

Fair Coflow Scheduling via Controlled Slowdown

Francesco De Pellegrini*, Vaibhav Kumar Gupta[†], Rachid El Azouzi*,
Serigne Gueye*, Cedric Richier* and Jeremie Leguay*

Abstract—The average coflow completion time (CCT) is the standard performance metric in coflow scheduling. However, standard CCT minimization may introduce unfairness between the data transfer phase of different computing jobs. Thus, while progress guarantees have been introduced in the literature to mitigate this fairness issue, the trade-off between fairness and efficiency of data transfer is hard to control.

This paper introduces a fairness framework for coflow scheduling based on the concept of slowdown, i.e., the performance loss of a coflow compared to isolation. By controlling the slowdown it is possible to enforce a target coflow progress while minimizing the average CCT. In the proposed framework, the minimum slowdown for a batch of coflows can be determined in polynomial time. By showing the equivalence with Gaussian elimination, slowdown constraints are introduced into primal-dual iterations of the CoFair algorithm. The algorithm extends the class of the σ -order schedulers to solve the fair coflow scheduling problem in polynomial time. It provides a 4-approximation of the average CCT w.r.t. an optimal scheduler. Extensive numerical results demonstrate that this approach can trade off average CCT for slowdown more efficiently than existing state of the art schedulers.

Index Terms—data transfer, coflow scheduling, fairness, progress, primal-dual scheduler.

I. INTRODUCTION

The coflow abstraction has been introduced in the seminal paper [1]. A coflow denotes a group of flows produced by data-intensive computing frameworks, e.g., Map Reduce, Giraph and Spark [2]–[4], during their data transfer phases. In the same fashion, the workflows of most data-intensive machine learning applications are based on data-transfers and operated by distributing training tasks over hundreds of individual compute nodes. A well-studied example of data transfer step is the shuffle phase of Hadoop MapReduce. The shuffle phase acts as a synchronization barrier since the next phase of the computation cannot be pursued until the end of the shuffle data transfer. Once each mapper node finishes running their job, partial results are fetched in the form of HTTP connections by peer reducer nodes and final results can be generated for the next computation phase or storage. At the application level, the shuffle phase has been shown to represent a significant part of the total computation time [5]. For this reason, efficient coflow scheduling has become a key aspect of traffic engineering in modern datacenters [6]–[10].

The standard metric to measure coflow scheduling performance is the weighted Coflow Completion Time (CCT), i.e., the average time by which the last flow of a coflow

is completed. Minimizing the average CCT is indeed the appropriate goal in order to increase the rate of computing jobs dispatched in a datacenter and improve the execution time of applications. The corresponding minimization problem is complicated by the fact that a shared datacenter fabric can be contended by hundreds of coflows at the same time. Thus, multiple congested links may appear under concurrent demands. In the last ten years, the problem of average CCT minimization has been addressed by several authors [1], [5], [8], [11], shading light on its complexity and devising several algorithmic solutions. The problem is found *NP*-hard by reduction to the open-shop problem [12], a mainstream operation research problem where jobs are scheduled on multiple machines. While inapproximability below a factor 2 has been proved [7], to date the best deterministic approximation ratio is 4 in the case when coflows are released at same time [8], [11], [13].

In the literature, maximizing network performance is known to entail potential unfairness among different flows [14]. Thus, the notion of per-flow fairness has been studied to balance the resource allocation, i.e., link utilization, among different flows. Max-min fairness and proportional fairness are reference concepts in this context and a series of fundamental works on utility-based fairness have showed that both the notions can be compounded under the larger concept of α -fairness [15].

Coflow fairness. Minimizing the average CCT suffers a similar fairness issue due to starvation [16] and fairness has been studied in the context of coflow scheduling as well [17] [18] [19] [20]. To obtain a resource allocation able to achieve a required trade-off between average CCT and fairness, it is tempting to directly use CCT as the variable in conventional fairness utility functions. However, a major obstacle to this approach exist. In fact, the definition of CCT entails a minimization, so that using CCT as argument of standard smooth convex functions such as the α -fair utility leads to minimization problems neither convex nor differentiable. The workaround is to allow for a constant rate allocation [21]. But, CCT minimization is a finite horizon problem for which optimizing stationary rates leads to suboptimal resource allocation.

Fairness in coflow scheduling relies on the notion of *coflow progress*, i.e., by guaranteeing a minimum resource allocation to coflows [18]–[20]. In particular, [19] studied the trade-off between performance and fairness by showing that it is possible to balance between minimum coflow progress and average CCT. In practice, since coflow links have intertwined dependencies, the coflow progress approach builds on the notion of DRF (Dominant Resource Fairness) [22] in order to allocate rates based on the maximum per port demand. The

arXiv:2208.06513v1 [cs.DC] 12 Aug 2022

*Laboratoire informatique d’Avignon (LIA), Avignon University, France,
[†]He was with LIA, Avignon University, France and now he is with Interdisciplinary Centre for Security, Reliability and Trust (SnT), University of Luxembourg, *Huawei Technologies, Paris Research Center, France.

progress of a coflow tracks the slowest flow – determining its CCT – on the related bottleneck port at each instant, while all remaining flows of the coflow are slowed down in order to improve network utilization without impairing the CCT. Finally, a max-min fair argument grants a bottleneck share per coflow in the form of a minimum rate allocation, i.e., a so-called minimum coflow progress guarantee. Apart from a few exceptions [10], any further performance optimization, e.g., minimization of average CCT or maximization of port utilization, is pursued on top of this baseline constant rate allocation [18]–[20], denoted as *static* progress in the rest of the paper. Note that, from the perspective of a coflow being served, this is just a lower bound on the average rate received before completion.

The notion of coflow fairness based on *static* progress has limitations as detailed formally in the next sections. In fact, it may not account for the structure of coflows and their conflicting demands. Furthermore, ensuring a constant minimal allocation corresponds to a non-preemptive definition of progress because it requires every coflow to send a minimum amount of data with no interruption over each engaged port.

Main contributions. The fairness framework introduced in this paper is based on a notion of progress which is the average rate granted to a coflow while in service. The coflow progress is tuned by controlling the inverse metric, e.g., the *slowdown* [20]. The slowdown measures the additional delay on the data transfer time of a coflow compared to isolation, i.e., when a coflow is scheduled alone in the datacenter fabric. An optimal scheduling problem is formulated under a generalized slowdown constraint, obtaining a notion of fairness more general than the notion of progress used in the literature. Similar to α -fair scheduling [15], the *slowdown* constraint acts as a single parameter able to strike the trade-off between efficiency, i.e., CCT performance, and fairness, i.e., the slowdown experienced by a coflow. Furthermore, a fully polynomial time algorithm, namely MPS (Minimum Primal Slowdown) provides an accurate estimation of the minimum feasible slowdown for a batch of coflows.

Finally, a scheduling algorithm CoFair extends the family of the primal-dual algorithms [8], [23]: by using the output of MPS, it produces a *primal-feasible* σ -order and an approximation factor of 4.

Actually, the target coflow progress is attained based on a per coflow prioritization, avoiding expensive rate-control mechanisms. CoFair is tested against Sincronia [8], near-optimal for CCT minimization, and against Utopia [10] for coflow fairness. CoFair performs closely to Sincronia in CCT and to Utopia in slowdown, thus improving the tradeoff, with low complexity since no rate control is required. To the best of the authors’ knowledge, the formal framework for fair scheduling in data-transfer based on controlled slowdown and the solutions provided in this paper are new contributions to the discussion on enforcing fairness in coflow scheduling.

The paper is organized as follows. Sec. II resumes existing works on fair coflow scheduling. Sec. III introduces the system model and the relations between CCT, coflow progress and

slowdown. Sec. IV describes the mathematical framework as a scheduling problem with slowdown constraints and its feasibility. Algorithmic solutions are proposed in Sec. VI in the set of primal-feasible σ -order schedulers. Sec. VII reports on numerical results and a concluding section ends the paper.

II. RELATED WORKS

In the coflow literature only a few papers deal with fairness issues. The seminal work [1] introduced the weighted CCT as objective function and later works pursued same approach, e.g., [24], even though the relation between coflow weights, CCT and coflow fairness remains elusive. In general, fairness metrics based either on weights or on constant rate allocation [18], [24] lead to sub-optimal scheduler and the usage of the CCT geometric mean [18] tends to privilege shorter coflows. In the scheduler Varys [12] coflow prioritization has been introduced as a means to mitigate starvation of coflows with low priority. Actually, the smallest-effective-bottleneck-first heuristic used in Varys is a pre-emptive scheduling which prioritizes greedily a coflow with the smallest remaining bottleneck’s completion time.

Thus, granting a constant fraction of port bandwidth per coflow became the accepted notion of coflow progress adopted later on in the literature [19], [20]. This coflow progress guarantee – already denoted static progress – has been formalized for the first time in [20]. However, in that context the progress is as a static rate constraint under the general objective of maximizing the fabric utilization. The idea of a progress constraint coupled to the CCT minimization appeared first in [19], where the (static) progress of each flow is to be set above certain target threshold.

An interesting attempt to develop a formal framework based on max-min fairness has been presented [17], based on lexicographic ordering for flow-level rate allocation [25]. However, due to its ease of implementation, the notion of fairness by static coflow progress guarantees has been preferred in the literature. Other authors [26] define the fairness degree of a data-transfer scheme with respect to standard schemes such as Hug [20] or DFR [22].

More recently, a *relative* fairness concept has been introduced in order to compare coflow schedulers to benchmarks such as Hug or DFR [26]. In particular, [10] has proposed a long term isolation guarantee by providing a bound on CCTs relative to DRF performance. While this type of bound resembles the slowdown constraint introduced in this work, the CCT-slowdown tradeoff cannot be controlled.

In the coflow deadline satisfaction (CDS) problem [27], [28] each coflow is subject to a completion deadline similar to a slowdown constraint. However, in that context the target is to operate joint *coflow admission control and scheduling* to maximize the number of admitted coflows which respect their deadlines.

The solution proposed in this paper is rooted on the analysis of a scheduling problem where the minimal slowdown constraint is connected to a primal-dual formulation. The resulting scheduler has no need to perform rate control and it can

Symbol	Meaning
\mathcal{L}	set of switch links $k = 1, \dots, 2M$
\mathcal{M}	coflow schedule
\mathcal{C}	set of coflows $\mathcal{C} = \{1, \dots, N\}$
\mathcal{F}_j	set of flows of coflow j ; $n_j = \mathcal{F}_j $ (coflow width)
$p_{\ell j}$	volume of coflow j on port ℓ
C_j	coflow completion time for coflow j
w_j	weight of coflow j
r_j	release time of coflow j
\bar{E}	slowdown constraint
$x_j^i(t)$	fraction of flow i of coflow j scheduled in slot t
$Y_j^i(t)$	progress of flow i of coflow j at time t
$Z_j(t)$	completed fraction of coflow j at time t
$Z_j^i(t)$	completed fraction of flow i of coflow j at time t
$x_j^i(t)$	fraction of flow i of coflow j at time t
$X_j^i(t)$	indicating variable for coflow j at slot t (binary)
$X_j^i(t)$	indicating variable for i of coflow j at slot t (binary)
B_ℓ	capacity of link ℓ
T	time horizon
Δ	time slot duration

TABLE I
LIST OF KEY NOTATIONS

be fully implemented by greedy rate allocation via priority queuing [8]. The key technical result is the equivalence of the weight adjustment step introduced in [23][Alg. 3.1] – later applied in [8] to coflow scheduling – with Gaussian elimination. Coupled with the novel results on primal-feasibility, this leads to an approximation solution based on a σ -order prioritization algorithm under generalized slowdown constraints.

In the rest of the paper, reference solutions are near-optimal Sincronia [8] for CCT minimization and Utopia [10] for coflow fairness.

III. MOTIVATIONS

A. System model

The most popular coflow scheduling model is based on a non blocking switch connection of the type reported in Fig. 1, often called the Big Switch model. This model is adequate because the bisection bandwidth of modern datacenters exceeds access capacity, so that congestion events occur at the inbound or outbound ports of top-of-rack switches. From now on the term link and port will be used interchangeably. The set of switch links is $\mathcal{L} = \{1, \dots, 2M\}$, where links $1 \leq \ell \leq M$ are ingress ports and $M+1 \leq \ell \leq 2M$ are egress ports. Each port has capacity B_ℓ .

Let $\mathcal{C} = \{1, \dots, N\}$ be an input set of coflows, i.e., a *batch*. Each coflow j is a set \mathcal{F}_j of n_j flows and a flow represents a shuffle connection over a pair of input-output ports. The release time $r_j \geq 0$ of coflow j is the time when its shuffle phase starts. A component flow $i \in \mathcal{F}_j$ has volume v_j^i . The *size* of coflow j , that is the total volume of coflow j , is denoted by $V_j = \sum_{i \in \mathcal{F}_j} v_j^i$. Let $\chi_j^i(\ell)$ indicate if flow $i \in \mathcal{F}_j$ is active on port ℓ . The volume of coflow j active on port ℓ is denoted

$$p_{\ell j} = \sum_{i \in \mathcal{F}_j} v_j^i \chi_j^i(\ell)$$

Tab. I summarizes all the important notations and their descriptions.

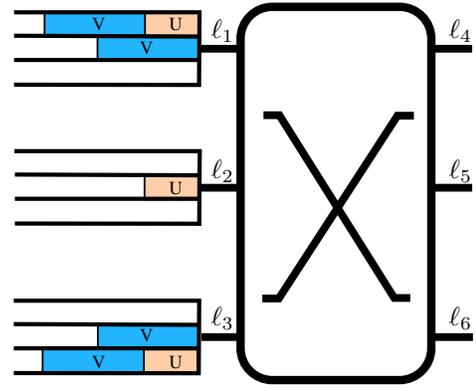


Fig. 1. Coflow scheduling over a 3×3 datacenter fabric with three ingress and three egress ports $\mathcal{L} = \{\ell_1, \dots, \ell_6\}$. Flows in ingress ports are organized by destinations and color-coded by coflows: $j = 1$ orange and $j = 2$ blue.

Let C_j be the CCT of coflow j : C_j is the epoch when the last flow in \mathcal{F}_j is fully transferred. Let C_j^0 be the coflow completion time *in isolation*, i.e., if all network resources serve solely the tagged coflow. It holds $C_j^0 =: r_j + \max_{\ell \in \mathcal{L}} \frac{1}{B_\ell} \sum_{i \in \mathcal{F}_j} v_j^i \chi_j^i(\ell)$. The standard coflow scheduling problem corresponds to determine a coflow schedule \mathcal{M} able to minimize the weighted sum of the CCTs under the capacity constraints imposed by the network fabric. A coflow schedule \mathcal{M} provides the set of rates which are assigned to each flow of each coflow per port. The large class of σ -order schedulers [7], for instance, prescribes to respect a static preemption priority σ and only requires monitoring of coflows at ingress/egress (I/E) ports of the Big Switch. As proved in [8], the efficiency loss w.r.t. an optimal scheduler respecting the given priority σ is at most a factor 2. In the numerical evaluations in Sec.VII, rates are assigned following a strict priority rule, namely the greedy rate allocation policy [8]. In practice, this rule can be implemented with a priority queuing scheduler with a number of queues equal to the number of coflows.

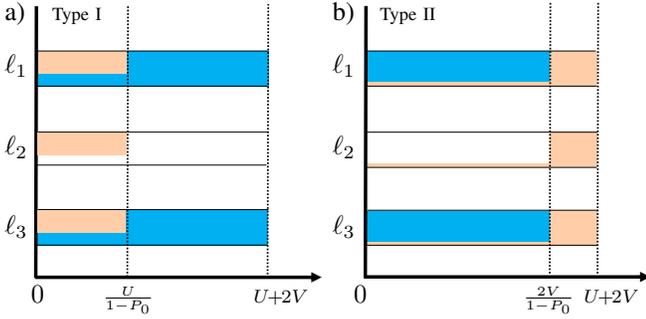
B. Static coflow progress

In the coflow literature, the standard definition of fairness [18], [20] is based on the notion of coflow progress. Given an input batch of coflows, the objective is to minimize the weighted CCT while ensuring that every coflow has a target minimum progress *per link*, e.g., per input or output port in the Big Switch model. Let $a_{\ell j}$ be the rate guaranteed to coflow j on port ℓ . The progress of coflow j is defined as

$$P_j = \min_{\ell \in \mathcal{L}} \left\{ \frac{a_{\ell j}}{d_{j\ell}} \right\} \quad (1)$$

where $d_{j\ell} := \frac{p_{\ell j}}{C_j^0 - r_j}$ is also called the rate demand of flow i of coflow j and it represents the lowest possible rate to dispatch $p_{\ell j}$ within C_j^0 . From the definition of C_j^0 indeed $P_j \leq 1$. Finally, coflow fairness ensures a per-coflow minimum progress

$$P_j \geq P_0, \quad \forall j \in \mathcal{C} \quad (2)$$



	slowdown ($\phi_j = 1$)		progress	
	E_1	E_2	R_1	R_2
Type I	$\frac{1}{1-P_0}$	$1 + \frac{U}{2V}$	$3(1-P_0)$	$\frac{4V}{2V+U}$
Type II	$1 + \frac{2V}{U}$	$\frac{1}{1-P_0}$	$\frac{3U}{2V+U}$	$2(1-P_0)$

Fig. 2. Example: a) schedule of Type 1 and b) schedule of Type 2.

It is worth remarking that, though equivalent to the definition appearing in the literature, (1) has a straightforward interpretation in terms of per port *bitrates*. In some works, the constant parameter P_0 is also called as isolation guarantee [20]. The next example is meant to outline why the notion of static progress does not fully capture the performance-fairness trade-off in coflow scheduling.

Example: The sample coflow batch of Fig. 1 is served on a fabric with $M = 3$ servers and unitary bandwidth links. Coflow 1 marked orange occupies all input and output ports of the fabric; it has 3 flows with volume U traffic units each. Coflow 2 marked blue consists of flows of volume V traffic units each, two on port ℓ_1 and two on port ℓ_3 . In Fig. 1 the index of the virtual queues at ingress ports indicate the output port. For instance the second flow of coflow 2 sends $v_2^2 = V$ traffic units to port ℓ_5 .

Let P_0 denote the static progress guarantee. By the symmetry of the problem, same rate guarantee is due on input ports for the two coflows, i.e., $P_1 = a_{\ell_{11}} = a_{\ell_{21}} = a_{\ell_{31}} \geq P_0$ and $P_2 = a_{\ell_{12}} = a_{\ell_{32}} \geq P_0$. Clearly, $P_0 \leq 1/2$ due to bandwidth constraints. As depicted in Fig. 2, based on the CCT order of coflows, two types of schedules are possible: Type I where coflow 1 finishes first or Type II where coflow 1 finishes last. For Type I schedules it holds $C_2 = U + V$, whereas $C_1 = U/(1 - P_0)$. For Type II, $C_1 = U + V$ and $C_2 = 2V/(1 - P_0)$. By letting $P_0 > 0$, for Type I the CCT of coflow 1 has been degraded with no gain for coflow 2. Conversely, for Type II schedules, the CCT of coflow 2 is degraded with no gain for coflow 1. Thus, while the static progress guarantees a minimum rate allocation per coflow, the CCT of some of the coflows can be degraded without any actual gain for the slowdown of the other ones.

C. Progress and slowdown

The progress of coflow j under coflow schedule \mathcal{M} is defined as its average rate attained while in the system, i.e., $R_j = \frac{V_j}{C_j - r_j}$. To control the fairness-performance tradeoff, it is convenient to use the coflow *slowdown*, which is in fact the

inverse quantity of the progress: a target progress guarantee can be enforced via a constraint on the maximum slowdown. The slowdown is the ratio between the progress in isolation R_j^0 and the progress under a given schedule \mathcal{M} . This plain notion of slowdown can be generalized to allow for a non-negative weight per coflow which is a function of the coflow geometry.

Definition 1 (Generalized coflow slowdown). *The generalized slowdown of coflow j under coflow schedule \mathcal{M} is the weighted ratio $E_j = \phi_j \frac{R_j^0}{R_j} = \phi_j \frac{C_j - r_j}{C_j^0 - r_j}$.*

The (generalized) slowdown captures the relative degradation of the data transfer duration when a coflow is scheduled within a batch of competing ones compared to its duration in isolation. In the numerical section, two variants of slowdown are considered:

- *plain* slowdown: $\phi_j = 1$, this is the standard definition of slowdown [10], [12], [29] where, between two coflows, a bound on the maximum slowdown prioritizes the coflow with smaller CCT in isolation;
- *slowdown with port-occupation*: $\phi_j = V_j$, with this definition, between two coflows with same volume, a bound on the maximum slowdown prioritizes the coflow with smaller port-occupation.

Example: The tradeoff between average CCT and progress for Type I and for Type II schedules is described in Fig. 2. For Type I schedules the slowdown of the coflow finishing first is $1/(1 - P_0)$. The progress guarantee as a function of the static progress guarantee P_0 writes $\bar{R}_I = \min\{3(1 - P_0), \frac{4V}{2V+U}\}$. In case of Type II schedules it writes $\bar{R}_{II} = \min\{2(1 - P_0), \frac{3U}{2V+U}\}$. Note that both functions are non increasing. Also, by letting $P_0 = 0$, both \bar{R}_I and \bar{R}_{II} attain their maximum and the attained average CCT is also minimized for $P_0 = 0$ for both Type I and Type II schedules.

Finally, the following cases hold. For $2V < U$ the Type II schedule attains the minimum average CCT, i.e., $(5V + U)/2$ and the progress guarantee $R^* = 3U/(2V + U)$ is also better off than Type I. Conversely for $4V > 3U$ a Type I schedule attains the minimum average CCT, i.e., $V + U$ and the best progress guarantee, i.e., $R^* = 4U/(2V + U)$.

It is interesting to remark that in the above two cases there is *no performance-fairness* tradeoff, because minimizing the average CCT is also ensuring the largest minimum progress.

Conversely, for $2U < 4V < 3U$, a performance-fairness tradeoff does exist: a Type II schedule attains better average CCT but worse progress than a Type I one; on the contrary, a Type I schedule attains worse average CCT but better progress than a Type II one.

From the example, the both Type I and Type II schedules do attain the best progress guarantee and are better off w.r.t. CCT figures for $P_0 = 0$ than for $P_0 > 0$. This fact is expressed as a general result in the next section, ruling out static progress guarantees in fair coflow scheduling.

IV. SCHEDULING UNDER SLOWDOWN CONSTRAINTS

The standard coflow scheduling optimization framework determines a coflow schedule \mathcal{M} by minimizing the weighted sum of the CCTs under capacity constraints imposed by the network fabric. It is often formulated as a Mixed Integer Linear problem (MILP) for a finite horizon where time is slotted into T time slots of duration Δ ; let $\Delta = 1$ without loss of generality. The decision variables $\{x_j^i(t)\}$ represent the fraction of flow $i \in \mathcal{F}_j$ of coflow j scheduled in slot t . The fraction of volume of flow i of coflow j already transmitted by slot t is $Z_j^i(t) = \sum_{v=1}^t x_j^i(v)$. The fraction of volume coflow j transmitted by time t writes $Z_j(t) := \frac{1}{V_j} \sum_{i \in \mathcal{F}_j} v_j^i Z_j^i(t)$; note that $C_j = \min\{t \mid Z_j(t) = 1\}$. Finally, the standard formulation with MILP representing the Progress Scheduling (PS) problem writes

$$\text{minimize: } \sum w_j C_j \quad (\text{PS})$$

$$\text{subj. to: } \sum_{i=1}^T x_j^i(t) = 1, \quad \forall j \in \mathcal{C}, i \in \mathcal{F}_j \quad (3)$$

$$X_j(t) \leq Z_j(t), \quad \forall j \in \mathcal{C}, \forall t \quad (4)$$

$$Z_j(t) \leq 0, \forall t \leq r_j \quad (5)$$

$$C_j \geq 1 + \sum_{z=1}^T (1 - X_j(z)), \quad \forall j \in \mathcal{C} \quad (6)$$

$$E_j := \frac{C_j - r_j}{C_j^0 - r_j} \leq \frac{E}{\phi_j}, \quad \forall j \in \mathcal{C} \quad (7)$$

$$\sum_{j \in \mathcal{C}} \sum_{i \in \mathcal{F}_i} v_j^i x_j^i(t) \chi_j^i(\ell) \leq B_\ell, \quad \forall t, \forall \ell \in \mathcal{L} \quad (8)$$

$$x_j^i(t) \in [0, 1], X_j(t) \in \{0, 1\}, \quad \forall t \quad (9)$$

where E is a non-negative parameter and weight w_j denotes the weight of coflow j . Let denote $\{C_j^*\}$ a solution of the PS and $\{R^*\}$ the corresponding progress vector.

The MILP makes use of coflow completion variables $Z_j(t)$ and binary variables $X_j(t)$. Note that, due to minimization, $X_j(t)$ behaves as the indicating variable of coflow completion, that is $X_j(t) = 1$ if $Z_j(t) = 1$ and $X_j(t) = 0$ if $Z_j(t) < 1$. Constraint (3) states that each flow should be completed by the end of the time horizon; (4) states that the completion of a coflow is subject to the transfer of the whole coflow volume. In (5) the release time of coflows imposes the constraint on the beginning of a coflow transmission. Constraint (6) binds coflow completion time and coflow completion variables, which appeared first in [7]. Finally, (8) provides the constraint on the port capacity; $Z_j(t)$ is defined implicitly in (4) and (5) for notation's sake. Finally, (7) is the slowdown constraint to enforce a target coflow progress, where constant $E \geq \max_j \phi_j$. Note that ϕ_j permits to express a preference: for instance, when $\phi_j = V_j$, if two coflows have same total volume, the one which occupies the fabric ports for shorter time is preferred.

Clearly, adding progress constraints leads to a deterioration of the average CCT with respect to the plain weighted CCT

minimization ($E = +\infty$). For the example in the previous section, when $\phi_j = V_j$, the constraint (7) writes $C_1/C_1^0 \leq E/3U$ and $C_2/C_2^0 \leq E/4V$. Letting $E = 2(U + 2V)$ attains the schedule of Type I. The schedule of Type II is attained when $2V < U$. In such case (7) is inactive, i.e., there is no tradeoff. However, Type I is to be preferred when the port occupation of coflow 2 is larger than that of coflow 1.

Next, the notion of static progress (1) and (2) existing in the literature [19] and that of progress based on coflow slowdown are compared. The Static Progress Scheduling (SPS) problem is described below:

$$\text{minimize: } \sum w_j C_j \quad (\text{SPS})$$

$$\text{subj. to: } (3)(4)(5)(6)(8) \quad (10)$$

$$\sum_{i \in \mathcal{F}_j} v_j^i x_j^i(t) \chi_j^i(\ell) \geq a_{\ell j} \Delta, \quad \forall t \forall \ell \in \mathcal{L} \quad (11)$$

$$x_j^i(t) \in [0, 1], X_j(t) \in \{0, 1\}, \quad \forall t \quad (12)$$

where the minimum rate allocations $\{a_{\ell j}\}$ obey to (1) and (2) for a given progress parameter P_0 . With some abuse of terminology, let denote $\{\bar{C}_j\}$ a solution of SPS.

From an analogous result for the weighted CCT minimization, it follows that both PS and SPS problems are NP-hard and inapproximable below a factor 2 [7].

Let \bar{R}_j the average progress guarantee offered by SPS for coflow j corresponding to allocation $\{a_{\ell j}\}$ so that $\bar{R}_j := \frac{1}{2} \sum_{\ell} \sum_{i \in \mathcal{F}_j} a_{\ell j}^i B_\ell$. As shown next, compared to (11), the slowdown constraint (7) allows more flexible rate allocation since the rate on some links can be increased by reducing or pausing the transmission on some other links used by the coflow. Formally, for every SPS scheduler there exists a slowdown value E and a PS scheduler which is better off.

Theorem 1. *Let consider \mathcal{C} and static progress defined for $P_0 > 0$ according to (1) and (2). There exists $E \geq 0$ such that a solution of SPS is feasible for PS. Furthermore, there exists a schedule such that*

$$\sum C_j^* < \sum \bar{C}_j, \quad \text{and} \quad R_j^* \geq \bar{R}_j \quad \forall j \in \mathcal{C}$$

Proof: i. Let \mathcal{M} be a schedule solving SPS for a given pair $\{a_{\ell j}\}$ and P_0 . Let $\{\bar{C}_j\}$ the corresponding result. It is possible to determine E so that such a schedule is a solution of PS as well. From (1) and (2)

$$\sum_{i \in \mathcal{F}_j} v_j^i x_j^i(t) \chi_j^i(\ell) \geq a_{\ell j} \Delta, \quad \frac{a_{\ell j} (C_j^0 - r_j)}{p_{\ell j}} \geq P_0$$

For every coflow we can write

$$\begin{aligned} \bar{C}_j - r_j &\leq \max_{\ell \in \mathcal{L}} \left\{ \frac{p_{\ell j}}{a_{\ell j}} \right\} = \frac{(C_j^0 - r_j)}{\min_{\ell \in \mathcal{L}} \frac{a_{\ell j} (C_j^0 - r_j)}{p_{\ell j}}} \\ &= \frac{C_j^0 - r_j}{P_j} \leq \frac{C_j^0 - r_j}{P_0}, \quad \forall j \in \mathcal{C} \end{aligned} \quad (13)$$

With the identification $E = \frac{\max V_j}{P_0}$, the schedule \mathcal{M} is feasible also for PS.

ii. The attained progress guarantees can be bounded by

$$\begin{aligned} \bar{R}_j &= \frac{1}{2} \sum_{\ell \in \mathcal{L}} a_{\ell j} B_\ell \leq \frac{1}{2} \frac{1}{\bar{C}_j - r_j} \sum_{t=r_j}^{\bar{C}_j} \sum_{\ell \in \mathcal{L}} \sum_{i \in \mathcal{F}_j} v_j^i x_j^i(t) \chi_j^i(\ell) \\ &= \frac{1}{2} \frac{2V_j}{\bar{C}_j - r_j} = R_j \end{aligned} \quad (14)$$

where the factor 2 appears since every flow is counted at the ingress at the egress port. By adding feasible constraints $C_j - r_j \leq V_j/\bar{R}_j$ to SPS and dropping (11) the thesis follows. ■

From Thm. 1, scheduling with slowdown constraints provides larger feasibility space and at least same progress guarantees than SPS. As seen in our example, the corresponding minimum weighted CCT which can be attained under same progress guarantees is strictly worse for $P_0 > 0$. However, the introduction of slowdown constraints on top of weighted CCT minimization poses the issue of feasibility, i.e., whether or not, for a given instance of the problem, a target slowdown constraint is feasible and what is the smallest feasible value of E .

V. SLOWDOWN FEASIBILITY

In this section the feasibility of PS is discussed and means to test it in polynomial time are presented. In particular, the largest possible average progress is obtained by minimizing E in (7) while ensuring the feasibility of (PS). Then, the class of σ -order preserving schedulers is extended to account for slowdown constraints.

Testing slowdown feasibility is possible in the form of a linear program (LP) feasibility test, from which the following result hold:

Theorem 2. *i. The feasibility of PS can be determined in polynomial time. ii. If an instance of PS is feasible for a given value if E where $r_j = 0$ for all $j \in \mathcal{C}$, then it is also feasible for all instances obtained by letting some of the r_{js} positive.*

The proof and the actual LP formulation are detailed in the Appendix. We denote E^* the minimum feasible parameter appearing in (7). It is interesting to remark that from the previous result the case of different release times is beneficial with respect to the coflow progress; the intuition is that for a later released coflow, it is convenient if some of earlier coflows have been already dispatched. In the next section we shall introduce a much faster approximated yet accurate algorithm to determine E^* (namely Alg. 1). From Thm. 2, E^* can be determined in pseudo-polynomial time by solving LP iteratively with a bisection search exploring the slowdown parameter E . By simple calculations this search is $O(\frac{(NV)^{3.5}}{\epsilon})$ where $V = \sum V_j$, and $\epsilon > 0$ is a tolerance, where we used the $O(n^{3.5})$ bound in the number of variables n .

A. σ -order schedulers

Even if the above result is encouraging, one fundamental problem with techniques based on MILPs in coflow schedul-

ing is that the number of flows per coflow maybe large¹. Furthermore, the number of variables involved in such time-indexed systems increases with the inverse of the time slot Δ . Hence, the number of per slot rate decision variables easily range in hundred of thousands for a time slot Δ in the order of seconds. Thus, scalability issues do arise even for approximation algorithms obtained by solving LPs from MILP relaxation [7], [11]. Same scalability issues occur solving the time-indexed LP to determine the minimal slowdown E^* .

For the sake of striking a fairness-performance tradeoff, solving PS exactly may not be feasible. Rather, it is possible to design a coflow scheduler to approximate a target average progress guarantee, and yet perform near-optimally with respect to the weighted CCT. We can regard this approach as a *soft* coflow fairness guarantee versus the *hard* coflow fairness guarantee provided by PS.

It is possible to do so using the class Σ of σ -order preserving schedulers (σ -order schedulers for short). These schedulers require to set the priority of coflows and use such order to determine rate allocation [8], [13]. Let order $\sigma : \mathcal{C} \rightarrow \mathcal{C}$ be a permutation of the coflow set: $j = \sigma(k)$ means that the k -th priority is assigned to coflow j .

Definition 2. *Fix an order $\sigma : \mathcal{C} \rightarrow \mathcal{C}$: the scheduler \mathcal{M} is a σ -order preserving scheduler if*

- *it is pre-emptive;*
- *it is work conserving;*
- *it respects the pre-emption order: no flow of coflow $\sigma(j)$ can be pre-empted by a flow of $\sigma(k)$, $k = j + 1, \dots, N$ on any port on which it has pending flows.*

The greedy rate allocation, for instance, blocks a coflow $\sigma(j)$ on some port iff it is occupied by some coflow $\sigma(k)$, $k < j$. As first proved in [8], the design of σ -order schedulers can be performed using the analogy of the Big Switch ports with correlated machines and using related results for open-shop scheduling [23]. Following the same idea, once E is fixed, schedulers in Σ can approximate the solution of PS problem. In the rest of the discussion, we consider the case $r_j = 0, \forall j \in \mathcal{C}$.

B. Extending σ -order schedulers: primal feasibility

The theory of σ -order scheduling is rooted on the work [23] for concurrent open shop problems where jobs are scheduled on parallel machines. Let define $p_{\ell j} = \sum v_j^i \chi_j^i(\ell)$ the total volume to be transferred by coflow j over port ℓ . Denote $V_\ell = \sum v_j^i \chi_j^i(\ell)$ the total volume on port ℓ and \mathcal{C}_ℓ is the set of coflows active on port ℓ . The equivalent slowdown constraint $D_j = \frac{EC_j^0}{V_j}$ is introduced for notation's sake. With the job scheduling terminology, $p_{\ell j}$ is the duration of the job j over

¹It may range in the order of tenths of thousands [5] in production datacenters

machine ℓ . The corresponding scheduling formulation writes [23]

$$\text{minimize: } \sum_{j \in \mathcal{C}} w_j C_j \quad (\text{Primal LP})$$

$$\text{subj. to: } \sum_{j \in A} p_{\ell j} C_j \geq f_{\ell}(A), \quad \forall A \subseteq \mathcal{C}, \ell \in \mathcal{L} \quad (15)$$

$$C_j \leq D_j, \quad \forall j \in \mathcal{C} \quad (16)$$

The set functions $f_{\ell} : \mathcal{P}(\mathcal{C}_{\ell}) \rightarrow \mathbb{R}$ are defined $f_{\ell}(S) = \frac{1}{2}[(\sum_{j \in S} p_{\ell j})^2 + \sum_{j \in S} p_{\ell j}^2]$ for $S \subset \mathcal{C}_{\ell}$. Constraints (15) are known as parallel inequalities. Let restrict to the set of schedulers Σ_b , that is the set of σ -order schedulers such that the permutation σ is *primal-feasible*, that is

$$C_{k\ell} := \frac{1}{B_{\ell}} \sum_{j=1}^k p_{\ell \sigma(j)} \leq D_k, \quad \forall \ell \in \mathcal{L}, \forall k \in \mathcal{C} \quad (17)$$

In fact, it is always possible to consider just σ -order schedulers

Theorem 3. *If Primal LP is feasible, it is solved by σ -order which is primal-feasible.*

Proof: Let consider the set of permutation schedules $\{\sigma_{\ell}\}$ which schedule coflows active on port ℓ according to the EDD rule (Earliest Due Date rule) [30]. Since Primal LP is primal-feasible, all such schedules $\{\sigma_{\ell}\}$ are feasible. Now let derive a feasible σ -order using the following procedure: select the coflow k for which C_k is maximum, and modify $\sigma_{\ell S}$ for all ports where $p_{\ell k} > 0$ by scheduling k last on all its active ports, and leaving the order of the others unchanged. After this operation it holds $C_{k\ell} \leq D_k$ and $C_{j\ell} \leq D_j$ for all $j \neq k$ and for all ports ℓ so that the resulting $\{\sigma_{\ell}\}$ are still feasible. Let $\sigma(k) = |\mathcal{C}|$, and then remove it from $\{\sigma_{\ell}\}$. Let $\mathcal{C} \leftarrow \mathcal{C} \setminus \{k\}$. By repeating this procedure, in N steps a feasible σ -order is obtained. ■

Let E^p be the minimal value such that Primal LP is feasible. From the proof of Thm. 4, a bisection search algorithm based on a per port EDD scheduling can determine the value E^p in pseudo-polynomial time. The next result provides a connection between the Primal LP and PS:

Theorem 4. *If PS is feasible, then Primal LP is feasible, i.e., $\Sigma \subseteq \Sigma_b$.*

Proof: It is sufficient to consider a schedule \mathcal{M} which is feasible for PS and then consider the σ -order which sorts coflows according to the CCTs produced by \mathcal{M} , denoted $\{\bar{C}_k\}$. In fact, let consider coflow $k \in \mathcal{C}$ and port ℓ where the last of its flows completes. Since coflow k completes on port ℓ after all coflows $j = 1, \dots, k-1$, the volume of all flows on ℓ which belong to $j = 1, \dots, k-1$ must be over by time \bar{C}_k . This implies that

$$D_k \geq \bar{C}_k \geq \sum_{j=1}^k p_{\ell j} \quad (18)$$

Algorithm 1 Pseudocode of MPS

Input: $\mathcal{C}, \{\phi_j\} \{p_{\ell,j}\}$

Output: E

```

1: Sort  $\mathcal{C}$  in decreasing order of  $\tilde{R}_j^0 = \frac{\phi_j V_j}{V_j C_j^0}$  % sorting by  $\mathcal{R}_j^0$ 
2:  $E_{\ell} = +\infty, \quad \forall \ell \in \mathcal{L}$  % initialise slowdown
3: for  $\ell = \ell_1, \dots, \ell_{2M}$  do
4:   for  $j \in \mathcal{C}_{\ell}$  do
5:      $Z_j = \tilde{R}_j^0 \sum_{k=1}^j v_k^{\ell}$  % per coflow estimate on  $\ell$ 
6:   end for
7:    $E_{\ell} = \max\{Z_1, \dots, Z_{|\mathcal{C}_{\ell}|}\}$  % slowdown estimate per port
8: end for
9:  $E = \max\{E_{\ell_1}, \dots, E_{\ell_{2M}}\}$  % for every port
10: return  $E$  % return the slowdown

```

so that (16) holds. Moreover, it holds

$$\sum_{j \in \mathcal{C}} p_{\ell,j} \bar{C}_j \geq \sum_{j \in \mathcal{C}} p_{\ell,j} \sum_{k=1}^j p_{\ell k} = f_{\ell}(\mathcal{C})$$

so that (15) is satisfied as well. By identifying $C_j := \sum_{k=1}^j p_{\ell k}$, the proof holds. ■

Thus, if there exists a solution for PS in Σ , then it is to be found in Σ_b . However, since primal-feasibility does not account for cross-dependencies among ports, it provides a weaker estimate of the completion time for the flows of a coflow according to a given σ -order.

Corollary 1. $E^p \leq E^*$.

Nevertheless, the experimental results reported in Sec VII indicate that the two values practically coincide. To this respect, the calculation of E^p provides a substantial computational advantage since it can be performed by a lightweight *fully polynomial time algorithm*. It is called MPS, which stands for Minimum Primal Slowdown.

The pseudocode of MPS is reported in Alg. 1. The value of E which satisfies the slowdown constraint (7) is minimum when the CCT of the coflow with the largest value $\frac{\phi_j C_j}{C_j^0}$ is minimized. The algorithm visits each port ℓ by sorting \mathcal{C}_{ℓ} in decreasing order of $\tilde{R}_j^0 = \frac{\phi_j V_j}{V_j C_j^0} = \frac{\phi_j}{V_j} R_j^0$. For each link ℓ , the iteration at step 7 provides a lower bound on E^p per coflow active on ℓ using (7). The maximum among such estimates is selected at step 9. Due to the sorting at line 1 and the N_{ℓ} products at line 5, it is easy to see that the complexity of Alg. 1 is in $O(MN + N \log(N))$. The next result proves that the output of MPS is correct.

Theorem 5. *MPS returns E^p*

Proof: From Primal LP, fixed a primal schedule, it holds $E \geq \tilde{R}_j^0 C_j$ for $\forall j \in \mathcal{C}$, and more precisely

$$E^p = \min_{\sigma \in \Sigma_b} \left\{ \max_{\ell \in \mathcal{L}} \max_{j \in \mathcal{C}_{\ell}} \tilde{R}_j^0 \sum_{i \in \mathcal{C}} p_{\ell \sigma(i)} \right\}$$

Let show that E^p is attained by a primal schedule which sorts coflows with decreasing \tilde{R}_j^0 . Let assume $\tilde{R}_1^0 \geq \dots \geq \tilde{R}_N^0$. Let

σ_ℓ be the permutation induced by σ on a tagged port $\ell \in \mathcal{L}$. For the last scheduled coflow on port ℓ it holds

$$E^P \geq \min_{j \in \mathcal{C}} \tilde{R}_j^0 \sum_{i \in \mathcal{C}} p_{\ell\sigma(i)} = \tilde{R}_{N_\ell}^0 \sum_{i=1}^{N_\ell} p_{\ell i} \quad (19)$$

Also, let $\sigma' = (\sigma(1), \dots, \sigma(N_\ell - 1))$, it holds

$$E^P \geq \max_{j \in \mathcal{C}_\ell \setminus \{N_\ell\}} \max \left\{ \tilde{R}_j^0 \sum_{\sigma'(i) \leq \sigma'(j)} p_{\ell\sigma'(i)} \right\}$$

because by eliminating the last scheduled coflow $\sigma(N)$, the maximum slowdown cannot increase. Now, it is possible to write

$$E^P \geq \max_{\sigma_\ell(N_\ell)} \left\{ \tilde{R}_{\sigma(N_\ell)}^0 \sum_{i=1}^{N_\ell} p_{\ell i}, \max_{j \in \mathcal{C}_\ell \setminus \{N_\ell\}} a_j \sum_{\substack{i \in \mathcal{C}_\ell \setminus \{N_\ell\} \\ \sigma(i) \leq \sigma'(j)}} p_{\ell\sigma'(i)} \right\} \\ \geq \max \left\{ \tilde{R}_{N_\ell}^0 \sum_{i=1}^{N_\ell} p_{\ell i}, \max_{j \in \mathcal{C}_\ell \setminus \{N_\ell\}} \sum_{\sigma'(i) \leq \sigma'(j)} p_{\ell\sigma'(i)} \right\}$$

Hence an optimal permutation writes $\sigma = (\sigma_1, \dots, \sigma_{N_\ell-1}, N_\ell)$. The same argument holds for σ' and $\mathcal{C}_\ell \setminus \{N_\ell\}$, so that an optimal permutation writes $\sigma = (\sigma_1, \dots, \sigma_{N_\ell-2}, N_{\ell-1}, N_\ell)$. After N_ℓ iterations, the theorem follows. ■

As a side result of the above proof it follows

Corollary 2. *Primal LP is feasible if and only if the schedule which sorts coflows with decreasing \tilde{R}_j^0 is feasible.*

On feasibility and primal feasibility.

Finally, the minimum value of the parameter E^P for which the formulation is primal-feasible is in general smaller than the minimum possible value E^* which renders PS feasible. The fact that they practically coincide – as seen in the numerical section – is to be ascribed to the fact that the largest slowdown corresponds invariably to the CCT of a coflow finishing last on some port. Hence, it is in fact determined by the total volume transferred on that port. This is actually the key step for the slowdown estimation performed by MPS at line (5).

A primal-feasible σ -order may not grant the target slowdown in PS for all coflows once rate allocation is performed. However, the objective is to provide guarantees on the slowdown of coflows by using primal-feasibility as a soft constraint. To this respect the feasibility of PS can be regarded as a hard constraint on the slowdown. In all numerical experiments described in Sec. VII the fraction of violations never exceeded a few percents.

VI. ALGORITHMIC SOLUTION

This section introduces a primal-dual algorithm, namely CoFair, which can determine a primal-feasible σ -order in polynomial time. It generalizes the primal-dual scheduling approach [23] for job scheduling over uncorrelated machines, a technique later applied in the context of coflow scheduling in

Algorithm 2 Pseudocode of CoFair

```

1: Input:  $\mathcal{C}, \{p_{\ell,j}\}, \{D_j\}, \{\alpha_j\}$ 
2: Output:  $\sigma, \{y_{\mu_k, F_k}\}, \{C_j\}$ 
3:  $\{y_{\ell,S}\} = 0$  for  $S \subseteq \mathcal{C}$  % initialise the dual variables
4:  $A_N = \mathcal{C}$  % initial unscheduled coflows
5:  $w_j^{(N+1)} = w_j + \alpha_j$  % initial weights
6: while  $A_k \neq \emptyset$  do
7:    $\mu_k \leftarrow \text{pivot\_bottleneck}(A_k)$  % pivot bottleneck
8:    $F_k \leftarrow F(A_k)$  % set of feasible tail coflows in  $A_k$ 
9:   if  $\mathcal{F}_k = \emptyset$  then % Primal LP non-feasible
10:    return  $\emptyset$ 
11:   end if
12:    $\sigma(k) \leftarrow \arg \min_{j \in F_k} \left\{ \frac{w_j^{(k)}}{p_{\mu_k, j}} \right\}$  % pivot coflow (Smith rule)
13:    $w_{\sigma(k)}^{(k)} \leftarrow \frac{w_{\sigma(k)}^{(k+1)}}{p_{\mu_k, j}}$  % update pivot coflow weight
14:    $y_{\mu_k, F_k} \leftarrow w_{\sigma(k)}^{(k)}$  % update dual problem solution
15:   for  $j \in A_k \setminus \{\sigma(k)\}$  do
16:     if  $j \in \mathcal{F}_k$  then
17:        $w_j^{(k)} \leftarrow w_j^{(k+1)} - w_{\sigma(k)}^{(k+1)} \frac{p_{\mu_k, j}}{p_{\mu_k, \sigma(k)}}$  % update weights
18:     else
19:        $w_j^{(k)} \leftarrow w_j^{(k+1)}$  % non tail-feasible coflows: unchanged
20:     end if
21:   end for
22:    $A_{k-1} \leftarrow A_k \setminus \{\sigma(k)\}$  % update unscheduled coflow set
23: end while
24:  $C_j := \max_{\ell \in \mathcal{L}} \sum_{h=1}^j p_{\ell, \sigma(h)}$ ,  $j \in \mathcal{C}$  % CCT lower bound
25: return  $\sigma, \{y_{\mu_k, F_k}\}, \{C_k\}$ 

```

[8] and [13]. The basic idea in [23] is to utilize the Smith rule over different machines and sort jobs accordingly: machines are mapped into ports, jobs into coflows and tasks into flows. Every port is treated as uncorrelated to the others and a primal-dual iteration performs a coflow selection and weight-adjustment technique [8], [23]

Nevertheless, it is not obvious whether or not it is possible to perform a similar primal-dual iteration while accounting for slowdown constraints and ultimately obtain a primal-feasible schedule. The rest of the section shows how to generate a primal-feasible σ -order if one exists. The pseudocode of the algorithm described in this section requires specific notation which are introduced next.

For the set of coflows $A \subseteq \mathcal{C}$, let $V_\ell(A) = \sum v_j^i \chi_j^i(\ell)$ the total volume transmitted over port ℓ ; if $A = \mathcal{C}$ the notation is simply V_ℓ . The set of ports engaged by coflows in A is $\mathcal{L}(A) = \{\ell \in \mathcal{L} | V_\ell(A) > 0\}$. $\mathcal{C}_\ell(A)$ is the set of coflows in A active over port ℓ . $T_\ell(A) = \frac{1}{B_\ell} V_\ell(A)$ is the minimum time required to complete the last coflow of set A over link ℓ . Coflow j is a feasible *tail coflow* of set A if $T_\ell \leq D_j$ for each link ℓ where $p_{\ell,j} > 0$. Formally, $F(A) = \{j \in \mathcal{C}_\ell(A) | T_\ell(A) \leq D_j, \forall \ell \in \mathcal{L}(A)\}$ is the set of feasible *tail coflows* of set A . With no loss of generality, in the rest of this section $B_\ell = 1$.

The CoFair pseudocode in Alg. 2 describes the primal-dual iterations to generate a σ -order which is primal-feasible. At each step k it tries to identify a *pivot* bottleneck-coflow pair (line 7 and 12). When this is not possible, the algorithm exits with returning a non-feasibility result (line 10).

Otherwise, first, a target bottleneck port μ_k is returned for the unscheduled coflows $A_k = \{\sigma(1), \sigma(2), \dots, \sigma(k)\}$ using a `pivot_bottleneck` procedure (line 7). Second, it identifies the coflow $\sigma(k)$ to be scheduled *last* over such port in the set of feasible tail coflows F_k . On the selected bottleneck, a feasible tail coflow is chosen according to the Smith rule [31]. Weight updates are only performed for the pivot coflow (line 8) and for feasible tail coflows (line 12). Afterwards, the algorithm eliminates the selected coflow from the set of unscheduled ones (line 17). Weights α_j have the meaning of multipliers and can be used in order to enforce different σ -orders: as it will be showed in the next section using such weights and `pivot_bottleneck` it is possible to attain all possible primal-feasible σ -orders.

The `pivot_bottleneck` procedure can be fully general, but in the numerical tests the selection is based on the heuristics that selects the most charged bottleneck. This is the baseline rule adopted in [8], [13], [23]. Direct calculations show that the computational complexity of CoFair is $O(N(M + N))$.

A. Correctness

We prove first that CoFair terminates in N steps: at each iteration it is possible to find at least one feasible tail coflow to be selected at step (8).

Lemma 1. *If Primal LP is feasible, then $F_k \neq \emptyset$ for $k = 1, \dots, N$.*

Proof: For $k = N$ the statement is true: $A_N = \mathcal{C} \neq \emptyset$ and $F_N \neq \emptyset$, otherwise Primal LP would not be feasible. Hence, Cor. 2 ensures that the σ -order which sorts coflows in order of decreasing values of R_j^0 is feasible. From the proof of Thm. 5 such schedule is feasible for any $A_k \subseteq \mathcal{C}$ so that $F_k \neq \emptyset$ for all $k = 1, \dots, N$. ■

From the proof of Lemma 1, we can observe two facts: i) the statement is true irrespective of the implementation of the `pivot_bottleneck` procedure; ii) it holds also true if we replaced the Smith selection rule with any other selection rule of coflows within set \mathcal{F}_k as well.

We now show that the output of the algorithm belongs to the set of solutions of Primal LP. Furthermore, the approximation properties of the proposed algorithmic solution rely on the dual formulation

$$\text{maximize: } \sum_{\ell \in \mathcal{L}} \sum_{A \subseteq \mathcal{C}} y_{\ell, A} f_{\ell}(A) - \sum_j \alpha_j D_j \quad (\text{Dual LP})$$

$$\text{subj. to: } \sum_{\ell \in \mathcal{L}} \sum_{A \subseteq \mathcal{C}: j \in S} y_{\ell, A} p_{\ell, j} = w_j + \alpha_j, \quad j \in \mathcal{C} \quad (20)$$

$$y_{\ell, A} \geq 0, \alpha_j \geq 0 \quad (21)$$

If Primal LP is feasible, the output of CoFair provides one solution for Primal LP and one solution for Dual LP.

Finally, the following result demonstrates that the primal-dual iteration [23] can be extended to account for primal feasibility, i.e., the output of CoFair is a primal-feasible σ -order scheduler if one exists. Otherwise, the algorithm provides a test of feasibility for Primal LP.

Theorem 6. *Let Primal LP be feasible. CoFair produces a primal-feasible solution $\{C_k\}$ for Primal LP and a feasible solution $\{y_{\ell, A}\}$ for Dual LP. Let Primal LP be not feasible, CoFair returns \emptyset .*

B. Output completeness

Under feasibility conditions, the algorithm can always produce a solution corresponding to a dual solution for $\alpha = 0$. Nevertheless, it is possible to use the multipliers of the dual problem in order to obtain any target primal-feasible σ -order.

Theorem 7. *For every primal-feasible σ^* , there exist multipliers $\{\alpha_j\}$, $\max \alpha_j = 1$ so that σ^* is the output of Alg. 2 under rescaled weights $\{\kappa \cdot w_j\}$, for some $0 \leq \kappa \leq 1$.*

C. Approximation factor

Let C_j be the output of Alg. 2. Let define C_j^{OPT} , $j \in \mathcal{C}$ a solution of PS and C_j^{LP} an solution of Primal LP. We also define \hat{C}_j the coflow completion times of a σ -order using a priority output of Alg. 2. C_j^{OPT} is defined accordingly. Let recall a few results:

- Lemma 2.** i. $\sum_{k=1}^N w_k C_k^{\text{LP}} \leq \sum_{k=1}^N w_k C_k^{\text{OPT}}$
ii. $\left(\sum_{j=1}^k p_{\mu_k \sigma(j)} \right)_2 \leq f_{\ell}(S)$, for $\ell \in \mathcal{L}$ and $S \subseteq \mathcal{C}$ [23]
Lemma 3.2
iii. $\hat{C}_j - r_j \leq 2(C_j - r_j)$ [8] Lemma 3

Theorem 8. *Alg. 2 produces a primal-feasible schedule such that for the corresponding σ -order coflow schedule σ^**

$$\sum w_j \hat{C}_j \leq 4 \sum w_j C_j^{\text{OPT}} + 4 \sum \alpha_j D_j \quad (22)$$

Assume $\alpha_j = 0$ and σ^* is feasible for PS, then it is a 4 approximation w.r.t. an optimal scheduler.

D. Smith rule is Gaussian elimination

The proof of Thm. 6 described in the Appendix offers a simple interpretation from basic linear algebra of the fact that the Smith rule is always leading to a near-optimal, primal-feasible σ -order. Ultimately, the iterative selection of a pivot bottleneck-coflow pair and the weight update operated according to the Smith rule can be seen as simultaneous Gaussian elimination and selection of a $N \times N$ sub-matrix over a suitable rectangular matrix of size $N \times M 2^N$ corresponding to the constraints of the dual problem. In turn the existence of a coflow to be selected at each step is granted by the feasibility of Primal LP, e.g., when using the output of MPS as input slowdown parameter E . This permits to extend the proof of the result in [23][Alg. 3.1] because under feasible slowdown constraints there is no need to update the weights of non tail-feasible coflows at each step. Furthermore, if one gives up the selection of the most charged bottleneck in the `pivot_bottleneck` selection and Smith rule, the output is still primal-feasible, but Thm. 8 may not hold.

VII. NUMERICAL RESULTS

This section reports on a set of extensive numerical experiments validating the proposed coflow fairness framework and the related algorithms. The numerical have been performed on a matlab[®] coflow scheduling simulator equipped with a coflow generator². For the algorithms based on σ -order, namely Sincronia and CoFair, the rate allocation is based on the greedy procedure. It assigns full priority to coflow k over coflow h on all ports where they are both active if $\sigma(k) \leq \sigma(k+1)$. This rate allocation has the advantage of simplicity of implementation, at the price of some performance loss against, for instance, the rate allocation adopted in Utopia which permits some rate back-filling for coflows in lower priority.

Sample Coflow batches. As indicated in the analysis of Facebook traces reported in [5], a typical coflow pattern observed in data centers may comprise a significant fraction of small coflows but also coflows with large *width* (i.e., number of flows). These type of instances are denoted *wide-narrow* coflows batches (WN). The first set of coflow batches considered has a fraction q of coflows whose width is drawn uniformly from $M/3, \dots, M$ (M being the number of ports) and a fraction $1 - q$ of coflows containing single flows. The second set of coflow batches are the Map-Reduce type (MR). They are classified according to the number of mappers and the number of reducers, which are drawn from a uniform distribution in $[1, m]$ and $[1, r]$, respectively. Each reducer fetches data from each mapper. For all coflow samples, individual flow volumes are exponentially distributed with average 10 and standard deviation 3 in normalized traffic volume units. Links are normalized with capacity 1 traffic volume units per second.

Performance and fairness metrics. The metrics employed for numerical validation are the CCT, the experimental slowdown and the stretch index (SI). The experimental slowdown for a tagged coflow j is $\hat{E}_j = \phi_j C_j / C_j^0$: it embraces the plain slowdown for $\phi_j = 1$ and the port-occupation one for $\phi = V_j$. Finally, for a tagged coflow j , the stretch index (SI), i.e., $SI := \sum \max(0, (\hat{E}_j / E) - 1)$, quantifies the relative magnitude of violations w.r.t. the target slowdown constraint parameter E .

A. Minimum slowdown estimation.

As described in Sec. IV, the minimum slowdown E^* corresponds to the largest average progress that can be ensured to all the coflows in a batch. Fig. 3a and 3b reported on the relative error which is obtained by approximating E^* using the value attained for the primal value E^P generated as output of MPS. Each point of those graphs is obtained by averaging the relative error obtained over 100 batch instances, where the LP solution corresponds to a time slot of 1 second.

Fig. 3a reports on the approximation results in the case of increasing values of q , $M = N = 10$, whereas Fig. 3b

²The simulator is Open Source and it will be made available on GitHub; to obtain the source code of the simulator the reader is encouraged to contact the authors of the manuscript.

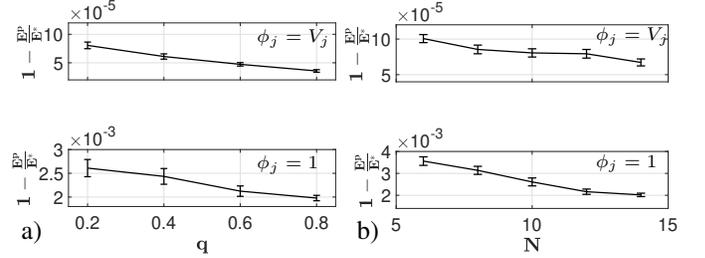


Fig. 3. Relative error on the minimum slowdown for $N = M = 10$: a) increasing q b) increasing N and increasing N for fixed $q = 0.2$.

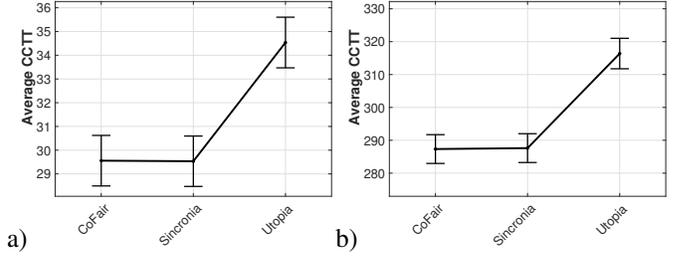


Fig. 4. Average CCT normalized against Sincronia for $M = 30$: a) $p = 0.2$ and $N = 30$ and b) $p = 0.8$ and $N = 100$; 95% confidence intervals are superimposed.

reports on the same experiment for increasing values of N . In both cases, the approximation improves for larger values of q and N , respectively. The important information is that the approximation improves the higher the fabric congestion, and it appears very tight, i.e., the relative error is bounded below 10^{-4} for $\phi = V_j$ and 3×10^{-3} for $\phi = 1$. These results confirm that the minimum slowdown of the primal problem E^P appears an extremely tight approximation compared against the exact value E^* obtained by solving the corresponding LP. For each set of 100 experiments we observed at most 1 outlier, i.e., a batch for which the relative error falls above 1%. By inspection, we found that this happens for very specific cases where a coflow batch has several bottlenecks of same size, and where a primal-feasible schedule may prioritize certain coflows on bottlenecks where in turn they slow down others, thus increasing the maximum slowdown value for that batch.

B. Plain slowdown: $\phi_j = 1$.

This set of experiments considers the plain slowdown, i.e., $\phi_j = 1$. With this choice, between two coflows with same volume, it is fair if the coflow with the lowest CCT in isolation finishes first. In each experiment, same coflow instances are processed using Sincronia, CoFair and Utopia. In all experiments, CoFair receives E^* as output of MPS. For the plain slowdown, the Jain index $J(\mathbf{R}) = \frac{1}{K} \frac{\sum_{j=1}^K R_j^2}{(\sum_{j=1}^K R_j)^2}$ can measure how fair is the progress distribution among coflows, i.e., the average rate R_j .

Fig. 4a, Fig. 5a and Fig. 5b refer to the performance figures measured for an experiment with 100 batches of coflows of NW type for $p = 0.2$ and $N = 30$. Performance figures reported in Fig. 4b and Fig. 5c and d refer to $p = 0.8$ and $N = 100$.

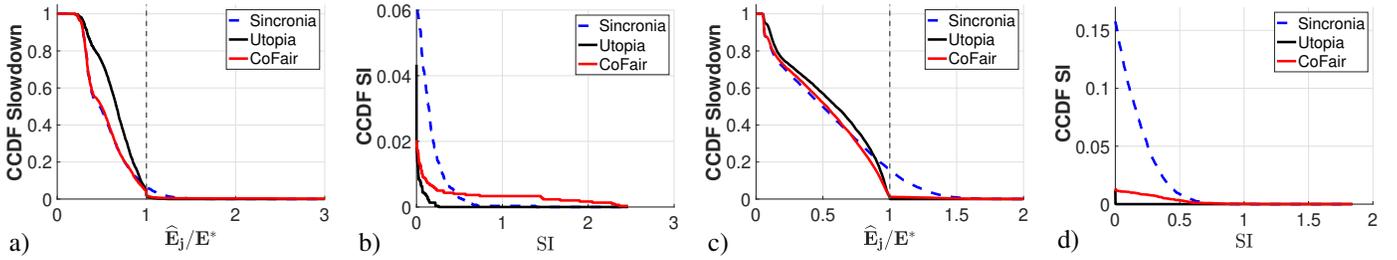


Fig. 5. WN traces: a) and c) CCDF of the normalized slowdown; b) and d) CCDF of the stretch index; setting: $M = 30$, $q = 0.2$ and $N = 30$ for a) and b); $q = 0.8$ and $N = 100$ for c) and d).

Fig. 4 reports on the average CCT attained by Sincronia, CoFair and Utopia: to this respect, the average CCT of CoFair and Sincronia coincide. On the other hand, the performance loss degradation of Utopia w.r.t. Sincronia and CoFair is significant, i.e., on the order of 17% for $p = 0.2$ and 9% for $p = 0.8$. For $p = 0.2$ the Jain index of CoFair is 0.60, for Sincronia is 0.59 whereas Utopia only scores 0.53. Conversely, for $p = 0.8$ the Jain index of CoFair is 0.61, for Sincronia is 0.60 whereas Utopia scores 0.58.

Fig. 5 illustrates the complementary cumulative distribution function (CCDF) of the slowdown and of the stretch index. The dashed vertical line seen in Fig. 5a and Fig. 5c separates the *inner interval* on the left where $\hat{E}_j \leq E^*$, i.e., no violation occurs, and the *outer interval* $\hat{E}_j > E^*$, i.e., where the slowdown does not meet the target. As observed in Fig. 5a, Sincronia and CoFair provide the best performance for $\hat{E}_j \leq E^*$. However, as expected, it suffers a larger number of violations compared to Utopia. Finally, CoFair matches the behavior of Sincronia for $\hat{E}_j \leq E^*$ and Utopia for $\hat{E}_j \geq E^*$, even though the rare violations, when they occur, appear relatively larger. Same behaviour is confirmed for $p = 0.8$ in Fig. 4a, Fig. 5a and Fig. 5b. In this case, due to the increased congestion, the violations of Sincronia are larger than in the previous experiments, and yet CoFair balances between a near optimal CCT minimization and slowdown guarantees.

The same experiments are repeated for MR coflow batches under different configurations on the number of mappers and reducers. These coflow batches generate severe congestion both at the input and at the output ports of the fabric. Fig. 6a and Fig. 6b refer to the case of 100 batches of coflows with $N = 30$ MR coflows where $m = 10$ and $r = 3$. As before, Sincronia and CoFair provide the best performance in the inner interval. However, Sincronia has a significant number of deviations, i.e., on the order of 10% for coflows, whereas Utopia provides a better performance in terms of deviation. This is confirmed also by the Jain index, where CoFair scores 0.89, Sincronia 0.88 and Utopia 0.87. Ultimately, CoFair matches Sincronia in the inner interval and Utopia in the outer one, providing a better tradeoff.

The experiment is hence performed for $m = 10$ and $r = 10$ and $N = 100$ coflows. Again, the slowdown CCDF of CoFair is practically superimposed to that of Sincronia in the inner interval and incurs very few violations ($< 10^{-3}$). With respect to the Jain index for the coflow progress, CoFair and Sincronia

score 0.84 whereas Utopia 0.80.

C. Slowdown with port occupation: $\phi_j = V_j$.

By letting $\phi_j = V_j$, the average port occupation is embedded in the slowdown definition: between two coflows with the same volume, it is fair if the one with lower average port occupation finishes first. Fig. 7a and Fig. 7b describe the effect of the slowdown constraint E onto the average CCT for a fabric with $M = 30$ ports and $N = 30, 50, 100$ coflows of type WN for $p = 0.2$. Results are averaged on hundred sample batches and normalized against the average CCT of the near optimal solution provided by Sincronia; 95% confidence intervals are superimposed. The average CCT decreases for increasing values of E : this behavior mimics the α parameter of classic flow α -fairness. As depicted in Fig. 7a, E can be tuned in order to attain a tradeoff between coflow slowdown and average CCT, where the loss in CCT gain over Sincronia tops at 40% for $p = 0.2$ and $N = 100$: this represents the loss in performance in order to grant all coflows the target slowdown. Same result is described in figure Fig. 4d for $p = 0.8$: in presence of many large coflows, the range of the tradeoff is smaller.

Fig. 8a/c and b/d report on the CCDF of the normalized slowdown and the corresponding stretch index, respectively, for two experiments on 100 coflow batches. By comparing the distributions of the different algorithms for $p = 0.2$, i.e., Fig. 5a and Fig. 5b, it is apparent that CoFair provides the best slowdown-efficiency tradeoff, whereas both baselines do not perform well both in terms of normalized slowdown and in terms of stretch index. The same behavior is observed for $p = 0.8$: CoFair and Utopia perform similar in the inner interval, but the slowdown CCDF of CoFair bends clearly in proximity of E^* thus attaining a negligible number of violations. The tail of the SI distribution shows that their magnitude is relatively larger than those of Utopia in the same region. In these experiments, both Utopia and Sincronia have a rate of violations on the order of 10 – 20%.

The same experiment is repeated for the MR batches in Fig. 9. As in the previous cases, Sincronia and CoFair provide the best performance in terms of slowdown. In this case, not only Sincronia but also Utopia has a significant number of deviations, i.e., on the order of 10% for coflows in the inner interval. CoFair provides the best tradeoff, outperforming Utopia in the outer interval. The experiment is hence performed for

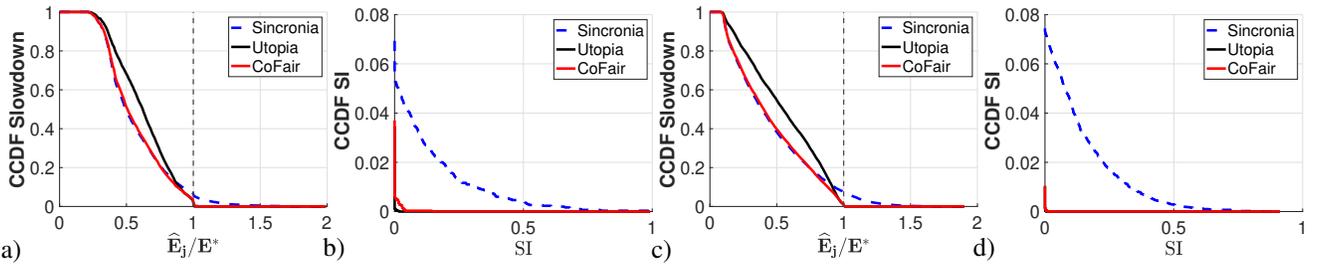


Fig. 6. MR traces: a) and c) CCDF of the normalized slowdown; b) and d) CCDF of the stretch index; setting: $M = 30$, and $N = 30$, $m = 10$ and $r = 3$ for a) and b); $M = 30$, and $N = 100$, $m = 10$ and $r = 10$ $q = 0.8$ and $N = 100$ for c) and d).

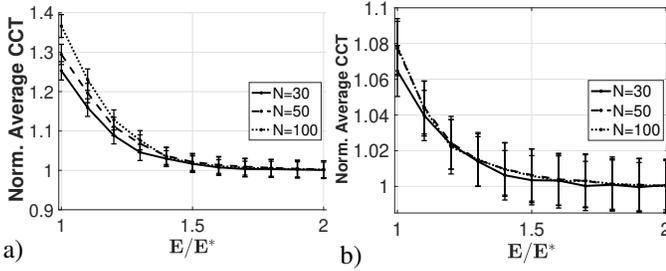


Fig. 7. Average CCT normalized against Sincronia for $M = 30$: a) $p = 0.2$ and $N = 30$; b) $p = 0.8$ and $N = 100$.

$m = 10$ and $r = 10$ and $N = 100$ coflows. In this case, the CCDF tail of CoFair in the outer interval is negligible, whereas Sincronia and Utopia show deviations larger than 15%.

Finally, in order to provide better understanding on the behavior of the algorithms, Fig. 10a and Fig. 10b show the scatter diagram of the normalized CCT and the slowdown for each coflow corresponding to the experiments of Fig. 8a and Fig. 8b. As seen in Fig. 10, in order to reduce the slowdown, CoFair tends to concentrate several coflows in the region where the slowdown is smaller at the price of higher CCT. In order to remark the tradeoff efficiency-fairness, the results for the runs of CoFair for $E = 1.1E^*$ have been added. As seen before in Fig. 5, the performance in terms of slowdown is very similar to Utopia. However, here parameter E acts as a fairness-efficiency knob as depicted in Fig. 9b. In fact, while matching the same CCT of Utopia, CoFair attains a significantly smaller number of violations. Fig. 9c and d report on the settings of Fig. 5c and Fig. 5d. Here, by letting on $E = 1.4E^*$ CoFair overlaps the slowdown performance of Utopia, and yet it provides average CCT figures very close to that of Sincronia.

VIII. CONCLUSIONS AND DISCUSSION.

The data transfer phase of modern computing frameworks requires to serve several concurrent coflows. This paper introduces a framework for coflow scheduling to trade off the average CCT for coflows' slowdown. It requires a single control parameter, that is the maximum slowdown for an input batch of coflows. Scheduling under slowdown constraints has larger feasibility set compared to the popular notion of progress based on minimum rate guarantees. The maximum

possible progress, i.e., the minimum possible slowdown, can be calculated with great accuracy in $O(MN + N \log(N))$.

A new scheduler, i.e., CoFair, has extended the framework of primal-dual σ -order schedulers to account for generalized slowdown constraints. The experimental results indicate that slowdown-aware σ -order generated by CoFair can adjust the priority of coflows with no apparent loss in performance with respect to near optimal benchmarks for average CCT minimization. Finally, based on the notion of generalized slowdown, CoFair and the proposed framework can perform coflow scheduling by accounting for specific features, e.g., the occupancy of the fabric ports.

REFERENCES

- [1] M. Chowdhury and I. Stoica, "Coflow: A networking abstraction for cluster applications," in *Proc. of ACM HotNets*, Redmond, Washington, 2012, p. 31–36.
- [2] J. Dean and S. Ghemawat, "Mapreduce: Simplified data processing on large clusters," in *Proc. of USENIX OSDI*, San Francisco, CA, 2004, pp. 137–150.
- [3] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, I. Stoica et al., "Spark: Cluster computing with working sets," in *Proc. of USENIX HotCloud*, vol. 10, no. 10-10, p. 95, 2010.
- [4] M. Han and K. Daudjee, "Giraph unchained: Barrierless asynchronous parallel execution in pregel-like graph processing systems," in *Proc. of VLDB*, vol. 8, no. 9, pp. 950–961, 2015.
- [5] N. M. K. Chowdhury, "Coflow: A networking abstraction for distributed data-parallel applications," Ph.D. dissertation, University of California, Berkeley, 2015.
- [6] Y. Chen and J. Wu, "Joint coflow routing and scheduling in leaf-spine data centers," *Elsevier Journal of Parallel and Distributed Computing*, vol. 148, pp. 83–95, 2021.
- [7] M. Zaharia et al., "Near optimal coflow scheduling in networks," in *Proc. of ACM SPAA*, Phoenix, AZ, USA, June 22-24 2019, pp. 123–134.
- [8] S. Agarwal, S. Rajakrishnan, A. Narayan, R. Agarwal, D. Shmoys, and A. Vahdat, "Sincronia: Near-optimal network design for coflows," in *Proc. of ACM SIGCOMM*, 2018, pp. 16–29.
- [9] R. Mao, V. Aggarwal, and M. Chiang, "Stochastic non-preemptive co-flow scheduling with time-indexed relaxation," in *Proc. of IEEE INFOCOM, DCPPerf Workshop*, Honolulu, Hawaii, US, April 16 2018, pp. 385–390.
- [10] L. Wang, W. Wang, and B. Li, "Utopia: Near-optimal coflow scheduling with isolation guarantee," in *Proc. of IEEE INFOCOM*, Honolulu, Hawaii, US, April 16 2018, pp. 891–899.
- [11] M. Shafiee and J. Ghaderi, "Scheduling coflows in datacenter networks: Improved bound for total weighted completion time," *SIGMETRICS Perform. Eval. Rev.*, vol. 45, no. 1, p. 29–30, Jun. 2017.
- [12] M. Chowdhury, Y. Zhong, and I. Stoica, "Efficient coflow scheduling with Varys," in *Proc. of ACM SIGCOMM*, Chicago, Illinois, USA, 2014, p. 443–454.
- [13] S. Ahmadi, S. Khuller, M. Purohit, and S. Yang, "On scheduling coflows," *Algorithmica*, vol. 82, no. 12, p. 3604–3629, dec 2020. [Online]. Available: <https://doi.org/10.1007/s00453-020-00741-3>

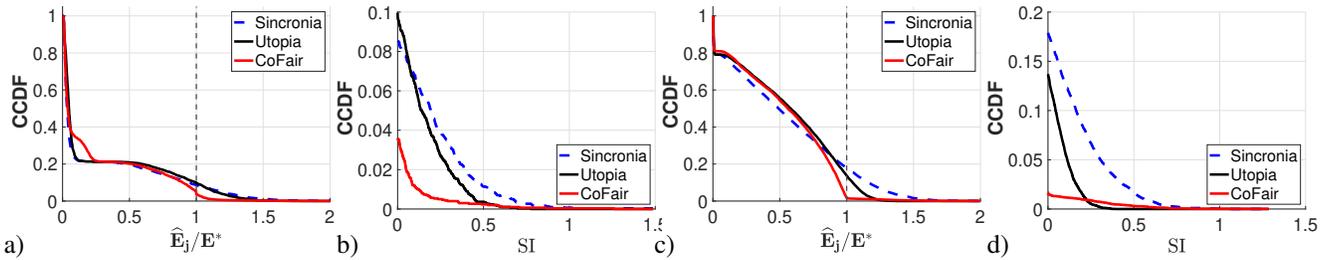


Fig. 8. WN traces, $\phi_j = V_j$: a) and c) CCDF of the normalized slowdown; b) and d) CCDF of the stretch index; setting: $M = 30$, $q = 0.2$ and $N = 30$ for a) and b); $q = 0.8$ and $N = 100$ for c) and d).

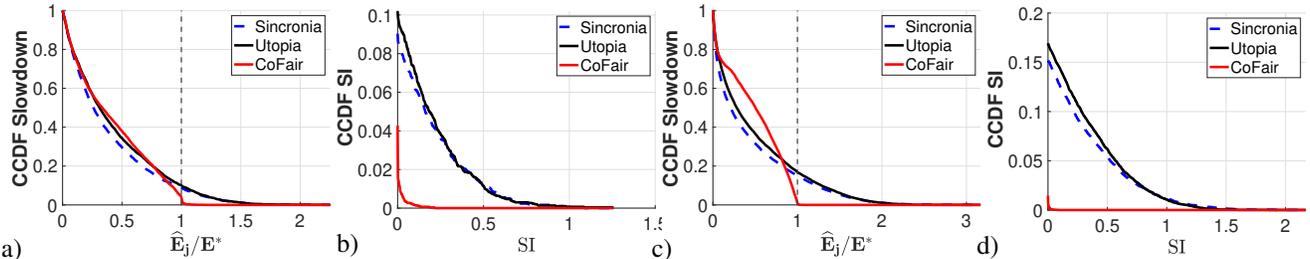


Fig. 9. MR traces, $\phi_j = V_j$: a) and c) CCDF of the normalized slowdown; b) and d) CCDF of the stretch index; setting: $M = 30$, and $N = 30$, $m = 10$ and $r = 3$ for a) and b); $M = 30$, and $N = 100$, $m = 10$ and $r = 10$ $q = 0.8$ and $N = 100$ for c) and d).

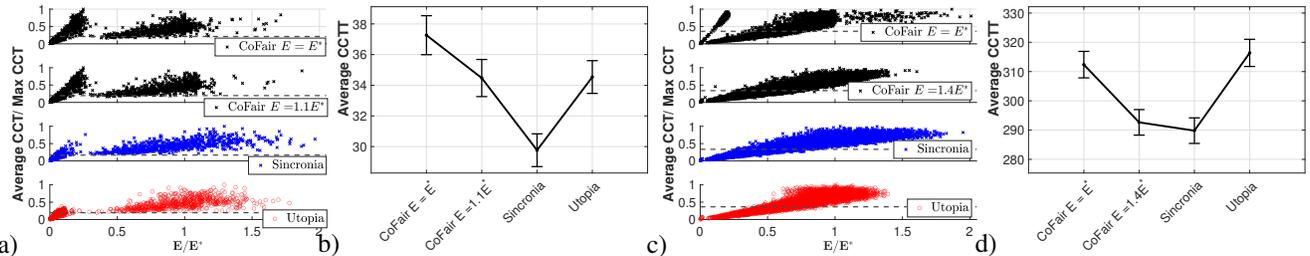


Fig. 10. WN traces, $\phi_j = V_j$: a) and c) scatter diagram and b) and d) corresponding average CCT; setting: $M = 30$, a) and b) $q = 0.2$ and $N = 30$ c) and d) $q = 0.8$ and $N = 100$.

[14] S. Shakkottai and R. Srikant, "Network optimization and control," *Foundations and Trends in Networking*, vol. 2, no. 3, pp. 271–379, 2008.

[15] J. Mo and J. C. Walrand, "Fair end-to-end window-based congestion control," *IEEE/ACM Trans. on Networking*, vol. 8, no. 5, pp. 556–567, 2000.

[16] M. Alizadeh, S. Yang, M. Sharif, S. Katti, N. McKeown, B. Prabhakar, and S. Shenker, "Pfabric: Minimal near-optimal datacenter transport," *SIGCOMM Comput. Commun. Rev.*, vol. 43, no. 4, p. 435–446, aug 2013. [Online]. Available: <https://doi.org/10.1145/2534169.2486031>

[17] L. Chen, W. Cui, B. Li, and B. Li, "Optimizing coflow completion times with utility max-min fairness," in *Proc. of IEEE INFOCOM*, April 2016, pp. 1–9.

[18] L. Chen, Y. Feng, B. Li, and B. Li, "Efficient performance-centric bandwidth allocation with fairness tradeoff," *IEEE TPDS*, vol. 29, no. 8, pp. 1693–1706, 2018.

[19] W. Wang, S. Ma, B. Li, and B. Li, "Coflex: Navigating the fairness-efficiency tradeoff for coflow scheduling," in *Proc. of IEEE INFOCOM*, Atlanta, GA, USA, 1–4 May 2017, pp. 1–9.

[20] M. Chowdhury, Z. Liu, A. Ghodsi, and I. Stoica, "HUG: Multi-resource fairness for correlated and elastic demands," in *Proc. of USENIX NSDI*, Santa Clara, CA, March 2016, pp. 407–424.

[21] L. Chen, Y. Feng, B. Li, and B. Li, "Efficient performance-centric bandwidth allocation with fairness tradeoff," *IEEE TPDS*, vol. 29, no. 8, pp. 1693–1706, 2018.

[22] A. Ghodsi, M. Zaharia, B. Hindman, A. Konwinski, S. Shenker, and I. Stoica, "Dominant resource fairness: Fair allocation of multiple resource types," in *Proc. of USENIX NSDI*, USA, 2011, p. 323–336.

[23] M. Mastrolilli, M. Queyranne, A. S. Schulz, O. Svensson, and N. A. Uhan, "Minimizing the sum of weighted completion times in a current open shop," *Operations Research Letters*, vol. 38, no. 5, pp. 390–395, 2010.

[24] T. Zhang, R. Shu, Z. Shan, and F. Ren, "Distributed bottleneck-aware coflow scheduling in data centers," *IEEE TPDS*, vol. 30, no. 7, pp. 1565–1579, July 2019.

[25] B. Radunovic and J.-Y. Le Boudec, "A unified framework for max-min and min-max fairness with applications," *IEEE/ACM Trans. on Networking*, vol. 15, pp. 1073 – 1083, 11 2007.

[26] H. Qu, R. Xu, W. Li, W. Qu, and X. Zhou, "OSTB: Optimizing fairness and efficiency for coflow scheduling without prior knowledge," in *Proc. of IEEE HPC*, 2019, pp. 1587–1594.

[27] S.-H. Tseng and A. Tang, "Coflow deadline scheduling via network-aware optimization," in *Proc. of the Annual Allerton Conference on Communications, Control and Computing*, 2018, pp. 829–833.

[28] Q. Luu, O. Brun, R. El Azouzi, F. De Pellegrini, B. J. Prabh, and C. Richier, "DCoflow: deadline-aware scheduling algorithm for coflows in datacenter networks," in *Proc. of IFIP Networking*, June 13–16 2022, pp. 1–9.

[29] W. Wang, S. Ma, B. Li, and B. Li, "Coflex: Navigating the fairness-efficiency tradeoff for coflow scheduling," in *Proc. of IEEE INFOCOM*, Atlanta, GA, USA, 1–4 May 2017, pp. 1–9.

[30] J. Cheriyan, R. Ravi, and M. Skutella, "A simple proof of the moore-hodgson algorithm for minimizing the number of late jobs," *Operations Research Letters*, vol. 49, no. 6, pp. 842–843, 2021.

[31] M. Queyranne, "Structure of a simple scheduling polyhedron," *Mathematical Programming*, vol. 58, pp. 263–285, 1993.

APPENDIX
PROOF OF THM. 2

Proof: The feasibility of the PS problem can be formulated with the following LP

$$\text{minimize: } 0 \quad (23)$$

$$\text{subj. to: } \sum_{t=1}^T x_j^i(t) = 1, \quad \forall j \in \mathcal{C}, i \in \mathcal{F}_j \quad (24)$$

$$Z_j(t) \leq 0, \forall t \leq r_j \quad (25)$$

$$Z_j(t) \geq 1, \forall t > \frac{E \cdot (C_j^0 - r_j)}{V_j} + r_j \quad (26)$$

$$\sum_{j \in \mathcal{C}} \sum_{i \in \mathcal{F}_i} v_j^i x_j^i(t) \chi_j^i(\ell) \leq B_\ell, \quad \forall t, \forall \ell \in \mathcal{L} \quad (27)$$

$$x_j^i(t) \in [0, 1], \quad \forall t \quad (28)$$

The above LP is obtained by considering any rate allocation for which the constraints of PS are respected: (26) has replaced the constraints (6) and (7) in PS. Since LP is in P this completes the proof of the first part of the statement. The second part statement follows by observing that from (26), $\frac{E \cdot (C_j^0 - r_j)}{V_j} + r_j \geq \frac{E \cdot C_j^0}{V_j}$ if and only if $E \geq V_j$ for all $j \in \mathcal{C}$, which holds true since $C_j \geq C_j^0$ for all $j \in \mathcal{C}$. ■

PROOF OF THM. 6

Proof: Let define a feasible solution for the Primal LP problem. For every pair $(\mu_k, \sigma(k))$ let define $C_{\sigma(k)} = \sum_{j \leq k} p_{\mu_k, \sigma(j)}$; since Primal LP is feasible, the slowdown constraint in the primal formulation is satisfied at each step as from Lemma 1. Parallel inequalities hold by applying portwise results in single machine scheduling [31].

For the Dual LP, let consider solutions of the type

$$y_{\ell, A} = \begin{cases} \theta_k > 0 & \ell = \mu_k \text{ and } A = F_k \text{ for some } k \\ 0 & \text{otherwise} \end{cases} \quad (29)$$

where μ_k is the bottleneck selected at step k and \mathcal{F}_k is the corresponding set of tail-feasible coflows. The corresponding form of the constraints of the dual problem

$$\sum_{k=j}^N q_{\mu_k, \sigma(j)} y_{\mu_k, F_k} = w_{\sigma(j)} + \alpha_{\sigma(j)}, \quad j \in \mathcal{C} \quad (30)$$

where $q_{\mu_k, j} = p_{\mu_k, j} \mathbf{1}\{j \in F_k\}$. If at each step $k = N, \dots, 2, 1$ the selection is

$$\sigma(k) = \operatorname{argmin}_{j \in F_k} \left\{ \frac{w_j^k}{p_{\mu_N, j}} \right\} \quad (31)$$

the algorithm produces a non-negative solution of the linear system (30). In fact, let first consider step $k = N$. The full rank system of the constraints of the dual system writes

$$\begin{pmatrix} q_{\mu_N \sigma(N)} & 0 & \dots & 0 & \left| \begin{array}{c} w_{\sigma(N)}^{(N+1)} \\ w_{\sigma(N)}^{(N+1)} \\ \vdots \\ w_{\sigma(N)}^{(N+1)} \end{array} \right. \\ q_{\mu_N \sigma(N-1)} & q_{\mu_{N-1} \sigma(N-1)} & \dots & 0 & \left| \begin{array}{c} w_{\sigma(N-1)}^{(N+1)} \\ \vdots \\ w_{\sigma(N-1)}^{(N+1)} \end{array} \right. \\ \vdots & \vdots & \vdots & \vdots & \left| \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right. \\ q_{\mu_N \sigma(1)} & q_{\mu_{N-1} \sigma(1)} & \dots & q_{\mu_1 \sigma(1)} & \left| \begin{array}{c} w_{\sigma(1)}^{(N+1)} \\ \vdots \\ w_{\sigma(1)}^{(N+1)} \end{array} \right. \end{pmatrix}$$

The Gaussian elimination using $\sigma(N)$ as pivot renders the new equivalent system

$$\begin{pmatrix} 1 & 0 & \dots & 0 & \left| \begin{array}{c} \frac{w_{\sigma(N)}^N}{p_{\mu_N, \sigma(N)}} \\ w_{\sigma(N-1)}^N - w_{\sigma(N)}^N \frac{p_{\mu_N \sigma(N-1)}}{p_{\mu_N \sigma(N)}} \\ \vdots \\ w_{\sigma(1)}^N - w_{\sigma(N)}^N \frac{p_{\mu_N \sigma(1)}}{p_{\mu_N \sigma(N)}} \end{array} \right. \end{pmatrix}$$

Observe that the first column is nullified irrespective of the chosen permutation $\sigma(j)$ for $j < N$. The system has been transformed in the equivalent one where Gaussian elimination can now be performed on the rightmost square sub-matrix. Let replace in the column on the right with the following values

$$w_j^{N-1} = \frac{w_{\sigma(N)}^N}{p_{\mu_N \sigma(N)}}, \quad j = \sigma(N)$$

$$w_j^{N-1} = w_j^N - w_{\sigma(N)}^N \frac{p_{\mu_N, j}}{p_{\mu_N \sigma(N)}} \geq 0, \quad j \in F_k \setminus \{\sigma(N)\}$$

$$w_j^{N-1} = w_j^N, \quad j \notin F_k$$

where $w_j^{N-1} \geq 0$ for all $j \in \mathcal{C}$. Indeed, $y_{\mu_N F_N} = \frac{w_{\sigma(N)}^N}{p_{\mu_N \sigma(N)}}$. The algorithm iterates the procedure on the remaining subsystem using the corner element of $\sigma(k)$ as pivot so that at each step $y_{\mu_k, F_k} = \theta_k = \frac{w_{\sigma(k)}^N}{p_{\mu_k \sigma(k)}} \geq 0$. ■

Note that unscheduled coflows which are non primal-feasible at a given step do not take part to gaussian elimination.

PROOF OF THM. 7

Proof: The sketch of the proof is as follows. First, by linearity rescaling weights $\{\kappa \cdot w_j\}$ is irrelevant. The Smith ratio at the k -th step for a candidate coflow

$$\frac{w_{\sigma(N-k)}^{(N-k)}}{P_{\mu_k \sigma(N-k)}} = \frac{w_{\sigma(N-k)}^{(N-k)} + \alpha_{\sigma(N-k)} - \sum_{h=1}^k w_{\sigma(N-k+h)}^{(N-k+h)} \frac{P_{\mu_h \sigma(N-k+h)}}{P_{\mu_h \sigma(N-k+h)}}}{P_{\mu_k \sigma(N-k)}}$$

Hence, to enforce the desired σ -order σ^* as output, let start by an input set $\alpha_j = 1$, for all j , and choose $\alpha_{\sigma^*(N-k)}$ such in a way that the Smith ratio of for $\sigma^*(N-k)$ be minimal among all candidate coflows, simply by rendering $(w_{\sigma(N-k)} + \alpha_{\sigma(N-k)})$ small enough as compared to the other coflows, possibly rescaling all weights w_j by κ . Using $\{\alpha_j\}$ as input to the algorithm produces the desired output σ^* . ■

PROOF OF THM. 8

Proof: For C_j , output of Alg. 2, it holds

$$\begin{aligned}
\sum_{j=1}^N w_j C_j &= \sum_{j=1}^N C_j \left(\sum_{\ell \in \mathcal{L}} \sum_{A \subseteq \mathcal{C}: j \in S} y_{\ell, A} p_{\ell, j} - \alpha_j \right) \\
&= \sum_{j=1}^N C_{\sigma(j)} \left(\sum_{k=j}^N p_{\mu_k, \sigma(j)} y_{\mu_k F_k} - \alpha_{\sigma(j)} \right) \\
&\stackrel{(i)}{\leq} \sum_{k=1}^N y_{\mu_k F_k} \sum_{j=1}^k p_{\mu_k, \sigma(j)} C_{\sigma(j)} - \sum_{j=1}^N \alpha_j C_j \\
&\stackrel{(ii)}{\leq} \sum_{k=1}^N y_{\mu_k F_k} \left(\sum_{j=1}^k p_{\mu_k, \sigma(j)} \right)^2 - \sum_{j=1}^N \alpha_j C_j \\
&\leq 2 \left(\sum_{k=1}^N y_{\mu_k F_k} f_{\mu_k}(F_k) - \alpha_k C_k \right) + \sum_{j=1}^N \alpha_j (2D_j - C_j) \\
&\stackrel{(iii)}{\leq} 2 \sum_{k=1}^N w_k C_k^{\text{LP}} + \sum_{j=1}^N \alpha_j (2D_j - C_j) \\
&\leq 2 \sum_{k=1}^N w_k C_k^{\text{OPT}} + 2 \sum_{j=1}^N \alpha_j D_j
\end{aligned}$$

Where the following arguments hold: (i) follows from the fact that in general $j \in F_k$ does not imply $j \in F_{k+1}$ and (ii) follows from the fact that $C_{\sigma(k)} = \sum_{j=1}^k p_{\mu_k, \sigma(j)} \geq C_{\sigma(j)}$ for all $k = 1, \dots, j$. According to Lemma. 2, it holds

$$\sum_{j=1}^N w_j \widehat{C}_j \leq 2 \sum_{j=1}^N w_j C_j \leq 4 \sum_{k=1}^N w_k C_k^{\text{OPT}} + 4 \sum_{j=1}^N \alpha_j D_j$$

which concludes the proof. \blacksquare