

Enhanced cluster computing performance through proportional fairness

Thomas Bonald^{a,*}, James Roberts^{b,1}

^a Télécom ParisTech, Paris, France

^b IRT-SystemX, Paris-Saclay, France

ARTICLE INFO

Article history:

Available online 10 July 2014

Keywords:

Cluster computing
Multi-resource sharing
Proportional fairness
Dominant resource fairness

ABSTRACT

The performance of cluster computing depends on how concurrent jobs share multiple data center resource types such as CPU, RAM and disk storage. Recent research has discussed efficiency and fairness requirements and identified a number of desirable scheduling objectives including so-called dominant resource fairness (DRF). We argue here that proportional fairness (PF), long recognized as a desirable objective in sharing network bandwidth between ongoing data transfers, is preferable to DRF. The superiority of PF is manifest under the realistic modeling assumption that the population of jobs in progress is a stochastic process. In random traffic the strategy-proof property of DRF proves unimportant while PF is shown by analysis and simulation to offer a significantly better efficiency–fairness tradeoff.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

We consider a cluster of servers, exclusively reserved for executing a certain class of jobs, as a set of distinct resource pools. The resources in question include CPU, RAM, disk space and I/O bandwidth, for instance. A job must execute a certain number of tasks, each task requiring a given quantity of each type of resource. An example of this type of computing environment and its workload are described by Reiss et al. [1]. Each task has its own particular requirements profile specifying how much of each resource must be reserved: some tasks require proportionally more CPU than RAM, others require proportionally more RAM than I/O bandwidth, etc. The issue at hand is how a central scheduler should initiate and run tasks in parallel to fairly and efficiently share cluster resources between jobs. We argue in this paper that such resource sharing should realize the objective of Proportional Fairness (PF), that is, it should maximize the sum over jobs of logarithms of their number of tasks.

The multi-resource sharing problem was formulated in a recent paper by Ghodsi et al. [2] under the simplifying assumption that resource pools are homogeneous and infinitely divisible between jobs. All tasks of a given job are assumed to have identical profiles and, rather than allocating resources to discrete tasks, a central scheduler is supposed to assign shares of each resource in proportion to the task profile. Sharing objectives are expressed in terms of a number of desirable properties. In particular, it is required that allocations be strategy-proof in the sense that the owner of a job cannot gain a bigger share by lying about its actual task profile. The authors of [2] show that this property is not satisfied by many otherwise intuitively appealing allocations. They proceed to define an original strategy-proof allocation called Dominant Resource Fairness (DRF) and explain how this can be realized.

* Corresponding author.

E-mail addresses: thomas.bonald@telecom-paristech.fr (T. Bonald), james.roberts@irt-systemx.fr (J. Roberts).

¹ The authors are members of the LINCOS, Paris, France. See www.lincos.fr.

Other authors have since come up with alternative strategies, based on different notions of fairness (e.g., [3,4]), or adapting DRF to account for more realistic cluster constraints (e.g., [5–7]). Gutman and Nisan [8] situate proposed fairness objectives like DRF [2] and “no justified complaints” [3] in a common economics framework and provide efficient polynomial time algorithms.

We believe the economics literature cited in [8] does not in fact constitute an appropriate background for modeling cluster resource sharing. This is because it completely ignores the dynamics of job arrivals and completions that characterize cluster workload [1]. These dynamics have an obvious impact on job completion times and interact closely with the applied allocation algorithm since a job finishes more or less quickly depending on the resources it is allocated. The notion of strategy-proofness must be revisited to account for the impact of false user requirements on the process of jobs in progress and their expected performance.

In advocating PF we draw on our understanding of bandwidth sharing performance in networks. PF was defined by Kelly et al. [9] as the bandwidth allocation that maximizes the sum of logs of flow rates, arguably thus realizing greater social welfare than alternative allocations like max–min fairness. We consider rather the impact of the sharing objective on flow completion times in dynamic traffic. In wired networks, it turns out that completion times are not highly dependent on the type of fairness imposed and DRF (equivalent here to max–min fairness) and PF have similar performance [10]. In wireless networks, on the other hand, the radio resource is measured in time slots per second rather than bit/s with the achievable bit rate per time slot depending significantly on user position. Sharing a channel equally in time, as realized in HDR/HSDPA systems for instance, actually corresponds to PF sharing in terms of bandwidth. In a mixed wired/wireless network, PF realizes a significantly more favorable efficiency–fairness tradeoff than max–min fairness [11].

To compare PF and DRF in dynamic traffic we adopt a simple Markovian traffic model assuming Poisson job arrivals and exponential job sizes. This simplifies analysis and clarifies the respective tradeoffs realized by the two allocation approaches. Simulation and analytical results confirm that PF performs significantly better than DRF. Given known insensitivity properties of fair resource sharing, we are confident that PF would be equally preferable under a more realistic traffic model [10].

In addition to discussing ideal sharing of infinitely divisible resources, Ghodsi et al. show how to approximately realize DRF in practice accounting for discrete, finite-size tasks. The scheduler algorithm preferentially launches tasks of the “most deprived” job. This is the job whose current share of its dominant resource is smallest. PF can be implemented similarly on re-defining “most deprived” in terms of the ideal shares determined for that allocation objective. The complexities of PF and DRF algorithms are similar and hardly constitute an implementation issue since the number of resource types to be shared is typically very small.

In the next section we define DRF and PF and demonstrate their respective sharing properties with respect to a static job population. We consider dynamic sharing in Section 3 assuming a fluid model where resources are infinitely divisible. The task-by-task implementations of DRF and PF are compared in Section 4. Finally, related work is discussed in Section 5 before we conclude.

2. Multi-resource sharing

We consider how multiple resources should be shared assuming a static population of jobs, each with a particular profile of per-task requirements. We adopt a fluid model where pools of resources are assumed infinitely divisible and compare two allocation strategies: dominant resource fairness (DRF) and proportional fairness (PF).

2.1. The fluid model

We consider J infinitely divisible pools of resources of respective capacities C_j , for $j = 1, \dots, J$, to be shared by n jobs indexed by i . Each task of job i requires A_{ij} units of resource j . Denoting by φ_i the number of ongoing tasks of job i , we have the capacity constraints:

$$\sum_{i=1}^n \varphi_i A_{ij} \leq C_j,$$

for $j = 1, \dots, J$. In the fluid model, we assume tasks are infinitesimally small and jobs can run a sufficiently large number of them to attain a given resource allocation. It then makes more sense to normalize resource capacities to 1 with $a_{ij} = A_{ij}/C_j$ representing fractional requirements. The (now) real numbers $\varphi_1, \dots, \varphi_n$ are then considered as *task volumes* satisfying capacity constraints:

$$\sum_{i=1}^n \varphi_i a_{ij} \leq 1, \tag{1}$$

for $j = 1, \dots, J$. The product $\varphi_i a_{ij}$ is the fraction of resource j allocated to job i . We also write capacity constraints (1) in matrix form $\varphi a \leq 1$, where inequality is understood component-wise.

We say that resource j is *saturated* if the corresponding capacity constraint is attained, i.e., if $\sum_i \varphi_i a_{ij} = 1$. Let a_i denote the vector whose j th component is a_{ij} . We say job i *needs* resource j if $a_{ij} > 0$ and we assume each job needs at least one resource. The *dominant resource* of job i is that for which the normalized requirement a_{ij} is largest.

2.2. Desirable properties

The following are desirable properties of a multi-resource allocation algorithm, as discussed by Ghodsi et al. [2]:

Pareto-efficiency. An algorithm is Pareto-efficient if each job needs some resource which is saturated.

Sharing-incentive. An algorithm satisfies sharing-incentive if each job gets no less than a $1/n$ share of its dominant resource.

Strategy-proof. An algorithm is strategy-proof if the resource shares of any job do not increase if this job modifies its vector of resource requirements.

Single-resource fairness. An algorithm is said to be single-resource fair if resource shares $\varphi_i a_{i1}$ are all equal to $1/n$ when $J = 1$.

It proves useful to introduce two additional properties:

Scale-invariance. An algorithm is scale-invariant if resource shares remain the same when the vector of resource requirements of any job is multiplied by a scaling factor.

Single-bottleneck fairness. An algorithm is said to be single-bottleneck fair if the shares $\varphi_i a_{ij}$ of resource j are all equal to $1/n$ whenever j is the only saturated resource.

Pareto-efficiency means no resources are needlessly left idle. Sharing-incentive means each job gets no less than it would under a strictly fair allocation where a $1/n$ share of each resource is allocated to every job. Strategy-proof means no job can get a better allocation by lying about its requirements. Scale-invariance means resource shares do not depend on the size of tasks (e.g., a scaling factor of 2 is equivalent to merging two tasks into one) and may be considered as a weaker form of strategy-proofness. Single-bottleneck fairness means the allocation reverts to simple fair sharing when there is only one saturated resource; this property and Pareto efficiency imply single-resource fairness.

Remark 1. It is straightforward to verify that single-resource fairness is realized by any Pareto-efficient algorithm that also satisfies any one of sharing-incentive, strategy-proof or scale-invariance. This means that if it is not single-resource fair, a Pareto efficient algorithm has none of the other properties.

2.3. Dominant resource fairness

We first note that simple max–min fairness is unsatisfactory. The max–min allocation is determined by the standard water-filling procedure: the φ_i are increased from zero together until some resource is saturated; the allocations of any job that needs this resource are frozen while the task volumes of the remainder increase together until another resource is saturated; the procedure continues until all job allocations are frozen.

The resulting allocation is uniquely characterized by the fact that, for each job i , there is some saturated resource j needed by job i such that φ_i is maximum over all jobs that need resource j . In particular, for $J = 1$, we have equal task volumes, $\varphi_i = 1/(\sum_{k=1}^n a_{k1})$, for $i = 1, \dots, n$. Resource shares $\varphi_i a_{i1}$ are thus proportional to requirements and max–min fairness is *not* single-resource fair. Since the algorithm is Pareto-efficient, Remark 1 shows it has neither sharing-incentive, strategy-proof nor scale-invariant properties. DRF is weighted max–min fairness with weights chosen to correct these failings. Weights are inversely proportional to the dominant resource requirements:

$$w_i = \frac{1}{\max_{j=1, \dots, J} a_{ij}},$$

for $i = 1, \dots, n$.

The allocation is derived from a modified water-filling procedure where filling rates are proportional to the per-job weights. It is the unique Pareto-efficient allocation such that, for each job i , there is some saturated resource j needed by job i such that φ_i/w_i is maximum over all jobs that need resource j . For $J = 1$, we get equal resource shares, $\varphi_i a_{i1} = 1/n$, for $i = 1, \dots, n$, so that DRF is single-resource fair. It is proved in [2] that DRF also satisfies the sharing-incentive and strategy-proof properties. For completeness, we provide concise proofs for these and the scale-invariant properties.

Proposition 1. DRF satisfies the sharing-incentive property.

Proof. For each job i , there is some saturated resource j needed by job i such that φ_i/w_i is maximum over all jobs that need resource j . In particular,

$$\varphi_i \max_{l=1, \dots, J} a_{il} = \frac{\varphi_i}{w_i} \geq \frac{1}{n} \sum_{k: a_{kj} > 0} \frac{\varphi_k}{w_k} \geq \frac{1}{n} \sum_{k: a_{kj} > 0} \varphi_k a_{kj} = \frac{1}{n}. \quad \square$$

Proposition 2. DRF is strategy-proof.

Proof. Assume that job i claims resource requirements a'_i instead of a_i . Denote by w'_i the corresponding weight and by φ'_i the resulting allocation. We compare the water-filling processes of both systems with the same filling rates for jobs other than

i. If job i is limited in the modified system first, then $\varphi'_i/w'_i \leq \varphi_i/w_i$ so that the dominant resource j of job i in the original system satisfies:

$$\varphi'_i a'_{ij} \leq \frac{\varphi'_i}{w'_i} \leq \frac{\varphi_i}{w_i} = \varphi_i a_{ij}.$$

If, on the other hand, job i in the original system is limited first on resource j , then all jobs except i that need this resource get a better allocation in the modified system. It follows that:

$$\varphi'_i a'_{ij} \leq 1 - \sum_{k \neq i} \varphi'_k a_{kj} \leq 1 - \sum_{k \neq i} \varphi_k a_{kj} = \varphi_i a_{ij}. \quad \square$$

Proposition 3. DRF is scale-invariant.

Proof. Assume that each job i claims resource requirements $a'_i = \alpha_i a_i$ instead of a_i , for some $\alpha_i > 0$. Let $w'_i = w_i/\alpha_i$ and $\varphi'_i = \varphi_i/\alpha_i$. The allocation φ' satisfies the new capacity constraints $\varphi' a' \leq 1$ and gives the same resource shares as the original allocation φ . It remains to show that φ' is DRF. We know that, for each job i , there exists some saturated resource j such that φ_i/w_i is maximum over all jobs that need resource j ; since $\varphi'_i/w'_i = \varphi_i/w_i$, this property is also satisfied by the allocation φ' which is therefore DRF. \square

Remark 2. Note that DRF also has the following somewhat surprising characteristic: the allocation may depend on resource requirements for non-saturated resources. Consider $J = 2$ resources, $n = 2$ jobs and $a_1 = (1, 1/3)$ and $a_2 = (1/2, 1)$. Then $\varphi = (2/3, 2/3)$ and resource 2 is not saturated. In the absence of resource 2 constraints, that is for $a_1 = (1, 0)$ and $a_2 = (1/2, 0)$, the allocation changes to $\varphi' = (1, 2)$ corresponding to equal shares of resource 1. The allocation depends on requirements for resource 2 even though this is not a bottleneck. This means, in particular, that DRF is *not* single-bottleneck fair.

2.4. Proportional fairness

An allocation that is well-known to have excellent properties in the context of sharing network bandwidth is proportional fairness (PF) [9]. The PF allocation is defined by:

$$\arg \max_{\varphi} \sum_{i=1}^n \log \varphi_i, \quad (2)$$

under capacity constraints (1). By construction, it is Pareto-efficient.

Let v be the vector of Lagrange multipliers associated with capacity constraints (1). The Karush–Kuhn–Tucker theorem implies that PF is the unique allocation φ such that:

$$\frac{1}{\varphi_i} = \sum_{j=1}^J a_{ij} v_j, \quad (3)$$

for $i = 1, \dots, n$, where

$$v_j \geq 0 \quad \text{and} \quad v_j \left(\sum_{i=1}^n \varphi_i a_{ij} - 1 \right) = 0, \quad (4)$$

for $j = 1, \dots, J$.

From (3) and (4) we deduce,

$$\sum_{j=1}^J v_j = \sum_{j=1}^J v_j \left(\sum_{i=1}^n \varphi_i a_{ij} \right) = \sum_{i=1}^n \varphi_i \left(\sum_{j=1}^J a_{ij} v_j \right) = n. \quad (5)$$

It follows that PF is single-bottleneck fair.

Proposition 4. PF satisfies the sharing-incentive property.

Proof. It follows from (3) and (5) that:

$$\varphi_i \max_j a_{ij} = \frac{\max_j a_{ij}}{\sum_{j=1}^J a_{ij} v_j} \geq \frac{1}{\sum_{j=1}^J v_j} = \frac{1}{n}. \quad \square$$

Proposition 5. PF is not strategy-proof.

Proof. It is sufficient to show this with an example. Consider $n = 2$ jobs and $J = 2$ resources with $a_1 = (1/2, 1)$ and $a_2 = (1, 1/2)$. We then have $v = (1, 1)$ and $\varphi = (2/3, 2/3)$. If job 1 claims $2/3$ units of resource 1 per task instead of $1/2$, we

have $a'_1 = (2/3, 1)$ and $a'_2 = (1, 1/2)$ yielding $v' = (2, 0)$ and $\varphi' = (3/4, 1/2)$. Job 1 receives a bigger share at the expense of job 2. \square

Remark 3. Unlike max–min fairness, it is not sufficient for a job to increase its resource requirements to increase its resource shares. Taking $a'_1 = (1, 1)$ and $a'_2 = (1, 1/2)$ in the above example yields $v' = (2, 0)$ and $\varphi' = (1/2, 1/2)$. Both job 1 and job 2 receive smaller resource shares.

Remark 4. A job can be sure to increase its resource shares only if it knows the resource requirements of the other jobs. Taking $a_1 = (1/2, 1)$ and $a_2 = (1, 1/3)$ yields $v = (1, 1)$ and $\varphi = (4/5, 3/5)$. If job 1 applies the same strategy as in the example of Proposition 5, we have $a'_1 = (2/3, 1)$ and $a'_2 = (1, 1/3)$ yielding $v' = (2, 0)$ and $\varphi' = (3/4, 1/2)$. Both jobs again receive fewer resources.

The above remarks suggest a user wishing to exploit PF would be hard-pressed to determine an optimal strategy. The impact of false resource claims can only be appraised by considering likely combinations of jobs in progress, as determined by a model of dynamic traffic. The following proposition shows that PF does have the scale-invariance guarantee.

Proposition 6. *PF is scale-invariant.*

Proof. Assume that each job i claims resource requirements $a'_i = \alpha_i a_i$ instead of a_i , for some $\alpha_i > 0$. Let $\varphi'_i = \varphi_i / \alpha_i$. The allocation φ' satisfies the new capacity constraints $\varphi' a' \leq 1$ and gives the same resource shares as the original allocation φ . Moreover, from (3) and (4), it can easily be verified that φ' is the PF allocation under capacity constraints $\varphi' a' \leq 1$. \square

2.5. Utility maximization

PF sharing maximizes overall utility under the assumption that the utility for each job i of allocation φ_i is proportional to $\log \varphi_i$. One might consider allocations defined in terms of more general, differentiable and strictly concave utility functions $U(\cdot)$, i.e.,

$$\varphi = \arg \max_{\varphi} \sum_{i=1}^n U(\varphi_i), \quad (6)$$

under capacity constraints (1).

For instance, $U = (\cdot)^\alpha / (1 - \alpha)$, for $\alpha \in (0, +\infty) \setminus \{1\}$, defines the family of α -fair allocations [12]. Proportional fairness corresponds to the limit $\alpha \rightarrow 1$ while max–min fairness is realized in the limit $\alpha \rightarrow +\infty$.

Let v be the vector of Lagrange multipliers associated with capacity constraints (1). The Karush–Kuhn–Tucker theorem states that a utility-based allocation is the unique allocation φ such that:

$$\forall i = 1, \dots, n, \quad U'(\varphi_i) = \sum_{j=1}^J a_{ij} v_j, \quad (7)$$

under constraints (4). Since utility maximization allocations are necessarily Pareto-efficient, the following property shows that no utility-based allocation except PF can be either sharing-incentive, strategy-proof or scale-invariant (cf. Remark 1).

Proposition 7. *PF is the only single-resource fair utility-based allocation.*

Proof. Take a single resource and two jobs with resource requirements $a_{11} = \frac{1}{2}$ and $a_{21} = \frac{1}{2x}$ for some $x > 0$. In view of (7), we have $U'(\varphi_1) = xU'(\varphi_2)$. If the allocation is single-resource fair, we have equal resource shares $a_{11}\varphi_1 = a_{21}\varphi_2 = \frac{1}{2}$ so that:

$$\forall x > 0, \quad U'(x) = \frac{U'(1)}{x}.$$

This implies that $U = \alpha \log(\cdot) + \beta$ for some constants $\alpha > 0$ and β , which is PF. \square

3. Behavior in dynamic traffic

Jobs handled by a data center cluster arrive over time and remain active for a finite duration that depends on the resource shares they are allocated. Algorithm performance, perceived in terms of job completion times, can only be realistically appraised under such dynamic traffic. We adopt a simple Markovian traffic model and retain the assumption of infinitely divisible resources to illustrate the relative performance of DRF and PF.

3.1. Traffic model

We consider the following dynamic traffic model. There are K distinct classes of jobs. Jobs of class k arrive at rate λ_k and require the completion of σ_k tasks on average. We denote by n_k the number of class- k jobs in progress and by $n = n_1 + \dots + n_K$

the total number of jobs. Each task of a class- k job requires a_{kj} units of resource j . The capacity constraints are:

$$\sum_{k=1}^K n_k \varphi_k a_{kj} \leq 1, \quad (8)$$

for $j = 1, \dots, J$. Let τ_k be the mean service time of one task of a class- k job. Then $\mu_k = 1/(\sigma_k \tau_k)$ is the completion rate of a class- k job with exactly one ongoing task and $\varphi_k \mu_k$ is the completion rate of a class- k job.

Under the fluid sharing model with exponential assumptions, i.e., Poisson job arrivals and exponentially distributed job sizes, the vector $\vec{n} = (n_1, \dots, n_K)$ is a Markov process with component- k birth rate λ_k and death rate $n_k \varphi_k \mu_k$, for $k = 1, \dots, K$.

Markovian assumptions are adopted for the sake of simplicity. We note however that insensitivity results for bandwidth sharing in networks suggest the performance of this simple model is indicative of much more general and realistic traffic models (e.g., see [10]). The impact of non-infinitely divisible resources is evaluated in the next section.

We assume allocations are instantaneously updated every time a new job starts or a job in progress completes. For DRF the allocation is determined by applying the modified water-filling algorithm. For PF we solve the associated convex optimization problem (2).

3.2. Traffic capacity

Tasks of class- k jobs arrive at rate $\lambda_k \sigma_k$ and last τ_k seconds on average. The corresponding traffic demand is thus $\rho_k = \lambda_k / \mu_k$. In dynamic traffic, the considered system may be unstable if demand is too high in the sense that the number of jobs in progress may then grow indefinitely. The set of loads for which the system remains stable defines a capacity region. The capacity region of DRF and PF (and alternative utility-based fairness allocations) can be directly deduced from known results for bandwidth sharing in networks [13]. It is defined by the following inequalities:

$$\sum_{k=1}^K \rho_k a_{kj} < 1, \quad (9)$$

for $j = 1, \dots, J$. In other words, the load on each resource must be less than its normalized capacity. In the following, we assume that this condition is satisfied and consider the system in steady state. We denote by π the stationary distribution of the Markov process \vec{n} .

3.3. Performance

Users experience system performance through the time it takes to complete their job. As this is job size dependent, we prefer to present class- k performance in terms of the mean service rate γ_k defined as the ratio of the mean job size σ_k to the mean completion time. By Little's law, we have the following expression:

$$\gamma_k = \frac{\lambda_k}{E(n_k)}, \quad (10)$$

where $E(n_k)$ is the mean number of ongoing class- k jobs in steady state.

We have evaluated the mean service rates γ_k for DRF and PF allocations by simulating some simple system configurations. We suppose each job requires quantities of two resources, CPU and RAM say, with relative per-task requirements $a_1 = (1, 1/10)$ and $a_2 = (1/10, 1)$. We set $\mu_1 = \mu_2 = 1$ and consider two demand patterns: *balanced* with $\rho_1 = \rho_2 = \rho/2$, where ρ is the total demand, and *unbalanced* with $\rho_1 = \frac{3}{4}\rho$ and $\rho_2 = \frac{1}{4}\rho$. For balanced demand the load of both resources is $\frac{\rho}{2}(1 + \frac{1}{10})$ whereas the unbalanced case yields a CPU load of $\frac{\rho}{4}(3 + \frac{1}{10})$ and a RAM load of $\frac{\rho}{4}(\frac{3}{10} + 1)$, nearly 3 times smaller.

Fig. 1 compares the mean relative service rate provided by DRF and PF allocations as a function of CPU load. Note that the service rate is 1 (in normalized units) at load zero, when never more than 1 job is in progress, and tends to zero as load attains system capacity. For a single resource, the service rate would in fact be $1 - \varrho$ for both DRF and PF, where ϱ is the resource load. This is so because allocations are single-resource fair and therefore realize a processor sharing system.

Both graphs in Fig. 1 present 4 curves representing γ_k for $k = 1, 2$ for DRF and PF. In the left hand plot the curves for both classes for each fairness objective coincide by symmetry while PF outperforms DRF, especially at high load. The right hand plot for unbalanced demand shows that $\gamma_2^{(PF)}$ is significantly higher than $\gamma_2^{(DRF)}$ while the service rates γ_1 for the CPU constrained class are practically the same. This occurs because of the bias exhibited by DRF but not PF when resource 2 (RAM) is not saturated, as noted in Remark 2.

Fig. 2 presents a case where PF also performs better than DRF at low load. To be precise, the service rate of class 2 is significantly greater at the cost of a slight decrease of that of class 1. The class profiles for this case are $a_1 = (1, 1/10)$ and $a_2 = (1/2, 1)$.

The above simple examples clearly do not prove that PF is always better than DRF. To gain insight we consider the light traffic and heavy traffic behavior of a general 2-class, 2-resource type system. Suppose, without loss of generality in view of scale invariance, job classes have task profiles $a_1 = (1, \alpha)$ and $a_2 = (\beta, 1)$ with $\alpha \leq \beta < 1$.

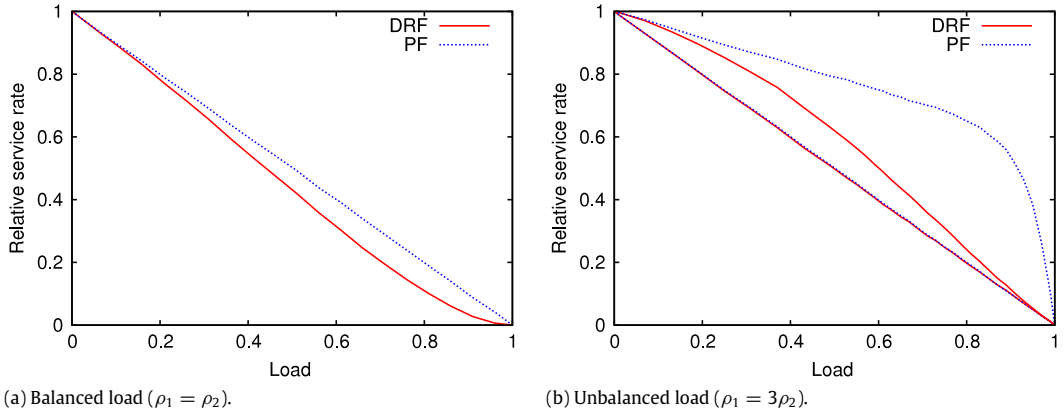


Fig. 1. Service rate against load of resource 1 (CPU) – $a_1 = (1, 1/10)$ and $a_2 = (1/10, 1)$.

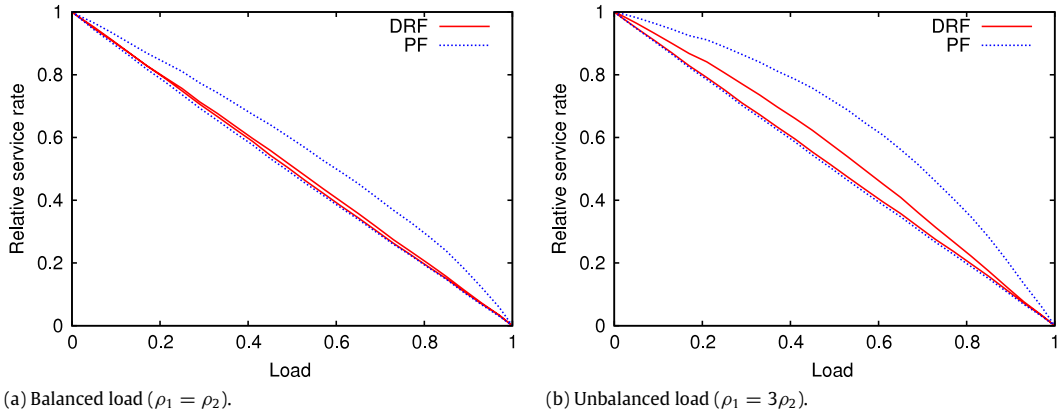


Fig. 2. Service rate against load of resource 1 (CPU) – $a_1 = (1, 1/10)$ and $a_2 = (1/2, 1)$.

Light traffic. When $\rho_1 \rightarrow 0$ and $\rho_2 \rightarrow 0$ with fixed positive ratio ρ_1/ρ_2 , the second-order approximation of the stationary distribution π is that of the Markov process \bar{n} truncated to states with no more than 2 ongoing jobs. We obtain:

$$\pi(0, 0) \approx 1 - \rho_1 - \rho_2 - \rho_1^2 - \rho_2^2 + 2 \frac{\varphi_1 + \varphi_2 - 1}{\varphi_1 + \varphi_2} \rho_1 \rho_2,$$

$$\pi(1, 0) \approx \rho_1 - \frac{2\varphi_1}{\varphi_1 + \varphi_2} \rho_1 \rho_2,$$

$$\pi(0, 1) \approx \rho_2 - \frac{2\varphi_2}{\varphi_1 + \varphi_2} \rho_1 \rho_2,$$

$$\pi(1, 1) \approx \frac{2}{\varphi_1 + \varphi_2} \rho_1 \rho_2,$$

$$\pi(2, 0) \approx \rho_1^2,$$

$$\pi(0, 2) \approx \rho_2^2,$$

where φ_1 and φ_2 denote the task volumes in state (1, 1) of jobs 1 and 2, respectively. Using (10), we deduce:

$$\gamma_1 \approx 1 - \rho_1 - 2 \frac{1 - \varphi_1}{\varphi_1 + \varphi_2} \rho_2, \quad \gamma_2 \approx 1 - \rho_2 - 2 \frac{1 - \varphi_2}{\varphi_1 + \varphi_2} \rho_1.$$

Under DRF, $\varphi_1 = \varphi_2 = \frac{1}{1+\beta}$ so that:

$$\gamma_1^{(\text{DRF})} \approx 1 - \rho_1 - \beta \rho_2 \quad \text{and} \quad \gamma_2^{(\text{DRF})} \approx 1 - \rho_2 - \beta \rho_1.$$

Under PF, $\varphi_1 = \max(\frac{1}{2}, \frac{1-\beta}{1-\alpha\beta})$, $\varphi_2 = \min(\frac{1}{2\beta}, \frac{1-\alpha}{1-\alpha\beta})$ so that:

$$\gamma_1^{(\text{PF})} \approx 1 - \rho_1 - 2 \min\left(\frac{\beta}{\beta+1}, \frac{\beta(1-\alpha)}{2-\alpha-\beta}\right) \rho_2 \quad \text{and} \quad \gamma_2^{(\text{PF})} \approx 1 - \rho_2 - 2 \max\left(\frac{2\beta-1}{\beta+1}, \frac{\alpha(1-\beta)}{2-\alpha-\beta}\right) \rho_1.$$

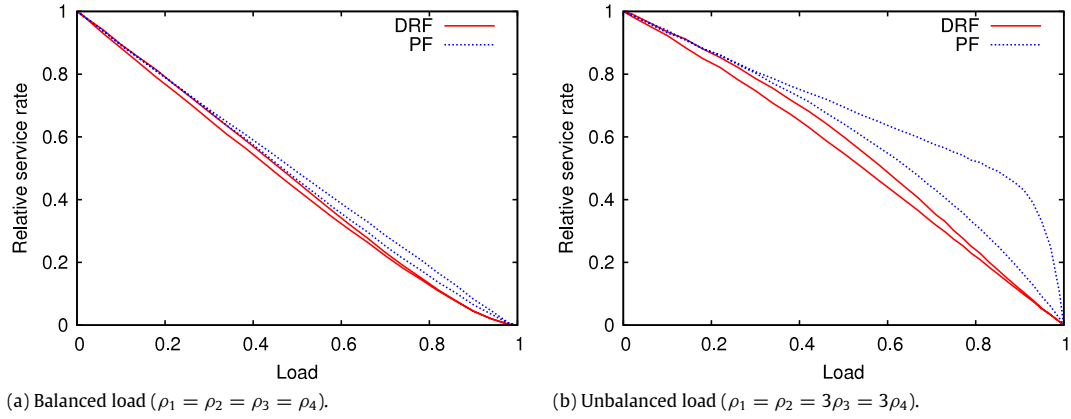


Fig. 3. Service rate against load of resource 1 (CPU).

We see that, in light traffic, PF and DRF yields the same performance when $\alpha = \beta$ while PF worsens the performance for class-1 jobs and improves the performance of class-2 jobs compared to DRF when $\alpha < \beta$. When $\alpha = 1/10$ and $\beta = 1/2$ for instance, corresponding to the results of Fig. 2, we get for balanced load:

$$\gamma_1^{(\text{DRF})} \approx \gamma_2^{(\text{DRF})} \approx 1 - \varrho,$$

where $\varrho = \rho_1 + \beta\rho_2$ is the load of resource 1, while:

$$\gamma_1^{(\text{PF})} \approx 1 - \frac{23}{21}\varrho, \quad \gamma_2^{(\text{PF})} \approx 1 - \frac{5}{7}\varrho.$$

Thus the performance improvement for class-2 jobs is significant, at the cost of a slight performance degradation for class-1 jobs. The conclusion is similar in case of unbalanced load.

Heavy traffic. Consider now the same profiles in heavy traffic with $\varrho = \rho_1 + \beta\rho_2 > \alpha\rho_1 + \rho_2$, so that resource 1 has higher load than resource 2, and $\varrho \rightarrow 1$. The service rate for PF is then determined only by the jobs' share of resource 1 yielding $\gamma_1^{(\text{PF})} \approx 1 - \varrho$ and $\gamma_2^{(\text{PF})} \approx (1 - \varrho)/\beta$. DRF, on the other hand, will equalize the service rates of both classes so that $\gamma_1^{(\text{DRF})} \approx \gamma_2^{(\text{DRF})} \approx 1 - \varrho$. This confirms the shapes of the curves in Figs. 1 and 2 for unbalanced load. Heavy traffic behavior for balanced load can be explained similarly, except in the symmetric case $\alpha = \beta$ where the performance DRF is much worse due to inefficient resource utilization.

In Fig. 3, we plot results for 4 job classes with profiles $a_1 = (1, 1/10)$, $a_2 = (1, 1/2)$, $a_3 = (1/2, 1)$ and $a_4 = (1/10, 1)$. The figures plot the service rates of classes 3 and 4 whose dominant resource is resource 2 (RAM). These results confirm the trends shown in previous cases are maintained when the traffic mix is more varied.

We have simulated several other configurations including cases with three types of resource. The results confirm the trends illustrated in the above plots: PF slightly outperforms DRF for balanced traffic; the gain is more significant for jobs whose dominant resource has relatively lower load; the gain comes at the cost of an imperceptible loss of service rate for jobs whose dominant resource is most heavily loaded. The latter is approximately $1 - \varrho$ for both DRF and PF where ϱ is the resource load, as predicted by the underlying processor sharing model, except for DRF under balanced load. DRF tends to behave poorly in heavy, balanced traffic since it penalizes efficiency in favor of strict fairness.

3.4. Strategy-proofness

The strategy-proof property in Section 2.2 relates to a static population of jobs. We demonstrate in this subsection that the absence of this property for PF is *not a disadvantage* in dynamic traffic. Even if some job knows some statistics about the resource requirements of competing jobs, it does not know how many jobs are active and thus can hardly define a consistent strategy to improve its service rate by lying about its resource requirements.

Consider the following example. There are $J = 2$ resources and $K = 2$ classes with resource requirements $a_1 = (1, \alpha)$ and $a_2 = (\alpha, 1)$ for some $\alpha \in (0, 1)$. There is a single, permanent class-1 job. The only way for this job to increase its service rate is to claim $\beta \in (\alpha, 1)$ units of resource 2 per task. Under PF, its service rate in the presence of n other jobs is:

$$\varphi_1(n) = \min \left(\frac{1 - \alpha}{1 - \alpha\beta}, \frac{1}{\beta(n + 1)} \right),$$

which is maximum for:

$$\beta = \frac{1}{1 + n(1 - \alpha)}.$$

Note that this value of β lies in $(\alpha, 1)$ as long as $n \leq 1/\alpha$, in which case job 1 gets a higher service rate and each class-2 job gets $\varphi_2(n) = \frac{1}{n}(1 - \beta\varphi_1(n))$. The number of class-2 jobs forms a birth-death process with birth rate λ_2 and death rate $\mu_2 n \varphi_2(n)$ in state n . This process is stable whenever $\rho_2 < 1$ and then has a certain stationary distribution π_2 yielding the mean service rate:

$$\gamma_1^{(\text{PF})} = \sum_{n \geq 0} \pi_2(n) \varphi_1(n).$$

Fig. 4 plots $\gamma_1^{(\text{PF})}$ against ρ_2 for $\alpha = \frac{1}{10}$ and three values of β , chosen to maximize the service rate $\varphi_1(n)$ assuming $n = 1, 2$ and 5, respectively. We observe that no benefit is gained from any of these static optimizations demonstrating that PF is indeed strategy proof in this scenario.

If, in an alternative scenario, the requirements of all jobs in a given class k were falsely declared with $a'_{kj} \geq a_{kj}$, for $j = 1, \dots, J$, and strict inequality for at least one resource, it is clear from (9) that the capacity region would be smaller for both DRF and PF. Simulation results shown in Fig. 5 confirm that service rates at stable loads are reduced for both DRF and PF. The figure plots the class 1 service rate for the configuration of Fig. 4 with balanced load. The dotted lines result from the same values of β used in Fig. 4 confirming there is no winning strategy. Traffic capacity is limited by $\beta\rho_1 + \rho_2 < 1$ instead of $\alpha\rho_1 + \rho_2 < 1$, corresponding to a maximum load of $\frac{1+\alpha}{1+\beta} \approx .72, .81, .93$, respectively. Of course, class 2 jobs also suffer from the false claims of class 1 in all these cases.

4. Task assignment algorithms

While the fluid models analyzed above are useful in appraising the relative merits of alternative fairness objectives, it is necessary in practice to account for the fact that tasks are not infinitely divisible. In this section we relax this assumption. We discuss the implementation of task-based versions of DRF and PF and compare their performance in simulation experiments.

4.1. “Serve the most deprived job”

We follow the proposal in [2] in seeking to realize the fairness objective by preferentially launching tasks of the most deprived job with respect to the objective fairness criterion. Let φ_i now denote the integer number of tasks in progress for job i . For DRF, the most deprived job is that for which $\varphi_i \max_j \{a_{ij}\}$ is minimal. For PF, the job to serve is that for which $\varphi_i \sum_{j=1}^J a_{ij} v_j$ is minimal, where the v_j are the Lagrange multipliers. These can be computed iteratively using the gradient descent method:

$$v_j \leftarrow \left(v_j + \gamma \left(\sum_{i=1}^n \psi_i a_{ij} - 1 \right) \right)^+,$$

where $\gamma > 0$ is a sufficiently small step size and $\psi_i = 1 / \sum_{j=1}^J a_{ij} v_j$ for $i = 1, \dots, n$.

Tasks are assigned when jobs arrive, if resources are available, or when other tasks end. If available resources are insufficient to accommodate a task of the most deprived job, allocations are frozen until a sufficient number of other tasks end or some other job takes over the most-deprived status. Note that we continue to assume resources are pooled though the assignment algorithm could be adapted to account for additional practical constraints such as assigning CPU and RAM on the same physical server, say, [6], or meeting other compatibility constraints [14].

4.2. Performance

We have simulated task-based allocations in the following configuration. Two classes of job require different amounts of CPU and RAM in the following proportions: $a_1 = (1, 1/10)$ and $a_2 = (1/10, 1)$, as previously considered in Section 3.3. Jobs of both classes require the execution of 500 tasks while CPU and RAM both have capacity 100. This means a single job of either class in an otherwise empty system will have 100 tasks in progress until 400 are completed and then 99, 98 and so on, until the last task finishes. This is a significant discrepancy with respect to the fluid model where the service rate would remain constant at 100 until job completion. With respect to previous notation, we have $\sigma_k = 500$, $\tau_k = 1/5$, $C = 100$ and therefore $\mu_k = C / (\sigma_k \tau_k) = 1$, for $k = 1, 2$.

Fig. 6 shows the simulation results for DRF and PF allocations under unbalanced load, $\rho_1 = 3\rho_2$. The left hand plots relate to exponentially distributed task sizes. Compared to the right hand plots of Fig. 1 there is a significant reduction in service rate for low loads due to the discrepancy noted above. The service time of an isolated job in the fluid model is the sum of 500 task inter-departure intervals at rate 100 compared, in the task-based system, to 401 intervals at rate 100 plus one at rate 99, one at rate 98, ..., and 1 at rate 1. The ratio of service times is thus

$$\frac{500}{100} \bigg/ \left(\frac{401}{100} + \frac{1}{99} + \frac{1}{98} + \dots + 1 \right) \approx 0.54,$$

corresponding to the zero load service rate seen in the figure.

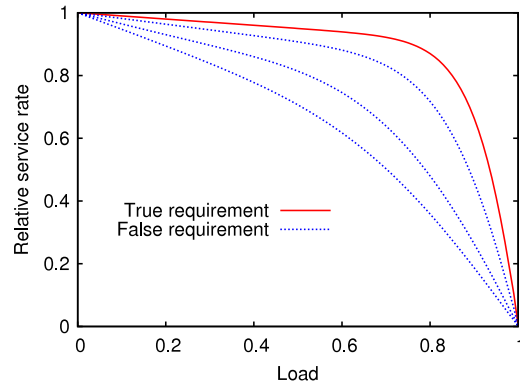


Fig. 4. Mean service rate of the test flow against CPU load with true and false requirements that are optimal under PF for 1, 2 and 5 other flows ($\beta \approx .53, .36, .18$, from bottom to top).

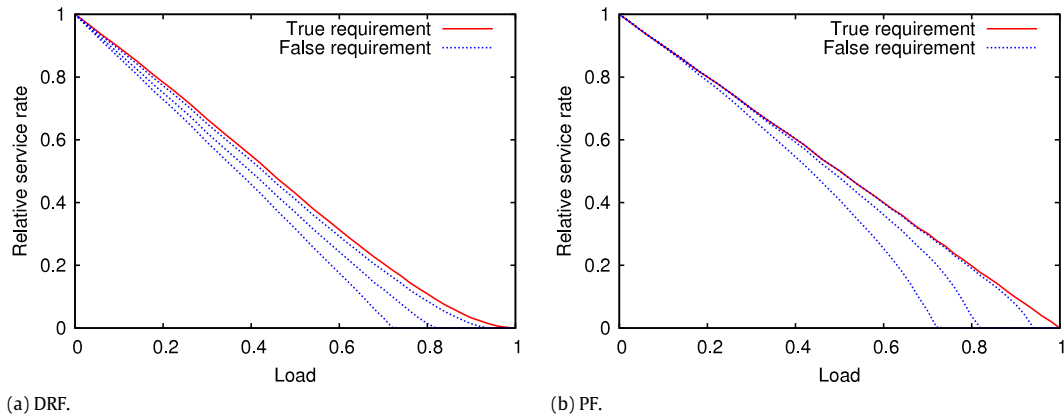


Fig. 5. Mean service rate of class 1 flows against CPU load with false class-1 requirements as in Fig. 4.

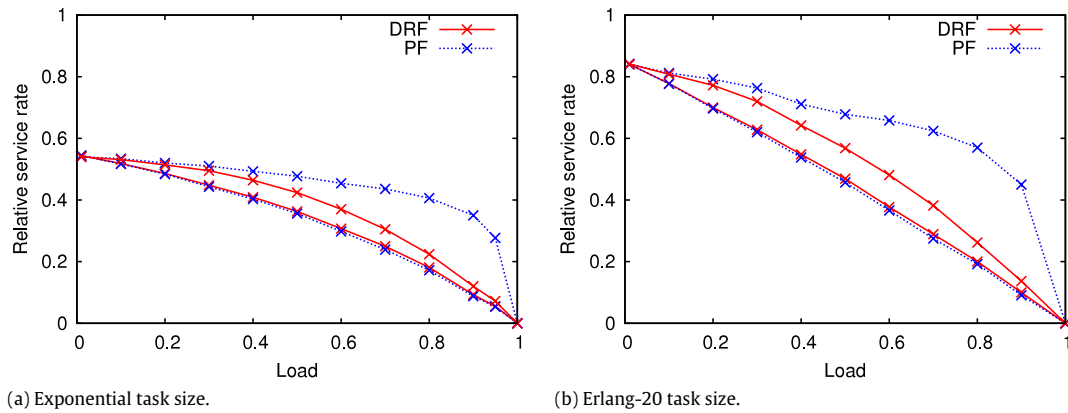


Fig. 6. Service rate against load of resource 1 (CPU) for unbalanced load ($\rho_1 = 3\rho_2$).

The discrepancy is smaller for task service times that are more deterministic. The right hand plot relates to an Erlang-20 distribution (i.e., the sum of 20 i.i.d. exponential random variables) with the same mean yielding service rates much closer to those of the fluid model. The fluid model is also a good fit when the number of tasks per job is much bigger than the resource capacity (results not shown).

Despite the discrepancy, it is clear that the relative advantage of PF over DRF, observed and analyzed in Section 3, is preserved under the task-based model. PF realizes a more favorable efficiency–fairness tradeoff, especially when demand for some resource is relatively low.

5. Related work

A number of authors have extended the DRF concept introduced by Ghodsi et al. [2]. Parkes et al. [5] generalize DRF notably to account for per-job sharing weights. Psomas and Schwartz [6] take account of the fact that task resources must fit into a single machine, bringing bin packing considerations to the allocation problem. Wang et al. [7] further generalize this approach by accounting for heterogeneous server capacity limits. Bhattacharya et al. [15] add the notion of hierarchical scheduling to DRF: cluster resource allocations account also for shares attributed to the departments that own particular jobs. It is noteworthy that DRF is now implemented in the Hadoop Next Generation Fair Scheduler². Lastly, we note the proposal to use DRF in a different application setting, namely sharing middlebox resources in software routers [16,17].

An alternative multi-resource fairness criterion is proposed by Dolev et al. [3]. The authors introduce the notion of *entitlement* corresponding to the share of each resource to which a job should minimally have access. The objective is then to define an allocation such that users do not receive less than their entitlement on at least one bottlenecked resource and therefore have “no justified complaints”. The resulting allocation, named bottleneck-based fairness (BBF), has been further developed by Zeldes and Feitelson [4] who propose an on-line algorithm. The multi-resource fairness problem is defined more formally by Gutman and Nisan who generalize the DRF and BBF definitions somewhat [8]. These authors identify similarities with prior work on fairness in the economics literature and propose polynomial complexity algorithms drawn from that source. Joe-Wong et al. [18] propose alternative generalized definitions of what constitutes a fair multi-resource allocation identifying notions of “Fairness on dominant shares” (FDS) and “Generalized fairness on jobs” (GFJ).

All the work cited in the above two paragraphs considers sharing between a fixed population of jobs and ignores the impact of traffic dynamics due to arrivals over time of finite size jobs. They do not therefore evaluate the performance of allocations in terms of realized job completion times which is arguably the most relevant criterion in appraising their effectiveness.

There is an abundant literature on sharing network bandwidth. PF and, more generally, the idea that bandwidth sharing between a fixed set of flows should maximize a sum of per-flow utilities were notions introduced by Kelly et al. in 1998 [9]. In that work, and subsequent contributions by many authors on network utility maximization, the population of flows in progress is assumed fixed. We believe this to be inappropriate, like the assumption that the population of jobs is fixed when evaluating cluster sharing. It does not make much sense to measure the utility of a flow as a function of its rate (e.g., the log function for PF) on recognizing that the population of flows in progress, and consequently the rate of any ongoing flow, changes rapidly and significantly throughout the lifetime of that flow.

The study of bandwidth sharing under a dynamic model of traffic also began at the end of the 90s [19]. The most significant results derived over an ensuing period of several years are summarized in the paper by Bonald et al. [10]. It turns out that PF brings significantly better flow completion time performance than max–min fairness for a network with a mix of wired and wireless links (e.g., see Fig. 17.14 in [11]). Flows in such networks require unequal resource shares, like jobs in the compute cluster, so that these bandwidth sharing results anticipate the excellent performance of PF demonstrated here.

6. Conclusions

Sharing multiple compute cluster resources like CPU and RAM is in many ways analogous to sharing bandwidth in a network where each link is considered as a distinct resource. The analogy is closest with respect to a wireless network where flows have different resource requirements on wired and radio links, respectively. We know proportional fairness (PF) performs well in such networks and have therefore sought to evaluate its suitability for sharing cluster resources. Our results confirm that PF is indeed preferable to so-called dominant resource fairness (DRF). DRF is the sharing criterion currently applied by the Next Generation Hadoop Fair Scheduler. On the basis of the present findings, it would be advisable to consider replacing this by PF.

DRF was developed to realize a set of sharing properties derived essentially from micro-economics considerations. In particular, DRF is claimed to be strategy proof in the sense that a job cannot gain better service by announcing false requirements. We have argued here that these properties must be revisited to account for the dynamics over time of the job population. In this context, an algorithm is strategy proof if the expected completion time of a job cannot be reduced by false announcements, given that the user is necessarily unaware of the precise mix of competing jobs. In the dynamic setting, PF is strategy proof to the same degree as DRF: it is not possible to improve an allocation by declaring false requirements.

PF realizes a better fairness–efficiency tradeoff than DRF over the range of possible loads. This is especially true when resource load is unbalanced with some resources remaining underloaded as others become saturated. DRF is then particularly deficient in overly limiting use of the underloaded resource in order to satisfy its particular notion of fairness.

PF can be realized with the same, low complexity as DRF by adapting its “serve-the-most-deprived-job” policy. When resources are freed as tasks complete, they are assigned to new tasks of the job whose current allocation is furthest from the PF ideal. This ideal allocation is easily determined as the result of a simple convex optimization problem.

It remains to adapt the PF allocation to account for practical placement constraints that occur in data centers. These include the requirement to allocate CPU and RAM on the same server, for example, or to use a specific type of server for some

² <http://hadoop.apache.org/docs/r2.3.0/hadoop-yarn/hadoop-yarn-site/FairScheduler.html>

jobs. It is also of interest to incorporate hierarchical scheduling criteria whereby the departments originating certain classes of jobs have overall service rate guarantees.

Available trace data suggests cluster resources are currently significantly over-provisioned so that optimized sharing is hardly an issue in practice. In future, however, one would expect tighter energy management leading to a requirement for more closely controlled performance. The key to this would be analytical modeling enabling one to predict performance given system capacity and load, by applying the robust engineering techniques developed over recent years for networks, for example.

References

- [1] C. Reiss, A. Tumanov, G.R. Ganger, R.H. Katz, M.A. Kozuch, Heterogeneity and dynamicity of clouds at scale: google trace analysis, in: Proceedings of the Third ACM Symposium on Cloud Computing, SoCC '12, ACM, New York, NY, USA, 2012, pp. 7:1–7:13.
- [2] A. Ghodsi, M. Zaharia, B. Hindman, A. Konwinski, S. Shenker, I. Stoica, Dominant resource fairness: Fair allocation of multiple resource types, in: Proceedings of the 8th USENIX Conference on Networked Systems Design and Implementation, USENIX Association, Berkeley, CA, USA, 2011, p. 24.
- [3] D. Dolev, D.G. Feitelson, J.Y. Halpern, R. Kupferman, N. Linial, No justified complaints: On fair sharing of multiple resources, in: Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, ITCS '12, ACM, New York, NY, USA, 2012, pp. 68–75.
- [4] Y. Zeldes, D.G. Feitelson, On-line fair allocations based on bottlenecks and global priorities, in: Proceedings of the 4th ACM/SPEC International Conference on Performance Engineering, ICPE '13, ACM, New York, NY, USA, 2013, pp. 229–240.
- [5] D.C. Parkes, A.D. Procaccia, N. Shah, Beyond dominant resource fairness: extensions, limitations, and indivisibilities, in: ACM Conference on Electronic Commerce, ACM, 2012, pp. 808–825.
- [6] C.-A. Psomas, J. Schwartz, Beyond beyond dominant resource fairness: indivisible resource allocation in clusters, Tech Report Berkeley, 2013.
- [7] W. Wang, B. Li, B. Liang, Dominant resource fairness in cloud computing systems with heterogeneous servers, in: Proceedings of INFOCOM 2014, 2014.
- [8] A. Gutman, N. Nisan, Fair allocation without trade, in: International Foundation for Autonomous Agents and Multiagent Systems, AAMAS '12, in: Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems, vol. 2, Richland, SC, 2012, pp. 719–728.
- [9] F.P. Kelly, A.K. Maulloo, D.K.H. Tan, Rate control for communication networks: shadow prices, proportional fairness and stability, J. Oper. Res. Soc. 49 (3) (1998) 237–252.
- [10] T. Bonald, L. Massoulié, A. Proutière, J. Virtamo, A queueing analysis of max–min fairness, proportional fairness and balanced fairness, Queueing Syst. Theory Appl. 53 (1–2) (2006) 65–84.
- [11] T. Bonald, A. Proutière, A queueing analysis of data networks, in: R. Boucherie, N. van Dijk (Eds.), Queueing Networks: A fundamental approach, Springer, 2011, pp. 729–765.
- [12] J. Mo, J. Walrand, Fair end-to-end window-based congestion control, IEEE/ACM Trans. Netw. 8 (5) (2000) 556–567.
- [13] H.-Q. Ye, Stability of data networks under an optimization-based bandwidth allocation, Automatic Control, IEEE Transactions on 48 (7) (2003) 1238–1242.
- [14] A. Ghodsi, M. Zaharia, S. Shenker, I. Stoica, Choosy: Max–min fair sharing for datacenter jobs with constraints, in: Proceedings of the 8th ACM European Conference on Computer Systems, EuroSys '13, ACM, New York, NY, USA, 2013, pp. 365–378.
- [15] A.A. Bhattacharya, D. Culler, E. Friedman, A. Ghodsi, S. Shenker, I. Stoica, Hierarchical scheduling for diverse datacenter workloads, in: Proceedings of the 4th Annual Symposium on Cloud Computing, SoCC '13, ACM, New York, NY, USA, 2013, pp. 4:1–4:15.
- [16] A. Ghodsi, V. Sekar, M. Zaharia, I. Stoica, Multi-resource fair queueing for packet processing, in: Proceedings of ACM SIGCOMM 2012, ACM, New York, NY, USA, 2012, pp. 1–12.
- [17] W. Wang, B. Li, B. Liang, Multi-resource round robin: A low complexity packet scheduler with dominant resource fairness, in: ICNP 2013, 2013.
- [18] C. Joe-Wong, S. Sen, T. Lan, M. Chiang, Multi-resource allocation: Fairness–efficiency tradeoffs in a unifying framework, in: INFOCOM, 2012 Proceedings IEEE, 2012, pp. 1206–1214.
- [19] L. Massoulié, J. Roberts, Bandwidth sharing and admission control for elastic traffic, Telecommunication Systems 15 (1–2) (2000) 185–201. <http://dx.doi.org/10.1023/A:1019138827659>.



Thomas Bonald is Professor in the department of Computer Science and Networking at Telecom ParisTech, France. He graduated as an Engineer from Ecole polytechnique in 1994 and from Telecom ParisTech in 1996. He received the Ph.D. degree in Applied Mathematics from Ecole polytechnique in 1999. From 1999 to 2009 he was a Research Engineer at Orange Labs, working on traffic modeling for wireline and wireless networks. His main research interests are in queueing theory and stochastic models of communication networks. Thomas Bonald is the author of 70 scientific papers, 10 patents and the book, Network Performance Analysis (Wiley, 2011). He received the Best Paper Award of ACM SIGMETRICS / IFIP Performance 2004. He is a member of the Board of Directors of ACM SIGMETRICS and has served on the Technical Program Committees of many conferences, including ACM SIGMETRICS IFIP Performance and IEEE INFOCOM. He has also served on the Editorial Boards of Queueing Systems, Theory and Applications and of the IEEE/ACM Transactions on Networking.



Jim Roberts very recently joined the French research institute IRT-SystemX to work on a project on Cloud computing and network architecture. He was previously with Inria from September 2009 after spending more than thirty years with France Telecom research labs. He received a degree in Mathematics from the University of Surrey in 1970 and a doctorate in computer science in 1987 from the University of Paris VI. His research is centered on the performance evaluation and design of traffic controls for communication networks. In a long career, he has published around 100 papers, chaired several program committees and been associate editor for a number of journals. He gave the Keynote at Infocom 2013. He is a Fellow of the Société des électriciens et électroniciens (SEE) and recipient of the Arne Jensen lifetime achievement award from the International Teletraffic Congress (ITC).