

# Fair Scheduling for Deadline-Driven, Resource-Constrained, Multi-Analytics Workloads

Stratos Dimopoulos, Chandra Krintz, Rich Wolski

Department of Computer Science  
University of California, Santa Barbara  
{stratos,ckrintz,rich}@cs.ucsb.edu

**Abstract**—In this paper, we analyze and empirically evaluate *Justice*, a fair-share, deadline-aware job scheduler for resource-constrained cloud deployments managed by big data resource negotiators. *Justice* provides admission control, which leverages historical traces and job deadlines to guide and adapt resource allocation decisions to changing cloud conditions. We evaluate *Justice* using different deadline types and production workloads. We find that it outperforms extant allocators in terms of fair allocation, deadline satisfaction, and useful work.

**Keywords**—cloud for IoT; scheduling; big-data; fairness;

## I. INTRODUCTION

Increasingly, cloud users deploy big data frameworks (e.g., Apache Hadoop and Spark) via resource negotiators such as Apache Mesos and YARN. Resource negotiators simplify deployment and enable multiple frameworks to execute concurrently using the same set of resources. They employ fair-share resource allocators [6, 8], which attempt to partition resources equally across frameworks.

In this paper, we investigate fair-share allocation for workloads with deadline and resource constraints. Deadline-driven workloads represent an important class of big data applications [12, 15, 20], which are unfortunately under supported in multi-analytic settings. Resource-limited deployments are those in which more resources (e.g., CPU, memory) cannot simply be added on-demand, in exchange for an additional charge, as they can in a public cloud. Such deployments include private clouds, IoT edge systems, and cloudlets in which data analytics is performed near where data is collected to provide low-latency (deadline-driven) actuation, control, data privacy, decision support, and to reduce bandwidth requirements [19]. Because modern resource negotiators and big data frameworks were not designed for this combination of constraints, their use can result in low utilization, poor performance, missed deadlines, and unfair sharing [3].

To address these limitations, we design and implement admission control for resource negotiators that satisfies deadlines while preserving fairness in resource-constrained environments. Our system, called *Justice*, is framework agnostic, so it can be part of existing open-source resource negotiators like Mesos, YARN, and Kubernetes and its focus on efficient resource usage makes it ideal for environments with limited resources. *Justice* uses historical job analysis and deadline information to assign the minimal fraction of resources required to meet a job’s deadline. *Justice* estimates this fraction from a

running tabulation of an expansion factor that it computes from an on-line, post-mortem analysis of all previous jobs executed.

We compare *Justice* to the baseline “fair” allocator employed by open-source systems like Mesos and YARN, to an extension of this allocator, and to an “oracle”, an allocator that cannot exist in practice, as it knows the exact minimum number of resources required for each job to meet its deadline (i.e., the “oracle” has access to information that real allocators can obtain only after jobs complete their executions). The metrics we use to do this comparison are two different shades of fairness, deadline satisfaction, productivity, wasted time and utilization. We use discrete-event, trace-driven simulation and deadline formulations from related work [7, 16, 20].

Our results show that *Justice* performs significantly better than the Mesos and YARN allocators and similarly to the oracle in terms of fairness, deadline satisfaction, and effective use of resources. This is because these resource negotiators attempt to preserve fairness without considering resource demand, which impedes performance when resources are constrained. We also find that *Justice* achieves greater productivity, wastes fewer resources, and has significantly better system utilization than its counter-parts for the workloads, deployment sizes, and deadlines that we consider.

The contributions of this paper are: a detailed analysis of job tracking, resource allocation, and admission control algorithms of *Justice*; an extended empirical evaluation of *Justice* against “fair-share” algorithms under different deadline types and cluster conditions on a YARN production trace that is 10x bigger in the number of submitted jobs and two times larger in terms of CPU cores, compared to the trace-based simulation performed in [4]; the introduction of a new fairness metric to distinguish “true fairness” from the traditional perception of fairness (“equality”). We define this metric using Jain’s fairness index [11] and compare the results of *Justice* and its competitive algorithms on our evaluation section.

## II. BACKGROUND AND RELATED WORK

In resource-constrained deployments encountered in the private cloud or the IoT, the fair-share policies [6, 8] fail to preserve fairness [3, 13]. Also, current fair allocators are deadline-agnostic. The kind of fairness they try to preserve is based on the resources requested by the job submitters and they ignore the actual demand of the job in order to meet its deadline. Moreover, to satisfy deadlines cluster administrators have to add extra resources, violate fairness by manually prioritizing one group of users over another, use solutions similar to

---

**Algorithm 1** *Justice* TRACK\_JOB Algorithm

---

```
1: function TRACK_JOB(requestedTasks, deadline, compTime,  
   numCPUsAllocd, success)  
2:   deadlineCPUs = compTime/deadline  
3:   maxCPUs = min(requestedTasks, cluster_capacity)  
4:   minReqRate = deadlineCPUs/maxCPUs  
5:   minReqRateList.add(minReqRate)  
6:   MinCPUFrac = min(minReqRateList)  
7:   MaxCPUFrac = max(minReqRateList)  
8:   LastCPUFrac = numCPUsAllocd/maxCPUs  
9:   LastSuccess = success  
10:  fractionErrorList.append(minReqRate - LastCPUFrac)  
11: end function
```

---

a capacity scheduler [1], or require users to reserve resources in advance [2, 17]. Such solutions are costly, inefficient, or impractical for resource constrained clusters as they further limit peak cluster capacity.

Moreover, in these multi-analytic settings, the intra-job allocators of frameworks like Hadoop and Spark greedily occupy resources even if they are not using them [3, 9, 22]. Authors in [9] attempt to address this issue by exploiting task-level resource requirements information and DAG dependencies. In contrast, *Justice* does not require job-repetitions and task-level information. Similarly to PYTHIA [5], *Justice* utilizes admission control to avoid wasting resources on infeasible jobs. But in addition, it monitors cluster conditions and adapts to changes in traffic patterns to avoid over-provisioning resources. This way it minimizes the amount of wasted resources and achieves “true fair-sharing” while it still satisfies more deadlines compared to fair-share approaches.

Much work [7, 18, 24] focuses on building performance profiles and scalability models offline or exploits historic and runtime information [10, 14, 20, 23, 24]. These approaches are not suitable for resource constrained, multi-analytics settings. Sampling, simulations, and extensive monitoring, impose overheads and additional cost. Also, trace analysis [4, 7] shows that some workloads have small ratio of repeated jobs with large execution time dispersion. Therefore, approaches based solely on past executions cannot predict with high statistical confidence for ad-hoc jobs or jobs that are not frequently repeated. Lastly, most of these approaches require task-level information, for the specific framework they target (e.g., Hadoop [10, 14, 23, 24] or Spark [21]) and consequently cannot be integrated into resource managers.

### III. JUSTICE ALGORITHM ANALYSIS

*Justice* can be conceptually separated into two main operations that it performs in parallel. The job tracking operation as described in Algorithm 1 and the admission control and resource allocation operation as described in Algorithm 2. Admission control and resource allocation depends on the online statistical model the tracking operation builds to estimate the resources jobs need to satisfy their deadlines (Function *alloc\_resources* in Algorithm 2) and to correct these estimations based on the Kalman filter mechanism described in Algorithm 3. The allocation operation depends on the quality of the statistical model for more accurate assignment of resources, but the two operations are decoupled so the allocator can still operate before there are sufficient historical statistics.

---

**Algorithm 2** Admission Control and Resource Allocation

---

```
1: function ADMISSION_CONTROL(RequesterJob)  
2:   for all  $j \in$  SubmittedJobs do  
3:     Feasible = True, TTD = Deadline - ElapsedTime  
4:     reqCpus = ESTIMATE_REQ( $j$ , TTD)  
5:     if reqCpus > min(taskCount, capacity) then  
6:       Feasible = False  
7:     end if  
8:     if Share( $j$ ) < reqCpus then  
9:       if Feasible == True then  
10:        priority = reqCpus/TTD, ADD2HEAP(priority,  $j$ )  
11:      else  
12:        DROP_JOB( $j$ )  
13:      end if  
14:    end if  
15:  end for  
16:  allocations = ALLOC_RESOURCES(heap)  
17:  if RequesterJob  $\notin$  allocations then  
18:    Add RequesterJob to queue  
19:  end if  
20: end function  
  
21: function ESTIMATE_REQ(Job)  
22:  maxCpus = min(tasks, capacity), reqCpus = maxCpus  
23:  if CompletedJobs > 1 then  
24:    fraction = CALCULATE_ALLOC_FRACTION()  
25:    fraction = CORRECT_ALLOC_FRACTION(fraction)  
26:    fraction = (deadline/(deadline - queue)) * fraction  
27:    reqCpus = max(ceil(fraction * maxCpus), 1)  
28:  end if  
29:  return reqCpus  
30: end function  
  
31: function ALLOC_RESOURCES(heap)  
32:  offers = CREATE_OFFERS(heap)  
33:  allocations = SEND_OFFERS(offers)  
34:  return allocations  
35: end function  
  
36: function CREATE_OFFERS(heap)  
37:  while availableCpus > 0 and heap not empty do  
38:    for all  $Job\ j \in$  heap do  
39:      offer = min(request( $j$ ), availableCpus)  
40:      if offer < request( $j$ ) then  
41:        offer = 0  
42:      else  
43:        availableCpus - = offer  
44:        offersDict[ $j$ ] = offer  
45:      end if  
46:    end for  
47:  end while  
48:  return offersDict  
49: end function
```

---

---

**Algorithm 3** Allocation Calculation and Correction

---

```
1: function CALCULATE_ALLOC_FRACTION  
2:   if LastSuccess then  
3:     CPUFrac = MinCPUFrac  
4:   else  
5:     CPUFrac = MaxCPUFrac  
6:   end if  
7:   fraction = (LastCPUFrac + CPUFrac)/2  
8:   return fraction  
9: end function  
10: function CORRECT_ALLOC_FRACTION(fraction)  
11:  correction = CALC_SMOOTHED_AVG(fractionErrorList)  
12:  correctedFraction = fraction + correction  
13:  correctedFraction = VALIDATE_FRACTION(correctedFraction)  
14:  return correctedFraction  
15: end function
```

---

This design allows *Justice* to perform these operations without delaying job scheduling as it can run on the background to calculate the new allocation fractions similarly to other window-based fairness algorithms. Its memory requirements scale linearly to the points used on the desired history window and, depending on the desired scheduling latency and estimation accuracy requirements, a limited amount of metrics could be stored in memory and the rest in database.

#### A. Job Metrics Tracking

*Justice* invokes Algorithm 1 every time a job completes its execution. The algorithm takes as inputs a number of metrics provided by the user at submission time and a number of metrics extracted by the job’s execution profile that is available in system logs after the job completes. In return, the algorithm produces a number of metrics that are used as inputs to the functions of Algorithm 2 and Algorithm 3. The inputs of the algorithm are the number of tasks a job has (`requestedTasks`), its deadline as defined in seconds by the job submitter, the total computation time of the job (`compTime`) expressed in `CPU*Seconds`, the number of CPUs allocated to the job (`numCPUsAllocd`), and a boolean `success` indicating whether the job was successful (completed its job before its deadline) or failed (exceeded its deadline or got dropped before completing).

Based on these five inputs, the algorithm derives a number of intermediate metrics for all jobs in order to produce its final results. These metrics are, the minimum number of CPUs a job would have needed to finish by its deadline (`deadlineCPUs`), the maximum parallelism of the job (`maxCPUS`), the minimum required rate (`minReqRate`) which is the fraction of resources the job would have needed compared to its maximum resources, in order to meet its deadline just in time. These rates are stored for all jobs encountered in the system (`minReqRateList`).

After producing these derived metrics, the algorithm calculates the desired results. These are `MinCPUFrac` and `MaxCPUFrac`, which correspond to the minimum and maximum request rates encountered across all jobs respectively and the `LastCPUFrac`, which is the last given rate observed for the job that completed and triggered the algorithm. It also stores whether the last job completed successful `LastSuccess`. Lastly, a historic fraction error is produced across all jobs `fractionErrorList` as the difference between the minimum required rate the job would have needed to meet its deadline and the fraction of resources *Justice* assigned to it. Note that `deadlineCPUs` cannot be greater than `maxCPU` (assuming feasible deadlines) and `MinCPUFrac` and `MaxCPUFrac` are always less or equal to 1.

#### B. Resource Estimation, Admission Control and Allocation

After a bootstrapping period in which the job tracking operation runs without any admission control in order to gather enough data and produce estimations with statistical significance, the admission control and resource allocation mechanism of *Justice*’s kick in. Algorithm 2 takes as input a submitted job (`RequesterJob`), and based on the metrics that the job tracking operation continuously produces and stores, as discussed on the previous section, it creates an

allocation for the job. This allocation should be sufficient to meet its deadline just in time or if, based on the cluster conditions, it estimates that it is impossible for the job to complete before its deadline, then, it drops the job.

1) *Resource Estimation*: To achieve this, it updates the deadlines for jobs in the queue, reducing each by the time that has passed since submission (line 3 in Algorithm 2). Then, it estimates the minimum amount of resources the job requires to meet its deadline (Function `estimate_req` in Algorithm 3). That happens by a subsequent call to Function `calc_alloc_frac` in Algorithm 3 that computes the CPU allocation fraction (`fraction`) for each newly submitted job as the average of the `LastCPUFrac` and either `MinCPUFrac` or `MaxCPUFrac`, as shown in Algorithm 3, depending on whether the last completed job met or missed its deadline, respectively. In other words, consecutive successes make *Justice* more aggressive, causing it to allocate smaller resource fractions (i.e., `fraction` converges to `MinCPUFrac`), while deadline violations make *Justice* more conservative, causing it to increase the fraction in an attempt to prevent future violations (`fraction` converges to `MaxCPUFrac`).

The fraction produced by Function `calc_alloc_frac` is further improved by a Kalman filter mechanism (Function `admission_control`). Every time a job completes its execution, *Justice* tracks the estimation error and uses it to correct the CPU allocation fraction. Estimation error is the difference between the allocation fraction and the ideal minimum fraction (`deadlineCPUs`). *Justice* calculates a weighted average of the historical errors (Function `correct_alloc_fraction`) and adds it to the allocation fraction. *Justice* can be configured to assign the same weights to all past errors or to use exponential smoothing (i.e., to weigh recent values higher than those that occurred in the distant past). Lastly, a validate function (that we do not include on the algorithm for brevity) ensures that the corrected fraction remains within allowable limits (no less than the minimum observed `MinCPUFrac` or greater than 1).

After *Justice* computes, corrects, and validates `allocCPUFrac`, it considers the time that the job has spent in the queue (line 26 in Function `estimate_req` of Algorithm 2). *Justice* multiplies `allocCPUFrac` by the number of tasks requested on job submission and uses this value as the number of CPUs to assign to the job (Function `estimate_req` in Algorithm 2).

2) *Admission Control*: *Justice* recomputes the CPU allocation of each enqueued job and, as part of its admission control policy, it either drops the ones with infeasible deadlines or keeps those that cannot be admitted but are still feasible (lines 12 and 18 respectively in Algorithm 2). *Justice* implements a *proactive* admission control mechanism to prevent jobs likely to miss their deadline from entering the system and consuming resources wastefully. This way, *Justice* attempts to maximize the number of jobs that meet their deadline even under severe resource constraints. *Justice* also tracks jobs that violate their deadlines and selectively drops some of them to avoid further waste of resources. It is selective in that it terminates jobs when their `requestedTasks` exceed a configurable threshold. Thus, it still able to collect statistics on “misses” to improve its estimations by letting the smaller violating jobs complete

their execution while at the same time it prevents the bigger violators from wasting resources.

The priority policy *Justice* uses is pluggable. In the evaluations of this paper we use a policy that aims to minimize the number of jobs that miss their deadlines. For this policy (line 10 in Algorithm 2), *Justice* prioritizes jobs with a small number of tasks and greatest time-to-deadline (TTD). However, all of the policies that we considered (including shortest time-to-deadline) perform similarly. Once *Justice* has selected a job for admission, it allocates the CPUs to the job and admits it to the system for execution. Once a job run commences, its CPU allocation does not change.

3) *Resource Allocation*: Finally, *Justice* allocates the calculated resources to jobs (Function `alloc_resources` in Algorithm 2) by creating offers according to job priorities. Function `offer_resources` creates offers for jobs until there are no other jobs to be scheduled or the available resources are exhausted. Lastly, *Justice* sends these offers to the frameworks (line 33 in Algorithm 2 - we omit Function `SEND_OFFERS` for brevity).

#### IV. EXPERIMENTAL METHODOLOGY

To evaluate *Justice* we use a discrete-event simulator written in Python with industry-provided production traces. The trace used in this paper is from a YARN cluster with 25 thousand CPU cores and more than 1 million jobs over a 3-month period. Approximately 60% of the jobs have a single task and 70-80% of the jobs have fewer than 10 tasks. 25% of the jobs repeat more than once and 16% of the jobs repeat more than 30 times. We omit further details of our test bed due to space limitations as they are similar to [4].

We compare *Justice* against the fair-share allocator implemented in open-source resource negotiators like Mesos and YARN. We refer to this allocator as `Baseline FS`. This allocator lacks any deadline information and therefore executes a job even after its deadline is exceeded. We also implement an extension of this allocator named `Reactive FS`, that enforces the same FS policy but reactively terminates a job that has exceeded its deadline. Lastly, we implement an “oracle” allocator that knows the exact amount of resources a job requires to meet its deadline without having, however, knowledge of the optimal schedule.

We evaluate the robustness of our approach by running experiments using deadline formulations from prior works [7, 16, 20] and variations on them. In particular, we assign deadlines that are multiples of the optimal execution time of a job as we extract it from our workload traces. We use two types of multiples: Fixed and variable.

**Fixed Deadlines**: With fixed deadlines, we use a deadline that is a multiple of the optimal execution time as described in [16]. Each deadline is expressed as  $D_i = x \cdot T_i$ , where  $T_i$  is the optimal runtime of the job and  $x \geq 1.0$  is some fixed multiplicative expansion factor. In our experiments, we use constant factors of  $x = 1$  and  $x = 2$ , which we refer to as *Fixed1x* and *Fixed2x* respectively.

**Variable Deadlines**: For variable deadlines, we compute deadline multiples by sampling distributions. *Jockey* deadlines pick randomly a factor  $x$  from two possible values as

described in [7]. In this work, we use the intervals from the sets with values (1, 2) and (2, 4) to choose  $x$  and, again, compute  $D_i = x \cdot T_i$ , where  $T_i$  is the minimum possible execution time. We refer to this variable deadline formulations as *Jockey1x2x* and *Jockey2x4x*. *90loose* are a variation of the *Jockey1x2x* deadlines, in which the deadlines take on the larger value (i.e. are loose) with a higher probability (0.9) while the other uses the smaller value. *Aria* deadlines are uniformly distributed in the intervals [1, 3] and [2, 4] as described in [20]; we refer to these deadlines as *Aria1x3x* and *Aria2x4x*, respectively.

#### V. RESULTS

We compare *Justice* in terms of fairness, deadline satisfaction, and effective resource utilization, for different resource-constrained cloud deployments, against fair share schedulers and an oracle using multiple deadline formulations, as described on Section IV.

##### A. Fairness Evaluation

Traditional fair-share allocators [6, 8] attempt to give an “equal” share of resources to concurrently executing jobs regardless of whether this share is sufficient to meet their deadlines. Herein, we will refer to this form of fairness as “equality”. To evaluate, the degree to which these allocators achieve this goal in resource-constrained settings, we use Jain’s fairness index  $\frac{|\sum_{i=1}^n F_i|^2}{n * \sum_{i=1}^n F_i^2}$  with  $F_i$  corresponding to the resource allocation of each job  $i$ .

To compute equality, we classify jobs based on their maximum demand. We then calculate the index for each job and the weighted average across indexes. Weights correspond to the number of jobs in each class (e.g., all jobs with demand of  $Y$  CPUs). We classify jobs in this way to avoid considering “unfair” (or “unequal”) allocations that correspond to different maximum demand classes because a job cannot be allocated more CPUs than it demands.

The top graphs of Figure 1 present equality results across all allocators and cluster capacities that we consider. The fairness index is averaged over 60-sec intervals on the lifetime of the workload. *Justice* achieves better fairness scores than the fair-share allocators by up to 23% and 17% for the two capacities. Even though the goal of *Justice* is not to preserve equality but instead to prioritize for what we consider actual fairness, it performs better than the fair-share allocators for two reasons. First, *Justice* keeps the system less utilized and therefore fewer jobs wait in the queue, which contributes negatively to equality. Second, due to constrained resources, *Justice* drops large jobs more frequently which provides opportunities for it to facilitate fairness at a finer grain across frameworks.

We argue that “equality” is not the desired property for deadline-driven workloads. Equality treats all jobs similarly regardless of their actual resource requirements. In practice, jobs have different priorities, max demands, and diverse deadline tightness. Instead, “true” fairness can be measured by using Jain’s fairness index with  $F_i$  corresponding to the fraction of demand of each job  $i$ . For each job  $i$ , among  $n$  total jobs, we define the fraction of demand as  $F_i = \frac{A_i}{D_i}$  where  $D_i$  is the resource request for job  $i$  and  $A_i$  is the allocation given to job  $i$ . When  $A_i \geq D_i$  the fraction is defined to be 1.

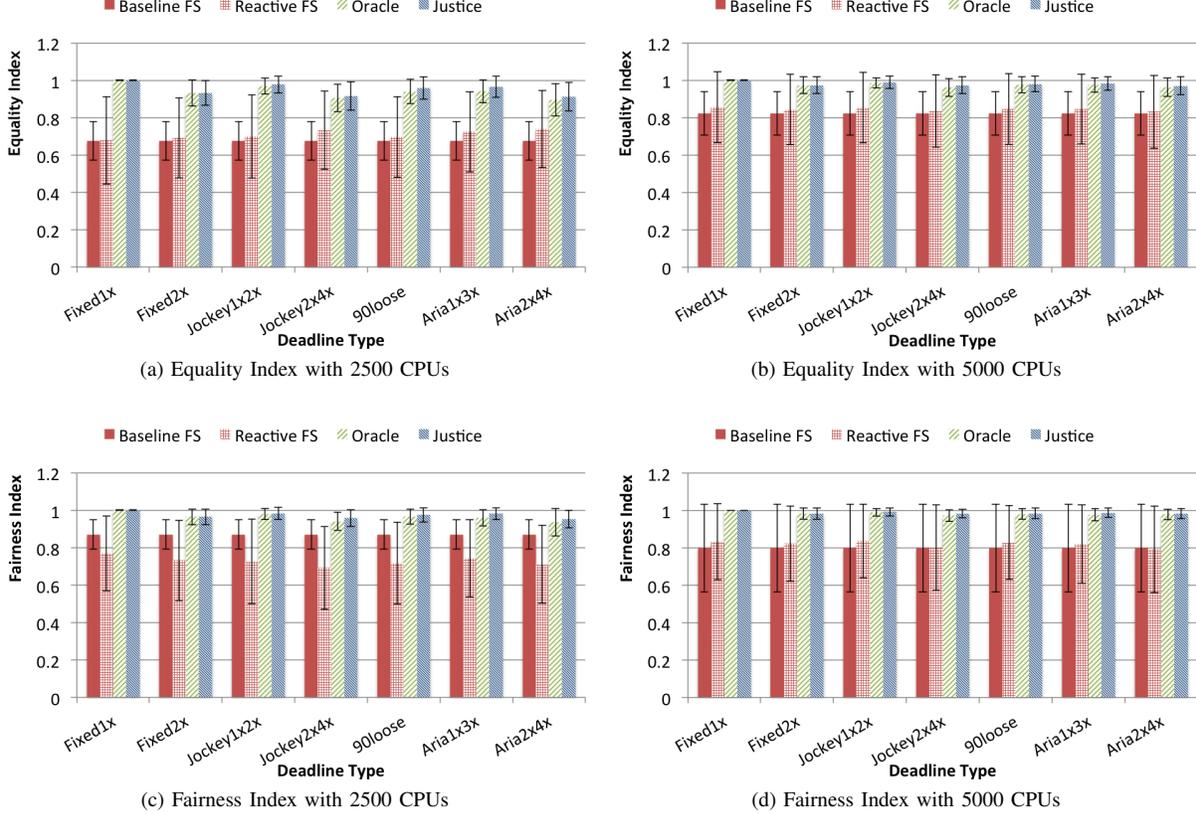


Fig. 1: **Equality Vs Fairness:** Average of Jain’s fairness index adapted for equality (top graphs) and fairness (bottom graphs) with highly constrained capacities (left graphs) and moderately constrained capacities (right graphs). Experiments denoted as “Fixed” have deadlines multiples of 1 and 2. Experiments denoted as “Jockey” have multiples picked randomly from a set with two values (1, 2) and (2, 4). Experiments denoted as “90loose” have 90% deadlines with a multiple of 2 and 10% deadlines with a multiple of 1. Experiments denoted as “Aria” have multiples drawn from uniformly distributed intervals [1, 3] and [2, 4]

By using this “true” fairness metric, *Justice* outperforms all “equality” allocators we evaluate in this study in clusters with constrained resources. It achieves this by applying admission control instead of greedily allocating resources to jobs and by predicting the amount of resources jobs require to meet their deadlines “just in time”. In contrast, the existing fair-share allocators cannot be fair under constrained resources as they cannot prevent larger jobs from taking over a significant portion of the cluster [3, 9, 22].

In addition, *Justice* outperforms the oracle for variable deadlines. This is an artifact of the oracle’s use of maximum job demand in the formula instead of the minimum required resources. Under this definition, our oracle is not a fairness oracle in terms of preserving fairness globally on the system. It is instead an oracle with respect to the minimum resource requirements needed to satisfy each job’s deadline.

### B. Deadline Satisfaction

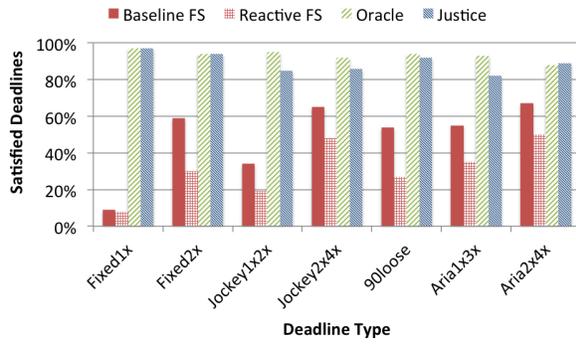
Being just fair in deadline-driven workloads is not enough. The main goal of a resource allocator in such settings is to satisfy deadlines. To investigate this, we compute the *Satisfied Deadline Ratio (SDR)* as the fraction of the jobs that complete before their deadline over the total number of submitted jobs.

Figures 2a and 2b show that fair-share allocators, fail to satisfy job deadlines as they lack deadline information and assign resources solely based on what they consider as “fair”. In contrast, *Justice* builds a statistical model based on previous job executions and assigns the amount of resources jobs need to satisfy their deadlines “just-in-time”. As a result, *Justice* achieves 80% of optimal allocation. Note that even the oracle doesn’t achieve a perfect SDR ratio, because it does not have knowledge of the perfect global schedule. Therefore, it also has to drop jobs that cannot achieve their deadlines.

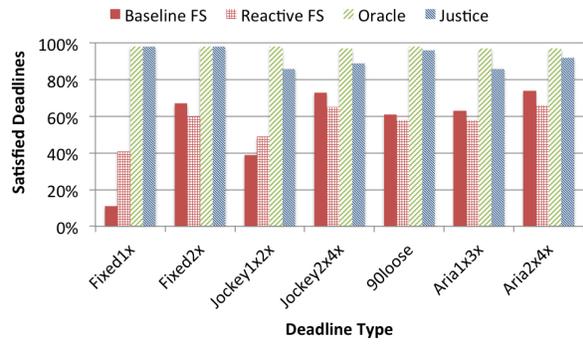
### C. Effective Resource Usage and Cluster Utilization

We next evaluate the resource allocators using three metrics that in combination show how effectively each utilizes cluster resources. For the set of submitted jobs  $J_1, J_2, \dots, J_n$  and their corresponding runtimes  $T_1, T_2, \dots, T_n$ , we consider the subset of  $m < n$  successful jobs  $J_1, J_2, \dots, J_m$  and the subset of  $k < n$  failed or dropped jobs  $J_1, J_2, \dots, J_k$  where  $n = m + k$ .

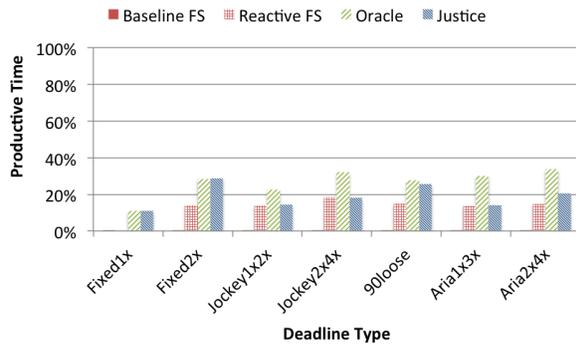
We define Productive Time Ratio (PTR) as  $\frac{\sum_{i=1}^m T_i}{\sum_{j=1}^n T_j}$  and Wasted Time Ratio WTR as  $\frac{\sum_{i=1}^k T_i}{\sum_{j=1}^n T_j}$ . Lastly, cluster utilization is  $\frac{busy}{idle+busy}$  where *busy* is the total busy time and *idle* is the total idle time across a workload.



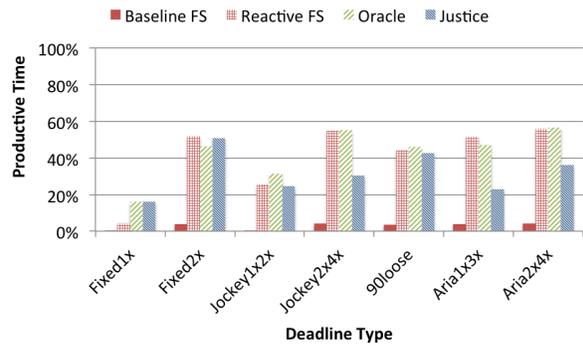
(a) Satisfied Deadlines with 2500 CPUs



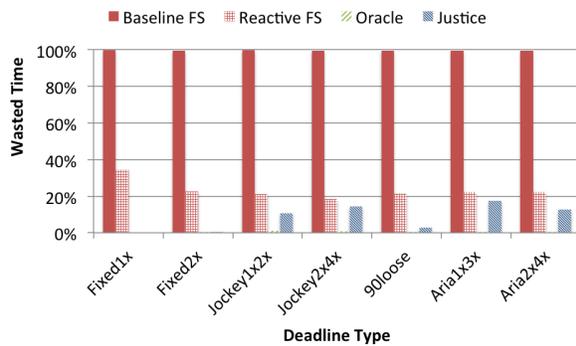
(b) Satisfied Deadlines with 5000 CPUs



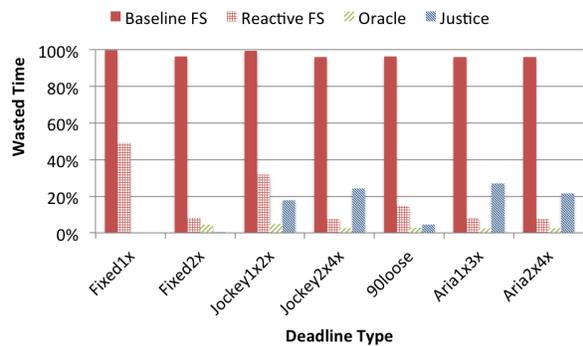
(c) Productive Time with 2500 CPUs



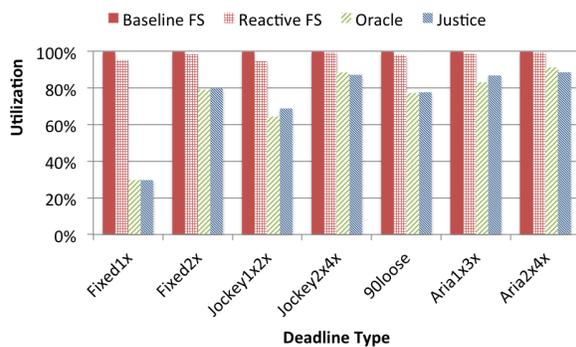
(d) Productive Time with 5000 CPUs



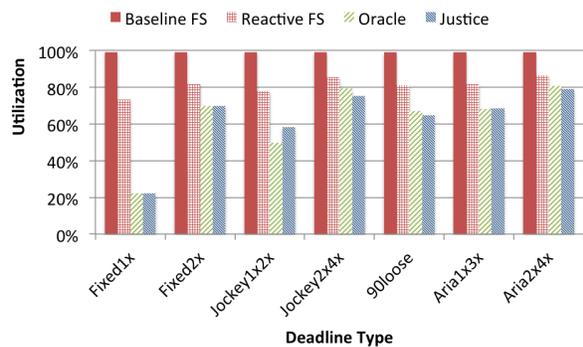
(e) Wasted Time with 2500 CPUs



(f) Wasted Time with 5000 CPUs



(g) Utilization with 2500 CPUs



(h) Utilization with 5000 CPUs

Fig. 2: **Deadline Satisfaction and Efficient Resource Utilization:** Satisfied Deadlines Ratio (SDR), Productive Time Ratio (PTR), Wasted Time Ratio (WTR), and cluster utilization with highly constrained cluster capacities (left graphs) and moderately constrained capacities (right graphs) for different deadline types.

The left side on Figure 2 shows that, under severe resource constraints, *Justice* spends more time productively and wastes fewer resources. It does so by dropping jobs that are likely to violate their deadlines according to its predictions. In contrast, fair-share policies attempt to share resources equally between smaller and bigger jobs. When resources are constrained, this share is insufficient for the bigger jobs to meet their deadlines. Moreover, Baseline FS wastes time on jobs that have already missed their deadlines, while Reactive FS avoids doing so by retroactively dropping such jobs.

*Justice* outperforms the existing fair-share allocators in all these metrics, while at the same time it satisfies more deadlines and achieves better fairness as discussed in Sections V-B and V-A. This means that the smaller utilization is not a by-product of added overhead but the result of effective admission control that filters out jobs that would not satisfy their deadlines under these constrained cluster resources.

The right side of Figure 2 shows that fair-share policies might be more suitable for optimizing productive work for clusters for which resource scarcity is not severe. In such conditions, and in combination with higher deadline variability, *Justice* might deny admission to some bigger jobs in order to preserve fairness and to satisfy deadlines for other (smaller) jobs. This effect is depicted both in a smaller PTR value (Figure 2d) and lower cluster utilization (Figure 2h).

## VI. CONCLUSIONS

*Justice* is a scheduler designed for clouds with constrained resources, commonly found in IoT and private clouds. It tracks job deadlines and runtime information to adapt its resource allocation and admission control mechanisms so it can achieve fairness and satisfy deadlines even when resource availability is scarce. We analyze *Justice*'s algorithm and empirically evaluate it using discrete-event simulation of deadline-driven, production workloads in resource-constrained clusters. We compare *Justice* to the existing fair-share allocator that ships with Mesos and YARN and find that *Justice* is able to achieve better "traditional" (equality) as well as "true" fairness, deadline satisfaction, and better resource utilization.

This work is supported by NSF (CNS-1703560, CCF-1539586, ACI-1541215) and AWS Cloud Credits for Research.

## REFERENCES

- [1] *YARN Capacity Scheduler*. <https://hadoop.apache.org/docs/r2.7.1/hadoop-yarn/hadoop-yarn-site/CapacityScheduler.html>.
- [2] C. Curino et al. Reservation-based Scheduling: If You're Late Don't Blame Us! In: *ACM Symposium on Cloud Computing*. 2014, pp. 1–14.
- [3] S. Dimopoulos, C. Krintz, and R. Wolski. Big Data Framework Interference In Restricted Private Cloud Settings. In: *IEEE International Conference on Big Data*. IEEE. 2016.
- [4] S. Dimopoulos, C. Krintz, and R. Wolski. Justice: A Deadline-aware, Fair-share Resource Allocator for Implementing Multi-analytics. In: *Cluster Computing (CLUSTER), 2017 IEEE International Conference on*. IEEE. 2017, pp. 233–244.
- [5] S. Dimopoulos, C. Krintz, and R. Wolski. PYTHIA: Admission Control for Multi-Framework, Deadline-Driven, Big Data Workloads. In: *International Conference on Cloud Computing*. IEEE. 2017.
- [6] *YARN Fair Scheduler*. <https://hadoop.apache.org/docs/r2.4.1/hadoop-yarn/hadoop-yarn-site/FairScheduler.html>.
- [7] A. D. Ferguson et al. Jockey: guaranteed job latency in data parallel clusters. In: *ACM European Conference on Computer Systems*. ACM. 2012, pp. 99–112.
- [8] A. Ghodsi et al. Dominant resource fairness: Fair allocation of multiple resource types. In: *NSDI*. 2011.
- [9] R. Grandl et al. Altruistic scheduling in multi-resource clusters. In: *USENIX Symposium on Operating Systems Design and Implementation*. 2016.
- [10] Z. Huang et al. RUSH: A RobUst ScHeduler to Manage Uncertain Completion-Times in Shared Clouds. In: *2016 IEEE 36th International Conference on Distributed Computing Systems (ICDCS)*. IEEE. 2016, pp. 242–251.
- [11] R. Jain, D.-M. Chiu, and W. R. Hawe. *A quantitative measure of fairness and discrimination for resource allocation in shared computer system*. Vol. 38. Eastern Research Laboratory, Digital Equipment Corporation Hudson, MA, 1984.
- [12] K. Kc and K. Anyanwu. Scheduling hadoop jobs to meet deadlines. In: *International Conference on Cloud Computing*. 2010, pp. 388–392.
- [13] J. Khamse-Ashari et al. An efficient and fair multi-resource allocation mechanism for heterogeneous servers. In: *IEEE Transactions on Parallel and Distributed Systems* 29.12 (2018), pp. 2686–2699.
- [14] P. Lama and X. Zhou. Aroma: Automated resource allocation and configuration of mapreduce environment in the cloud. In: *ACM International Conference on Autonomic Computing*. 2012, pp. 63–72.
- [15] S. Li et al. WOHA: deadline-aware map-reduce workflow scheduling framework over hadoop clusters. In: *Distributed Computing Systems (ICDCS), 2014 IEEE 34th International Conference on*. IEEE. 2014, pp. 93–103.
- [16] J. Liu, H. Shen, and H. S. Narman. CCRP: Customized Cooperative Resource Provisioning for High Resource Utilization in Clouds. In: *IEEE International Conference on Big Data*. IEEE. 2016.
- [17] A. Tumanov et al. TetriSched: global rescheduling with adaptive plan-ahead in dynamic heterogeneous clusters. In: *European Conference on Computer Systems*. 2016, p. 35.
- [18] S. Venkataraman et al. Ernest: efficient performance prediction for large-scale advanced analytics. In: *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*. 2016, pp. 363–378.
- [19] T. Verbelen et al. Cloudlets: bringing the cloud to the mobile user. In: *ACM workshop on Mobile cloud computing and services*. ACM. 2012.
- [20] A. Verma, L. Cherkasova, and R. H. Campbell. ARIA: automatic resource inference and allocation for mapreduce environments. In: *ACM International Conference on Autonomic Computing*. 2011, pp. 235–244.
- [21] K. Wang and M. M. H. Khan. Performance Prediction for Apache Spark Platform. In: *2015 IEEE 17th International Conference on High Performance Computing and Communications (HPCC)*. IEEE. 2015, pp. 166–173.
- [22] Y. Yao et al. Admission control in YARN clusters based on dynamic resource reservation. In: *IEEE International Symposium on Integrated Network Management*. 2015, pp. 838–841.
- [23] N. Zafeilas and V. Kalogeraki. Real-time scheduling of skewed mapreduce jobs in heterogeneous environments. In: *11th International Conference on Autonomic Computing (ICAC 14)*. 2014, pp. 189–200.
- [24] W. Zhang et al. Mimp: Deadline and interference aware scheduling of hadoop virtual machines. In: *IEEE Cluster, Cloud and Grid Computing*. 2014, pp. 394–403.