

Per-Server Dominant-Share Fairness (PS-DSF): A Multi-Resource Fair 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—Users of cloud computing platforms pose different types of demands for multiple resources on servers (physical or virtual machines). Besides differences in their resource capacities, servers may be additionally heterogeneous in their ability to service users - certain users' tasks may only be serviced by a subset of the servers. We identify important shortcomings in existing multi-resource fair allocation mechanisms - Dominant Resource Fairness (DRF) and its follow up work - when used in such environments. We develop a new fair allocation mechanism called Per-Server Dominant-Share Fairness (PS-DSF) which we show offers all desirable sharing properties that DRF is able to offer in the case of a single "resource pool" (i.e., if the resources of all servers were pooled together into one hypothetical server). We evaluate the performance of PS-DSF through simulations. Our evaluation shows the enhanced efficiency of PS-DSF compared to the existing allocation mechanisms. We argue how our proposed allocation mechanism is applicable in cloud computing networks and especially large scale data-centers.

I. 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. Resource allocation and especially fair sharing in such shared computing system is particularly challenging because of the following reasons: a) heterogeneity of servers, b) placement constraints, c) dealing with multiple types of resources, and d) diversity of workloads and demands.

Real world data-centers are comprised of heterogeneous machines/servers with different configurations, where some machines might be incompatible for some processing purposes/tasks. Furthermore, each user may have specific requirements which further restrict the set of servers that the tasks of the user may run on. For example, a user may require a machine with a public IP address, particular kernel version, special hardware such as GPUs, or large amounts of memory, and might be unable to run on machines which lack such requirements. For instance, it has been observed that over 50%

of tasks at Google clusters have strict constraints about the machines they can run on [1], [2].

Besides placement constraints, users present diversity over the amount of resources they need for executing one task. For instance, the tasks of some users might be CPU intensive while for others memory or I/O bandwidth might be a bottleneck. Dominant Resource Fairness is the first allocation mechanism which describes a notion of fairness when allocating multiple types of resources [3]. With DRF users receive a fair share of their *dominant resource*. 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 expressed as a fraction of the overall resource capacity spread across all available servers. There are several other works investigating DRF allocation in case that different resources are distributed over heterogeneous servers but there are no placement constraints [4], [5], [6], [7].

There are some recent works investigating max-min fair allocation/scheduling for one type of resource while respecting placement constraints [2], [8], [9], [10], [11], [12]. These schedulers could be useful in a multi-resource setting only when one of the resources serves as the bottleneck for all users, otherwise they might result in poor resource utilization [3], [2]. There are limited works in the literature investigating multi-resource fair allocation while respecting placement constraints [13], [14], [15], [16]. In this case, it is unclear how to globally identify the dominant resource as well as the dominant share for different users, as each user may have access only to a subset of servers. [14], [15] present an elementary extension of DRF which identify the share of each user by ignoring the placement constraints and applying the same ideas as the unconstrained setting. We show that this approach does not achieve fairness even in the specific case that one of the resources serves as a bottleneck (Further discussions could be found in Section II-B).

Our Contributions: We propose a new allocation mechanism called *Per-Server Dominant Share Fairness*. We show that PS-DSF achieves all the desirable properties offered by DRF for a single resource pool: sharing incentive, strategy proofness, envy freeness, Pareto optimality, bottleneck fairness and single

resource fairness (A detailed description of these properties can be found in Section II-A). In fact, PS-DSF reduces to DRF when one considers a single server system.

The intuition behind PS-DSF is to compare and weigh the allocated resources to each user from the perspective of each server. PS-DSF identifies a dominant resource and 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 (w.r.t.) server i describes the fraction of the dominant resource which should be allocated to user n from server i as if all user n 's tasks were allocated resources solely from server i . Each server may then use this localized metric to decide whether to increase or decrease the allocated tasks to each user without the need to identify a global dominant share. Besides its enhanced performance, PS-DSF is the first (to our knowledge) principled allocation mechanism which could be intrinsically implemented in a distributed manner.

The rest of this paper is organized as follows: In Section II, after describing the model, we give the necessary background and discuss insufficiency of the existing multi-resource allocation mechanisms, especially in case of heterogeneous servers with placement constraints. After presenting our proposed allocation mechanism in Section III, we investigate different sharing properties that it satisfies and present a distributed algorithm to realize it. We present some numerical experiments in Section IV, and finally we draw conclusions in Section V.

II. BACKGROUND AND MODEL

Consider a set \mathcal{K} of $K = |\mathcal{K}|$ heterogeneous servers (resource pools) each containing M types of resources. We denote by $c_{i,r}$ the capacity of resource r on server i , where $c_{i,r} \geq 0$. Let \mathcal{N} denote the set of active users, where $N = |\mathcal{N}|$. Let $d_n = [d_{n,r}]$ denote the per task *demand vector* for user $n \in \mathcal{N}$, that is the amount of each resource required for executing one task for user n . Let $\phi_n > 0$ denote the weight associated with user n . The weights reflect the priority of users with respect to each other.

Due to heterogeneity of users and servers, each user may be restricted to get service only from a subset of servers. For example, each user may have some special hardware/software requirements (e.g., public IP address, a particular kernel version, GPU, etc.) which restrict the set of servers that the tasks of the user may run on. Besides such explicit placement constraints, users may not run their tasks on servers which lack some required resources.

For instance, consider the example in Figure 1, where three types of resources, CPU, memory, and network bandwidth are available over two servers in the amounts of $\mathbf{c}_1 = [9 \text{ cores}, 12\text{GB}, 100\text{Mb/s}]$ and $\mathbf{c}_2 = [12 \text{ cores}, 12\text{GB}, 0\text{Mb/s}]$, where no communication bandwidth is available over the second server. Consider three users with the weights $\phi_1 = \phi_2 = 1$, $\phi_3 = 2$, whose demand vectors are $\mathbf{d}_1 = [1, 2, 10]$, $\mathbf{d}_2 = [1, 2, 1]$ and $\mathbf{d}_3 = [1, 2, 0]$. Accordingly, users 1 and 2 are restricted to get service only from the first server, while user 3

may get service from both servers. In summary, let $\delta_{n,i} = 1$ if the tasks of user n can run on server i , and otherwise $\delta_{n,i} = 0$.

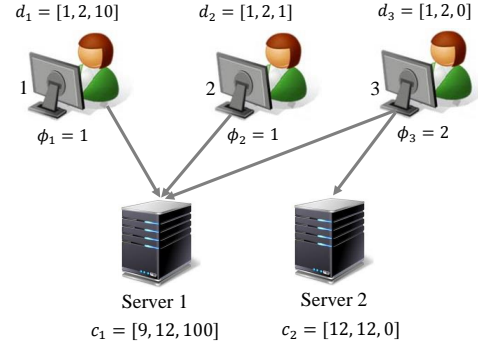


Fig. 1: A heterogeneous multi-resource system with three users and two servers.

A. Dominant Resource Fairness

The problem of multi-resource fair allocation was originally studied in [3] 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 \mathcal{A} . 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:

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

where, x_n is a non-negative real number. [3] argues that the following important properties must be satisfied by a multi-resource allocation mechanism:

- *Sharing Incentive*: Consider a generic *uniform allocation* where every user n is allocated $\phi_n / \sum_m \phi_m$ portion of each resource. An allocation is said to provide sharing incentive, when each user is able to run more tasks compared to the uniform allocation.
- *Envy freeness*: A user should not prefer the allocation of another user when adjusted according to their weights, i.e., $U_n(\mathbf{a}_n) \geq U_n(\frac{\phi_n}{\phi_m} \mathbf{a}_m)$, for all m .
- *Pareto Optimality*: It should not be possible to increase x_n for any user n , without decreasing x_m for some user m .
- *Strategy Proofness*: Users should not be able to increase their utilization by lying about their resource demands.

Sharing incentive provides 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. Pareto optimality results in maximizing system utilization. Finally, strategy proofness prevents users from gaming the allocation mechanism. The reader is referred to [3] or [17] for further details.

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 [3]:

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

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 *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., $\mathbf{a}_n = x_n \mathbf{d}_n$, $\forall n$. In this case, an allocation $\{x_n\}$ is feasible when:

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

Definition 1. It is said that $\{x_n\}$ satisfies DRF, if it is feasible and the normalized dominant share for each user, s_n/ϕ_n cannot be increased while maintaining feasibility without decreasing s_m for some user m with $s_m/\phi_m \leq s_n/\phi_n$ [3].

DRF is a restatement of max-min fairness in terms of *dominant shares*. What make it appealing are desirable sharing properties which are satisfied under this allocation mechanism. Besides the above-mentioned essential properties, DRF also satisfies the following simple but essential properties [3].

- *Single Resource Fairness*: When there is only one resource type, the allocation satisfies max-min fairness.
- *Bottleneck Fairness*: If there is one resource which is dominantly requested by each user, then the allocation satisfies max-min fairness for that resource.

B. Challenges with Heterogeneous Resource-Pools and Placement Constraints

The notion of DRF has been extended to the case of heterogeneous servers, when all types of resources are available within each server and there are no placement constraints [7]. In this case, DR for user n is readily identified as the resource whose greatest portion is required for execution of one task as if all resources were integrated at resource pool. That is, DR for user n could be identified according to (2), where $c_r := \sum_i c_{i,r}$ is the total capacity of resource r . 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 here is the total number of tasks which are allocated to user n from different servers, that is $x_n := \sum_i x_{n,i}$. In [7] it is proposed to find $\{x_{n,i}\}$ such that max-min fairness is achieved in terms of global dominant shares. This mechanism, which is referred to as DRFH, has been shown to achieve Pareto optimality, strategy proofness, envy freeness and bottleneck fairness. However, it fails to provide sharing incentive [7].

When there are placement constraints, it is unclear how to define a single system-wide DR for a user similar to that in [7]. A natural first thought may be to identify the DR over the set of eligible servers for each user. However, in this case users may have an incentive to misreport the set of eligible servers [14]. A strategy-proof approach is to identify the DR for each user as if there were no placement constraints and all

resources were integrated at one resource pool. We argue that this approach, which we refer to as C-DRFH, does not result in a fair allocation as it does not satisfy bottleneck fairness.

To appreciate this shortcoming of C-DRFH, consider the example in Figure 1, where the second resource (RAM) is dominantly requested by every user from its eligible servers. If we allocate the available RAM proportionate to the weights, 6GB is allocated to the first two users and 12 GB is allocated to the third user. Accordingly, each user is allocated $x_1 = x_{1,1} = 3$, $x_2 = x_{2,1} = 3$, $x_3 = x_{3,2} = 6$ tasks (this allocation follows from our proposed allocation mechanism - see Section III). However, C-DRFH would instead identify bandwidth as the dominant resources for the first user and identifies RAM as the dominant resource for the second and third users. Hence, if we allocate global dominant shares in a weighted fair manner, each user is allocated $x_1 = x_{1,1} = 2.609$, $x_2 = x_{2,1} = 3.130$, and $x_3 = x_{3,1} + x_{3,2} = 6 + 0.261 = 6.261$ tasks respectively, which obviously violates fairness on the bottleneck resource.

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

Property	C-DRFH	TSF	PS-DSF
SI		✓	✓
EF	✓	✓	✓
SP	✓	✓	*
PO	✓	✓	*
BF			✓

Yet another extension of DRF that also considers heterogeneous servers, all containing all types of resources without any placement constraints, is CDRF [4]. Specifically, let $\gamma_n := \sum_i \gamma_{n,i}$ be defined as the number of tasks which are allocated to user n when monopolizing the whole cluster (i.e., if n were the only user running on the cluster). An allocation is said to satisfy CDRF¹, when x_n/γ_n satisfies max-min fairness. In case of one server, x_n/γ_n gives dominant share for each user n . As a result, CDRF reduces to DRF in case of one server. In case of multiple heterogeneous servers with no placement constraints, CDRF is shown to satisfy Pareto optimality, strategy proofness, envy freeness and sharing incentive properties [4].

In [14] CDRF has been extended to address the placement constraints. Specifically, let $\gamma_n := \sum_i \gamma_{n,i}$ be (re)defined as the number of tasks which are allocated to user n from different servers when monopolizing all servers as if there were no placement constraints [14]. An allocation is said to satisfy Task Share Fairness (TSF), when x_n/γ_n satisfies max-min fairness. TSF is shown to satisfy Pareto optimality, strategy proofness, envy freeness and sharing incentive properties in case of heterogeneous servers with placement constraints [14]. However, we argue that this mechanism is not essentially fair as it does not satisfy bottleneck fairness.

For instance, consider again the example in Figure 1. The number of tasks that each user may run in the whole cluster

¹Containerized DRF

is $\gamma_1 = \gamma_2 = 6$, and $\gamma_3 = 12$ tasks, respectively. Hence, each user is allocated $x_1 = x_{1,1} = 2$, $x_2 = x_{2,1} = 2$ and $x_3 = x_{3,1} + x_{3,2} = 6 + 2 = 8$ tasks according to TSF mechanism, which is completely different and far from the fair allocation. Table I summarizes different sharing properties which could be satisfied under different allocation mechanisms. Shortcomings of the existing allocation mechanisms in case of heterogeneous servers with placement constraints motivates us to develop a new allocation mechanism.

III. PER-SERVER DOMINANT SHARE FAIRNESS

In this section we describe PS-DSF, an extension of DRF that is applicable for heterogeneous resource-pools in the presence of placement constraints. As discussed in the previous section, in the case of heterogeneous servers and in the presence of placement constraints, it is unclear how to globally identify one DR and the corresponding dominant share for each user. The intuition behind PS-DSF is to define a *virtual dominant share* for every user w.r.t. each server. Towards this, we first define the DR for every user n w.r.t. each server i as:

$$\rho(n, i) := \arg \max_r \frac{d_{n,r}}{c_{i,r}}. \quad (6)$$

Let $\gamma_{n,i}$ denote the number of tasks which could be executed by user n when monopolizing server i :

$$\gamma_{n,i} := \delta_{n,i} \min_r \frac{c_{i,r}}{d_{n,r}} = \delta_{n,i} \frac{c_{i,\rho(n,i)}}{d_{n,\rho(n,i)}}. \quad (7)$$

We say that server i is *eligible* to serve user n when $\gamma_{n,i} > 0$ or equivalently $\delta_{n,i} = 1$. Without loss of generality we restrict ourselves to non-wasteful allocations, that is $\mathbf{a}_{n,i} = x_{n,i} \mathbf{d}_n$, where $\mathbf{a}_{n,i} = [a_{n,i,r}]$ is the vector of allocated resources to user n from server i and $x_{n,i} \in \mathbb{R}^+$ is the number of allocated tasks from the same server.

Definition 2. The *Virtual Dominant Share (VDS)* for user n w.r.t. server i , $s_{n,i}$, is defined as:

$$s_{n,i} = \frac{x_n}{\gamma_{n,i}}, \quad (8)$$

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).

Intuitively, $s_{n,i}$ gives the fraction² of the dominant resource for user n w.r.t. server i which should be allocated to it as if x_n tasks were allocated to it entirely from server i . When the available resources over each server are arbitrarily divisible, we have the following condition on $\{x_{n,i}\}$ to be feasible.

Definition 3. An allocation, $\{x_{n,i}\}$, is said to satisfy *Resource Division Multiplexing (RDM)* constraint, when:

$$\sum_n x_{n,i} d_{n,r} \leq c_{i,r}, \quad \forall i, r. \quad (9)$$

For a data-center comprising of a plurality of servers, it is sometimes of more practical interest to assume that servers

²The reader may note that $s_{n,i}$ could be possibly greater than 1, as some tasks could be allocated to user n from other servers.

may not be divided to finer partitions [2]. Accordingly, the hypervisor may only time-share servers among different users. In this case, we have the following condition on $\{x_{n,i}\}$ to be feasible.

Definition 4. An allocation, $\{x_{n,i}\}$, is said to satisfy *Time Division Multiplexing (TDM)* constraint, when³:

$$\sum_n x_{n,i} / \gamma_{n,i} \leq 1, \quad \forall i. \quad (10)$$

It can be observed that TDM constraint is more stringent than RDM constraint, as (10) implies (9):

$$1 \geq \sum_n \frac{x_{n,i}}{\gamma_{n,i}} = \sum_n \frac{x_{n,i} d_{n,\rho(n,i)}}{c_{i,\rho(n,i)}} \geq \frac{\sum_n x_{n,i} d_{n,r}}{c_{i,r}}, \quad \forall i, r. \quad (11)$$

We investigate our proposed allocation mechanism under both of these feasibility conditions.

Definition 5. An allocation $\{x_{n,i}\}$ satisfies *Per-Server Dominant-Share Fairness*, 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 $s_{m,i} / \phi_m \leq s_{n,i} / \phi_n$.

A. An Example

Consider again the heterogeneous servers from our earlier example this time serving four equally weighted users whose demand vectors are $\mathbf{d}_1 = [1.5, 1, 10]$, $\mathbf{d}_2 = [1, 2, 10]$, $\mathbf{d}_3 = [0.5, 1, 0]$, $\mathbf{d}_4 = [1, 0.5, 0]$. We show this in Figure 2. Note the placement constraints for users 1 and 2 whose tasks may only run on the first server. We show the PS-DSF allocation (based on RDM) in Figure 3. The allocated tasks to each user are $x_1 = x_{1,1} = 3.6$, $x_2 = x_{2,1} = 3.6$, $x_3 = x_{3,2} = 8$, $x_4 = x_{4,2} = 8$, respectively, where no tasks are allocated to users 3 and 4 from the first server. Specifically, the VDS for user 3 (and user 4 respectively) w.r.t. the first server is $s_{3,1} = 8/12$ ($s_{4,1} = 12/12 = 1$), while the VDS of users 1 and 2 w.r.t. this server is $s_{1,1} = s_{2,1} = 0.6$. The VDS of users 3 and 4 w.r.t. the second server is $s_{3,2} = s_{4,2} = 8/12$. The reader may verify that for each server i the allocated tasks to each user may not be increased without decreasing the allocated tasks of some user with less VDS.

B. The Properties of the PS-DSF Allocation Mechanism

Before examining different sharing properties satisfied under PS-DSF allocation mechanism, we describe a necessary and sufficient condition to achieve PS-DSF.

Definition 6. Given a feasible allocation $\{x_{n,i}\}$ based on RDM, we say that r is a *bottleneck resource* for user n w.r.t. an eligible server i if $d_{n,r} > 0$, $\sum_m x_{m,i} d_{m,r} = c_{i,r}$ (i.e. r is saturated), and

$$\frac{s_{n,i}}{\phi_n} \geq \frac{s_{m,i}}{\phi_m}, \quad \forall m \text{ such that } x_{m,i} d_{m,r} > 0. \quad (12)$$

³Considering resources such as CPU, BW, \dots , which are attributed a processing speed per time-unit, $x_{n,i} / \gamma_{n,i}$ represent the percentage of time-unit that server i is allocated to user n .

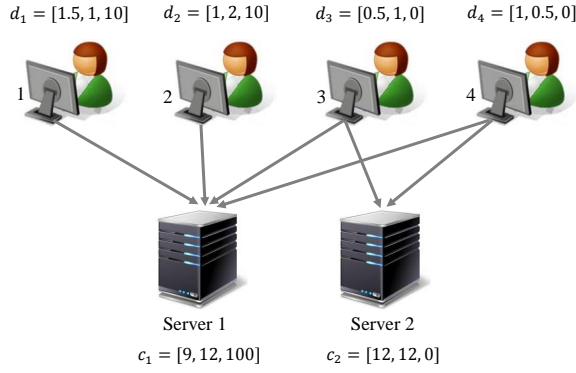


Fig. 2: A heterogeneous multi-resource system with four users and two servers.

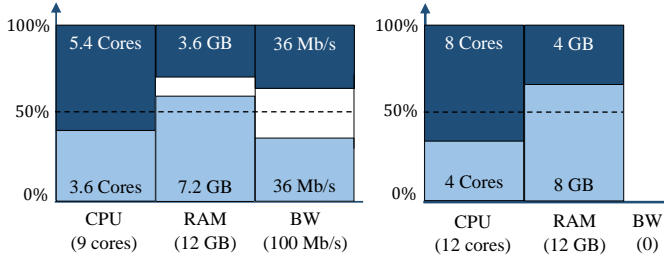


Fig. 3: PS-DSF allocation for the example in Figure 2.

Theorem 1. A feasible allocation $\{x_{n,i}\}$ based on RDM satisfies PS-DSF if and only if there exists a bottleneck resource for every user w.r.t. every eligible server.

Theorem 2. A feasible allocation $\{x_{n,i}\}$ based on TDM satisfies PS-DSF if and only if (10) holds with equality, and

$$\frac{s_{n,i}}{\phi_n} \geq \frac{s_{m,i}}{\phi_m}, \quad \forall n \text{ and } \forall m \text{ such that } x_{m,i} > 0. \quad (13)$$

The proofs are given in the appendix. These conditions will be useful in determining a PS-DSF allocation (see Section III-C). In the following we examine different sharing properties that are satisfied under PS-DSF. In case of a heterogeneous system with placement constraints, we will need to extend the notion of *Sharing Incentive*, *Strategy Proofness* and *Bottleneck Fairness*. Other properties, *Envy Freeness*, *Pareto Optimality* and *Single Resource Fairness*, will follow the same definitions as described in Section II-A.

The generalization of sharing incentive property is straightforward. We consider a *uniform allocation* which allocates $\phi_n / \sum_m \phi_m$ portion of the resources on each server (whether this server is eligible or not) to each user n . An allocation is said to satisfy *sharing incentive*, when each user is able to run more tasks compared to such uniform allocation.

For the strategy proofness property, we may note that we assume each user to declare its demand vector and also the set of eligible servers. We say that an allocation satisfies *strategy proofness* when users may not increase their utilization by lying about their resource demands or the set of eligible servers.

Finally, a resource is considered as a bottleneck in the whole system when it is dominantly requested by each user from every eligible server. If there is a bottleneck resource, then the allocation should satisfy max-min fairness w.r.t. such resource.

Theorem 3. PS-DSF allocation mechanism (whether based on RDM or TDM) satisfies single resource fairness, bottleneck fairness, envy freeness, and sharing incentive properties. It also satisfies Pareto optimality and strategy proofness in case of TDM.

The proof is given in the appendix. Unfortunately PS-DSF does not satisfy Pareto optimality in case of RDM. This is the reason why strategy proofness is *not generally satisfied* in case of RDM. The following lemma describes the behaviour of PS-DSF allocation mechanism from this respect.

Lemma 1. Assume that all users demand all type of resources, that is $d_{n,r} > 0$, $\forall n, r$. Under the PS-DSF allocation mechanism with RDM, each user cannot decrease the utilization of other users by lying about its resource demands or the set of eligible servers, without decreasing its own utilization.

For the proof refer to the appendix.

C. PS-DSF Allocation Algorithm

In this subsection, we present an algorithm which realizes the PS-DSF allocation in the case of RDM⁴. According to Theorem 1, an allocation satisfies PS-DSF when every user has a bottleneck resource w.r.t. every eligible server. Let \mathcal{N}_i denote the set of users for which $\gamma_{n,i} > 0$. The following corollary describes a condition to check whether a saturated resource serves as a bottleneck for user n w.r.t. server i .

Corollary 1. If r is saturated at server i and

$$n \in \arg \min_{m \in \mathcal{N}_i} \left\{ \frac{s_{m,i}}{\phi_m} \mid d_{m,r} > 0 \right\} \quad (14)$$

Then, r is a bottleneck resource for user n at server i when:

$$\frac{s_{m,i}}{\phi_m} > \frac{s_{n,i}}{\phi_n}, \quad d_{m,r} > 0 \Rightarrow x_{m,i} = 0. \quad (15)$$

To find a PS-DSF allocation, we may apply an iterative algorithm beginning with an initial allocation. Assume that servers are indexed from 1 to K . Starting from the first server, the proposed algorithm sequentially updates the allocation for different servers, so that at the end a bottleneck resource is identified for every user w.r.t. every eligible server. In the following we describe the procedure for updating the allocation at each server.

Specifically, for each server i let \mathcal{N}_i initially denote the set of users for which $\gamma_{n,i} > 0$. Given a feasible allocation, $\{x_{n,i}\}$, find S_i^* as the *minimum VDS* at server i :

$$S_i^* := \min_{m \in \mathcal{N}_i} \left\{ \frac{s_{m,i}}{\phi_m} \right\}. \quad (16)$$

The set of users achieving the minimum in (16) is denoted by \mathcal{N}_i^* . Let \mathcal{R}_i^* denote the set of saturated resources at server i

⁴A simplified version of this algorithm can be used in the case of TDM.

for which $d_{n,r} > 0$ for some user $n \in \mathcal{N}_i^*$. These resources are the potential bottleneck resources for users $n \in \mathcal{N}_i^*$. If the condition in Corollary 1 is satisfied for some resource $r^* \in \mathcal{R}_i^*$, then this resource serves as the bottleneck for users $n \in \mathcal{N}_i^*$ with $d_{n,r^*} > 0$. In this case, we restrict our attention to the users for which no bottleneck resource is identified w.r.t. server i . Specifically, \mathcal{N}_i is updated to:

$$\mathcal{N}_i = \mathcal{N}_i - \{n \mid d_{n,r^*} > 0\}. \quad (17)$$

When the condition in Corollary 1 is not satisfied for any resource $r \in \mathcal{R}_i^*$, the algorithm updates the allocation for server i . Specifically, for every resource $r \in \mathcal{R}_i^*$, a user n_r is chosen such that:

$$n_r \in \arg \max_{n \in \mathcal{N}_i} \left\{ \frac{s_{n,i}}{\phi_n} \mid x_{n,i} d_{n,r} > 0 \right\}. \quad (18)$$

If we release the whole allocated resources to these users from server i , the maximum potential increase in S_i^* is given by z^* (see the *Update-Allocation* subroutine in Algorithm II). To make sure that S_i^* is monotonically increasing, $\beta \in (0, 1]$ is chosen such that $S_i^* + \beta z^*$ remains less than or equal to the updated VDS w.r.t. server i for all users n_r , $r \in \mathcal{R}_i^*$.

For each server i , the above procedure is repeated until \mathcal{N}_i becomes empty. At the end of this procedure, a bottleneck resource is identified for every user eligible to be served by server i . However, the subsequent updates for the next servers, may violate this condition for server i and the previous servers. Hence, we repeat the whole process for all servers, until no more update is possible for any of the servers⁵. This process is described in Algorithm I.

D. Distributed Implementation

One of the advantages of the PS-DSF allocation mechanism is that it locally identifies the dominant resource for every user w.r.t. each server, without any knowledge of the available resources on the other servers, as opposed to existing allocation mechanisms which need to *globally* identify dominant resource and/or dominant share for each user. This is of great importance from a practical point of view, as we may develop a distributed algorithm to find the PS-DSF allocation.

Specifically, consider the inner while-loop in the main sub-routine of Algorithm I which we refer to as “*server procedure*”. According to this procedure the allocated tasks to different users from each server i are updated only based on the knowledge of the available resources on server i and the total allocated tasks to each user. Accordingly, we may come up with a distributed version of Algorithm I where each server individually (and even asynchronously) executes the server procedure every T seconds. When T is chosen sufficiently smaller than period of changes in a cluster (like changes in the set of active users and/or servers), such distributed algorithm may dynamically achieve the PS-DSF allocation. We implement this algorithm in our experiments in Section V.

⁵Convergence properties of this algorithm will be studied in our future work.

Algorithm I: PS-DSF Allocation Algorithm

Initialization

Initially allocate available resources by applying DRF individually to each server.

The main subroutine

```

while (1)
  Last-round-flag := 1
  for ( $i = 1; i \leq K; i++$ )
     $\mathcal{N}_i := \{n \in \mathcal{N} \mid \gamma_{n,i} > 0\}$ .
    while ( $\mathcal{N}_i \neq \emptyset$ )
      Find  $S_i^*$  according to (16).
      Identify  $\mathcal{N}_i^*$  as the set of users achieving the minimum
      in (16).
      Identify  $\mathcal{R}_i^*$  as the set of saturated resources at server  $i$ 
      for which
         $d_{n,r} > 0$  for some  $n \in \mathcal{N}_i^*$ .
      If ( $S_i^* = \max_{n \in \mathcal{N}_i^*} \{ \frac{s_{n,i}}{\phi_n} \mid x_{n,i} d_{n,r^*} > 0 \}$ , for  $r^* \in \mathcal{R}_i^*$ )
        Update  $\mathcal{N}_i = \mathcal{N}_i - \{n \mid d_{n,r^*} > 0\}$ 
      else
        Last-round-flag = 0
        Call Update-Allocation( $\mathbf{x}, i$ ).
      If (Last-round-flag = 1)
        break;
  Update-Allocation( $\mathbf{x}, i$ ) subroutine
    Identify  $\mathbf{f}_i = [f_{i,r}]$  as the amount of unallocated resources under  $\mathbf{x}$ .
    for ( $r \in \mathcal{R}_i^*$ )
      Choose  $n_r \in \arg \max_{n \in \mathcal{N}_i^*} \{ \frac{s_{n,i}}{\phi_n} \mid x_{n,i} d_{n,r} > 0 \}$ .
      Update  $\mathbf{f}_i = \mathbf{f}_i + x_{n_r,i} \mathbf{d}_n$ .

    Find  $D_i^* := \sum_{n \in \mathcal{N}_i^*} \phi_n \gamma_{n,i} \mathbf{d}_n$ .
    Find  $z^* := \min_r \frac{f_{i,r}}{D_{i,r}^*}$ .
    Choose  $\beta \in (0, 1]$  such that:  $S_i^* + \beta z^* \leq \frac{x_{n_r} - \beta x_{n_r,i}}{\phi_{n_r} \gamma_{n_r,i}}, \forall r \in \mathcal{R}_i^*$ .
    Update  $x_{n,i} = x_{n,i} + \beta \phi_n \gamma_{n,i} z^*, \forall n \in \mathcal{N}_i^*$ .
    Set  $x_{n_r,i} = (1 - \beta) x_{n_r,i}, \forall r \in \mathcal{R}_i^*$ .

```

IV. EXTENSIONS

In this section we present an extension which directly follows from our proposed approach in Section III. Specifically, consider the case where the effective capacity of resources on each server may vary for different users. In this case, that is unclear how to define a global dominant resource for each user even when there are no placement constraints. However, according to the formulation in Section III, we may define $\gamma_{n,i}$ as the number of tasks which could be executed by user n when monopolizing server i . We may also define the VDS for every user w.r.t. each server in the same way, and then find an allocation which satisfies PS-DSF. To gain more intuition, we consider two specific example scenarios in the following.

Example scenario 1: First consider a simple scenario where only one type of resource, notably bandwidth, is available on different servers. Specifically, we may consider different servers as different frequency channels which are subject to multi-user diversity in a wireless system. For instance, consider the example in Figure 4 where two users share three wireless channels. Without the insight of our proposed approach, that is unclear how to allocate the capacity of servers

among different users in a fair manner, as we may not weigh different servers with respect to each other.

Let define the utility of each user, x_n , as the number of bits which are given service in one second (i.e. the service rate). Also define $\gamma_{n,i}$ as the achievable service rate by user n when monopolizing server i . For the example in Figure 4 PS-DSF results in allocating the first channel (the third channel respectively) to the first user (second user), while the second channel is equally shared between the two users. Accordingly, user 1 gets a service rate of 1.5Mb/s while user 2 gets 1Mb/s. It can be observed that x_n can not be increased for any user without decreasing $x_{m,i}$ while $x_m/\gamma_{m,i} \leq x_n/\gamma_{n,i}$.

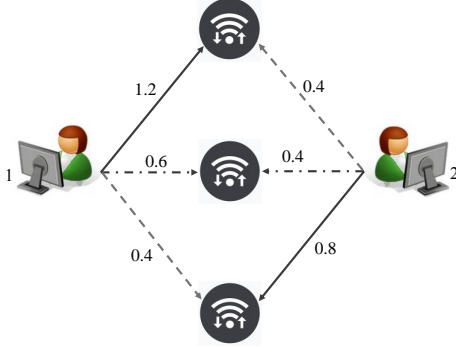


Fig. 4: An example with two equally weighted users sharing three frequency channels. The achievable service rates by each user over different channels are shown beside the arrows (in Mb/s).

Example scenario 2: Consider a set of heterogeneous servers in a computing cluster, where each server consists of different types of resources, such as CPU, RAM, bandwidth, etc. Although the CPU on each server has a fixed physical capacity, different users may experience different effective processing capacities when specific co-processors are available at a server. Coprocessors are supplementary processing units which are specialized for specific arithmetics or other processing purposes. As a result, they might be useful only for some users, for which they accelerate processing performance.

Our approach in Section III could be readily extended to incorporate the effect of coprocessors. Assume that $\mathbf{d}_n = [d_{n,r}]$ describes the demand of user n from each type of resource when no co-processor is utilized. Let $\gamma_{n,i}$ denote the maximum number of tasks which could be executed by user n when monopolizing server i and utilizing any available co-processor. Given the insight of PS-DSF, we may find an allocation, $\{x_{n,i}\}$, such that for every user, x_n cannot be increased while maintaining feasibility without decreasing $x_{m,i}$ for some user m and sever i with $x_m/\phi_m\gamma_{m,i} \leq x_n/\phi_n\gamma_{n,i}$. This problem will be studied in more details in our future work.

V. NUMERICAL RESULTS

In this section we evaluate performance of the PS-DSF allocation mechanism through some numerical experiments. In our simulations, we consider a cluster with four different

classes of servers (120 servers in total), where the configuration of servers are drawn from the distribution of Google cluster servers [18]. It is assumed that the available resources over each server can be partitioned in any arbitrary way. We consider four users where the last two users may run their tasks only by the last two classes of servers (see Figure 5).

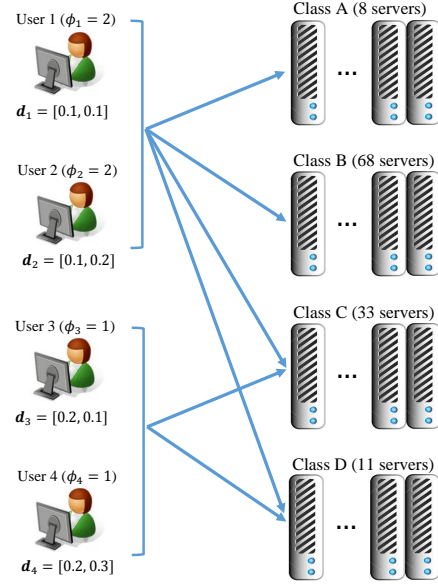


Fig. 5: A cluster with four classes of servers (120 servers in total) and four users. The weight of the first two users is twice the weight of the last two users. The configurations of resources (CPU and memory respectively) for servers of each class are as follows: $C_A = [1, 1]$, $C_B = [0.5, 0.5]$, $C_C = [0.5, 0.25]$, $C_D = [0.5, 0.75]$, where CPU and memory units for each server are normalized w.r.t. the servers of the first class.

The number of tasks that each user may run when monopolizing each class of servers are given in Table III. Assume that all users are active. The PS-DSF (based on RDM) and the TSF allocations in this case are given in Table IV. Under both allocations the servers of the first two classes (the second two classes respectively) are allocated to the first (the last) two users. According to the PS-DSF allocation, the servers of the third class (the fourth class respectively) are entirely allocated to the third user (the fourth user), which results in maximizing the minimum VDS w.r.t. these servers.

Intuitively, PS-DSF tries to allocate each server to the most efficient users. Therefore, we expect that PS-DSF results in greater utilization for different resources of a server compared to other allocation mechanisms such as TSF and C-DRFH. To observe this, we have executed these algorithms over the interval (0, 300) sec for the cluster in Figure 5. For the PS-DSF, we start with an initial allocation and update the allocation every second according to the *servers' procedure* (see our discussions in Section III-D on distributed implementation). For TSF and C-DRFH mechanisms we precisely find these allocations every second.

It is assumed that all users except User 4 are continuously active during the simulation interval. User 4 is inactive during interval (100, 250) sec, and is active elsewhere. The utilization

TABLE III: The total number of tasks that each user may run when monopolizing each class of servers.

$\gamma_{n,i}$	Class A	Class B	Class C	Class D
User 1	80	340	82.5	55
User 2	40	170	41.25	41.25
User 3	0	0	82.5	27.5
User 4	0	0	27.5	27.5

TABLE IV: The total number of tasks allocated to each user from each class of servers under PS-DSF and TSF allocations.

PS-DSF	Class A	Class B	Class C	Class D
User 1	40	170	0	0
User 2	20	85	0	0
User 3	0	0	82.5	0
User 4	0	0	0	27.5

TSF	Class A	Class B	Class C	Class D
User 1	35	170	0	0
User 2	22.5	85	0	0
User 3	0	0	58.33	0
User 4	0	0	8.05	27.5

that is achieved under any of these allocation mechanisms for the CPU at the third and the fourth classes of servers are shown respectively in Figure 6 (The CPU on the first two classes of servers and also the memory on all servers are fully utilized under any of the allocation mechanisms). It can be observed that the PS-DSF allocation mechanism results in greater utilization compared to the two other mechanisms in this example. Furthermore, it may be observed that the distributed version of the PS-DSF allocation algorithm promptly converges when changes occur in the set of active users.

VI. CONCLUSION

In summary, we studied the problem of multi-resource fair allocation for heterogeneous servers while respecting placement constraints. We identified important shortcomings in existing multi-resource fair allocation mechanisms when used in such environments. Hence, we proposed a new allocation mechanism, called PS-DSF. We discussed how our proposed allocation mechanism achieves different sharing properties which are satisfied under DRF in the case of one resource-pool/server. Furthermore, we discussed how PS-DSF could be implemented in a distributed manner. The performance of the PS-DSF allocation mechanism was compared against the existing allocation mechanisms and its enhanced performance was demonstrated through the numerical experiments. Further studies are under way and they will appear in future work.

REFERENCES

- [1] V. Chudnovsky, R. Rifaat, J. Hellerstein, B. Sharma, and C. Das, "Modeling and synthesizing task placement constraints in google compute clusters," in *Symposium on Cloud Computing*, 2011.
- [2] 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.
- [3] 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.
- [4] E. Friedman, A. Ghodsi, and C.-A. Psomas, "Strategyproof allocation of discrete jobs on multiple machines," in *Proceedings of the ACM conference on Economics and computation*. ACM, 2014, pp. 529–546.

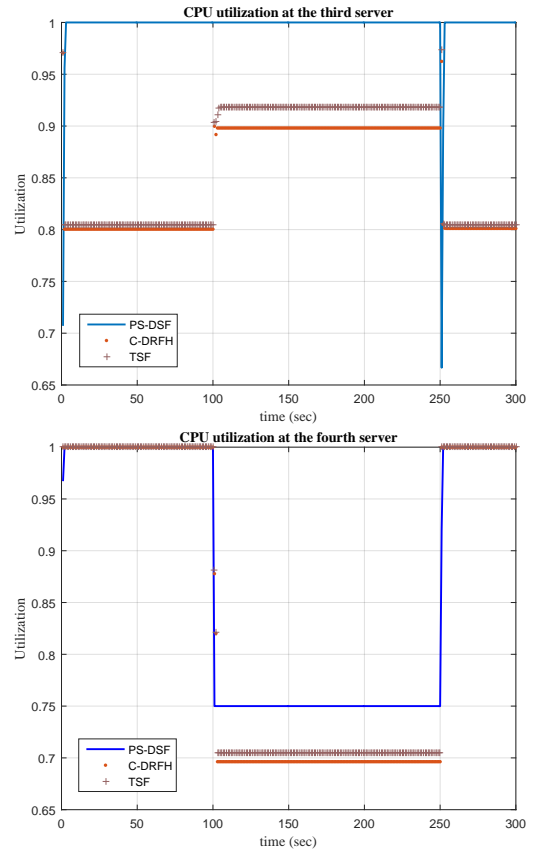


Fig. 6: The utilization that is achieved for the CPU at the third and the fourth classes of servers under PS-DSF, TSF and C-DRFH allocation mechanisms.

- [5] 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.
- [6] M. Chowdhury, Z. Liu, A. Ghodsi, and I. Stoica, "Hug: Multi-resource fairness for correlated and elastic demands," in *Proc. NSDI*, Mar 2016.
- [7] W. Wang, B. Liang, and B. Li, "Multi-resource fair allocation in heterogeneous cloud computing systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 10, pp. 2822–2835, Oct 2015.
- [8] 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.
- [9] J. Khamse-Ashari, I. Lambadaris, and Y. Q. Zhao, "Constrained multi-user multi-server max-min fair queuing," <http://arxiv.org/abs/1601.04749>, 2016. [Online]. Available: <http://arxiv.org/abs/1601.04749>
- [10] J. Khamse-Ashari, G. Kesidis, I. Lambadaris, B. Urgaonkar, and Y. Zhao, "Max-min fair scheduling of variable-length packet-flows to multiple servers by deficit round-robin," in *Proceedings of the CISS, Princeton*, Mar 2016.
- [11] —, "Constrained max-min fair scheduling of variable-length packet-flows to multiple servers," in *Proc. IEEE Globecom*, Dec 2016.
- [12] —, "Efficient and Fair Scheduling of Placement Constrained Threads on Heterogeneous Multi-Processors," in *preprint*, Sept. 2016.
- [13] 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.
- [14] W. Wang, B. Li, B. Liang, and J. Li, "Multi-resource fair sharing for datacenter jobs with placement constraints," *SC 2016*.
- [15] —, "Towards multi-resource fair allocation with placement constraints," in *Proc. ACM SIGMETRICS*, Antibes, France, 2016.
- [16] G. Kesidis, Y. Wang, B. Urgaonkar, J. Khamse-Ashari, and I. Lambadaris, "Fair Scheduling of Multiple Resource Types over Multiple and Heterogeneous Resource Pools," CSE Dept, PSU, Tech. Rep. CSE-16-

- [17] D. Parkes, A. Procaccia, and N. Shah, "Beyond dominant resource fairness: Extensions, limitations, and indivisibilities," in *Proc. ACM EC*, Valencia, Spain, June 2012.
- [18] C. Reiss, J. Wilkes, and J. L. Hellerstein, "Google cluster-usage traces," 2011, <http://code.google.com/p/googleclusterdata/>.

APPENDIX

Proof of Theorem 1. First consider a feasible allocation $\{x_{m,i}\}$ for which there exists a bottleneck resource for every user w.r.t. every eligible server. Let $b(n,i)$ denote the bottleneck resource for user n w.r.t. server i . That is, $b(n,i)$ is saturated, $d_{n,b(n,i)} > 0$, and

$$\frac{s_{n,i}}{\phi_n} \geq \frac{s_{m,i}}{\phi_m}, \quad \forall m \text{ such that } x_{m,i} d_{m,b(n,i)} > 0. \quad (19)$$

Given that $b(n,i)$ is saturated, it is not possible to increase $x_{n,i}$, unless decreasing $x_{m,i}$ for some user m with $x_{m,i} d_{m,b(n,i)} > 0$. On the other hand, (19) implies that $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$ for any user m with $x_{m,i} d_{m,b(n,i)} > 0$. Hence, we may not increase the allocated tasks to user n from any server i unless decreasing $x_{m,i}$ for some user m with $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$. This implies that $\{x_{m,i}\}$ satisfies PS-DSF.

Now consider an allocation $\{x_{m,j}\}$ which satisfies PS-DSF. Let $\mathcal{R}_{n,i}$ denote the set of demanded resources by user n which are saturated at an eligible server i under the allocation $\{x_{m,j}\}$, that is:

$$\mathcal{R}_{n,i} := \{r \mid d_{n,r} > 0 \text{ and } \sum_m x_{m,i} d_{m,r} = c_{i,r}\}. \quad (20)$$

This set includes the *potential* bottleneck resources for user n w.r.t. server i . First we prove that $\mathcal{R}_{n,i}$ may not be empty under a PS-DSF allocation. By contradiction, assume that $\mathcal{R}_{n,i} = \emptyset$, that is none of the demanded resources by user n are saturated at server i . In this case, we may increase $x_{n,i}$ by:

$$z_{n,i} := \min_{r: d_{n,r} > 0} \frac{c_{i,r} - \sum_m x_{m,i} d_{m,r}}{d_{n,r}} > 0, \quad (21)$$

without decreasing $x_{m,i}$ for any user m . However, this contradicts to the fact that $\{x_{m,j}\}$ satisfies PS-DSF.

Next, we show that there exists some resource $r \in \mathcal{R}_{n,i}$ which serves as a bottleneck for user n w.r.t. server i . By contradiction, assume that none of the resources in $\mathcal{R}_{n,i}$ is a bottleneck. That is, for any resource $r \in \mathcal{R}_{n,i}$ we may find some user p with $x_{p,i} d_{p,r} > 0$ such that $s_{n,i}/\phi_n < s_{p,i}/\phi_p$. Hence, we can increase $x_{n,i}$ by decreasing $x_{p,i}$ for some user(s) p with $s_{n,i}/\phi_n < s_{p,i}/\phi_p$. That is, $x_{n,i}$ could be increased without decreasing $x_{m,i}$ for any user m with $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$. However, this contradicts to the fact that $\{x_{m,j}\}$ satisfies PS-DSF. \square

Proof of Theorem 2. Consider a feasible allocation $\{x_{n,i}\}$ for which (10) holds with equality, and the condition in (13) is established. Since (10) holds with equality, it is not possible to increase $x_{n,i}$, unless decreasing $x_{m,i}$ for some user m with $x_{m,i} > 0$. On the other hand, (13) implies that $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$ for any user m with $x_{m,i} > 0$. Therefore, we may

not increase the allocated tasks to any user n from any server i unless decreasing $x_{m,i}$ for some user m with $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$. This implies that $\{x_{n,i}\}$ satisfies PS-DSF.

Now assume that $\{x_{n,i}\}$ satisfies PS-DSF. By contradiction, assume that (10) holds with inequality for some server i . In this case we may increase $x_{n,i}$ by:

$$z_{n,i} = \gamma_{n,i} \left[1 - \sum_m \frac{x_{m,i}}{\gamma_{m,i}} \right] > 0, \quad (22)$$

without decreasing $x_{m,i}$ for any user m . However, this contradicts to the fact that $\{x_{n,i}\}$ satisfies PS-DSF. Next, we show that (13) is established under the PS-DSF allocation. By contradiction, assume that we can find some user p with $x_{p,i} > 0$ such that $s_{p,i}/\phi_p > s_{n,i}/\phi_n$. Hence, we can increase $x_{n,i}$ by decreasing $x_{p,i}$ for user p with $s_{p,i}/\phi_p > s_{n,i}/\phi_n$. That is, $x_{n,i}$ could be increased without decreasing $x_{m,i}$ for any user m with $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$. However, this contradicts to the fact that $\{x_{n,i}\}$ satisfies PS-DSF. \square

Proof of Theorem 3. We prove single resource fairness, bottleneck fairness, envy freeness, and sharing incentive properties for the more complicated case of RDM. The proofs of these properties follow the same line of arguments in case of TDM, so we do not repeat them here.

Single resource fairness: When there is only one type of resource, then $d_n = d_{n,1}$, $\forall n$ and $\gamma_{n,i} = \delta_{n,i} c_{i,1} / d_{n,1}$, $\forall n, i$. As a result:

$$s_{n,i} = \frac{x_n}{\gamma_{n,i}} = \frac{x_n d_{n,1}}{c_{i,1} \delta_{n,i}} = \frac{a_n}{c_{i,1} \delta_{n,i}}, \quad (23)$$

where a_n is the allocated resource to user n from all servers. According to the PS-DSF allocation, we may not increase x_n (or equivalently a_n) while maintaining feasibility without decreasing $x_{m,i}$ for some user m with $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$ (or $a_m/\phi_m \leq a_n/\phi_n$). Therefore, the allocated resource to different users, $\{a_n\}$ satisfies (constrained) weighted max-min fairness.

Bottleneck fairness: Assume that there is one resource, say r^* , which is dominantly requested by every user from every eligible server. By definition, r^* is considered as the *dominant resource* for every user n w.r.t. every eligible server i . Accordingly, $\gamma_{n,i}$ is given by $\gamma_{n,i} = \delta_{n,i} c_{i,r^*} / d_{n,r^*}$, and the VDS for user n w.r.t. server i is given by:

$$s_{n,i} = \frac{x_n}{\gamma_{n,i}} = \frac{x_n d_{n,r^*}}{c_{i,r^*} \delta_{n,i}} = \frac{a_{n,r^*}}{c_{i,r^*} \delta_{n,i}}. \quad (24)$$

where a_{n,r^*} is the amount of the bottleneck resource allocated to user n from all servers. According to the PS-DSF allocation, we may not increase x_n (or equivalently a_{n,r^*}) while maintaining feasibility without decreasing $x_{m,i}$ for some user m with $s_{m,i}/\phi_m \leq s_{n,i}/\phi_n$ (or $a_{m,r^*}/\phi_m \leq a_{n,r^*}/\phi_n$). Hence, the allocated bottleneck resource to different users, $\{a_{n,r^*}\}$ satisfies (constrained) weighted max-min fairness.

Envy freeness: Given $U_n(\mathbf{a})$ as the utility function for user n , $U_n(\mathbf{a}_n \phi_n / \phi_m)$ gives the utility of user n of the allocated

resources to user m (i.e., $\mathbf{a}_m = x_m \mathbf{d}_m$), when adjusted according to their weights. According to (1):

$$U_n\left(\frac{\phi_n}{\phi_m} \mathbf{a}_m\right) = \frac{\phi_n}{\phi_m} \min_r \frac{a_{m,r}}{d_{n,r}} = \frac{\phi_n x_m}{\phi_m} \min_r \frac{d_{m,r}}{d_{n,r}}. \quad (25)$$

We show that:

$$\min_r \frac{d_{m,r}}{d_{n,r}} \leq \frac{d_{m,\rho(n,i)}}{d_{n,\rho(n,i)}} = \frac{\gamma_{n,i} d_{m,\rho(n,i)}}{c_{i,\rho(n,i)}} \leq \frac{\gamma_{n,i}}{\gamma_{m,i}}, \quad \forall i. \quad (26)$$

As a result:

$$U_n\left(\frac{\phi_n}{\phi_m} \mathbf{a}_m\right) \leq x_m \frac{\phi_n \gamma_{n,i}}{\phi_m \gamma_{m,i}}, \quad \forall i. \quad (27)$$

According to Theorem 1, there exists a bottleneck resource for every user w.r.t. every eligible server. Consider server i for which $x_{m,i} > 0$. Let $b(n,i)$ denote the bottleneck resource for user n w.r.t. server i . For $U_n(\mathbf{a}_m \phi_n / \phi_m)$ to be greater than zero (see (25)), we need $d_{m,b(n,i)} > 0$. Given that $b(n,i)$ is the bottleneck for user n and $x_{m,i} d_{m,b(n,i)} > 0$, it follows that:

$$\frac{x_m}{\phi_m \gamma_{m,i}} \leq \frac{x_n}{\phi_n \gamma_{n,i}}. \quad (28)$$

This along with (27) results in:

$$U_n\left(\frac{\phi_n}{\phi_m} \mathbf{a}_m\right) \leq x_n. \quad (29)$$

Sharing Incentive: Without loss of generality assume that the demand vector for every user n , \mathbf{d}_n , is normalized by $\sum_i \gamma_{n,i}$ (the number of tasks which could be executed by user n if the whole system is allocated to it). In this case, x_n and $\gamma_{n,i}$ will be normalized by the same factor and the VDS for user n w.r.t. different servers and also the resulting PS-DSF allocation won't be changed. Specifically, define:

$$\hat{x}_n := \frac{x_n}{\sum_j \gamma_{n,j}} \quad (30)$$

$$\hat{\gamma}_{n,i} := \frac{\gamma_{n,i}}{\sum_j \gamma_{n,j}} \quad (31)$$

For the uniform allocation:

$$\hat{x}_n^{unif} = \frac{\phi_n}{\sum_m \phi_m} \sum_i \hat{\gamma}_{n,i} = \frac{\phi_n}{\sum_m \phi_m},$$

where the second equality follows from the fact that $\sum_i \hat{\gamma}_{n,i} = 1$. We assert that \hat{x}_n / ϕ_n is greater than or equal to $1 / \sum_m \phi_m$ for all users under the PS-DSF allocation.

The proof is by induction on the number of users, N . Specifically, for $N = 2$ consider two users, n and m . To consider the worst-case, assume that both users have the same bottleneck w.r.t. each server. Assume that servers are indexed in increasing order of $\hat{\gamma}_{n,j} / \hat{\gamma}_{m,j}$. Let j_0 denote the least indexed server from which some tasks are allocated to user n under the PS-DSF allocation, that is $\hat{x}_{n,j_0} > 0$. Given that $x_{n,j_0} > 0$ and both users have the same bottleneck w.r.t. server j_0 , it follows that (see Definition 6):

$$s_{m,j_0} / \phi_m \geq s_{n,j_0} / \phi_n. \quad (32)$$

Without loss of generality assume that $\hat{\gamma}_{n,j} / \hat{\gamma}_{m,j} > \hat{\gamma}_{n,j_0} / \hat{\gamma}_{m,j_0}$, for $j > j_0$. For these servers it follows that $s_{m,j} / \phi_m > s_{n,j} / \phi_n$. This along with the assumption that user m has the same bottleneck as user n imply that $\hat{x}_{m,j} = 0$ for $j > j_0$. Therefore, server j_0 is the only server for which $\hat{x}_{n,j} \hat{x}_{m,j}$ could be greater than zero. Accordingly:

$$\hat{x}_n = \alpha_{n,j_0} \hat{\gamma}_{n,j_0} + \sum_{j=j_0+1}^K \hat{\gamma}_{n,j}, \quad (33)$$

$$\hat{x}_m = \alpha_{m,j_0} \hat{\gamma}_{m,j_0} + \sum_{j=1}^{j_0-1} \hat{\gamma}_{m,j}, \quad (34)$$

where α_{n,j_0} (α_{m,j_0} respectively) denotes the portion of the DR for user n (user m) w.r.t. server j_0 that is allocated to it under PS-DSF. Substituting \hat{x}_n and \hat{x}_m from (33) and (34) into (32) results in:

$$\begin{aligned} \frac{\alpha_{m,j_0}}{\phi_m} + \sum_{j=1}^{j_0-1} \frac{\hat{\gamma}_{m,j}}{\phi_m \hat{\gamma}_{m,j_0}} &\geq \frac{\alpha_{n,j_0}}{\phi_n} + \sum_{j=j_0+1}^K \frac{\hat{\gamma}_{n,j}}{\phi_n \hat{\gamma}_{n,j_0}} \\ &\geq \frac{1 - \alpha_{m,j_0}}{\phi_n} + \sum_{j=j_0+1}^K \frac{\hat{\gamma}_{n,j}}{\phi_n \hat{\gamma}_{n,j_0}}, \end{aligned}$$

where the second inequality follows from the fact that $\alpha_{m,j_0} + \alpha_{n,j_0} \geq 1$. After some manipulations, it follows that $\alpha_{m,j_0} \geq (A + \phi_m) / (\phi_m + \phi_n)$, where:

$$A := \sum_{j=j_0+1}^K \phi_m \frac{\hat{\gamma}_{n,j}}{\hat{\gamma}_{n,j_0}} - \sum_{j=1}^{j_0-1} \phi_n \frac{\hat{\gamma}_{m,j}}{\hat{\gamma}_{m,j_0}}. \quad (35)$$

Applying the lower bound of α_{m,j_0} into (34) and after some manipulations, it follows that:

$$\begin{aligned} \frac{\hat{x}_m}{\phi_m} &\geq \frac{\sum_{j=1}^{j_0} \hat{\gamma}_{m,j} + \sum_{j=j_0+1}^K \frac{\hat{\gamma}_{n,j}}{\hat{\gamma}_{n,j_0}} \hat{\gamma}_{m,j_0}}{\phi_m + \phi_n} \\ &\geq \frac{\sum_{j=1}^K \hat{\gamma}_{m,j}}{\phi_m + \phi_n} = \frac{1}{\phi_m + \phi_n} \end{aligned} \quad (36)$$

where the second inequality follows from the fact that $\hat{\gamma}_{n,j} / \hat{\gamma}_{m,j} \geq \hat{\gamma}_{n,j_0} / \hat{\gamma}_{m,j_0}$, $j \geq j_0$, and the last equality follows from the fact that $\sum_j \hat{\gamma}_{m,j} = 1$. The lower bound in (36) could be shown for \hat{x}_n / ϕ_n in the same way.

Assume that the statement is established for $N = N_0$ users. For the case that the set of all users, \mathcal{N} , consists of $N_0 + 1$ users, we may assume that \mathcal{N} is comprised of a subset \mathcal{N}_0 of N_0 users with the total weight of $\Phi_0 := \sum_{n \in \mathcal{N}_0} \phi_n$ and a singular user n_0 with the weight of ϕ_0 . We assume that user n_0 is chosen arbitrarily. We may consider $\phi_n / (\Phi_0 + \phi_0)$ portion of the resources on every server as the share of each user n . Assume that user n_0 does not share its resources with others. In this case, $\hat{x}_{n_0} = \phi_0 / (\Phi_0 + \phi_0)$.

User n_0 would prefer to exchange all or part of its allocated resources from server j with all or part of the allocated resources to user m from server i , if

$$\frac{\hat{\gamma}_{n_0,i}}{\hat{\gamma}_{m,i}} > \frac{\hat{\gamma}_{n_0,j}}{\hat{\gamma}_{m,j}}. \quad (37)$$

In this case, the number of allocated tasks to both of them could be increased compared to the generic uniform allocation. We may repeat the same process, exchanging the allocated resources to user n_0 by that for other users, until no more exchange is possible and \hat{x}_{n_0} cannot be further increased. After that, we may freeze the allocated resources to user n_0 and allocate the remaining resources among other users according to PS-DSF. Given that sharing incentive is provided by PS-DSF for the set \mathcal{N}_0 with N_0 users, it follows that:

$$\hat{x}_n \geq \frac{\Phi_0}{\Phi_0 + \phi_0} \frac{\phi_n}{\Phi_0} = \frac{\phi_n}{\Phi_0 + \phi_0}, \quad \forall n \in \mathcal{N}_0. \quad (38)$$

If we allocate the whole resources among all users $n \in \mathcal{N}$ according to PS-DSF allocation, the number of allocated tasks to users $n \in \mathcal{N}_0$ may not be decreased compared to the above-described allocation (because there is no reservation for user n_0 in this case). That is, (38) is established under the PS-DSF allocation. Since n_0 is chosen arbitrarily, we may repeat the same discussions by choosing a different set \mathcal{N}'_0 which includes n_0 . Hence, we may conclude that the lower bound in (38) is established for all users.

Pareto optimality: Consider an allocation, $\{x_{n,i}\}$, satisfying PS-DSF based on TDM. For such an allocation, Theorem 2 implies that (10) holds with equality for each server i . Hence, we may not increase $x_{n,i}$ without decreasing $x_{m,i}$ for some user m with $x_{m,i} > 0$. Furthermore, according to Theorem 2, for any user m and server i with $x_{m,i} > 0$:

$$x_{m,i} > 0 \Rightarrow s_{m,i}/\phi_m = \min_n s_{n,i}/\phi_n. \quad (39)$$

In fact, PS-DSF allocation mechanism maximizes x_m for each user m subject to (39). The following lemma shows that this condition is not restricting, as we may not increase x_n without decreasing x_m for some user m , even when violating the condition in (39). This means that PS-DSF allocation is Pareto optimal in case of TDM.

Lemma 2. Assume that $\{x_{n,i}\}$ satisfies PS-DSF based on TDM. Consider two arbitrary users, n and m , for which $x_{n,i} > 0$ and $x_{m,j} > 0$. If user n exchanges all or part of its allocated tasks from server i with all or part of the allocated tasks to user m from server j , then the allocated tasks to at least one of them will be decreased compared to the PS-DSF allocation.

Proof. Given that $x_{m,j} > 0$ and $x_{n,i} > 0$, we may decrease $x_{m,j}$ and $x_{n,i}$, and increase $x_{n,j}$ and $x_{m,i}$. Let $\Delta x_{n,i}$, $\Delta x_{m,i}$ and $\Delta x_{n,j}$, $\Delta x_{m,j}$ denote a feasible change in the number of allocated tasks to users n and m while (10) holds with equality for both servers i and j . For (10) to hold with equality we have:

$$\Delta x_{n,j} = -\Delta x_{m,j} \frac{\gamma_{n,j}}{\gamma_{m,j}} \quad (40)$$

$$\Delta x_{m,i} = -\Delta x_{n,i} \frac{\gamma_{m,i}}{\gamma_{n,i}} \quad (41)$$

Assume that $\Delta x_{n,i} + \Delta x_{n,j} > 0$ or $-\Delta x_{n,i} < \Delta x_{n,j}$. This along with (40) and (41) results in:

$$\Delta x_{m,i} < \Delta x_{n,j} \frac{\gamma_{m,i}}{\gamma_{n,i}} \quad (42)$$

$$< -\Delta x_{m,j} \frac{\gamma_{n,j}}{\gamma_{m,j}} \frac{\gamma_{m,i}}{\gamma_{n,i}}. \quad (43)$$

The fact that $x_{m,j} > 0$ and $x_{n,i} > 0$ along with (13) result in:

$$\frac{\gamma_{m,i}}{\gamma_{n,i}} \leq \frac{\phi_n}{\phi_m} \frac{x_m}{x_n} \quad (44)$$

$$\frac{\gamma_{n,j}}{\gamma_{m,j}} \leq \frac{\phi_m}{\phi_n} \frac{x_n}{x_m} \quad (45)$$

Combining (43) with (44) and (45) results in $\Delta x_{m,i} < -\Delta x_{m,j}$. This means that the number of allocated tasks to user m is decreased compared to PS-DSF allocation. \square

Strategy proofness: Let $A := \{a_{m,i}\}$ (and $A' := \{a'_{m,i}\}$, respectively) denote the resulting PS-DSF allocation when user n trustfully declares \mathbf{d}_n and $\delta_n = [\delta_{n,i}]$ (non-trustfully declares \mathbf{d}'_n and δ'_n). Users other than n take the same actions in both cases (whether trustful or non-trustful). Hence, $\gamma'_{m,i} = \gamma_{m,i}$ for $m \neq n$. The number of tasks that user n may actually execute under the allocation A' (i.e., by using $\mathbf{a}'_n = x'_n \mathbf{d}'_n$) is given by:

$$U_n(\mathbf{a}'_n) = \min_r \frac{a'_{n,r}}{d_{n,r}} = x'_n \min_r \frac{d'_{n,r}}{d_{n,r}}. \quad (46)$$

As in (26), we can show that $\min_r \frac{d'_{n,r}}{d_{n,r}} \leq \frac{\gamma_{n,i}}{\gamma'_{n,i}}$, $\forall i$. Hence:

$$U_n(\mathbf{a}'_n) = x'_n \min_r \frac{d'_{n,r}}{d_{n,r}} \leq x'_n \frac{\gamma_{n,i}}{\gamma'_{n,i}}. \quad (47)$$

For strategy proofness we need to show that $U_n(\mathbf{a}'_n) \leq x_n$. By contradiction, assume that $U_n(\mathbf{a}'_n) > x_n$. It follows that:

$$s_{n,i} = \frac{x_n}{\gamma_{n,i}} < \frac{U_n(\mathbf{a}'_n)}{\gamma_{n,i}} \leq \frac{x'_n}{\gamma'_{n,i}} = s'_{n,i}, \quad \forall i. \quad (48)$$

That is, the VDS for user n is increased w.r.t. all servers under the allocation A' compared to the allocation A , provided that $U_n(\mathbf{a}'_n) > x_n$. Let \mathcal{U} denote the set of users for which $U_m(\mathbf{a}'_m) > x_m$. For users $m \in \mathcal{U}$, $m \neq n$, it follows that

$$s_{m,i} = \frac{x_m}{\gamma_{m,i}} < \frac{U_m(\mathbf{a}'_m)}{\gamma_{m,i}} = \frac{x'_m}{\gamma'_{m,i}} = s'_{m,i}, \quad \forall i. \quad (49)$$

We define:

$$s_i := \min_n \frac{s_{n,i}}{\phi_n} \quad (50)$$

as the Virtual Dominant Share Level, VDSL, at server i under the allocation A . In the same way, we define s'_i as the VDSL at server i under the allocation A' . For any user $m \in \mathcal{U}$, Theorem 2 implies that $s'_i = s'_{m,i}/\phi_m$ provided that $x'_{m,i} > 0$. It follows that:

$$s'_i = \frac{s'_{m,i}}{\phi_m} > \frac{s_{m,i}}{\phi_m} \geq s_i. \quad (51)$$

That is, the VDSL is increased at all servers for which $x'_{m,i} > 0$ for some $m \in \mathcal{U}$. Let define:

$$\mathcal{S} := \{i \mid x'_{m,i} > 0 \text{ for some } m \in \mathcal{U}\}. \quad (52)$$

Accordingly, no tasks are allocated under the allocation A' from servers $j \notin \mathcal{S}$ to users $m \in \mathcal{U}$, i.e., $x'_{m,j} = 0$, $m \in \mathcal{U}$, $j \notin \mathcal{S}$. Hence, VDSL at servers $j \notin \mathcal{S}$ may not be decreased under the allocation A' compared to allocation A , that is $s'_j \geq s_j$ for servers $j \notin \mathcal{S}$. Therefore, $s'_i \geq s_i$, $\forall i$, which in turn implies that $U_m(\mathbf{a}'_m) = x'_m \geq x_m$ for $m \neq n$. This along with the assumption that $U_n(\mathbf{a}'_n) > x_n$ contradict to Pareto optimality of the allocation A . \square

Proof of Lemma 1. The proof follows the same line of arguments as the proof of strategy proofness in case of TDM. Specifically, when all users demand all types of resources, the same resource serves as the bottleneck for all users w.r.t. each server. Hence, we may define the VDSL at each server i as in (50). With the same line of arguments we may conclude that the VDSL is not decreased at any server i under the allocation A' compared to allocation A , provided that $U_n(a'_n) \geq x_n$. That is, $s'_i \geq s_i$, which in turn implies that $U_m(\mathbf{a}'_m) \geq x_m$ $\forall m$. Therefore, user n may not decrease the utilization of other users, by lying about its resource demands or the set of eligible servers, unless decreasing its own utilization. \square