to know the allocation parameters at the other servers. To achieve an efficient distributed implementation, we propose a simple heuristic algorithm which directly applies to the NCP in (28)-(29). First, assume that the Lagrange multipliers,  $\{\lambda_{i,r}\}$ , are known, so that the complementary conditions in (28) are satisfied. Starting with a feasible allocation  $\mathbf{x}^1$ , in each iteration  $h \geq 1$ , one may update  $x_{n,i}^h$  in the following direction:

$$v_{n,i}^h = [-f_{n,i}(\mathbf{x}^h, \lambda)]_{x_{n,i}^h}^+, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (55)$$

$$x_{n,i}^{h+1} = [x_{n,i}^h + \kappa_i v_{n,i}^h]^+, \quad (56)$$

where,  $[v]_z^+ = v$  if  $z > 0$ , and otherwise  $[v]_z^+ = \max\{v, 0\}$ . The step size,  $\kappa_i > 0$ , is chosen independently by each server. According to (55) and (56),  $x_{n,i}^h$  is decreased if it is positive and  $f_{n,i}(\mathbf{x}^h, \lambda_i) > 0$ . Otherwise, it will be increased when  $f_{n,i}(\mathbf{x}^h, \lambda_i) < 0$ . The above dynamic converges/terminates when the complementary conditions in (29) are satisfied. The only issue is to find the Lagrange multipliers. As also can be seen in [27], the Lagrange multipliers could be approximated by:

$$\lambda_{i,r}(\mathbf{x}^h) = \left[ \sum_m x_{m,i}^h d_{m,r} - c_{i,r} + \varepsilon \right]^+ / \varepsilon^2, \quad (57)$$

where a better approximation is achieved when  $\varepsilon \rightarrow 0^8$ . The major advantage of the above dynamic is that the allocation at each server could be updated based on the local information on the available resources at each server, without any knowledge of the available resources at other servers. The only information that each server requires is the *total* number of tasks that are allocated to eligible users, which is assumed to be updated whenever an update is made by any of the servers. The convergence of such a heuristic algorithm is shown through numerical experiments in the next section.

8. This is equivalent to relaxing the capacity constraint in (13) for every resource  $r$ , and adding a quadratic barrier function,  $(\sum_m x_{m,i} d_{m,r} - c_{i,r} + \varepsilon)^2 / 2\varepsilon^2$ , to the objective function in (12).

## 6 TRACE DRIVEN SIMULATION

In this section we evaluate the performance of the  $\alpha$ PF-VDS allocation mechanism through several experiments driven by real-world traces. Among the publicly available traces, Google cluster-usage data-set is the most extensive one which reports the resource usage for different tasks of different users (Google engineers and services) running on a cluster of 12000 servers over a period of one month. The resource usage for each task has been measured at 1 second intervals, however, its average value is reported in the data-set with a period of 5 minutes. If a task is terminated during the five minutes period, its resource usage is reported over a shorter interval. Specifically, the data-set is given in a table format where each row reports the start and the end of each measurement period (which is typically 5 minutes), the job ID and task index, and finally the usage of CPU and RAM for the specified task. The reader may refer to [32] for further details. Despite the detailed information that is provided by this data-set, it does not report the usage of other resources such as network/IO bandwidth.

To do experiments with resource demand vectors of higher dimensionality, we also use a data-set provided by Bitbrain IT services incorporation which gives cloud service to users with business-critical workloads [33]. This data-set reports the resource usage (of CPU, RAM, network and storage bandwidth) for different virtual machines each giving service to one user. Again, the measurements have been reported every 5 minutes for a period of 4 months. There are only a few other traces which are publicly available (such as those from Yahoo or Facebook). However, we do not use such traces herein as they do not provide the information that we need on the usage of different resources. Furthermore, some of them pertain to specific processing tasks, such as Map-reduce, which may not represent the input to a real world data-center [33].

### 6.1 Experimental setup

To do experiments with Google traces, we consider a cluster consisting of four different classes of servers (120 servers in total as shown in Fig. 4), where the configuration of servers are drawn from the distribution of Google cluster servers [32]. For the input workload, we randomly sample 2% of users from the Google traces, so that the cluster is heavily loaded. The jobs for different users belong to 4 different scheduling classes (specified in a table for different jobs), where the last two classes are more latency sensitive. We classify users/jobs<sup>9</sup> into two different groups,  $\mathcal{U}_1$  and  $\mathcal{U}_2$ , where the users in  $\mathcal{U}_1$  are less latency sensitive. Servers in classes A and B are assumed to be public (available to all users), while classes C and D are only available to delay sensitive users,  $\mathcal{U}_2$ .

9. While each user in practice may submit several jobs at the same time, for the sake of simplicity in presentation we assume that each user only submits one job. So, we may use jobs or users interchangeably.

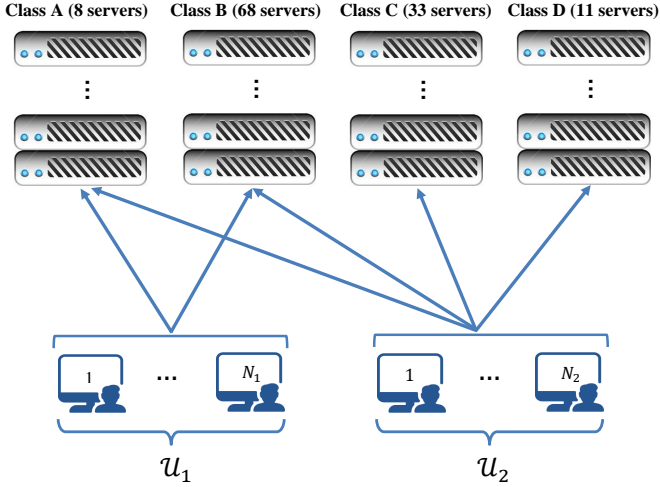


Fig. 4: A cluster with four classes of servers (120 servers in total) and two classes of users. Users are assumed to be equally weighted. The configurations of resources (CPU and memory respectively) for servers of each class are as follows:  $C_A = [1, 1]$ ,  $C_B = [0.5, 0.5]$ ,  $C_C = [0.5, 0.25]$ ,  $C_D = [0.5, 0.75]$ , where CPU and memory units for each server are normalized with respect to the servers of the first class.

To do experiments with the Bitbrain data-set, first we need to choose a set of servers from which the resources are allocated to different users. For this purpose, we find the overall resource usage by active users at different instants of time, so we provide enough capacity of each resource to meet the maximum overall usage. It is worth noting that users demand resources to meet their *maximum* usage. Hence, the overall demands might be more than the overall resource usages at any of instant of time. Given the required resource capacities, we assume that CPU and RAM resources are provided by three types of servers, where there are 75 servers of type 1 with 4 GHz of CPU and 12 GBytes of RAM, 100 servers of type 2 with 8 GHz of CPU and 8 GBytes of RAM, and 75 servers of type 3 with 16 GHz of CPU and 4 GBytes of RAM. It is assumed that servers are distributed over 3 different locations as shown in Fig. 5. Each type of servers at each location is connected to a storage device and also is equipped with a network connection (see Fig. 5). It is assumed that users are uniformly distributed over different locations. Each user may get service from the servers at the same, or nearby location.

## 6.2 Adjusting the resource utilization

As discussed in Section 3.2 and shown by an illustrative example in Fig. 3, the resource utilization is improved as the parameter  $\alpha$  in the  $\alpha$ PF-VDS allocation mechanism reduces. In this subsection we study this effect when applying this mechanism to real-world workloads.

In the Bitbrain workload, users become active/inactive with a low churn. So, the resources could be allocated to different users/virtual-machines in a semi-static manner. In this case, we do several

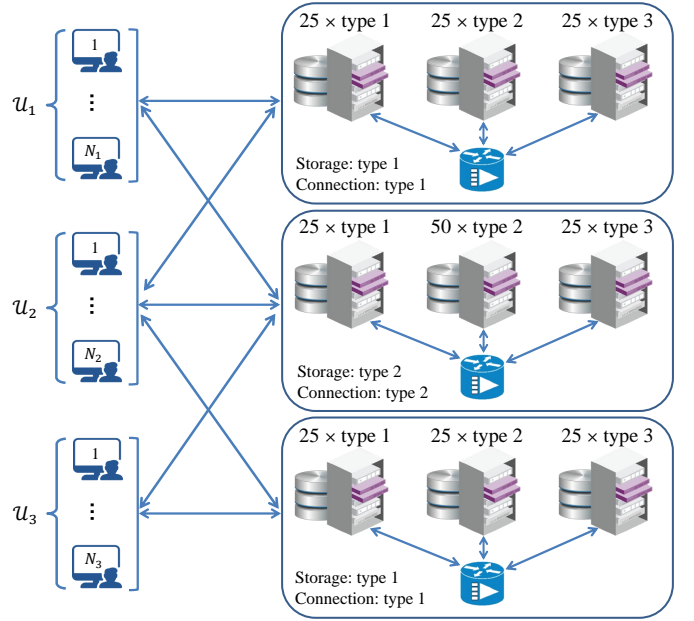


Fig. 5: A data-center distributed over three different locations. There are three types of servers, where the configuration of resources (CPU and memory respectively) for each type of server is as follows: (4 GHz, 12 GBytes) for type 1, (8 GHz, 8 GBytes) for type 2, (16 GHz, 4 GBytes) for type 3. The storage devices are of two different types, where the read/write bandwidth for type 1, used at the first and the last locations, is 32 MB/s, and for type 2, used at the second location, is 100 MB/s. Finally, there are two types of broadband connections, where the first type provides a bandwidth of 100 Mb/s (and 1 Gb/s respectively) to send (receive) data, while the second type provides a bandwidth of 1 Gb/s (and 2 Gb/s respectively) to send (receive) data.

experiments for different sets of active users chosen at random instants of time. In particular, for each set of active users we find the  $\alpha$ PF-VDS allocation for  $\alpha = 1$ ,  $\alpha = 3$  and  $\alpha = \infty$ . In case of  $\alpha = 1$  and  $\alpha = 3$  we employ the distributed iterative algorithm proposed in Section 5.2. For  $\alpha = \infty$ , we employ the customized algorithm proposed in [12] to implement PS-DSF. Fig. 6 shows how our proposed distributed algorithm converges to an optimal solution to Problem 3 where  $\Psi(\mathbf{x}, \lambda) = 0$ . In Fig. 7 we report the overall resource utilization that is achieved on average over different servers and over 100 runs, for different variants of  $\alpha$ PF-VDS. As expected, the  $\alpha$ PF-VDS results in a greater utilization of different resources for smaller values of  $\alpha$ . The improvement in utilization could be significant when  $\alpha$  ranges from  $\infty$  to 1.

For the Google workload, we allocate resources in a (semi)dynamic manner. In particular, consider the computing cluster in Fig. 4, where 2% of users from the Google traces are randomly chosen as the input workload. In such a setting, we (re)allocate resources from the servers to demanding jobs (at least) every 5 minutes. Specifically, given the resource usage for different tasks of each job by the Google traces, we may find the demand vector for each job (at the beginning of each

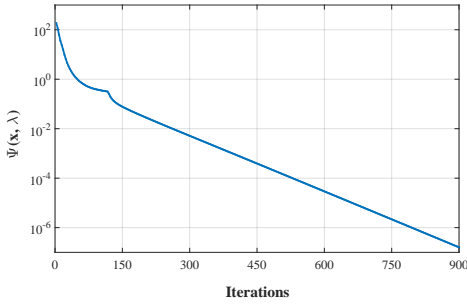


Fig. 6: The convergence of the proposed distributed algorithm to find the  $\alpha$ PF-VDS allocation for the computing cluster in Fig. 5 when  $\alpha = 3$ . The algorithm shows a linear convergence.

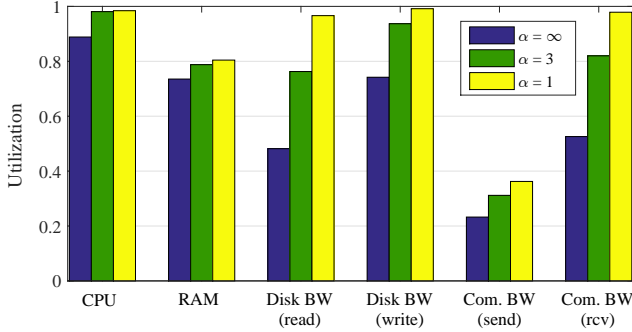


Fig. 7: The overall resource utilization, averaged over different servers and over different runs, for different variants of  $\alpha$ PF-VDS:  $\alpha = 1$  (proportional fairness),  $\alpha = 3$ ,  $\alpha = \infty$  (PS-DSF).

5 minutes interval) as the summation of the resource usage for different tasks (different tasks of the same job usually have proportional demands [32]). Given the total demand for each job,  $\mathbf{d}_n = [d_{n,r}]$ , we define a **quantum** for job  $n$  as a block of resources in the amount of  $\tilde{\mathbf{d}}_n := \mathbf{d}_n / \max_r d_{m,r}$  that is allocated to job  $n$  for 1 second. Accordingly, job  $n$  requires  $q_n := 300 \max_r d_{m,r}$  execution quanta for the next 5 minutes interval. We use the normalized demand vectors,  $\{\tilde{\mathbf{d}}_n\}$ , as the input to the  $\alpha$ PF-VDS mechanism in order to find the number of tasks that are allocated to each job under this mechanism. Given the allocated tasks to each job, the completion time for job  $n$  is given by  $q_n/x_n$ . If a job completes/leaves the system during the 5 minutes period, the released resources are reallocated among the remaining jobs.

The required execution quanta for different jobs, and also their activity duration, span a very wide range. Fig. 8 plots the cumulative density function (CDF) of the activity duration (in terms of the number of 5 min intervals) and also the CDF of the required quanta for different jobs in the Google traces over an interval of 24 hours. As can be seen in Fig. 8, around 38% of jobs are completed within a 5 min period, while 16% of them are active more than 24 hours. Moreover, 53% of jobs require less than 10 execution quanta, while 5% of them require more than 10000 quanta. The variety of jobs and the

existence of such intensive ones indicate the necessity for an efficient and fair allocation mechanism, which from one hand prevents intensive jobs starving others, and on the other hand results in an efficient resource utilization.

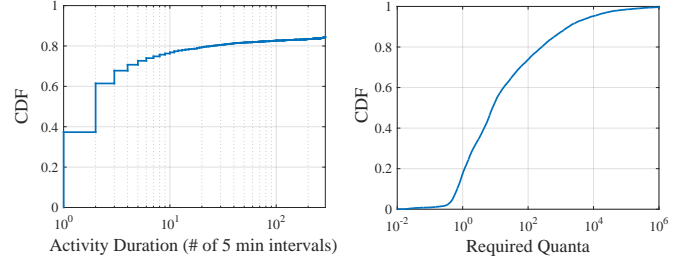


Fig. 8: The CDF of the activity duration and the required quanta for different jobs over an interval of 24 hours.

Fig. 9 compares the overall resource utilization (averaged over different servers) that is achieved on average over an interval of *one hour*, when allocating resources using different variants of  $\alpha$ PF-VDS. Again, it can be observed that the proportional fair allocation ( $\alpha = 1$ ) results in a greater resource utilization, while the PS-DSF allocation ( $\alpha = \infty$ ) is less efficient compared to other variants. It is worth noting that the achieved increase in utilization for smaller values of  $\alpha$  comes at the price of compromising on fairness. Intuitively, PS-DSF strives to balance the (weighted) VDS for all users with  $x_{n,i} > 0$ , at each server  $i$ . Specifically, the PS-DSF allocation results in  $\tilde{s}_{n,i} = \min_m \tilde{s}_{m,i}$ , for all users with  $x_{n,i} > 0$ ,  $\forall i$ , provided that  $d_{n,r} > 0$ ,  $\forall n, r$ . For an arbitrary allocation  $\mathbf{x}$ , we define:

$$D_n(\mathbf{x}) := \sum_i \frac{x_{n,i}}{x_n} \left( \frac{\tilde{s}_{n,i} - \min_m \tilde{s}_{m,i}}{\min_m \tilde{s}_{m,i}} \right), \quad (58)$$

as a measure for *deviation* of each user  $n$  from the fair allocation (given by PS-DSF). In Fig. 10 we report the average deviation,  $\sum_n \phi_n D_n(\mathbf{x}) / \sum_n \phi_n$ , and also the maximum deviation among different users,  $\max_n D_n(\mathbf{x})$ , for our previous experiment in Fig. 9. It can be observed that a larger deviation is experienced for smaller values of  $\alpha$ . This in turn results in variations in the quality of service (e.g. per-quantum delay) experienced by different jobs. In Fig. 11 we report the average per-quantum delay and its standard deviation which are experienced by *short-duration jobs* (those requiring less than one execution quantum) under different variants of the  $\alpha$ PF-VDS mechanism. Again a larger deviation is experienced for smaller values of  $\alpha$ .

### 6.3 Comparison with existing mechanisms

In this subsection, we compare our proposed multi-resource allocation mechanism in terms of resource utilization against the existing multi-resource allocation mechanisms. Specifically, we compare the PS-DSF allocation (which is the least efficient variant of  $\alpha$ PF-VDS), against the DRFH and TSF allocation mechanisms

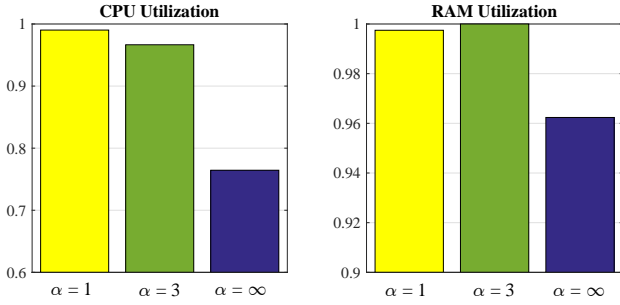


Fig. 9: The overall resource utilization that is achieved (on average) over different servers during a 1 hour period, when allocating resources using  $\alpha$ PF-VDS with:  $\alpha = 1$  (proportional fairness),  $\alpha = 3$ ,  $\alpha = \infty$  (PS-DSF).

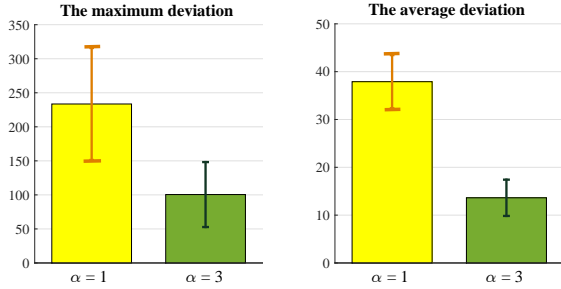


Fig. 10: The average and the maximum deviation (averaged over a 1 hour period), for different variants of the  $\alpha$ PF-VDS allocation mechanism.

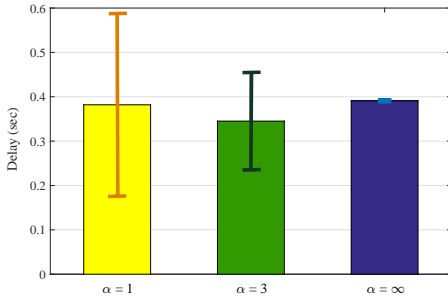


Fig. 11: The average and the standard deviation of the per-quantum delay under different variants of the  $\alpha$ PF-VDS allocation mechanism.

(which all are applicable to heterogeneous servers in the presence of placement constraints). First, consider the computing cluster in Fig. 5 fedded by the Bitbrain workload. We employ each of the aforementioned mechanisms to allocate resources of the servers in Fig. 5 to different sets of active users chosen at random instants of time. The overall utilization that is achieved by of each of these mechanisms for different resources, is shown in Fig. 12, when averaged over different servers and over 100 runs. It can be observed that the PS-DSF allocation mechanism outperforms the two other mechanisms in terms of the achieved utilization for different resources. In particular, the resource utilization is enhanced by the PS-DSF mechanism for up to 20% for some resources.

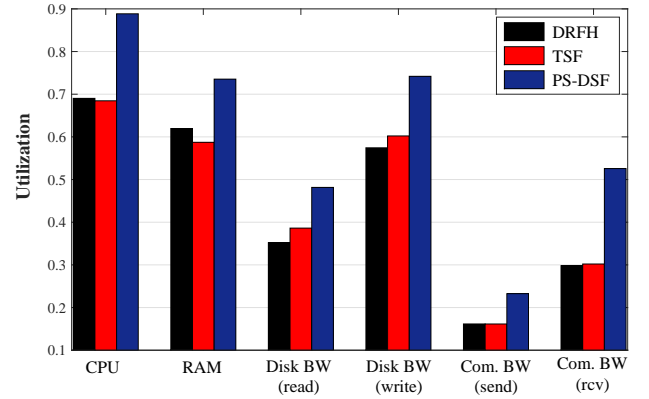


Fig. 12: The overall resource utilization, averaged over different servers and over different runs, for different allocation mechanisms.

We make similar observations with the Google traces. Specifically, consider again the computing cluster in Fig. 4, where 2% of jobs in the Google traces are randomly chosen as the input workload. We employ each of the PS-DSF, DRFH and TSF allocation mechanisms to allocate resources of the servers in Fig. 4 to demanding jobs over an interval of 24 hours. Fig. 13 compares the overall resource utilization (averaged over different servers) that is achieved by different allocation mechanisms. It can be observed that PS-DSF is again more efficient in utilizing different resources, compared to the DRFH and TSF allocation mechanisms, while the achieved resource utilization by DRFH and TSF mechanisms is almost the same<sup>10</sup>. The overall resource utilization that is achieved on average over the 24 hour period is shown in Fig. 14 for different allocation mechanisms. The resource utilization over the last two classes of servers is also shown in Fig. 14. It can be observed that the PS-DSF allocation improves the resource utilization over the last two classes of servers more significantly.

Intuitively, the PS-DSF allocation mechanism allocates resources at each server based on the *per-server virtual dominant shares*. So, at each server it gives more priority to users which may run more tasks (c.f. (9)). Hence, one may expect that the PS-DSF allocation mechanism results in a greater resource utilization compared to the DRFH and TSF mechanisms, especially when the resources are heterogeneously distributed over different servers. That is the reason why a more significant increase in utilization is achieved by the PS-DSF allocation over the last two classes of servers, where the resources are more heterogeneously distributed (the available resources over the first two classes of servers in Fig. 4 are almost proportional to the overall resource capacities). This is also consistent with our observation in the first experiment (with the Bitbrain workload), where the variety of re-

<sup>10</sup> DRFH is just the extension of DRF for multiple heterogeneous servers. There are also other mechanisms, such as those in [18], which approximate DRFH. As their utilization is the same as, or inferior to the DRFH allocation, we do not consider them herein.

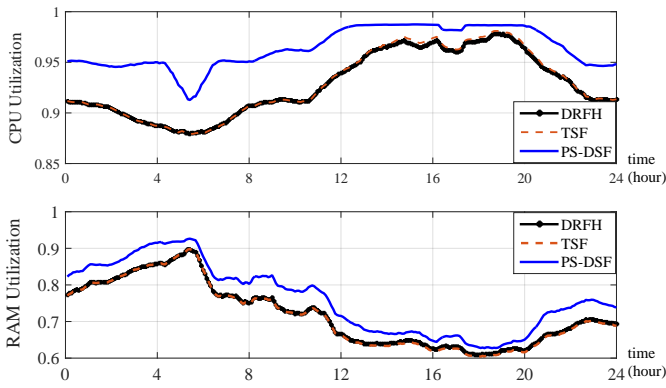


Fig. 13: The overall resource utilization (averaged over different servers) that is achieved by different allocation mechanisms during an interval of 24 hours. To get a better view, a moving average with a window size of 1 hour is applied to all graphs.

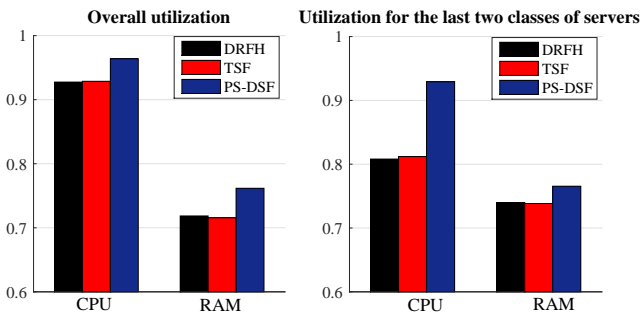


Fig. 14: The overall resource utilization that is achieved on average over an interval of 24 hours.

sources along with the heterogeneity of servers results in a significant outperformance by the PS-DSF mechanism.

## 7 CONCLUSION

We studied efficient and fair allocation of *multiple types* of resources in an environment of *heterogeneous servers* in the presence of *placement constraints*. We identified potential limitations in the existing multi-resource allocation mechanisms, DRF and its follow up work, when used in such environments. In certain occasions, they may not succeed in satisfying all of the essential fairness-related properties, may not readily be implementable in a distributed fashion, and may lead to inefficient resource utilization. We proposed a *new server-based* approach to efficiently allocate resources while capturing server heterogeneity. We showed how our proposed mechanism could be parameterized (by  $\alpha$ ) to adjust the trade-off between efficiency and fairness. Our resource allocation mechanism was shown to satisfy all of the essential fairness-related properties, i.e., sharing incentive, envy-freeness and bottleneck fairness, for every  $\alpha \geq 1$ , and Pareto optimality for the case  $\alpha = 1$ . The amenability to distributed implementation usually comes at the price of degrading the performance. Our proposed mechanism

not only is amenable to distributed implementation, but also results in an enhanced resource utilization compared to the existing mechanisms. We carried out extensive simulations, driven by real-world traces, to demonstrate the performance improvements of our resource allocation mechanism.

## ACKNOWLEDGMENT

The Bitbrain data-set, used for numerical experiments, was graciously provided by Bitbrains IT Services Inc. This data-set is publicly available at <http://gwa.ewi.tudelft.nl>.

## REFERENCES

- [1] A. Ghodsi, M. Zaharia, B. Hindman, A. Konwinski, S. Shenker, and I. Stoica, "Dominant resource fairness: Fair allocation of multiple resource types," in *Proc. NSDI*, June 2011.
- [2] C. Joe-Wong, S. Sen, T. Lan, and M. Chiang, "Multi-resource allocation: Fairness-efficiency tradeoffs in a unifying framework," *IEEE/ACM Trans. Networking*, vol. 21, no. 6, Dec. 2013.
- [3] R. Grandl, G. Ananthanarayanan, S. Kandula, S. Rao, and A. Akella, "Multi-resource packing for cluster schedulers," *SIGCOMM Rev.*, vol. 44, no. 4, pp. 455–466, Aug. 2014.
- [4] T. Bonald and J. Roberts, "Enhanced cluster computing performance through proportional fairness," *Performance Evaluation*, vol. 79, pp. 134–145, 2014.
- [5] D. Bertsekas and R. Gallager, *Data networks*. Prentice Hall, 1992.
- [6] E. Friedman, A. Ghodsi, and C.-A. Psomas, "Strategyproof allocation of discrete jobs on multiple machines," in *Proc. ACM Conf. on Economics and Computation*, 2014.
- [7] W. Wang, B. Liang, and B. Li, "Multi-resource fair allocation in heterogeneous cloud computing systems," *IEEE TPDS*, vol. 26, no. 10, pp. 2822–2835, Oct 2015.
- [8] M. Chowdhury, Z. Liu, A. Ghodsi, and I. Stoica, "Hug: Multi-resource fairness for correlated and elastic demands," in *Proc. NSDI*, Mar 2016.
- [9] Q. Zhu and J. C. Oh, "An approach to dominant resource fairness in distributed environment," in *Proc. IEA-AIE*. Springer, May 2015.
- [10] Y. Tahir, S. Yang, A. Kolioussis, and J. McCann, "Udrf: Multi-resource fairness for complex jobs with placement constraints," in *GLOBECOM*, Dec 2015, pp. 1–7.
- [11] W. Wang, B. Li, B. Liang, and J. Li, "Multi-resource fair sharing for datacenter jobs with placement constraints," *SC 2016*.
- [12] J. Khamse-Ashari, I. Lambadaris, G. Kesidis, B. Urgaonkar, and Y. Zhao, "Per-server dominant-share fairness (ps-dsf): A multi-resource fair allocation mechanism for heterogeneous servers," in *Proc. ICC*, May, 2017.
- [13] J. Bredin, R. T. Maheswaran, C. Imer, T. Başar, D. Kotz, and D. Rus, "A game-theoretic formulation of multi-agent resource allocation," in *Proc. Autonomous Agents*, 2000, pp. 349–356.
- [14] V. Jalaparti and G. D. Nguyen, "Cloud resource allocation games," *Tech. Rep.*, 2010.
- [15] G. Wei, A. V. Vasilakos, Y. Zheng, and N. Xiong, "A game-theoretic method of fair resource allocation for cloud computing services," *The journal of supercomputing*, vol. 54, no. 2, 2010.
- [16] X. Xu and H. Yu, "A game theory approach to fair and efficient resource allocation in cloud computing," *Mathematical Problems in Engineering*, vol. 2014, 2014.
- [17] Q. Zhu and J. C. Oh, "Learning fairness under constraints: A decentralized resource allocation game," in *Proc. IEEE ICMLA*. IEEE, 2016, pp. 214–221.

- [18] —, "Equality or efficiency: A game of distributed multi-type fair resource allocation on computational agents," in *Proc. IEEE/ACM WI-IAT*, vol. 2, 2015, pp. 139–142.
- [19] A. Ghodsi, M. Zaharia, S. Shenker, and I. Stoica, "Choosy: Max-min fair sharing for datacenter jobs with constraints," in *Proc. ACM EuroSys*, 2013, pp. 365–378.
- [20] K. Yap, T. Huang, Y. Yiakoumis, S. Chinchali, N. McKeown, and S. Katti, "Scheduling packets over multiple interfaces while respecting user preferences," in *Proc. ACM coNEXT*, Dec. 2013.
- [21] J. Khamse-Ashari, I. Lambadaris, and Y. Q. Zhao, "Constrained multi-user multi-server max-min fair queuing," <http://arxiv.org/abs/1601.04749>, Jan. 2016.
- [22] J. Khamse-Ashari, G. Kesidis, I. Lambadaris, B. Urgaonkar, and Y. Zhao, "Efficient and fair scheduling of placement constrained threads on heterogeneous multi-processors," in *Proc. DCPeef*, Atlanta, USA, May 2017.
- [23] —, "Constrained max-min fair scheduling of variable-length packet-flows to multiple servers," *Annals of Telecom.*, Aug 2017.
- [24] D. Parkes, A. Procaccia, and N. Shah, "Beyond dominant resource fairness: Extensions, limitations, and indivisibilities," in *Proc. ACM EC*, Valencia, Spain, June 2012.
- [25] J. Mo and J. Walrand, "Fair end-to-end window-based congestion control," *IEEE/ACM Trans. Networking*, vol. 8, no. 5, Oct 2000.
- [26] J. B. Rosen, "Existence and uniqueness of equilibrium points for concave n-person games," *Econometrica: Journal of the Econometric Society*, pp. 520–534, 1965.
- [27] F. P. Kelly, A. K. Maulloo, and D. K. Tan, "Rate control for communication networks: shadow prices, proportional fairness and stability," *Journal of the Operational Research*, 1998.
- [28] A. Dreves, F. Facchinei, C. Kanzow, and S. Sagratella, "On the solution of the kkt conditions of generalized nash equilibrium problems," *SIAM Journal on Opt.*, vol. 21, no. 3, 2011.
- [29] A. Fischer, "New constrained optimization reformulation of complementarity problems," *Journal of Optimization Theory and Applications*, vol. 97, no. 1, pp. 105–117, 1998.
- [30] C. Geiger and C. Kanzow, "On the resolution of monotone complementarity problems," *Computational Optimization and Applications*, vol. 5, no. 2, pp. 155–173, 1996.
- [31] F. Facchinei, "Structural and stability properties of p0 nonlinear complementarity problems," *Math. Oper. Res.*, vol. 23, no. 3, pp. 735–745, Mar. 1998. [Online]. Available: <http://dx.doi.org/10.1287/moor.23.3.735>
- [32] C. Reiss, J. Wilkes, and J. L. Hellerstein, "Google cluster-usage traces," 2011, <http://code.google.com/p/googleclusterdata/>.
- [33] S. Shen, V. van Beek, and A. Iosup, "Statistical characterization of business-critical workloads hosted in cloud datacenters," in *Proc. 15th IEEE/ACM CCGrid*, May 2015, pp. 465–474.
- [34] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.

## APPENDIX

**Proof of Theorem 1:** With  $g'_i(z) = z^{-\alpha}$ , the partial derivative of  $U_i(\mathbf{x})$  in (12) with respect to  $x_{n,i}$

$$\frac{\partial U_i(\mathbf{x})}{\partial x_{n,i}} = \frac{g'_i(\tilde{s}_{n,i})}{\gamma_{n,i}} = \frac{1}{\gamma_{n,i} \tilde{s}_{n,i}^\alpha}, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (59)$$

and  $\partial U_i / \partial x_{n,i} = 0$  if  $n \notin \mathcal{N}_i$ . Hence,

$$\nabla_{\mathbf{x}_i} U_i(\mathbf{x})^T (\mathbf{y}_i - \mathbf{x}_i) = \sum_{n \in \mathcal{N}_i} \frac{(y_{n,i} - x_{n,i}) / \gamma_{n,i}}{\tilde{s}_{n,i}^\alpha}. \quad (60)$$

In case that  $\mathbf{x}$  is a solution to Problem 1, the inequality in (16) follows directly from the first order optimality condition. Specifically,

$$\sum_{n \in \mathcal{N}_i} \frac{(y_{n,i} - x_{n,i}) / \gamma_{n,i}}{\tilde{s}_{n,i}^\alpha} = \nabla_{\mathbf{x}_i} U_i(\mathbf{x})^T (\mathbf{y}_i - \mathbf{x}_i) \leq 0, \quad (61)$$

for every feasible  $\mathbf{y}_i = [y_{n,i}]$ , and for all servers. Now assume that (16) is established for a feasible allocation  $\mathbf{x}$ . The equality in (60) implies that  $\nabla_{\mathbf{x}_i} U_i(\mathbf{x})^T (\mathbf{y}_i - \mathbf{x}_i) \leq 0$ . This along with concavity of  $U_i(\mathbf{x})$  in terms of  $\mathbf{x}_i$  imply that:

$$\begin{aligned} U_i(\mathbf{y}_i, \mathbf{x}_{-i}) &\leq U_i(\mathbf{x}_i, \mathbf{x}_{-i}) + \nabla_{\mathbf{x}_i} U_i(\mathbf{x})^T (\mathbf{y}_i - \mathbf{x}_i) \\ &\leq U_i(\mathbf{x}_i, \mathbf{x}_{-i}), \end{aligned}$$

for every feasible  $\mathbf{y}_i$ , and for all servers. This in turn implies that  $\mathbf{x}$  is a solution (Nash equilibrium) to Problem 1.  $\square$

**Proof of Theorem 2:** For  $\alpha = 1$ , we may sum (16) over different servers. Then, for every feasible allocation  $\mathbf{y}$ :

$$\sum_n \phi_n \frac{y_n - x_n}{x_n} \leq 0, \quad (62)$$

which implies that  $\mathbf{x}$  satisfies (weighted) proportional fairness. In case that  $\alpha \rightarrow \infty$ , the proof would be similar to that for Lemma 5 in [25]. Let  $\mathbf{x}(\alpha)$  denote an allocation satisfying  $\alpha$ PF-VDS. Since the feasible region, described by (13)-(15), is a compact set, we can find a sequence of  $\alpha$ ,  $\{\alpha_l, l \geq 1\}$ , such that  $\lim_{l \rightarrow \infty} \alpha_l = \infty$  and  $\{\mathbf{x}(\alpha_l)\}$  converges to some feasible  $\mathbf{x}(\infty)$  as  $l \rightarrow \infty$ . By definition, if  $\mathbf{x}(\alpha_l)$  satisfies  $\alpha$ PF-VDS, then for every server  $i$  and for all feasible allocations,  $\mathbf{y}_i$ ,

$$\sum_{m \in \mathcal{N}_i} \frac{(y_{m,i} - x_{m,i})}{\gamma_{m,i} [\tilde{s}_{m,i}(\alpha_l)]^{\alpha_l}} \leq 0,$$

where  $\tilde{s}_{m,i}(\alpha_l) := x_m(\alpha_l) / \phi_m \gamma_{m,i}$  is the weighted VDS for user  $m$  with respect to server  $i$ . Consider an arbitrary user  $n$  and server  $i$  for which  $y_{n,i} \neq x_{n,i}$ . It follows that:

$$\Delta_{n,i} := \frac{(y_{n,i} - x_{n,i})}{\gamma_{n,i} [\tilde{s}_{n,i}(\alpha_l)]^{\alpha_l}} \leq - \sum_{m \neq n, m \in \mathcal{N}_i} \frac{(y_{m,i} - x_{m,i})}{\gamma_{m,i} [\tilde{s}_{m,i}(\alpha_l)]^{\alpha_l}},$$

Dividing both sides of the above inequality by  $\Delta_{n,i} \neq 0$ ,

$$\begin{aligned} 1 &\leq \sum_{m \neq n, m \in \mathcal{N}_i} h_{m,i} \left[ \frac{\tilde{s}_{n,i}(\alpha_l)}{\tilde{s}_{m,i}(\alpha_l)} \right]^{\alpha_l} \\ &\leq \sum_{m \neq n, h_{m,i} > 0} h_{m,i} \left[ \frac{\tilde{s}_{n,i}(\alpha_l)}{\tilde{s}_{m,i}(\alpha_l)} \right]^{\alpha_l}, \quad (63) \end{aligned}$$

where  $h_{m,i} := -\frac{\gamma_{n,i}(y_{m,i} - x_{m,i})}{\gamma_{m,i}(y_{n,i} - x_{n,i})} < \infty$ .

Unless there exists some user  $p$  with  $h_{p,i} > 0$  and  $\tilde{s}_{p,i}(\infty) \leq \tilde{s}_{n,i}(\infty)$ , the right hand side of (63) approaches zero as  $\alpha_l \rightarrow \infty$ . Hence, for the inequality in (63) to hold, there must exist some user  $p$  with  $h_{p,i} > 0$  such that  $\tilde{s}_{p,i}(\infty) \leq \tilde{s}_{n,i}(\infty)$ . Note that  $h_{p,i} > 0$  implies



that if  $y_{n,i} > x_{n,i}(\infty)$ , then  $y_{p,i} < x_{p,i}(\infty)$ . In other words, for any feasible allocation  $\mathbf{y} \neq \mathbf{x}(\infty)$ , we may not increase the allocated tasks to user  $n$  from server  $i$ ,  $y_{n,i} > x_{n,i}(\infty)$ , unless decreasing the allocated tasks from server  $i$ ,  $y_{p,i} < x_{p,i}(\infty)$ , for some user  $p$  with  $\tilde{s}_{p,i}(\infty) \leq \tilde{s}_{n,i}(\infty)$ . This implies that  $\mathbf{x}(\infty)$  satisfies PS-DSF.  $\square$

**Proof of Theorem 3:** For  $\alpha = 1$ , the  $\alpha$ PF-VDS allocation reduces to a proportional fair allocation, maximizing a *global objective*,  $\sum_n \phi_n \log(x_n)$ . So, the resulting allocation is Pareto optimal. To prove sharing incentive and envy freeness properties, first we need to derive some inequalities.

As in Section 4, let  $\lambda_{i,r}$ , and  $\nu_{n,i}$ , denote the Lagrange multipliers corresponding to the constraints in (13), and (14), respectively. Given the Lagrange multipliers, the KKT conditions for Problem 1 are described by (23)-(25). Specifically, the first order optimality condition in (25) implies that:

$$\gamma_{n,i}^{-1} g'_i(\tilde{s}_{n,i}) - \sum_r \lambda_{i,r} d_{n,r} + \nu_{n,i} = 0, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (64)$$

where  $\tilde{s}_{n,i} = x_n / \phi_n \gamma_{n,i}$ . Since  $\nu_{n,i} \geq 0$ ,

$$\gamma_{n,i}^{-1} g'_i(\tilde{s}_{n,i}) \leq \sum_r \lambda_{i,r} d_{n,r}, \quad n \in \mathcal{N}_i, \quad \forall i. \quad (65)$$

By definition,  $\gamma_{n,i} = \min_r \{c_{i,r} / d_{n,r}\}$ . Hence,

$$g'_i(\tilde{s}_{n,i}) \leq \sum_r \lambda_{i,r} \gamma_{n,i} d_{n,r} \leq \sum_r \lambda_{i,r} c_{i,r}, \quad n \in \mathcal{N}_i, \quad \forall i. \quad (66)$$

Multiplying both sides of (64) by  $x_{n,i}$ , and summing over different users:

$$\begin{aligned} \sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}} g'_i(\tilde{s}_{n,i}) &= \sum_{n \in \mathcal{N}_i} \left[ \sum_r \lambda_{i,r} x_{n,i} d_{n,r} + x_{n,i} \nu_{n,i} \right] \\ &= \sum_r \sum_{n \in \mathcal{N}_i} \lambda_{i,r} x_{n,i} d_{n,r}, \quad \forall i, \end{aligned} \quad (67)$$

where the second equality follows from the fact that  $x_{n,i} \nu_{n,i} = 0$ ,  $n \in \mathcal{N}_i$ ,  $\forall i$  (c.f. (24)). Moreover, the complementary condition in (23) implies that:

$$\sum_r \sum_{n \in \mathcal{N}_i} \lambda_{i,r} x_{n,i} d_{n,r} = \sum_r \lambda_{i,r} c_{i,r}, \quad \forall i. \quad (68)$$

From (66), (67) and (68) it follows that:

$$g'_i(\tilde{s}_{n,i}) \leq \sum_{m \in \mathcal{N}_i} \frac{x_{m,i}}{\gamma_{m,i}} g'_i(\tilde{s}_{m,i}), \quad n \in \mathcal{N}_i, \quad \forall i. \quad (69)$$

Now, we are ready to prove the sharing incentive property. Let  $n_i^* := \arg \min_n \tilde{s}_{n,i}$ , and multiply both sides of the inequality in (69) for server  $i$  and user  $n_i^*$  by  $(\tilde{s}_{n_i^*,i})^{\alpha-1}$ . For  $g'_i(z) = z^{-\alpha}$  and  $\alpha \geq 1$ , it follows that:

$$\begin{aligned} (\tilde{s}_{n_i^*,i})^{-1} &\leq \sum_{m \in \mathcal{N}_i} \phi_m \frac{x_{m,i}}{x_m} \left[ \frac{\tilde{s}_{n_i^*,i}}{\tilde{s}_{m,i}} \right]^{\alpha-1} \\ &\leq \sum_{m \in \mathcal{N}} \phi_m \frac{x_{m,i}}{x_m}, \quad \forall i, \end{aligned} \quad (70)$$

where the second inequality follows from the facts that  $\tilde{s}_{n_i^*,i} \leq \tilde{s}_{m,i}$ ,  $\forall m$ , and  $x_{m,i} = 0$ ,  $n \notin \mathcal{N}_i$ . Furthermore,

$$(\tilde{s}_{n,i})^{-1} \leq (\tilde{s}_{n_i^*,i})^{-1} \leq \sum_m \phi_m \frac{x_{m,i}}{x_m}, \quad \forall n, i. \quad (71)$$

Summing (71) over different servers,

$$\frac{\phi_n}{x_n} \sum_i \gamma_{n,i} \leq \sum_m \sum_i \phi_m \frac{x_{m,i}}{x_m} = \sum_m \phi_m, \quad (72)$$

or,

$$x_n \geq \frac{\phi_n}{\sum_m \phi_m} \sum_i \gamma_{n,i}. \quad (73)$$

To show envy freeness, (by contradiction) assume that user  $n$  envies user  $m$ 's allocation vector, when adjusted according to their weights. That is,  $U_n(\frac{\phi_n}{\phi_m} \mathbf{a}_m) > U_n(\mathbf{a}_n)$ , where  $\mathbf{a}_m = x_m \mathbf{d}_m$  and  $\mathbf{a}_n = x_n \mathbf{d}_n$ . It follows that (c.f. (1)):

$$\frac{x_n d_{n,r}}{\phi_n} < \frac{x_m d_{m,r}}{\phi_m}, \quad \forall r. \quad (74)$$

Consider some server  $j$  for which  $x_{m,j} > 0$ , so that  $\nu_{m,j} = 0$ . Then, (64) implies that:

$$\sum_r \lambda_{j,r} d_{m,r} = \gamma_{m,j}^{-1} g'_j(\tilde{s}_{m,j}). \quad (75)$$

If we multiply both sides of (65) by  $x_n / \phi_n$ , then it follows from (74) and (75) that:

$$\begin{aligned} \tilde{s}_{n,j} g'_j(\tilde{s}_{n,j}) &\leq \sum_r \frac{\lambda_{j,r} d_{n,r} x_n}{\phi_n} \\ &< \sum_r \frac{\lambda_{j,r} d_{m,r} x_m}{\phi_m} \\ &= \tilde{s}_{m,j} g'_j(\tilde{s}_{m,j}) \end{aligned} \quad (76)$$

In case that  $g'_j(z) = z^{-\alpha}$  and  $\alpha > 1$ , it follows from (76) that  $\tilde{s}_{n,j} > \tilde{s}_{m,j}$ , and therefore,

$$\frac{x_n d_{n,\rho(n,j)}}{\phi_n c_{j,\rho(n,j)}} = \frac{x_n}{\phi_n \gamma_{n,j}} > \frac{x_m}{\phi_m \gamma_{m,j}} \geq \frac{x_m d_{m,r}}{\phi_m c_{j,r}}, \quad \forall r, \quad (77)$$

which contradicts (74) for  $r = \rho(n,j)$ . In case that  $g'_j(z) = z^{-1}$ , (76) requires that  $1 < 1$ , which is a contradiction.  $\square$

**Proof of Theorem 4:** A resource, say  $\rho(i)$ , is identified as a bottleneck resource at server  $i$ , if it is dominantly requested by all eligible users for server  $i$ . That is  $\rho(n,i) = \rho(i)$ ,  $n \in \mathcal{N}_i$ . Hence,

$$\sum_{n \in \mathcal{N}_i} x_{n,i} \frac{d_{n,r}}{c_{i,r}} \leq \sum_{n \in \mathcal{N}_i} x_{n,i} \frac{d_{n,\rho(i)}}{c_{i,\rho(i)}} = \sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}}, \quad \forall i, r. \quad (78)$$

The above inequality implies that all of the capacity constraints for server  $i$  could be replaced by:

$$\sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}} \leq 1. \quad (79)$$

Let  $\lambda_i$  denote the Lagrange multiplier corresponding to the constraint in (79). Hence, the KKT conditions for Problem 1 (given by (23)-(25)) can be simplified/rewritten as:

$$g'_i\left(\frac{x_n}{\phi_n \gamma_{n,i}}\right) - \lambda_i + \nu'_{n,i} = 0, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (80)$$

$$0 \leq \lambda_i \perp \left(1 - \sum_{n \in \mathcal{N}_i} x_{n,i} / \gamma_{n,i}\right) \geq 0, \quad \forall i, \quad (81)$$

$$0 \leq \nu'_{n,i} \perp x_{n,i} \geq 0, \quad n \in \mathcal{N}_i, \quad \forall i, \quad (82)$$

where,  $\nu'_{n,i} := \nu_{n,i} \gamma_{n,i}$ .

First assume that  $\mathbf{x}$  is a solution (Nash equilibrium) for Problem 1, so that the KKT conditions in (80)-(82) are satisfied. We know that  $g'(\cdot) > 0$ , owing to the assumption that  $g(\cdot)$  is strictly increasing and differentiable. This along with (80) imply that  $\lambda_i > \nu'_{n,i}$ ,  $n \in \mathcal{N}_i \forall i$ , and therefore,  $\lambda_i > 0, \forall i$ . According to (81),

$$\sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}} = 1, \quad \forall i. \quad (83)$$

The assumption that  $g(\cdot)$  is *strictly concave* implies that  $g'(\cdot)$  is strictly decreasing, and therefore is invertible. The inverse of  $g'(\cdot)$  would be also a strictly decreasing function which we denote it by  $h(\cdot)$ . It follows from (80) that:

$$\tilde{s}_{n,i} := \frac{x_n}{\phi_n \gamma_{n,i}} = h(\lambda_i, -\nu'_{n,i}), \quad n \in \mathcal{N}_i \quad \forall i. \quad (84)$$

Consider two users  $n, m \in \mathcal{N}_i$ , where  $x_{n,i} > 0$ , so that  $\nu'_{n,i} = 0$  (c.f. (82)). It follows that,

$$\tilde{s}_{n,i} = h(\lambda_i) \leq h(\lambda_i - \nu'_{m,i}) = \tilde{s}_{m,i}, \quad (85)$$

where the above inequality follows from the fact that  $h(\cdot)$  is a decreasing function. According to (83) and (85), we may not increase  $x_{m,i}$  for any user  $m$ , unless decreasing  $x_{n,i}$  for some user  $n$  with  $\tilde{s}_{n,i} \leq \tilde{s}_{m,i}$ . So, by definition,  $\mathbf{x}$  satisfies PS-DSF.

Now consider an allocation  $\mathbf{x}$  satisfying PS-DSF. For such an allocation, there exists (at least) one resource at each server  $i$  for which the capacity constraint holds with equality. This along with (78) imply that the constraint in (79) holds with equality at each server  $i$ . Hence, (81) is established for all servers. For each server  $i$ , consider some user  $n$  for which  $x_{n,i} > 0$ . By definition,  $\tilde{s}_{p,i} \geq \tilde{s}_{n,i}$ ,  $\forall p$ , and  $\tilde{s}_{p,i} = \tilde{s}_{n,i}$ ,  $\forall p : x_{p,i} > 0$ . Hence, we may choose  $\lambda_i := g'(\tilde{s}_{n,i})$ , and

$$\nu'_{m,i} = \lambda_i - g'_i(\tilde{s}_{m,i}) \geq 0, \quad m \in \mathcal{N}_i, \quad \forall i, \quad (86)$$

so that the first order optimality condition in (80) is established. Furthermore, (86) implies that  $\nu'_{m,i} = 0$  when  $x_{m,i} > 0$ . So, all of the conditions in (80)-(82) are established in conjunction with the chosen multipliers. In other words, given a PS-DSF allocation,  $\mathbf{x}$ , we may find a set of multipliers such that the KKT conditions in (80)-(82) are established. This implies that  $\mathbf{x}$  is a solution to Problem 1.  $\square$

**Proof of Theorem 5:** The proof follows exactly the same line of arguments as Theorem 4. Specifically, let  $\lambda_i$ , and  $\nu_{n,i}$ , denote the Lagrange multipliers corresponding to the constraints in (19), and (20). The Lagrangian function for Problem 2 at server  $i$  is given by:

$$\begin{aligned} \mathcal{L}_i(\mathbf{x}, \lambda, \nu) := & \sum_{n \in \mathcal{N}_i} \left[ \phi_n g_i\left(\frac{x_n}{\phi_n \gamma_{n,i}}\right) + \sum_n \nu_{n,i} x_{n,i} \right] \\ & + \lambda_i \left[ 1 - \sum_{n \in \mathcal{N}_i} \frac{x_{n,i}}{\gamma_{n,i}} \right]. \end{aligned} \quad (87)$$

It can be observed that the KKT conditions for this problem are exactly the same as those given by (80)-(82) in Theorem 4. Hence, the same arguments apply here.  $\square$

**Proof of Corollary 2:** As in Theorem 3, let  $n_i^* := \arg \min_n \tilde{s}_{n,i}$ . Then, divide both sides of the inequality in (69) for server  $i$  and user  $n_i^*$  by  $g'_i(\tilde{s}_{n_i^*,i}) \tilde{s}_{n_i^*,i}$ . It follows that:

$$(\tilde{s}_{n_i^*,i})^{-1} \leq \sum_{m \in \mathcal{N}_i} \phi_m \frac{x_{m,i}}{x_m} \left[ \frac{\tilde{s}_{m,i}}{\tilde{s}_{n_i^*,i}} \frac{g'_i(\tilde{s}_{m,i})}{g'_i(\tilde{s}_{n_i^*,i})} \right]. \quad (88)$$

To simplify the right hand side of (88), let define:

$$q_i(z) := \frac{1}{z g'_i(z)} = \frac{\left(\frac{z}{A_i}\right)^{\alpha_i - 1}}{A_i + B_i \left(\frac{z}{A_i}\right)^{\alpha_i - 1}}. \quad (89)$$

It can be observed that  $q_i(z)$  is an increasing function. Hence,

$$\begin{aligned} (\tilde{s}_{n,i})^{-1} \leq (\tilde{s}_{n_i^*,i})^{-1} & \leq \sum_{m \in \mathcal{N}_i} \phi_m \frac{x_{m,i}}{x_m} \left[ \frac{q_i(\tilde{s}_{n_i^*,i})}{q_i(\tilde{s}_{m,i})} \right] \\ & \leq \sum_{m \in \mathcal{N}_i} \phi_m \frac{x_{m,i}}{x_m}, \quad \forall n, i, \end{aligned} \quad (90)$$

where the first and the last inequalities follow from the fact that  $\tilde{s}_{n_i^*,i} \leq \tilde{s}_{m,i}$ ,  $\forall m$ . As in Theorem 3, we may sum (90) over different servers to get the required result (c.f. (72) and (73)).  $\square$

**Proof of Corollary 3:** Consider a sequence,  $\{\alpha_l, l \geq 1\}$ , which converges to infinity,  $\lim_{l \rightarrow \infty} \alpha_l = \infty$ , and  $\{\mathbf{x}(\alpha_l)\}$  converges to some feasible  $\mathbf{x}(\infty)$  as  $l \rightarrow \infty$ . When  $\mathbf{x}(\alpha_l)$  is a solution to Problem 1, then for every feasible allocation,  $\mathbf{y}$ , it can be shown that (c.f. Theorem 1):

$$\sum_m g'_i(\tilde{s}_{m,i}(\alpha_l)) \frac{y_{m,i} - x_{m,i}}{\gamma_{m,i}} \leq 0, \quad \forall i.$$

Now, consider an arbitrary user  $n$  and server  $i$  for which  $y_{n,i} \neq x_{n,i}$ . As in Theorem 2, it can be shown that:

$$1 \leq \sum_{m \neq n, h_{m,i} > 0} h_{m,i} \frac{g'_i(\tilde{s}_{m,i}(\alpha_l))}{g'_i(\tilde{s}_{n,i}(\alpha_l))}, \quad (91)$$

where  $h_{m,i} := -\frac{\gamma_{n,i}(y_{m,i} - x_{m,i})}{\gamma_{m,i}(y_{n,i} - x_{n,i})} < \infty$ . On the other hand,

$$\lim_{l \rightarrow \infty} \frac{g'_i(\tilde{s}_{m,i}(\alpha_l))}{g'_i(\tilde{s}_{n,i}(\alpha_l))} = \lim_{l \rightarrow \infty} \left[ \frac{\tilde{s}_{n,i}(\alpha_l)}{\tilde{s}_{m,i}(\alpha_l)} \right]^{\alpha_l}, \quad (92)$$

provided that  $\tilde{s}_{n,i}/A_i \leq 1$ ,  $n \in \mathcal{N}_i$ . The limit in (92), and therefore the right hand side of (91), approaches zero as  $\alpha_l \rightarrow \infty$ , unless there exists some user  $p$  with  $h_{p,i} > 0$  and  $\tilde{s}_{p,i}(\infty) \leq \tilde{s}_{n,i}(\infty)$ . In other words, for the inequality in (91) to be established, there must exist some user  $p$  with  $h_{p,i} > 0$  such that  $\tilde{s}_{p,i}(\infty) \leq \tilde{s}_{n,i}(\infty)$ . As in Theorem 2, this implies that  $\mathbf{x}(\infty)$  satisfies PS-DSF.  $\square$

**Proof of Theorem 6:** Let  $\nabla_{\mathbf{x}}\Psi = [(\nabla_{\mathbf{x}}\Psi)_{n,i}]_{NK \times 1}$  denote the gradient of  $\Psi$  with respect to  $\mathbf{x}$ , and  $\nabla_{\lambda}\Psi = [(\nabla_{\lambda}\Psi)_{i,r}]_{MK \times 1}$  denote the gradient of  $\Psi$  with respect to  $\lambda$ . In particular,

$$(\nabla_{\mathbf{x}}\Psi)_{n,i} := \frac{\partial\Psi}{\partial x_{n,i}} = \left[ \frac{\partial\psi_{n,i}}{\partial x_{n,i}} + \sum_j \frac{\partial\psi_{n,j}}{\partial f_{n,j}} \frac{\partial f_{n,j}}{\partial x_n} \right], \quad (93)$$

$$(\nabla_{\lambda}\Psi)_{i,r} := \frac{\partial\Psi}{\partial \lambda_{i,r}} = \sum_m \frac{\partial\psi_{m,i}}{\partial f_{m,i}} d_{m,r}, \quad (94)$$

where  $\psi_{n,i} := \psi(x_{n,i}, f_{n,i}(\mathbf{x}, \lambda_i))$ ,  $n \in \mathcal{N}_i$ , and  $\psi_{n,i} = 0$ ,  $n \notin \mathcal{N}_i$ . Let

$$\mathbf{d}_b\psi := \left[ \frac{\partial\psi_{n,i}}{\partial f_{n,i}} \frac{\partial f_{n,i}}{\partial x_n} \right]_{NK \times 1}. \quad (95)$$

$(\mathbf{x}, \lambda)$  is a stationary point of  $\Psi$  when  $\nabla\Psi = [\nabla_{\mathbf{x}}\Psi; \nabla_{\lambda}\Psi] = 0$ . If we pre-multiply  $\nabla_{\mathbf{x}}\Psi$  by  $\mathbf{d}_b\psi$ , it follows that (c.f. (93)):

$$\begin{aligned} (\mathbf{d}_b\psi)^T \nabla_{\mathbf{x}}\Psi &= \sum_{n,i} \frac{\partial\psi_{n,i}}{\partial x_{n,i}} \frac{\partial\psi_{n,i}}{\partial f_{n,i}} \frac{\partial f_{n,i}}{\partial x_n} \\ &+ \sum_n \left( \sum_i \frac{\partial\psi_{n,i}}{\partial f_{n,i}} \frac{\partial f_{n,i}}{\partial x_n} \right)^2 = 0, \end{aligned} \quad (96)$$

where,

$$\frac{\partial f_{n,i}}{\partial x_n} := \frac{-1}{\phi_n \gamma_{n,i}^2} g''_i \left( \frac{x_n}{\phi_n \gamma_{n,i}} \right) > 0, \quad (97)$$

owing to strict concavity of  $g_i(\cdot)$ . For the complementary function  $\psi(a, b)$  defined in (36), it is shown that [30]:

$$\frac{\partial\psi}{\partial a} \frac{\partial\psi}{\partial b} \geq 0, \quad \forall (a, b) \in \mathbb{R}^2, \quad (98)$$

$$\frac{\partial\psi}{\partial a} = \frac{\partial\psi}{\partial b} = 0 \iff \psi(a, b) = 0. \quad (99)$$

Hence, the right hand side of (96) is strictly positive, unless  $\partial\psi_{n,i}/\partial x_{n,i} = \partial\psi_{n,i}/\partial f_{n,i} = 0$ ,  $\forall n, i$ , or equivalently  $\psi(x_{n,i}, f_{n,i}(\mathbf{x}, \lambda_i)) = 0$ ,  $\forall n, i$ , which is the case only if  $\mathbf{x}$  is a solution to Problem 1. Conversely, for every solution to Problem 1,  $\psi(x_{n,i}, f_{n,i}(\mathbf{x}, \lambda_i)) = 0$ ,  $\forall n, i$ , which implies that  $\partial\psi_{n,i}/\partial x_{n,i} = \partial\psi_{n,i}/\partial f_{n,i} = 0$ ,  $\forall n, i$ , and therefore  $\nabla\Psi = 0$ .  $\square$

**Proof of Theorem 7:** By contradiction, assume that there exists another solution  $(\mathbf{x}', \lambda')$ , arbitrarily close to  $(\mathbf{x}, \lambda)$ , for which  $x'_m \neq x_m$ , for some user  $m$ . Since  $(\mathbf{x}', \lambda')$  is arbitrarily close to  $(\mathbf{x}, \lambda)$ , by continuity we conclude that  $f_{n,i}(\mathbf{x}', \lambda'_i) > 0$  if  $f_{n,i}(\mathbf{x}, \lambda_i) > 0$ . Hence,  $x'_{n,i} > 0$  only if  $f_{n,i}(\mathbf{x}, \lambda_i) = 0$  (c.f. (29)). In the same way,  $\lambda'_{i,r} > 0$  only if  $f_{i,r}(\mathbf{x}_i) = 0$ . On the other hand, the assumption that

$(\mathbf{x}', \lambda')$  is a solution to Problem 3, implies that  $f_{i,r}(\mathbf{x}'_i) = 0$ , for every server  $i$  and resource  $r$  with  $\lambda'_{i,r} > 0$  (c.f. (28)). Hence,

$$\Delta f_{i,r} := f_{i,r}(\mathbf{x}'_i) - f_{i,r}(\mathbf{x}_i) = 0, \quad \forall i, r : \lambda'_{i,r} > 0. \quad (100)$$

Furthermore, (29) implies that  $f_{n,i}(\mathbf{x}', \lambda'_i) = 0$ , for every user  $n$  and server  $i$  with  $x'_{n,i} > 0$ . Hence,

$$\Delta f_{n,i} := f_{n,i}(\mathbf{x}', \lambda'_i) - f_{n,i}(\mathbf{x}, \lambda_i) = 0, \quad \forall n, i : x'_{n,i} > 0 \quad (101)$$

Let  $\delta x_{n,i} := x'_{n,i} - x_{n,i}$  and  $\delta \lambda_{i,r} := \lambda'_{i,r} - \lambda_{i,r}$ . Since  $(\mathbf{x}', \lambda')$  could be arbitrarily close to  $(\mathbf{x}, \lambda)$ , it follows from (100) and (101) that:

$$\Delta f_{i,r} = - \sum_n d_{n,r} \delta x_{n,i} = 0, \quad \forall i, r : \lambda'_{i,r} > 0, \quad (102)$$

$$\begin{aligned} \Delta f_{n,i} &\simeq \frac{\alpha_i \delta x_n}{x_n} \left[ \frac{\phi_n^{\alpha_i} \gamma_{n,i}^{\alpha_i-1} A_i^{\alpha_i}}{x_n^{\alpha_i}} \right] + \frac{B_i \delta x_n}{x_n^2} \\ &+ \sum_r d_{n,r} \delta \lambda_{i,r} = 0, \quad \forall n, i : x'_{n,i} > 0, \end{aligned} \quad (103)$$

where  $\delta x_n := \sum_i \delta x_{n,i}$ . Let  $r^*(i)$  denote the only resource at server  $i$  for which  $\lambda_{i,r^*(i)} > 0$ . The fact that  $f_{n,i}(\mathbf{x}, \lambda_i) = 0$  for every user  $n$  with  $x'_{n,i} > 0$ , implies that:

$$d_{n,r^*(i)} \lambda_{i,r^*(i)} = \frac{\phi_n^{\alpha_i} \gamma_{n,i}^{\alpha_i-1} A_i^{\alpha_i}}{x_n^{\alpha_i}} + \frac{B_i}{x_n}. \quad (104)$$

Hence, we may rewrite (103) as:

$$d_{n,r^*(i)} \delta \lambda_{i,r^*(i)} + \omega_{n,i} \delta x_n = 0, \quad \forall n, i : x'_{n,i} > 0, \quad (105)$$

where,  $\omega_{n,i}$  is defined as:

$$\omega_{n,i} := \left[ \frac{\alpha_i d_{n,r^*(i)} \lambda_{i,r^*(i)}}{x_n} + \frac{B_i(1 - \alpha_i)}{x_n^2} \right]. \quad (106)$$

In the following, we check the possibility for existence of  $(\mathbf{x}', \lambda')$  such that (102) and (105) are established. In general, one may partition the set of users, and the set of servers respectively, into  $L$  partitions,  $\mathcal{N} = \{\mathcal{N}_1, \mathcal{N}_2, \dots, \mathcal{N}_L\}$  and  $\mathcal{K} = \{\mathcal{K}_1, \mathcal{K}_2, \dots, \mathcal{K}_L\}$ , such that  $x'_{n,i} = 0$  for each user  $n \in \mathcal{N}_l$  and for every server  $i \notin \mathcal{K}_l$ . Here, without loss of generality, we assume that  $L = 1$  is the greatest number of partitions which could be found [Otherwise, we should consider each partition separately].

By assumption,  $\sum_i \delta x_{m,i} = \delta x_m \neq 0$  for user  $m$ . Without loss of generality assume that  $\delta x_m > 0$ . From (105) it follows that  $\delta \lambda_{i,r^*(i)} < 0$ , for every server  $i$  for which  $x'_{m,i} > 0$ . Furthermore, (105) implies that  $\delta x_n > 0$  for every user  $n$ , for which  $x'_{m,j} x'_{n,j} > 0$  for some server  $j$ . Given that all users and all servers reside in the same partition, it follows that  $\delta x_n > 0$  for all users, and  $\delta \lambda_{i,r^*(i)} < 0$  for all servers.

Let  $\mathcal{S}$ ,  $\mathcal{S} := |\mathcal{S}|$ , denote the set of servers for which  $\delta x_i = [\delta x_{n,i}] \neq 0$ . For every server  $i \in \mathcal{S}$ , (102) implies that there exists some user  $p$  for which  $\delta x_{p,i} < 0$ . On the other hand, for each user  $n$  there exists some server  $i$  for

which  $\delta x_{n,i} > 0$ . Hence,  $x'_{n,i} > 0$  for at least  $N+S$  pairs of users and servers, for which (105) should be established.

To present the system of equations in (105) in a matrix form, let define  $\chi_i := \{n \mid x'_{n,i} > 0\}$ . Also, define  $W_i = \text{diag}([\omega_{n,i}])$ , and  $V_i = W_i(\{n \in \chi_i, \cdot\})$  as a sub-matrix of  $W_i$ , consisting of a subset of rows in  $W_i$  which corresponds to users  $n \in \chi_i$ . Then, the system of equations in (105) can be written as:

$$\begin{bmatrix} \mathbf{d}_1^* & 0 & \cdots & 0 & V_1 \\ 0 & \mathbf{d}_2^* & \cdots & 0 & V_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \mathbf{d}_S^* & V_S \end{bmatrix} \begin{bmatrix} \delta\lambda_{1,r^*(1)} \\ \vdots \\ \delta\lambda_{S,r^*(S)} \\ \delta\mathbf{x} \end{bmatrix} = 0, \quad (107)$$

where  $\mathbf{d}_i^* := [d_{n,r^*(i)} \mid n \in \chi_i]$ , and  $\delta\mathbf{x} := [\delta x_n \mid n \in \mathcal{N}]$ . Here, without loss of generality, we have assumed that servers in  $\mathcal{S}$  are indexed from 1 to  $S$ . As a short hand notation, we denote the coefficient matrix in (107) by  $[D \mid V]$ . In the following we show that the coefficient matrix in (107) has a column rank of  $N + S$ .

For the matrix  $V$ , we know that exactly one element is non-zero in each row. Furthermore, there exists at least a non-zero element in each column, owing to the fact that for each user  $n$ ,  $x'_{n,i} > 0$  for at least one server  $i$ . Hence, we may find  $N$  linear independent rows in  $V$ , or equivalently  $N$  linear independent columns, which form this matrix. It means that  $V$  is a full rank matrix which has a column rank of  $N$ . Furthermore, from (105) and (106), it can be observed that none of the columns (or group of columns) in  $V$  can be expressed as a linear combination of columns of  $D$ , provided that  $B_i > 0, \forall i$ . Also, none of the columns (or group of columns) in  $D$  can be expressed in terms of columns in  $V$ , owing to the assumption that all users and all servers reside in the same partition. Accordingly, the coefficient matrix,  $[D, V]$ , has a column rank of  $N + S$ . Hence, the only solution to (107) is  $\delta\lambda_i = 0, \forall i$ , and  $\delta\mathbf{x} = 0$ . However, this contradicts the assumption that  $\delta x_m \neq 0$  for user  $m$ .  $\square$

**Proof of Lemma 3:**  $(\mathbf{v}_x^h)_i$  is the projection of  $-\nabla_{\mathbf{x}_i} \Psi$  onto  $\Omega_i$ , when:

$$(\mathbf{v}_x^h)_i = \arg \min_{\mathbf{v}_i \in \Omega_i} \frac{1}{2} \|\mathbf{v}_i + \nabla_{\mathbf{x}_i} \Psi\|^2. \quad (108)$$

According to the constrained optimality theorem,  $(\mathbf{v}_x^h)_i$  is a solution to minimization in (108) if and only if [34]:

$$[(\mathbf{v}_x^h)_i + \nabla_{\mathbf{x}_i} \Psi]^T [\mathbf{v}_i - (\mathbf{v}_x^h)_i] \geq 0, \quad (109)$$

for every  $\mathbf{v}_i \in \Omega_i$ . Substituting  $\mathbf{v}_i = 0$  results in:

$$(\mathbf{v}_x^h)_i^T \nabla_{\mathbf{x}_i} \Psi \leq -\|(\mathbf{v}_x^h)_i\|^2 \quad (110)$$

which implies that  $\mathbf{v}_x^h$  is a strictly descent direction, unless  $\mathbf{v}_x^h = 0$ . On the other hand, for every  $i, r$  with

$(\mathbf{v}_\lambda^h)_{i,r} \neq 0$ , it follows from (47) and (48) that:

$$\begin{aligned} \frac{\partial \Psi}{\partial \lambda_{i,r}} \times (\mathbf{v}_\lambda^h)_{i,r} &= - \left( \frac{\partial \Psi}{\partial \lambda_{i,r}} \right)^2 \\ &+ \beta_{i,r}^h \frac{\partial \Psi}{\partial \lambda_{i,r}} \cup \left( -\beta_{i,r}^h \frac{\partial \Psi}{\partial \lambda_{i,r}} \right) < 0, \quad (111) \end{aligned}$$

which implies that  $\mathbf{v}_\lambda^h$  is a strictly descent direction, unless  $\mathbf{v}_\lambda^h = 0$ .  $\square$

**Proof of Lemma 4:** By contradiction, assume that  $\nabla_{\mathbf{x}} \Psi(\mathbf{x}^h, \lambda^h) \neq 0$  and  $\mathbf{v}^h = [\mathbf{v}_x^h; \mathbf{v}_\lambda^h] = 0$ . Consider some server  $i$  for which  $\nabla_{\mathbf{x}_i} \Psi(\mathbf{x}^h, \lambda^h) \neq 0$ . The projection of  $-\nabla_{\mathbf{x}_i} \Psi$  onto  $\Omega_i$  can be found by solving the optimization in (108). The problem in (108) is to minimize a convex function over an affine feasible region (see the definition of  $\Omega_i$  in (46)). The KKT conditions for this problem imply that:

$$0 \leq \varsigma_{i,r} \perp -(\mathbf{v}_x^h)_i^T \mathbf{d}_r \geq 0, \quad r \in \mathcal{R}_i^h : \lambda_{i,r} = 0, \quad (112)$$

$$(\mathbf{v}_x^h)_i^T \mathbf{d}_r = 0, \quad r \in \mathcal{R}_i^h : \lambda_{i,r} > 0, \quad (113)$$

$$(\mathbf{v}_x^h)_i + \nabla_{\mathbf{x}_i} \Psi + \sum_{r \in \mathcal{R}_i^h} \varsigma_{i,r} \mathbf{d}_r = 0. \quad (114)$$

Here, without loss of generality, we assume that the vectors in  $\{\mathbf{d}_r \mid r \in \mathcal{R}_i^h\}$  are linear independent [Otherwise we may restrict the conditions in  $\Omega_i$  to a subset  $\mathcal{A} \subseteq \mathcal{R}_i^h$  for which  $\{\mathbf{d}_r \mid r \in \mathcal{A}\}$  are linear independent.]. The projection of  $-\nabla_{\mathbf{x}_i} \Psi$  onto  $\Omega_i$  is zero,  $(\mathbf{v}_x^h)_i = 0$ , when:

$$-\nabla_{\mathbf{x}_i} \Psi = \sum_{r \in \mathcal{R}_i^h} \varsigma_{i,r} \mathbf{d}_r, \quad (115)$$

for some scalars  $\varsigma_{i,r} \in \mathbb{R}$  satisfying (112). In other words,  $\mathbf{v}_x^h = 0$  if and only if there exists some scalars  $\varsigma_{i,r} \in \mathbb{R}$ , so that (115) is established for all servers with  $\nabla_{\mathbf{x}_i} \Psi \neq 0$ .

On the other hand,  $\mathbf{v}_\lambda^h = 0$  requires that  $(\tilde{\mathbf{v}}_\lambda^h)_{i,r} = 0$ , for every  $r \in \mathcal{R}_i^h$  with  $\lambda_{i,r} > 0$  (c.f. (48)). However,  $(\tilde{\mathbf{v}}_\lambda^h)_{i,r} = 0$  only when  $\partial \Psi / \partial \lambda_{i,r} = 0$  and  $\beta_{i,r}^h = -\mathbf{d}_r^T \nabla_{\mathbf{x}_i} \Psi = 0$  (c.f. (47)). it follows from (115) that:

$$\beta_{i,r}^h = -\mathbf{d}_r^T \nabla_{\mathbf{x}_i} \Psi = \mathbf{d}_r^T \sum_{r' \in \mathcal{R}_i^h} \varsigma_{i,r'} \mathbf{d}_{r'}. \quad (116)$$

Since the vectors  $\{\mathbf{d}_r \mid r \in \mathcal{R}_i^h\}$  are linear independent, It follows from (116) that  $\beta_{i,r}^h = 0$  only when  $\varsigma_{i,r} = 0$ . Hence,  $\varsigma_{i,r} = 0$  for every  $r \in \mathcal{R}_i^h$  with  $\lambda_{i,r} > 0$ . So, the right hand side of (115) would be zero under Assumption 1. That is,  $\nabla_{\mathbf{x}_i} \Psi = 0, \forall i$ . This, however, violates the contradictory assumption.

Now let Assumption 1 do not hold, so that for server  $i$  with  $\nabla_{\mathbf{x}_i} \Psi \neq 0$  there exists some resource  $r \in \mathcal{R}_i^h$  with  $\lambda_{i,r} = 0$  and  $\varsigma_{i,r} > 0$  (c.f. (112)). Let  $\bar{\mathcal{R}}_i^h := \{r \in \mathcal{R}_i^h \mid \lambda_{i,r} = 0 \text{ and } \varsigma_{i,r} > 0\}$ . Then, (115) implies that  $\partial \Psi / \partial x_{n,i} < 0$ , for every user  $n$  with  $d_{n,r} > 0, r \in \bar{\mathcal{R}}_i^h$ . On the other hand,  $\partial \Psi / \partial x_{n,i} = 0$  for every user  $n \notin \bar{\mathcal{N}}_i$  (c.f. (93)). Hence, (115) may not be established if there exists some user  $n \notin \bar{\mathcal{N}}_i$  with  $d_{n,r} > 0, r \in \bar{\mathcal{R}}_i^h$ , which is the case under Assumption 2. In other words, (115) may not be established under Assumption 2.  $\square$