

Sincronia: Near-Optimal Network Design for Coflows

Saksham Agarwal*
Cornell University

Shijin Rajakrishnan*
Cornell University

Akshay Narayan
MIT

Rachit Agarwal
Cornell University

David Shmoys
Cornell University

Amin Vahdat
Google

ABSTRACT

We present Sincronia, a near-optimal network design for coflows that can be implemented on top on any transport layer (for flows) that supports priority scheduling. Sincronia achieves this using a key technical result — we show that given a “right” ordering of coflows, any per-flow rate allocation mechanism achieves average coflow completion time within $4\times$ of the optimal as long as (co)flows are prioritized with respect to the ordering.

Sincronia uses a simple greedy mechanism to periodically order all unfinished coflows; each host sets priorities for its flows using corresponding coflow order and offloads the flow scheduling and rate allocation to the underlying priority-enabled transport layer. We evaluate Sincronia over a real testbed comprising 16-servers and commodity switches, and using simulations across a variety of workloads. Evaluation results suggest that Sincronia not only admits a practical, near-optimal design but also improves upon state-of-the-art network designs for coflows (sometimes by as much as $8\times$).

CCS CONCEPTS

• **Networks** → **Network protocol design**; • **Theory of computation** → *Scheduling algorithms*;

KEYWORDS

Coflow, Datacenter Networks, Approximation Algorithms

ACM Reference Format:

Saksham Agarwal, Shijin Rajakrishnan, Akshay Narayan, Rachit Agarwal, David Shmoys, and Amin Vahdat. 2018. Sincronia: Near-Optimal Network Design for Coflows. In *SIGCOMM '18: ACM SIGCOMM 2018 Conference, August 20–25, 2018, Budapest, Hungary*. ACM, New York, NY, USA, 19 pages. <https://doi.org/10.1145/3230543.3230569>

1 INTRODUCTION

Traditionally, networks have used the abstraction of a “flow”, that captures a sequence of packets between a single source and a single destination. This abstraction has been a mainstay for decades and for a good reason — network designs were optimized for latency and/or throughput for a point-to-point connection, precisely the performance metrics important to traditional applications (e.g., file transfers, web access, etc.). However, distributed applications running across datacenter networks use programming models (e.g., bulk synchronous programming and partition-aggregate model) that require optimizing performance for a collection of flows rather than individual flows. The network still optimizes the performance of individual flows, leading to a fundamental mismatch between performance objectives of applications and the optimization objectives of network designs.

The coflow abstraction [7] mitigates this mismatch, allowing distributed applications to more precisely express their performance objectives to the network fabric. For instance, many distributed services with stringent performance constraints must essentially block until receiving all or almost all responses from hundreds or even thousands of remote servers (§2). Such services can specify a collection of flows as a coflow. The network fabric now optimizes for average Coflow Completion Time (CCT) [7, 10, 12], where the CCT of a coflow is defined as the time when some percentage, perhaps 100%, of flows in the coflow finish. Several recent evaluations show that optimizing for average CCT can significantly improve application-level performance [7, 10, 12].

There has been tremendous recent effort on network designs for coflows, both in networking [7–10, 12] and in the theory community [2, 20, 23]. However, prior designs require a centralized coordinator to perform complex per-flow rate allocations, with rate allocated to a flow being dependent on the rate allocated to other flows in the network. Such

*The first two authors contributed equally to the paper.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGCOMM '18, August 20–25, 2018, Budapest, Hungary
© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-5567-4/18/08...\$15.00
<https://doi.org/10.1145/3230543.3230569>

centralized inter-dependent per-flow rate allocation make it hard to realize these designs in practice for several reasons. First, per-flow rate allocation naturally requires knowledge about location of congestion in the network and paths taken by each (co)flow, making it hard to use these designs when congestion is in the fabric core and/or changes dynamically. Second, since rates allocated to flows are correlated, arrival or departure of even one coflow may result in reallocation of rate for each and every flow in the network. Such reallocation is impractical in datacenters where thousands of coflows may arrive each second. As a result, a practical near-optimal network design for coflows still remains elusive.

This paper presents Sincronia, a new datacenter network design for coflows that achieves near-optimal average CCT without any explicit per-flow rate allocation mechanism. The high-level design of Sincronia can be summarized as:

- Time is divided into epochs;
- In each epoch, a subset of unfinished coflows are selected and “ordered” using a simple greedy algorithm;
- Each host independently sets a priority for its flows (based on the corresponding coflow’s ordering), and offloads the flow to underlying priority-enabled transport mechanism;
- Coflows that arrive between epoch boundaries are greedily scheduled for work conservation.

Sincronia’s minimalistic design is based on a key technical result — with a “right” ordering of coflows, it is possible to achieve average CCT within $4\times$ of the optimal¹ as long as (co)flow scheduling is “order-preserving” — if coflow C is ordered higher than coflow C' , flows/packets in C must be prioritized over those in C' . From a practical perspective, this result is interesting because it shows that as long as (co)flow scheduling is order-preserving, *any* per-flow rate allocation mechanism results in average CCT within $4\times$ of the optimal. Using this result, Sincronia admits a practical, near-optimal network design for coflows (§3) — a simple greedy algorithm periodically orders the set of unfinished coflows; each host, without any explicit coordination with other hosts, sets priorities for its flows based on corresponding coflow’s ordering, and offloads scheduling of and rate allocation to individual flows to the underlying priority-enabled transport layer.

Sincronia thus overcomes the aforementioned practical challenges in existing network designs for coflows by avoiding per-flow rate allocation and by being agnostic to the underlying transport layer. First, by avoiding per-flow rate

allocation, Sincronia design is independent of underlying network topology, location of congestion in the fabric, and paths taken by each (co)flow. This also allows Sincronia to transparently respond to network failures. Second, coflow arrivals and departures do not require explicit rate reallocation for existing flows, leading to a much more scalable design. Third, Sincronia design using simple priority mechanisms enables coexistence of flows and coflows (§4), supporting backward compatibility. Finally, by being transport-agnostic, Sincronia admits efficient implementation on top of any existing transport mechanism that supports priority scheduling including TCP (using DiffServ [6] for priority scheduling), pHost [13] and pFabric [3]. We discuss implementation of Sincronia on top of several transport mechanisms in §4.3.

We have implemented Sincronia on top of TCP, with DiffServ [6] for priority scheduling. Our implementation is work conserving, efficiently handles online arrival of coflows, and allows coexistence of flows and coflows. We evaluate Sincronia implementation on a 16-server testbed interconnected with a FatTree topology comprising 20 commodity switches. We have also incorporated Sincronia into existing coflow simulators; we use these to perform sensitivity analysis of Sincronia performance against variety of workloads, number of coflows, network load, transport mechanisms, etc. Our implementation and simulation results show that Sincronia not only provides near-optimal average CCT but also outperforms state-of-the-art network designs for coflows across all evaluated workloads (sometimes by as much as $8\times$).

2 SINCRONIA OVERVIEW

In this section, we briefly recall the coflow abstraction (§2.1) and formally define our optimization objective for network design for coflows (§2.2). We also review some of the known results in coflow scheduling in §2.2.

2.1 The Coflow Abstraction

Existing distributed programming frameworks [7, 9, 11, 18, 21, 22, 29, 30] often have a communication stage that is structured and takes place between successive computation stages. In these frameworks, execution of a task (or even an entire computation stage) cannot begin until all flows in the preceding communication stage have finished. A coflow [7, 10] is a collection of such flows, with a shared performance goal (e.g., minimizing the completion time of the last flow in a coflow). Figure 1 shows an example.

We assume that coflows are defined such that flows within a coflow are independent; that is, the input of a flow does not depend on the output of another flow within that coflow. Similar to most existing designs [10, 12, 14, 20], we focus on a clairvoyant design that assumes information about a coflow (set of flows, and corresponding sources, destinations and sizes) is known at coflow’s arrival time but no earlier.

¹Under the standard assumption that the fabric core can sustain 100% throughput. That is, only the ingress and the egress access links are potential bottlenecks. This is the same assumption made in the big switch model for traditional abstraction of flows [3, 13, 16]. For the case when the fabric core cannot sustain 100% throughput, no approximation algorithm is known even for the traditional abstraction of flows. While we use this assumption for our theoretical bounds, our implementation makes no such assumption and adapts well to in-network congestion.

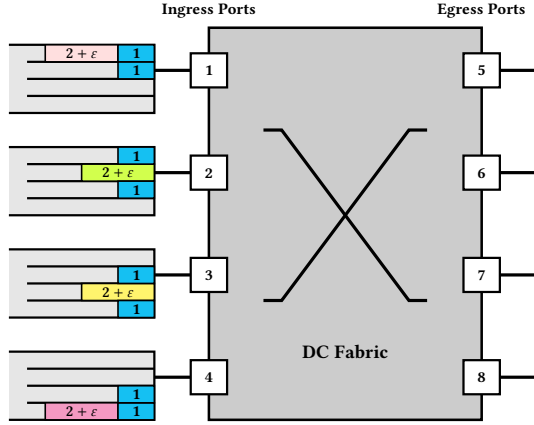


Figure 1: An instance of a coflow scheduling problem, used as a running example in the paper. The datacenter has 4 ingress and egress ports, with each ingress port having a “virtual output queue” for each egress port (see §2.2 for detailed model description). The example has 5 coflows. Coflow C1 has eight flows, with each ingress port sending unit amount of data to two egress ports; coflows C2, C3, C4 and C5 have one flow each sending $2 + \epsilon$ amount of data to one of the egress ports (the ϵ amount is only for breaking ties consistently). Coflow C2, C3, C4 and C5 being single flow coflow is only for simplicity; the ϵ amount of flow could be sent to any egress port without changing the results in §3.

2.2 Problem Statement and Prior Results

We now describe the network model used for our theoretical bounds, and the network performance objective.

Conceptual Model (for theoretical bounds). Similar to near-optimal network designs for traditional abstraction of network flows [3, 13, 16] and coflows [2, 10, 12, 14, 20, 23], we will abstract out the datacenter network fabric as one big switch that interconnects the servers. In such a big switch model, the ingress queues correspond to the NICs and the egress queues to the last-hop TOR switches. The model assumes that the fabric core can sustain 100% throughput and only the ingress and egress queues are potential congestion points. Under this model, each ingress port has flows from one or more coflows to various egress ports (see Figure 1). For ease of exposition, we organize the flows in virtual output queues at the ingress ports. We use this abstraction to simplify our theoretical analysis and algorithmic description, but we do not enforce it in our design and experiments.

Performance Objective. Formally, we assume that the network is a big switch comprising m ingress ports $\{1, 2, \dots, m\}$ and m egress ports $\{m + 1, m + 2, \dots, 2m\}$. Unless mentioned otherwise, all ports have the same bandwidth. We are given a collection of n coflows $\mathbb{C} = \{1, 2, \dots, n\}$, indexed using c . Each coflow c may be assigned a weight w_c (default weight is 1), has an arrival time a_c and comprises a set of flows \mathbb{F}_c .

Table 1: Notation used in the paper.

| | |
|-----------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------|
| w_c | weight of coflow c (default = 1) |
| d_c^{ij} | data sent by coflow c between ingress port i & egress port j |
| d_c^p | Total data sent by coflow c at port p |
| $= \begin{cases} \sum_i d_c^{pi} & \text{if } p \text{ is an ingress port} \\ \sum_j d_c^{jp} & \text{if } p \text{ is an egress port} \end{cases}$ | |

The source, the destination and the size for each flow in the coflow is known at time a_c . The total amount of data sent by coflow c between ingress port i and egress port j is denoted by d_c^{ij} (see Table 1 for notation).

The completion time of a coflow (CCT_c) is the difference between the time when the last of its flows finishes and its arrival time a_c . The average CCT for \mathbb{C} is the average of individual completion times of all coflows $\sum_c \text{CCT}_c / n$. The weighted average CCT is defined as $(\sum_c w_c \times \text{CCT}_c) / n$. Given this formulation, prior work has established that (detailed discussion of related work in §6):

NP-Hardness [10]. Even when all coflows arrive at time 0 and their sizes are known in advance, the problem of minimizing average CCT is NP-hard (via reduction from concurrent open-shop scheduling problem [25]). Thus, the best we can hope for is an approximation algorithm.

Lower Bounds [4, 26]. Even under the big switch model, the only known lower bound is a natural generalization of the lower bound for flows — under a complexity-theoretic assumption somewhat stronger than $P \neq NP$, it is impossible to minimize (weighted) average CCT within a factor of $2 - \epsilon$.

Necessity for Coordination [8]. There exists an instance of coflow scheduling problem, where a scheduling algorithm that does not use any coordination will achieve average CCT $\Omega(\sqrt{n})$ of the optimal. Thus, at least some coordination is necessary to achieve any meaningful approximation.

3 SINCRONIA DESIGN

In this section, we present the core of Sincronia design — an offline algorithm for scheduling coflows for the case when all coflows arrive at time 0; next section describes how Sincronia incorporates this algorithm into an end-to-end network design that achieves near-optimal performance while scheduling coflows in an online and work conserving manner.

Our offline algorithm has two components. The first component is a combinatorial primal-dual greedy algorithm, Bottleneck-Select-Scale-Iterate, for ordering coflows (§3.1); the second component shows that any per-flow rate allocation mechanism that is work conserving, preemptive and schedules flows in order of corresponding coflow ordering achieves average CCT within $4\times$ of the optimal (§3.2).

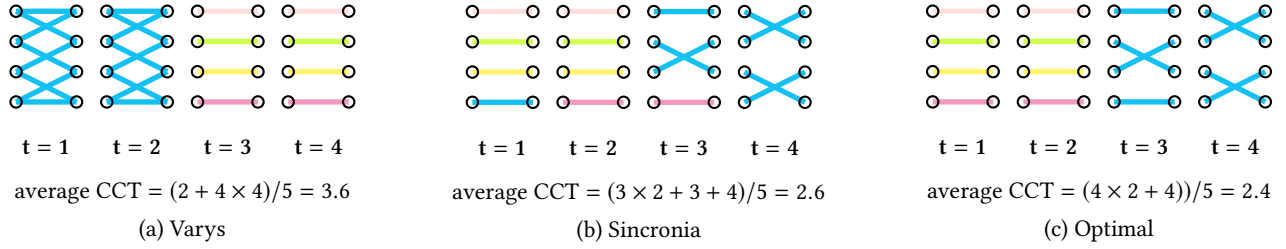


Figure 2: Comparison of Sincronia against Varys and Optimal for the example of Figure 1; corresponding average CCTs are shown for $\varepsilon = 0$. Each figure shows a “matching” between ports at each time slot; a link indicates data being sent between ports in that time slot, with multiple links at a port indicating port bandwidth being equally allocated to each of the links in that time slot. For instance, (a) shows that Varys sends data from first ingress port to the first two egress ports in first two time steps (with each outgoing link getting equal rate allocation), and from first ingress port to first egress port in third and fourth time steps (with the outgoing link getting full rate). The orderings produced by Varys and Sincronia are $\{C1, C2, C3, C4, C5\}$ and $\{C2, C3, C4, C1, C5\}$, respectively; the optimal schedule requires ordering $\{C2, C3, C4, C5, C1\}$ (modulo permutations within coflows $C2, C3, C4$ and $C5$).

Algorithm 1 Bottleneck-Select-Scale-Iterate Algorithm

$\mathbb{C} = [n]$ ▷ Initial set of unscheduled coflows

procedure ORDER-COFLAWS(J)

for $k = n$ to 1 **do** ▷ Note ordering is from last to first

▷ **Find the most bottlenecked port**

$b \leftarrow \arg \max_p \sum_{c \in \mathbb{C}} d_c^p$

▷ **Select weighted largest job to schedule last**

$\sigma(k) \leftarrow \arg \min_{c \in \mathbb{C}} (w_c / d_c^b)$

▷ **Scale the weights**

$w_c \leftarrow w_c - w_{\sigma(k)} \times \frac{d_c^b}{d_{\sigma(k)}^b} \quad \forall c \in \mathbb{C} \setminus \{\sigma(k)\}$

▷ **Iterate on updated set of unscheduled jobs**

$\mathbb{C} \leftarrow \mathbb{C} \setminus \{\sigma(k)\}$

return σ ▷ **Output the coflow ordering**

3.1 Coflow Ordering

Sincronia uses a primal-dual based greedy algorithm — Bottleneck-Select-Scale-Iterate (BSSI) — for ordering coflows (Algorithm 1). BSSI generalizes a near-optimal flow scheduling mechanism, “Shortest Remaining Processing Time” first (SRPT-first) [3], to the case of coflows. The main challenge in achieving such a generalization is to capture how scheduling a coflow impacts the completion time of other coflows in the network. BSSI achieves this using a novel weight scaling step that is derived based on the primal-dual framework.

BSSI operates in four steps — bottleneck, select, scale, iterate. In its first two steps, BSSI generalizes SRPT-first policy for flows to the case of coflows. Intuitively, it does so using an alternative view of SRPT — Largest Remaining Processing Time last (LRPT-last). In particular, the first step finds the most bottlenecked ingress or egress port, say b , defined as the port that sends or receives the most amount of data across all

Table 2: Execution of Algorithm 1 on Figure 1 example. In this example, we break ties in favor of ingress port with the largest index. The final ordering produced by the algorithm is $\{C2, C3, C4, C1, C5\}$.

| k | b | $\sigma(k)$ | $\{w_1, w_2, w_3, w_4, w_5\}$ | \mathbb{C} |
|---|---|-------------|-------------------------------------------------|---------------------|
| — | — | — | $\{1, 1, 1, 1, 1\}$ | $\{1, 2, 3, 4, 5\}$ |
| 5 | 4 | C5 | $\{\varepsilon/(2 + \varepsilon), 1, 1, 1, 0\}$ | $\{1, 2, 3, 4\}$ |
| 4 | 3 | C1 | $\{0, 1, 1, 1 - \varepsilon/2, 0\}$ | $\{2, 3, 4\}$ |
| 3 | 3 | C4 | $\{0, 1, 1, 0, 0\}$ | $\{2, 3\}$ |
| 2 | 2 | C3 | $\{0, 1, 0, 0, 0\}$ | $\{2\}$ |
| 1 | 1 | C2 | $\{0, 0, 0, 0, 0\}$ | \emptyset |

unordered coflows; the second step then implements LRPT-last: it chooses the coflow with largest remaining weighted processing time at port b and places this coflow the last among all unordered coflows. The third step in BSSI scales the weights of all unordered coflows to capture how ordering the coflow chosen in the second step impacts the completion time of all remaining coflows. The final step is to simply iterate on the set of unordered flows until all coflows are ordered.

An Example. Table 2 shows execution of Algorithm 1 on example of Figure 1. In first iteration ($k = 5$), the algorithm chooses bottleneck port $b=4$ and LRPT coflow $\sigma(5) = C5$. The algorithm then scales the weights — in this example, it ends up reducing the weight of coflow C1 while keeping other weights unchanged. This reduction in weight allows C1 to be selected as the LRPT coflow in next iteration ($k = 4$). Figure 2 compares the performance of Sincronia against Varys [10] and optimal for this example. It is not very hard to show that that the average CCT of Varys can be made arbitrarily worse compared to Sincronia (and optimal) by adding more ports and corresponding coflows in the above example (Figure 4).

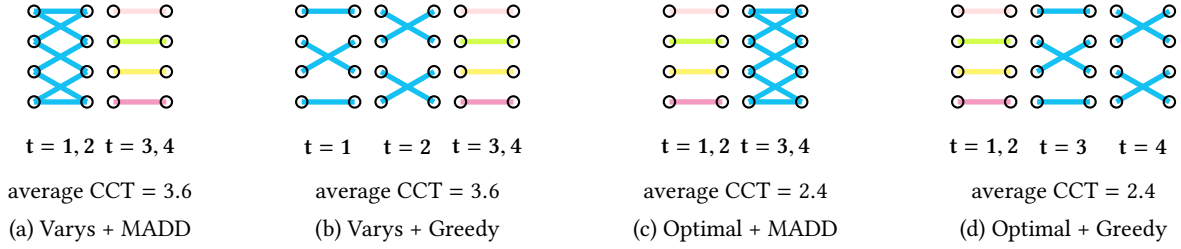


Figure 3: Intuition behind our results in §3.2 using the example of Figure 1. We use two per-flow rate allocation mechanisms in this example. The first one is weighted fair sharing proposed in Varys [10] (in this example, it simply allocates equal rates to all flows at any ingress or egress port). The second one is a greedy rate allocation mechanism that simply chooses one flow from the currently highest ordered coflow at each port, and assigns it the full rate (see Algorithm 2 in §4). The example shows that irrespective of the per-flow rate allocation mechanism used, both Varys and optimal achieve the same average CCT. We do not show Sincronia in this example because MADD ends up allocating non-equal rates to flows in the first coflow and it is hard to depict it pictorially.

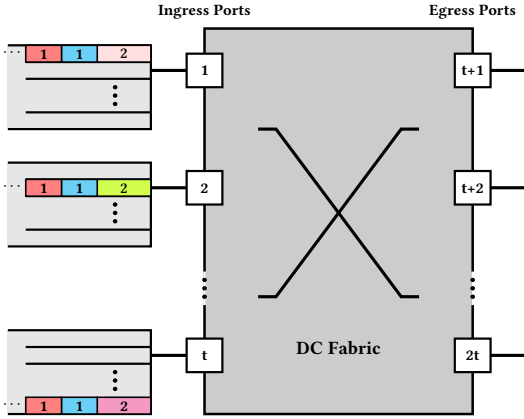


Figure 4: An example to show Varys can achieve unbounded average CCT compared to optimal. The example has t ingress and egress ports and a total of $t + N$ coflows. First t coflows have single flow of two units of data each, with every coflow occupying a different ingress and egress port. Rest N coflows are identical (with two of them shown in the figure), having t unit-sized flows each, with every flow within a coflow occupying different ingress and egress port as shown. The optimal schedule is obtained when all t coflows having single flow each are scheduled in parallel, followed by remaining N coflows in arbitrary order. Varys, however first schedules the N identical coflows first, followed by the t single-flow coflows in parallel. It can be easily shown that for the case when $t > N^{1+\epsilon}$ for any $\epsilon > 0$, the average CCT of Varys compared to the optimal can grow arbitrarily worse as N becomes very large.

3.2 Per-Flow Rate Allocation is Irrelevant

As discussed earlier, prior network designs for coflows require a centralized coordinator to perform complex per-flow rate allocation, where rate allocated to a flow is dependent on rates allocated to other flows in the network; for instance, Varys [10] allocates rates to flows in proportion to their sizes. Such centralized inter-dependent per-flow rate allocations make it hard to realize these designs in practice since

changes in location of congestion in the network, transient failures, and arrival or departure of even one coflow may result in reallocation of rate for each and every flow in the network. Such reallocations are impractical in large datacenters where thousands of coflows may arrive each second. We discuss how Sincronia completely offloads the rate allocation to and scheduling of individual flows to the underlying priority-enabled transport layer. We start with an intuitive discussion, followed by a more formal statement of the result.

High-level idea. Given a coflow ordering produced by our BSSI algorithm, we show that it is sufficient for Sincronia to schedule flows in an order-preserving manner; that is, at any time, a flow from coflow C is blocked if and only if either its ingress port or its egress port is busy serving a flow from a different coflow C' that is ordered before C . The reason this is sufficient is that once a strict ordering between coflows has been established, the proof simply requires finishing each coflow as soon as possible while following the ordering. The main insight is that if there are multiple flows within a coflow starting at an ingress port or ending at an egress port, sharing the link bandwidth does not improve the completion time of this coflow. For instance, in example of Figure 2(a), if we would have given full rate to one flow at each ingress port in the first time step and to the other flow in the other time step, the completion time of coflow $C1$ would not have changed (and so, the same is true for the overall average CCT). Figure 3 demonstrates the irrelevance of per-flow rate allocation for both Varys and the optimal.

Formal statement of results. We now formally state the result regarding the irrelevance of per-flow rate allocation. The detailed proofs for these results are in §9; we give a high-level idea of the proofs in §7.

Definition 1. Let $\sigma : [n] \mapsto [n]$ be an ordering of coflows. A flow scheduling mechanism M is said to be σ -order-preserving if M blocks a flow f from coflow $\sigma(k)$ if and only

if either its ingress port or its egress port is busy serving a flow from a coflow $\sigma(i)$, $i < k$ (preemption is allowed).

Theorem 1. Consider a set of coflows \mathbb{C} , all of which arrive at time 0. Let \mathcal{O} be the ordering of coflows produced by the Bottleneck-Select-Scale-Iterate algorithm for \mathbb{C} and consider **any** work-conserving, pre-emptive and \mathcal{O} -order-preserving flow rate allocation scheme used to schedule \mathbb{C} in Sincronia. Then, under the big switch model, Sincronia achieves average coflow completion time within $4\times$ of the optimal average coflow completion time for \mathbb{C} .

It turns out that if work conservation is desired, preemption is necessary to achieve bounded average CCT:

Claim 1. There exists a set of coflows \mathbb{C} for which the average coflow completion time using any work-conserving, **non-preemptive** flow rate allocation scheme can be arbitrarily worse than the optimal average coflow completion time for \mathbb{C} .

Definition 2. Let σ be an ordering of coflows. For any coflow $\sigma(k)$, let the **ordered load** for coflow $\sigma(k)$ on port p with respect to σ be: $\sum_{i \leq k} d_{\sigma(i)}^p$. Furthermore, let $\hat{p}(k)$ be the port with highest ordered load for $\sigma(k)$. That is

$$\hat{p}(k) \leftarrow \arg \max_p \sum_{i=1}^k d_{\sigma(i)}^p$$

Definition 3. Let σ be an ordering of coflows. Let $\mathcal{A}(\sigma)$ be the class of flow scheduling algorithms for which the completion time of each coflow $\sigma(k)$ is no earlier than its ordered load at port $\hat{p}(k)$ divided by the bandwidth of port $\hat{p}(k)$.

$\mathcal{A}(\sigma)$ is indeed a large class of flow scheduling mechanisms — the only condition is that the last bit of coflow $\sigma(k)$ is sent no earlier than the ordered load at port $\hat{p}(k)$. Let $\text{OPT}(\mathcal{A}(\sigma))$ be the optimal average coflow completion time for the set of coflows across all flow scheduling mechanisms in $\mathcal{A}(\sigma)$.

Theorem 2. Suppose all coflows arrive at time 0 and let σ be an ordering of coflows. Then **any** flow rate allocation scheme that is work-conserving, is pre-emptive and is σ -order-preserving achieves average coflow completion time within $2\times$ of $\text{OPT}(\mathcal{A}(\sigma))$. This bound is tight.

4 SINCRONIA IMPLEMENTATION

We now provide details on the Sincronia implementation. We start with a description of the end-to-end system (§4.1). The remainder of the section focuses on three important aspects of the implementation. First, the BSSI algorithm from §3.1 assumes that all coflows arrive at time 0. In §4.2, we discuss how Sincronia implementation incorporates the BSSI algorithm to efficiently handle coflows arriving at arbitrary times. Second, we showed in §3.2 that Sincronia decouples coflow ordering from rate allocation to and scheduling of individual flows. In §4.3, we discuss how this result enables

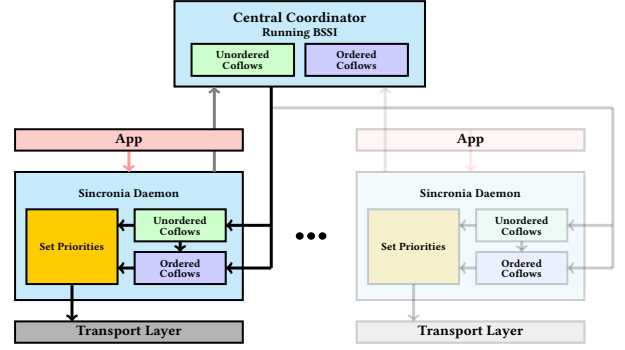


Figure 5: Sincronia end-to-end architecture. See §4.1 for description of various components.

Sincronia to be efficiently integrated with existing datacenter transport mechanisms, including TCP, pHost and pFabric. Finally, we discuss how Sincronia implementation achieves work conservation (§4.4), enables co-existence of flows and coflows, and resolves various other practical issues (§4.5).

4.1 End-to-end system design

We have implemented Sincronia in C++ using just 3000 lines of code. This includes a central coordinator and at each server, a shim layer that sits between the application and the transport layers (see Figure 5).

Applications, upon arrival of a coflow, inform Sincronia daemon about the coflow (coflow ID, flows, and corresponding sources, destinations, sizes, etc.); the daemon adds this coflow to the list of “unordered coflows” (to be used for work conservation) and then uses the above information to register the coflow to the coordinator. The daemons also maintain a list of “ordered coflows”, ones that have been assigned an ordering by the coordinator (as discussed below). When the ongoing flow finishes, the daemon picks one flow from the currently highest ordered coflow (if one exists) or from list of unordered coflows (if no ordered coflow exists), assigns the flow an appropriate priority, and sends it to the underlying transport layer. The daemon also unregisters the finished (co)flow from the coordinator. The priorities assigned to both ordered and unordered coflows depend on the underlying transport mechanism and are discussed in more depth in §4.3.

The coordinator performs the following tasks. It divides the time into epochs. The coordinator maintains a list of “unordered coflows”, those that have been registered but not yet ordered (ordering is done only at the start of each epoch). At the start of each epoch, the coordinator selects a subset of coflows from the unordered coflow list and uses the offline algorithm to order these coflows. We discuss several strategies for deciding the epoch size and for selecting the

subset of coflows at the starting of epoch in §4.2. Once computed, the coordinator removes the ordered coflows from the unordered coflow list, and sends the “ordered coflow” list to all the servers that have unfinished coflows; we use several optimizations here such that the coordinator only informs the servers of the “delta”, changes in ordered list, rather than resending the entire ordered list. Note that each server is oblivious to epochs maintained at the coordinator; hence, servers are not required to be time synchronized with the coordinator.

As discussed in §2, there is a lower bound of $\Omega(\sqrt{n})$ on achievable approximation for average CCT for mechanisms that do not use coordination [8]; thus, some coordination is necessary. However, Sincronia admits much simpler coordinator than existing network designs for coflows — in contrast to existing designs that require the central coordinator to perform per-flow rate allocation, Sincronia requires it to just order the coflows using its BSSI algorithm.

4.2 From Offline to Online

We now discuss how Sincronia implementation incorporates the BSSI algorithm into a system that can efficiently handle the online case, where coflows may arrive at arbitrary times. An obvious way to incorporate BSSI algorithm into an online design is to use the approach taken by prior works (e.g., Varys [10]) — run the offline algorithm upon each coflow arrival. However, in datacenter networks where thousands of coflows may arrive within each second, this may lead to high complexity (although we do use this algorithm as a baseline in our simulations).

Sincronia avoids this high complexity approach using a recently proposed framework [20] along with its own BSSI algorithm from the previous section. We provide a high-level description of the framework to keep the paper relatively self-contained. The framework works in three steps. First, time horizon is divided into exponentially increasing sized epochs. At the start of each epoch, the framework runs an approximation algorithm to select a subset of unfinished coflows. This approximation algorithm, introduced by Garg et. al. [14], works as follows: formulate an integer program with a decision variable for each coflow to indicate whether (or not) it should be assigned to complete within that epoch; we add constraints to enforce that the total work required for the selected subset does not exceed the amount of data that can be transferred within the epoch, and we aim to maximize the total weight of the coflows assigned to this epoch. We solve the linear programming relaxation of this integer program, and select each coflow that is at least “half-assigned” by this optimal solution. It is easy to understand why this process loses a factor of 8 in the approximation

guarantee. Once the subset of coflows is selected, any α -approximate offline algorithm can be used to order coflows arriving over time while providing $(8 + \alpha)$ -competitive ratio².

Sincronia implementation of this framework uses the BSSI algorithm (in the last step) to order coflows arriving in an online manner; we set smallest epoch size to be 100ms, with every subsequent epoch being $2\times$ larger. However, Sincronia makes two modifications in its implementation of the framework. The first modification is to incorporate a work conservation step; we describe this in more detail in §4.4. The second modification Sincronia implementation makes in using this framework is to avoid performance degradation due to large epoch sizes. Indeed, if epoch sizes grow arbitrarily large, increasingly larger number of coflows arrive within an epoch and have to wait until the starting of the next epoch to be scheduled (despite work conservation, as discussed in §4.4); small coflows, in particular, observe poor performance. Sincronia thus bounds the maximum number of epochs and once this number is reached, it “resets the time horizon”. Since BSSI algorithm provides 4-approximation guarantees, using the above formula, Sincronia system achieves a competitive ratio of 12 for the online case.

We emphasize that exponentially-increasing sized epochs are needed only for polynomial-time complexity of the online algorithm and for theoretical guarantees [20]. If these were not the goals, Sincronia could simply use suitably chosen fixed-sized epochs. We evaluate this in §5 and show that Sincronia performance with exponentially-increasing sized epochs is very similar to that with fixed-sized epochs.

4.3 Sincronia + Existing Transport Layers

We now discuss Sincronia implementation on top of several existing transport layer mechanisms for flows. We already described Sincronia coordinator and daemon functionalities in §4.1; these remain unchanged across Sincronia implementation on top of various transport mechanisms. We focus on the only difference across implementations — assignment of priorities to individual flows before Sincronia daemon offloads the flows to the underlying transport mechanism.

Sincronia + TCP. If the underlying network fabric supports infinite priorities, implementing Sincronia on top of TCP is straightforward — each server daemon, for any given flow, simply assigns it a priority equal to the order of the corresponding coflow, sets the priority using the priority bits in DiffServ [6], and sends the data over TCP. However, in practice, the underlying fabric may only support a fixed number of priorities due to hardware limitations, or due to use of some priority levels for other purposes (e.g., fault tolerance).

²As with all online algorithms, Sincronia guarantees for the online version are in terms of competitive ratio.

With finite number p of priorities, Sincronia implementation on top of TCP (with DiffServ) approximates the ideal Sincronia performance using a simple modification: the daemon sets the priority of a given flow to be the order of the corresponding coflow if the current order of the coflow is less than $p - 1$, else it assigns priority p to the flow. As the list of active coflows is updated, the priority used for a flow is also updated accordingly. When using unordered coflows for work conservation, the daemon also sets priority p . While not ideal, our experimental results over real testbeds that support 8 priority levels (§5) show that Sincronia achieves significant improvements in average CCT even with small number of priority levels.

Sincronia + pHost. We now discuss Sincronia implementation on top of pHost, a receiver-driven transport layer mechanism. This is particularly interesting because pHost handles incast and outcast traffic patterns efficiently, precisely a use case for coflows. We extend the pHost implementation to support special prioritized scheduling required for Sincronia implementation (at the receiver as well as at the source and at intermediate network elements). We assume familiarity with pHost design. In our implementation, the pHost receiver sends a **token** for packets in flows in order of corresponding coflow ordering. The server daemons, among all the received **tokens**, use the one for a flow in the currently highest ordered coflow. In-network priorities are not necessary (congestion is at the edge due to per-packet scheduling and packet spraying) but can be supported using TCP style priority assignment. Interestingly, pHost implementation of Sincronia converges to the greedy algorithm shown in Algorithm 2 that, in a converged state, assigns a single flow the full outgoing access link rate at any given point of time.

Algorithm 2 Greedy Rate Allocation Algorithm

```

All access links have uniform bandwidth
 $\mathbb{C} = \sigma$  is the input coflow ordering
procedure GREEDYFLOWSCHEULING( $\sigma$ )
  while  $\mathbb{C}$  is not empty do
    for  $i = 1$  to  $|\mathbb{C}|$  do
      for  $j = 1$  to  $|C_i.flows|$  do
        if Ingress port of  $C_i.flows(j)$  free then
          if Egress port of  $C_i.flows(j)$  free then
            Allocate entire BW to  $C_i.flows(j)$ 
        Update flow sizes and available link bandwidth
      GreedyFlowScheduling( $\mathbb{C}$ )
  return

```

Sincronia + pFabric. Sincronia implementation on top of pFabric admits an even simpler design. Since pFabric already

supports infinite priority levels, we only need a minor change in pFabric priority assignment mechanism: each flow is now assigned a priority equal to the ordering of its coflow, rather than the size of the flow as in original pFabric paper [3]. A special “minimum priority level” is assigned to unordered coflows for work conservation purposes. Since Sincronia generates a total ordering across coflows, this implementation results in ideal Sincronia performance.

4.4 Prioritized Work Conservation

Sincronia coordinator runs the BSSI algorithm for coflow ordering at the starting of each epoch. Thus, coflows that arrive in the middle of an epoch (referred to as “orphan coflows”) may not be assigned an ordering until the starting of the next epoch. While our flow scheduling mechanisms from the previous section are naturally work-conserving (because underlying transport layer mechanisms are work conserving), our preliminary implementation highlighted a potential performance issue. Orphan coflows, since unordered, end up fair sharing the bandwidth. For “short” orphan coflows, such fair sharing could lead to a long tail.

Sincronia uses a simple optimization for work conservation that alleviates this problem. Each orphan coflow is assigned an ordering among all the orphan coflows based on its “oracle completion time (OCT)”, which is defined as the time the coflow would take to finish if it were the only coflow in the network. Note that the OCT of a coflow c is simply $\max_p(d_c^p / b_p)$, where b_p is the bandwidth at port p . Specifically, the orphan coflows are chosen to be scheduled for work conservation in increasing order (smaller first) of the following metric:

$$\frac{\text{OCT}}{\text{current_time} + \text{max_epoch_size} - \text{arrival_time}}$$

The reason this metric is interesting for work conservation is that it finds the right balance between coflows that are “small” (that is, have small OCTs) and coflows that are large but have been waiting for a long while (that is, coflows that are being starved) while selecting coflows to use for work conservation. For instance, consider a time instant when there are three coflows in the system; two of these coflows arrived at time 20 and have OCTs equal to 1 and 10, respectively, and the third coflow arrived at time 0 and has OCT equal to 25. Suppose the `max_epoch_length` is 4. Then, when choosing a coflow for work conservation at time 21, the Sincronia daemon will choose the coflow with OCT equal to 1 first (the metric value being $1/5$), then choose the coflow with OCT equal to 25 that has been waiting for too long (the metric value being 1) and then finally send the coflow with OCT equal to 10. Thus, Sincronia is able to find the right balance between

small coflows and starving coflows to choose from for work conservation purposes.

4.5 Other Practical Considerations

Finally, we discuss a few other techniques incorporated within Sincronia implementation to handle practical issues that may arise in real-world scenarios.

Co-existence of flows and coflows. In most datacenter networks, multiple applications co-exist, some of which may require the network fabric to optimize for coflow-based metrics while others may care about the performance of each individual flow. Prior results on network design for coflows enable coexistence of such applications by treating each individual flow as a coflow. This may be restrictive since flow based applications may have different performance goals compared to coflow based applications. Sincronia partially handles such scenarios using coflow weights.

Specifically, while our discussion so far has focused on Sincronia design and implementation for achieving near-optimal performance in terms of average CCT, Sincronia achieves something much more general — it optimizes for “weighted” average CCT, as defined in §2. That is, it allows network operators to set a weight for each individual coflow. Network operators can use different weights for different applications (e.g., those that require optimizing for flows and those that require optimizing for coflows), and the BSSI algorithm computes ordering of coflows that optimize for the weighted average CCT. Finding the right mechanism to set weights depends on the applications and is beyond the scope of this paper.

Achieving other performance objectives. Sincronia also allows achieving several other performance objectives using the coflow weights. For instance, there has been quite a bit of work in the community on deadline-aware scheduling for the traditional abstraction of network flows [3, 17, 27]. One way to handle deadlines in Sincronia is to assign coflow weights inversely proportional to their deadlines. Sincronia also supports admission control by assigning zero weights to coflows which would definitely miss their deadlines, hence scheduling them after other coflows. Designing coflow scheduling algorithms that provide provable guarantees for the case of deadlines is an interesting open problem (§6).

Starvation Freedom. Minimizing average completion time necessitates starvation [3, 10]. However, as discussed in §4.4, the prioritized work conservation mechanism used in Sincronia alleviates starvation to some extent. Intuitively, since the choice of an unordered coflow for work conservation depends inversely on the “waiting time” of the coflow (difference between current time and the arrival time), as the

waiting time of the coflow increases, its chances of being selected for work conservation purposes also improve.

5 EVALUATION

We now present evaluation results for Sincronia³. We start by describing the workloads and performance metrics used in our evaluation (§5.1). We then evaluate Sincronia using simulations (§5.2); the main goal here is to understand the performance of Sincronia against the state-of-the-art [10], to understand the envelope of workloads for which Sincronia performs well, and to perform sensitivity analysis of Sincronia performance over a variety of parameters including workloads, number of coflows, network load and epoch sizes. We then evaluate the Sincronia implementation with TCP/DiffServ on a 16-server testbed interconnected with a FatTree topology comprising 20 commodity switches (§5.3).

5.1 Workloads and Performance Metrics

We use two workloads in our evaluation. The first workload is the one used in all prior network designs [10]: a 526-coflow trace obtained from one-hour long run of MapReduce jobs on a 3000-machine cluster at Facebook. Unfortunately, the Facebook trace makes several simplifying assumptions, is limited in number of coflows and imposes low network load. The second workload is a collection of traces generated using a coflow workload generator [1] that allows upsampling the Facebook trace to desired number of coflows, network load, etc., while keeping workload characteristics similar to the original Facebook trace. Unless mentioned otherwise, all results for these “custom” traces use the baseline of 2000 coflows, 0.9 network load and time horizon reset after 8 epochs. Varys does not support weighted coflow scheduling; thus, to make a fair comparison, we use unit weights for all coflows. We provide other setup details (e.g., topology, link bandwidth, etc.) for simulations and implementation in respective subsections.

Metrics. For the offline algorithm, we evaluate the performance in terms of CCT, on an average and at high percentiles. For the online algorithm, we evaluate the performance of Sincronia using the slowdown metric; for a given coflow, the slowdown is defined as the ratio of its CCT and its OCT (recall from §4.4, the OCT of a coflow is its completion time when its the only coflow in the network). In presence of other coflows, CCT of a coflow may be much larger than its OCT; thus, our evaluation results for the online version of Sincronia are against the best possible baseline.

³Sincronia implementation and simulator are publicly available at: <https://github.com/sincronia-coflow>.

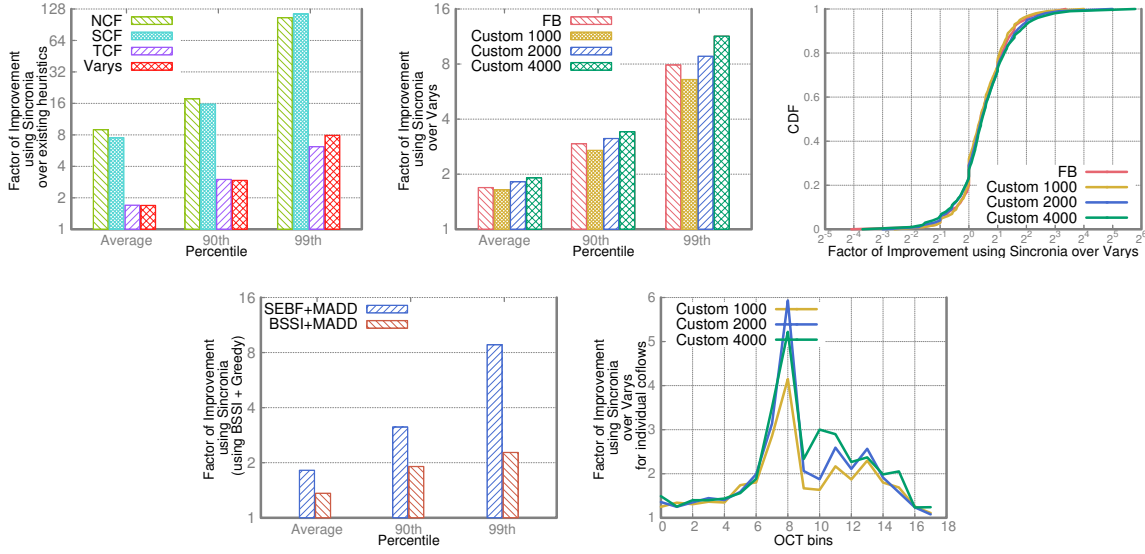


Figure 6: Sincronia evaluation results for the offline case. Sincronia achieves consistently better performance compared to existing heuristics for the original Facebook trace (top left) and maintains these improvements across a variety of workloads (top center, top right). Most of Sincronia’s improvements are due to its BSSI scheduling algorithm from §3 (bottom left, bottom right). More discussion in §5.2.

5.2 Simulation Results

We now evaluate Sincronia against the state-of-the-art [10] using simulations, and perform sensitivity analysis of Sincronia performance over a variety of parameters.

Setup. Varys uses a flow-level simulator in its evaluation. For a fair comparison, we incorporate Sincronia within Varys simulator, with BSSI for coflow ordering (Algorithm 1, §3.1) and greedy algorithm for flow scheduling (Algorithm 2, §4.3). We use the same setup used in Varys, modeling the network using a non-blocking 150-port switch, where each port corresponds to a top-of-the-rack switch in 3000-machine cluster.

Offline Algorithm. We first evaluate Sincronia against existing network designs for the offline case, that is, when all coflows arrive at time 0. While not a practical scenario, this allows us to tease out the performance benefits of Sincronia’s BSSI algorithm against those used in prior network designs.

Figure 6(a)-6(c) show that Sincronia significantly improves upon state-of-the-art (Varys [10] and other heuristics) across all evaluated workloads. Specifically, Figure 6(a) shows that, when compared to TCF and Varys, Sincronia improves CCT by 1.7× on an average and by 7.9× at 99-th percentile. In addition, Figure 6(b) shows that Sincronia improvements over Varys increase with larger number of coflows in the trace, both on an average and at high percentiles. Since both Varys and Sincronia are work-conserving, CCT for some coflows has to degrade for improvement in CCT for other coflows. Figure 6(c) shows the CDF for distribution of improvement

across all coflows; interestingly, Sincronia degrades the CCT for less than 20% of the coflows while improving the CCT for more than 65% of the coflows.

Figure 6(d)-6(e) provide more insights on how Sincronia achieves the observed improvements. Figure 6(d) shows that most of the improvements come from Sincronia’s ordering algorithm, Bottleneck-Select-Scale-Iterate (BSSI) — even when BSSI is used with MADD, the rate allocation mechanism from Varys, we observe significant CCT improvements on an average and at tail. Figure 6(e) provides more insights about BSSI’s performance. The figure plots the average factor of improvement measured across coflows binned by their OCTs normalized by the smallest OCT; that is, an OCT bin i contains the set of coflows whose OCT is within $[2^i, 2^{i+1})$ of the OCT of the smallest coflow. The figure shows that, when compared to Varys, Sincronia results in CCT improvements by making better scheduling decisions for coflows that are neither too small (low OCTs) nor too large (high OCTs). Indeed, scheduling decisions for small and large coflows are rather straightforward, and both Sincronia and Varys end up making similar decisions; it is precisely the medium size coflows where the weight scaling step in BSSI helps Sincronia in making the “right” decisions.

Online Algorithm. We now evaluate Sincronia for the case when coflows arrive over time. We focus on custom traces only since the Facebook trace has very low load and does not provide meaningful insights. As discussed in §5.1, we use the absolute best possible baseline here — ratio of coflow

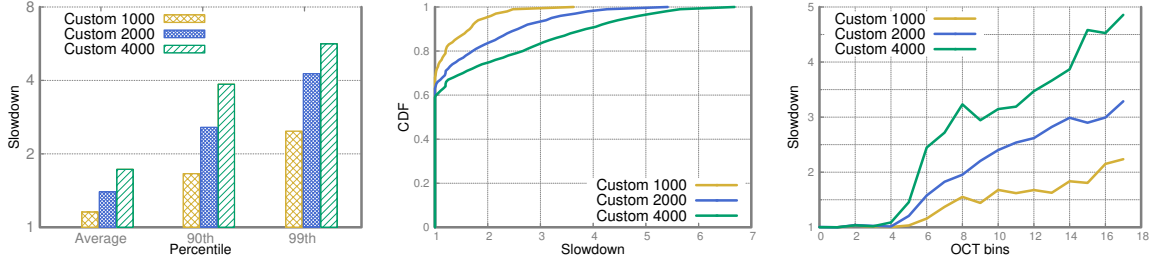


Figure 7: Sincronia evaluation results for the online case. (left, center) For the evaluated workloads, Sincronia achieves performance within $1.7\times$ on an average and $5.7\times$ at high percentiles when compared to an unloaded network; (right) Coflows with larger sizes observe larger slowdown in online version of Sincronia. More discussion in §5.2.

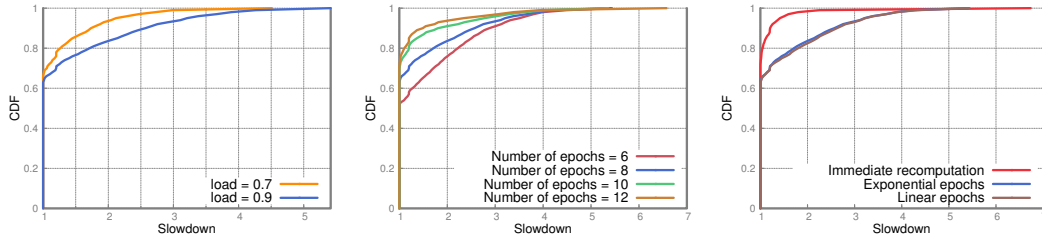


Figure 8: Sensitivity analysis of Sincronia performance over varying network load (left), varying number of epochs before resetting the time horizon (center) and varying epoch mechanisms (right). More discussion in §5.2.

CCT and its OCT. Figure 7 shows that Sincronia handles online arrival of coflows very well; when compared to an unloaded network (slowdown = 1), Sincronia at load 0.9 performs within $1.7\times$ on an average and within $5.7\times$ at high percentiles for the Custom 4000 trace.

Figure 7(b)-7(c) provide more insights into Sincronia’s performance. The former shows that around 60% of the coflows achieving slowdown 1, that is, are optimally scheduled by Sincronia. As the size of the trace increases, contention for network resources increases resulting in higher slowdowns. Figure 7(c) shows that larger slowdowns are caused mainly by larger coflows; this is not only because smaller coflows should indeed be prioritized but also that the work conservation heuristic in Sincronia (§4.2) prioritizes smaller coflows until long coflows have been waiting for a long time.

Sensitivity Analysis. Figure 6 and Figure 7 already show impact of trace sizes on Sincronia performance. We perform additional sensitivity analysis of Sincronia performance against network load, number of epochs before resetting the time horizon, and various epoch mechanisms in Figure 8. Figure 8(a) shows that, as expected, Sincronia performance improves as the network load decreases from 0.9 to 0.7. We see in Figure 8(b) that Sincronia performance varies with number of epochs before resetting the time horizon — as the number of epochs increase, larger fraction of coflows observe completion times close to their OCTs but the tail

slowdown also worsens. Intuitively, as the number of epochs increase, the largest epoch size also increases; this, in turn, increases the chance of a small coflow arriving between epoch boundaries and being blocked by a large “ordered” coflow, one that was ordered at the start of the epoch and hence has higher priority than the small coflow. Finally, we compare the performance of Sincronia with varying epoch mechanisms. We use exponentially-increasing epoch sizes, fixed-epoch sizes (with size equal to the longest epoch in exponentially-increasing epoch size mechanism) and the “immediate recomputation” mechanism that executes the BSSI algorithm upon each coflow arrival and departure. Figure 8(c) shows that, over the evaluated workloads, the first two mechanisms perform comparably and are only marginally worse when compared to the (impractical) immediate recomputation mechanism.

5.3 Implementation Results

We now present evaluation results for Sincronia implementation on top of TCP (using Diffserv [6] for priority scheduling). Our implementation runs on a topology shown in Figure 10. We compare Sincronia performance against coflow-agnostic TCP implementation⁴.

⁴Despite trying for several weeks, we were not able to run Varys [10] over our topology; the code for Varys is not maintained and the software stack for which this code was written has evolved. Nevertheless, we would like

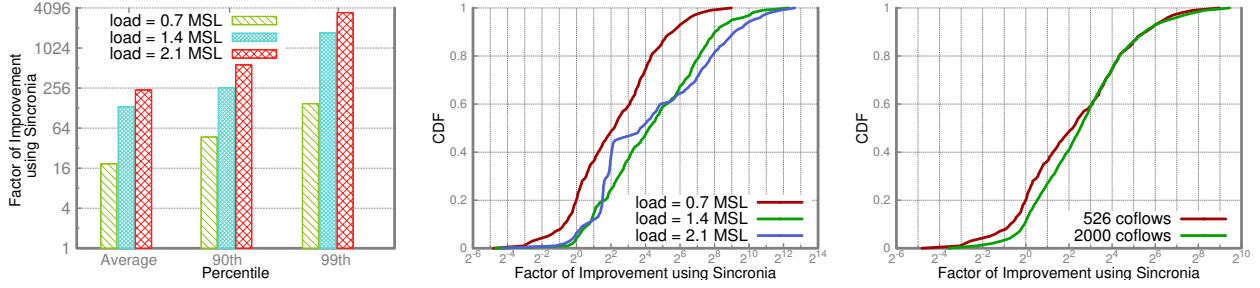


Figure 9: Evaluation results for Sincronia implementation on top of TCP over our testbed in Figure 10. Sincronia when implemented on top of TCP/DiffServ significantly improves coflow performance when compared to coflow-agnostic TCP both on an average and at high percentiles (left), and maintains these improvements across varying network loads (center) and trace sizes (right). More discussion in §5.3.

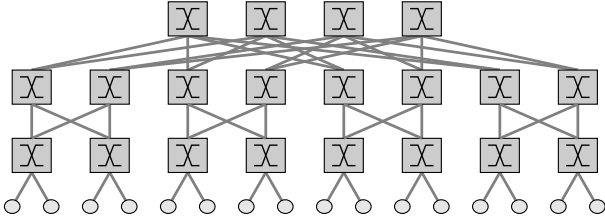


Figure 10: Testbed used in our experiments has 16 servers, each with 1Gbps access link bandwidth, interconnected with a FatTree topology comprising 20 PICA8 switches. Each switch supports priority queueing with 8 distinct levels.

The main challenge in running these experiments is that TCP does not handle incast traffic very well, precisely the scenario for coflows (multiple sources sending data to the same destination). Thus, we compute the expected load on a destination (by using network load along with expected number of sources sending data to a destination in the Facebook trace), and refer to it as Maximum Sustainable Load (MSL) for the cluster. We then generate the workloads for the cluster using the workload generator for 16 ingress/egress ports, with 1Gbps access link bandwidth and varying loads ($0.7 \times \text{MSL}$, $1.4 \times \text{MSL}$ and $2.1 \times \text{MSL}$).

Figure 9 shows that for a workload comprising 526 coflows at load $0.7 \times \text{MSL}$, Sincronia implementation on top of TCP improves the CCT by $18.61 \times$ on an average, by $46.95 \times$ at 90-th percentile, and $149.05 \times$ at 99-th percentile when compared to coflow-agnostic TCP.

As expected, the improvements are even more significant for larger loads. Figure 9(b) shows that Sincronia achieves these improvements by slowing down less than 10% of the coflows by more than $2 \times$ (and even smaller fraction for higher loads). Finally, Figure 9(c) shows that Sincronia

achieves similar benefits for larger traces. This evaluation is a bit unfair — TCP is neither designed for coflows nor for minimizing average completion times. However, these results give us some indication on performance of Sincronia against existing network designs; for instance, Varys [10] reports to improve upon TCP by a factor of $1.85 \times$ on an average, and roughly by the same number at 95 percentiles.

6 RELATED WORK & OPEN PROBLEMS

There has been tremendous recent effort on network design for coflows, both in networking and in theory communities. In this section, we compare and contrast Sincronia against related work from these communities, and discuss a number of problems that remain open in the context of network design for coflows.

Gap between lower and upper bounds. Existing network designs for coflows from networking community [7, 9, 10, 12] can perform arbitrarily worse than the optimal in terms of average CCT (§3). Qiu et. al. [23] provided the first coflow scheduling algorithm that achieves approximation guarantees of $64/3$ for the offline case with zero release dates. This result was improved by Ahmadi et. al. [2] to 4-approximation for zero release dates and to 5-approximation with release dates. However, all these results require a central coordinator to either solving complex linear programs or to perform complex per-flow rate allocations. Sincronia achieves state-of-the-art approximation guarantees, both for zero release dates and, using a simple extension of the BSSI algorithm, for the case of release dates (§9-10). However, Sincronia shows that these guarantees can be achieved without any complex per-flow rate allocation, thus overcoming the limitations of all existing results from the networking and the theory community. Despite all the above results, there is a gap between known lower bounds for coflow scheduling (impossibility of better than 2-approximation, §2) and the known upper

to thank the authors of Varys [10] to share their implementation with us and helping us to try and run their code.

bounds (4-approximation for offline problem with zero release dates) [2]. It remains an intriguing problem to bridge this gap.

Improved competitive ratio for online coflow scheduling. Khuller et. al. [20] were the first to present an online algorithm that achieves a competitive ratio of 12. Sincronia achieves similar guarantees. Is it possible to design a coflow scheduling algorithm with a better competitive ratio? Note that an offline algorithm that achieves a better approximation factor for the case of zero release dates will already lead to improvements in the online algorithm by using the framework from Khuller et. al. [20].

Necessity of centralized solutions. As discussed in §2, there is a strong lower bound of $\Omega(\sqrt{n})$ on achievable approximation for average CCT using algorithms that do not use any coordination [8]. However, the “amount of coordination” required to mitigate this lower bound is unclear. For instance, similar to [10, 31, 32], Sincronia requires a centralized controller to implement BSSI. While the amount of computation required by Sincronia’s centralized scheduler is much lower than that of prior solutions (since Sincronia scheduler does not need to do per-flow rate allocation), it remains an open problem to understand the amount of computation needed to be done at the centralized controller.

Coflows with paths. Sincronia, similar to [10, 12], assumes that routing of flows within and across coflows are decided by the network layer (either via specific routing mechanisms, or via packet spraying). However, it is a priori conceivable that better performance may be achievable by co-designing routing along with scheduling of coflows. This problem has been studied in a few recent papers [19, 32]; however, designing optimal algorithms for this problem remains open.

Extensions to Non-clairvoyant scheduler. Sincronia, similar to [10, 12, 32], assumes that the information about a coflow is available at the arrival time (and no earlier). For many applications, this is indeed the case (see [10] and the discussion therein). Moreover, recent work has shown that it is possible to identify coflows and their properties within reasonable estimate for many applications [31]. However, designing a near-optimal non-clairvoyant scheduler remains an interesting future direction.

Extensions to other performance metrics. Finally, similar to most of the existing literature on coflows [8, 10, 31, 32], Sincronia is designed to optimize for average (weighted) completion time. It remains an interesting question to extend our results to the case of deadline-sensitive coflows and for minimizing other metrics such as tail coflow completion time. An intriguing question in this direction is also to define a notion of fairness for coflows.

7 PROOF OUTLINE FOR THEOREM 1

In this section, we provide a high-level idea for our results in §3. The detailed proofs for these results are in §9–10.

At a high-level, we use a linear programming relaxation to obtain a lower bound on the optimal value of the average CCT, and obtain our approximation guarantee by comparing the average CCT of our algorithm to this lower bound. Our algorithm is purely combinatorial, in that it does not require solving an LP.

In our algorithm, we decouple the problem of obtaining a feasible schedule into two parts: we first obtain an ordering of coflows, and then obtain a feasible schedule by using a greedy rate allocation scheme that maintains this order.

To obtain an ordering of coflows, we also relax the problem by ignoring the dependencies between the input and output ports; more specifically, we consider an instance of a concurrent open shop problem where there are $2m$ machines corresponding to the m ingress and the m egress ports, and one job corresponding to each of the n coflows. The processing requirement for job c on machine k is the total load at port k due to coflow c , and the weight of job c is the weight associated with coflow c . Observe that the optimal value of this concurrent open shop instance is a lower bound on the optimal value of the original coflow scheduling instance since any feasible solution to the latter can be viewed as a feasible solution for the former with the same objective function value. Next, observe that there is an optimal solution for the concurrent open shop input in which the order of the jobs processed on each machine is the same: given any optimal solution, if we consider the last job on the most heavily loaded machine we can push that job to the end on each other machine while maintaining that none of the job completion times increases. We can repeat this to obtain a solution where the jobs are processed in the same order on each machine without increasing the objective function value. So we have reduced our problem to one of just finding an ordering of coflows, since given a ordering, determining a feasible schedule is straightforward.

We use a primal-dual algorithm to compute an ordering such that the weighted completion time of the jobs is at most twice the optimal value for the concurrent open shop instance; hence, by our lower bound argument, this weighted completion time is at most twice the optimal value for the coflow scheduling problem as well.

For the second part of the problem, we present a greedy rate-allocation scheme that maintains the order returned by the primal-dual algorithm. More specifically, we ensure that a flow from input port i to output port j on a coflow c is scheduled only after all flows between this pair of ports from coflows of a higher priority have completed. So, if we consider the last flow processed for some coflow c , say from

input port i to output port j , this property, coupled with the fact that the algorithm is work-conserving, implies that at least one of the ports i or j was busy for at least half the completion time of coflow c with processing flows from coflows with an equal or higher priority than coflow c . Since the total amount of work from higher priority coflows that can be done on these ports is at most the completion time of coflow c in the concurrent open shop instance, we arrive at the conclusion that the completion time of coflow c in our greedy rate-allocation scheme is at most twice the completion time of coflow c in the concurrent open shop instance. This result in fact extends to any rate allocation scheme that is preemptive, work-conserving, and maintains the ordering of the coflows. Combining the above two results directly yields the 4-approximation result.

8 CONCLUSION

We have presented Sincronia, a network design for coflows that provides near-optimal performance and can be implemented on top on any transport layer mechanism for flows that supports priority scheduling. Sincronia achieves this using a key technical result — we show that given a “right” ordering of coflows, any per-flow rate allocation mechanism achieves average coflow completion time within $4\times$ of the optimal as long as (co)flows are prioritized with respect to the ordering. This allows Sincronia to use a simple greedy mechanism to “order” all unfinished coflows; all flows within and across coflows can then be greedily scheduled using any transport mechanism that supports priority scheduling (without any per-flow rate allocation mechanism). Evaluation results suggest that Sincronia not only admits a practical, near-optimal design but also improves upon state-of-the-art network designs for coflows.

ACKNOWLEDGMENTS

This work was supported in part by NSF grants CCF-1526067, CMMI-1537394, CCF- 1522054, and CCF-1740822.

REFERENCES

- [1] 2018. Sincronia Repository. <https://github.com/sincronia-coflow>.
- [2] Saba Ahmadi, Samir Khuller, Manish Purohit, and Sheng Yang. 2017. On scheduling coflows. In *MOS IPCO*.
- [3] Mohammad Alizadeh, Shuang Yang, Milad Sharif, Sachin Katti, Nick McKeown, Balaji Prabhakar, and Scott Shenker. 2013. pFabric: Minimal Near-optimal Datacenter Transport. In *ACM SIGCOMM*.
- [4] Nikhil Bansal and Subhash Khot. 2010. Inapproximability of hypergraph vertex cover and applications to scheduling problems. In *EATCS ICALP*.
- [5] Soumen Chakrabarti, Cynthia A Phillips, Andreas S Schulz, David B Shmoys, Cliff Stein, and Joel Wein. 1996. Improved scheduling algorithms for minsum criteria. In *EATCS ICALP*.
- [6] Kwok Ho Chan, Jozef Babiarz, and Fred Baker. 2006. Configuration Guidelines for DiffServ Service Classes. <https://tools.ietf.org/html/rfc4594>.
- [7] Mosharaf Chowdhury and Ion Stoica. 2012. Coflow: A networking abstraction for cluster applications. In *ACM HotNets*.
- [8] Mosharaf Chowdhury and Ion Stoica. 2015. Efficient coflow scheduling without prior knowledge. In *ACM SIGCOMM*.
- [9] Mosharaf Chowdhury, Matei Zaharia, Justin Ma, Michael I Jordan, and Ion Stoica. 2011. Managing data transfers in computer clusters with orchestra. In *ACM SIGCOMM*.
- [10] Mosharaf Chowdhury, Yuan Zhong, and Ion Stoica. 2014. Efficient coflow scheduling with varys. In *ACM SIGCOMM*.
- [11] Jeffrey Dean and Sanjay Ghemawat. 2004. MapReduce: simplified data processing on large clusters. In *USENIX OSDI*.
- [12] Fahad R Dogar, Thomas Karagiannis, Hitesh Ballani, and Antony Rowstron. 2014. Decentralized task-aware scheduling for data center networks. In *ACM SIGCOMM*.
- [13] Peter X Gao, Akshay Narayan, Gautam Kumar, Rachit Agarwal, Sylvia Ratnasamy, and Scott Shenker. 2015. phost: Distributed near-optimal datacenter transport over commodity network fabric. In *ACM CoNEXT*.
- [14] Naveen Garg, Amit Kumar, and Vinayaka Pandit. 2007. Order scheduling models: hardness and algorithms. In *IARCS FSTTCS*.
- [15] Leslie A Hall, Andreas S Schulz, David B Shmoys, and Joel Wein. 1997. Scheduling to minimize average completion time: Off-line and on-line approximation algorithms. In *Mathematics of operations research*, 22(3): 513–544. INFORMS.
- [16] Mark Handley, Costin Raiciu, Alexandru Agache, Andrei Voinescu, Andrew Moore, Gianni Antichi, and Marcin Wojcik. 2017. Re-architecting datacenter networks and stacks for low latency and high performance. In *ACM SIGCOMM*.
- [17] Chi-Yao Hong, Matthew Caesar, and P Godfrey. 2012. Finishing flows quickly with preemptive scheduling. In *ACM SIGCOMM*.
- [18] Michael Isard, Mihai Budiu, Yuan Yu, Andrew Birrell, and Dennis Fetterly. 2007. Dryad: distributed data-parallel programs from sequential building blocks. In *ACM EuroSys*.
- [19] Hamidreza Jahanjou, Erez Kantor, and Rajmohan Rajaraman. 2017. Asymptotically Optimal Approximation Algorithms for Coflow Scheduling. In *ACM SPAA*.
- [20] Samir Khuller, Jingling Li, Pascal Sturm, Kevin Sun, and Prayaag Venkat. 2018. Select and Permute: An Improved Online Framework for Scheduling to Minimize Weighted Completion Time. In *LATIN*.
- [21] Yucheng Low, Danny Bickson, Joseph Gonzalez, Carlos Guestrin, Apo Kyröla, and Joseph M Hellerstein. 2012. Distributed GraphLab: a framework for machine learning and data mining in the cloud. *Proceedings of the VLDB Endowment*, 5(8): 716–727.
- [22] Grzegorz Malewicz, Matthew H Austern, Aart JC Bik, James C Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. 2010. Pregel: a system for large-scale graph processing. In *ACM SIGMOD*.
- [23] Zhen Qiu, Cliff Stein, and Yuan Zhong. 2015. Minimizing the total weighted completion time of coflows in datacenter networks. In *ACM SPAA*.
- [24] Maurice Queyranne. 1993. Structure of a simple scheduling polyhedron. In *Mathematical Programming*, 58(1): 263–285. Springer.
- [25] Thomas A. Roemer. 2006. A note on the complexity of the concurrent open shop problem. In *Journal of Scheduling*, 9(4): 389–396. Springer.
- [26] Sushant Sachdeva and Rishi Saket. 2013. Optimal inapproximability for scheduling problems via structural hardness for hypergraph vertex cover. In *IEEE CCC*.
- [27] Christo Wilson, Hitesh Ballani, Thomas Karagiannis, and Ant Rowstron. 2011. Better never than late: Meeting deadlines in datacenter networks. In *ACM SIGCOMM*.
- [28] L. A. Wolsey. 1990. Formulating single machine scheduling problems with precedence constraints. In *Economic Decision-Making: Games, Econometrics and Optimisation*, 473 – 484. North-Holland.

- [29] Yuan Yu, Michael Isard, Dennis Fetterly, Mihai Budiu, Úlfar Erlingsson, Pradeep Kumar Gunda, and Jon Currey. 2008. DryadLINQ: A System for General-Purpose Distributed Data-Parallel Computing Using a High-Level Language.. In *USENIX OSDI*.
- [30] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J Franklin, Scott Shenker, and Ion Stoica. 2010. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *USENIX NSDI*.
- [31] Hong Zhang, Li Chen, Bairen Yi, Kai Chen, Mosharaf Chowdhury, and Yanhui Geng. 2016. CODA: Toward automatically identifying and scheduling coflows in the dark. In *ACM SIGCOMM*.
- [32] Yangming Zhao, Kai Chen, Wei Bai, Minlan Yu, Chen Tian, Yanhui Geng, Yiming Zhang, Dan Li, and Sheng Wang. 2015. Rapier: Integrating routing and scheduling for coflow-aware data center networks. In *IEEE INFOCOM*.

9 OFFLINE APPROXIMATION ALGORITHM

In this section, we prove that the algorithm described in section 3 is a 5-approximation algorithm for the general case and a 4-approximation algorithm when there are no release dates. We do this by dividing the problem into two parts - obtaining an ordering of the coflows to be scheduled, and a rate allocation algorithm that respects this ordering. We provide algorithms for each part, with a combined analysis. We first obtain an ordering of coflows by formulating a linear programming relaxation, from which we derive a primal-dual algorithm with provable guarantees. We then use a greedy algorithm for rate allocation, while maintaining the order derived from the primal-dual algorithm.

The algorithms are based on a linear program whose constraints are similar to those introduced by Wolsey [28] and Queyranne [24] for the problem of scheduling n jobs on a single machine: $\sum_{j \in S} p_j C_j \geq \frac{1}{2} \left[\left(\sum_{j \in S} p_j \right)^2 + \sum_{j \in S} p_j^2 \right]$ for each subset $S \subseteq [n]$, where C_j , $j \in [n]$ are variables corresponding to the job completion time and p_j is the processing time of job j , for $j = 1, \dots, n$. Although there are an exponential number of constraints, Queyranne [24] describes a polynomial-time separation oracle for the set of constraints.

The linear program we consider is not a relaxation of coflow scheduling (in the sense that a feasible schedule directly corresponds to a feasible LP solution), but we will prove that the optimal value of this LP is a lower bound on the weighted completion time of any feasible schedule. We focus on the total flow demands at each input and output port and introduce notation for the total flow at each ingress/egress port p : given a coflow k and the associated traffic matrix D_k , we define the total flow at port p

$$v_k^p = \begin{cases} \sum_j v_k(p, j), & \text{if } p \in \{1, \dots, m\}; \\ \sum_i v_k(i, p), & \text{if } p \in \{m+1, \dots, 2m\}. \end{cases}$$

We also define $\rho(D_k) = \max_p \{v_k^p\}$ to be the maximum row or column sum in D_k .

The LP has a variable C_k for the completion time of each coflow k , and a variable C_k^p that corresponds to the time that coflow k is done sending or receiving flows on the ingress/egress port p :

$$\min \sum_{k=1}^n w_k C_k \quad [\text{LP-primal}]$$

$$\text{s.t. } C_k \geq C_k^p, \quad \forall p \in [2m], k \in [n]; \quad (1)$$

$$C_k^p \geq r_k, \quad \forall p \in [2m], k \in [n]; \quad (2)$$

$$\sum_{k \in S} v_k^p C_k^p \geq \frac{1}{2} \left[\left(\sum_{k \in S} v_k^p \right)^2 + \sum_{k \in S} (v_k^p)^2 \right], \quad \forall p \in [2m], S \subseteq [n]. \quad (3)$$

We prove that the optimal value of LP-primal provides a lower bound for the coflow scheduling problem.

Lemma 1. *Let OPT be an optimal schedule for the set of coflows K and let $(C_k^{\text{OPT}})_{k \in [n]}$ be the completion times of the coflows in this optimal schedule. Furthermore, let C_1^*, \dots, C_n^* be the optimal values of C_1, \dots, C_n in LP-primal. Then $\sum_{k=1}^n w_k C_k^* \leq \sum_{k=1}^n w_k C_k^{\text{OPT}}$.*

Proof 1. *To prove this, we establish a mapping from feasible schedules of $\text{Coflow}|r_k| \sum w_k C_k$ to feasible solutions of LP-primal.*

Consider a feasible schedule and let \hat{C}_k be the completion time of coflow k , $k = 1, \dots, n$. Furthermore, assume that the coflows are re-indexed so that $\hat{C}_1 \leq \hat{C}_2 \leq \dots \leq \hat{C}_n$.

We now construct a feasible solution to LP-primal. For each coflow k , we set $C_k^p = \hat{C}_k$ for each $p \in [2m]$ and set $C_k = \hat{C}_k$ as well. This satisfies constraints (1). Constraints (2) follow since no coflow can complete before its release time. To see that constraints (3) hold, consider any set $S \subseteq [n]$ and $p \in [2m]$; then, since coflow k finishes after all of the coflows $\{1, \dots, k-1\}$ have finished, it must be true that $\hat{C}_k \geq \sum_{\kappa=1}^k v_\kappa^p$. Then,

$$\begin{aligned} \sum_{k \in S} v_k^p C_k^p &= \sum_{k \in S} v_k^p \hat{C}_k \geq \sum_{k \in S} v_k^p \sum_{\kappa \in S, \kappa \leq k} v_\kappa^p \\ &= \frac{1}{2} \left[\left(\sum_{k \in S} v_k^p \right)^2 + \sum_{k \in S} (v_k^p)^2 \right]. \end{aligned}$$

To complete the proof, we note that the objective value of LP-primal is $\sum_{k=1}^n w_k C_k = \sum_{k=1}^n w_k \hat{C}_k$ and so every feasible schedule can be mapped to a feasible solution to LP-primal with the same value of the objective function. This proves that there exists a feasible solution to LP-primal whose value is equal to $\sum_{k=1}^n w_k C_k^{\text{OPT}}$, and clearly $\sum_{k=1}^n w_k C_k^$ is at most this value.*

We now consider the dual of the linear program described above to analyze the combinatorial algorithm used to schedule coflows. We introduce a dual variable α_k^p for each coflow k and each port p , and a variable β_S^p for each port p and each subset of coflows $S \subseteq [n]$, to obtain the following linear

program:

$$\begin{aligned} \max \quad & \sum_p \sum_k \gamma_k^p r_k + \frac{1}{2} \sum_S \sum_p \beta_S^p \left[\left(\sum_{j \in S} v_k^p \right)^2 + \sum_{j \in S} \left(v_k^p \right)^2 \right] \\ \text{s.t.} \quad & \sum_p \alpha_k^p \leq w_k, \quad \forall k \in [n]; \end{aligned} \quad (4)$$

$$\begin{aligned} & \sum_{S \ni k} \beta_S^p v_k^p \leq \alpha_k^p - \gamma_k^p, \quad \forall p \in [2m], k \in [n]; \\ & \alpha_k^p \geq 0, \quad \forall p \in [2m], k \in [n]; \\ & \beta_S^p \geq 0, \quad \forall p \in [2m], S \subseteq [n]. \end{aligned} \quad (5)$$

The primal-dual algorithm is described in Algorithm 3.

Algorithm 3 Greedy algorithm for scheduling coflows with release dates

$J = [n]$ ▷ Initial set of unscheduled coflows
 ▷ $v_k(i, j)$ is the size of flow between ports i and j and belonging to coflow k with weight w_k
 ▷ Flow sum at ingress ports for coflow k
 $v_k^p = \sum_j v_k(i, j) \quad \forall p \in \{1, 2, \dots, m\}$
 ▷ Flow sum at egress ports for coflow k
 $v_k^p = \sum_i v_k(i, j) \quad \forall p \in \{m+1, m+2, \dots, 2m\}$
procedure COFLOWPERMUTATION(J)
for $k = n$ **to** 1 **do**
 ▷ Find the bottleneck port
 $b_k \leftarrow \arg \max_p \sum_{j \in J} v_j^p$
 ▷ Find the latest release date
 $r_max \leftarrow \max_{j \in J} r_j$
 if $r_max \leq \frac{1}{2} \sum_{j \in J} v_j^{b_k}$ **then**
 ▷ Select weighted largest coflow to schedule
 $\sigma(k) \leftarrow \arg \min_{j \in J} \frac{w_j - \sum_{\ell > k} \beta_\ell v_j^{b_\ell}}{v_j^{b_k}}$
 $\beta_k \leftarrow \frac{w_{\sigma(k)} - \sum_{\ell > k} \beta_\ell v_{\sigma(k)}^{b_\ell}}{v_{\sigma(k)}^{b_k}}$
 else ▷ Coflow with large release date exists
 ▷ Schedule latest released coflow
 $\sigma(k) \leftarrow \arg \max_{j \in J} r_j$
 $\gamma_{\sigma(k)} = w_{\sigma(k)} - \sum_{\ell > k} \beta_\ell v_{\sigma(k)}^{b_\ell}$
 ▷ Update set of unscheduled jobs
 $J \leftarrow J \setminus \sigma(k)$
return σ ▷ Output the coflow permutation

For ease of analysis, assume that the coflows are numbered so that $\sigma(k) = k, \forall k \in [n]$. If $(C_k^{\text{OPT}})_{k \in [n]}$ are the optimal completion times given the set of coflows $\{1, \dots, n\}$, then we first prove bounds for the permutation produced by the above algorithm.

Lemma 2. Algorithm 3 produces a permutation of coflows σ such that

$$\sum_{k=1}^n w_k \left(\max_{j \leq k} r_j + 2 \sum_{j \leq k} v_j^{b_k} \right) \leq 5 \sum_k w_k C_k^{\text{OPT}}.$$

Following this, we prove an upper bound on the completion times of the coflows on using a rate allocation scheme that satisfies the properties of being work conserving (no idle time) and preemptive.

Lemma 3. Let $\hat{C}_1, \dots, \hat{C}_n$ be the completion times of the coflows upon using any work conserving and flow-level preemptive rate allocation algorithm that ensures that it maintains the order returned by the primal-dual algorithm on each set of input-output ports. Then,

$$\hat{C}_k \leq r_k + 2 \sum_{j \leq k} v_j^{b_k}.$$

Since a greedy rate allocation scheme on the permutation order given by the first algorithm satisfies the conditions of the Lemma 3, it is easy to see that by combining Lemma 2 and Lemma 3, we obtain Theorem 1.

Proof 2. Lemma 2 We analyze the primal-dual algorithm by constructing a feasible dual solution and comparing the objective function value to the quantity of interest on the left hand side, $\sum_{k=1}^n w_k (\max_{j \leq k} r_j + 2 \sum_{j \leq k} v_j^{b_k})$.

First set the new dual variables

$$\gamma_k^p = \begin{cases} \gamma_k, & \text{if } p = b_k; \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

Note that the γ_k^p variables are non-zero if the **else** condition in the algorithm is triggered, and 0 otherwise. Furthermore, in the case that there is a coflow with a large release date, the variable corresponding to it is non-zero only for the port that has the most load at time of scheduling.

We set the dual variables

$$\alpha_k^p = \begin{cases} \sum_{\ell > k | b_\ell = p} \beta_\ell v_k^{b_\ell}, & \text{if } p \neq b_k; \\ w_k - \sum_{\ell > k | b_\ell \neq b_k} \beta_\ell v_k^{b_\ell}, & \text{if } p = b_k. \end{cases} \quad (7)$$

We set the variables

$$\beta_S^p = \begin{cases} \beta_k, & \text{if } S = [k] \text{ and } p = b_k; \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

These variables are now non-zero only if S represents the set of unscheduled coflows J at some iteration, and in that iteration, the **if** condition was satisfied in the algorithm, i.e., all the release dates were sufficiently small.

We observe some crucial properties of these variables

Lemma 4. The dual variables α_k^p, β_S^p , and γ_k^p as defined in (7), (8), and (6), respectively, have the following properties:

- $\alpha_k^p \geq 0, \forall k \in [n], p \in [2m];$
- $\beta_S^p \geq 0, \forall S \subseteq [n], p \in [2m];$
- $\gamma_k^p \geq 0, \forall k \in [n], p \in [2m];$
- the dual constraints (4) are satisfied with equality;

- (e). the dual constraints (5) are satisfied with equality;
 (f). when γ_k^p is non-zero, $r_k \geq \frac{1}{2} \sum_{j \leq k} v_j^p$ and when $\beta_{[k]}^p$ is non-zero, $\max_{j \leq k} r_j \leq \frac{1}{2} \sum_{j \leq k} v_j^p$.

Proof 3. (a,b). We prove this by induction on the iteration number k . When $k = n$, we have $\alpha_n^{b_n} = w_n \geq 0$ and $\alpha_n^p = 0, \forall p \neq b_n$. Furthermore, $\beta_{[n]}^{b_n} = \beta_n = \frac{w_n}{v_n^{b_n}} \geq 0$ and $\beta_{[n]}^p = 0, \forall p \neq b_n$.

When $k < n$, it is clear that $\alpha_k^p \geq 0, \forall p \neq b_k$, since by induction we know that $\beta_\ell \geq 0, \forall \ell > k$. We also have $\beta_{[k]}^p = 0, \forall p \neq b_k$. Now for the remaining variables $\alpha_k^{b_k}$ and $\beta_{[k]}^{b_k} (= \beta_k)$, we know that coflow $k + 1$ was chosen as the minimizer of the modified weight-to-processing time ratio at iteration $k + 1$, which gives us

$$\frac{w_k - \sum_{\ell=k+2}^n \beta_\ell v_k^{b_\ell}}{v_k^{b_{k+1}}} \geq \frac{w_{k+1} - \sum_{\ell=k+2}^n \beta_\ell v_{k+1}^{b_\ell}}{v_{k+1}^{b_{k+1}}}.$$

Noting that the right hand side is just β_{k+1} and rearranging yields $w_k - \sum_{\ell=k+1}^n \beta_\ell v_k^{b_\ell} \geq 0$. Dividing this by $v_k^{b_k}$ leads to $\beta_{[k]}^{b_k} \geq 0$, and from the induction hypotheses, since $\beta_\ell \geq 0, \forall \ell > k$, dropping terms from the summation only increases its value, and therefore $\alpha_k^{b_k} = w_k - \sum_{\ell > k | b_\ell \neq b_k} \beta_\ell v_k^{b_\ell} \geq w_k - \sum_{\ell=k+1}^n \beta_\ell v_k^{b_\ell} \geq 0$.

- (c). This follows from the fact that $w_k \geq \sum_{i > k} \beta_i v_k^{b_i}$, as noted in the proof of part (a, b).
 (d). The definition of α_k^p immediately yields the result $\sum_p \alpha_k^p = w_k, \forall k$.
 (e). Recall that β_S^p is non-zero only when $S = [k]$ and $p = b_k$, for each $k \in [n]$; furthermore, $\gamma_k^p = 0$ when $p \neq b_k$. Therefore, when $p \neq b_k$, we have $\gamma_k^p + \sum_{S \ni k} \beta_S^p v_k^p = \sum_{\ell > k | b_\ell = p} \beta_\ell v_k^{b_\ell} = \alpha_k^p$ by definition. If $p = b_k$ and $\gamma_k = 0$, then we have,

$$\begin{aligned} \sum_{S \ni k} \beta_S^{b_k} v_k^{b_k} &= \beta_k v_k^{b_k} + \sum_{\ell > k | b_\ell = b_k} \beta_\ell v_k^{b_\ell} \\ &= w_k - \sum_{\ell > k} \beta_\ell v_k^{b_\ell} + \sum_{\ell > k | b_\ell = b_k} \beta_\ell v_k^{b_\ell} \\ &= w_k - \sum_{\ell > k | b_\ell \neq b_k} \beta_\ell v_k^{b_\ell} \\ &= \alpha_k^{b_k}, \end{aligned}$$

where the second equality follows from the definition of β_k .

For the remaining case, we look at port $p = p_k$ and assume that in iteration k , we have $\gamma_k \neq 0$. Then we

know that $\beta_k = 0$, and so

$$\begin{aligned} \gamma_k^{b_k} + \sum_{S \ni k} \beta_S^{b_k} v_k^{b_k} &= \gamma_k^{b_k} + \sum_{\ell > k | b_\ell = b_k} \beta_\ell v_k^{b_\ell} \\ &= w_k - \sum_{\ell > k} \beta_\ell v_k^{b_\ell} + \sum_{\ell > k | b_\ell = b_k} \beta_\ell v_k^{b_\ell} \\ &= w_k - \sum_{\ell > k | b_\ell \neq b_k} \beta_\ell v_k^{b_\ell} \\ &= \alpha_k^{b_k}, \end{aligned}$$

where the second equality follows from the definition of $\gamma_k (= \gamma_k^{b_k})$.

- (f). This is evident from the **if-else** condition in the algorithm. (Note that γ_k^p and $\beta_{[k]}^p$ are zero for all $p \neq b_k$).

Now note that

$$\begin{aligned} &\sum_{j=1}^n w_k (\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j}) \\ &\stackrel{\text{Lemma 4(d)}}{=} \sum_{j=1}^n \left(\sum_p \alpha_j^p \right) (\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j}) \\ &\stackrel{\text{Lemma 4(e)}}{=} \sum_j \left(\sum_p \gamma_j^p + \sum_p \sum_{S \ni j} \beta_S^p v_j^p \right) (\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j}). \end{aligned} \quad (9)$$

We deal with the two terms separately. First,

$$\begin{aligned} \sum_j \sum_p \gamma_j^p (\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j}) &\stackrel{(a)}{=} \sum_j \sum_p \gamma_j^p (r_j + 2 \sum_{i \leq j} v_i^{b_j}) \\ &\stackrel{(b)}{\leq} 5 \sum_j \sum_p \gamma_j^p r_j, \end{aligned} \quad (10)$$

where (a) follows from the fact that when γ_k is non-zero (i.e., the **else** condition in the algorithm is triggered), the coflow scheduled in that iteration is the one with the maximum release date amongst the set of coflows yet to be scheduled, and (b) follows from Lemma 4(f).

Now looking at the second term,

$$\begin{aligned} &\sum_j \left(\sum_p \sum_{S \ni j} \beta_S^p v_j^p \right) \left(\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j} \right) \\ &= \sum_j \left(\sum_{k \geq j} \beta_{[k]}^{b_k} v_j^{b_k} \right) \left(\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j} \right) \\ &= \sum_k \beta_{[k]}^{b_k} \sum_{j \leq k} v_j^{b_k} \left(\max_{i \leq j} r_i + 2 \sum_{i \leq j} v_i^{b_j} \right) \\ &\leq \sum_k \beta_{[k]}^{b_k} \sum_{j \leq k} v_j^{b_k} \left(\max_{i \leq k} r_i + 2 \sum_{i \leq k} v_i^{b_k} \right) \\ &\stackrel{(a)}{\leq} \frac{5}{2} \sum_k \beta_{[k]}^{b_k} \sum_{j \leq k} v_j^{b_k} \sum_{i \leq k} v_i^{b_k} = \frac{5}{2} \sum_k \beta_{[k]}^{b_k} \left(\sum_{j \leq k} v_j^{b_k} \right)^2 \end{aligned}$$

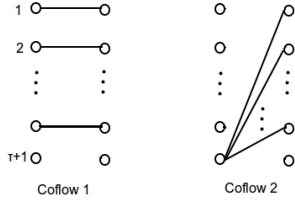


Figure 11: Coflow Instance

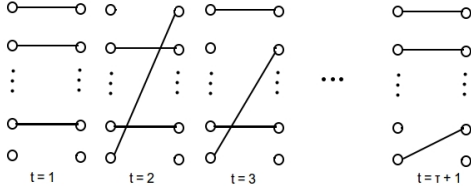


Figure 12: Optimal rate allocation

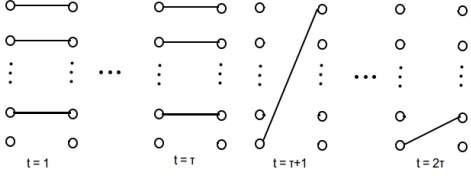


Figure 13: Poor rate allocation

$$\begin{aligned}
 &\leq \frac{5}{2} \sum_k \beta_{[k]}^{b_k} \left(\sum_{j \leq k} [v_k^{b_k}]^2 + \left[\sum_{j \leq k} v_k^{b_k} \right]^2 \right) \\
 &= \frac{5}{2} \sum_p \sum_s \beta_s^p \left(\sum_{j \leq k} [v_k^{b_k}]^2 + \left[\sum_{j \leq k} v_k^{b_k} \right]^2 \right), \quad (11)
 \end{aligned}$$

where (a) follows from Lemma 4(f).

Combining (10) and (11) with (9) and noting from weak duality that

$$\sum_p \sum_k \gamma_k^p r_k + \frac{1}{2} \sum_s \sum_p \beta_s^p \left[\left(\sum_{j \in S} v_k^p \right)^2 + \sum_{j \in S} \left(v_k^p \right)^2 \right] \leq \sum_k w_k C_k^{\text{OPT}}$$

for any feasible dual solution, we get the desired result.

Proof 4. Lemma 3

Consider any rate allocation algorithm \mathcal{A} that is work conserving and follows the priorities established by the primal-dual algorithm. More specifically, given any flow from some input port p_i to output port p_o in a coflow k , no flow (p_i, p_o) from a later coflow $\ell : \ell > k$ is scheduled before all of the packets in this flow for coflow k have completed.

Now consider any coflow k , and let \hat{C}_k be its completion time. Since the release time of this coflow is r_k , the flows of this coflow can only be scheduled after r_k . Since the coflow completes at time \hat{C}_k , there exists a packet from port p_i to port p_o that is sent at time \hat{C}_k . It follows from work conservation that this

occurs because the ports p_i and p_o were never simultaneously free during the time interval $[r_k, \hat{C}_k)$. In fact, because of the strict adherence to the priorities as mentioned above, we can see that the ports p_i and p_o are neither free nor serving a flow from a coflow $\ell > k$ during this time interval (since otherwise, the corresponding flow from coflow k would have been picked by the algorithm).

Given the above statements, it immediately follows that at least one of the ports among p_i and p_o is busy serving flows from the set of coflows $\{1, \dots, k\}$ for at least $\frac{1}{2}(\hat{C}_k - r_k)$ amount of time in the interval $[r_k, \hat{C}_k)$. Now noting that the total time a port p can be busy serving flows from a set S is $\sum_{j \in S} v_j^p$ and thus we get $\frac{1}{2}(C_k - r_k) \leq \max\{\sum_{j \leq k} v_j^{p_i}, \sum_{j \leq k} v_j^{p_o}\}$. Combining with the fact that b_k was the most heavily loaded port and thus $\max\{\sum_{j \leq k} v_j^{p_i}, \sum_{j \leq k} v_j^{p_o}\} \leq \sum_{j \leq k} v_j^{b_k}$ and rearranging completes the proof.

We can improve this result for the case when there are no release dates, leading to a 4-approximation algorithm.

If $(C_k^{\text{OPT}})_{k \in [n]}$ are the optimal completion times given the set of coflows $\{1, \dots, n\}$, then

Lemma 5. When there are no release dates, Algorithm 3 produces a permutation of coflows σ such that

$$\sum_{k=1}^n w_k \left(\sum_{j \leq k} v_j^{b_k} \right) \leq 2 \sum_k w_k C_k^{\text{OPT}}$$

PROOF. Noting that all the γ_k^p dual variables are 0, the proof of Lemma 2 yields this result. \square

Combining this with the fact that the completion time of coflow k under the greedy rate allocation scheme is at most $2 \sum_{j \leq k} v_j^{b_k}$ (follows from Lemma 3 by setting $r_k = 0, \forall k$), we obtain the 4-approximation guarantee.

An important point to note is that the only properties we used about the rate allocation scheme was that it was work conserving and that it preserved the order of the coflows returned by the primal-dual algorithm while scheduling. Thus any rate allocation which maintains these two properties provides, together with the primal-dual algorithm to provide an ordering, a 4-approximation algorithm to the coflow scheduling problem.

We now show that the bound given in Theorem 2 is tight.

Let $\tau = \lceil \frac{3}{\epsilon} \rceil$ and consider the coflow instance shown in Figure 11 on $\tau + 1$ ports and the permutation being that coflow 1 has higher priority. Let the weight of the first coflow be $w_1 = 1$ and that of the second coflow be $w_2 = \tau$. Now consider the optimum rate allocation scheme that is work conserving, pre-emptive and order preserving, shown in Figure 12. The weighted completion time for this scheme is $(1 + \tau)^2$. Now consider another rate allocation scheme within

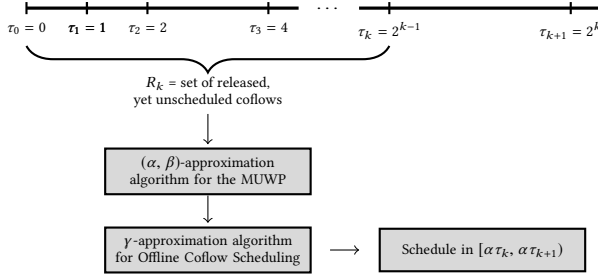


Figure 14: Algorithm to provide a $(2\alpha\beta + \gamma)$ -competitive online algorithm, as described in Khuller et al.

the same class, shown in Figure 13. The weighted completion time of this scheme is $\tau + 2\tau^2$. The ratio of this value to the optimal completion time is more than $2 - \varepsilon$.

10 ONLINE APPROXIMATION ALGORITHM

For the online setting, we use the generalization of the framework of Hall et al. [15] by Khuller et al. [20] for a 12-competitive online algorithm and a 9.78-competitive randomized algorithm for this problem. Their framework makes use of a γ -approximation algorithm to the offline variant of the same problem and a (α, β) -approximation to the Minimum Unscheduled Weight Problem (MUWP) to obtain a $2\alpha\beta + \gamma$ -competitive algorithm.

The MUWP with respect to the scheduling environment of coflows is as follows.

Definition 4. [MINIMUM UNSCHEDULED WEIGHT PROBLEM] *Given a set of coflows available at time 0, a weight for each coflow, and a deadline D , find a subset of coflows that can be scheduled to complete by time D while minimizing the weight of unscheduled coflows.*

Definition 5. *An algorithm is an (α, β) -approximation algorithm to the MUWP if it returns a set of coflows such that*

they can be scheduled to complete by time αD and ensures that the weight of the unscheduled coflows is at most β times the optimum solution to the MUWP with deadline D .

We next briefly describe the algorithm used. As in Figure 14, divide the time horizon into a sequence of intervals whose lengths increase geometrically; let $\tau_0 = 0$, and for $k \geq 1$, $\tau_k = 2^{k-1}$. For each $k \geq 0$, let R_k be the set of coflows that have been released by τ_k but not yet scheduled. Then at each time τ_k , run the (α, β) -approximation algorithm for MUWP on the set R_k . We then run the γ -approximation algorithm on the resulting set of coflows and schedule them to run in the interval $[\alpha\tau_k, \alpha\tau_{k+1})$. Khuller et al. prove that this algorithm is $(2\alpha\beta + \gamma)$ -competitive, with an added term of αW where W is the total weight of the coflows.

Using the $(2, 2)$ -approximation for the MUWP in the setting of order scheduling of Garg et al. [14] and the 4-approximation for the offline scheduling problem in this paper leads to a 12-competitive algorithm for online coflow scheduling.

The randomized algorithm achieves a better bound by modifying the framework so that the intervals lengths are not chosen deterministically. To be more specific, they take $\tau_k = \eta 2^k$ where $\eta = 2^{-X}$ and X is chosen uniformly from $(0, 1]$. If B_j denotes the start of the interval in which coflow j is scheduled, then their framework produces a schedule with the weighted completion time being at most $2\alpha\beta \sum_j w_j B_j + \gamma \sum_j w_j C_j^{\text{OPT}}$ where C_j^{OPT} is the completion time of coflow j in an optimal schedule.

Now using the result by Chakrabarti et al. [5] that $E[B_j] = \frac{1}{2 \ln 2} C_j^{\text{OPT}}$ when the intervals are chosen as specified above, and using linearity of expectation, Khuller et al. arrive at an algorithm that is $\left(\frac{1}{\ln 2} \alpha\beta + \gamma\right)$ -competitive in expectation. Using the values of α, β and γ as described above leads to a 9.78-competitive randomized algorithm for online coflow scheduling.