



HAL
open science

Revisiting I/O bandwidth-sharing strategies for HPC applications

Anne Benoit, Thomas Herault, Lucas Perotin, Yves Robert, Frédéric Vivien

► **To cite this version:**

Anne Benoit, Thomas Herault, Lucas Perotin, Yves Robert, Frédéric Vivien. Revisiting I/O bandwidth-sharing strategies for HPC applications. RR-9502 v3, INRIA. 2023, pp.57. hal-04038011v3

HAL Id: hal-04038011

<https://inria.hal.science/hal-04038011v3>

Submitted on 15 Aug 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Inria

Revisiting I/O bandwidth-sharing strategies for HPC applications

Anne Benoit, Thomas Herault, Lucas Perotin, Yves Robert, Frédéric Vivien

**RESEARCH
REPORT**

N° 9502

March 2023

Project-Team ROMA

ISRN INRIA/RR--9502--FR+ENG

ISSN 0249-6399



Revisiting I/O bandwidth-sharing strategies for HPC applications

Anne Benoit*, Thomas Herault†, Lucas Perotin‡, Yves Robert§,
Frédéric Vivien¶

Project-Team ROMA

Research Report n° 9502 — version 3 — initial version March 2023 —
revised version August 2023 — 57 pages

Abstract: This work revisits I/O bandwidth-sharing strategies for HPC applications. When several applications post concurrent I/O operations, well-known approaches include serializing these operations (FCFS) or fair-sharing the bandwidth across them (FAIRSHARE). Another recent approach, I/O-Sets, assigns priorities to the applications, which are classified into different sets based upon the average length of their iterations. We introduce several new bandwidth-sharing strategies, some of them simple greedy algorithms, and some of them more complicated to implement, and we compare them with existing ones. Our new strategies do not rely on any a-priori knowledge of the behavior of the applications, such as the length of work phases, the volume of I/O operations, or some expected periodicity. We introduce a rigorous framework, namely *steady-state windows*, which enables to derive bounds on the competitive ratio of all bandwidth-sharing strategies for three different objectives: minimum yield, platform utilization, and global efficiency. To the best of our knowledge, this work is the first to provide a quantitative assessment of the online competitiveness of any bandwidth-sharing strategy. This theory-oriented assessment is complemented by a comprehensive set of simulations, based upon both synthetic and realistic traces. The main conclusion is that two of our simple and low-complexity greedy strategies significantly outperform FCFS, FAIRSHARE and I/O-Sets, and we recommend that the I/O community would implement them for further assessment.

Key-words: I/O, bandwidth sharing, scheduling strategy, HPC applications

* Univ Lyon, EnsL, UCBL, CNRS, Inria, LIP, F-69342, LYON Cedex 07, France.

† University of Tennessee, Knoxville, TN, USA

‡ Univ Lyon, EnsL, UCBL, CNRS, Inria, LIP, F-69342, LYON Cedex 07, France.

§ Univ Lyon, EnsL, UCBL, CNRS, Inria, LIP, F-69342, LYON Cedex 07, France and University of Tennessee, Knoxville, TN, USA

¶ Univ Lyon, EnsL, UCBL, CNRS, Inria, LIP, F-69342, LYON Cedex 07, France.

RESEARCH CENTRE
GRENOBLE – RHÔNE-ALPES

Inovallée
655 avenue de l'Europe Montbonnot
38334 Saint Ismier Cedex

Stratégies de partage de la bande passante d'entrée-sortie entre applications de calcul haute performance

Résumé : Ce travail revisite les stratégies de partage de la bande passante d'entrée-sortie entre applications de calcul haute performance. Quand plusieurs applications postent simultanément des opérations d'entrée-sortie, les approches classiques incluent la sérialisation de ces opérations (FCFS) et le partage équitable de la bande passante (FAIRSHARE). Une approche récente, *I/O-Sets*, attribue des priorités aux applications, qui sont classées en différents ensembles basés sur la longueur moyenne de leurs itérations.

Nous introduisons plusieurs stratégies nouvelles de partage de la bande passante, de simples heuristiques gloutonnes et des stratégies plus compliquées à implémenter, et nous les comparons aux solutions préexistantes. Nos stratégies n'utilisent aucune connaissance *a priori* du comportement des applications, telle que la longueur des phases de calcul, le volume des entrées-sorties, ou leur périodicité.

Nous introduisons un cadre rigoureux, les fenêtres de régime permanent, qui permet de définir des bornes sur le facteur de compétitivité de toutes les stratégies de partage de la bande passante, et ce pour trois objectifs: le *yield* (rendement) minimal, l'utilisation de la plateforme, et l'efficacité globale. À notre connaissance, ce travail est le premier à proposer une évaluation quantitative du facteur de compétitivité de stratégies de partage de la bande passante. Cette évaluation théorique est complétée par un ensemble de simulations, utilisant des traces synthétiques et d'autres réalistes. La principale conclusion de ces simulations est que deux de nos simples heuristiques gloutonnes de faible complexité obtiennent de bien meilleures performances que FCFS, FAIRSHARE et *I/O-Sets*. Nous recommandons donc que la communauté les implémente afin de permettre leur évaluation en pratique.

Mots-clés : Entrées-sorties, partage de bande passante, ordonnancement, calcul haute performance

Contents

1	Introduction	4
2	Related Work	6
3	Framework	8
3.1	Applications	9
3.1.1	Application Characteristics	9
3.1.2	Bandwidth Allocation	10
3.2	Steady-State Windows	11
3.2.1	Interaction with the Batch Scheduler	11
3.2.2	Cost Model for Steady-State Windows	12
3.3	Objectives	13
4	Bandwidth-Sharing Strategies	15
4.1	Greedy Strategies	15
4.2	SET-10 Strategy	16
4.3	Maximizing the Minimum Yield at the Next Event	17
5	Lower Bounds on Competitive Ratios	19
5.1	Example 1	20
5.2	Example 2	21
5.3	Example 3	22
5.4	Example 4	23
5.5	Example 5	23
5.6	Example 6	25
5.7	Tight Bounds	27
5.7.1	The competitive ratio of FAIRSHARE is exactly m for MINYIELD	27
5.7.2	The competitive ratios of FCFS, SET-10, GREEDY YIELD, LOOKAHEAD-GREEDY YIELD, and PERIODICGREEDY YIELD are exactly m for EFFICIENCY	27
5.7.3	FAIRSHARE can be arbitrarily better than BESTNEXT EVENT	28
6	Performance Evaluation	28
6.1	I/O Pressure	28
6.2	Synthetic Traces	28
6.2.1	Framework	28
6.2.2	Results for Synthetic Traces	29
6.2.3	Synthesis of the Evaluation on Synthetic Scenarios	33
6.3	Evaluation on APEX workloads	33
6.3.1	Apex Traces [17]	33
6.3.2	MINYIELD of FAIRSHARE on APEX Scenarios	34
6.3.3	MINYIELD of All Strategies on APEX NERSC Scenarios	36
6.3.4	MINYIELD of All Strategies on APEX TRILAB Scenarios	36
6.3.5	EFFICIENCY of All Strategies on APEX Scenarios	38
6.3.6	UTILIZATION of All Strategies on APEX Scenarios	39
6.3.7	Computation Time of All Strategies on APEX Scenarios	40
6.3.8	Synthesis of the Evaluation on APEX Scenarios	40
7	Conclusion	41

A Detailed description of BESTNEXTEVENT	45
A.1 Maximizing the minimum yield in an interval	46
A.2 Maximizing the minimum yield in the whole execution window	50
A.3 Maximizing the minimum yield lexicographically at time $t + u$	50
B Additional Simulation Results	54
B.1 Synthetic Traces	54
B.2 APEX Traces	55

1 Introduction

HPC applications do not share computing resources: all the nodes assigned to a given application are dedicated to that application throughout its execution. Such a mode of operation is enforced to guarantee a sustained level of performance to all applications that execute concurrently on the platform. However, concurrent applications do share both the interconnexion network and the parallel file system. When several applications request to perform an I/O operation simultaneously, they have to share the resource, which leads to interferences and performance degradation.

Several researchers have already identified and addressed this problem (see [7, 12, 25, 2, 28, 27] among others). Performance degradation due to I/O is already significant for current state-of-the-art platforms and is expected to worsen due to the faster increase in processing speed than in I/O bandwidth [22]. The problem can be partially mitigated by reducing the volume of data transfers, e.g., via compression or in-situ processing. But the main question remains: given several applications executing concurrently and competing for I/O resources, how to orchestrate I/O operations? In other words, scheduling strategies must be designed and evaluated to dynamically assign a fraction of the total I/O bandwidth to individual application transfers. Well-known strategies are FCFS, which gives exclusive I/O access to the first pending I/O operation, and FAIRSHARE, which assigns bandwidth proportionally to application transfers.

From a scheduling perspective, which fraction goes to which application at any given time depends upon the optimization metric, such as application progress rate (minimum or average) or platform utilization. When targeting fairness across concurrent applications, a classical objective is to maximize the minimum *yield*, where the yield of an application is the ratio of its actual progress rate over the progress rate that would have been achieved if the application was executing with a dedicated I/O system and always granted the total available bandwidth. We discuss optimization metrics in detail in Section 3.3.

This work focuses on I/O bandwidth-sharing scheduling strategies for HPC applications, revisiting existing strategies and introducing new ones. Our major contributions are described in the following four paragraphs.

General framework We provide and assess online scheduling strategies that are agnostic of the characteristics of the concurrent applications in terms of processing time and I/O requests. In particular, we do not assume any periodic behavior; several applications execute concurrently and alternate phases of work and phases of I/O operations, whose lengths are not known a priori. Instead, we discover the timing and size of I/O transfers on the fly, as each application posts its operations. We allow for interrupting and resuming on-going I/O operations dynamically, and launching newly posted ones.

Novel strategies We introduce novel I/O bandwidth-sharing strategies that aim at allocating a fraction of the bandwidth to each application as a function of the current progress of all

applications. The main motivation is to maximize the minimum yield that can be achieved each time a scheduling decision is made. These novel heuristics come in several flavors, from simple greedy algorithms to sophisticated decision mechanisms.

Competitiveness analysis We provide a rigorous framework by focusing on a *steady-state* time window, defined with the following three rules. Throughout the window: (i) several applications, each with a processing history, execute concurrently; (ii) none of them terminates; and (iii) no new application can start. Thus, the window corresponds to a steady-state mode of behavior where each application progresses at the rate enforced by the I/O bandwidth-sharing strategy. Focusing on such a window is key to assess performance. Otherwise, say if some application would terminate before the end of the window, the batch scheduler would likely launch a new application, whose starting time and progress up to the end of the window would depend on all previous scheduling decisions. The same holds if a new application is launched in the middle of the window. Getting rid of the interaction with the batch scheduler, we provide the first complexity results on the performance of several I/O bandwidth-sharing strategies, some old and some new, and for various optimization objectives.

Comprehensive simulation campaign We compare existing and novel I/O bandwidth-sharing strategies on an extensive set of application scenarios, some generated from realistic traces derived from the APEX workflows report [17], and some with synthetic parameters. A key parameter is the I/O pressure W , defined for a steady-state window $[T_{begin}, T_{end}]$, as the ratio $\frac{V}{B(T_{end}-T_{begin})}$, where (i) V is the total I/O volume (accumulated for all applications) to transfer during the window; and (ii) B is the total I/O bandwidth (see Section 6.1 for details). In a nutshell, if this ratio is close to 1 or even exceeds 1, the set of I/O operations saturate the I/O system, and many I/O operations will have to be delayed. We study how rapidly the performance of each strategy degrades for high I/O pressures, thereby paving the way for a fair bandwidth allocation on future platforms. We point out that simulations are a first but mandatory step to assess the limitations and strengths of all the I/O bandwidth-sharing strategies. Our extensive set of experiments corresponds to several months of platform usage and would have been impossible to deploy on a large-scale platform, even if we had both permission and budget to conduct them. Our extensive set of experiments corresponds to several months of platform usage and would have been impossible to deploy on a large-scale platform, even if we had both permission and budget to conduct them. The main conclusion is that two simple and low-complexity greedy strategies significantly outperform FCFS, FAIRSHARE and I/O-Sets.

We conclude this section by stating the main limitations of this work (see Section 3 for a detailed list of application and platform parameters):

- We model the storage system as a black box: we assume that it is a monolithic block and that it can offer a fixed bandwidth to all applications regardless of their I/O patterns and placement on the network. In practice this is a simplifying assumption because: (1) an HPC storage system is made of hundreds (sometimes thousands) of storage nodes in a complex topology; (2) each compute node does not have the same bandwidth, latency, and number of hops to reach each storage node; and (3) applications do not usually access all the storage nodes, and for those nodes that they do access, they do not necessarily access them in the same manner.
- We assume the existence of an I/O controller, i.e. a centralized entity that: (1) can see all the I/O traffic; and (2) can make split-microsecond decisions on how to handle it.

Both assumptions are introduced to make the problem tractable. As for the first assumption, the recent survey [3] states that *the multi-layered software and hardware HPC I/O stack is complex. To access data in HPC systems, applications issue requests that, while traversing the I/O stack,*

are reshaped via a series of data transformations. These originate from distinct abstractions and mappings between the data models used in each layer combined with optimization techniques applied before reaching the file system and, eventually, the storage hardware. Assessing the performance of scheduling algorithms in the framework of the full I/O stack is impossible without extensive experiments conducted on a variety of platforms. As for the second assumption, although an application-level I/O controller is not generally available everywhere, projects like [12] and [7] provide initial implementations on which such a system can be built. In our evaluation, we consider the cost of taking the scheduling decisions, and this cost partly drives our final recommendations. Altogether, the design, analysis and comparison in simulation of all the scheduling algorithms introduced in this work lay the foundations of I/O bandwidth-sharing strategies, and represent a preliminary but mandatory first step before further assessment by the community.

The paper is organized as follows. We first survey related work in Section 2. Then, we detail the application and platform framework in Section 3, together with the optimization objectives. We detail well-known bandwidth-sharing strategies, and introduce new ones, in Section 4. Complexity results are stated in Section 5 in the form of lower bounds for competitive ratios. The experimental evaluation in Section 6 presents extensive simulation results comparing all the strategies. Finally, we conclude and provide hints for future work in Section 7.

2 Related Work

We discuss related work in this section. We survey existing approaches before pointing to a related problem in the scheduling literature.

CALCioM [7] This pioneering paper introduces and experimentally compares three policies to manage cross-application coordination of I/O operations: (i) *Interference* (called FAIRSHARE in this paper), where the total bandwidth is shared equally¹ among all concurrent operations; (ii) *FCFS-based serialization* (called FCFS in this paper), where I/O operations are serialized based upon an FCFS priority; and (iii) *Interruption-based serialization*, where I/O operations are serialized but preemptive, allowing for another operation B to interrupt the current operation A, which resumes only after the completion of B. Some examples are given to explain when to favor a given strategy, but no general approach is explored. In particular, interruption-based serialization would require to set priorities among applications, which are not detailed in the paper. Altogether, this work presents one of the first comparisons of bandwidth-sharing strategies, and we build upon their ideas to cast a general framework and introduce new strategies.

CLARISSE [12] This paper introduces a middleware designed to enhance data-staging coordination and control in the HPC software storage I/O stack. Among many other contributions, the CLARISSE middleware enables to directly compare the *no-scheduling* strategy (called FAIRSHARE in this paper) with FCFS and reports performance gains for the latter. Intuitively, the superiority of FCFS can be expected as it comes from a classic result in parallel computing: when scheduling two identical communications that can each make use of the full bandwidth, better serialize them than execute them concurrently. Indeed, with serialization, the first communication ends at time t and the second one at time $2t$ (for a duration t , assuming a start at time 0), while in parallel, both communications end at time $2t$. However, our analysis and experiments reveal that this intuition can be misleading and that (i) FAIRSHARE prevails over FCFS

¹More precisely, FAIRSHARE shares the total bandwidth in proportion to the size of the concurrent applications, see Section 4.1 for details.

in many practical scenarios and (ii) more sophisticated policies that account for past history to set priority-based bandwidth assignments perform even better.

I/O-Cop [25] I/O-Cop is a prototype system aimed at exploring access control mechanisms to manage the shared Parallel File System (PFS) of the platform. This work is motivated by revealing the contention incurred when several applications aim at performing I/O transfers simultaneously. The I/O-Cop prototype is limited to the case when the access controller to the Parallel File System (PFS) provides exclusive access to a single application at a given time, and without allowing for preemption of ongoing I/O operations.

QoS-based and reward-based approaches In [26], the authors also advocate controlling accesses to the PFS in order to achieve some Quality of Service (QoS) for each application. They envision a system with several I/O storage devices (disks, SSDs or NVRAMs) and aim at load-balancing I/O requests across all storage types to minimize contention. In [24], the authors consider several applications that execute concurrently and post I/O requests. They partition all the I/O requests into several queues, one per application, and aim at establishing priorities across the applications. The idea is that after completing some I/O transfer, a given application could be granted access for its next I/O transfer before all the other applications would have completed one I/O transfer themselves. At each time-step, the progress of each application is monitored as the number of I/O transfers that have been granted so far. In a related paper [11], the authors survey I/O capabilities of state-of-the-art supercomputers and enforce QoS constraints for I/O transfers by implementing a token-based bucket algorithm that works similarly to that of [24]. Finally, the authors of [23] target a system with several I/O sub-systems (OST, which stands for Object Storage Target, typically a RAID array of disks). For each application, they aim at the same share of available bandwidth on each OST, because it balances transfers (one needs to wait for the last node to complete its transfer before resuming work). The allocation of nodes (hence applications) to the different OSTs is given by some external mechanism. Then, on a given OST, some application may benefit from an increased bandwidth, which is done by throttling another application. The throttled application is issued a coupon, to be redeemed later. They do not deal with the interplay of successive I/O operations and work phases, and no comparison is made with other strategies. In contrast, our work restricts to a single OST but provides a comprehensive comparison of several bandwidth-sharing strategies

Periodic applications A series of papers [9, 1, 2, 6, 13] focus on periodic applications that consist of work phases followed by I/O operations. More precisely, each application repeats a two-phase period with a fixed computing length followed by an I/O of volume. The CPU lengths and I/O volume depend upon the application, but remain the same from one period to the next. The major goal of these works is to orchestrate a global periodic scheme where I/O transfers are meticulously shaped to fill up the smallest possible rectangle that will repeat. While the problem of finding the minimum size rectangle is shown to be NP-complete in the initial work [9], several interesting heuristics have been developed in the subsequent papers. The approach is quite flexible, with I/O transfers possibly split into different sub-transfers, each with a different bandwidth. The main limitation is of course the assumed periodicity of each application. An extension is provided by other authors in [28], where applications still consist of phases with work followed by I/O transfers, but now CPU phases have stochastic lengths taken from some probability distribution, while I/O phases have constant length. As a motivation, for CPU phases, we can think of a constant amount of flops to perform, with some system-dependent or data-dependent noise, while for I/O transfers, we can think of a fixed-size checkpoint operation. In contrast, our approach does not assume any a priori knowledge of the concurrent applications.

I/O-Sets [27] This recent work can be viewed as an interesting extension of the work in [2] for periodic applications. Each application consists of several iterations, which as above are work phases followed by I/O operations. Periodicity is no longer assumed. Instead, for each application, they determine the value of ω , which is the average length of an iteration so far. In [27], CPU lengths and I/O volumes are sampled from some probability distributions (that differ for each application), which enables to compute ω with the expectations of these distributions, but one could envision to acquire the value of ω on the fly, as the application progresses. Then, the applications are partitioned into I/O-sets: two applications belong to the same set if they have the same value for $\lceil \log_{10} \omega \rceil$. Each I/O-set is assigned a priority. The I/O bandwidth-sharing strategy is described in detail in Section 4.2. In a nutshell, FCFS is enforced within each I/O-set; hence, at most one application per I/O-set is competing for bandwidth at any time-step. Then some priority-based sharing is enforced across I/O-sets. The motivation for using a mixture of FCFS and FAIRSHARE (or more precisely a priority-based variant of sharing) is very interesting: small and large applications (characterized by different orders of magnitude for ω) should not be treated equally by the scheduler. The I/O-sets strategy has several parameters, and we use the same instantiation as in [2], with the same name SET-10. We use SET-10 as a competitor for our novel strategies.

A note on the painter problem In the scheduling literature, the *painter* problem, a.k.a the *scheduling with delays* problem, is the following: (i) several chains of tasks are to be scheduled on a single machine; (ii) for each chain, there is a minimal *delay* to be enforced between the completion of a task and the start of its successor. As for the analogy with a painter: the painter is the machine and has several rooms to paint on its agenda, each with several paint layers (a task is the application of a paint layer); for each room (each chain), there is a delay between the end of a layer and the next one. The tasks are not preemptive. Release times can be simulated by adding delays from a fake source task. This is an offline problem where the objective is to minimize either the makespan (maximum completion time of a task) or the total flow (unweighted or weighted sum of all completion times). The analogy with the I/O problem is clear: the machine is the I/O resource, the task chains are the applications, the tasks are the I/O operations, and the delays are the computing phases between two consecutive I/O operations. The main differences with the I/O scheduling problem are the following:

1. Execution is not preemptive in the painter problem, while one can pause an on-going I/O operation;
2. A single task is executed at any time step while several I/O operations (from different applications) can share the I/O bandwidth;
3. All chain parameters (task lengths and delay values) are known at the beginning of the execution while the lengths of work phases and the volumes of I/O operations are discovered on the fly in the I/O scheduling problem.

Particular instances of the painter problem have been shown to have polynomial complexity. We refer to the interested reader to [20, 21, 5, 18, 8] for details. A survey of recent results and extensions is available in [16].

3 Framework

In this section, we describe the framework. We start with application characteristics and detail rules for I/O operations and bandwidth allocation in Section 3.1. We discuss the interaction with the batch scheduler and explain why we restrict to steady-state time windows in Section 3.2. We conclude with optimization objectives in Section 3.3. Main notations are summarized in Table 1.

3.1 Applications

3.1.1 Application Characteristics

We consider a very general framework where applications are submitted online to the batch scheduler. Each application \mathcal{A}_i requests p_i nodes and starts executing as soon as the batch scheduler has been able to allocate that many nodes. Thus, each application executes on a dedicated set of nodes throughout its execution, which is the standard approach on large-scale HPC platforms. However, all applications execute I/O transfers (reads and writes) through the I/O controller and share the bandwidth of the I/O system. Our approach is agnostic of the nature of the storage (SSDs, NVRAMs, disks or tapes), and of the organization of the PFS (Parallel File System).

Each application \mathcal{A}_i executes an alternating sequence of work phases and I/O operations, which we represent as follows:

$$\mathcal{A}_i \equiv v_i^{(0)}, w_i^{(1)}, v_i^{(1)}, w_i^{(2)}, \dots, v_i^{(n_i-1)}, w_i^{(n_i)}, v_i^{(n_i)}, \dots$$

where $v_i^{(j)}$ stands for I/O volumes, and $w_i^{(j)} > 0$ stands for (parallel) work units. Because computing nodes are dedicated to the application, we can assume w.l.o.g. that one unit of work lasts one second, so that the $w_i^{(j)}$ represent the duration of the work phases; more precisely, within w_i seconds, each of the p_i nodes performs w_i work units. However, because we do not know the bandwidth of I/O operations in advance, we have to express them in volume (amount of bytes) rather than in duration. We detail rules for bandwidth allocation in Section 3.1.2. We will discuss rules for posting and managing I/O operations in Section 3.2.2.

As stated before, the length $w_i^{(j)}$ of each work phase is not known until it terminates, and the volume $v_i^{(j)}$ of each I/O operation is not known until the operation is posted to the I/O controller. Similarly to the related work surveyed in Section 2, we also assume that I/O operations are blocking and coordinated between the different nodes of the application, and that the application does not overlap I/O operations with some work phase. This is typical of HPC applications using a synchronous global interface like MPI-IO [19, 14], which also provides the I/O controller with critical information like the volume of data to transfer.

m	number of applications
p_i	size (number of nodes) of application \mathcal{A}_i
τ_i	release time of application \mathcal{A}_i
$w_i^{(j)}$	duration of work phase number j for application \mathcal{A}_i
$v_i^{(j)}$	volume of I/O operation number j for application \mathcal{A}_i
B	total bandwidth of the I/O system
b	bandwidth of each platform node
b_i	maximal bandwidth of application \mathcal{A}_i : $b_i = \min(p_i b, B)$
α_i^j	fraction of bandwidth assigned to \mathcal{A}_i for I/O operation j (can vary over time)
$[T_{begin}, T_{end}]$	steady-state window
$y_i(t)$	yield of application \mathcal{A}_i at time t

Table 1: Summary of main notations.

3.1.2 Bandwidth Allocation

Consider an application \mathcal{A}_i executing on p_i nodes and initiating an I/O operation of volume $v_i^{(j)}$. What are the bandwidth allocation rules for this operation? We let b be the bandwidth of the network card (of interface card) of each node, and B be the total bandwidth of the I/O system.

First, assume for simplicity that the I/O operation is not interrupted, and is granted the same bandwidth from start to completion. The maximal bandwidth that can be granted by the I/O controller is

$$b_i = \min(p_i b, B). \quad (1)$$

Note that Equation (1) implicitly assumes that each node of \mathcal{A}_i has to transfer (approximately) the same volume of data to/from the PFS. If transfers are unbalanced from one node to another, we should redefine $v_i^{(j)}$ as $v_i^{(j)} = p_i v_{i,\max}^{(j)}$, where $v_{i,\max}^{(j)}$ is the maximum volume of data to be transferred by any of the p_i nodes of \mathcal{A}_i . The main rule of the game for the scheduler is to assign a fraction $\alpha_i^{(j)}$ of the maximal bandwidth b_i to the I/O operation $v_i^{(j)}$. The duration of the I/O operation will then be

$$d_i^{(j)} = \frac{v_i^{(j)}}{\alpha_i^{(j)} b_i}. \quad (2)$$

Of course, if no I/O operation has been posted by another application, the scheduler will enforce $\alpha_i^{(j)} = 1$ to ensure fastest possible completion. In that case, we use the notation

$$d_{i,\min}^{(j)} = \frac{v_i^{(j)}}{b_i} \quad (3)$$

to denote the minimal possible duration of the I/O operation. On the contrary, in the presence of several concurrent I/O operations, the scheduler will resort to some bandwidth-sharing strategies, like the ones studied in this paper.

We are ready to discuss the general case, which will require some additional notations. Intuitively, a given I/O operation will NOT be granted the same bandwidth fraction throughout execution. At any time-step t , some I/O operations that were posted before are granted some bandwidth and executing, while some others may be pending (that is to say their fraction is currently 0). A new I/O operation may be posted at time t , which the scheduler can account for by granting it some bandwidth, at the price of reducing the fraction of other applications. On the contrary, some on-going I/O operation may complete at time t , thereby opening the possibility of a larger fraction to be granted to some applications. We see that bandwidth fractions are granted only for some duration, which we call the *horizon*. Decisions are taken at specific instants, which we call *events*. Typically, an event corresponds to the posting of a new I/O operation, or to the termination of an on-going one. But an event can also be triggered by the I/O scheduler, e.g., for a strategy where additional events are created periodically, say every 10 seconds. The I/O controller takes a new decision at every event, as explained below. The constraints on the number of events, and the cost of bandwidth-sharing strategies, will be detailed in Section 3.2.2.

Consider an event at time t , and let $S(t)$ be the index set of *active* applications, i.e., applications that have posted an I/O operation before time t which is not yet completed, or applications that post a new I/O operation exactly at time t . Among the applications with incomplete I/O-operations, some may be transferring data at some bandwidth fraction and some may be kept waiting. Each active application \mathcal{A}_i , $i \in S(t)$, is allotted a bandwidth $\alpha_i^t b_i$ (with some α_i^t possibly 0) so that

$$\sum_{i \in S(t)} \alpha_i^t b_i \leq B. \quad (4)$$

This bandwidth allocation remains valid until the next event at time $t+h$, where h is the horizon. The bandwidth allocation depends upon the bandwidth-sharing strategy, whose inputs are the volume of data that must still be transferred for each on-going I/O operation, the knowledge of the progress of all active applications so far, and the optimization objective.

We stress that the horizon h is unknown at time t . The next event is triggered either by a new post or a completion, or again by an external decision given to the I/O controller. At time t , after having granted bandwidth fractions to active applications, we only know that h is greater than the time needed to complete the shortest on-going I/O operation, given that no new event (new post or external) will happen before that.

When the next event takes place at time $t+h$, we update the set of active applications, leaving out I/O operations that have completed and including new posts, if any. We also update the remaining volume of data still to be transferred for each active application. The I/O controller applies the bandwidth-sharing strategy for this new set of parameters.

3.2 Steady-State Windows

In this section, we recall the management of HPC applications by the batch scheduler, and explain why we need to restrict to steady-state time windows to assess the performance of bandwidth-sharing strategies.

3.2.1 Interaction with the Batch Scheduler

HPC applications are submitted to the batch scheduler. Each application \mathcal{A}_i has a release time τ_i , a size p_i and a wall-time res_i (length of the reservation slot). Upon release, application \mathcal{A}_i is put in the queue of the batch scheduler and will be allocated resources at time $t_i^{alloc} \geq \tau_i$, which means that p_i nodes are dedicated to the application during the interval $[t_i^{alloc}, t_i^{alloc} + res_i)$. The p_i nodes are released as soon as the application completes its execution or its deadline is reached, whichever comes first.

Each application has dedicated nodes but all applications that execute concurrently share the I/O system. I/O operations are posted by the applications and managed by the I/O controller. If an application posts an I/O operation while another I/O operation has already been granted access, several scenarios can happen, depending upon the bandwidth-sharing policy implemented by the I/O controller. We have already discussed the FCFS and FAIRSHARE strategies in Section 2, and will introduce other strategies in Section 4. Whenever the I/O controller makes a decision according to its bandwidth-sharing policy, this decision has an impact on the progress of all active applications. Altogether, the bandwidth-sharing policy will change the termination time of all applications. In theory, some applications may even fail to complete before the end of their reservation due to the bandwidth-sharing strategy being disadvantageous to them. On the contrary, some applications may benefit from the strategy and complete early, thereby releasing their resources early. In summary, the opportunities for decisions of the batch scheduler to allocate new applications will depend upon the bandwidth-sharing strategy applied to the applications that are currently executing. Furthermore, any decision of the batch scheduler changes the mix of applications that run concurrently and possibly compete for I/O resources. This, in turn, changes the scope and impact of the decisions of the bandwidth-sharing policy implemented by the I/O controller. Altogether, the interplay between the batch scheduler and the decisions of the I/O controller is hard to comprehend.

To the best of our knowledge, none of the papers surveyed in Section 2 has dealt with this difficulty. Instead, these papers consider a fixed number of applications that execute concurrently (each on a dedicated set of nodes) and compete for I/O access. This amounts to consider an

execution window $[T_{begin}, T_{end}]$ where all applications start executing at time T_{begin} and do not complete execution before time T_{end} , regardless of the I/O policy that is implemented. In other words, the platform operates in steady-state mode during the window $[T_{begin}, T_{end}]$ with no application terminating nor no new application launched throughout the window. This assumption is never stated in recent papers. Again, the reason why it is assumed that applications do not complete before the end of the window is the following: if an application terminates at time $T < T_{end}$, the batch scheduler might launch another application right after the completion. Because T depends on the bandwidth-sharing strategy that is enforced, it becomes impossible to assess the performance of the strategy by itself.

In this paper, we use a steady-state execution window $[T_{begin}, T_{end}]$ and assume that m applications \mathcal{A}_i ($1 \leq i \leq m$) execute concurrently throughout the window. To eliminate side effects and deal with a general scenario, we do not assume that the applications start executing at time T_{begin} : on the contrary, the applications may have been launched earlier and have been executing for some time. The history of the applications will be taken into account when evaluating the objective function (see Section 3.3).

3.2.2 Cost Model for Steady-State Windows

Given a steady-state execution window $[T_{begin}, T_{end}]$, assume that m applications \mathcal{A}_i ($1 \leq i \leq m$) execute concurrently throughout the window. Each application \mathcal{A}_i will execute a series of work phases followed by I/O transfers. If the application \mathcal{A}_i was alone on the platform, all I/O transfers would be granted maximal bandwidth b_i . Let $N_{op}(i)$ be the number of I/O operations that would be initiated from time T_{begin} until time T_{end} , assuming such a dedicated mode.

In concurrent mode, we introduce two *events* for each I/O operation, one when it is posted, and one when it completes. The total number of events due to I/O operations is upper bounded by

$$E = \sum_{i=1}^m 2N_{op}(i). \quad (5)$$

Indeed, no application will perform more I/O operations by the end of the window than in dedicated mode, hence the number of events for each application \mathcal{A}_i never exceeds $2N_{op}(i)$, regardless of the bandwidth-sharing strategy.

The value of E is a key parameter to the size of the problem (other parameters include the binary encoding of work lengths and I/O volumes). We enforce that all bandwidth strategies have a cost polynomial in E , meaning that the number of bandwidth-sharing decisions remains polynomial in E . For instance, if the I/O controller enforces periodic decisions every h seconds, where h is a fixed horizon, the number of additional events $E^{(+)} = \lfloor \frac{T_{end} - T_{begin}}{h} \rfloor$ must remain polynomial in E . We use $E^{(+)} = E$ in the simulations to add equi-spaced decisions across the steady-state window. Note that triggering an external event every second would lead to $T_{end} - T_{begin}$ external events, which is exponential in the problem size (we use a logarithmic encoding for all parameters).

To the best of our knowledge, none of the papers surveyed in Section 2 has discussed how frequently decisions should be taken, nor has included the cost of the bandwidth-sharing strategy each time a decision is taken. We could easily include that cost into the assessment of the performance of the strategies. We do not, because the cost is inherent to the strategy and independent of the actual length of the work phase and I/O operations: if we multiply the latter quantities (and the window size) by a factor 10 or 100, the cost of the strategy remains the same and becomes negligible in front of the execution time of the applications.

3.3 Objectives

In this section, we define the *yield* of an application. The major objective of our novel bandwidth-sharing strategies is MINYIELD, the maximization of the minimum yield over all applications executing within the steady-state window $[T_{begin}, T_{end}]$. However, we also report performance for two other objectives, UTILIZATION and EFFICIENCY, which we describe at the end of this section.

Consider an application \mathcal{A}_i that is released at time $\tau_i = 0$. Consider a steady-state window $[T_{begin}, T_{end}]$. At any time $t \geq T_{begin}$, we want to monitor the progress of \mathcal{A}_i in terms of work done and data volume transferred. Recall that \mathcal{A}_i executes an alternating sequence of work phases (work) and I/O operations:

$$\mathcal{A}_i \equiv v_i^{(0)}, w_i^{(1)}, v_i^{(1)}, w_i^{(2)}, \dots, v_i^{(n_i-1)}, w_i^{(n_i)}, v_i^{(n_i)}, \dots$$

We have assumed unit speed for work phases, and we normalize I/O volumes by the maximal possible bandwidth $b_i = \min(p_i b, B)$. Letting $d_{i,\min}^{(j)} = \frac{v_i^{(j)}}{b_i}$ be the minimum duration for I/O operation number j of volume $v_i^{(j)}$, we rewrite \mathcal{A}_i as

$$\mathcal{A}_i \equiv d_{i,\min}^{(0)}, w_i^{(1)}, d_{i,\min}^{(1)}, w_i^{(2)}, \dots, d_{i,\min}^{(n_i-1)}, w_i^{(n_i)}, d_{i,\min}^{(n_i)}, \dots$$

The *ideal progress* of \mathcal{A}_i at time t is the amount of work plus the volume of data transferred since its release time τ_i and up to time t , when all I/O operations have taken place with no delay and at the maximal possible bandwidth b_i . This corresponds to \mathcal{A}_i progressing at maximal rate, which happens if it executes in dedicated mode on the platform. By definition, at time t , the ideal progress is equal to $t - \tau_i$.

In a concurrent execution, the *actual progress* of \mathcal{A}_i at time t is the amount of work plus the volume of data transferred since its release time τ_i and up to time t . While work phases still progress at full (unit) speed, I/O operations are slowed down by interferences. For any time $t \in [T_{begin}, T_{end}]$, let $W_i^{(done)}(t)$ be the total amount of work done up to time t , and $V_i^{(transferred)}(t)$ be the total volume of data transferred up to time t . The *yield* of \mathcal{A}_i at time t is defined as the ratio of the actual progress over the ideal progress, namely

$$y_i(t) = \frac{W_i^{(done)}(t) + \frac{V_i^{(transferred)}(t)}{b_i}}{t - \tau_i}. \quad (6)$$

As a side note, we show how to compute the value of $V_i^{(transferred)}(t)$ as the concurrent execution goes. We do this computation incrementally, one work phase or I/O operation after another. Consider the I/O operation number j and assume that it has occurred during the interval $[start_i^{(j)}, end_i^{(j)}]$ ($end_i^{(j)}$ is equal to the completion time of this I/O operation and $start_i^{(j)}$ to the completion time of the previous work phase). Let $\alpha_i^{(j)}(u)b_i$ be the bandwidth granted at time $u \in [start_i^{(j)}, end_i^{(j)}]$, where $0 \leq \alpha_i^{(j)}(u) \leq 1$ (and let $\alpha_i^{(j)}(u) = 0$ for u outside this interval). If the I/O operation number j is not complete at time t , i.e., if $t \in [start_i^{(j)}, end_i^{(j)})$, the amount of data volume $V_i^{(j)}(t)$ transferred up to time t is

$$\int_{start_i^{(j)}}^t \alpha_i^{(j)}(u)b_i du = V_i^{(j)}(t). \quad (7)$$

In fact, the integral is a discrete sum of at most E components, since we change bandwidth allocation only when a new event takes place. Note that if $t \geq end_i^{(j)}$, we obtain $V_i^{(j)}(t) = v_i^{(j)}$.

Equation (7) enables us to compute the actual progress incrementally, from one work phase or I/O operation to the next. Of course, the actual progress depends upon the bandwidth-sharing strategy through the choice of the fractions $\alpha_i^{(j)}(u)$ of the maximal bandwidth b_i allotted at every instant u .

We are ready to state the optimization objectives, together with their initial motivation. Consider a steady-state window $[T_{begin}, T_{end}]$ and m applications. Each application \mathcal{A}_i has a yield $y_i(T_{begin})$ when entering the window. The three target objectives are **MINYIELD**, **UTILIZATION** and **EFFICIENCY**.

MINYIELD The objective is to maximize the minimum yield at the end of the window:

$$\text{MAXIMIZE } \min_{1 \leq i \leq m} y_i(T_{end}). \quad (8)$$

This objective aims at enforcing fairness among all the applications, regardless of their characteristics. The intuition is that all applications suffer from the same slowdown factor if they achieve the same yield. As discussed in Sections 1 and 2, previous work has shown the limitations of FCFS and FAIRSHARE, which give priority to some applications and severely slow down other ones. **MINYIELD** will guide bandwidth-sharing decisions so that all applications exit the window with balanced yields. An application entering the window with a very low yield will be granted more bandwidth to catch up.

UTILIZATION The objective is to maximize platform utilization throughout the window:

$$\text{MAXIMIZE } \frac{\sum_{1 \leq i \leq m} p_i \left(W_i^{(done)}(T_{end}) - W_i^{(done)}(T_{begin}) \right)}{(T_{end} - T_{begin}) \sum_{1 \leq i \leq m} p_i}. \quad (9)$$

The work $W_i^{(done)}(T_{end}) - W_i^{(done)}(T_{begin})$ done by each application \mathcal{A}_i within the window is weighted by its size p_i . This objective is the classical performance objective from the perspective of the administrator or owner of the platform, because it measures the fraction of time where computing nodes have been used for actual application work. Hence, this objective is natural for HPC applications that perform no or little I/O transfers. However, it may seem ill-suited in a framework focusing on I/O transfers, because it is very sensitive to the ratio of work over data volumes (normalized by maximal bandwidth). For instance, if we multiply all data volumes by, say, 10, platform utilization will plummet, even if we keep the same bandwidth-sharing strategy. This observation leads to introducing the objective **EFFICIENCY**.

EFFICIENCY The objective is to maximize the sum of the actual progress of all applications throughout the window:

$$\text{MAXIMIZE } \frac{\sum_{1 \leq i \leq m} p_i \left(W_i^{(done)}(T_{end}) - W_i^{(done)}(T_{begin}) + \frac{V_i^{(transferred)}(T_{end}) - V_i^{(transferred)}(T_{begin})}{b_i} \right)}{(T_{end} - T_{begin}) \sum_{1 \leq i \leq m} p_i}. \quad (10)$$

Comparing Equations (9) and (10), we see that I/O operations are taken into account with **EFFICIENCY**: this objective aims at optimizing the combined progress of all applications. It can be viewed as a measure of how efficiently platform resources (both compute nodes and the I/O system) are used.

4 Bandwidth-Sharing Strategies

We describe bandwidth-sharing strategies in this section. We start by recalling a few notations and introducing new ones. Consider a steady-state window $[T_{begin}, T_{end}]$ with m applications executing concurrently. Consider an event at time t and let $S(t)$ be the index set of active applications at time t . Note that applications that are not active are engaged in work phases at time t and progress independently of the decisions made by the I/O controller.

Each active application \mathcal{A}_i , $i \in S(t)$, has posted an I/O operation at time $R_i \leq t$ that is not complete at time t . Let \mathcal{V}_i denote the remaining volume still to be transferred for the I/O operation. Each active application is allotted a fraction α_i^t (with some α_i^t possibly 0) of its maximum possible bandwidth $b_i = \min(p_i b, B)$. The bandwidth-sharing strategy consists in determining α_i^t for each active application \mathcal{A}_i . Finally, let $\mathcal{BW}_i(t', y)$ denote the bandwidth that should be allotted to application \mathcal{A}_i for it to achieve a yield of at least y at time t' .

We start with some simple greedy strategies, some old and some new, in Section 4.1. Then in Section 4.2, we detail the recent SET-10 strategy proposed in [4]. Finally, in Section 4.3, we sketch an elaborate strategy whose aim is to compute the best horizon for maximizing the minimum yield.

4.1 Greedy Strategies

We discuss below six greedy strategies. The first three strategies do not rely on any (tentative) horizon, while the last two aim at taking some future events into account. Finally, the sixth strategy re-evaluates the current bandwidth allocation at periodic time-steps.

- **FAIRSHARE**: each active application \mathcal{A}_i with $i \in S(t)$ is allocated $\alpha_i = \min(1, \frac{B}{\sum_{j \in S(t)} b_j})$. Therefore, each application will either saturate its maximal bandwidth b_i , or it will receive a fair share (proportional to its size p_i) of the total bandwidth B . This is the de-facto strategy implemented by the parallel filesystems available in most HPC centers. This strategy does not need to consider what application is requesting the I/O operation, but just how many I/O operations are currently concurrent.
- **FCFS**: greedily allocate the bandwidth to active applications sorted by non-decreasing R_i . More precisely, up to some re-ordering, let $S(t) = \{1, 2, \dots, k\}$ with $R_i \leq R_{i+1}$ for $1 \leq i < k$. \mathcal{A}_1 is granted its maximum bandwidth b_1 (hence, $\alpha_1 = 1$), then \mathcal{A}_2 is granted $\alpha_2 b_2 = \min(b_2, B - \alpha_1 b_1)$, and so on until no more bandwidth is available.
- **GREEDY YIELD**: greedily allocate the bandwidth to active applications sorted by non-decreasing yields $y_i(t)$. The greedy allocation process is the same as for FCFS but with a different criterion, current minimum yield instead of oldest posting time. This strategy gives priority to applications with low yield, so that they can catch up.
- **GREEDY COM**: greedily allocate the bandwidth to the applications sorted by non-decreasing ratio \mathcal{V}_i/b_i , i.e., by the remaining time to complete the pending I/O operation at maximum possible bandwidth. This strategy gives priority to completing shorter transfers, with the goal of freeing the I/O system as fast as possible and/or give more bandwidth to forthcoming I/O operations.
- **LOOKAHEAD GREEDY YIELD**: for each active application \mathcal{A}_i , compute the minimum yield Z_i that can be achieved (over all active applications) if \mathcal{A}_i is given priority and allocated the maximum possible bandwidth b_i , and where the remaining bandwidth $B - b_i$ is allocated following **GREEDY YIELD** for the other applications in $S(t)$. Then, we retain the allocation

that maximizes the minimum yield Z_i obtained with these $|S(t)|$ possible priority choices. The rationale for LOOKAHEADGREEDYIELD is to look ahead and maximize the minimum yield not at time t , but at time $t+h$, where the horizon h is (tentatively) computed as the end of one ongoing I/O operation.

- PERIODICGREEDYIELD (δ): this strategy is a variant of GREEDYIELD where I/O decisions are triggered by external (periodic) events submitted to the I/O controller every δ seconds, in addition to the regular events that correspond to posting and completion of I/O operations. As discussed in Section 3.2.2, we must restrict to a polynomial number of external events. With the notations of Section 3.2.2, we use $E^{(+)} = E$ in the simulations, which leads to choosing $\delta = \frac{T_{end} - T_{begin}}{E^{(+)}}$. At every event, external or regular, bandwidth-sharing decisions are the same as for GREEDYIELD. The rationale for adding periodic events is to avoid the risk that GREEDYIELD would apply a bad decision for too long: with several concurrent I/O operations lasting for a long time, greedy decisions are updated every δ seconds, instead of waiting for the first completion of one of these I/O operations.

4.2 SET-10 Strategy

This section provides a description of SET-10, the I/O-sets bandwidth-sharing strategy from [4].

Determination of I/O-sets With the notations of Section 3.3, consider an application \mathcal{A}_i composed of operations

$$v_i^{(0)}, w_i^{(1)}, v_i^{(1)}, w_i^{(2)}, \dots, v_i^{(n_i-1)}, w_i^{(n_i)}, v_i^{(n_i)}, \dots$$

Assume that \mathcal{A}_i has just completed the I/O operation $v_i^{(j)}$. Then, the current value of ω^i , the average length of an iteration for \mathcal{A}_i , is defined as

$$\omega^i = \frac{1}{j} \sum_{k=1}^j (w_i^{(k)} + d_{i,\min}^{(k)}),$$

where $d_{i,\min}^{(j)} = \frac{v_i^{(j)}}{b_i}$, and $b_i = \min(p_i b, B)$. Note that we neglect the initial I/O operation $v_i^{(0)}$ to match the specification of [4]. Then, \mathcal{A}_i is assigned to I/O-set \mathcal{S}_n , where $n = \lceil \log_{10} \omega^i \rceil$, and $\lceil x \rceil$ denotes the nearest integer to x . Note that an application \mathcal{A}_i may be dynamically reassigned to another I/O-set depending upon the duration of its next work phases and I/O operations. In [4], I/O-set \mathcal{S}_n , where $n = \lceil \log_{10} \omega^i \rceil$, receives a priority $q_n = 10^{-n}$.

Bandwidth assignment Consider an event occurring at time t , and let $S(t)$ denote the index set of active applications that have a pending I/O transfer at time t . Each participating application \mathcal{A}_i , $i \in S(t)$, is allotted a bandwidth $\alpha_i b_i$ computed via the following algorithm [4]:

1. Assume that the applications in $S(t)$ belong to s different I/O-sets $\mathcal{S}_{n_1}, \mathcal{S}_{n_2}, \dots, \mathcal{S}_{n_s}$.
2. Within each I/O-set, a single application is granted access to the I/O system. In other words, there is exclusive access within sets. If several applications in $S(t)$ belong to the same I/O set, the one with the smallest value of R_i (FCFS, the one that posted its request first) is selected.
3. Now, we have a subset of s applications, one per I/O subset, which will be granted some bandwidth. The intuition is to partition the bandwidth according to the priorities defined

above. For simplicity, let us renumber the applications so that \mathcal{A}_j is the application chosen from set \mathcal{S}_{n_j} , for $1 \leq j \leq s$. Then, each application \mathcal{A}_j should be granted the fraction $\alpha_j = \frac{q_{n_j}}{\sum_{1 \leq k \leq s} q_{n_k}}$ of the total bandwidth B .

4. As usual, this bandwidth assignment remains valid until the next event.

However, this bandwidth-sharing algorithm implicitly assumes that each application can use the whole system bandwidth: $b_i = B$ for each application \mathcal{A}_i . To cope with general scenarios where this is not the case, we have to extend the algorithm. The natural idea is allocate bandwidth to several applications in the same I/O subset, rather than one, while still enforcing the priorities. More precisely, the fraction $\frac{q_{n_j}}{\sum_{1 \leq k \leq s} q_{n_k}}$ of the total bandwidth B is now assigned to several applications from \mathcal{S}_{n_j} , chosen greedily in FCFS order. Here is the extended algorithm for bandwidth-sharing:

1. Assume that the applications in $S(t)$ belong to s different I/O-sets $\mathcal{S}_{n_1}, \mathcal{S}_{n_2}, \dots, \mathcal{S}_{n_s}$.
2. For each I/O-set \mathcal{S}_{n_j} , compute the maximum bandwidth fraction that it can receive, namely $\beta_j = \frac{\sum_{k \in \mathcal{S}_{n_j}} b_k}{B}$. As before, let $\alpha_j = \frac{q_{n_j}}{\sum_{1 \leq k \leq s} q_{n_k}}$.
3. We partition the s I/O sets into two categories, those that can receive the fraction α_j and those that are limited by their maximal bandwidth fraction β_j . Let \mathcal{C} be the set of I/O sets of the latter category, i.e., such that $\beta_j \leq \alpha_j$.
4. All the applications \mathcal{A}_k in an I/O set belonging to \mathcal{C} receive their maximal bandwidth b_k .
5. We compute the remaining bandwidth $B_{left} = (1 - \sum_{\mathcal{S}_{n_j} \in \mathcal{C}} \beta_j)B$.
6. We repeat the whole procedure with the remaining I/O-sets and B_{left} , until either there is no I/O-set left, or all remaining I/O-sets have a larger maximal bandwidth than their priority share: $\beta_j \geq \alpha_j$. In the final step, the remaining I/O-sets are granted the fraction α_j of the remaining bandwidth B_{left} . Within each of these I/O sets, bandwidth is allotted greedily in FCFS order.

Remark on Framework The I/O-sets strategy [4] does not assume that the total volume of an I/O operation is known when that operation is posted. Instead, they assume that this volume is unknown until the I/O operation ends. They rely on the knowledge of the average length of an iteration for each application, which is acquired from past behavior traces. In our simulations of SET-10, we acquire information on average iteration length on the fly as execution progresses.

As stated in Section 3.1.1, we do assume that the total volume of each I/O operation is known when posted. This knowledge is necessary for GREEDYCOM, LOOKAHEADGREEDYIELD (described in Section 4.1) and BESTNEXTEVENT (described below in Section 4.3). However, GREEDYIELD and LOOKAHEADGREEDYIELD (also described in Section 4.1) do not need any information at all on the applications, they only need to compute application yields on the fly. And of course FAIRSHARE and FCFS do not need any information either.

4.3 Maximizing the Minimum Yield at the Next Event

Given an event at time $t \in [T_{begin}, T_{end}]$, the aim of strategy BESTNEXTEVENT is to find the *best predictable* event in the remainder of the window $]t, T_{end}]$. A predictable event is either the end of the execution window (at time T_{end}) or the first time one of the currently on-going I/O

operations is completed, whichever comes first. The best predictable event is the predictable event at which point the minimum yield will be maximized. Of course, if an unpredictable event, such as the posting of a new I/O operation, surges before the best predictable event, the bandwidth-sharing strategy will account for it and recompute the best predictable event from that time on.

A priori, there are infinitely many dates in the interval $[t, T_{end}]$, at which the next predictable event can happen; hence, we cannot test each and every one of them. Instead, we partition the interval $[t, T_{end}]$ into a polynomial (in practice, quadratic) number of sub-intervals. The extremities of these sub-intervals will be either the earliest date at which an I/O operation can complete, or the time at which the characteristic yield functions of two applications intersect (see below for details; the characteristic yield function of an application will be, for instance, its maximum achievable yield at time t' , or its yield at time t' if it is allocated no bandwidth, etc.).

Let $t = t_1 \leq t_2 \leq \dots \leq t_{n_{int}} = T_{end}$ be the extremities of these sub-intervals. For each sub-interval $[t_i, t_{i+1}]$, we will consider each application \mathcal{A}_k that can define an event in (t_i, t_{i+1}) (hence, each application \mathcal{A}_k such that $t_i \geq \frac{\mathcal{V}_k}{b_k}$). Then, we search for the event defined by \mathcal{A}_k that maximizes the minimum yield in $[t_i, t_{i+1}]$. For that purpose, we start by looking for the best solution at time t_i . Once we have identified that solution, we determine the largest interval $[t_i, t'_i] \subset [t_i, t_{i+1}]$ such that for any $t' \in [t_i, t'_i]$ the optimal solution at time t' has the same structure (which applications are allocated bandwidth, which application is allocated its maximal bandwidth, etc.) as the one at time t_i . If $t_i = t_{i+1}$ we conclude. Otherwise, we call recursively the algorithm on the interval $[t'_i, t_{i+1}]$.

Because application \mathcal{A}_k is defining an event at time t_i , it receives the bandwidth $\frac{\mathcal{V}_k}{t_i - t}$, where \mathcal{V}_k is the remaining volume at time t (hence, the I/O operation completes at time t_i). The remaining bandwidth $B - \frac{\mathcal{V}_k}{t_i - t}$ must be distributed among the other applications. We first compute an upper-bound, y^{UB} , on the maximum minimum yield: y^{UB} is the minimum, over all applications, of the maximum yield achievable by each application at time t_i . We then check whether this upper-bound can be achieved without exceeding the total bandwidth B . In the following, let $y_j(t', b)$ denote the yield of application j at time t' if it is allocated a bandwidth of b during the interval $[t, t']$.

Assume first that y^{UB} is not achievable. Then, let $y^{opt}(t_i)$ denote the maximum minimum yield at t_i . The goal is to find the value of $y^{opt}(t_i)$ and compute the set \mathcal{I} of applications to which some bandwidth must be allocated. This is exactly the set of applications whose yield ($y_j(t_i, 0)$) is strictly lower than $y^{opt}(t_i)$ if they were not allocated any bandwidth. \mathcal{I} can be computed by checking the total bandwidth required for all applications to achieve a yield of at least $y_j(t_i, 0)$, for any application \mathcal{A}_j in $S(t) \setminus \{k\}$. We can show that all applications receiving bandwidth must achieve the same yield. Then, knowing \mathcal{I} , we can compute the value of $y^{opt}(t_i)$. Hence, we know how to maximize the optimal minimum yield at time t_i (under our hypothesis, namely that \mathcal{A}_k defines an event at that time). We have built a solution: \mathcal{A}_k defines an event and the remaining bandwidth is distributed among the applications in \mathcal{I} which all achieve the same yield. We can write the yield achieved by this solution has a function of $t' - t_i$ for $t' \in [t_i, t_{i+1}]$. This function is of the form: $a \frac{b+(t'-t_i)}{c+t'-t_i}$. Hence, it is monotonic. If it is non-increasing, we conclude that the optimal is found at time t_i . If it is increasing, we compute the latest time $t'_i \in [t_i, t_{i+1}]$ for which our solution defines the optimal solution. t'_i is the last time at which all the conditions defining the solution hold, namely, for all $t' \in [t_i, t'_i]$:

- only applications in \mathcal{I} receive bandwidth: $\forall j \in S(t) \setminus \mathcal{I}, j \neq k \Rightarrow y_j(t', 0) \geq y^{opt}(t')$;
- the bandwidth limits of all applications are satisfied: $\forall j \in \mathcal{I}, \mathcal{B}\mathcal{W}_j(t', y^{opt}(t')) \leq b_j$;
- $y^{opt}(t')$ is not greater than the yield of \mathcal{A}_k : $y^{opt}(t') \leq y_k(t', \frac{\mathcal{V}_k}{t'})$.

t'_i is computed by solving a set of second degree polynomials. Then, if $t'_i = t_{i+1}$, the optimum is achieved at time t_{i+1} . Otherwise, the algorithm is called recursively on the interval $[t'_i, t_{i+1}]$.

Now, assume that y^{UB} is achievable. If the upper-bound is achieved by application \mathcal{A}_k , because its yield is decreasing on $[t_i, t_{i+1}]$, the optimum is achieved at time t_i . Otherwise, let \mathcal{A}_j be an application whose maximum achievable yield is minimal throughout $[t_i, t_{i+1}]$ (which implies $y_j^{\max}(t_i) = y^{UB}$). Therefore, \mathcal{A}_j achieves its maximum achievable yield at time t_i . Two cases can happen.

- \mathcal{A}_j I/O operation ends at time t_i . Then, the optimum is achieved at time t_i because the maximum achievable yield of \mathcal{A}_j is then decreasing on $[t_i, t_{i+1}]$.
- \mathcal{A}_j is allocated its maximum bandwidth b_j . Then, the yield of \mathcal{A}_j is increasing over $[t_i, t_{i+1}]$ (as long as we can allocate it a bandwidth of b_j). Once again, we compute the set \mathcal{I} of applications to which bandwidth must be allocated. Then, we compute the latest time $t'_i \in [t_i, t_{i+1}]$ for which our solution defines the optimal solution. t'_i is the last time at which all the conditions defining the solution hold, namely:
 - only applications in \mathcal{I} receive bandwidth: $\forall l \in S(t) \setminus \mathcal{I}, l \neq k \Rightarrow y_l(t', 0) \geq y_j(t', b_j)$;
 - the total bandwidth is not exceeded: $\sum_{l \in \mathcal{I}} \mathcal{BW}_l(t', y_j(t', b_j)) \leq B - \frac{y_k}{t'}$;
 - the yield of \mathcal{A}_j is not greater than the yield of \mathcal{A}_k : $y_j(t', b_j) \leq y_k(t', \frac{y_k}{t'})$.

t'_i is computed by solving a set of second degree polynomials. Then, if $t'_i = t_{i+1}$, the optimum is achieved at time t_{i+1} . Otherwise, the algorithm is called recursively on the interval $[t'_i, t_{i+1}]$.

This concludes the high-level description of BESTNEXTEVENT. All details and algorithms are available in Section A of the Appendix. Altogether, BESTNEXTEVENT is quite complicated, and admittedly too complicated for practical use. But it will serve as a reference to help us assess the quality of the (simpler) greedy strategies of Section 4.1.

	MINYIELD	EFFICIENCY	UTILIZATION
FAIRSHARE [7, 12]	$m^{(1)}$	$\frac{m}{4}^{(2)}$	$\infty^{(2)}$
FCFS [7, 12]	$\infty^{(2)}$	$m^{(3)}$	$\infty^{(2)}$
SET-10 [27]	$\infty^{(2)}$	$m^{(3)}$	$\infty^{(2)}$
GREEDY YIELD	$\infty^{(2)}$	$m^{(3)}$	$\infty^{(2)}$
GREEDY COM	$\infty^{(2)}$	$\frac{m}{4}^{(2)}$	$\infty^{(2)}$
LOOKAHEAD GREEDY YIELD	$\infty^{(2)}$	$m^{(3)}$	$\infty^{(2)}$
PERIODIC GREEDY YIELD ($\delta \rightarrow 0$)	$2^{(4)}$	$m^{(3)}$	$\infty^{(2)}$
BESTNEXTEVENT	$\frac{m}{2} - 4^{(6)}$	$m^{(3)}$	$\infty^{(2)}$
Any strategy	$\frac{3}{2}^{(5)}$	$\frac{m}{4}^{(2)}$	$\infty^{(2)}$

Table 2: Lower bounds for the competitive ratios of bandwidth-sharing strategies.

5 Lower Bounds on Competitive Ratios

This section provides lower bounds for the performance of the bandwidth-sharing strategies. The results are summarized in Table 2. For instance, the first entry $m^{(1)}$ in the table means that FAIRSHARE has a competitive ratio not better than m , and that the proof of this result is given by Example 1. An entry ∞ means that the strategy does not have a ρ competitive ratio for any

value of $\rho \geq 1$. Sections 5.1 to 5.6 deal with the examples that fill the lower bounds in Table 2. Finally, we give some tight bounds in Section 5.7.

5.1 Example 1

We consider a window $[T_{begin}, T_{end}] = [T, T + 1]$ with $T \gg 1$. The first $m - 1$ applications are released at time 0, and have a yield of 1 at the beginning of the window. The m -th application is released at time T . Each application \mathcal{A}_i verifies $b_i = B = 1$ and $p_i = 1$, and posts an I/O operation of volume 1 at time T . FAIRSHARE allocates a bandwidth of $\frac{1}{m}$ to all applications, resulting in a yield of $\frac{T+\frac{1}{m}}{T+1}$ for the first $m - 1$ applications, and a yield of $\frac{1}{m}$ for the last one. Therefore, the minimum yield is $\frac{1}{m}$.

However, if we had allocated all the bandwidth to the last application, its yield would be 1 and the minimum yield would be $\frac{T}{T+1}$. By taking T large enough, we see that FAIRSHARE has not a competitive ratio smaller than m for MINYIELD. This gives the first entry in Table 2.

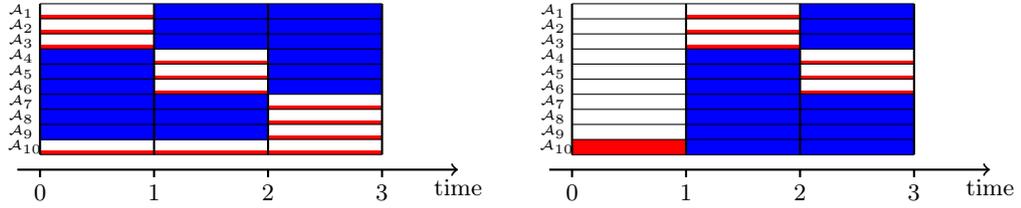


Figure 1: Illustration of the example for FAIRSHARE when all applications are released at time T_{begin} , with $K = 3$ (hence, $m = 10$). The schedule achieved by FAIRSHARE is depicted on the left, and the alternative schedule on the right. Blue areas represent computations, red ones I/Os, and white ones idle time. The height of a red area shows whether the I/O uses the whole bandwidth (like \mathcal{A}_{10} during $[0, 1]$ in the second schedule) or which fraction of it (one fourth for all other I/Os).

When all applications are released at time T_{begin} We now provide another example where each application \mathcal{A}_i is released at the beginning of the window: $\tau_i = T_{begin}$ for all i . This example establishes another negative result for the competitiveness of FAIRSHARE for MINYIELD: there is no constant competitive ratio even when all the applications are *fresh* when entering the window.

For this example, we consider a window $[T_{begin}, T_{end}] = [0, K]$ with $K \geq 3$. The $m = K^2 + 1$ applications are all released at time 0. Each application \mathcal{A}_i verifies $b_i = B = 1$ and $p_i = 1$. The applications are as follows:

- $\forall i \in [1, K]$, \mathcal{A}_i consists of the phases $v_i^{(0)} = \frac{1}{K+1}$, $w_i^{(1)} = K$;
- $\forall j \in [1, K - 1]$, $\forall i \in [Kj + 1, K(j + 1)]$, \mathcal{A}_i consists of the phases $v_i^{(0)} = 0$, $w_i^{(1)} = j$, $v_i^{(1)} = \frac{1}{K+1}$, $w_i^{(2)} = K$;
- \mathcal{A}_{K^2+1} consists of the phases $v_{K^2+1}^{(0)} = 1$, $w_{K^2+1}^{(1)} = K$.

Thus, the first K applications and the last one post an I/O operation at time 0, while the others start with a work phase. The schedule of FAIRSHARE is presented on the left of Figure 1; it is as follows:

- During the interval $[0, 1]$, applications $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_K$ and \mathcal{A}_{K^2+1} receive a bandwidth $\frac{1}{K+1}$. Applications $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_K$ finish their I/O operation at time 1. The other applications work. Applications $\mathcal{A}_{K+1}, \mathcal{A}_{K+2}, \dots, \mathcal{A}_{2K}$ finish their work at time 1.
- During the interval $[1, 2]$, applications $\mathcal{A}_{K+1}, \mathcal{A}_{K+2}, \dots, \mathcal{A}_{2K}$ and \mathcal{A}_{K^2+1} receive a bandwidth $\frac{1}{K+1}$. Applications $\mathcal{A}_{K+1}, \mathcal{A}_{K+2}, \dots, \mathcal{A}_{2K}$ finish their I/O operation at time 2. The other applications work. Applications $\mathcal{A}_{2K+1}, \mathcal{A}_{2K+2}, \dots, \mathcal{A}_{3K}$ finish their work at time 2.
- ...
- During the interval $[K-1, K]$, applications $\mathcal{A}_{K^2-K+1}, \mathcal{A}_{K^2-K+2}, \dots, \mathcal{A}_{K^2}$ and \mathcal{A}_{K^2+1} receive a bandwidth $\frac{1}{K+1}$. $\mathcal{A}_{K^2-K+1}, \mathcal{A}_{K^2-K+2}, \dots, \mathcal{A}_{K^2}$ finish their I/O operation at time K . The other applications work.

In the end, the minimum yield is the yield of application \mathcal{A}_{K^2+1} which is $\frac{1}{K+1}$.

We consider the following alternative schedule (presented on the right of Figure 1):

- During the interval $[0, 1]$, application \mathcal{A}_{K^2+1} executes its I/O operation. The other applications do nothing (even the applications that could have processed their work).
- During the interval $[1, 2]$, $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_K$ receive a bandwidth $\frac{1}{K+1}$. $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_K$ finish their I/O operation at time 2. The other applications work. Applications $\mathcal{A}_{K+1}, \mathcal{A}_{K+2}, \dots, \mathcal{A}_{2K}$ finish their work at time 2.
- During the interval $[2, 3]$, applications $\mathcal{A}_{K+1}, \mathcal{A}_{K+2}, \dots, \mathcal{A}_{2K}$ receive a bandwidth $\frac{1}{K+1}$. Applications $\mathcal{A}_{K+1}, \mathcal{A}_{K+2}, \dots, \mathcal{A}_{2K}$ finish their I/O operation at time 3. The other applications work. Applications $\mathcal{A}_{2K+1}, \mathcal{A}_{2K+2}, \dots, \mathcal{A}_{3K}$ finish their work at time 3.
- ...
- During the interval $[K-1, K]$, applications $\mathcal{A}_{K^2-2K+1}, \mathcal{A}_{K^2-2K+2}, \dots, \mathcal{A}_{K^2-K}$ receive a bandwidth $\frac{1}{K+1}$. $\mathcal{A}_{K^2-2K+1}, \mathcal{A}_{K^2-2K+2}, \dots, \mathcal{A}_{K^2-K}$ finish their I/O operation at time K . The other applications work. Applications $\mathcal{A}_{K^2-K+1}, \mathcal{A}_{K^2-K+2}, \dots, \mathcal{A}_{K^2}$ finish their work at time K .

In $[1, K]$, all applications either work or communicate. They communicate during a time at most 1, therefore they work during at least $K-2$ units of time. The minimum yield is therefore larger than $\frac{K-2}{K}$. This shows that FAIRSHARE has not a ρ competitive ratio, where $\rho \leq \frac{(K-2)(K+1)}{K}$. But $\frac{(K-2)(K+1)}{K} > K-2 = \sqrt{m-1} - 2 > \sqrt{m} - 3$; hence, the result that FAIRSHARE has not a $\sqrt{m} - 3$ competitive ratio when all applications start execution at the beginning of the window.

5.2 Example 2

We consider a window $[T_{begin}, T_{end}] = [0, 1]$. $m \geq 4$ applications are released at time 0. Each application \mathcal{A}_i verifies $b_i = B = 1$ and $p_i = 1$. We assume that m is even. We suppose that $m/2$ applications are in category A , i.e., have an I/O operation of volume $\frac{2}{m}$, followed by a work phase of length 1. The other $\frac{m}{2}$ applications are in category B , i.e., have an I/O operation of volume $\frac{2}{m}$ followed by a work phase of length $\alpha = \frac{\varepsilon}{m/2-1}$, where $\varepsilon > 0$ is a small number, and by another I/O operation phase of volume 1.

There are m I/O operations of volume $\frac{2}{m}$ posted at time 0. With a total bandwidth $B = 1$, it is impossible to complete more than $m/2$ of them by time 1. Because the m applications

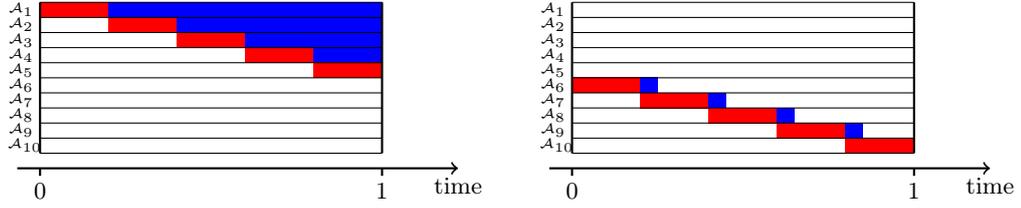


Figure 2: Illustration of Example 2 for the case $m = 10$ and $\epsilon = 1/5$. The best schedule for MINYIELD that completes the initial I/O of each application of category A (resp. of category B) is depicted on the left (resp. on the right).

are not distinguishable at time 0, the adversary might force the scheduler to complete only I/O operations of applications of category B at time 1, and have no application of category A having completed its I/O operation by the end of the window. Proceeding with this scenario, only applications of category B may have executed some work at time 1. In fact, the most efficient scenario (which is illustrated on the right side of Figure 2) is to grant the full bandwidth to each application in category B one after the other, so that $m/2 - 1$ of them can complete their work phase by the end of the window; indeed, it is impossible for all $m/2$ applications in category B to terminate their work phase by $t = 1$, and the best, in order to maximize the work done, is to schedule the I/O operations without sharing. Finally, no application of category B can complete its second I/O transfer by $t = 1$. The efficiency at time $t = 1$ is therefore upper bounded as

$$\mathcal{E} = \frac{\sum_{i \in A \cup B} V_i^{(transferred)} + \sum_{i \in B} W_i^{(done)}}{m} \leq \frac{1}{m} + \frac{(m/2 - 1)\alpha}{m} \leq \frac{1 + \epsilon}{m}.$$

A strategy that would process the I/O operations of the jobs in category A without sharing is illustrated on the left side of Figure 2 and would reach an efficiency

$$\mathcal{E}' = \frac{\sum_{i=1}^{m/2} \frac{i}{m/2}}{m} = \frac{m/2 + 1}{2m} > \frac{m}{4} \mathcal{E}$$

for ϵ small enough. Therefore, there is no competitive ratio lower than $\frac{m}{4}$ for EFFICIENCY for any strategy.

Now, if we consider the UTILIZATION objective function, we get $u = \frac{\epsilon}{m}$ with the first scenario and $u' = \frac{\sum_{i=1}^{m/2-1} \frac{i}{m/2}}{m} = \frac{m-2}{4m}$ with the second scenario. Therefore, we can get a competitive ratio arbitrarily large by choosing ϵ small enough.

Finally, for the MINYIELD objective, the best strategy is sharing I/O bandwidth equally among the $m = 2K$ applications, with gives a minimum yield of $\frac{1}{m}$. Any heuristic that serializes the I/O operations reaches a minimum yield of 0. For this example, this includes FCFS, GREEDY YIELD, GREEDY COM, LOOKAHEAD GREEDY YIELD and SET-10 (if additionally we assume that all applications belong to the same I/O-set when starting execution at time $t = 0$).

5.3 Example 3

We consider a window $[T_{begin}, T_{end}] = [1, 2]$. m applications are released at time 0. Each application \mathcal{A}_i verifies $b_i = B = 1$ and $p_i = 1$. The first $m - 1$ applications have a yield of 1 at time 1 and consists of an I/O operation of volume $\frac{\epsilon}{m}$ followed by a work phase of length 1. The last application has an initial yield of 0 at time $t = 1$, and consist of an I/O operation of volume 1.

A schedule that executes the I/O operations of the first $m-1$ applications in parallel completes these operations in time less than ϵ . These first $m-1$ applications then work for at least $1-\epsilon$ units of time, while the last application executes a fraction of its I/O operation. This schedule has an efficiency EFFICIENCY larger than $1-\epsilon$.

However, GREEDY YIELD , $\text{LOOK AHEAD GREEDY YIELD}$, $\text{PERIODIC GREEDY YIELD}$ and BEST-NEXT EVENT would all allocate the whole bandwidth to the last application, because of its low initial yield. Thus they obtain an efficiency of $\frac{1}{m}$ (although they optimize the minimum yield). The same is true for FCFS , which could allocate the whole bandwidth to the last application (since all I/O operations are posted at time $t=1$ and applications are indistinguishable). Finally, if additionally we assume all applications belong to the same I/O-set when starting execution at time $t=1$, SET-10 might do the same and select the last application. Altogether, all these heuristics have a competitive ratio of at least m .

5.4 Example 4

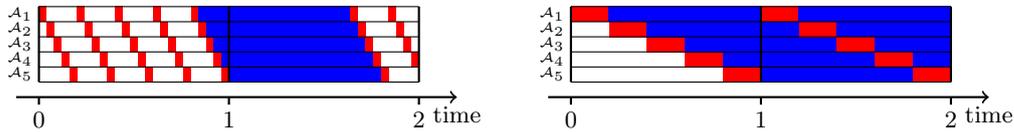


Figure 3: Illustration of Example 4 for the case $m=5$ and $T=2$. The schedule of $\text{PERIODIC GREEDY YIELD} (\frac{1}{m^2})$ is depicted on the left, and an alternative schedule on the right.

We consider a window $[T_{begin}, T_{end}] = [0, T]$. m identical applications are released at time 0. Each application \mathcal{A}_i verifies $b_i = B = 1$ and $p_i = 1$. Each application \mathcal{A}_i repeats the same cycle indefinitely, namely an I/O operation of volume $\frac{1}{m}$ followed by a work phase of length $1 - \frac{1}{m}$.

We study $\text{PERIODIC GREEDY YIELD} (\delta)$ with $\delta = \frac{1}{m^2}$. This strategy will interleave the I/O operations of all applications in a cyclic fashion, moving from one application to the next every δ units of time. This is illustrated on the left-hand side of Figure 3. Without loss of generality, \mathcal{A}_i will be the i -th application to complete its first I/O operation, at time $\frac{m-1}{m} + \frac{i}{m^2}$. The same cyclic sequence of I/O operations may well repeat after the first work phase, and so on until the end of the window. The minimum yield is that of \mathcal{A}_m , and it reaches its highest value at the end of each work phase (and then decreases during the next I/O operation). This highest value is $\frac{1}{2 - \frac{1}{m}}$. Indeed, by induction, \mathcal{A}_m finishes its k -th phase of work at time $k(1 + 1 - \frac{1}{m})$, while its progress at the end of this k -th phase is equal to k .

However, a strategy that processes the I/O operations in sequence would result in a schedule with perfect yield after time 1. Such a strategy is illustrated on the right-hand side of Figure 3. In such a strategy, application \mathcal{A}_i , with $1 \leq i \leq m$ would perform I/O in the intervals $[j + \frac{i-1}{m}, j + \frac{i}{m}]$ and work in the intervals $[j + \frac{i}{m}, j + 1 + \frac{i-1}{m}]$ for all $j \geq 0$. Therefore the yield of each application will be at least $\frac{T-1}{T}$ with that strategy.

Letting T large enough, this shows that $\text{PERIODIC GREEDY YIELD} (\delta)$ has a competitive ratio at least $2 - \frac{1}{m}$ for MIN YIELD . Letting m large enough, this shows that $\text{PERIODIC GREEDY YIELD} (\delta)$ has a competitive ratio at least 2 for MIN YIELD .

5.5 Example 5

This example provides a general lower bound for the competitive ratio of any bandwidth-sharing strategy for the MIN YIELD objective:

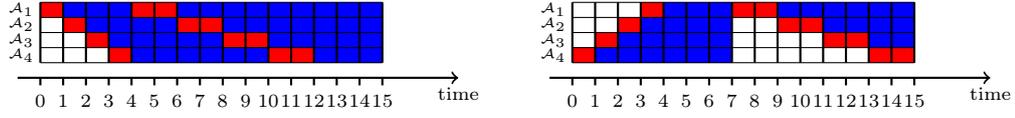


Figure 4: Illustration of Example 5 for the case $m = 4$. The schedule on the left maximizes MINYIELD. The first phase (first I/O and first computation of each application) of the schedule on the right is the best possible schedule when the first set of I/Os is completed in the worst possible order. In the second phase of the schedule, the only thing that matters is not to share the bandwidth (the order of I/O completions has no impact).

Lemma 1. *Given $\varepsilon > 0$, there does not exist any $\frac{3}{2} - \varepsilon$ competitive algorithm for MINYIELD.*

Proof. We consider a set of m applications $\mathcal{A}_1, \dots, \mathcal{A}_m$ and the window $[T_{begin}, T_{end}] = [0, 4m - 1]$. The characteristics of application \mathcal{A}_i , $1 \leq i \leq m$, are as follows:

- $b_i = B = 1$ and $p_i = 1$;
- \mathcal{A}_i starts with an I/O operation of volume 1;
- \mathcal{A}_i then has a work phase of duration $m - 2 + i$;
- \mathcal{A}_i then has an I/O operation of volume 2;
- \mathcal{A}_i then has a work phase of duration $3m$.

A possible schedule for each application \mathcal{A}_i is the following (it is illustrated on the left-hand side of Figure 4):

- \mathcal{A}_i waits during the time interval $[0, i - 1]$;
- \mathcal{A}_i performs its initial I/O operation during the time interval $[i - 1, i]$;
- \mathcal{A}_i computes during the time interval $[i, m + 2i - 2]$;
- \mathcal{A}_i performs its second I/O operation during the time interval $[m + 2i - 2, m + 2i]$;
- \mathcal{A}_i computes during the time interval $[m + 2i, 4m + 2i]$ (note that $4m + 2i \geq 4m - 1$).

All the initial I/O operations are scheduled one after the other. The same holds for all second I/O operations. Moreover, the last initial I/O operation (by \mathcal{A}_m) ends at time m , while the first second I/O operation (by \mathcal{A}_1) starts at time $m + 2 - 2 = m$. Hence, the $2m$ I/O operations are sequentialized and performed at maximal bandwidth. The application that suffers from the largest slowdown is \mathcal{A}_m . Its yield is $y_m = \frac{(4m-1)-(m-1)}{4m-1} = \frac{3m}{4m-1}$.

Because all m applications start with an I/O operation of volume 1, an arbitrary bandwidth-sharing strategy has no way to differentiate them. Hence, an adversary can decide that the initial I/O operations are completed in reverse order of application indices: the first completed I/O operation is that of \mathcal{A}_m , the second completed I/O operation is that of \mathcal{A}_{m-1} , and so on. Such a scenario is illustrated on the right-hand side of Figure 4. Then, the initial I/O operation of application \mathcal{A}_i cannot finish before time $m - i + 1$, because the total volume transferred when \mathcal{A}_i completes its I/O operation is $m - i + 1$. Then \mathcal{A}_i cannot finish its work phase before time $(m - i + 1) + (m - 2 + i) = 2m - 1$. Therefore, the second I/O operation of all applications starts at or after time $2m - 1$. Because the total volume of all second I/O operations is $2m$, at least

one application, say \mathcal{A}_{i_0} , will not finish its second I/O operation before time $4m - 1$, which is the end of the window. In total, this application \mathcal{A}_{i_0} will have executed at most 3 units of time of I/O operation and $m - 2 + i_0 \leq 2m - 2$ work units. Hence, the yield of \mathcal{A}_{i_0} satisfies $y_{i_0} \leq \frac{2m+1}{4m-1}$.

We can then derive a bound on the competitive ratio ρ of any strategy:

$$\rho \geq \frac{\frac{3m}{4m-1}}{\frac{2m+1}{4m-1}} = \frac{3m}{2m+1} \xrightarrow{m \rightarrow +\infty} \frac{3}{2}.$$

□

5.6 Example 6

We consider $m > 10$ applications and a window $[T_{begin}, T_{end}] = [m, \frac{m^2}{2} - 2m + 2]$. Each application \mathcal{A}_i is released at time 0, verifies $b_i = B = 1$, $p_i = 1$ and has an initial yield of 1 at time T_{begin} . The first application, \mathcal{A}_1 , repeats the same cycle indefinitely, namely an I/O operation of volume $\frac{1}{m}$, followed by a work phase of length $\frac{\epsilon}{m^3}$. The other $m - 1$ applications are identical and consist of an I/O operation of volume $1 + \epsilon$, where $\epsilon > 0$ is a small number, followed by a work phase of length m^2 .

A possible schedule would first share the bandwidth among the I/O operations of the last $m - 1$ applications up to time $t = m + (m - 1)(1 + \epsilon)$. The first application would then remain idle up to time t . After time t , only the first application \mathcal{A}_1 is posting I/O operations, so there is no further interference until the end of the window. The minimum yield would then be that of \mathcal{A}_1 and would verify:

$$y_1^* = \frac{T_{end} - t}{T_{end} - T_{begin}} = \frac{\frac{m^2}{2} - 2m + 2 - (m - 1)(1 + \epsilon)}{\frac{m^2}{2} - 3m + 2} \geq \frac{\frac{m^2}{2} - 3m}{\frac{m^2}{2} - 2m + 2}$$

if ϵ is small enough (say $\epsilon < \frac{1}{m}$).

We now study the performance of BESTNEXTEVENT. Intuitively, the parameters of the example have been chosen so that: (i) each time \mathcal{A}_1 posts an I/O operation, BESTNEXTEVENT assigns it the whole bandwidth immediately; and (ii) during each work phase of \mathcal{A}_1 , BESTNEXTEVENT assigns the whole bandwidth to the same application (the one that was granted the whole bandwidth during the first work phase of \mathcal{A}_1). We now prove these facts by induction on the number of events, with time $T_{begin} = m$ corresponding to event number 1.

Lemma 2. *All events taking place at a time $t < \frac{m^2}{2} - 2m$ are triggered by \mathcal{A}_1 (and maybe other applications). Odd-numbered events correspond to \mathcal{A}_1 posting an I/O operation, and even-numbered events correspond to \mathcal{A}_1 completing an I/O operation. Up to time t , BESTNEXTEVENT will never share the I/O bandwidth, and it will assign it as follows:*

- \mathcal{A}_1 is granted the whole bandwidth at every odd-numbered event;
- the same application, say \mathcal{A}_2 , is granted the whole bandwidth at every even-numbered event.

Therefore, only \mathcal{A}_1 and \mathcal{A}_2 make some progress up to time $\frac{m^2}{2} - 2m$.

Proof. We proceed by induction on events. We assume that we are currently at the i -th event and that all previous decisions have fulfilled the conditions of the lemma so far.

Case 1: $i = 2k + 1$ for $k \geq 0$. Since the lemma is true so far, $i = 2k + 1$ means application \mathcal{A}_1 is ready to perform an I/O operation, and $t = m + k \left(\frac{1}{m} + \frac{\epsilon}{m^3} \right) \leq \frac{m^2}{2} - 2m$, thus $k < m^3$, and the $m - 1$ last applications \mathcal{A}_i for $i \geq 2$ have more than 1 unit of I/O volume remaining (by induction, because \mathcal{A}_1 has been working for at most k phases of length $\frac{\epsilon}{m^3}$, hence for a duration at most ϵ). We distinguish two cases:

Case 1.1: the next event happens in $\delta \geq 1$ units of time. In this case, we consider the yield of the application that has been assigned, in the interval $[t, t + \delta]$, the least bandwidth, b , among the last $m - 2$ applications that have not started their first I/O operation and have a current progress of m . We have $b \leq \frac{1}{m-2}$. The minimum yield at the next event would verify $y \leq f(\delta) = \frac{m+b\delta}{t+\delta}$. Using $b \leq \frac{1}{m-2}$ and $t \leq \frac{m^2}{2} - 2m$, we differentiate and find that f is non increasing. Hence, we can safely replace δ by 1 to get $y \leq \frac{m+b}{t+1} \leq \frac{m+\frac{1}{m-2}}{t+1}$.

Case 1.2: the next event happens in $\delta < 1$ units of time. In this case, the next event is defined by the completion of the I/O operation of \mathcal{A}_1 (the end of the window is in $T_{end} - t > 1$ time units; therefore, the next event cannot be the completion of the I/O operation of one the last $m - 1$ applications). Letting x be the fraction of bandwidth assigned to \mathcal{A}_1 , we have $\delta = \frac{1}{mx}$. We consider the yield of the application that has the least bandwidth, b , allocated during the interval $[t, t + \delta]$, among the $m - 2$ applications whose processing have not yet started (by induction): $b \leq \frac{1-x}{m-2}$.

Clearly the minimum yield will not exceed $y = \frac{m+b\delta}{t+\delta} \leq \frac{m+\frac{1}{mx} \frac{1-x}{m-2}}{t+\frac{1}{mx}} = \frac{(m-2)m^2x+1-x}{(m-2)(1+mtx)}$.

We let $f(x) = \frac{(m-2)m^2x+1-x}{(m-2)(1+mtx)}$. We differentiate and get $f'(x) = \frac{m^3-2m^2-mt-1}{(m-2)(1+mtx)^2}$, which is positive if m is large enough, since $t \leq \frac{m^2}{2} - 2m$. Therefore, we can safely replace x by 1 and get $y \leq \frac{m}{t+\frac{1}{m}}$. This upper bound on the minimum yield can actually be achieved by allocating all the bandwidth to \mathcal{A}_1 ; indeed, \mathcal{A}_1 would then have a yield of 1 and every other application a yield at least equal to $\frac{m}{t+\frac{1}{m}}$.

To conclude, we need to show that Case 1.2 results in a better minimum yield than Case 1.1, i.e., that $\frac{m+\frac{1}{m-2}}{t+1} \leq \frac{m}{t+\frac{1}{m}}$. This last inequality is equivalent to $m \geq 1 + \frac{t}{m-2} + \frac{1}{m(m-2)}$, and is true as $t \leq \frac{m^2}{2} - 2m < \frac{(m-2)^2}{2}$. As a consequence, the best decision is to allocate all the bandwidth to \mathcal{A}_1 , which concludes the induction step for $i = 2k + 1$.

Case 2: $i = 2k$ for $k > 0$. Since the lemma is true so far, the event $i = 2k$ occurs at time $t = m + (k-1) \left(\frac{1}{m} + \frac{\epsilon}{m^3} \right) + \frac{1}{m} \leq \frac{m^2}{2} - 2m$. Hence, again, $k < m^3$ and \mathcal{A}_1 is starting a work phase. If $k = 1$ we will show that the whole bandwidth is assigned to one application, say \mathcal{A}_2 . If $k > 1$, by induction \mathcal{A}_2 has transferred a volume $\frac{(k-1)\epsilon}{m^3} < \epsilon$ so far. Regardless of the value of k , by induction, any application \mathcal{A}_i with $i \geq 3$ has not started its execution.

Case 2.1: The next event is the end of the window. In this case, let b be the minimum bandwidth allocated to an application other than \mathcal{A}_2 (and \mathcal{A}_1) during the interval $[t, t + \delta]$. Thus, $b \leq \frac{1}{m-2}$ and $\delta \geq T_{end} - t > 2$ does not depend on b . We get

$y \leq f(\delta) = \frac{m+\frac{1}{m-2}\delta}{t+\delta}$. We obtain $f'(\delta) = \frac{-m^2+2m+t}{(m-2)(\delta+t)^2} < 0$ (because $t < \frac{m^2}{2} - 2m$). We

replace δ by 2 and get $y \leq \frac{m+\frac{2}{m-2}}{t+2}$.

Case 2.2: The next event is not the end of the window. If the next event is not the end of the window, let \mathcal{A}_{i_0} be the application defining the next event at time $t + \delta$ and let x be the bandwidth allocated to \mathcal{A}_{i_0} (we will eventually show that $\mathcal{A}_{i_0} = \mathcal{A}_2$). Let \mathcal{A}_{i_1} be the application with smallest bandwidth, b , allocated to it during $t, t + \delta]$, among the applications that are not \mathcal{A}_1 , \mathcal{A}_2 and \mathcal{A}_{i_0} . Therefore, $b \leq \frac{1-x}{m-3}$ and by assumption \mathcal{A}_{i_1} has not started. The minimum yield will not exceed the yield of \mathcal{A}_{i_1} , which is $y = \frac{m+b\delta}{t+\delta}$. We distinguish two last sub-cases:

Case 2.2.1: $i_0 = 2$. We have $\delta = \frac{1+\epsilon-\frac{(k-1)\epsilon}{m}}{x}$, and $y \leq y(\mathcal{A}_{i_1}) \leq f(x) = \frac{m+\frac{1-x}{m-3} \frac{1+\epsilon-\frac{(k-1)\epsilon}{m}}{x}}{t+\frac{1+\epsilon-\frac{(k-1)\epsilon}{m}}{x}}$.

We obtain $f'(x) = \frac{(\epsilon(-k+m+1)+m)(\epsilon(k-m-1)+m(m^2-3m-t-1))}{(m-3)(\epsilon(-k+m+1)+mtx+m)^2}$. If ϵ is small enough $f'(x)$ will have the sign of $m^2(m^2-3m-t-1)$ which is positive as $t \leq \frac{m^2}{2} - 2m$. We can safely bound $f(x)$ by $f(1)$ and the minimum yield will verify $y \leq \frac{m}{t+1+\epsilon - \frac{(k-1)\epsilon}{m}}$. We point out that this bound on the minimum yield is achievable by allocating all the bandwidth to application \mathcal{A}_2 ; hence, we have an equality.

Case 2.2.2: $i_0 \neq 2$. In this case, we have $\delta = \frac{1+\epsilon}{x}$, and $y \leq y(\mathcal{A}_{i_1}) \leq f(x) = \frac{m + \frac{1-x}{m-3} \frac{1+\epsilon}{x}}{t + \frac{1+\epsilon}{x}}$. We obtain $f'(x) = \frac{(\epsilon+1)(m^2-3m-t-1-\epsilon)}{(m-3)(\epsilon+tx+1)^2}$. If ϵ is small enough, $f'(x)$ will have the sign of $m^2-3m-t-1$, which is strictly positive as $t \leq \frac{m^2}{2} - 2m$. We can safely bound $f(x)$ by $f(1)$ and the minimum yield will verify $y \leq \frac{m}{t+1+\epsilon}$. Note that, once again, this bound on the minimum yield is achievable by allocating all bandwidth to application \mathcal{A}_{i_0} ; hence we have an equality.

A quick computation shows that if $k = 1$, cases 2.2.1 and 2.2.2 are equivalent and better than 2.1 because $\frac{m}{t+1+\epsilon} > \frac{m + \frac{m^2-2}{t+2}}{t+2} \Leftrightarrow m(m-2) > \frac{2(t+1+\epsilon)}{1-\epsilon}$, which is true if ϵ is small enough. Therefore, if $k = 1$, the algorithm allocates all the bandwidth to an arbitrarily chosen application. When $k > 1$, Case 2.2.1 achieves a strictly better yield than the other two cases. Hence, BESTNEXTEVENT will always chose to allocate bandwidth to the same application.

This concludes the proof of the lemma. \square

A direct consequence of Lemma 2 is that each application \mathcal{A}_i with $i \geq 3$ will not have any progress in the interval $[m, \frac{m^2}{2} - 2m]$. Therefore, the minimum yield is at most $y' = \frac{m+2}{\frac{m^2}{2} - 2m+2} = \frac{2(m+2)}{m^2-4m+4}$. Therefore, $\frac{y^*}{y'} \geq \frac{m^2-6m}{2(m+2)} > \frac{m}{2} - 4$.

5.7 Tight Bounds

This section provides additional results on the performance of the bandwidth-sharing strategies.

5.7.1 The competitive ratio of FAIRSHARE is exactly m for MINYIELD

Section 5.1 has shown that the competitive ratio of FAIRSHARE for the objective MINYIELD is at least m . In fact, this competitive ratio is exactly m . Intuitively, with FAIRSHARE, each application progresses at full rate when computing, and at least at the fraction $\frac{1}{m}$ of the optimal rate when performing an I/O operation.

To see this formally, whenever an application \mathcal{A}_i posts an I/O operation at time t , it receives either its maximal bandwidth b_i or the fair fraction $\frac{b_i}{\sum_{j \in S(t)} b_j} \geq \frac{b_i}{m}$. Therefore, if \mathcal{A}_i was released at time τ_i with an initial yield $y_i(T_{begin})$ at the beginning of the window $[T_{begin}, T_{end}]$, we get

$$y_i(T_{end}) \geq \frac{y_i(T_{begin}) \times (T_{begin} - \tau_i) + \frac{T_{end} - T_{begin}}{m}}{T_{end} - \tau_i} \geq \frac{y_i(T_{begin}) \times (T_{begin} - \tau_i) + (T_{end} - T_{begin})}{m(T_{end} - \tau_i)} \geq \frac{y_i^{opt}}{m},$$

where y_i^{opt} is the best achievable yield for \mathcal{A}_i .

5.7.2 The competitive ratios of FCFS, SET-10, GREEDY YIELD, LOOKAHEADGREEDY YIELD, and PERIODICGREEDY YIELD are exactly m for EFFICIENCY

Any bandwidth-sharing strategy that always allocates the whole bandwidth that can be allocated (regardless of the details of the allocation) does achieve a competitive ratio of m for EFFICIENCY.

Indeed, if there is no I/O operation at time t , the efficiency is 1, and if there is some I/O operation, the efficiency is at least $\frac{1}{m}$. The five strategies listed here do allocate the whole possible bandwidth at each event. The lower bounds come from Section 5.3.

5.7.3 FAIRSHARE can be arbitrarily better than BESTNEXTEVENT

In the example of Section 5.6, the I/O operations of applications \mathcal{A}_i for $i > 1$ will always receive a fraction of the total bandwidth greater than $\frac{1}{m}$ with FAIRSHARE. These applications will complete their I/O operations at time at most $m + m(1 + \epsilon) < 3m$. After this point, there would not be further interference, and the minimum yield y^* for FAIRSHARE verifies

$$y^* \geq \frac{\frac{m^2}{2} - 5m}{\frac{m^2}{2} - 2m + 2} > \left(\frac{m}{2} - 6\right) \frac{2(m+2)}{(m-2)^2} \geq \left(\frac{m}{2} - 6\right) y_N,$$

where y_N is the minimum yield achieved by BESTNEXTEVENT, whose upper bound $\frac{2(m+2)}{(m-2)^2}$ is given in Section 5.6.

6 Performance Evaluation

We first formally define the main parameter for the experiments, namely the I/O pressure, in Section 6.1. Then, we detail the simulations conducted with synthetic traces in Section 6.2 before discussing results for the APEX workloads in Section 6.3.

6.1 I/O Pressure

For a given a steady-state window $[T_{begin}, T_{end}]$ with m applications, we compute the volume V_i that each application \mathcal{A}_i would be able to transfer if it was executed in dedicated mode throughout the window. The total I/O volume to transfer during the window is $V = \sum_{i=1}^m V_i$. The I/O pressure W is then

$$W = \frac{V}{B(T_{end} - T_{begin})}. \quad (11)$$

The I/O pressure W is the ratio of this total volume V over the maximum volume that could have been transferred during the window, assuming that it consists of a single block of data available at T_{begin} . Of course, if W exceeds 1, some transfers will necessarily be delayed. But even if W is lower than 1 but high, say 0.8, it is likely that I/O interferences and delays due to work phases will prevent to transfer the whole data volume V before the end of the steady-state window.

The I/O pressure W is a key parameter for the simulations: most bandwidth-strategies are expected to perform well when W is low, but we aim to assess how much their performance drops when W is high.

6.2 Synthetic Traces

6.2.1 Framework

The synthetic traces follow the methodology of [4] and consist of $m = 60$ applications, each of them being able to saturate the bandwidth (we have $b_i = B = 1$), with an approximate horizon of $h = 2,000,000$. For a given aimed pressure W^{GOAL} , each application \mathcal{A}_i ($1 \leq i \leq m$) is defined by the three parameters $(\mu_i, \sigma'_i, \nu_i)$: μ_i and σ'_i represents expectation and standard deviation

and impact the length of the repetitions for each applications and ν_i determines how much the application differs from one iteration to another. More precisely,

- We generate an iteration duration ω_i for \mathcal{A}_i , which corresponds to the sum of a work phase and an I/O phase if the application was alone on the platform. This duration is generated using the two parameters μ_i and σ_i' : ω_i is drawn from the normal distribution $\mathcal{N}(\mu_i, \sigma_i')$, truncated so that we consider only positive results.
- The number of iterations of application \mathcal{A}_i is $n_i = \left\lceil \frac{h}{\omega_i} \right\rceil$ so its total completion time if it were alone on the platform is close to h .
- All applications are released at time T_{begin} : $\tau_i = T_{begin}$ for each application \mathcal{A}_i . In other words, all applications are *fresh* when entering the window and have the same yiled (equal to 0). To avoid having all applications synchronized, we add a work phase $w_i^{(0)}$ whose length is generated in $\mathcal{U}[0, \omega_i]$, so that application \mathcal{A}_i effectively starts at time $w_i^{(0)}$. To simulate SET-10, we put all applications in the same I/O-set with highest priority initially, and hence process them in FCFS order at the beginning of the execution. After that each application has completed its first I/O operation, the duration of each iteration is updated on the fly, and applications get classified into different I/O-sets.
- Next, for each application, we fix the time spent on I/Os vs. on computing, so that the total pressure is around W^{GOAL} . This is done by drawing a value u_k uniformly at random in $\mathcal{U}[0, 1]$ for each application \mathcal{A}_k ($1 \leq k \leq m$), and then by defining the fraction of I/O for application \mathcal{A}_i as $\phi_i = \frac{u_i W^{GOAL}}{\sum_{k=1}^m u_k}$. This guarantees that the I/O pressure W is around W^{GOAL} . Indeed, ϕ_i allows us to define the average duration of computing phases $t_{i,cpu} = (1 - \phi_i)\omega_i$ and the average volume of I/O phases: $t_{i,io} = \phi_i\omega_i$. Thus

$$W \approx \frac{\sum_{i=1}^m t_{i,io} n_i}{B(T_{end} - T_{begin})} = \frac{\sum_{i=1}^m \phi_i \omega_i n_i}{B(T_{end} - T_{begin})} \approx \frac{\sum_{i=1}^m \phi_i T_{end}}{B(T_{end} - T_{begin})} = \frac{T_{end} W^{GOAL}}{T_{end}} = W^{GOAL}.$$

We point out that we cannot enforce exactly $W = W^{GOAL}$ due to the randomness in the generation of instances.

- Finally, for each application \mathcal{A}_i , we consider a noise parameter ν_i to generate iterations of different lengths. For all $j \leq n_i$, we draw two variables $\gamma_{cpu}^{(j)}$ and $\gamma_{io}^{(j)}$ from a uniform distribution $\mathcal{U}[-\nu_i, \nu_i]$ and let $w_i^{(j)} = (1 + \gamma_{cpu}^{(j)})t_{i,cpu}$ and $v_i^{(j)} = (1 + \gamma_{io}^{(j)})t_{i,io}$.

6.2.2 Results for Synthetic Traces

Still following the methodology of [4], the experiments are conducted by varying four different key parameters for the 60 applications. For the application length, we consider 20 applications of medium size, and then a proportion of smaller and larger applications, as determined by the parameter n_{small} (number of small applications). The standard deviation is dictated by parameter σ , the noise is set to ν , and the pressure is W^{GOAL} . Overall, the applications are as follows:

- n_{small} *small* applications with parameters $(\mu = 1\,000, \sigma' = \mu\sigma, \nu)$;
- 20 *medium* applications with parameters $(\mu = 10\,000, \sigma' = \mu\sigma, \nu)$;
- 40 - n_{small} *big* applications with parameters $(\mu = 100\,000, \sigma' = \mu\sigma, \nu)$.

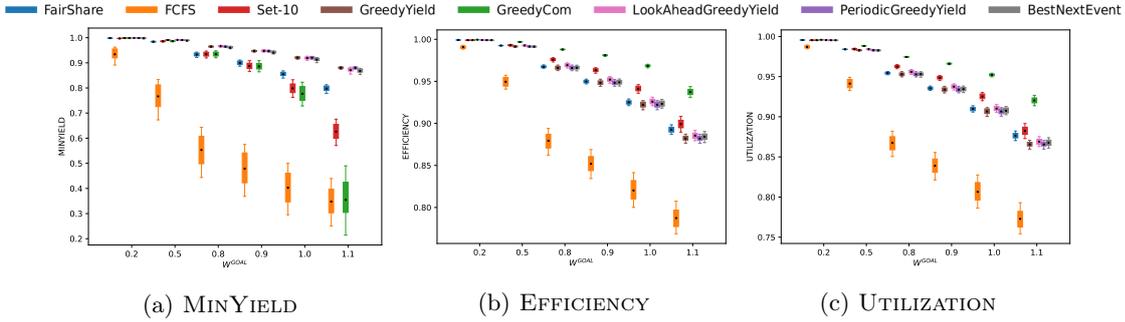
The time window is defined as $[T_{begin} = 0, T_{end} \approx h]$, where T_{end} is the smallest time required to complete an application when it is running alone on the platform. Each application is generated in such a way that T_{end} is approximately equal to $h = 2,000,000$. For each set of experiments, we study the results of all the heuristics for the three objectives (MINYIELD, EFFICIENCY, UTILIZATION).

Finally, for each set of parameters, we generate $K = 200$ instances on which we test all the heuristics presented in Section 4 (including the reference heuristics FAIRSHARE, FCFS and SET-10). In the following sections, we vary the parameters one by one and present the results on different figures. Each set of instances is represented by a boxplot of a color associated with the studied heuristic. In these boxplots, the 25th and 75th percentiles of the K instances delimit the box, and the 10th and 90th percentiles are at the end of the whiskers. Finally, the boxplots are connected by a line passing through their means.

Impact of the target I/O pressure (W^{GOAL}). We first set $\nu = \sigma = 0.5$, and $n_{small} = 20$ (20 applications of each category), and we present the results of the experiments for all values of the aimed I/O pressure $W^{GOAL} \in [0.2, 0.5, 0.8, 0.9, 1.0, 1.1]$ on Figure 5. As soon as the pressure increases, we see that the state-of-the-art strategies FAIRSHARE, FCFS, and SET-10, along with the new GREEDYCOM, fail to keep a minimum yield close to 1. The other newly proposed strategies, which all focus on the yield, successfully maintain a very high minimum yield, and achieve a similar performance which very slowly degrades when the I/O pressure increases. LOOKAHEADGREEDYIELD and BESTNEXTEVENT achieve a very slightly worse performance than GREEDYIELD and PERIODICGREEDYIELD. This counter-intuitive result may be explained by the fact that an I/O phase is always followed by a computation phase during which the progress rate of an application is perfect. Hence, what heuristics GREEDYIELD and PERIODICGREEDYIELD may lose in terms of application yield during an I/O phase may be made up later on in the subsequent computation phase. The fact that GREEDYIELD, PERIODICGREEDYIELD and LOOKAHEADGREEDYIELD achieve a minimum yield no worse than that of BESTNEXTEVENT, a costly strategy which exhaustively looks for the best solution, strongly validates these three low-cost strategies.

The classical FCFS strategy also has very poor results in terms of efficiency and utilization, while GREEDYCOM is actually the best for these objective functions since it will complete short I/Os first, with a risk of starvation for applications with long I/Os. This explains the poor performance of GREEDYCOM for MINYIELD for higher values of W^{GOAL} . The yield-based strategies tend to balance the yield of all applications, which optimizes the MINYIELD. However, not allowing any application to starve requires prioritizing some long I/Os that saturate the bandwidth, which can negatively impact both EFFICIENCY and UTILIZATION. The underlying tradeoff explains why heuristics achieving a significantly better performance than FAIRSHARE for the MINYIELD usually achieve slightly worse performance than FAIRSHARE for EFFICIENCY and UTILIZATION. However, the performance degradation in terms of either EFFICIENCY or UTILIZATION is quite small (under 5%) and only happens for the largest value of I/O pressure. For all but the largest value of W^{GOAL} , LOOKAHEADGREEDYIELD even achieves better EFFICIENCY and UTILIZATION than FAIRSHARE.

Impact of iteration size (ω) and I/O fraction (ϕ). We investigate the impact of ω and ϕ on the yield of each individual application for a fixed set of parameters: $\nu = 0.5$, $\sigma = 0.5$, $n_{small} = 20$, and $W^{GOAL} \in \{0.5, 0.8, 1.1\}$. More precisely, for each experiment \mathcal{E}_k , we define two permutations on the index set $\{1, 2, \dots, 60\}$ to sort the applications by increasing values of ω (permutation π_ω^k) or of ϕ (permutation π_ϕ^k). For each value of W^{GOAL} , we compute the average

Figure 5: Impact of the aimed I/O pressure (W^{GOAL}).

yield of applications in each position i under each permutation, denoted as $\bar{y}_i^{(\omega)}$ (resp. $\bar{y}_i^{(\phi)}$), for $i \in \{1, 2, \dots, 60\}$. It is computed as follows:

$$\bar{y}_i^{(\omega)} = \frac{1}{K} \sum_{k=1}^K y_{\pi_{\omega}^k(i)} \quad \text{and} \quad \bar{y}_i^{(\phi)} = \frac{1}{K} \sum_{k=1}^K y_{\pi_{\phi}^k(i)}.$$

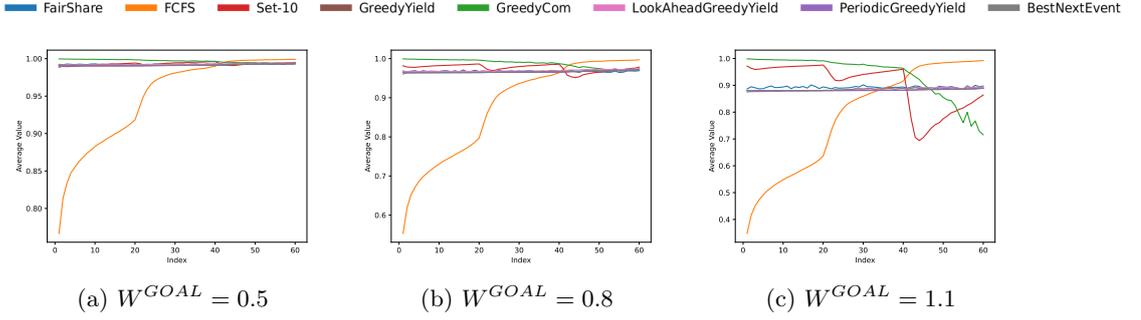
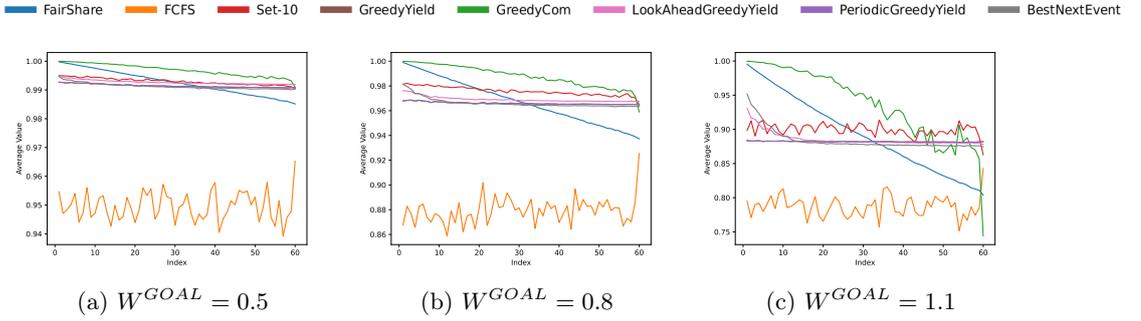
We then plot the value of $\bar{y}_i^{(\omega)}$ for i varying in $[1, 60]$ on Figure 6 and the value of $\bar{y}_i^{(\phi)}$ on Figure 7. Therefore, the leftmost point on Figure 6 (respectively, on Figure 7) corresponds to the average yield of the application with the smallest value of ω (resp., of ϕ), while the rightmost point corresponds to the average yield of the application with the largest value of ω (resp., of ϕ).

Impact of iteration size (ω). On Figure 6, we observe that the differences between the heuristics are more pronounced when W^{GOAL} increases. This is because the increase in W^{GOAL} increases the I/O interferences. For this reason, we now focus on the figure on the right (case $W^{GOAL} = 1.1$). First, we can observe that the variation of ω has little impact on the yields achieved by FAIRSHARE, GREEDY YIELD, LOOKAHEADGREEDY YIELD, PERIODICGREEDY YIELD, and BESTNEXTEVENT. GREEDY YIELD, LOOKAHEADGREEDY YIELD PERIODICGREEDY YIELD, and BESTNEXTEVENT tend to balance the yield of the different applications, resulting in a constant function. For FAIRSHARE, there seems to be no correlation between ω and the yield. This can be explained by the fact that there is no correlation between ω and ϕ in the generated instances.

This figure is more enlightening for the other heuristics. First, the yield seems to be positively correlated with ω for FCFS. This is because ϕ is not correlated with ω . Hence, a small value of ω corresponds to short I/O phases. For FCFS, the longest I/Os will saturate the bandwidth more often. Indeed, a single application can saturate the bandwidth, so when a long I/O is executed, all the other applications wanting to perform some I/O are stopped. For an application with a small I/O, the waiting time may be very long compared to its size, and the next waiting phase may also come quickly if some long I/O is posted between two of its I/O phases. Therefore, applications with short I/Os, i.e., a small value of ω , will spend a large part of their time waiting.

We observe the opposite behavior for the GREEDYCOM strategy since, this time, small I/Os are given priority. As previously mentioned, a low value of ω induces short I/Os; hence, the yield decreases with ω .

Finally, this figure perfectly illustrates the behavior of SET-10. Indeed, we can clearly distinguish the three steps corresponding to the three priority categories in these synthetic traces.

Figure 6: Yields sorted by iteration size (ω).Figure 7: Yields sorted by I/O fraction (ϕ).

Moreover, within each of these steps, we see that the yield increases with ω , just like for FCFS. This is because SET-10 behaves like FCFS inside each of these categories.

Impact of I/O fraction (ϕ). Figure 7 may appear a bit more cluttered, but illustrates some interesting behaviors. Once again, we only focus on the figure on the right, that is, on the case $W^{GOAL} = 1.1$. First, we can see a difference between GREEDY YIELD (hidden under PERIODIC GREEDY YIELD) and LOOK AHEAD GREEDY YIELD for small values of ϕ , showing that the best immediate choice is not always the best choice in the long term. We can also see that BEST NEXT EVENT favors applications with smaller I/Os so that the next event arrives as soon as possible and the yield do not have the time to significantly decrease (because of I/O interference). FCFS is erratic because an application with a large ω but a small ϕ will still have larger I/O volumes per phase than an application with a small ω but a large ϕ . The same argument also explains the non-monotonic behavior of GREEDY COM when ϕ becomes large. The only heuristic that is strongly (negatively) correlated with ϕ is FAIR SHARE. Indeed, the larger ϕ , the longer the application will spend performing I/Os, and the lower the yield will be, whereas in a working phase, the instantaneous yield is 1. The linear shape of this curve is related to the uniform distribution of ϕ .

Impact of the other parameters. In Section B.1 of the Appendix, we report on experiments detailing the impact of the other parameters, namely the number of small applications (n_{small}), the standard deviation (σ), and the noise (ν). The number of small applications (n_{small}) and the noise (ν) do not impact the performance of the strategies. High values of the standard deviation (σ) only impacts the MIN YIELD achieved by FCFS and SET-10, and it impacts them negatively.

6.2.3 Synthesis of the Evaluation on Synthetic Scenarios

For the MINYIELD objective, the greedy strategies GREEDYIELD, LOOKAHEADGREEDYIELD and PERIODICGREEDYIELD achieve comparable performance, and much better performance than the competitors FCFS, FAIRSHARE and SET-10. Furthermore, we stress that the complicated strategy BESTNEXTEVENT does *not* turn out to be superior to the simpler ones, which is good news: GREEDYIELD, LOOKAHEADGREEDYIELD and PERIODICGREEDYIELD are all simple to implement and use. Finally, for the EFFICIENCY and UTILIZATION objectives, GREEDYCOM is the best, FCFS is the worst, and the other strategies achieve close performance in between.

6.3 Evaluation on APEX workloads

6.3.1 Apex Traces [17]

We use the workload and platform described in [17] to evaluate the bandwidth-sharing strategies on realistic scenarios. The table in Figure 3 of [17] describes two very different workloads: the NERSC workload and the TRILAB workload. The NERSC workload contains a large number of small applications (e.g., a single pipeline of the SkySurvey workflow runs over 24 cores for 4 hours, but the set of SkySurvey workflows represents 12% of the overall core-hours used by the workload on the machine), some large applications (GTS spans over 16,512 cores, or 1/8 of the platform, for 48 hours), and some very long running applications (CESM applications run for 10 days over 8,000 cores). The TRILAB workload contains a more homogeneous set of applications (4096 to 32768 cores), and all applications run for a significantly longer time (64 hours for the smallest duration, and up to 12 days for the longest). From this table, we take the application walltime, its number of cores, and the data information to build a possible schedule on the target machine. The table reports how much input, output, and checkpoint data each application uses. The trace does not provide fine-grain information on how the data is consumed or produced. To simulate the schedules, we assume that all inputs happen at the beginning of the application, which then does periodic checkpoints, and eventually outputs all its output data just before its completion. As is often the practice in HPC centers [10], we use a fixed period of 1 hour for the checkpoint interval.

Based on this information, we generate machine schedules using the first-fit strategy. We consider independently NERSC or TRILAB workloads and, for a given workload, we randomly pick applications from this workload, and place them on the schedule, until two conditions are met: 1) the schedule follows the application workload distribution described in the APEX table, and 2) the schedule represents at least 3 months of machine use. For each target machine considered (see below), we generate 100 schedules for the TRILAB workload and 100 schedules for the NERSC workload. In each schedule, we then find the 20 longest windows during which no application is joining or leaving the machine, to fit the analysis conditions with steady-state windows described in Sections 3 and 4. We then assume that each application joined the system at the window start ($\tau_i = T_{begin}$ for each application \mathcal{A}_i).

On the Celio system², both the NERSC and TRILAB workloads represent a small I/O pressure (about 0.15 on average). However, I/O pressure is a metric that tends to increase as we consider larger platforms and newer systems. In [15], the authors look at the architectural trends and system balance of the top 500 supercomputers. The Parallel File System (PFS) bandwidth is studied for systems that existed between 2009 and 2018. The authors compare the PFS bandwidth with the aggregated memory bandwidth. The different systems have a ratio of aggregated

²Celio is the platform used for the NERSC and TRILAB workloads [17].

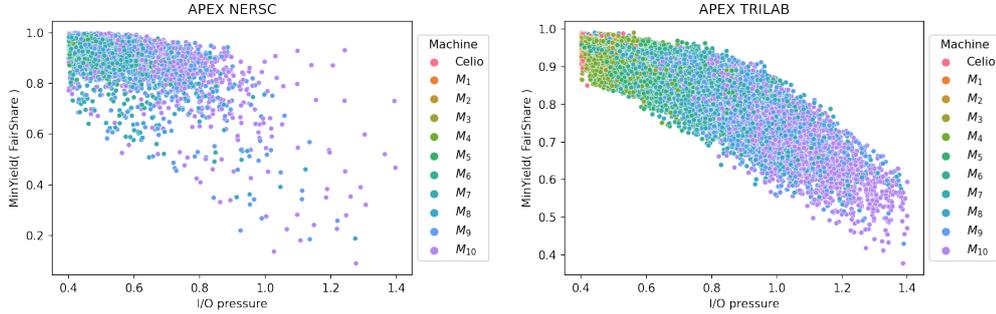


Figure 8: MINYIELD of the FAIRSHARE strategy for the NERSC and TRILAB workloads, as a function of the I/O pressure and of the target platform.

memory bandwidth by PFS bandwidth between 50 and 17000, with an average of 13,353, without a clear trend in time.

The ratio of aggregated memory bandwidth per computing performance, however, shows a clear diminishing trend. As an example, this ratio decreased by a factor 9 between the No. 1 machine in 2009 and the No. 1 machine in 2018. As a consequence, the ratio between the PFS bandwidth and the computing performance also has a clearly decreasing trend. In [4], the authors note that this ratio has decreased by a factor 24.8 over 20 years. Over long periods of time, it looks like the trend of the PFS bandwidth progresses more slowly than the computing power by a linear factor.

To study how the different algorithms behave with higher values of the I/O pressure, we have considered a set of target machines that are scaled versions of the Celio system. Let C_c and C_{bw} be respectively the total number of cores and system bandwidth of Celio, and let t represent the passing time. The system M_t has $C_c \times 2^{\frac{t}{\alpha}}$ cores (representing a doubling of computing power every α time units, in accordance of the observed progression in [15]), and $C_{bw} \times 2^{\frac{t}{\alpha}/t}$ system bandwidth, following the observation above. M_y , $y > 0$, represents machines built y time units later than the Celio machine, and for each target machine, we compute the schedules and corresponding windows for both workloads. We thus obtain a range of I/O pressures between 0.15 and 1.4, and simulate the behavior of the bandwidth-sharing strategies in each window, to evaluate our metrics as a function of the I/O pressure.

6.3.2 MINYIELD of FAIRSHARE on APEX Scenarios

We use the FAIRSHARE strategy as the basis for our evaluation, so we study first how FAIRSHARE behaves as a function of the I/O pressure. Figure 8 presents the MINYIELD obtained by the FAIRSHARE strategy within each of the 2,000 windows obtained during the simulation, as a function of the I/O pressure observed inside each window. The color of points denote on which target platform this I/O pressure and MINYIELD were observed.

On the NERSC workload, we see that the MINYIELD stays above 0.8 when the I/O pressure is low (0.4), and the distribution tends to decrease as the I/O pressure increases, with some scenarios that obtain a MINYIELD under 0.5 when the I/O pressure is 1, and the number of runs that have a low MINYIELD continue to increase as the I/O pressure continues to increase. The machine scale has some impact on the I/O pressure inside the various windows, but most of the runs present a relatively low I/O pressure, and a MINYIELD of 1 for FAIRSHARE is observed for some runs with high I/O pressures (up to 1.4). We conjecture that this is a consequence

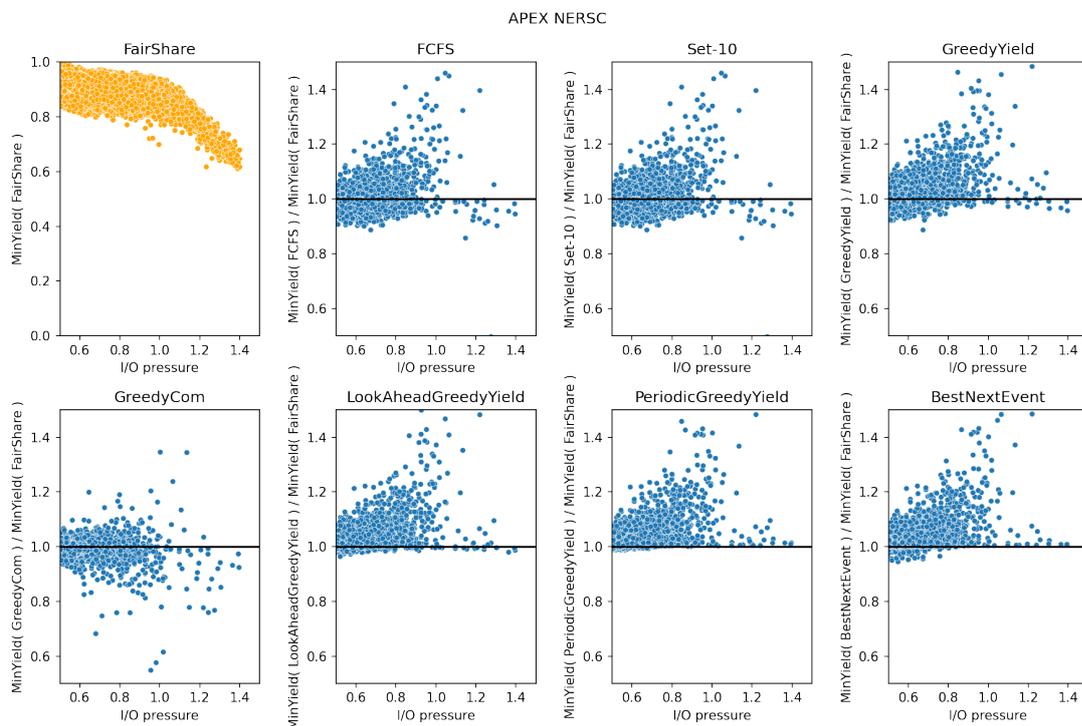


Figure 9: MINYIELD of all strategies, as a ratio of the MINYIELD with the FAIRSHARE strategy for the same experiment, for the NERSC workload, as a function of the I/O Pressure.

of the relatively small windows for the NERSC workload. Small scale, short lived applications constitute the bulk of many windows of the NERSC workload. These applications only do I/O at the beginning and end of their execution, limiting the opportunities for interferences. These I/O are also small (even relative to the short duration of the application), so when they interfere (which is unavoidable when the I/O pressure is higher than 1), they still reduce the MINYIELD by only a fraction. Only on windows that feature the few larger applications and those with costly checkpoints can we observe a measurable decrease of MINYIELD for FAIRSHARE.

This conjecture is corroborated by the measurement of the TRILAB workload. The same trends for this workload are more clearly marked: the larger the machine, the higher the I/O pressure, and the higher the I/O pressure, the lower MINYIELD for FAIRSHARE. Although there are no scenarios where MINYIELD goes under 0.4, there are also no scenarios with a MINYIELD close to 1 when the I/O pressure is above 1. The windows are much longer in the TRILAB experiments, and applications have time to checkpoint regularly during these windows. As a consequence, interferences between applications that have overlapping I/O create slowdowns that reduce the MINYIELD. We note from the left graph of Figure 8 that no NERSC scenario on the Celio platform obtains an I/O pressure of at least 0.5, while some scenarios of the TRILAB workload can saturate the I/O bandwidth. We explore the characteristics of the windows duration, size and utilization in more details in Section B.2 of the Appendix.

6.3.3 MINYIELD of All Strategies on APEX NERSC Scenarios

Figure 9 presents all the scenarios used in Figure 8 for the NERSC workload, and considers the MINYIELD of each strategy as a ratio of MINYIELD for FAIRSHARE, with an independent graph per strategy. As a reference, the MINYIELD of FAIRSHARE is also presented in a different color. A value of 1 of the ratios means that the target strategy obtains the exact same MINYIELD as FAIRSHARE for the scenario, while a value higher than 1 means that a higher MINYIELD than FAIRSHARE is obtained for this scenario, and a value lower than 1 that on this scenario, the strategy obtains a lower MINYIELD than FAIRSHARE.

There are three classes of graphs in this figure. The strategy GREEDYCOM presents on average a ratio distributed approximately uniformly between 0.9 and 1.1. This means that this strategy fails to reliably improve the MINYIELD in at least half of the scenarios. The second set of graphs show that FCFS, SET-10, and GREEDYIELD have a non-negligible set of scenarios where they decrease MINYIELD compared to FAIRSHARE, but as the I/O pressure increases they tend to behave slightly better than MINYIELD on average (with still a risk of significant performance degradation for all I/O pressures). When they experience gains, the gains are more pronounced for high I/O pressures. SET-10 and FCFS behave strictly identically over the NERSC workload, and this is because the NERSC workload featuring very small windows with typically at most one phase for many applications, SET-10 does not have time to learn the phases, and thus puts all the applications in the same set, behaving as FCFS.

The third set of graphs include LOOKAHEADGREEDYIELD, PERIODICGREEDYIELD, and BESTNEXTEVENT. These three strategies have a very high probability of increasing MINYIELD compared to FAIRSHARE, and that performance increase tends to be higher as the I/O pressure increases. BESTNEXTEVENT is the strategy of this set that features the highest risk of decreasing MINYIELD (although the decrease is limited to 95% of the MINYIELD of FAIRSHARE in the worst scenario), while PERIODICGREEDYIELD has almost no scenario with a MINYIELD lower than FAIRSHARE.

6.3.4 MINYIELD of All Strategies on APEX TRILAB Scenarios

Figure 10 presents the same evaluation, for the TRILAB workload (relative to the experiments shown in the right graph of Figure 8). With this workload, the ratio of MINYIELD behaves differently than with the NERSC workload. Overall, all strategies tend to behave better (with relatively less scenarios presenting a ratio lower than 1), and the gains over FAIRSHARE are on average higher for all strategies at low I/O pressure and for most strategies at high I/O pressure.

GREEDYCOM presents better behaviors than over the NERSC workload, with only a few scenarios underperforming FAIRSHARE, until the I/O pressure reaches a ratio of 1, i.e., until the system reaches saturation of the communication system. Then, the performance of GREEDYCOM quickly drops dramatically, with eventually all scenarios obtaining a lower MINYIELD than FAIRSHARE.

FCFS, SET-10 and GREEDYIELD continue to behave similarly, but the trend is more clear, with a significant risk of MINYIELD degradation for low I/O pressures, but significant gains as the I/O pressure, and consistent gains at I/O saturation (when the I/O pressure is higher than 1). FCFS and SET-10 continue to behave identically. However, this time this is not due to a lack of time to learn the periodicity of the applications: in the TRILAB workload, each application has a minimum of 5 phases during the window, which is long enough to converge on the phase duration and categorize the application in the appropriate set. Because all applications checkpoint with the same approximate checkpointing period, only the duration of the checkpoint operation can define different categories of phases. The checkpoint duration of the different applications can vary by an order of magnitude or more in the TRILAB workload, but the duration of the

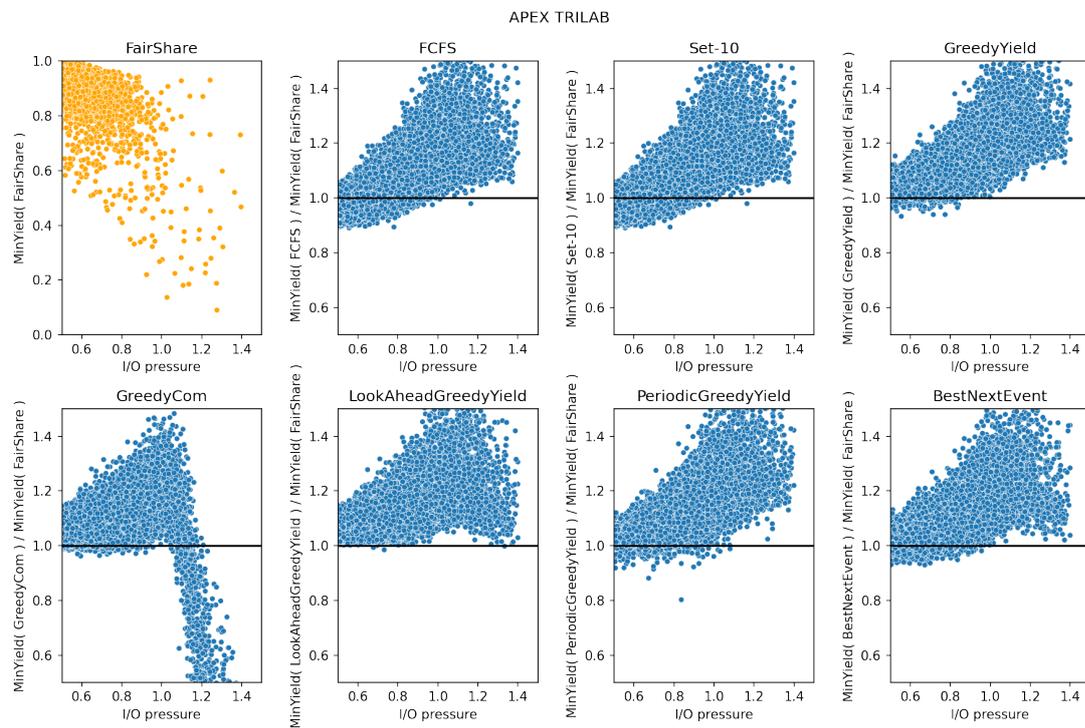


Figure 10: MINYIELD of all strategies, as a ratio of the MINYIELD with the FAIRSHARE strategy for the same experiment, for the TRILAB workload, as a function of the I/O pressure.

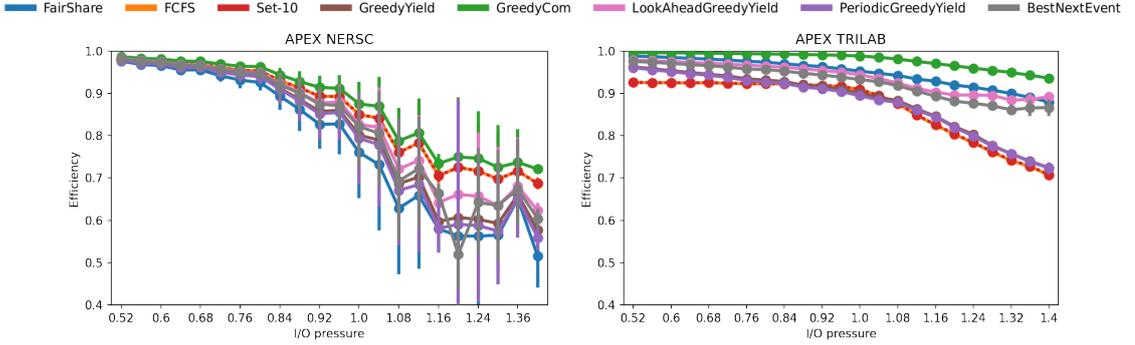


Figure 11: EFFICIENCY of all strategies for the NERSC and TRILAB workloads, as a function of the I/O pressure.

slowest checkpointing operation still remains small compared to the checkpointing period. As a consequence, the SET-10 strategy tends to put all applications in the same category, and falls back to applying the FCFS strategy.

Among the three winning strategies for the NERSC workload, LOOKAHEADGREEDYIELD, PERIODICGREEDYIELD and BESTNEXTEVENT, the trends observed for the NERSC workload are enforced: until the system reaches I/O saturation, PERIODICGREEDYIELD and BESTNEXTEVENT feature a few scenarios where MINYIELD can be slightly decreased compared to FAIRSHARE, but in most scenarios (and in almost all scenarios for LOOKAHEADGREEDYIELD), these strategies improve MINYIELD, and that improvement becomes higher as the I/O pressure increases. Contrarily to NERSC traces, GREEDYIELD performs similarly to PERIODICGREEDYIELD and BESTNEXTEVENT on the TRILAB traces.

On these longer windows, the I/O pressure seems to have a more significant impact than on the smaller windows of the NERSC workload, and as the I/O pressure increases, the gains relative to FAIRSHARE tend to increase (see Section B.2 in the Appendix for detailed results). When the I/O pressure is higher than 1, interference is unavoidable, and the I/O scheduling strategy becomes critical to the performance of applications. Naive strategies, or strategies that are not well suited for the irregular nature of the applications present in these workloads, have then a higher risk of taking the wrong decision and performing worse than FAIRSHARE.

6.3.5 EFFICIENCY of All Strategies on APEX Scenarios

Figure 11 presents the mean and standard deviation of the EFFICIENCY metric for each strategy as a function of the I/O pressure. To synthesize these graphs, we split the I/O pressure domain in 25 intervals and compute the mean EFFICIENCY value and its standard deviation for all scenarios with an I/O pressure that falls in this interval. The point is presented at the middle of the interval.

The NERSC and TRILAB workloads present both some commonalities and some significantly different features. In the NERSC workload, EFFICIENCY quickly drops as the I/O pressure increases for all strategies, while each strategy seems to hold its EFFICIENCY until the system reaches saturation (I/O pressure of 1) in the TRILAB workload. Once the I/O pressure is above 1, EFFICIENCY drops with the I/O pressure for both workloads, but this drop is more pronounced, and becomes chaotic, for the NERSC workload, while the EFFICIENCY with the TRILAB workload remains stable and supports higher I/O pressures for all strategies.

EFFICIENCY measures the sum of actual progress of all applications throughout the window.

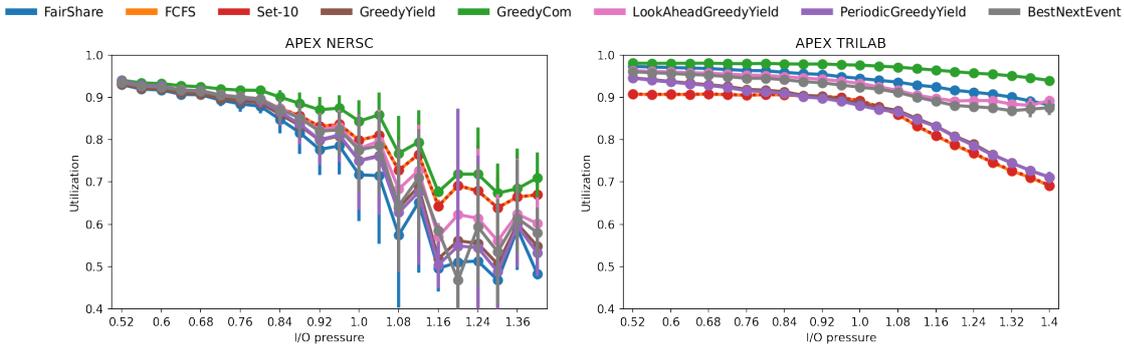


Figure 12: UTILIZATION of all strategies for the NERSC and TRILAB workloads, as a function of the I/O Pressure.

As NERSC has on average much smaller windows than TRILAB, the effect of a few bad I/O schedule decisions can be much more impactful on a small window than on a large one. This explains the chaotic EFFICIENCY measurement on the NERSC workload compared to TRILAB.

TRILAB is also a workload on which it is easier, for all strategies, to maintain a high EFFICIENCY compared to NERSC, because the windows feature a lower number of long and large-scale applications, where the I/O is close to periodic per application (mostly driven by fixed-period checkpointing), allowing many opportunities to overlap I/O operations and computation. However, at high I/O pressure, we observe three groups of strategies on the TRILAB workload: GREEDYCOM, which targets a balance of I/O operation progress, remains the most efficient; FAIRSHARE, LOOKAHEADGREEDYIELD and BESTNEXTEVENT provide a similar EFFICIENCY, slightly under GREEDYCOM; and in the third group, SET-10, FCFS (hidden by SET-10 in the figure), GREEDYIELD (hidden by PERIODICGREEDYIELD in the figure), and PERIODICGREEDYIELD present the worst EFFICIENCY. As the I/O pressure is above 1, contentions are unavoidable, and the strategies that pursue too eagerly an optimization of MINYIELD fail at providing a good EFFICIENCY. FCFS and SET-10 take I/O scheduling decisions that are detrimental to EFFICIENCY because the I/O that is favored is arbitrary.

In the NERSC workload, the metric is too chaotic at high I/O pressure to define a clear order, but GREEDYCOM remains the strategy with the highest EFFICIENCY, which is expected as GREEDYCOM targets this metric.

6.3.6 UTILIZATION of All Strategies on APEX Scenarios

Figure 12 presents the mean and standard deviation of the UTILIZATION metric for each strategy as a function of the I/O pressure. We used the same binning approach as for Figure 11 to present trends from individual scenarios.

UTILIZATION is overall lower in the NERSC workload than in the TRILAB workload. This is corroborated by the window characteristics detailed in Section B.2 of the Appendix: windows in the NERSC workload have on average a lower UTILIZATION than for the TRILAB workload, even without considering I/O interferences.

As the I/O pressure increases and in the saturated domain in particular, I/O interferences reduce even more UTILIZATION, for all strategies and in all scenarios. GREEDYCOM, which targets a balance of I/O operation progress, shines with this metric as well as for EFFICIENCY, at the cost of a worst MINYIELD as illustrated in Figures 9 and 10. On these practical scenarios, EFFICIENCY and UTILIZATION seem to behave very similarly.

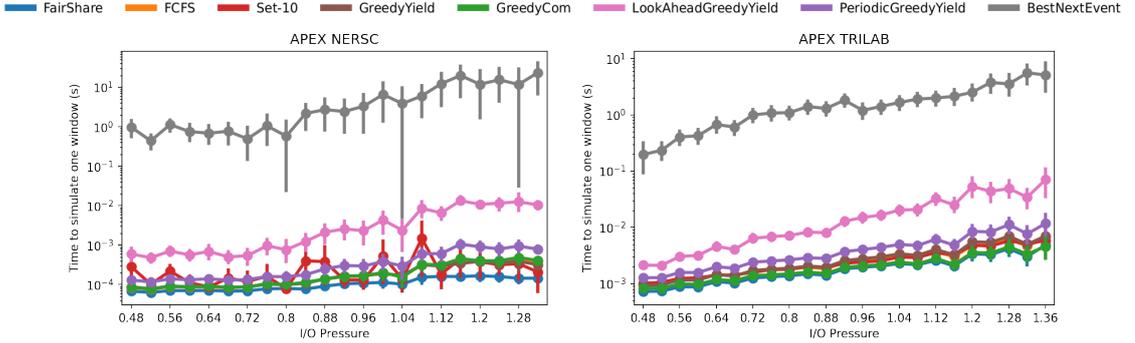


Figure 13: Simulation time of all strategies for the NERSC and TRILAB workloads, as a function of the I/O Pressure.

6.3.7 Computation Time of All Strategies on APEX Scenarios

Last, we look at the computation time of the different strategies. Each strategy decides to take a scheduling decision at different times, and each scheduling decision impacts the order of events and when the next scheduling decision will happen. To compare the computation time of the different strategies in a practical setup, we have thus chosen to measure the entire simulation time of a given window, for a given strategy. This time includes the simulation, but also, for each scheduling event, the cost of computing the decision, as an implementation of the strategy would have to do.

The mean and standard deviation of the time to simulate each of the windows is presented in Figure 13. We used the same binning approach as for Figures 11 and 12, in order to present trends. As the different strategies exhibit very different simulation computation times, the time axis is using a logarithmic scale.

BESTNEXTEVENT is the only strategy that requires significant computation time, with a few seconds (and never more than 60 seconds) needed to simulate an entire window for both the NERSC and TRILAB workloads. The second highest demanding strategy, LOOKAHEADGREEDYIELD, only requires 10s of milliseconds to simulate an entire window, and all the other strategies are yet an order of magnitude faster.

Although BESTNEXTEVENT is the most demanding strategy in terms of computational complexity, its runtime remains small enough to be considered in practice. The PERIODICGREEDYIELD strategy, which needs to re-compute regularly the entire schedule, can be called with a very small period (seconds to milliseconds), as its computational demand on realistic scenarios is achieved in a fraction of this time.

6.3.8 Synthesis of the Evaluation on APEX Scenarios

Overall, LOOKAHEADGREEDYIELD is the strategy that shows the best performance for the MINYIELD metric on the most variety of scenarios, closely followed by PERIODICGREEDYIELD and BESTNEXTEVENT. PERIODICGREEDYIELD requires to re-compute goals at a higher frequency, namely twice the frequency of the other greedy strategies with our choice for the periodicity of external events; but LOOKAHEADGREEDYIELD remains more costly, because each decision requires a set of goal computations, one per active application, and BESTNEXTEVENT, with its exhaustive search, is far more computationally demanding. GREEDYCOM is a strategy

that would perform the best on the UTILIZATION and EFFICIENCY metrics, and its MINYIELD remains reasonable on the TRILAB workload, as long as the I/O pressure is not saturated, but it is a risky choice for the NERSC workload.

7 Conclusion

This work has revisited I/O bandwidth-sharing strategies for concurrent applications. Our main contributions are two-fold. On the theoretical side, we have provided the first competitive ratios for such strategies, owing to a rigorous framework based upon steady-state windows. These competitive ratios are mostly negative. In particular, the lower bound for MINYIELD is as high as the (order of) number of applications for all strategies except PERIODICGREEDYIELD. These results bring new insights on the hardness of the problem, and lay the foundations for the study of its complexity. On the practical side, we have introduced several new greedy heuristics and have compared them to well-established strategies such as FCFS, FAIRSHARE and SET-10. We have used a comprehensive set of experiments, some based upon synthetic traces and some based upon an extended version of APEX traces. In both cases, the well-established strategies perform worse, and often much worse, than the new heuristics. As a global conclusion, although there is no absolute winner for all scenarios and objectives, we recommend using LOOKAHEADGREEDYIELD, which achieves an excellent performance for MINYIELD on all scenarios, achieves better EFFICIENCY and UTILIZATION than FAIRSHARE for the NERSC workload and comparable ones for the TRILAB and synthetic workloads. LOOKAHEADGREEDYIELD requires knowing the volume of an I/O operation when it is posted. If such an information is not available, one can use PERIODICGREEDYIELD: it achieves very good MINYIELD on all scenarios, achieves better EFFICIENCY and UTILIZATION than FAIRSHARE for the NERSC workload, comparable ones for the synthetic workload, but worse ones for the TRILAB workload. We acknowledge that these results have been obtained by simulation, and we recommend that the I/O community would implement LOOKAHEADGREEDYIELD and PERIODICGREEDYIELD for further assessment.

References

- [1] G. Aupy, A. Gainaru, and V. Le Fèvre. Periodic I/O scheduling for super-computers. In S. A. Jarvis, S. A. Wright, and S. D. Hammond, editors, *PMBS workshop*, volume 10724 of *Lecture Notes in Computer Science*, pages 44–66. Springer, 2017.
- [2] G. Aupy, A. Gainaru, and V. Le Fèvre. I/O scheduling strategy for periodic applications. *ACM Trans. Parallel Comput (TOPC)*, 6(2):1–26, 2019.
- [3] J. L. Bez, S. Byna, and S. Ibrahim. I/O access patterns in HPC applications: A 360-degree survey. *ACM Computing Surveys*, jul 2023.
- [4] F. Boito, G. Pallez, L. Teylo, and N. Vidal. IO-SETS: Simple and efficient approaches for I/O bandwidth management. working paper or preprint, May 2022.
- [5] P. Brucker, S. Knust, and C. Oğuz. Scheduling chains with identical jobs and constant delays on a single machine. *Mathematical Methods of Operations Research*, 63(1):63–75, 2006.
- [6] J. Carretero, E. Jeannot, G. Pallez, D. E. Singh, and N. Vidal. Mapping and scheduling HPC applications for optimizing I/O. In *ICS: International Conference on Supercomputing*. ACM, 2020.

- [7] M. Dorier, G. Antoniu, R. Ross, D. Kimpe, and S. Ibrahim. CALCioM: Mitigating I/O Interference in HPC Systems through Cross-Application Coordination. In *Proc. 2014 IEEE 28th International Parallel and Distributed Processing Symposium, IPDPS '14*, page 155–164. IEEE Computer Society, 2014.
- [8] K. Ecker and M. Tanaś. Complexity of scheduling of coupled tasks with chains precedence constraints and any constant length of gap. *Journal of the Operational Research Society*, 63(4):524–529, 2012.
- [9] A. Gainaru, G. Aupy, A. Benoit, F. Cappello, Y. Robert, and M. Snir. Scheduling the I/O of HPC applications under congestion. In *IPDPS'2015, the 29th IEEE International Parallel and Distributed Processing Symposium*. IEEE Computer Society Press, 2015.
- [10] T. Herault and Y. Robert, editors. *Fault-Tolerance Techniques for High-Performance Computing*, Computer Communications and Networks. Springer Verlag, 2015.
- [11] Y. Hua, X. Shi, H. Jin, W. Liu, Y. Jiang, Y. Chen, and L. He. Software-defined qos for I/O in exascale computing. *CCF Trans. High Perform. Comput.*, 1(1):49–59, 2019.
- [12] F. Isaila, J. Carretero, and R. Ross. CLARISSE: A Middleware for Data-Staging Coordination and Control on Large-Scale HPC Platforms. In *16th IEEE/ACM Int. Symp. on Cluster, Cloud and Grid Computing (CCGrid)*, pages 346–355. IEEE, 2016.
- [13] E. Jeannot, G. Pallez, and N. Vidal. Scheduling periodic I/O access with bi-colored chains: models and algorithms. *J. Scheduling*, 24(5):469–481, 2021.
- [14] Q. Kang, S. Lee, K. Hou, R. Ross, A. Agrawal, A. Choudhary, and W.-k. Liao. Improving MPI collective i/o for high volume non-contiguous requests with intra-node aggregation. *IEEE Transactions on Parallel and Distributed Systems*, 31(11):2682–2695, 2020.
- [15] A. Khan, H. Sim, S. S. Vazhkudai, A. R. Butt, and Y. Kim. An analysis of system balance and architectural trends based on top500 supercomputers. In *The International Conference on High Performance Computing in Asia-Pacific Region, HPC Asia 2021*, page 11–22, New York, NY, USA, 2021. Association for Computing Machinery.
- [16] Y. Kuo, S.-I. Chen, and Y.-H. Yeh. Single machine scheduling with sequence-dependent setup times and delayed precedence constraints. *Operational Research*, 20(2):927–942, 2020.
- [17] LANL, NERSC, SNL. APEX workflows. Technical report, Los Alamos National Laboratory (LANL), National Energy Research Scientific Computing Center (NERSC), Sandia National Laboratory (SNL)., 2016. Available online at <http://www.nersc.gov/assets/apex-workflows-v2.pdf>.
- [18] L. Lozano, M. J. Magazine, and G. G. Polak. Decision diagram-based integer programming for the paired job scheduling problem. *IIEE Transactions*, 53(6):671–684, 2021.
- [19] Message Passing Interface Forum. MPI: A message-passing interface standard, version 4.0. <https://www.mpi-forum.org/>, 2021.
- [20] A. Munier, M. Queyranne, and A. S. Schulz. Approximation bounds for a general class of precedence constrained parallel machine scheduling problems. In *Integer Programming and Combinatorial Optimization*, pages 367–382. Springer, 1998.
- [21] A. Munier and F. Sourd. Scheduling chains on a single machine with non-negative time lags. *Mathematical Methods of Operations Research*, 57(1):111–123, 2003.

-
- [22] S. Oral, S. S. Vazhkudai, F. Wang, C. Zimmer, C. Brumgard, J. Hanley, G. Markomanolis, R. Miller, D. Leverman, S. Atchley, and V. V. Larrea. End-to-End I/O Portfolio for the Summit Supercomputing Ecosystem. In *Proc. Int. Conf. for High Performance Computing, Networking, Storage and Analysis*, SC '19. ACM, 2019.
 - [23] T. Patel, R. Garg, and D. Tiwari. GIFT: A Coupon Based Throttle-and-Reward Mechanism for Fair and Efficient I/O Bandwidth Management on Parallel Storage Systems. In *Proc. 18th USENIX Conf. on File and Storage Technologies*, page 103–120. USENIX Association, 2020.
 - [24] Z. Tan, L. Du, D. Feng, and W. Zhou. Eml: An i/o scheduling algorithm in large-scale-application environments. *Future Generation Computer Systems*, 78:1091–1100, 2018.
 - [25] S. Thapaliya, P. Bangalore, J. Lofstead, K. Mohror, and A. Moody. IO-Cop: Managing concurrent accesses to shared parallel file system. In *43rd International Conference on Parallel Processing (ICPP) Workshops*, pages 52–60. IEEE, 2014.
 - [26] Y. Yang, X. Shi, W. Liu, H. Jin, Y. Hua, and Y. Jiang. DDL-QoS: a dynamic I/O scheduling strategy of QoS for HPC applications. *Concurrency and Computation: Practice and Experience*, 33(7), 2021.
 - [27] F. Zanon Boito, G. Pallez, L. Teylo, and N. Vidal. IO-SETS: Simple and efficient approaches for I/O bandwidth management. <https://hal.inria.fr/hal-03648225>, May 2022.
 - [28] B. Zha and H. Shen. Adaptively Periodic I/O Scheduling for Concurrent HPC Applications. *Electronics*, 11(9), 2022.

A Detailed description of BESTNEXTEVENT

In this appendix we detail Algorithm BESTNEXTEVENT which was sketched in Section 4.3.

The aim of algorithm BESTNEXTEVENT is to maximize the minimum yield lexicographically at the next *predictable* event, that is, either at the end of the execution window or the first time one of the current I/O requests is completed, whichever comes first. BESTNEXTEVENT is called each time an event occurs, either the completion of an I/O request, the release of a new I/O request, or the beginning of a new execution window. Then BESTNEXTEVENT defines a *constant* bandwidth allocation that will be applied up to the next event.

The algorithm itself is the combination of three algorithms. Algorithm BESTNEXTEVENT itself partitions the whole execution window in a set of what we call “simple intervals”. Searching for an event that maximizes the minimum yield is relatively easy in a “simple interval” because the peculiar events which may change the nature of the optimal solution can only happen at the extremities of a simple interval. Algorithm BESTNEXTEVENT is described in Section A.2. Once BESTNEXTEVENT has partitioned the whole execution window it searches, in each simple interval, and for each application that can define the next event in that interval, the solution maximizing the minimum yield by calling algorithm MINYIELDININTERVAL (described in Section A.1). Then, among all the solutions maximizing the minimum yield, BESTNEXTEVENT picks the one that maximizes the yield lexicographically by calling Algorithm LEXICOMINYIELD (Section A.3).

Before describing in detail algorithm MINYIELDININTERVAL we describe its working principle.

Algorithm principle. Even if algorithm MINYIELDININTERVAL looks complicated because of the many cases it has to deal with, its principle is rather simple. Let us consider a time $t + u$ and an application \mathcal{A}_k defining an event at that time. Because application \mathcal{A}_k defines an event at time $t + u$ its I/O ends precisely at that time, it uses a bandwidth $\frac{v_k}{u}$ and its yield is $y_k(t + u, \frac{v_k}{u})$. The remaining bandwidth is used to maximize the minimum yield achieved by the other applications. This minimum yield is defined either by the maximal bandwidth of an application ($\min_i y_i(t + u, b_i)$) or by the amount of remaining bandwidth. In the latter case, it turns out that all applications receiving some bandwidth achieve the same yield.

So far, we have described the algorithm principle for a given value of u . However, whatever the type of solution found (yield defined by the event-defining application, the maximal bandwidth of an application, etc.), the solution is valid in a neighbourhood of u . Therefore we study whether the minimum yield increases in this neighbourhood. Finally the neighbourhood itself is defined by the set of conditions defining the solution: which applications must receive bandwidths (which depends on the rankings and thus the intersection of the functions $\min_i y_i(t + u, 0)$), etc. In many cases we will have to determine the extent of this neighbourhood.

Preliminary remark. If the I/O requests can not saturate the available bandwidth, that is, if $\sum_{i \in S(t)} b_i \leq B$ then, obviously, each application is allotted its maximum bandwidth and there is nothing to discuss. Otherwise, there is nothing we can do to optimize the yield of a computing application and the yield of an application is increasing while it computes. Hence, in the following, we only consider applications which have posted I/O operations.

Let t denote the date of the considered event. For each application \mathcal{A}_i which has an I/O operation pending at time t (i.e., $i \in S(t)$), let \mathcal{T}_i be the (minimum) time it would have required application \mathcal{A}_i to progress as much as it did by time t if it was the sole application running on the platform. If \mathcal{A}_i is the only application running on the platform then $\mathcal{T}_i = t - \tau_i$.

A.1 Maximizing the minimum yield in an interval

Let us consider a sub-interval $[t + u_{\min}, t + u_{\max}]$ of the window and the case where application \mathcal{A}_k defines the next event. We will later loop over all applications to find the overall best solution in $[t + u_{\min}, t + u_{\max}]$. As stated in the introduction of this section, in order to simplify the study and the computations, we assume that the sub-interval $[t + u_{\min}, t + u_{\max}]$ is “simple”, that is, that it does not contain any “peculiar” events, except maybe at its extremities ($t + u_{\min}$ and/or $t + u_{\max}$). We will define these “peculiar” events and show how to compute them during the algorithm walkthrough. Algorithm MINYIELDININTERVAL (Algorithm 1) shows how to find the solution maximizing the minimum yield in the considered interval.

Algorithm MINYIELDININTERVAL starts by studying the situation at time $t + u_{\min}$. That is, we consider the case where the next-event happens at time $t + u = t + u_{\min}$. Because \mathcal{A}_k is the event-defining application, its I/O completes at time $t + u_{\min}$ and it is allocated a bandwidth $\frac{\mathcal{V}_k}{u_{\min}}$. Hence, $\alpha_k = \frac{\mathcal{V}_k}{u_{\min} b_k}$ and the remaining bandwidth will be distributed among the other applications (Step 1).

The algorithm then computes an upper bound on the achievable minimum yield. Let us consider any application \mathcal{A}_i . The maximum bandwidth allocatable to \mathcal{A}_i is b_i . Hence, the earliest time \mathcal{A}_i 's I/O request could complete is at time $t + \frac{\mathcal{V}_i}{b_i}$. This time is a peculiar time and thus we add the set

$$\left\{ \frac{\mathcal{V}_i}{b_i} \mid i \in S(t) \right\}$$

to the set of peculiar times (Step 3 of Algorithm 2).

By hypothesis, if the next event happens at time $t + u$, no I/O request can complete in the interval $(t, t + u)$. Therefore, because the yield of an application at time $t + u$ is a decreasing function of its allocated bandwidth, the maximum achievable yield for application \mathcal{A}_i at time $t + u$ is

$$y_i^{\max}(t + u) = \begin{cases} y_i(t + u, b_i) & \text{if } u \leq \frac{\mathcal{V}_i}{b_i} \\ y_i(t + u, \frac{\mathcal{V}_i}{u}) & \text{otherwise.} \end{cases} \quad (12)$$

MINYIELDININTERVAL starts by computing (at Step 2) an upper bound on the maximal minimum yield, that is, the minimum, over all applications performing I/Os, of their maximum achievable yield, and of the yield of the event-defining application \mathcal{A}_k . We then check whether there is enough bandwidth overall to achieve this upper bound (Step 3).

The yield upper-bound is not achievable. We first consider the case where there is *not* enough bandwidth overall to achieve the upper bound. We want to compute the optimal yield y^{opt} . If no bandwidth is allocated to an application \mathcal{A}_i , it achieves a yield of $y_i(t + u_{\min}, 0)$ at time $t + u_{\min}$. Therefore, the only applications to which bandwidth is allocated are those such that $y^{opt} > y_i(t + u_{\min}, 0)$. At Step 4, MINYIELDININTERVAL builds and sorts the set \mathcal{Y}_0 of the minimum yields achieved by the different applications (application \mathcal{A}_k excepted) if they are not allocated any bandwidth. Let us denote by l the index of the element of \mathcal{Y}_0 such that $\mathcal{Y}_0[l] \leq y^{opt}$ and, either $y^{opt} < \mathcal{Y}_0[l + 1]$ or $l = |\mathcal{Y}_0|$ (i.e., $\mathcal{Y}_0[l]$ is the largest element in \mathcal{Y}_0 smaller than or equal to y^{opt}). l is computed at Step 5 either in linear time through an exhaustive search or, more cleverly, in logarithmic time through a binary search. Then, bandwidth should be allocated to all applications if $l = |\mathcal{Y}_0|$ and, otherwise, bandwidth should only be allocated to applications such that $y_i(t + u_{\min}, 0) \leq \mathcal{Y}_0[l]$. The yield of each of these applications is then equal to y^{opt} . Indeed, by definition of y^{opt} the yield of an application cannot be smaller. Because we are in the case where the upper-bound on the yield is not achievable, each application (if allocated enough bandwidth) can achieve a yield strictly larger than y^{opt} at time $t + u_{\min}$ without violating its bandwidth limit b_i . Then, if the yield of one application receiving some bandwidth was strictly

Algorithm 1: MINYIELDININTERVAL($\mathcal{A}_k, u_{\min}, u_{\max}, \mathcal{S}$)

```

1  $RB \leftarrow B - \mathcal{V}_k / u_{\min}$  /* Remaining bandwidth */
2  $y^{UB} \leftarrow \min \left\{ \min_{i \in \mathcal{S}} y_i^{\max}(t + u_{\min}), y_k \left( t + u_{\min}, \frac{\mathcal{V}_k}{u_{\min}} \right) \right\}$ 
3 if  $\sum_{i \in \mathcal{S}} \mathcal{B}\mathcal{W}_i(t + u_{\min}, y^{UB}) > RB$  then
4    $\mathcal{Y}_0 \leftarrow \text{sort}(\{y_i(u_{\min}, 0) \mid i \in \mathcal{S}\})$ 
5   Find  $l$  such that  $\sum_{i \in \mathcal{S}} \mathcal{B}\mathcal{W}_i(u_{\min}, \mathcal{Y}_0[l]) \leq RB$  and  $l = |\mathcal{Y}_0|$  or
      $\sum_{i \in \mathcal{S}} \mathcal{B}\mathcal{W}_i(u_{\min}, \mathcal{Y}_0[l+1]) > RB$ 
6    $\mathcal{I} \leftarrow \{i \in \mathcal{S} \mid y_i(u_{\min}, 0) \leq \mathcal{Y}_0[l]\}$ 
7    $y_{\mathcal{I}}^{\text{opt}}(u) = \frac{uB - \mathcal{V}_k + \sum_{i \in \mathcal{I}} b_i \tau_i}{\sum_{i \in \mathcal{I}} b_i (t + u - \tau_i)}$ 
8   if  $y_{\mathcal{I}}^{\text{opt}}(u)$  is a non-increasing function then return  $(y_{\mathcal{I}}^{\text{opt}}(u_{\min}), u_{\min})$ 
9   If  $l < \mathcal{Y}_0$  let  $\mathcal{A}_j$  be such that  $y_j(u_{\min}, 0) = \mathcal{Y}_0[l+1]$  and  $y_j(u_{\max}, 0)$  is minimal
10  Let  $u_{\text{intersect}}$  be the first intersection (if it exists) in  $[u_{\min}, u_{\max}]$  of  $y_{\mathcal{I}}^{\text{opt}}(u)$  with
     either  $y_k(t + u, \frac{\mathcal{V}_k}{u})$  or  $y_j(t + u, 0)$  (if the later exists) or with  $y_m(t + u, b_m)$  for
     some  $m \in \mathcal{I}$ 
11  if  $u_{\text{intersect}}$  does not exist then return  $(y_{\mathcal{I}}^{\text{opt}}(u_{\max}), u_{\max})$ 
12  else if  $y_{\mathcal{I}}^{\text{opt}}(u_{\text{intersect}}) = y_k \left( t + u_{\text{intersect}}, \frac{\mathcal{V}_k}{u_{\text{intersect}}} \right)$  then return
      $\left( y_k \left( t + u_{\text{intersect}}, \frac{\mathcal{V}_k}{u_{\text{intersect}}} \right), u_{\text{intersect}} \right)$ 
13  else return MINYIELDININTERVAL( $\mathcal{A}_k, u_{\text{intersect}}, u_{\max}, \mathcal{S}$ )
14 if  $y^{UB} = y_k \left( t + u_{\min}, \frac{\mathcal{V}_k}{u_{\min}} \right)$  then
15   return  $\left( y_k \left( t + u_{\min}, \frac{\mathcal{V}_k}{u_{\min}} \right), u_{\min} \right)$ 
16 else
17   Let  $\mathcal{A}_j$  be such that  $y_j^{\max}(t + u_{\min}) = y^{UB}$  and  $y_j^{\max}(t + u_{\max})$  is minimal
18   if  $\mathcal{B}\mathcal{W}_j(u_{\min}, y_j^{\max}(t + u_{\min})) = \frac{\mathcal{V}_j}{u_{\min}}$  then return  $(y_j^{\max}(t + u_{\min}), u_{\min})$ 
19    $\mathcal{I} \leftarrow \{i \in \mathcal{S} \mid y_i(u_{\min}, 0) \leq y^{UB}\}$ 
20    $y_{\mathcal{I}}^{\text{opt}}(u) = \frac{uB - \mathcal{V}_k + \sum_{i \in \mathcal{I}} b_i \tau_i}{\sum_{i \in \mathcal{I}} b_i (t + u - \tau_i)}$ 
21   If  $\mathcal{I} \neq \mathcal{S}$  then let  $\mathcal{A}_l$  be an application in  $\mathcal{S} \setminus \mathcal{I}$  such that  $y_l \left( \frac{u_{\min} + u_{\max}}{2}, 0 \right)$  is minimal
22   Let  $u_{\text{intersect}}$  be the first intersection (if it exists) in  $[u_{\min}, u_{\max}]$  of  $y_j^{\max}(t + u)$  with
      $y_{\mathcal{I}}^{\text{opt}}(u)$ ,  $y_k \left( t + u, \frac{\mathcal{V}_k}{u} \right)$ , or  $y_l(t + u, 0)$ 
23   if  $u_{\text{intersect}}$  does not exist then return  $y_j^{\max}(t + u_{\max})$ 
24   if  $y_j^{\max}(t + u_{\text{intersect}}) = y_k \left( t + u_{\text{intersect}}, \frac{\mathcal{V}_k}{u_{\text{intersect}}} \right)$  then return
      $\left( y_k \left( t + u_{\text{intersect}}, \frac{\mathcal{V}_k}{u_{\text{intersect}}} \right), u_{\text{intersect}} \right)$ 
25   else return MINYIELDININTERVAL( $\mathcal{A}_k, u_{\text{intersect}}, u_{\max}, \mathcal{S}$ )

```

larger than y^{opt} , some of its bandwidth could be distributed to the other applications to increase the value of y^{opt} which would contradict the optimality of y^{opt} .

Finally, the total bandwidth should be distributed among the applications. Hence, the constraints on the distribution of bandwidth are:

$$\forall i \in \mathcal{I}, y_{\mathcal{I}}^{opt}(u) = \frac{\mathcal{T}_i + \alpha_i u}{t + u - \tau_i} \quad \text{and} \quad \sum_{i \in \mathcal{I}} \alpha_i b_i = B - \frac{\mathcal{V}_k}{u}$$

where \mathcal{I} is the set of the applications performing I/Os which require bandwidth, as defined at Step 6. We wrote the constraints for an undefined variable u rather than just for the value u_{\min} for which we know there are true. This is because we are next going to study the evolution of the defined solution in a neighbourhood of $t + u_{\min}$.

From the first equation we obtain:

$$(t + u - \tau_i) y_{\mathcal{I}}^{opt}(u) = \mathcal{T}_i + \alpha_i u \quad \Leftrightarrow \quad \frac{1}{u} ((t + u - \tau_i) y_{\mathcal{I}}^{opt}(u) - \mathcal{T}_i) = \alpha_i.$$

Then the second equation can be rewritten:

$$\sum_{i \in \mathcal{I}} \frac{b_i}{u} ((t + u - \tau_i) y_{\mathcal{I}}^{opt}(u) - \mathcal{T}_i) = RB \quad \Leftrightarrow \quad y_{\mathcal{I}}^{opt}(u) = \frac{uRB + \sum_{i \in \mathcal{I}} b_i \mathcal{T}_i}{\sum_{i \in \mathcal{I}} b_i (t + u - \tau_i)}$$

Hence:

$$y_{\mathcal{I}}^{opt}(u) = \frac{uB - \mathcal{V}_k + \sum_{i \in \mathcal{I}} b_i \mathcal{T}_i}{\sum_{i \in \mathcal{I}} b_i (t + u - \tau_i)}. \quad (13)$$

We then study the variations of $y_{\mathcal{I}}^{opt}(u)$:

$$\begin{aligned} y_{\mathcal{I}}^{opt}(u) &= \frac{uB - \mathcal{V}_j + \sum_{i \in \mathcal{I}} b_i \mathcal{T}_i}{\sum_{i \in \mathcal{I}} b_i (t + u - \tau_i)} \\ &= \frac{B}{\sum_{i \in \mathcal{I}} b_i} + \frac{1}{\sum_{i \in \mathcal{I}} b_i} \frac{(\sum_{i \in \mathcal{I}} b_i) (\sum_{i \in \mathcal{I}} b_i \mathcal{T}_i) - (\sum_{i \in \mathcal{I}} b_i (B(t - \tau_i) + \mathcal{V}_j))}{\sum_{i \in \mathcal{I}} b_i (t + u - \tau_i)}. \end{aligned}$$

Therefore, the yield is increasing with u on the considered interval if and only if the expression

$$\left(\sum_{i \in \mathcal{I}} b_i \right) \left(\sum_{i \in \mathcal{I}} b_i \mathcal{T}_i \right) - \left(\sum_{i \in \mathcal{I}} b_i (B(t - \tau_i) + \mathcal{V}_j) \right) \quad (14)$$

is negative. Hence, if this expression is non-negative, the yield is non-increasing over the interval and the optimum is found for $u = u_{\min}$ (Step 8).

If the yield is increasing, things are more complicated. Equation (13) defines the optimum yield for a given set \mathcal{I} of applications among which the remaining bandwidth should be distributed. In turn, the definition of \mathcal{I} depends on the considered time $t + u$ and on the amount of remaining bandwidth which is also a function of u . Let us assume that the interval $(t + u_{\min}, t + u_{\max})$ is such that no two curves $u \mapsto y_i(t + u, 0)$ intersects in this interval (except if the two curves are identical). Therefore, the intersection of two curves $u \mapsto y_i(t + u, 0)$ defines a ‘‘peculiar event’’ and the set of these intersections is computed at Step 9 of Algorithm 2. Then, if $l = |\mathcal{Y}_0|$, there is no curve $u \mapsto y_i(t + u, 0)$ above the curve $u \mapsto y_{\mathcal{I}}^{opt}(u)$. Otherwise, let application \mathcal{A}_j be an application whose curve is the first one above the curve $u \mapsto y_{\mathcal{I}}^{opt}(u)$. In practice, let application \mathcal{A}_j be an application such that $y_j(t + u_{\min}, 0) = \mathcal{Y}_0[l + 1]$ and such that $y_j(t + u_{\max}, 0)$ is minimal (Step 9). When the yield is increasing, it may intersects the

curve $y_j(t+u, 0)$, changing the definition of the set \mathcal{I} of applications requiring bandwidth. If this happens at a date $t + u_{\text{intersect}}$, then we know that the maximal min yield on the interval $[u_{\text{min}}, u_{\text{intersect}}]$ is achieved for $u_{\text{intersect}}$. Then, we recursively call Algorithm MINYIELDININTERVAL on the interval $[u_{\text{intersect}}, u_{\text{max}}]$ (Step 13).

Another potential problem when the yield is increasing is that the (minimum) yield of the applications receiving bandwidth, $y_{\mathcal{I}}^{\text{opt}}(u)$, becomes equal to the yield of the event-defining applications, \mathcal{A}_k , at some date $t + u_{\text{intersect}}$. Then, we know that the maximal min yield on the interval $[u_{\text{min}}, u_{\text{intersect}}]$ is achieved for $u_{\text{intersect}}$. Furthermore, because the function $y_k(t+u, \frac{v_k}{u})$ is non-increasing, the minimum yield is non increasing on the interval $[u_{\text{intersect}}, u_{\text{max}}]$. Hence, the maximum minimum yield is obtained at time $t + u_{\text{intersect}}$ (Step 12).

The last problem that may happen, when u is increasing, is that the bandwidth allocated to an application \mathcal{A}_m may increase so much that it reaches its limit b_m , at some date $t + u_{\text{intersect}}$. At that point, the nature of the solution changes. We know that the maximal min yield on the interval $[u_{\text{min}}, u_{\text{intersect}}]$ is achieved for $u_{\text{intersect}}$ and, once again, we recursively call Algorithm MINYIELDININTERVAL on the interval $[u_{\text{intersect}}, u_{\text{max}}]$ (Step 13).

All the “problems” we just highlighted are defined by the equality of two yield functions. All these yield functions are of the form $\frac{\alpha u + \beta}{\gamma u + \delta}$. Therefore, looking for the equality of two such functions requires to solve a second degree polynomial and to see whether it has roots in the interval $[u_{\text{min}}, u_{\text{max}}]$ and, if there are two of them, which one is the smallest. Therefore, the algorithm can easily check, at Step 10, whether any of these potential problems occurs and if this is the case, which one happens the earliest. If no problem occurs, $y_{\mathcal{I}}^{\text{opt}}(u)$ defines a valid, increasing, solution throughout the interval and the best solution is found at time $t + u_{\text{max}}$ (Step 11). Otherwise, the algorithm applies the relevant case, as defined above.

The yield upper-bound is achievable. We now consider the case where there is enough bandwidth overall to achieve the upper bound. We first check whether the application defining the upper-bound is also the event-defining application \mathcal{A}_k (Step 14). If this is the case, because the yield of the event-defining application is a non-increasing function, $y_k(t+u, \frac{v_k}{u})$, the optimal solution on the interval $[t + u_{\text{min}}, t + u_{\text{max}}]$ is obtained in u_{min} (Step 15).

Otherwise, we identify (Step 17) an application, say \mathcal{A}_j , which defines the upper bound on the minimum yield at time $t + u_{\text{min}}$. If there are several candidates we pick one which also defines this bound at the end of the interval; this is possible because we assume the interval includes no *peculiar* events. Therefore we add to the set of peculiar events the times at which two curves $u \mapsto y_i^{\text{max}}(t+u)$ intersect (Steps 5, 6, and 7 of Algorithm 2). If there are still several candidates we pick one arbitrarily as they all have the same maximum yield throughout the interval.

The yield achieved by application \mathcal{A}_j , namely $y_j^{\text{max}}(t+u)$, is an increasing function of u . This yield defines the optimal minimum yield as long as two conditions hold: 1) it is not greater than the yield $y_k(t+u, \frac{v_k}{u})$ of application \mathcal{A}_k ; 2) the bandwidth required for all applications to achieve a yield at least equal to $y_j^{\text{max}}(t+u)$ does not exceed the total available bandwidth. Algorithm MINYIELDININTERVAL first identifies (at Step 19) to which applications bandwidth should be allocated for all applications to have a yield of at least $y_j^{\text{max}}(t+u)$. Then it computes the maximal yield $y_{\mathcal{I}}^{\text{opt}}(u)$ achieved when the total remaining bandwidth is distributed among these applications. Note that this equation is only meaningful when one must distribute at least the total remaining bandwidth to achieve a yield at least equal to $y_j^{\text{max}}(t+u)$. Finally, the algorithm computes the latest time, in the interval $[t + u_{\text{min}}, t + u_{\text{max}}]$, at which the two above conditions hold (Step 22): it computes the earliest time $t + u_{\text{intersect}}$ in $[t + u_{\text{min}}, t + u_{\text{max}}]$ (if it exists) when the curves $y_j^{\text{max}}(t+u)$ and $y_k(t+u, \frac{v_k}{u})$ intersect or when the curves $y_j^{\text{max}}(t+u)$

and $y_{\mathcal{I}}^{opt}(u)$ intersect.³

If $u_{\text{intersect}}$ does not exist, the optimal solution is $y_j^{\max}(t + u_{\max})$ (Step 23). Otherwise, the best solution in $[t + u_{\min}, t + u_{\text{intersect}}]$ is reached at time $t + u_{\text{intersect}}$ at which time the nature of the solution changes. If $y_j^{\max}(t + u_{\text{intersect}}) = y_k\left(t + u_{\text{intersect}}, \frac{v_k}{u_{\text{intersect}}}\right)$ then, like previously, the best solution in $[t + u_{\min}, t + u_{\max}]$ is reached in $t + u_{\text{intersect}}$ (Step 24). Otherwise, Algorithm `MINYIELDININTERVAL` is recursively called on the interval $[u_{\text{intersect}}, u_{\max}]$ (Step 25).

A.2 Maximizing the minimum yield in the whole execution window

`BESTNEXTEVENT` starts by partitioning the whole (remaining) execution window $[t, T_{\text{end}}]$ based on the peculiar events identified in the previous section. This is done at Step 1 and at Steps 5 through 10.

We compute at Step 1 the earliest time each I/O can complete. If none of these dates happens before the end of the window, the next event happens at the end of the window (Step 2). If there is enough total bandwidth for all I/Os to be allocated their maximum bandwidth, the next event happens the first time an I/O can complete (Step 3). Otherwise, `BESTNEXTEVENT` scans one by one the intervals defined by the peculiar events. For each of these intervals, and for each application whose I/O can complete in the studied interval, it calls `MINYIELDININTERVAL` to find the best possible solution. The best solution, among all the identified candidates, is eventually selected at Step 20.

The test at Step 18 enables us to implement an optimization: once a candidate solution is identified where the event-defining application is also the application with minimum yield then, because the yield $y_k\left(t, \frac{v_k}{t}\right)$ is a decreasing function, and it is the highest yield application \mathcal{A}_k can achieve at time t , we know that no better solution can be found for larger values of t and the search can stop.

A.3 Maximizing the minimum yield lexicographically at time $t + u$

`LEXICOMINYIELD` (see Algorithm 3) is an algorithm which builds a bandwidth allocation that lexicographically maximizes the minimum yield at the date $t + u$ (if no event occurs in the interval $(t, t + u)$) and an event defined by application \mathcal{A}_k happens at time $t + u$. The algorithm principle is straightforward. It starts, if an event-defining application \mathcal{A}_k is designed, to allocate it the required bandwidth (Steps 3 through 5). Then, `LEXICOMINYIELD` calls, at Step 10, `MINYIELD` to obtain a bandwidth allocation maximizing the minimum yield. By definition, the yield of an optimal solution cannot be increased. The yield of a solution cannot be increased because achieving such a yield requires either to saturate the total available bandwidth, or to allocate one application its maximum allocatable bandwidth, or for an I/O request to end at time $t + u$. In the first case, the yield of each application performing an I/O is equal to the maximum minimum yield and an optimal solution has been built (Steps 12 and 14). For the other cases (which are not excluding), we fix the bandwidth of the applications whose I/O ends at time $t + u$ and of those whose allocated bandwidth is equal to the maximum allocatable one (Step 18). We compute the total remaining bandwidth (Step 19) and the minimum yield maximization is redone on the remaining applications (defined at Step 20) with the remaining bandwidth.

Algorithm `LEXICOMINYIELD` has complexity of $O(m^2 \log(m))$, because of the complexity of `MINYIELD` and because the while loop is executed at most m times.

³It may happen (especially after a recursive call), that $y_j^{\max}(t + u_{\min}) = y_{\mathcal{I}}^{opt}(u_{\min})$. In such a case, the algorithm determines which of the two functions $y_j^{\max}(t + u)$ and $y_{\mathcal{I}}^{opt}(u)$ achieves the smallest yield on $[t + u, t + u + \epsilon]$, that is, when u is infinitesimally greater than u_{\min} . If it is $y_j^{\max}(t + u)$, it is used instead of $y_{\mathcal{I}}^{opt}(u)$ to compute $u_{\text{intersect}}$.

Algorithm 2: BESTNEXTEVENT(t, T_{end})

```

1  $\mathcal{S} \leftarrow \left\{ \frac{v_i}{b_i} \mid i \in S(t) \right\} \cap [0, T_{end} - t]$ 
2 if  $\mathcal{S} = \emptyset$  then return LEXICOMINYIELD( $t, T_{end}, \emptyset$ )
3 if  $\sum_{i \in S(t)} b_i \leq B$  then return LEXICOMINYIELD( $t, \min \mathcal{S}, \emptyset$ )
4  $\mathcal{D} \leftarrow \{(h, \emptyset)\}$ 
5 JointEvents  $\leftarrow$ 
    $\left\{ u \mid \exists i, j \in S(t), T_{end} - t \geq u \geq 0, u \geq \frac{v_i}{b_i}, u \geq \frac{v_j}{b_j}, y_i\left(u, \frac{v_i}{u}\right) = y_j\left(u, \frac{v_j}{u}\right), \mathcal{T}_i + \frac{v_i}{b_i} \neq \mathcal{T}_j + \frac{v_j}{b_j} \right\}$ 
6  $\mathcal{M}_1 \leftarrow \left\{ u \mid 0 \leq u \leq T_{end} - t, i, j \in S(t), i \neq j, y_i\left(u, b_i\right) = y_j\left(u, \frac{v_j}{u}\right), u \leq \frac{v_i}{b_i}, u \geq \frac{v_j}{b_j} \right\}$ 
7  $\mathcal{M}_2 \leftarrow \left\{ u \mid 0 \leq u \leq T_{end} - t, i, j \in S(t), i \neq j, y_i\left(u, b_i\right) = y_i\left(u, b_j\right), u \leq \frac{v_i}{b_i}, u \leq \frac{v_j}{b_j} \right\}$ 
8  $\mathcal{X} \leftarrow$ 
    $\left\{ u \mid 0 \leq u \leq T_{end} - t, i, j \in S(t), i \neq j, T_{end} - t \geq u, u \geq \frac{v_j}{b_j}, y_i(u, 0) = y_j\left(u, \frac{v_j}{u}\right) \right\}$ 
9  $\mathcal{Z} \leftarrow \{u \mid 0 \leq u \leq T_{end} - t, i, j \in S(t), i \neq j, 0 \leq u \leq T_{end} - t, y_i(u, 0) = y_j(u, 0)\}$ 
10  $\mathcal{U} \leftarrow \text{sort}(\text{JointEvents} \cup \mathcal{S} \cup \mathcal{X} \cup \mathcal{M}_1 \cup \mathcal{M}_2 \cup \{T_{end} - t\})$ 
11 DominantSolutionNotFound  $\leftarrow$  true
12  $i \leftarrow 1$ 
13 while  $i \leq |\mathcal{U}| - 1$  and DominantSolutionNotFound do
14   for  $k \in S(t)$  do
15     if  $u_i \geq \frac{v_k}{b_k}$  then
16        $(y, u) \leftarrow \text{MINYIELDININTERVAL}(k, u_i, u_{i+1}, S(t) \setminus \{k\})$ 
17        $\mathcal{D} \leftarrow \mathcal{D} \cup \{u, k\}$ 
18       if  $y = y_k\left(u, \frac{v_k}{u}\right)$  then DominantSolutionNotFound  $\leftarrow$  false
19    $i \leftarrow i + 1$ 
20 return  $\max_{(u,k) \in \mathcal{D}} \text{LEXICOMINYIELD}(t, u, k)$ 

```

Algorithm MINYIELD is pretty straightforward. It finds the minimum yield that can be achieved among the set of applications whose bandwidth has not already been fixed, knowing the amount of bandwidth that remains to be allocated.

Algorithm 3: LEXICOMINYIELD(t, u, k)

```

1  $\forall i \in S(t) \ b_i \leftarrow 0$  /* Current bandwidth allocation */
2 if  $k \neq \emptyset$  then
3    $b_k \leftarrow \frac{\mathcal{V}_k}{u}$ 
4    $RB \leftarrow B - b_k$  /* Remaining bandwidth */
5    $\mathcal{F} \leftarrow S(t) \setminus \{k\}$  /* Applications whose bandwidth allotment is not yet fixed */
6 else
7    $RB \leftarrow B$  /* Remaining bandwidth */
8    $\mathcal{F} \leftarrow S(t)$  /* Applications whose bandwidth allotment is not yet fixed */
9 while  $\mathcal{F} \neq \emptyset$  and  $RB > 0$  do
10   $y \leftarrow \text{MINYIELD}(t, u, RB, \mathcal{F})$ 
11  if  $\sum_{i \in \mathcal{F}} \mathcal{BW}_i(u, y) = RB$  then /* The whole remaining bandwidth is used */
12     $\forall i \in \mathcal{F}, \ b_i \leftarrow \mathcal{BW}_i(u, y)$ 
13     $\mathcal{F} \leftarrow \emptyset$ 
14     $RB \leftarrow 0$ 
15  else
16    for  $i \in \mathcal{F}$  do
17      if  $\mathcal{BW}_i(u, y) = b_i$  or  $\mathcal{BW}_i(u, y)u = \mathcal{V}_i$  then
18        /* The application bandwidth cannot be increased */
19         $b_i \leftarrow \mathcal{BW}_i(u, y)$ 
20         $RB \leftarrow RB - b_i$ 
21         $\mathcal{F} \leftarrow \mathcal{F} \setminus \{i\}$ 
21 if  $RB > 0$  then Error

```

Algorithm 4: MINYIELD(t, u, B, \mathcal{S})

- 1 $y^{\max} \leftarrow \min_{i \in \mathcal{S}} y_i^{\max}(u)$
 - 2 **if** $\sum_{i \in \mathcal{S}} \mathcal{BW}_i(u, y^{\max}) \leq B$ **then return** y^{\max}
 - 3 $\mathcal{Y}_0 \leftarrow \text{sort}(\{y_i(u, 0) \mid i \in \mathcal{S}\})$
 - 4 Find l such that $\sum_{i \in \mathcal{S}} \mathcal{BW}_i(u_{\min}, \mathcal{Y}_0[l]) \leq B$ and $l = |\mathcal{Y}_0|$ or
 $\sum_{i \in \mathcal{S}} \mathcal{BW}_i(u_{\min}, \mathcal{Y}_0[l+1]) > RB$
 - 5 $\mathcal{I} \leftarrow \{i \mid i \in \mathcal{S}, y_i(u, 0) \leq \mathcal{Y}_0[l]\}$
 - 6 **return** $\frac{uB + \sum_{i \in \mathcal{I}} \mathcal{T}_i b_i}{\sum_{i \in \mathcal{I}} (t + u - \tau_i) b_i}$
-

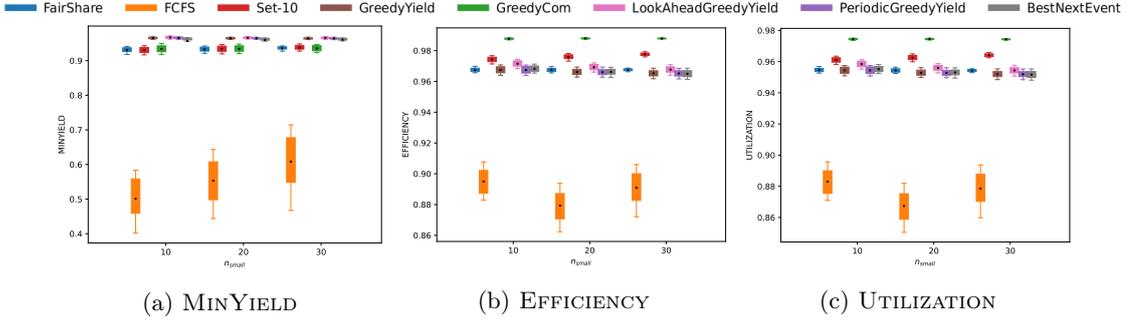


Figure 14: Impact of the number of small tasks (n_{small}) with low I/O pressure ($W^{GOAL} = 0.8$).

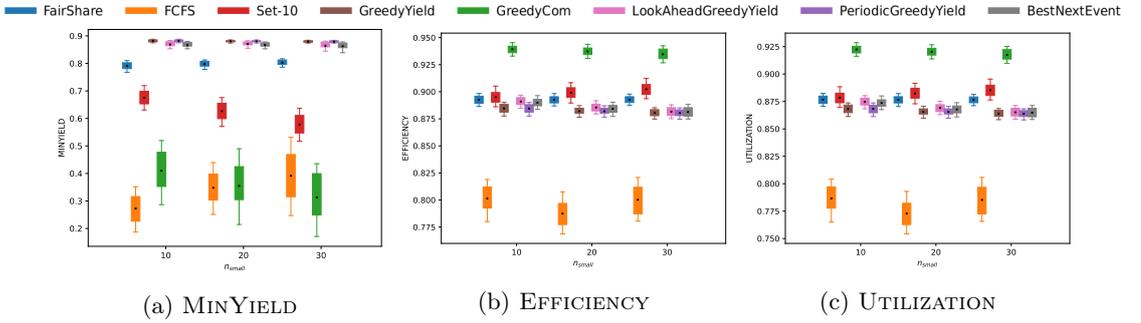


Figure 15: Impact of the number of small tasks (n_{small}) with high I/O pressure ($W^{GOAL} = 1.1$).

B Additional Simulation Results

B.1 Synthetic Traces

Impact of the number of small applications n_{small} . We once again let $\nu = \sigma = 0.5$, and run experiments for all values of $n_{small} \in [10, 20, 30]$, both in the low I/O pressure scenario (Figure 14, with $W^{GOAL} = 0.8$) and in the high I/O pressure scenario (Figure 15, with $W^{GOAL} = 1.1$). The results are pretty similar for all values of n_{small} , and we draw the same conclusions as in Section 6.2.2), in particular when the I/O pressure is high. Indeed, with a low I/O pressure, almost all heuristics succeed to achieve a high value of MINYIELD.

Impact of the standard deviation σ . We set $\nu = 0.5$, $n_{small} = 20$, and consider $\sigma \in [0, 0.25, 0.5, 0.75, 1]$, both in the low I/O pressure scenario (Figure 16, with $W^{GOAL} = 0.8$) and in the high I/O pressure scenario (Figure 17, with $W^{GOAL} = 1.1$). We observe that the novel heuristics (except maybe GREEDYCOM when the I/O pressure is high) are not affected by an increase in the standard deviation σ , while FCFS suffers from high standard deviation. In the scenario with a high I/O pressure, the MINYIELD achieved by SET-10 significantly decreases when σ increases.

Impact of the noise. In these experiments, we set $\sigma = 0.5$, $n_{small} = 20$, and we consider values $\nu \in [0, 0.25, 0.5, 0.75, 1]$, both in the low I/O pressure scenario (Figure 18, with $W^{GOAL} = 0.8$) and in the high I/O pressure scenario (Figure 19, with $W^{GOAL} = 1.1$). We see that the

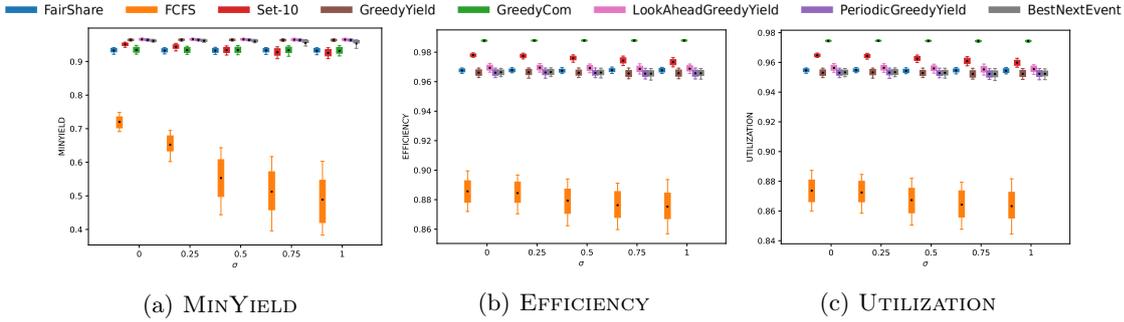


Figure 16: Impact of the standard deviation (σ) with low I/O pressure ($W^{GOAL} = 0.8$).

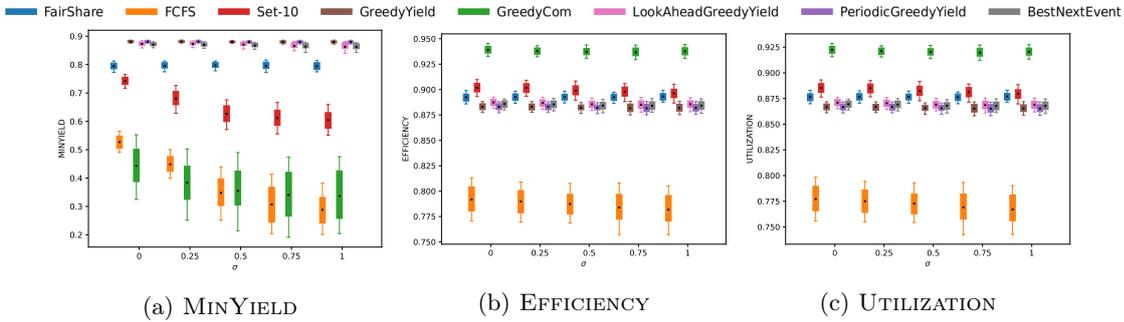


Figure 17: Impact of the standard deviation (σ) with high I/O pressure ($W^{GOAL} = 1.1$).

noise does not affect any of the heuristics, which are resilient to variations in the lengths of the working and I/O phases.

B.2 APEX Traces

To understand the results obtained on the APEX traces, we look at the characteristics of the TRILAB and NERSC workloads on the different projected machines. Figure 20 presents the distribution of the window duration for both workloads, and for 11 projected machines (Celio, and M_1 to M_{10}). We observe that window duration is an order of magnitude longer for the TRILAB workload compared to the NERSC workload, for all machines. This is easily explained by the difference of jobs between the workloads: the NERSC workload features a large number of short-lived applications, it is thus hard to find long time windows during which no application completes or starts, compared to the TRILAB workload. We also observe that, as the platform becomes larger, the average window duration increases for both workloads. Problem size, and thus amount of I/O, is defined as a function of the aggregated memory in the APEX report. When we scale up the number of nodes to project a future machine, we maintain a memory of 1GB per core, but as the number of cores per node increases, the amount of memory per application increases. This results in longer-running applications, and thus in longer windows.

These observations are corroborated by Figure 21. This figure observes how many applications belong to a given window, for the different machines. The TRILAB workload presents windows that have an order of magnitude fewer applications than the NERSC workload, which corresponds well to the relative workload distribution of the different workloads.

Figure 22 presents the distribution of platform usage during a simulation window for both

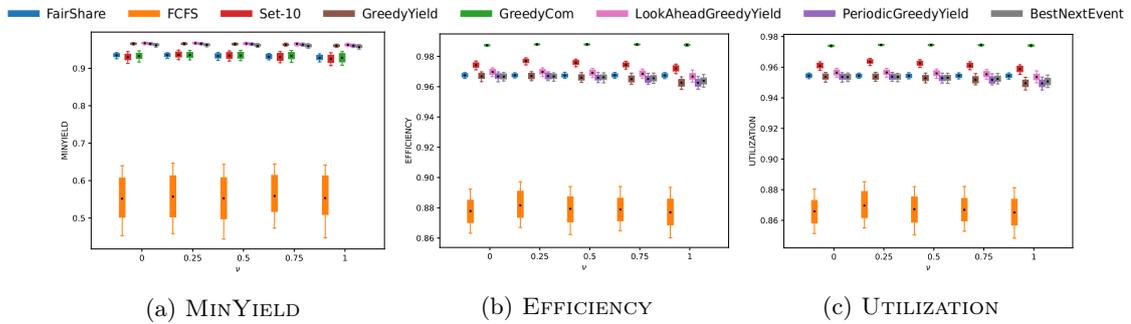
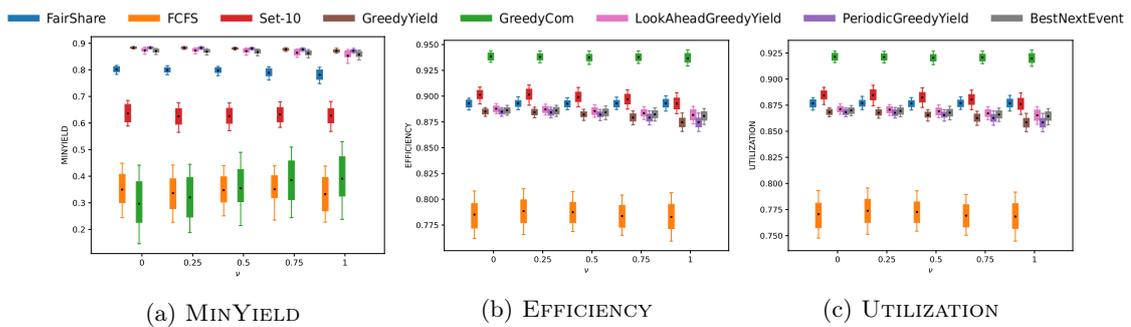
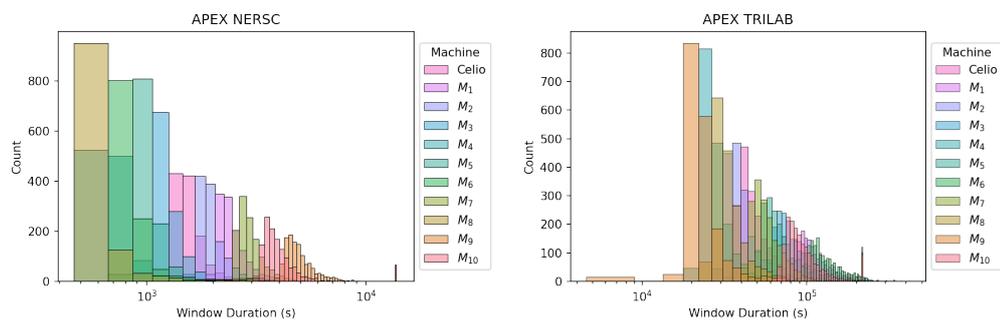
Figure 18: Impact of the noise (ν) with low I/O pressure ($W^{GOAL} = 0.8$).Figure 19: Impact of the noise (ν) with high I/O pressure ($W^{GOAL} = 1.1$).

Figure 20: Distribution of the window durations for the APEX campaign on the NERSC and TRILAB workloads, as function of the projected machine.

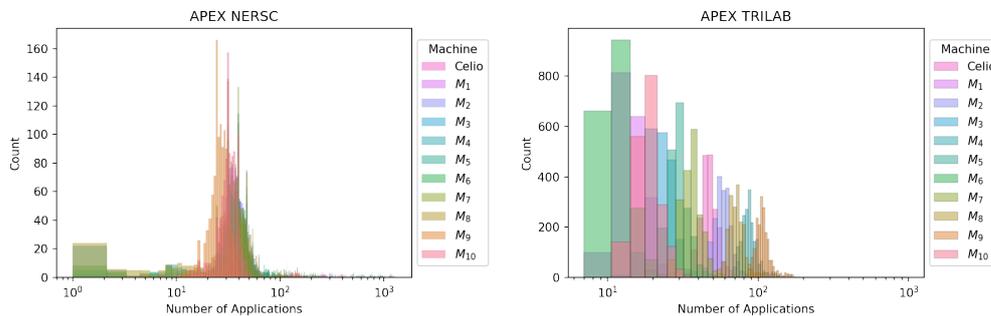


Figure 21: Distribution of the number of applications within each window for the APEX campaign on the NERSC and TRILAB workloads, as function of the projected machine.

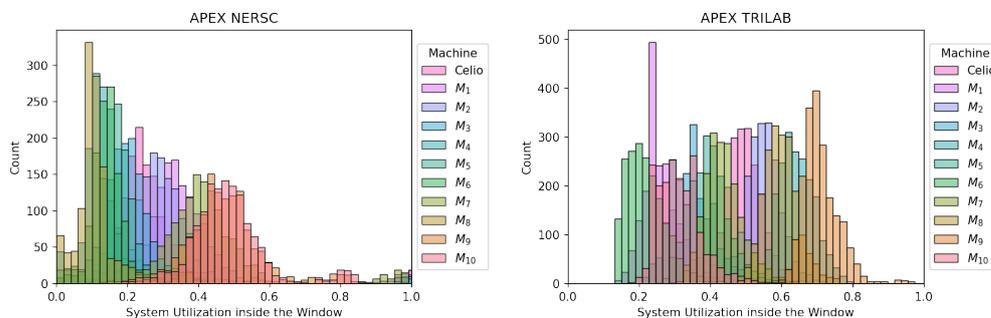


Figure 22: Distribution of the number of applications within each window for the APEX campaign on the NERSC and TRILAB workloads, as function of the projected machine.

workloads. This metric varies significantly, but is clearly higher for the TRILAB workload on average than what is observed for the NERSC platform. Again, the small duration of some applications in the NERSC workload introduces this imbalance: selecting long windows in order to provide a statistically accurate result favors windows that have only a few occurrences of the short-lived application, which increases the probability that some nodes are left idle by the first-fit scheduler. This low utilization translates in a lower I/O pressure, which explains why the NERSC workload shows less points in the high I/O pressure end of the figures.

Inria

**RESEARCH CENTRE
GRENOBLE – RHÔNE-ALPES**

Inovallée
655 avenue de l'Europe Montbonnot
38334 Saint Ismier Cedex

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399