# Efficient Inter-Datacenter Bulk Transfers with Mixed Completion Time Objectives[☆]

Mohammad Noormohammadpour

*University of Southern California*

Srikanth Kandula

*Microsoft*

Cauligi S. Raghavendra

*University of Southern California*

Sriram Rao

*Facebook*

**Abstract**

Bulk transfers from one to multiple datacenters can have many different completion time objectives ranging from quickly replicating some $k$ copies to minimizing the time by which the last destination receives a full replica. We design an SDN-style wide-area traffic scheduler that optimizes different completion time objectives for various requests. The scheduler builds, for each bulk transfer, one or more multicast forwarding trees which preferentially use lightly loaded network links. Multiple multicast trees are used per bulk transfer to insulate destinations that have higher available bandwidth and can hence finish quickly from congested destinations. These decisions–how many trees to construct and which receivers to serve using a given tree–result from an optimization problem that minimizes a weighted sum of transfers' completion time objectives and their bandwidth consumption. Results from simulations and emulations on Mininet show that our scheduler, Iris, can improve different completion time objectives by about $2.5\times$.

*Keywords:* Replication, Bulk Multicast, Traffic Engineering, Inter-Datacenter

*May 30, 2022*

arXiv:1905.01749v3 [cs.DC] 15 Sep 2019

## 1. Introduction

A wide range of distributed applications replicate content and data to increase end-users' quality of experience [1, 2, 3, 4, 5, 6, 7, 8] which results in inter-datacenter bulk multicast transfers with a given set of receivers. For a variety of applications, objects may be replicated to at least four datacenters and for some applications potentially to tens of datacenters [9, 10]. Moreover, an analysis of Baidu's traffic [11] across 30 datacenters showed that over 90% of the traffic is multicast and over 90% of the multicast transfers are destined to at least 60% of datacenters.

A variety of approaches can be used to perform bulk multicast transfers. We can model a bulk multicast transfer as multiple independent unicast bulk transfers [12, 13, 14, 4] which wastes network capacity and can increase the transfer completion times. Standard internet multicasting [15] builds multicast trees incrementally as receivers join the multicast session without considering the distribution of traffic load across network links. Therefore, generated multicast trees can be considerably larger than necessary with highly unbalanced network load distribution. Overlay multicasting [16] builds application layer multicast trees which may be far from the optimal due to limited visibility of network link level status and little control of traffic routing at the network layer. Peer-to-peer file distribution [17, 18] aims to maximize throughput per receiver in a decentralized fashion and greedily, which can be far from global optimization. Centralized multicast routing approaches allow for better multicast tree selection by incorporating a global view of network status. Some centralized methods, such as [19, 20], target the regular and structured topologies of networks inside datacenters, which cannot be directly applied to inter-datacenter networks. Many other centralized techniques, such as [21, 22, 23, 24, 25, 26, 27], do not consider the optimization of receiver completion times, especially in an online scenario with many concurrent bulk multicast transfers. Finally, very recently, several centralized proposals aim to optimize the completion times of bulk multicast receivers [9, 28]. Our work in this paper builds on these techniques.

When receivers of a bulk multicast transfer have very different network bandwidth available on paths from the sender, the slowest receiver dictates the completion time for all receivers. Recent work suggests using multiple multicast trees to separate the faster receivers

which will improve the average receiver's completion time [28]. However, each additional tree consumes more network bandwidth and at the extremum, this idea devolves to one tree per receiver. We aim to answer the following questions:

1. What is the right number of trees per transfer?
2. Which receivers should be grouped in each tree?

We analyze a relaxed version of this partitioning problem where each partition is a subset of receivers attached to the sender with a separate forwarding tree. We first propose a partitioning technique that reduces the average receiver completion times of receivers by isolating slow and fast receivers. We study this approach in the relaxed setting of having a congestion-free network core, i.e., links in/out of the datacenters are the capacity bottlenecks, and considering max-min fair rate allocation from the underlying network. We then develop a partitioning technique for real-world inter-datacenter networks, without relaxations, and inspired by the findings from studying the relaxed scenario. The partitioning technique operates by building a hierarchy of valid partitioning solutions and selecting the one that offers the best average receiver completion times. Our evaluation of this partitioning technique on real-world topologies, including ones with bottlenecks in the network core, show that the technique yields completion times that are close to a lower bound and hence nearly optimal.

Back-end geo-distributed applications running on datacenters can have different requirements on how their objects are replicated to other datacenters. Hence, inter-datacenter traffic is usually a mix of transfers with various completion time objectives. For example, while reproducing $n$ copies of an object to $n$ different datacenters/locations, one application may want to transfer $k$ copies quickly to any among $n$ given receivers, and another application may want to minimize the time when the last copy finishes. In the former case, grouping the slower $n - k$ receivers into one partition consumes less bandwidth and this spare bandwidth could be used to speed up the other transfers. In the latter case, by grouping all receivers except the slowest receiver together to use one tree, we can isolate them from the slowest receiver with minimal bandwidth consumption. Minimizing the completion times of all receivers is another possible objective. Our technique takes as input a binary objective vector whose $i^{\text{th}}$ element expresses interest in the completion time of the $i^{\text{th}}$ fastest receiver;

it aims to minimize the completion times of receivers whose rank is set to one in this objective vector. It is easy to see that following values of the objective vector achieve the goals discussed so far; when $k = 1$, $n = 3$, $\{1, 0, 0\}$, $\{0, 0, 1\}$ and $\{1, 1, 1\}$ aim to minimize the completion time of the fastest $k$ out of $n$ receivers, the slowest receiver, and all receivers, respectively.

We have built a system called Iris which combines the proposed partitioning technique and the application/user supplied objective vectors. It operates in a logically centralized manner, receives bulk multicast transfer requests from end-points, and computes receiver partitions along with their multicast forwarding trees. We create forwarding trees using group tables [29]. Iris uses a RESTful API to communicate with the end-points allowing them to specify their transfer properties and requirements (i.e., objective vectors) using which it computes and installs the required rules in the forwarding plane. We believe our techniques are easily applicable in today's inter-datacenter networks [4, 30, 2]. Our contributions can be summarized as follows:

- We propose a partitioning approach that reduces the effect of slow receivers by isolating them and attaching them using independent paths. We discuss various scenarios where this approach offers different levels of performance compared to the optimal partitioning on relaxed network topology and given max-min fair rate allocation.

- We incorporate binary objective vectors which allow applications to indicate transfer-specific objectives for receivers' completion times. Using the application-provided objective vectors, we can optimize for mixed completion time objectives based on the trade-off between total network capacity consumption and the receivers' average completion times.

- We present the Iris heuristic, which computes a partitioning of receivers for every transfer given a binary objective vector. Iris aims to minimize the completion time of receivers whose rank is indicated by applications/users with a one in the objective vector while saving as much bandwidth as possible by grouping receivers whose ranks are indicated with consecutive zeros in the objective vector.

4

- We perform extensive simulations and Mininet emulations with Iris using synthetic and real-world Facebook inter-datacenter traffic patterns over large WAN topologies. Simulation results show that Iris speeds up transfers to a small number of receivers (e.g., $\geq 8$ receivers) by $\geq 2\times$ on the average completion time while the bandwidth used is $\leq 1.13\times$ compared to state-of-the-art. Transfers with more receivers receive larger benefits. For transfers to at least 16 receivers, 75% of the receivers complete at least $5\times$ faster and the fastest receiver completes $2.5\times$ faster compared to state-of-the-art. Compared to performing multicast as multiple unicast transfers with shortest path routing, Iris reduces mean completion times by about $2\times$ while using $0.66\times$ of the bandwidth. Finally, Mininet emulations show that Iris reduces the maximum group table entries needed by up to $3\times$.

## 2. Background and related Work

**Point to Multipoint (P2MP) Bulk Transfers:** Similar to bulk multicast, P2MP transfers push data and content from one location to multiple locations. The transfer size, the source and the receivers are known and fixed prior to initiation of the transfer. Load-aware forwarding trees [9, 28, 31, 27], flexible source selection [26], and store-and-forward [32, 33] have been used in recent work to save bandwidth and speed up transfer completion. Several works around this subject focus on admission control for multicast transfers with deadlines [26, 27] which is orthogonal to our work in this paper. Other works focus on transferring large volumes of data with minimum increase in the network bandwidth cost [32, 33]. We build our work in this paper on top of recent work on optimizing the completion times of P2MP transfers [9, 28].

**Network-layer Multicast:** A vast variety of network-layer multicast solutions have been proposed [15, 34, 35, 36, 37, 38, 39]. In general, these solutions consider the dynamic scenarios where receivers may join or leave at any time; hence, they greedily adapt the multicast distribution tree as the receiver set evolves. The problem considered in this paper differs in the following key ways: we assume a known and fixed transfer size and set of receivers, we assume SDN-style visibility and control on the network routes, and we support for general completion time objectives. Compared to general multicast solutions that build multicast

trees incrementally and greedily as new receivers join, given a known and fixed set of receivers, we can use the global knowledge of network load distribution and topology to select further optimized multicast trees that reduce the number of bottlenecks and improve the receiver completion times.

**End-system Multicast:** These works form multicast trees in the application layer among the participating end-points [40, 16, 41, 42, 43, 11]; they have limited visibility and control of the underlying network, i.e., they cannot easily change routes, identify available bandwidth along network paths, etc. Therefore, depending on the network topology and the underlying routing approach, these techniques may offer solutions far from the optimal in minimizing the completion times of receivers.

**Datacenter Multicast:** Recently, some works propose using multicast trees within and between datacenters [20, 19, 44]. These approaches, however, operate specifically on datacenter networks that have regular and structured topologies whereas inter-datacenter networks are usually neither structured nor regular. Besides, these solutions do not aim at optimizing the completion times of multicast receivers.

**Reliable transport protocols for multicast:** Reliable multicast schemes ensure that data is completely and correctly received by all receivers using a variety of techniques such as FEC codes [45, 46, 36, 47, 48, 49], retransmissions [50, 44, 36], etc. Lost data can be detected using ACKs or NACKs. Iris is orthogonal to and can use any multicast transport protocol.

**Bit Indexed Explicit Replication:** A recent proposal, BIER [51], encodes forwarding state in the packet headers which simplifies network forwarding and improves scalability on large networks. BIER allows changes to multicast trees with a small overhead. Iris can adopt BIER for creating and updating multicast trees to reduce the cost of rule installations and updates from SDN Group Tables.

**Store-and-Forward (SnF):** SnF techniques [33, 52, 32, 53] have been proposed for unicast delay tolerant transfers to avoid periods of network congestion; a recent work called BDS uses SnF for bulk multicast over unicast TCP connections that connect receivers in a line [11] to increase network utilization given diurnal patterns of available capacity on backbone links.

SnF, however, increases the protocol complexity and can incur additional bandwidth and storage costs on intermediate datacenters. Besides, formulating the bulk multicast problem using SnF will result in a completely new model as the data transmission rate across the edges of a multicast tree may be different which can change the nature of the problem in several ways. First, additional optimization metrics should be considered, among which are the total network storage and maximum per-node storage budget over multicast trees. Next, in case there is no limit on how much data can be stored per-node in a multicast tree, the slow receivers will automatically be isolated from the fast receivers over the tree. However, given that we most likely will have storage budgets per intermediate node on a multicast tree, this may create a complex relationship between the download speeds of fast and slow receivers. Therefore, modeling bulk multicasting with SnF can generate highly complex linear programs solving which may be slow. Application of SnF for bulk multicasting to optimize receiver completion times is considered part of the future work.

**Peer-to-Peer File Distribution:** These techniques [54, 18, 17, 55] function locally and greedily and cannot take advantage of SDN-style visibility and control of the inter-datacenter networks for global optimization of transfers among datacenters.

## 3. System Model

We focus on proprietary inter-datacenter networks that connect geographically dispersed datacenters such as Microsoft Global WAN [30], Facebook Express Backbone [2] and Google GScale [4]. These networks are managed by one organization and their forwarding state can be managed in a logically centralized manner. A traffic engineering server that runs Iris algorithm decides how traffic is forwarded in-network similar to other related work [4, 14, 12, 13].

Cross datacenter traffic can, in general, be categorized as high priority user traffic that is highly sensitive to latency and internal traffic (also known as elastic or background traffic [14, 12]) that is more resilient to latency. Internal traffic constitutes the majority of cross datacenter traffic, consists of huge volumes of replicated data and content that generate long-running transfers, and is growing at a much faster pace than user-generated traffic [2]. By forwarding such traffic according to transfer properties (i.e., end-points and volume)

and network topology (i.e., connectivity and available bandwidth) we can optimize some network-wide utility. Wide-area traffic management is a complex problem and in general a variety of metrics can be considered for optimization [56]. We focus on internal traffic that is a result of data and content replication which can be modeled as bulk multicast transfers. These transfers are processed in an online manner as they arrive with the main objective of optimizing completion times. Also, forwarding entries, which are installed for every transfer upon arrival, are fixed until transfers' completion and are only updated in case of failures. Finally, we assume that all multicast transfers are of the same priority. Extension of the proposed solutions to a scenario where transfers are of different value to the operator/client(s) is considered as part of the future work.

We consider max-min fair [57] rate allocation across multicast forwarding trees. Traffic is transmitted with the same rate from the source to all the receivers attached to a forwarding tree. To reach max-min fair rates, such rates can either be computed centrally over specific time periods, i.e., timeslots, and then be used for end-point traffic shaping or end-points can gradually converge to such rates in a distributed fashion in a way similar to TCP [44] (fairness is considered across trees). In our evaluations, we will consider the former approach for increased network utilization. Using a fair sharing policy addresses the starvation problem (such as in SRPT policy) and prevents larger transfers from blocking edges (such as in FCFS policy). Recent work has also shown that such conditions can worsen over trees [28].

We use the notion of **objective vectors** to allow applications to define transfer-specific requirements which in general can improve overall system performance and reduce bandwidth consumption. An objective vector for a transfer is a vector of zeros and ones which is the same size as the number of receivers of that transfer. From left to right, the binary digit $i$ in this vector is associated with the $i^{\text{th}}$ fastest receiver. A one in the objective vector indicates that we are specifically interested in the completion time of the receiver associated with that rank in the vector. By assigning zeros and ones to different receiver ranks, it is possible to respect different applications' preferences or requirements while allowing the system to optimize bandwidth consumption further. The application/user, however, needs not be aware of the mapping between the downlink speeds (rank in the objective vector) and the receivers themselves.

Table 1: Behavior of Several Objective Vectors

| Objective Vector ($\omega$) | Outcome (given $n$ receivers) |
|---|---|
| $\{1, \ldots, 1\}$ over $n$ | Interested in completion times of all individual receivers |
| $\{1, \ldots, 1, 0, \ldots, 0\}$ over $k$ and $n\text{-}k$ | Interested in completion times of the top $k$ receivers (groups the bottom $n - k$ receivers to save bandwidth) |
| $\{0, \ldots, 0, 1, 0, \ldots, 0\}$ over $k\text{-}1$ and $n\text{-}k$ | Interested in the completion time of the $k^{\text{th}}$ receiver (groups the top $k - 1$ receivers into a fast partition, and the bottom $n - k$ receivers into a slow one to save bandwidth) |
| $\{0, \ldots, 0\}$ over $n$ | Not interested in the completion time of any specific receiver (all receivers form a single partition) |

Table 1 offers several examples. For instance, an objective vector of $\{0, 0, 0, 1, 0, 0, 0, 0\}$ indicates the application's interest in the fourth fastest receiver. To respect the application's objective, we initially isolate the fourth receiver and do not group it with any other receiver. The first three fastest receivers can be grouped into a partition to save bandwidth. The same goes for the four slowest receivers. However, we do not group all receivers indicated with zeros into one partition initially (i.e., the top three receivers and the bottom four) to avoid slowing some of them down unnecessarily (in this case, the top three receivers). This forms the basis for the partitioning technique proposed in §5.4 that operates by building a hierarchy with multiple layers, where each layer is a valid partitioning solution, and selects the layer that gives the smallest average receiver completion times.

Although the objective vector can, in general, be any binary string of zeros and ones, it is worth noting that not all such combinations lead to meaningful objectives for datacenter applications. For example, an objective vector of $\{0, 1, 0, 1, 0, 1, 0, 1\}$ may be unlikely to be used by an application. Having the ability to define and enforce any objective function though makes the system highly configurable and adaptable. Operators may come up with a set of rules based on which they can decide whether the objective vector proposed by an application is meaningful, or propose changes to a submission that is not deemed useful.

**Problem Statement:** Given an inter-datacenter topology with known available bandwidth per link, the traffic engineering server is responsible for *partitioning receivers* and *selecting a forwarding tree per partition* for every incoming bulk multicast transfer. A bulk multicast transfer is specified by its source, set of receivers and volume of data to be delivered. The primary objective is minimizing average receiver completion times. In case an objective vector is specified, we want to minimize average completion times of receivers whose ranks are indicated with a 1 in the vector as well as receivers indicated with consecutive 0's in the vector together as groups (receivers noted with consecutive 0's use the same forwarding tree and will have the same completion times). Minimizing bandwidth consumption, which is directly proportional to the size of selected forwarding trees, is considered a secondary objective.

*3.1. Online Greedy Optimization Model*

The online bulk multicast partitioning and forwarding tree selection problem can be formulated using Eq. 1-3 with the added constraint that our rate allocation is max-min fair across forwarding trees for any selection of the partitions and the trees. Table 2 lists the variables used in the formulation below.

The set $\mathbf{R}$ includes both the new transfer $R_N$ and all the ones already in the system for which we already have the partitions and forwarding trees. The optimization objective of Eq. 1 is to minimize the weighted sum of completion times of receivers of all requests $R \in \mathbf{R}$ according to their objective vectors, and the total bandwidth consumption of $R_N$ by partitioning its receivers and selecting their forwarding trees (indicated by the term $\sum_{P \in \mathbf{P}_{R_N}} \mathcal{V}_P |T_P|$). Operators can choose the non-negative coefficient $\epsilon$ according to the overall system objective to give a higher weight to minimizing the weighted completion time of receivers than reducing bandwidth consumption. Eq. 2 shows the demand constraints which state that the total sum of transmission rates over every tree for future timeslots is equal to the remaining volume of data per partition (each partition uses one tree). Eq. 3 presents the capacity constraints which state that the total sum of transmission rates per timeslot for all trees that share a common edge has to not go beyond its available bandwidth.

Table 2: Definition of Variables

| Variable | Definition |
|---|---|
| $t_{now}$ | Current timeslot |
| $e$ | A directed edge |
| $C_e$ | Capacity of $e$ in bytes/second |
| $B_e(t)$ | Available bandwidth on edge $e$ at timeslot $t$ after setting aside usage of high priority user traffic |
| $B_e$ | Average available bandwidth on edge $e$ |
| $G$ | The directed inter-datacenter graph |
| $T$ | A directed Steiner tree |
| $\mathbf{V}_G$ and $\mathbf{V}_T$ | Set$\langle\rangle$ of vertices of $G$ and $T$ |
| $\mathbf{E}_G$ and $\mathbf{E}_T$ | Set$\langle\rangle$ of edges of $G$ and $T$ |
| $r_T(t)$ | The transmission rate over tree $T$ at $t$ |
| $\delta$ | Duration of a timeslot |
| $R$ | A bulk multicast transfer request |
| $S_R$ | Source datacenter of request $R$ |
| $A_R$ | Arrival time of request $R$ |
| $\mathcal{V}_R$ | Original volume of request $R$ |
| $\mathbf{D}_R$ | Set$\langle\rangle$ of destinations of request $R$ |
| $\mathbf{R}$ | Set$\langle\rangle$ of ongoing transfers |
| $P$ | A receiver partition of some request |
| $\mathbf{P}_R$ | Set$\langle\rangle$ of partitions of some request $R$ |
| $T_P$ | The forwarding tree of partition $P$ |
| $\mathcal{V}_P^{[res]}$ | Current residual volume of partition $P$ of request $R$ |
| $\kappa_P$ | Estimated minimum completion time of partition $P$ |
| $L_e$ | Edge $e$'s total outstanding load (see §5.1) |
| $\omega_R$ | Objective vector assigned to request $R$ |
| $\omega_R^\star$ | Weighted completion time vector computed from $\omega_R$ by replacing the last zero in a pack of consecutive zeros with the number of consecutive zeros in that pack (e.g., $\omega_R = \{0,0,0,1,0,0\} \to \omega_R^\star = \{0,0,3,1,0,2\}$) |
| $t_{\mathbf{D}_R}$ | Vector of completion times of receivers of request $R$ sorted from fastest to slowest |

$$\textbf{min} \quad \sum_{R \in \mathbf{R}} \left( t_{\mathbf{D}_R} \cdot \omega_R^\star \right) + \epsilon \sum_{P \in \mathbf{P}_{R_N}} \mathcal{V}_P |T_P| \qquad (1)$$

$$\textbf{S.t.} \quad \sum_t r_{T_P}(t) = \mathcal{V}_P^{[res]} \qquad\qquad \forall P \in \mathbf{P}_R, R \in \mathbf{R} \qquad (2)$$

$$\sum_{\{P | e \in T_P\}} r_{T_P}(t) \le B_e(t) \qquad \forall t, e, P \in \mathbf{P}_R, R \in \mathbf{R} \qquad (3)$$

This online discrete optimization problem is highly complex as it is unclear how receivers should be partitioned into multiple subsets to reduce completion times and there is an exponential number of possibilities. Selection of forwarding trees to minimize completion times is also a hard problem. In §5, we will present a heuristic that approximates a solution to this optimization problem inspired by the findings in §4.

## 4. Partitioning of Receivers on a Relaxed Topology

Due to the high complexity of the partitioning problem as a result of physical topology, we first study a relaxed topology where every datacenter is attached with a single uplink/downlink to a network with infinite core capacity and so the network core cannot become a bottleneck. As shown in Figure 1, the sender has a maximum uplink rate of $r_s$ and transmits to a set of $n$ receivers with different maximum downlink rates of $r_i, \forall i \in \{1, \ldots, n\}$. In §5.1, we discuss a load-balancing forwarding tree selection approach that aims to distribute load across the network to minimize the effect of bottlenecks within the network core.

Without loss of generality, let us also assume that the receivers in Figure 1 are sorted by their downlink rates in descending order. The sender can initiate multicast flows to any partition, i.e., a subset of receivers, given that every receiver appears in exactly one partition. All receivers in a partition will have the same multicast rate that is the rate of the slowest receiver in the partition. To compute rates at the uplink, we consider the max-min fair rate allocation policy as stated earlier in §3. In this context, we would like to compute the number of partitions as well as the receivers that should be grouped per partition to minimize mean completion times.
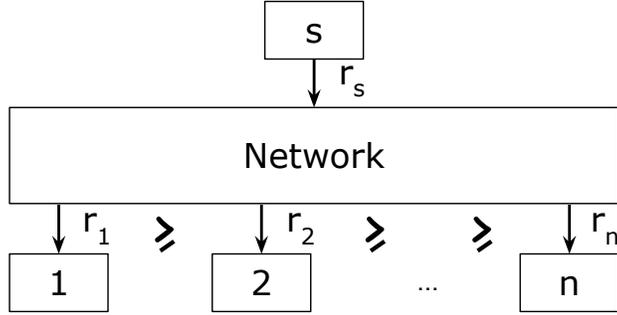
Figure 1: A relaxed topology with infinite core capacity, and uplink and downlink capacities of $r_s$ and $r_1 \geq \cdots \geq r_n$.

**Theorem 1.** *Given receivers sorted by their downlink rates, partitioning that groups consecutive receivers is pareto-optimal with regards to minimizing completion times.*

*Proof.* We use proof by contradiction. Let us assume a partitioning where non-consecutive receivers are grouped together, that is, there exist two partitions $P_1$ and $P_2$ where part of partition $P_1$ falls in between receivers of $P_2$ or the other way around. Let us call the slowest receivers of $P_1$ and $P_2$ as $j_1$ and $j_2$, respectively. Across $j_1$ and $j_2$, let us pick the fastest and call it $f(j_1, j_2)$. If $f(j_1, j_2) = j_1$ (i.e., in the non-decreasing order of downlink speed from left to right, $P_2$ appears before $P_1$ as in $P_2\{\ldots\}\ P_1\{\ldots, j_1\}\ P_2\{\ldots, j_2\}\ \ldots$), then by swapping the fastest receiver in $P_2$ and $j_1$, we can improve the rate of $P_1$ while keeping the rate of $P_2$ the same. If $f(j_1, j_2) = j_2$, then by swapping the fastest receiver in $P_1$ and $j_2$, we can improve the rate of $P_2$ while keeping the rate of $P_1$ the same. This can be done in both cases without changing the number of partitions, or number of receivers per partition across all partitions. Since the new partitioning has a higher or equal achievable rate for one of the partitions, the total average completion times will be less than or equal to that of original partitioning, which means the original partitioning could not have been optimal.

*4.1. Our Partitioning Approach*

Based on Theorem 1, the number of possible partitioning scenarios that can be considered for minimum average completion times is the number of compositions of integer $n$, that is, $2^{n-1}$ ways which can be a large space to search. To reduce complexity, we propose to isolate slow receivers from the rest of receivers to minimize their effect. In other words, given an
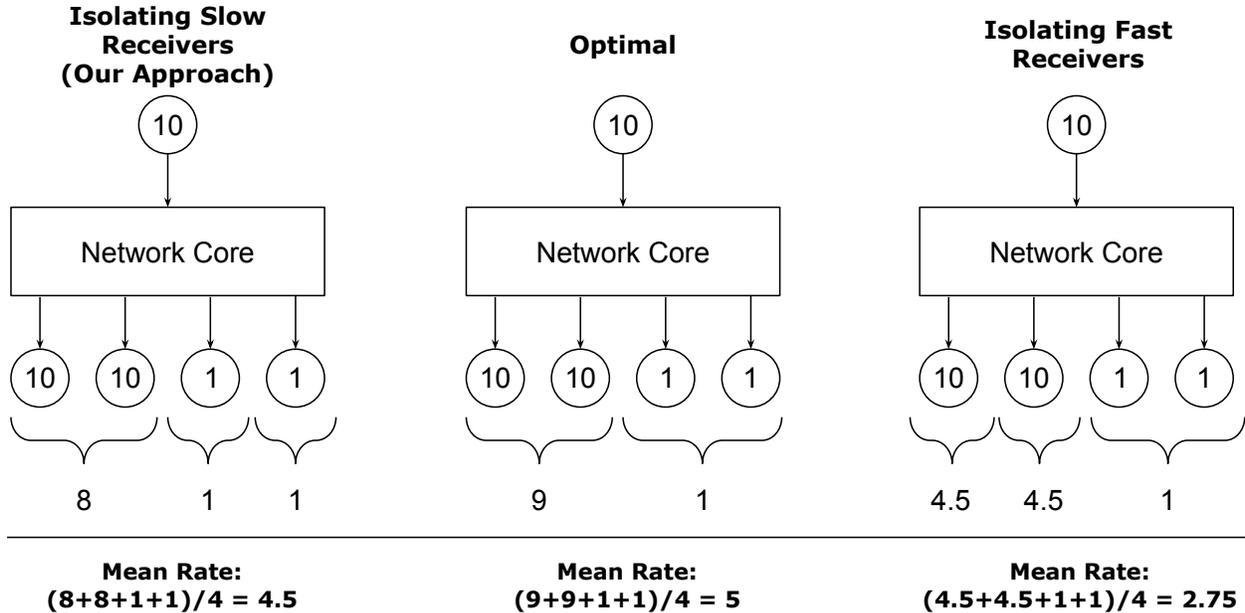
Figure 2: Various partitioning solutions for a scenario with four receivers. Numbers show the downlink and uplink speeds of nodes and curly brackets indicate the partitions where all nodes in a partition receive data at the same rate. The objective is to maximize the average rate of receivers given the max-min fairness policy.

integer $1 \leq M \leq n$, we propose to group the first $n-M+1$ fastest receivers into one partition and the rest of the receivers as separate 1-receiver partitions ($M-1$ in total). Since we do not know the value of integer $M$, we will try all possible values, that is, $n$ in total which will help us find the right threshold for the separation of fast and slow receivers. In particular, we compute the total average downlink rate of all receivers for the given transfer for every value of $M$ and select the $M$ that maximizes the average rate.[1] As shown in Figure 1, the uplink at the sender has a rate of $r_s$ which will be divided across all the multicast flows that deliver data to the receivers. Isolating a slow receiver only takes a small fraction of the sender's uplink which is why this technique is effective as we will later see in evaluations. An example of this approach and how it compares with the optimal solution is shown in Figure 2 where our solution selects $M = 3$ partitions isolating the two slow receivers.

A main determining factor in the effectiveness of this approach is how $r_s$ compares with

---

[1]Or alternatively minimizes the average completion times of receivers.

$\sum_{1 \le i \le n} r_i$. If $r_s$ is larger, then simply using $n$ partitions will offer the maximum total rate to the receivers. The opposite is when $r_s \ll \sum_{1 \le i \le n} r_i$ in which case using a single partition offers the highest total rate. In other cases, given the partitioning approach mentioned above, the worst-case scenario happens when there are many slow receivers and only a handful of fast receivers. An example has been shown in Figure 3. In the scenario on the left, our approach groups all the receivers into one partition where they all receive data at the rate of one. That is because by isolating slow receivers we can either get a rate of one or less than one if we isolate more than nine slow receivers, which means using one partition is enough. The optimal case, however, groups all the slow receivers into one partition. In general, scenarios like this rarely happen as the number of slow receivers over inter-datacenter networks is usually small, i.e., most datacenters are connected using high capacity links with large available bandwidth.[2] In general, since we consider all values of $M$ from 1 to $n$ partitions, the solution obtained from our partitioning approach cannot be worse than the two baseline approaches of using a single multicast tree for all receivers and unicasting to all receivers using separate paths.

*4.2. Incorporating Objective Vectors*

We allow users to supply an objective vector along with their multicast transfers to better optimize the network performance, that is, total network capacity consumption and receiver completion times. We incorporate the objective vectors by grouping receivers with consecutive ranks that are indicated with zeros in the objective vector and treating them as one partition in the whole process. That is because the users have indicated no interest in the completion times of those receivers, so we might as well reduce the network capacity usage by grouping them from the beginning. Figure 4 shows an example of building possible solutions by isolating slow receivers and incorporating the user-supplied objective vector, which we refer to as the partitioning hierarchy. Please note that this hierarchy moves in the reverse direction, that is, instead of isolating slow receivers, it merges fast receivers from

---

[2]We have deduced this by looking at many WAN topologies available on the Internet Topology Zoo [58]. We found that in most topologies, a small fraction of nodes are connected using significantly slower links while the variation of downlink/uplink capacity for the rest of the nodes is not significant.
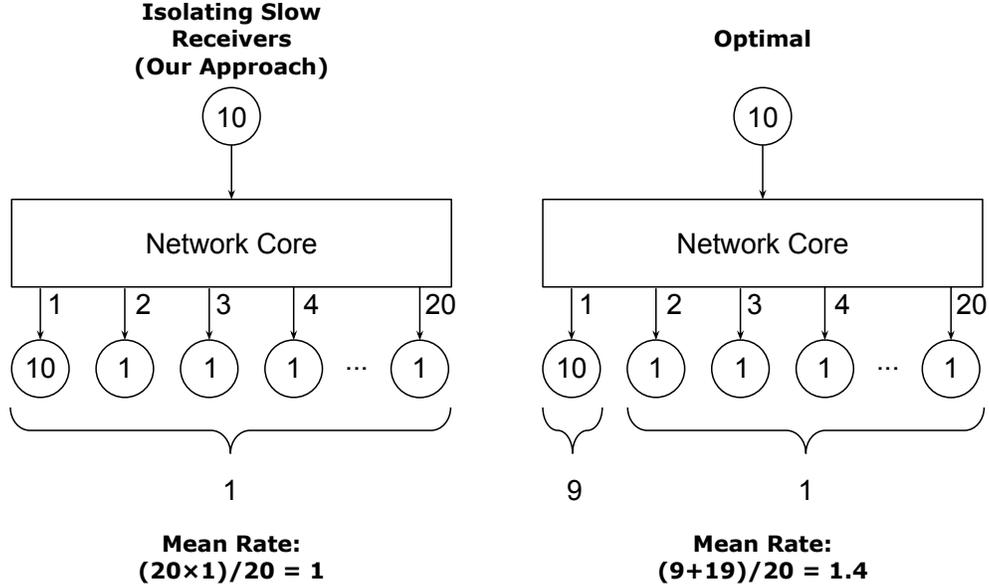
Figure 3: A worst-case scenario for the proposed partitioning scenario. Numbers within the nodes show the downlink and uplink speeds of nodes and curly brackets indicate the partitions where all nodes in a partition receive data at the same rate. The objective is to maximize the average rate of receivers given the max-min fairness policy.

bottom to the top.

Each layer in this hierarchy, labeled as $\mathbf{P}_i, 1 \leq i \leq 5$, represents a valid partitioning solution.[3] We see that receivers indicated with consecutive zeros in $\omega_R$ are merged into one big partition at the base layer or $\mathbf{P}_5$. Also, we see that as we move up, the two fastest partitions at each layer are merged, which reduces total bandwidth consumption. For each layer, we compute the average completion time of receivers and then select the layer that offers the least value, in this case, $\mathbf{P}_3$ was chosen.

## 5. Iris

We apply the partitioning technique discussed in the previous chapter to real-world inter-datacenter networks. We develop a heuristic for partitioning receivers on real-world topologies without relaxations of §4. We will generate multiple valid partitioning solutions in the

---

[3]The associated network topology is not shown.

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $\mathbf{P}_1$ | 8, 5, 3, 9, 7, 10, 1, 4, 2, 6 | | | | | | | | | | |
| $\mathbf{P}_2$ | 8, 5, 3, 9, 7, 10, 1, 4, 2 | | | | | | | | | | 6 |
| $\mathbf{P}_3$ | 8, 5, 3, 9, 7, 10, 1, 4 | | | | | | | | | 2 | 6 |
| $\mathbf{P}_4$ | 8, 5 | | 3, 9, 7, 10, 1, 4 | | | | | | | 2 | 6 |
| $\mathbf{P}_{base}$ ($\mathbf{P}_5$) | 8 | 5 | 3, 9, 7, 10, 1, 4 | | | | | | | 2 | 6 |
| $\mathbf{D}_R$ | 8 | 5 | 3 | 9 | 7 | 10 | 1 | 4 | | 2 | 6 |
| $\Psi$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | | 9 | 10 |
| $\omega_R$ | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | | 1 | 1 |

*Partitioning Hierarchy (right side label). Receiver IDs → $\mathbf{D}_R$. Receiver Ranks → $\Psi$. Objective Vector → $\omega_R$.*
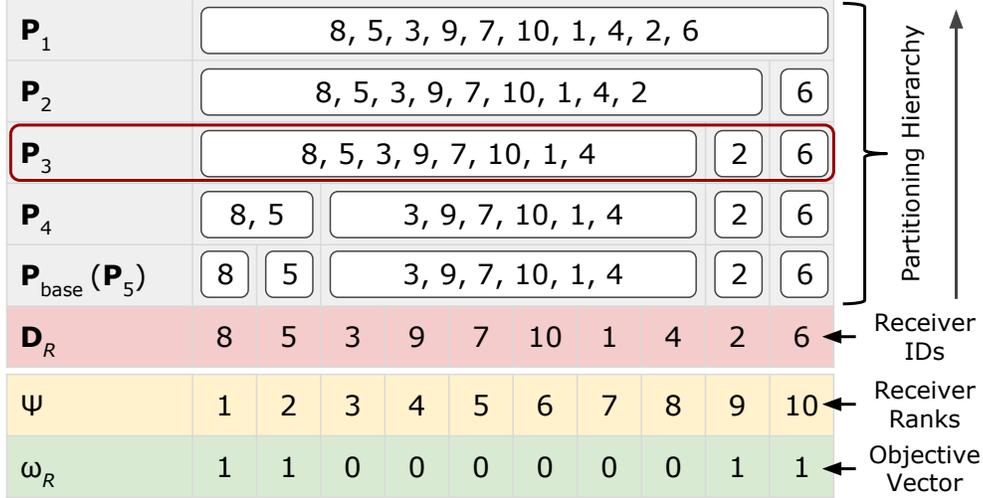
Figure 4: Example of a partitioning hierarchy for a transfer with 10 receivers (the topology not shown).

form of a hierarchy where layers of the hierarchy present feasible partitioning solutions and each layer is formed by merging the two fastest partitions of the layer below.[4]

We present Iris, a heuristic that runs on the traffic engineering server to manage bulk multicast transfers.[5] When a bulk multicast transfer arrives at an end-point, it will communicate the request to the traffic engineering server which will then invoke Iris. It uses the knowledge of physical layer topology, available bandwidth on edges after deducting the share of high priority user traffic and other running transfers to compute partitions and forwarding trees. The traffic engineering server pulls end-points' actual progress periodically to determine their exact remaining volume across transfers to compute the total outstanding load per edge for all edges. Iris consists of four modules as shown in Figure 5 which we discuss in the following subsections. Iris aims to find an approximate solution to the optimization problem of Eq. 1 assuming $\epsilon \ll 1$ to prioritize minimizing completion times over minimizing

---

[4]In general, it is not possible to offer optimality guarantees due to the highly varying factors of network topology, transfer arrivals, and the distribution of transfer volumes. However, our extensive simulations in §6 show that our approach can offer significant improvement on other approaches over various topologies and traffic patterns. Also, as a result of building a hierarchy of partitioning options and selecting the best one, our solution will be at least as good as either using a single multicast tree or using unicasting to all receivers.

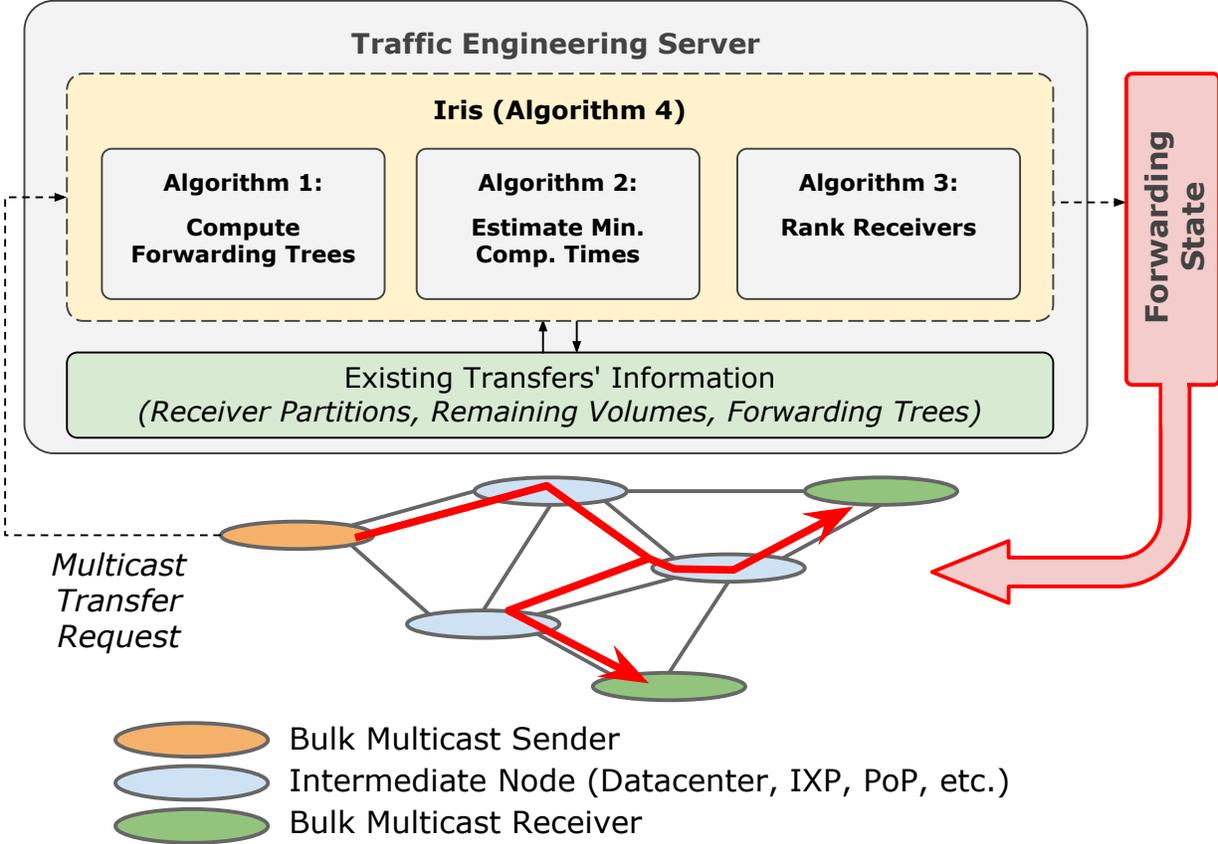[5]Unicast transfers are a special case with a single receiver.

Figure 5: Pipeline of Iris.

bandwidth consumption. We will empirically evaluate Iris by comparing it to recent work and a lower bound in §6.

## 5.1. Choosing Forwarding Trees

Load aware forwarding trees are selected given the link capacity information on the topology and according to other ongoing bulk multicast transfers across the network to reduce the completion times by mitigating the effect of bottlenecks. Tree selection should also aim to keep bandwidth consumption low by minimizing the number of edges per tree where an edge could refer to any of the links on the physical topology. To select a forwarding tree, a general approach that can capture a wide range of selection policies is to assign weights to edges of the inter-datacenter graph $G$ and select a minimum weight Steiner tree [59]. Per edge $e \in \mathbf{E}_G$, we assume a virtual queue that increases by volume of every transfer scheduled on that edge and decreases as traffic flows through it. Since edges differ in capacity,

---
**Algorithm 1:** Compute A Forwarding Tree
---
  **Input:** Steiner tree terminal nodes $\boldsymbol{\Gamma}$, request $R$

  **Output:** Edges of a tree

  **CompForwardingTree** $(\boldsymbol{\Gamma}, R)$

    | To every edge $e \in \mathbf{E}_G$, assign weight $W_e = (L_e + \frac{\mathcal{V}_R}{B_e})$;

    | **return** *A minimum weight Steiner tree that connects the nodes in set $\boldsymbol{\Gamma}$ (we used a*

    | *hueristic [61, 60]);*
---

completing the same virtual queue size may need significantly different times for different links. We extend the metric used in a recent work [28] that is called load $L_e$ as follows.

$$L_e = \frac{1}{B_e} \sum_{\{P \in \mathbf{P_R} | e \in T_P\}} \mathcal{V}_P^{[res]} \tag{4}$$

In Eq. 4, $\mathbf{P_R}$ is the set of partitions of receivers of all ongoing transfers. This equation sums up the remaining volumes of all trees that use a specific edge (total virtual queue size) and divides that by the average available bandwidth on that edge to compute the minimum possible time it takes for all ongoing transfers on that edge to complete. In the tree selection process, to keep completion times low, we need to avoid edges for which this value is large.

With this metric available, to select a forwarding tree given a sender and several receivers, we will first assign an edge weight of $W_e = L_e + \frac{\mathcal{V}_R}{B_e}$ to all edges and then select a minimum weight Steiner tree as shown in Algorithm 1. With this edge weight, compared to edge utilization which has been extensively used in literature for traffic engineering, we achieve a more stable measure of how busy a link is expected to be in the near future on average. We considered the second term in edge weight to reduce total bandwidth use when there are multiple trees with the same weight. It also leads to the selection of smaller trees for larger transfers which decreases the total bandwidth consumption of Iris further in the long run.

**Complexity:** To compute a minimum weight Steiner tree we use a heuristic that is called GreedyFLAC [60] which given the set of terminal nodes $\Gamma$, has a guaranteed polynomial running time of $\mathcal{O}(|\mathbf{V}_G||\mathbf{E}_G| + |\mathbf{V}_G|^2 log(|\mathbf{V}_G|)|\Gamma| + |\mathbf{V}_G|^2|\Gamma|^3)$.

*5.2. Estimating Minimum Completion Times*

The purpose of this procedure is to estimate the minimum completion time of different partitions of a given transfer considering available bandwidth over the edges and applying max-min fair rate allocation when there are shared links across forwarding trees. Algorithms 3 and 4 then use the minimum completion time per partition to rank the receivers (i.e., faster receivers have an earlier completion time) and then decide which partitions to merge. Computing the minimum completion times is done by assuming that the new transfer request has access to all the available bandwidth and compared to computing the exact completion times is much faster. Besides, calculating the exact completion times is not particularly more effective due to the continuously changing state of the system as new transfer requests arrive. Since available bandwidth over future timeslots is not precisely known, we can use estimate values similar to other work [33, 12, 62]. Algorithm 2 shows how the minimum completion times are computed.

**Complexity:** For a new request $R$ with $|\mathbf{P}|$ partitions, this algorithm calls Algorithm 1, $|\mathbf{P}|$ times. It then computes max-min fair rates per partition for timeslots until all partitions finish. Computing max-min fair rates per timeslot has a complexity of $\mathcal{O}(|\mathbf{P}||\mathbf{E}_G|)$. This process continues for $\mathcal{O}(\frac{|\mathbf{P}| \, \mathcal{V}_R}{\min_{e,t} B_e(t)})$ iterations. Therefore, the complexity of this algorithm is $\mathcal{O}(|\mathbf{P}| \, (\mathcal{C}_{Algorithm \ 1} + |\mathbf{E}_G| \, \frac{|\mathbf{P}| \, \mathcal{V}_R}{\min_{e,t} B_e(t)}))$.

*5.3. Assigning Ranks to Receivers*

Algorithm 3 assigns ranks to individual receivers according to their minimum completion times taking into account available bandwidth over edges as well as edges' load in the path selection process. This ranking is used along with the provided objective vector later to partition receivers.

**Complexity:** This algorithm calls Algorithm 2 over all receivers as separate 1-node partitions, then sorts the nodes which gives a complexity of $\mathcal{O}(\mathcal{C}_{Algorithm \ 2} + |\mathbf{D}_R| \, log(|\mathbf{D}_R|))$.

*5.4. The* Iris *Algorithm*

The Iris algorithm computes receiver partitions using hierarchical partitioning and assigns each partition a multicast forwarding tree. The partitioning problem is solved per transfer

**Algorithm 2:** Minimum Completion Times

**Input:** A set of partitions $\mathbf{P}$, request $R$

**Output:** Completion time of every partition in $\mathbf{P}$

**MinimumCompletionTimes** $(\mathbf{P}, R)$

> $\mathbf{f} \leftarrow \emptyset$, $t \leftarrow t_{now} + 1$;
>
> $\gamma_P \leftarrow \mathcal{V}_R$, $\forall P \in \mathbf{P}$;
>
> $T_P \leftarrow$ `CompForwardingTree`$(P, R)$, $\forall P \in \mathbf{P}$;
>
> **while** $|\mathbf{f}| < |\mathbf{P}|$ **do**
>
> > Compute $r_P(t), \forall P \in \{\mathbf{P} - \mathbf{f}\}$, max-min fair rate [57] allocated to tree $T_P$ at
> > timeslot $t$ given available bandwidth of $B_e(t)$ on every edge $e \in \mathbf{E}_G$;
> >
> > $\gamma_P \leftarrow \gamma_P - \delta \times r_P(t), \ \forall P \in \mathbf{P}$;
> >
> > **foreach** $P \in \{\mathbf{P} - \mathbf{f}\}$ **do**
> >
> > > **if** $\gamma_P = 0$ **then**
> > >
> > > > $\kappa_P \leftarrow t, \mathbf{f} \leftarrow \mathbf{f} \cup P$;
> >
> > $t \leftarrow t + 1$;
>
> **return** $\kappa_P, \forall P \in \mathbf{P}$

and determines the number of partitions and the receivers that are grouped per partition. Iris uses a partitioning technique inspired by the findings of §4 that is computationally fast, significantly improves receiver completion times, and operates only relying on network topology and available bandwidth per edge.[6] We refer to our approach as hierarchical partitioning as it builds a hierarchy of different partitioning solutions.

We build a partitioning hierarchy with numerous layers and examine the various number of partitions from bottom to the top of the hierarchy while looking at the average of minimum completion times. Given that each node in a real-world topology may have multiple interfaces, by building a hierarchy, we consider the discrete nature of forwarding tree

---

[6]The available bandwidth per edge is computed by deducting the quota for higher priority user traffic from the link capacity.

---
**Algorithm 3:** Assign Receiver Ranks
---
**Input:** Request $R$

**Output:** $\psi_r$, i.e., rank of receiver $r \in \mathbf{D}_R$

**AssignReceiverRanks** $(R)$

   /* Every receiver is treated as a separate partition */

   $\{\kappa_r, \ \forall r \in \mathbf{D}_R\} \leftarrow$ `MinimumCompletionTimes`$(\mathbf{D}_R, R)$;

   $\psi_r \leftarrow$ Position of receiver $r$ in the list of all receivers sorted by their estimated minimum completion times (fastest receiver is assigned a rank of 1), $\forall r \in \mathbf{D}_R$;

   **return** $\psi_r, \forall r \in \mathbf{D}_R$;

---

selection on the physical network topology. The process consists of two steps as follows.

Algorithm 4 illustrates how Iris partitions receivers with an objective vector. We first use the receiver ranks from Algorithm 3 and the objective vector to create the base of partitioning hierarchy, $\mathbf{P}_{base}$. We first sort the receivers by their ranks from fastest to slowest and then group them according to the weights in the objective vector. For any receiver whose rank in the objective vector has a value of 1, we consider a separate partition (single node partition) which allows the receiver to complete as fast as possible by not attaching it to any other receiver. Next, we group receivers with consecutive ranks that are assigned a value of 0 in the objective vector into partitions with potentially more than one receiver, which allows us to save as much bandwidth as possible since the user has not indicated interest in their completion times.

Now that we have a set of base partitions $\mathbf{P}_{base}$, a heuristic creates a hierarchy of partitioning solutions with $|\mathbf{P}_{base}|$ layers where every layer $1 \leq l \leq |\mathbf{P}_{base}|$ is made up of a set of partitions $\mathbf{P}_l$. Each layer is created by merging two partitions from the layer below going from the bottom to the top of hierarchy. At the bottom of the hierarchy, we have the base partitions. Also, at any layer, any partition $P$ is attached to the sender using a separate forwarding tree $T_P$. We first compute the average of minimum completion times of all receivers at the bottom of the hierarchy. We continue by merging the two partitions that hold receivers with highest ranks. When merging two partitions, the faster partition is slowed down to the speed of slower partition. A new forwarding tree is computed for the resulting

---
**Algorithm 4:** Compute Receiver Partitions and Trees (Iris)

---

**Input:** Request $R$, binary objective vector $\omega_R$

**Output:** Partitions of request $R$ and their forwarding trees

**CompPartitionsAndTrees** $(R, \omega_R)$

/* Initial partitioning using the objective vector $\omega_R$ */

$\{\psi_r, \forall r \in \mathbf{D}_R\} \leftarrow$ `AssignReceiverRanks`$(R)$;

$\mathbf{D}_R^s \leftarrow$ Receivers $r$ sorted by $\psi_r, \forall r \in \mathbf{D}_R$ ascending;

$\mathbf{P}_{base} \leftarrow \{$Any receiver $r \in \mathbf{D}_R$ for which $\omega_R < \psi_r >$ is 1 as a separate partition$\} \cup \{$Group receivers that appear consecutively on $\mathbf{D}_R^s$ for which $\omega_R < \psi_r >$ is 0, each group forms a separate partition$\}$;

/* Building the partitioning hierarchy for $\mathbf{P}_{base}$ */

$\mathbf{P}_{|\mathbf{P}_{base}|} \leftarrow \mathbf{P}_{base}$;

**for** $l = |\mathbf{P}_{base}|$ **to** $l = 1$ **by** $-1$ **do**

    $\{\kappa_P, \ \forall P \in \mathbf{P}_l\} \leftarrow$ `MinimumCompletionTimes`$(\mathbf{P}_1, R)$;

    $\kappa_l \leftarrow \frac{\sum_{P \in \mathbf{P}_l}(|P| \ \kappa_P)}{|\mathbf{D}_R|}$;  /* Compute the best average completion times */

    Assuming receivers are sorted from left to right by increasing order of rank, merge the two partitions on the left, $P$ and $Q$, to form $PQ$;

    $\mathbf{P}_{l-1} \leftarrow \{PQ\} \cup \{\mathbf{P}_l - \{P, Q\}\}$;

Find $l_{min}$ for which $\kappa_{l_{min}} \leq \min_{1 \leq l \leq |\mathbf{P}_{base}|} \kappa_l$, if multiple layers have the same $\kappa_l$, choose the layer with minimum total weight over all of its forwarding trees, i.e., select $l_{min}$ to $\min \sum_{P \in \mathbf{P}_{l_{min}}} (\sum_{e \in T_P} W_e)$;

/* Compute forwarding trees for the partitions */

**foreach** $P \in \mathbf{P}_{l_{min}}$ **do**

    $T_P \leftarrow$ `CompForwardingTree`$(P, R)$;

    **foreach** $e \in T_P$ **do**

        $L_e \leftarrow L_e + \frac{\mathcal{V}_R}{B_e}, \ W_e \leftarrow W_e + \frac{\mathcal{V}_R}{B_e}$;

**return** $(P, \ T_P), \ \forall P \in \mathbf{P}_{l_{min}}$;

---

Table 3: Various topologies and traffic patterns used in evaluation. One unit of traffic is equal to what can be transmitted at the rate of the fastest link over a given topology per timeslot.

|  | Name | Description |
|---|---|---|
| Topology | GEANT | Backbone and transit network across Europe with 34 nodes and 52 links. Link capacity from 45 Mbps to 10 Gbps. |
|  | UNINETT | Backbone network across Norway with 69 nodes and 98 links. Most links have a capacity of 1, 2.5 or 10 Gbps. |
| Traffic Pattern | Light-tailed | Based on Exponential distribution with a mean of 20 units per transfer. |
|  | Heavy-tailed | Based on Pareto distribution with the minimum of 2 units, the mean of 20 units, and the maximum capped at 2000 units per transfer. |
|  | Hadoop | Generated by geo-distributed data analytics over Facebook's inter-datacenter WAN (distribution mean of 20 units per transfer). |
|  | Cache-follower | Generated by geo-distributed cache applications over Facebook's inter-datacenter WAN (distribution mean of 20 units per transfer). |

partition using the forwarding tree selection heuristic of Algorithm 1 to all receivers in that partition, and the average of minimum completion times for all receivers are recomputed. This process continues until we reach a single partition that holds all receivers. In the end, we select the layer at which the average of minimum completion times across all receivers is minimum, which gives us the number of partitions, the receivers that are grouped per partition, and their associated forwarding trees. If there are multiple layers with the minimum average completion times, the one with minimum total forwarding tree weight across its forwarding trees is chosen which on average leads to better load distribution.

**Complexity:** This algorithm first calls Algorithm 3, then it calls Algorithm 2 up to $|\mathbf{D}_R|$ times. At the end, it also runs Algorithm 1 up to $|\mathbf{D}_R|$ times. Therefore, this algorithm has a complexity of $\mathcal{O}(\mathcal{C}_{Algorithm\ 3} + |\mathbf{D}_R|\ (\mathcal{C}_{Algorithm\ 2} + \mathcal{C}_{Algorithm\ 1}))$.

## 6. Evaluation

We considered various topologies and transfer size distributions as shown in Table 3. We selected two research topologies with given capacity information on edges from the Internet Topology Zoo [63]. We could not use other commercial topologies as the exact connectivity and link capacity information were not publicly disclosed. We also considered multiple transfer volume distributions including synthetic (light-tailed and heavy-tailed) and real-world Facebook inter-datacenter traffic patterns (Hadoop and Cache-follower) [64]. Transfer arrival pattern was according to Poisson distribution with a rate of $\lambda$ per timeslot. For simplicity, we assumed an equal number of receivers for all bulk multicast transfers per experiment. We performed simulations and Mininet emulations to evaluate Iris.

We compare Iris with multiple baseline techniques and QuickCast [28] which also focuses on partitioning receivers into groups for improved completion times. We were unable to evaluate our work against another recent work called BDS [11], which has been developed by Baidu, due to source code unavailability. BDS takes advantage of store-and-forward and operates at the application layer. We qualitatively compare Iris with BDS. Iris can offer lower bandwidth consumption since BDS uses paths instead of trees, and can more effectively exercise physical links as it manages traffic at the network layer. On the other hand, since BDS uses all available overlay paths (possibly many routes to any specific receiver), under the lightly loaded regime, BDS may offer higher throughput (at the cost of considerably higher bandwidth consumption). Extension of Iris to use parallel trees is considered part of future work.

### 6.1. Computing a Lower Bound

We develop a technique to compute a lower bound on receiver completion times by creating an aggregate topology from the actual topology. As shown in Figure 6, to create the aggregate topology, we combine all downlinks and uplinks with rates $r_i^{[node]}$ for all interfaces $i$ per node to a single uplink and downlink with their rates set to the sum of rates of physical links. Also, the aggregate topology connects all nodes in a star topology using their uplinks and downlinks and so assumes no bottlenecks within the network. Since this topology is a relaxed version of the physical topology, any solution that is valid for the physical topology
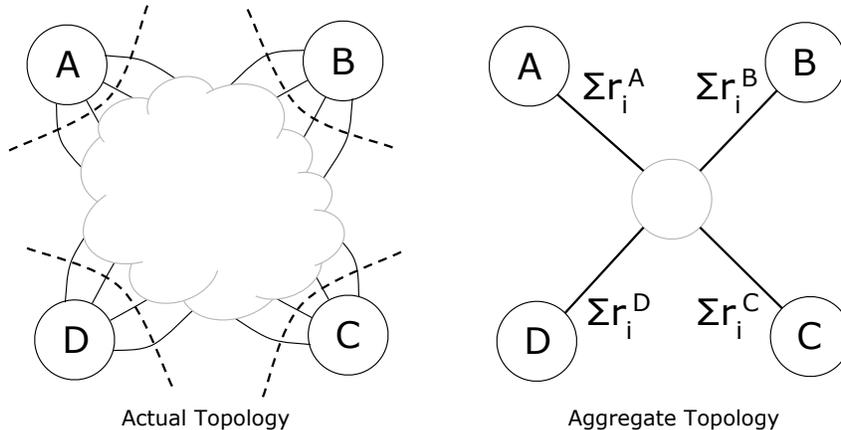
Figure 6: The physical topology, and the aggregate topology to compute a lower bound on receiver completion times. The aggregate topology is used only for evaluation purposes and it does not play any part in the design of Iris.

is valid on this topology as well. Therefore, the solution to the aggregate topology is a lower bound that can be computed efficiently but may be inapplicable to the actual physical topology. We will use this approach in §6.2.1 for evaluation of Iris.

### 6.2. Simulations

In simulations, we focus on computing gains and therefore assume no dropped packets and accurate max-min fair rates. We normalized link capacities by maximum link rate per topology and fixed the timeslot length to $\delta = 1.0$.

**Effect of User Traffic:** We account for the effect of higher priority user traffic in the simulations. The amount of available bandwidth per edge per timeslot, i.e., $B_e(t)$, is computed by deducting the rate of user traffic from the link capacity $C_e$. Recent work has shown that this rate can be safely estimated [12, 33]. For evaluations, we assume that user traffic can take up to 30% of a link's capacity with a minimum of 5% and that its rate follows a periodic pattern going from low to high and to low again. Per link, we consider a random period in the range of 10 to 100 timeslots that is generated and assigned per experiment instance.

### 6.2.1. Minimizing Average Completion Times

This is when the objective vector is made of all ones. The partitioning hierarchy then begins with all receivers forming their 1-receiver partitions. This is a highly general objective

and can be considered as the default approach when the application/user does not specify an objective vector. We discuss multiple simulation experiments.

In Figure 7, we measure the completion times (mean and tail) as well as bandwidth consumption by the number of receivers (tail is $99.9^{th}$ percentile). We consider two baseline cases: unicast shortest path and static single tree (i.e., minimum edge Steiner tree) routing. The shortest path routing is the unicast scenario that uses minimum bandwidth possible. The minimum edge Steiner tree routing uses minimum bandwidth possible while connecting all receivers with a single tree. The first observation is that using unicast, although leads to highest separation of fast and slow receivers, does not lead to the fastest completion as it can lead to many shared bottlenecks and that is why we see long tail times. Iris offers the minimum completion times (mean and tail) across all scenarios. Also, its completion times grow much slower compared to others as the number of receivers (and so overall network load) increases. This is while Iris uses only up to 35% additional bandwidth compared to the static single tree (unicast shortest path routing uses up to $2.25\times$). Compared to QuickCast, Iris offers up to 26% lower tail times and up to $2.72\times$ better mean times while using up to 13% extra bandwidth.

In Figure 8, we show the completion times speedup of receivers by their rank. As seen, gains depend on the topology, traffic pattern, and receiver's rank. The dashed line is the baseline, i.e., no-partitioning case. Compared to QuickCast [28], the fastest node always completes faster and up to $2.25\times$ faster with Iris. Also, the majority of receivers complete significantly faster. In case of four receivers, the top 75% receivers complete between $2\times$ to $4\times$ faster than baseline and with sixteen receivers, the top 75% receivers complete at least $8\times$ faster than baseline. This is when QuickCast's gain drops quickly to one after the top 25% of receivers.

In Figure 9, we measure the CDF of completion times for all receivers. As seen, tail completion times are two to three orders of magnitude longer than median completion times which is due to variable link capacity and transfer volumes. We evaluate the completion times of QuickCast and Iris and compare them with a lower bound which considers the aggregate topology (see §6.1) and applies Theorem 2 directly. It is likely that no feasible solution exists that achieves this lower bound. Under low arrival rate (light load), we see
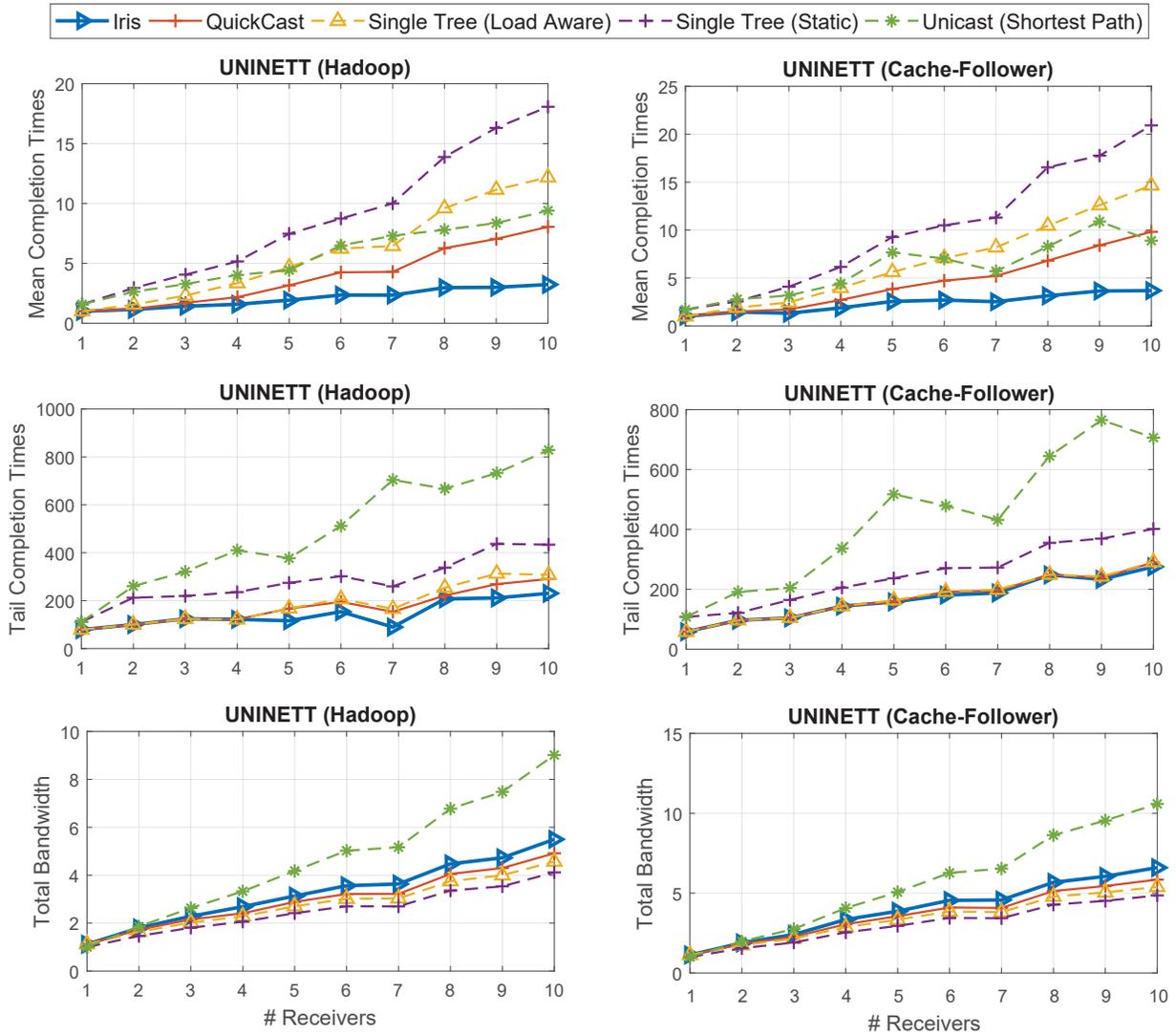
Figure 7: Comparison of various techniques by number of multicast receivers. Plots are normalized by the minimum data point (mean and tail charts are normalized by the same minimum), $\lambda = 1$, and lower values are better.
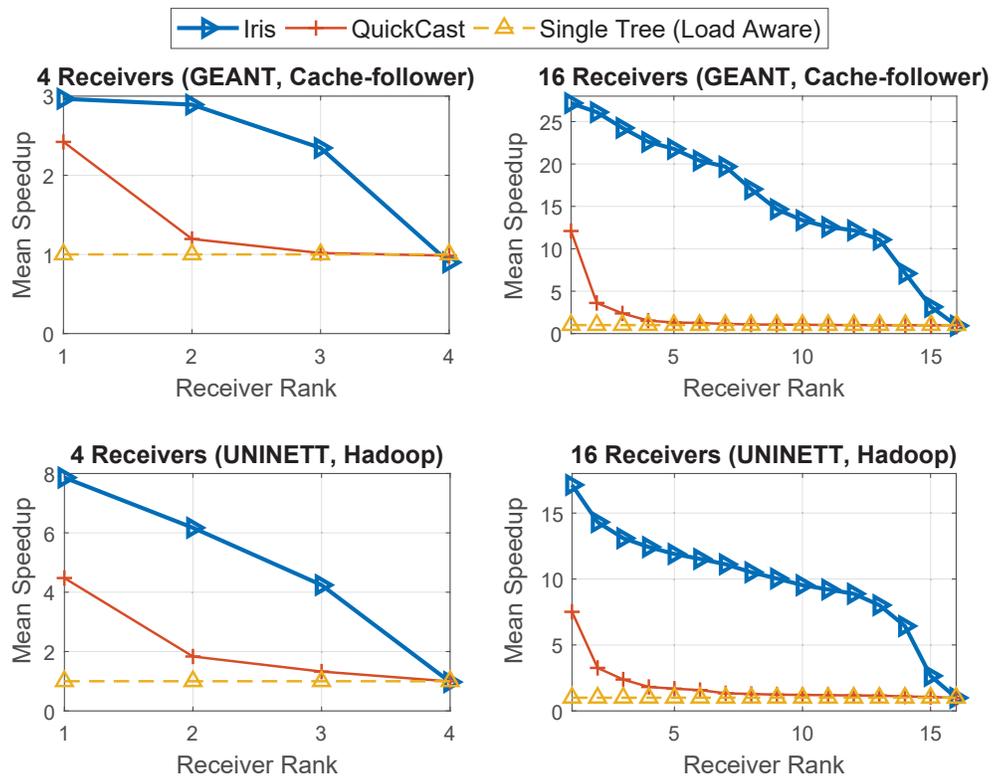
Figure 8: Mean completion time speedup (larger is better) of receivers normalized by no partitioning (load aware single tree) case given their rank from fastest to slowest, every node initiates equal number of transfers, receivers were selected according to uniform distribution from all nodes, and we considered $\lambda$ of 1.

that Iris tracks the lower bound nicely with a marginal difference. Under high arrival rate (heavy load), Iris stays close to the lower bound for lower and higher percentiles while not far from it for others.

### 6.2.2. Other Objective Vectors

We discuss four different objective vectors of $A$, $B$, $C$ and $D$ as shown in Figure 10. This figure shows the mean speedup of receivers given their ranks, and the bandwidth consumption associated with each vector. In $A$, we aim to finish one copy quickly while not being concerned with completion times of other receivers. We see a gain of between $9\times$ to $18\times$ across the two topologies considered for the first receiver. We also see that this approach uses much less extra bandwidth compared to when we have a vector with more ones (e.g., case $B$). In $B$, we aim to speed up the first four receivers (we care about each one) while in $C$, we want to speed up the fourth receiver not directly concerning ourselves with the top three receivers. As can be seen, $B$ offers increasing speedups for the top three receivers while $C$'s speedup is flatter. Also, $C$ uses less bandwidth compared to $B$ by grouping the top three receivers into one partition at the base of the hierarchy. Finally, $D$'s vector specifies that the application/user only cares about the completion time of the last receiver which means that receiver will be put in a separate partition at the base of the hierarchy while other receivers will be grouped into one partition. Since the slowest receiver is usually limited by its downlink speed, this cannot improve its completion time. However, with minimum extra bandwidth, this speeds up all receivers except the slowest by as much as possible. Except for the slowest, all receivers observe a speedup of between $3\times$ to $6\times$ while using $8\%$ to $16\%$ less bandwidth compared to $B$. A tradeoff is observed, that is, $D$ offers lower speedup but consistent gain for more receivers with less bandwidth use compared to $B$.

### 6.3. Mininet Emulations

We used Mininet to build and test a prototype of Iris and compare it with QuickCast and set up the testbed on CloudLab [65]. We used OpenvSwitch (OVS) 2.9 in the OpenFlow 1.3 compatibility mode along with the Floodlight controller 1.2 connecting them to a control network. We assumed fixed available bandwidth over edges according to GEANT topology
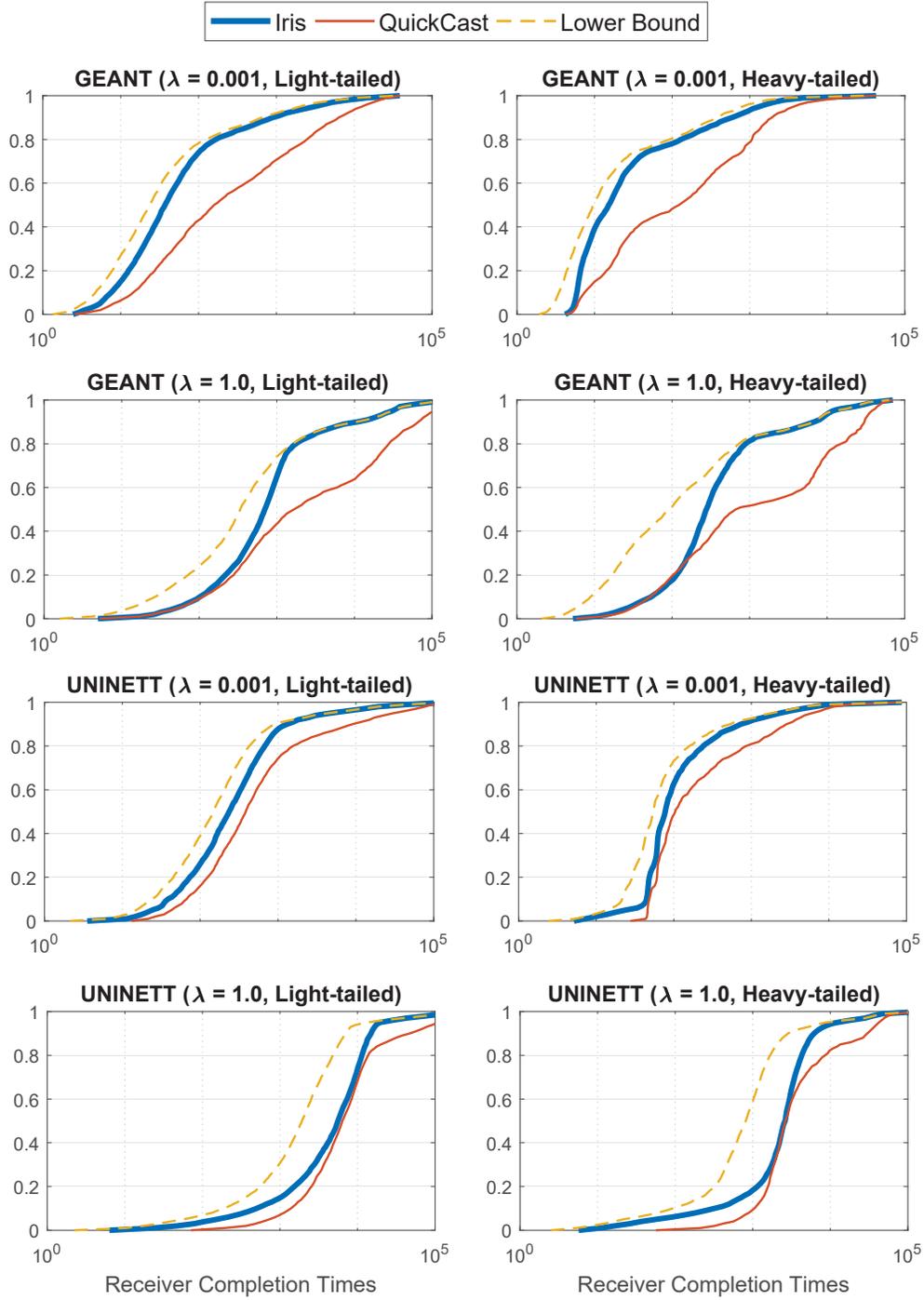
Figure 9: CDF of receiver completion times. Every transfer has 8 receivers selected uniformly across all nodes. "Lower Bound" is computed by finding the aggregate topology and applying Theorem 2.
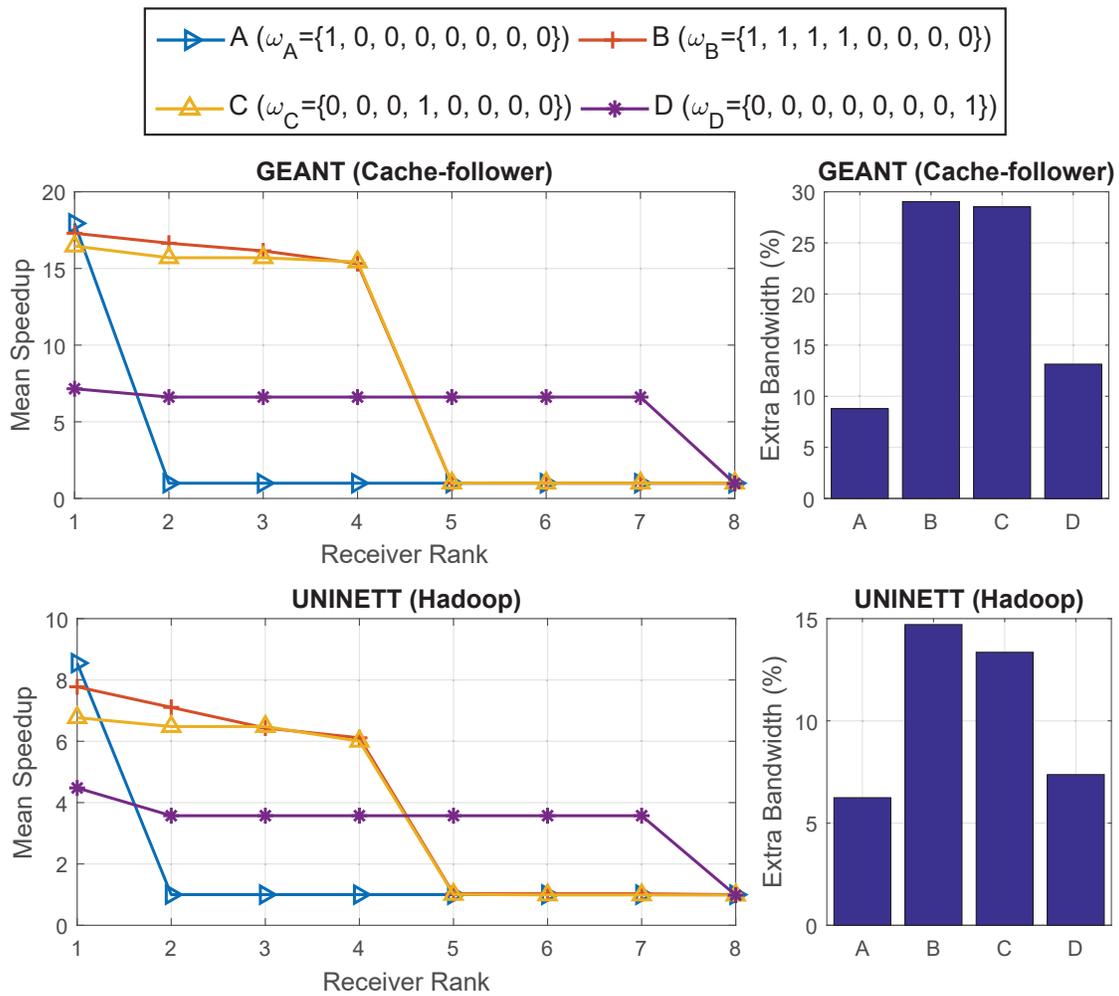
Figure 10: Gain by rank for different receivers per transfer averaged over all transfers for four different objective vectors. We set $\lambda = 0.1$ and there are 8 receivers.

[66] while scaling downlinks' capacity so that the maximum is 500 Mbps.[7] We did this to reduce the CPU overhead of traffic shaping over TCLink Mininet modules. Our traffic engineering program communicated with end-points through a RESTful API. We used NORM [67] for multicast session management along with its rate-control module. To increase efficiency, we computed max-min fair rates centrally at the traffic engineering program and let the end-points shape their traffic using NORM's rate control module. The experiment was performed using twelve trace files generated according to Facebook traffic patterns concerning transfers' volumes [64], and each trace file had 200 requests in total with an arrival rate of one request per timeslot based on Poisson distribution. We also considered timeslots of one second, a minimum transfer volume of 5 MB and limited the maximum transfer volume to 500 MB.[8] We considered three schemes of Iris, QuickCast and a single tree approach (no partitioning). The total emulation time was about 24 hours. Figure 11 shows our emulation results. To allow comparison between the tail, i.e., 95[th] percentile, and mean values, we have normalized both plots by the same minimum in each row. Also, the group table usage plots are not normalized and show the actual average and actual maximum across all switches. The reason why data points jump up and down is the randomness of generated traces that comes from transfers (volume, source, receivers, arrival pattern, etc).

**Completion Times and Bandwidth:** Iris offers up to $2.5\times$ speed up in mean completion times compared to QuickCast and $4\times$ compared to using a single multicast tree per transfer. We also see that compared to using a single multicast tree, Iris consumes at most $25\%$ extra bandwidth.

**Forwarding Plane:** We see that Iris uses up to about $4\times$ less group table entries at the

---

[7]In general, inter-datacenter link capacities may go beyond tens of Gbps. Due to the limitations of our emulation server, we had to use 500 Mbps as the maximum link capacity. That is, although we used a server with 56 logical CPU cores, even with a maximum link rate of 500 Mbps, the machine ran at close to full CPU utilization across most cores. Using a higher rate would have led to inaccurate emulation results due to the timing inaccuracies caused at high CPU utilization. The high CPU utilization in Mininet is caused mostly by the traffic shapers that Mininet uses to model links' capacities. Using a lower rate, however, should have a negligible effect on the validity of results as a proof of concept.

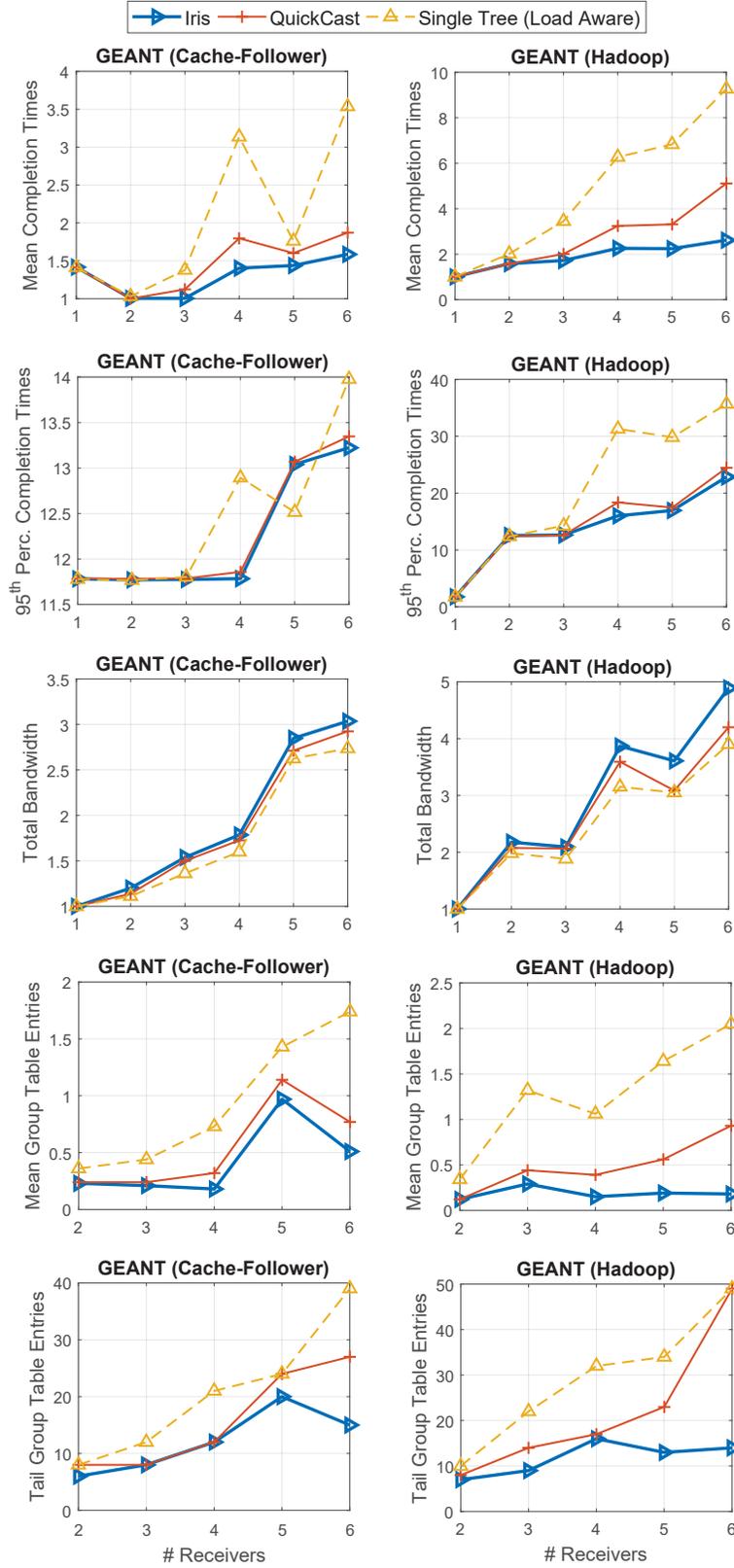[8]These parameters also match the distribution of YouTube video sizes [68].

Figure 11: Mininet Emulation Results

switches where the maximum number of entries were exhausted which allows more parallel transfers across the same network. Iris achieves this by allowing a larger number of partitions per transfer whenever it does not hurt the completion times. By allowing more partitions, each tree will branch less times on average reducing the number of group table entries.

**Running Time:** Across all experiments, the computation time needed to run Iris to calculate partitions and forwarding trees, i.e., running Algorithm 4, stayed below 5 ms per request.

## 6.4. Practical Concerns

New challenges, such as increased communication latency across network elements and failures, may arise while deploying Iris on a real-world geographically distributed network. Communication latency may not affect the performance considerably as we focus on long-running internal transfers that are notably more resilient to latency overhead of scheduling and routing compared to interactive user traffic. Failures may affect physical links or the traffic engineering server. Loss of a physical link can be addressed by rerouting the affected transfers reactively either by the controller or by use of SDN fast failover mechanisms. Endpoints may be equipped with distributed congestion control (similar to [44]) which they can fall back to in case the centralized traffic engineering fails.

## 7. Conclusions and Future Work

Replication of content and data across geographically dispersed datacenters creates a large volume of multicast traffic that needs to be managed for increased performance. A bulk multicast transfer can be indicated with a source, set of receivers and total transfer volume. In this paper, we focused on the problem of grouping receivers into multiple partitions to minimize the effect of receiver downlink speed discrepancy on completion times of receivers. We analyzed a relaxed version of this problem and proposed a partitioning that reduces mean completion times of multicast receivers given max-min fair rates. We also set forth the idea of applications/users expressing their requirements in the form of binary objective vectors which allows us to optimize resource consumption and performance further. We then proposed Iris, a system that computes partitions and forwarding trees for incoming bulk multicast transfers as they arrive given objective vectors. We showed that Iris could

significantly reduce mean completion times with a small increase in bandwidth consumption and can fulfill the requirements expressed using objective vectors while saving bandwidth whenever possible. It is worth noting that performance of any partitioning and forwarding tree selection algorithm rests profoundly on the network topology and transfer properties. Study of rate allocation policies besides max-min fairness and handling failures are among future directions.

## References

[1] Chunqiang Tang, Thawan Kooburat, et al. Holistic configuration management at facebook. In *Proceedings of the 25th Symposium on Operating Systems Principles*, SOSP '15, pages 328–343. ACM, 2015.

[2] Building express backbone: Facebook's new long-haul network. `https://code.facebook.com/posts/1782709872057497/building-express-backbone-facebook-s-new-long-haul-network/`. visited on September 30, 2017.

[3] Ken Florance. How netflix works with isps around the globe to deliver a great viewing experience. `https://goo.gl/72CbM7`, 2016.

[4] Sushant Jain, Alok Kumar, et al. B4: Experience with a globally-deployed software defined wan. *SIGCOMM*, 43(4):3–14, 2013.

[5] Y. Xia, T. S. E. Ng, and X. S. Sun. Blast: Accelerating high-performance data analytics applications by optical multicast. In *2015 IEEE Conference on Computer Communications (INFOCOM)*, pages 1930–1938, April 2015.

[6] Kirill Bogdanov, Miguel Peón-Quirós, Gerald Q. Maguire, Jr., and Dejan Kostić. The Nearest Replica Can Be Farther Than You Think. In *Proceedings of the Sixth ACM Symposium on Cloud Computing*, SoCC '15, pages 16–29, New York, NY, USA, 2015. ACM.

[7] Sarah Wassermann, John P Rula, et al. Anycast on the Move: A Look at Mobile Anycast Performance. *Network Traffic Measurement and Analysis Conference*, 2018.

[8] Felipe Huici, Mohamed Ahmed, Sofia Nikitaki, and Saverio Niccolini. Efficient caching in content delivery networks based on popularity predictions, May 9 2017. US Patent 9,648,126.

[9] M. Noormohammadpour, C. S. Raghavendra, S. Rao, and S. Kandula. Dccast: Efficient point to multipoint transfers across datacenters. In *HotCloud*. USENIX Association, 2017.

[10] Volker Stocker, Georgios Smaragdakis, William Lehr, and Steven Bauer. The growing complexity of content delivery networks: Challenges and implications for the Internet ecosystem. *Telecommunications Policy*, 41(10):1003 – 1016, 2017. Celebrating 40 Years of Telecommunications Policy A Retrospective and Prospective View.

[11] Yuchao Zhang, Junchen Jiang, Ke Xu, Xiaohui Nie, Martin J. Reed, Haiyang Wang, Guang Yao, Miao Zhang, and Kai Chen. BDS: A Centralized Near-optimal Overlay Network for Inter-datacenter Data Replication. In *Proceedings of the Thirteenth EuroSys Conference*, EuroSys '18, pages 10:1–10:14. ACM, 2018.

[12] Srikanth Kandula, Ishai Menache, Roy Schwartz, and Spandana Raj Babbula. Calendaring for wide area networks. *SIGCOMM*, 44(4):515–526, 2015.

[13] Xin Jin, Yiran Li, Da Wei, Siming Li, Jie Gao, Lei Xu, Guangzhi Li, Wei Xu, and Jennifer Rexford. Optimizing bulk transfers with software-defined optical wan. In *SIGCOMM*, pages 87–100. ACM, 2016.

[14] Chi-Yao Hong, Srikanth Kandula, Ratul Mahajan, et al. Achieving high utilization with software-driven wan. In *SIGCOMM*, pages 15–26. ACM, 2013.

[15] M. Cotton, L. Vegoda, and D. Meyer. IANA guidelines for IPv4 multicast address assignments. Internet Requests for Comments, 2010.

[16] Suman Banerjee, Bobby Bhattacharjee, and Christopher Kommareddy. Scalable application layer multicast. In *SIGCOMM*, pages 205–217. ACM, 2002.

[17] R. Sherwood, R. Braud, and B. Bhattacharjee. Slurpie: a cooperative bulk data transfer protocol. In *INFOCOM*, volume 2, pages 941–951, 2004.

[18] Johan Pouwelse, PawełGarbacki, Dick Epema, and Henk Sips. The bittorrent p2p file-sharing system: Measurements and analysis. In *Proceedings of the 4th International Conference on Peer-to-Peer Systems*, IPTPS'05, pages 205–216, Berlin, Heidelberg, 2005. Springer-Verlag.

[19] Aakash Iyer, Praveen Kumar, and Vijay Mann. Avalanche: Data center multicast using software defined networking. In *COMSNETS*, pages 1–8. IEEE, 2014.

[20] J. Cao, C. Guo, G. Lu, Y. Xiong, Y. Zheng, Y. Zhang, Y. Zhu, C. Chen, and Y. Tian. Datacast: A scalable and efficient reliable group data delivery service for data centers. *IEEE Journal on Selected Areas in Communications*, 31(12):2632–2645, 2013.

[21] S. H. Shen, L. H. Huang, D. N. Yang, and W. T. Chen. Reliable multicast routing for software-defined networks. In *INFOCOM*, pages 181–189, April 2015.

[22] L. H. Huang, H. C. Hsu, S. H. Shen, D. N. Yang, and W. T. Chen. Multicast traffic engineering for software-defined networks. In *INFOCOM*, pages 1–9. IEEE, 2016.

[23] A. Nagata, Y. Tsukiji, and M. Tsuru. Delivering a file by multipath-multicast on open-flow networks. In *International Conference on Intelligent Networking and Collaborative Systems*, pages 835–840, 2013.

[24] K. Ogawa, T. Iwamoto, and M. Tsuru. One-to-many file transfers using multipath-multicast with coding at source. In *IEEE International Conference on High Performance Computing and Communications*, pages 687–694, 2016.

[25] M. Noormohammadpour and C. S. Raghavendra. DDCCast: Meeting Point to Multi-point Transfer Deadlines Across Datacenters using ALAP Scheduling Policy. *Technical Report, Department of Computer Science, University of Southern California*, Report No. 17-972, 2017.

[26] Long Luo, Hongfang Yu, and Zilong Ye. Deadline-guaranteed Point-to-Multipoint Bulk Transfers in Inter-Datacenter Networks. *ICC*, 2018.

[27] Long Luo, Klaus-Tycho Foerster, Stefan Schmid, and Hongfang Yu. Dartree: Deadline-aware multicast transfers in reconfigurable wide-area networks. In *27th IEEE/ACM International Symposium on Quality of Service (IWQoS 2019)*, June 2019.

[28] Mohammad Noormohammadpour, Cauligi S. Raghavendra, Srikanth Kandula, and Sriram Rao. QuickCast: Fast and Efficient Inter-Datacenter Transfers using Forwarding Tree Cohorts. *INFOCOM*, 2018.

[29] Ben Pfaff, Bob Lantz, Brandon Heller, et al. Openflow switch specification, version 1.1.0 implemented (wire protocol 0x02). `http://archive.openflow.org/documents/openflow-spec-v1.1.0.pdf`, 2011.

[30] How microsoft builds its fast and reliable global network. `https://azure.microsoft.com/en-us/blog/how-microsoft-builds-its-fast-and-reliable-global-network/`. visited on September 30, 2017.

[31] Mohammad Noormohammadpour, Cauligi S. Raghavendra, Srikanth Kandula, and Sriram Rao. Fast and efficient bulk multicasting over dedicated inter-datacenter networks, 2018.

[32] Nikolaos Laoutaris, Georgios Smaragdakis, Rade Stanojevic, Pablo Rodriguez, and Ravi Sundaram. Delay-tolerant bulk data transfers on the internet. *IEEE/ACM TON*, 21(6), 2013.

[33] Nikolaos Laoutaris, Michael Sirivianos, Xiaoyuan Yang, and Pablo Rodriguez. Inter-datacenter bulk transfers with netstitcher. In *SIGCOMM*, pages 74–85. ACM, 2011.

[34] Srinivasan Keshav and Sanjoy Paul. Centralized multicast. In *International Conference on Network Protocols*, pages 59–68. IEEE, 1999.

[35] S. Liang and D. Cheriton. Tcp-smo: extending tcp to support medium-scale multi-cast applications. In *Proceedings.Twenty-First Annual Joint Conference of the IEEE Computer and Communications Societies*, volume 3, pages 1356–1365, 2002.

[36] Brian Adamson, Carsten Bormann, Mark Handley, and Joe Macker. Nack-oriented reliable multicast (norm) transport protocol, 2009.

[37] M. J. Donahoo, M. H. Ammar, and E. W. Zegura. Multiple-channel multicast scheduling for scalable bulk-data transport. In *INFOCOM '99. Eighteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, volume 2, pages 847–855 vol.2, Mar 1999.

[38] S. Bhattacharyya, J. F. Kurose, et al. Efficient rate-controlled bulk data transfer using multiple multicast groups. *TON*, 11(6):895–907, 2003.

[39] Meeyoung Cha, W. Art Chaovalitwongse, Jennifer Yates, Aman Shaikh, and Sue Moon. Efficient and scalable provisioning of always-on multicast streaming services. *Computer Networks*, 53(16):2825 – 2839, 2009.

[40] D. Li, M. Xu, M. c. Zhao, C. Guo, Y. Zhang, and M. y. Wu. Rdcm: Reliable data center multicast. In *2011 Proceedings IEEE INFOCOM*, pages 56–60, 2011.

[41] A. Rodriguez, D. Kostic, and A. Vahdat. Scalability in adaptive multi-metric overlays. In *International Conference on Distributed Computing Systems*, pages 112–121, 2004.

[42] Karthik Nagaraj, Hitesh Khandelwal, Charles Killian, and Ramana Rao Kompella. Hierarchy-aware distributed overlays in data centers using dc2. In *COMSNETS*, pages 1–10. IEEE, 2012.

[43] Miguel Castro, Peter Druschel, Anne-Marie Kermarrec, Animesh Nandi, Antony Rowstron, and Atul Singh. Splitstream: High-bandwidth multicast in cooperative environments. In *SOSP*, pages 298–313. ACM, 2003.

[44] Tingwei Zhu, Fang Wang, Yu Hua, Dan Feng, et al. Mctcp: Congestion-aware and robust multicast tcp in software-defined networks. In *International Symposium on Quality of Service*, pages 1–10, June 2016.

[45] Michael Luby, Lorenzo Vicisano, Jim Gemmell, et al. The use of forward error correction (fec) in reliable multicast, 2002.

[46] M. Luby, J. Gemmell, L. Vicisano, L. Rizzo, and J. Crowcroft. Asynchronous layered coding (alc) protocol instantiation, 2002.

[47] Christos Gkantsidis, John Miller, and Pablo Rodriguez. Comprehensive view of a live network coding p2p system. In *IMC*, pages 177–188. ACM, 2006.

[48] A. Shokrollahi. Raptor codes. *IEEE Transactions on Information Theory*, 52(6):2551–2567, 2006.

[49] John W. Byers, Michael Luby, Michael Mitzenmacher, and Ashutosh Rege. A digital fountain approach to reliable distribution of bulk data. In *SIGCOMM*, pages 56–67. ACM, 1998.

[50] K. Jeacle and J. Crowcroft. Tcp-xm: unicast-enabled reliable multicast. In *ICCCN*, pages 145–150, 2005.

[51] IJsbrand Wijnands, Eric C. Rosen, Andrew Dolganow, Tony Przygienda, and Sam Aldrin. Multicast Using Bit Index Explicit Replication (BIER). RFC 8279, November 2017.

[52] Sen Su, Yiwen Wang, Sujuan Jiang, Kai Shuang, and Peng Xu. Efficient algorithms for scheduling multiple bulk data transfers in inter-datacenter networks. *International Journal of Communication Systems*, 27(12), 2014.

[53] Yiwen Wang, Sen Su, et al. Multiple bulk data transfers scheduling among datacenters. *Computer Networks*, 68:123–137, 2014.

[54] M. Hefeeda, A. Habib, B. Botev, et al. Promise: Peer-to-peer media streaming using collectcast. In *MULTIMEDIA*, pages 45–54. ACM, 2003.

[55] K. Suh, C. Diot, J. Kurose, L. Massoulie, C. Neumann, D. Towsley, and M. Varvello. Push-to-Peer Video-on-Demand System: Design and Evaluation. *IEEE Journal on Selected Areas in Communications*, 25(9):1706–1716, December 2007.

[56] Srinivas Narayana, Joe Wenjie Jiang, Jennifer Rexford, and Mung Chiang. Distributed wide-area traffic management for cloud services. *SIGMETRICS Perform. Eval. Rev.*, 40(1):409–410, June 2012.

[57] Dimitri Bertsekas and Robert Gallager. Data networks, 1987.

[58] The Internet Topology Zoo. `http://www.topology-zoo.org/`.

[59] FK Hwang and Dana S Richards. Steiner tree problems. *Networks*, 22(1):55–89, 1992.

[60] Dimitri Watel and Marc-Antoine Weisser. *A Practical Greedy Approximation for the Directed Steiner Tree Problem*, pages 200–215. Springer International Publishing, Cham, 2014.

[61] Evaluation of approximation algorithms for the directed steiner tree problem. `https://github.com/mouton5000/DSTAlgoEvaluation`. visited on Apr 27, 2017.

[62] Hong Zhang, Kai Chen, Wei Bai, et al. Guaranteeing deadlines for inter-datacenter transfers. In *EuroSys*, page 20. ACM, 2015.

[63] The internet topology zoo (dataset). `http://topology-zoo.org/dataset.html`.

[64] Arjun Roy, Hongyi Zeng, Jasmeet Bagga, George Porter, and Alex C. Snoeren. Inside the social network's (datacenter) network. In *SIGCOMM*, pages 123–137. ACM, 2015.

[65] CloudLab. `https://www.cloudlab.us/`.

[66] The Internet Topology Zoo (GEANT). `http://www.topology-zoo.org/files/Geant2009.gml`.

[67] NACK-Oriented Reliable Multicast (NORM). `https://www.nrl.navy.mil/itd/ncs/products/norm`.

[68] Phillipa Gill, Martin Arlitt, Zongpeng Li, and Anirban Mahanti. Youtube traffic characterization: A view from the edge. In *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement*, IMC '07, pages 15–28, New York, NY, USA, 2007. ACM.