

# Datacast: A Scalable and Efficient Reliable Group Data Delivery Service for Data Centers

Jiaxin Cao, Chuanxiong Guo, Guohan Lu, Yongqiang Xiong, Yixin Zheng, Yongguang Zhang, Yibo Zhu, Chen Chen, and Ye Tian

**Abstract**—Reliable Group Data Delivery (RGDD) is a pervasive traffic pattern in data centers. In an RGDD group, a sender needs to reliably deliver a copy of data to all the receivers. Existing solutions either do not scale due to the large number of RGDD groups (e.g., IP multicast) or cannot efficiently use network bandwidth (e.g., end-host overlays).

Motivated by recent advances on data center network topology designs (multiple edge-disjoint Steiner trees for RGDD) and innovations on network devices (practical in-network packet caching), we propose *Datacast* for RGDD. *Datacast* explores two design spaces: 1) *Datacast* uses multiple edge-disjoint Steiner trees for data delivery acceleration. 2) *Datacast* leverages in-network packet caching and introduces a simple soft-state based congestion control algorithm to address the scalability and efficiency issues of RGDD.

Our analysis reveals that *Datacast* congestion control works well with small cache sizes (e.g., 125KB) and causes few duplicate data transmissions (e.g., 1.19%). Both simulations and experiments confirm our theoretical analysis. We also use experiments to compare the performance of *Datacast* and BitTorrent. In a BCube(4, 1) with 1Gbps links, we use both *Datacast* and BitTorrent to transmit 4GB data. The link stress of *Datacast* is 1.01, while it is 1.39 for BitTorrent. By using two Steiner trees, *Datacast* finishes the transmission in 16.9s, while BitTorrent uses 52s.

**Index Terms**—Multicast, congestion control, content distribution

## I. INTRODUCTION

**R**ELIABLE Group Data Delivery (RGDD) is widely used in cloud services (e.g., GFS [15] and MapReduce [5]) and applications (e.g., social networking, Search, scientific computing). In RGDD, we have a group which contains one data source and a set of receivers. We need to reliably deliver the same copy of bulk data from the source to all the receivers.

Existing solutions for RGDD can be classified into two categories: 1) Reliable IP multicast. IP multicast suffers from scalability issues, since it is hard to manage a large number of group states in the network. Adding reliability is also

challenging, due to the ACK implosion problem [13]. 2) End-host based overlays. Overlays are scalable, since devices in the network do not maintain group states. Reliability is easily achieved by using TCP in overlays. However, overlays do not use network bandwidth efficiently. The same copy of data may traverse the same link several times, resulting in high link stress. For example, ESM [19] reported that the average and worst-case link stresses are 1.9 and 9, respectively.

Motivated by the recent progresses on data center network (DCN) topologies and network devices, we explore new opportunities in supporting RGDD for DCN: 1) Recently proposed DCN topologies have multiple edge-disjoint Steiner trees<sup>1</sup>, which has not been well studied before. These multiple Steiner trees may enable full utilization of DCN bandwidth. 2) There is a clear technical trend that network devices are providing powerful packet processing abilities by integrating CPUs and large memory. This makes in-network packet caching practical. By leveraging in-network packet caching, we can address the scalability and bandwidth efficiency issues of RGDD.

However, it is challenging to take advantage of these opportunities. The multiple edge-disjoint Steiner trees problem has been studied for decades. Unfortunately, existing algorithms [6] cannot generate enough edge-disjoint Steiner trees within a short time, even in well structured data center networks. Although network devices are becoming capable of in-network packet caching, the resource is not unlimited. We need to use as small caches as possible for each group to maximize the number of simultaneously supported groups. At the same time, we need to increase bandwidth efficiency by reducing duplicate packets transmitted in the network.

In this paper, we design *Datacast* to address the above challenges. Leveraging the properties of the DCN topologies, *Datacast* introduces an efficient algorithm to calculate multiple edge-disjoint Steiner trees, and then distributes data among them. In each Steiner tree, *Datacast* leverages the concept of CCN [14]. To help *Datacast* achieve high bandwidth efficiency with small cache size in intermediate nodes, we design a rate-based congestion control algorithm, which follows the classical Additive Increase and Multiplicative Decrease (AIMD) approach. *Datacast* congestion control leverages a key observation: the receiving of a duplicate packet request at the source can be interpreted as a congestion signal. Different from previous work (e.g., TFMCC [27] and pgmcc [22]),

<sup>1</sup>In this paper, we define a Steiner tree as a tree whose root is the data source, and spans all the receivers.

Manuscript received November 29, 2012; revised May 21, 2013.

J. Cao, C. Guo, G. Lu, Y. Xiong, and Y. Zhang are with Microsoft Research Asia (e-mail: {jjiaoc, chguo, gulv, yqx, ygz}@microsoft.com). J. Cao is also with University of Science and Technology of China.

Y. Zheng is with Tsinghua University (e-mail: zhengyx12@mails.tsinghua.edu.cn).

Y. Zhu is with the University of California, Santa Barbara (e-mail: yibo@cs.ucsb.edu).

C. Chen is with the University of Pennsylvania (e-mail: chenche@seas.upenn.edu).

Y. Tian is with the University of Science and Technology of China (e-mail: yetian@ustc.edu.cn).

Digital Object Identifier 10.1109/JSAC.2013.131205.

which uses explicit information exchanges between the source and receivers, Datacast is much simpler. To understand the performance of Datacast, we build a fluid model. By analyzing the model, we prove that Datacast works at the full rate when the cache size is greater than a small threshold (e.g., 125KB), and also derive the ratio of duplicate data sent by the data source (e.g., 1.19%). We have built Datacast in NS3, and also have implemented it with the ServerSwitch [8] platform. Simulations and experiments verify our theoretical results, which suggest that Datacast achieves both scalability and high bandwidth efficiency.

This paper makes the following contributions: 1) We design a simple and efficient multicast congestion control algorithm, and build a fluid model to understand its properties. 2) We propose a low time-complexity algorithm for multiple edge-disjoint Steiner trees calculation. 3) We implement Datacast with the ServerSwitch platform, and validate its performance.

## II. BACKGROUND

### A. Reliable group data delivery

In data center applications and services, Reliable Group Data Delivery (RGDD) is a pervasive traffic pattern. The problem of RGDD is, *given a data source, Src, and a set of receivers,  $R_1, R_2, \dots, R_n$ , how to reliably transmit bulk data from Src to all the receivers*. A good RGDD design should be scalable and achieve high bandwidth efficiency. The following cases are typical RGDD scenarios.

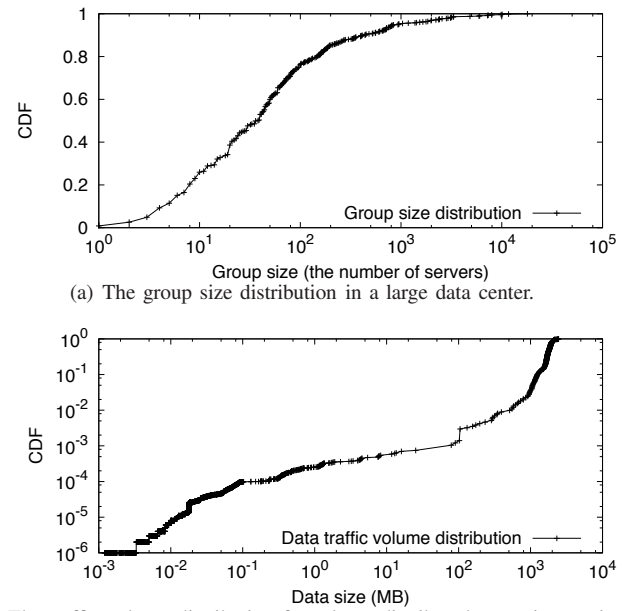
Case 1: In data centers, servers are typically organized as physical clusters. During bootstrapping or OS upgrading, the same copy of the OS image needs to be transferred to all the servers in the same cluster. A physical cluster is further divided into sub-clusters of different sizes. A sub-cluster is assigned to a service. All the servers in the same sub-cluster may need to run the same set of applications. We need to distribute the same set of program binaries and configuration data to all the servers in the sub-cluster.

Case 2: In distributed file systems, e.g., GFS [15], a chunk of data is replicated to several (typically three) servers to improve reliability. The sender and receivers form a small replication group. A distributed file system may contain tens of Peta bytes using tens of thousands machines. Hence the number of replication groups is huge. In distributed execution engine, e.g., Dryad [20], a copy of data may need to be distributed to many servers for JOIN operations.

Case 3: In Amazon EC2 or Windows Azure, a tenant may create a set of virtual machines. These virtual machines form an isolated computing environment dedicated to that tenant. When setting up the virtual machines, customized virtual machine OSes and application images need be delivered to all the physical servers that host these virtual machines.

Figure 1(a) and 1(b) show the group size and traffic volume distributions for a RGDD service in a large production data center. We use these two figures to show the challenges in supporting RGDD.

**The system should be scalable.** As we have mentioned in the above scenarios, we need to support a large number of RGDD groups in large data centers. Figure 1(a) further shows that the group size varies from several servers to thousands of servers



(a) The group size distribution in a large data center.

(b) The traffic volume distribution for a large distributed execution engine.

Fig. 1. RGDD groups and traffics in data centers.

and even more. The large number of groups and the varying group sizes pose scalability challenges, since maintaining a large number of group states in the network is hard (as demonstrated by IP multicast).

**Bandwidth should be efficiently and fully used.** Figure 1(b) shows the traffic volume distribution for group communications. It shows that the groups transmitting more than 550MB data contribute 99% RGDD data traffic volume. Due to the large number of groups and the large data sizes, RGDD contributes a significant amount of traffic. This requires that RGDD uses network bandwidth efficiently. On the other hand, the new DCN topologies (e.g., BCube [7] and CamCube [9]) provide high network capacity with multiple data delivery trees. An RGDD design should take full advantage of these new network topologies to speedup data delivery.

In what follows, we introduce recent technology progresses on DCN topologies and network devices, which we leverage to address the above challenges.

### B. New opportunities

**Multiple edge-disjoint Steiner trees.** Different from the Internet, DCNs are owned and operated by a single organization. As a result, DCN topologies are known in advance, and we can assume that there is a centralized controller to manage and monitor the whole DCN. Leveraging such information, we can improve RGDD efficiency by building efficient data delivery trees. Furthermore, several recently proposed DCNs (e.g., BCube [7] and CamCube [9]) have multiple edge-disjoint Steiner trees which can be used to further accelerate RGDD.

**In-network packet caching becomes practical.** Recently, we observe a clear technical trend for network devices (switches and routers). First, powerful CPUs and large memory are being included in network devices. The new generation of devices are equipped with multi-core X64 CPUs and several

GB memory, e.g., Arista 7504 has 2 AMD Athlon X64 Dual-Core CPUs and 4GB DRAM. Second, the merchant switching ASIC, CPU and DRAM can be connected together by using the state-of-the-art PCI-E interface, as demonstrated by research prototype (e.g., ServerSwitch [8]) and products (e.g., Force10 S7000 [21]). With the new abilities of network devices, many in-network packet processing operations (e.g., in-network packet caching) become practical. In this paper, we explore in-network packet caching. By turning hard-states for group managements in intermediate network devices into soft-states based packet caching, we address the scalability and efficiency issues of RGDD.

However, technical challenges exist to take advantage of these opportunities. First, given the network topology, calculating one single Steiner tree with minimal cost is NP-hard [16]. What is more challenging is that we have to calculate multiple Steiner trees, and the calculation has to be fast enough (otherwise it may be more time consuming than data dissemination). Second, we have a large number of RGDD groups to support and have limited resources in intermediate network devices. How to use as few resources as possible to support more RGDD groups is a challenge.

We design Datacast to explore the new design spaces provided by the new opportunities. The design goal of Datacast is *to achieve scalability and also high bandwidth efficiency*. In what follows, we first introduce the architecture of Datacast, then describe how Datacast addresses the above technical challenges.

### III. DATACAST OVERVIEW

Figure 2 shows the architecture of Datacast. There are five components in Datacast: Fabric Manager, Master, data source, receivers, and intermediate devices (IMD). Fabric Manager is a centralized controller, which maintains a global view of the network topology. When we need to start an RGDD group, we first start a Master. The Master will get topology information from Fabric Manager and then calculate multiple edge-disjoint Steiner trees. After that, the Master will send the tree information and other signalling messages (e.g., which file to fetch) to receivers via a signalling protocol. Then data transmission begins. When transmitting data, the data source will run our congestion control algorithm. During the whole process, intermediate devices do not interact with Fabric Manager, Master, the source or any receivers. These devices just cache and service data based on their local decisions.

To deliver signalling messages efficiently, we have built a signalling protocol, which uses a hierarchical transmission tree structure (generated by the Breadth First Search algorithm) to transmit signalling messages. It encodes the transmission tree into the message. Each node in the transmission tree decodes the signalling message, splits the tree into subtrees and forwards each subtree to its corresponding children. When the signalling messages reach the leaves, ACKs are generated and aggregated along the paths from leaves to the root. Using the message split and aggregation, signalling messages can be reliably and efficiently delivered.

In large data centers, failures are inevitable. Different from BitTorrent [4], which achieves fault tolerant in a distributed

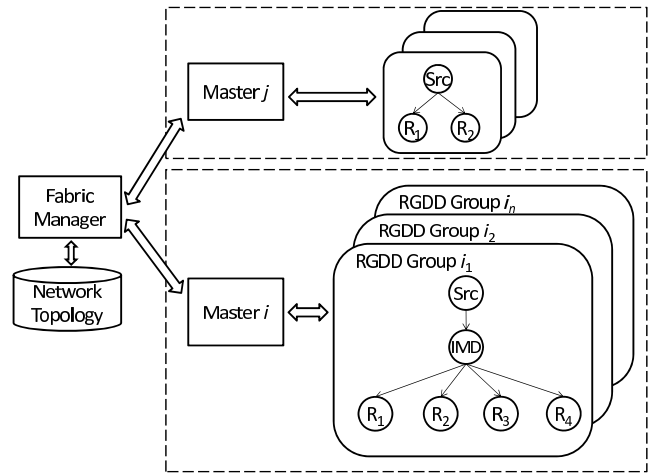


Fig. 2. The architecture of Datacast.

way, Datacast handles network failures in a centralized manner. In Datacast, Fabric Manager monitors the network status in real time. When network failures happen, Fabric Manager will send the new topology information to all the Masters, and each Master will recalculate the Steiner trees and notify the affected receivers accordingly.

To monitor the network status in real time, LSAs (Link State Advertisement) are used. A network device sends LSAs to all its direct neighbors under two conditions: 1) A network device sends LSAs periodically (e.g., 5s). 2) A network device sends LSAs when it detects link state changes (e.g., a link encounters a failure). To detect link state changes, each network device uses a simple heartbeat protocol. When a network device receives a new LSA, it forwards the LSA to all its ports except the incoming one. Fabric Manager uses the latest received LSAs to decide the real time network status and construct the spanning tree for signaling delivery.

In the following sections, we will present two key designs of Datacast: the fast calculation of multiple edge-disjoint Steiner trees, and the Datacast congestion control protocol which helps Datacast achieve scalability and high bandwidth efficiency.

### IV. MULTIPLE EDGE-DISJOINT STEINER TREES IN DCN

In this section, we first present the algorithm on multiple Steiner trees calculation, then discuss how to use these multiple Steiner trees for data delivery.

#### A. Calculation of multiple Steiner trees

It has been known that using multiple Steiner trees can improve the transmission efficiency [3]. However, constructing multiple edge-disjoint Steiner trees in a given (data center) topology has not been investigated before. The problem is, for a given network  $G(V, E)$ , where  $V$  is the set of nodes and  $E$  is the set of edges, and a group  $D$  containing one source and a set of receivers, how to calculate the maximum number of edge-disjoint Steiner trees. This is the well known multiple edge-disjoint Steiner trees problem, which has been studied for decades. Unfortunately, calculating Steiner trees is NP-hard [16].

```

// G is the DCN network, D is the Datacast group.
CalcSteinerTrees(G, D):
  // 1) construct multiple spanning trees
  SPTSet = G.CalcSpanningTrees(D.src);

  // 2) prune each spanning trees
  foreach (SPT in SPTSet)
    SteinerTree = Prune(SPT, D);
    SteinerTreeSet.add(SteinerTree);

  // 3) repair Steiner trees if they are broken
  foreach (SteinerTree in SteinerTreeSet)
    if (SteinerTree has broken links)
      if (RepairSteinerTree(SteinerTree, G) == false)
        Release(SteinerTree);
        SteinerTreeSet.remove(SteinerTree);
  return SteinerTreeSet;

```

Fig. 3. The algorithm for multiple edge-disjoint Steiner trees calculation.

We therefore turn our attention to heuristic algorithms. One reasonable approach is as follows. There are algorithms for calculating multiple edge-disjoint spanning trees (e.g., [6]). We can first find the multiple edge-disjoint spanning trees, and then prune the unneeded edges and nodes to get the Steiner trees.

However, the generic multiple spanning trees algorithms do not work well for our case. First, the time complexity of calculating the spanning trees is high. The best algorithm we know is Po's algorithm [25]. Its time complexity is  $O((k')^2|V||E|)$ , which is too high for RGDD (we will see that in Section VI-A1). Second, the depths of the spanning trees generated by the generic algorithm can be very large. For example, the average and worst-case depths of the trees for RGDDs in BCube can be 1000+ and 2000+ hops, whereas the network diameter is only 8.

Fortunately, we observe that DCNs, e.g., Fattree, BCube and multi-dimensional Torus, are well structured topologies. These topologies are also well studied. Multiple spanning trees construction algorithms for these topologies are already known (e.g., [7], [23]), and these spanning trees have good qualities, e.g., small tree depths. However, network failures (e.g., link failures) are common in real networks. Without reorganizing the spanning trees, network failures could possibly break all the trees generated by these algorithms. In order to solve the problem, we propose a multiple edge-disjoint Steiner trees algorithm, which is shown in Figure 3. The algorithm contains three parts.

The first part of this algorithm uses specific algorithms to construct spanning trees for specific DCN topologies (without considering network failures). For example, in Fattree [1], Breadth First Search (BFS) can generate a spanning tree, and the spanning trees algorithms for BCube and Torus are proposed in [7] and [23]. The time complexity of these algorithms are  $O(k|V|)$ , where  $k$  is the number of edge disjoint Spanning trees.

The second part prunes the links that are not used in data transmissions. To prune the spanning tree, we calculate the paths from the receivers to the source in the spanning tree. Then the set of links involved in the paths form a Steiner tree. The time complexity of pruning all the spanning trees is  $O(|E|)$ , since each link will only be traversed once.

The third part tries to repair the broken trees affected by link failures. The core idea of repairing a Steiner tree is: *we first release the broken tree, and then try to use BFS to traverse the free and active links to construct a new Steiner tree*. The repairing algorithm applies this idea to the broken trees one by one as shown in Figure 3. Although this idea is simple, it has the following benefits: 1) It guarantees at least one Steiner tree if all the receivers are connected. 2) The depth of the tree is locally minimized due to the use of BFS. The time complexity of repairing all the trees is  $O(k'|E|)$ , where  $k'$  is the number of Steiner trees to be repaired.

Our multiple Steiner trees calculation algorithm is fast. The time complexity of the algorithm is  $O(k|V|) + O(|E|) + O(k'|E|)$ , which contains the construction and pruning of spanning trees and the repairing of Steiner trees. Our algorithm also has good performance (in terms of the number of Steiner trees) and is fault tolerant. Even if there are network failures, we can still create a number of Steiner trees. We have derived an upper bound of the number of Steiner trees, and found that the number of Steiner trees generated by our algorithm is very close to the upper bound (details will be shown in Section VI-A2).

### B. Data distribution among multiple Steiner trees

To use multiple Steiner trees for data delivery, we first split the data into blocks, and then feed each tree with a block. When a Steiner tree finishes transmitting the last data packet of the current block, we know that the transmission of the current block is finished. Then the data source will use our signalling protocol to deliver the information of the next block to be transferred, e.g., the name of the block, to the receivers. After that the Steiner tree will start to transmit the next block. This process repeats until all the blocks are successfully delivered.

## V. DATACAST TRANSPORT PROTOCOL

In this section, we introduce in-network packet caching in Datacast, present the Datacast congestion control algorithm and discuss the cache management mechanism. By building a fluid model for the congestion control, we also derive the condition under which Datacast operates at the full rate, and its efficiency.

### A. Data transmission with in-network caching

In-network packet caching has been used in many previous works, including Active Networking [24], RE (redundancy elimination) [2], and CCN [14]. Datacast is built on top of CCN. In CCN, every single packet is assigned a unique, hierarchical name. A user needs to explicitly send an interest packet to ask for the data packet. Any intermediate device that has the requested data along the routing path can respond with the data packet. The network devices along the reverse routing path then cache the data packet in their content stores for later uses. CCN therefore turns group communication into in-network packet caching.

Datacast improves CCN as follows: 1) Datacast introduces a congestion control algorithm to achieve scalability and high bandwidth efficiency. 2) Datacast only caches data packets at

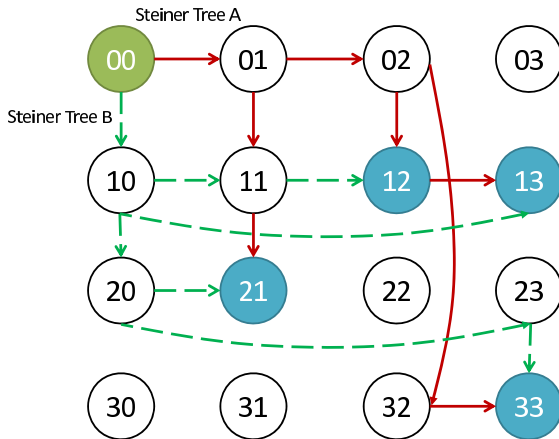


Fig. 4. An illustration of in-network caching.

branching nodes, which helps the whole system save memory. 3) Datacast uses source routing to enforce routing paths, so no Forwarding Information Base (FIB) management is needed at the intermediate devices.

Figure 4 shows an example of data delivery with in-network caching supports. The green node, 00, is the data source. The blue nodes, 12, 13, 21 and 33, are the receivers. Two Steiner trees calculated by the algorithm proposed in Section IV are shown in solid lines and dashed lines separately. The transmission in Steiner tree A could take the following steps: 1) Node 21 sends an interest packet to node 00 through the path {21, 11, 01, 00}. Node 00 sends the requested data back along the reverse path. Then the data packet is cached at the branch node 01. 2) Node 12 sends an interest packet along the path {12, 02, 01, 00} asking for the same data. When the interest arrives at node 01, node 01 finds that it has already cached the data packet, so it terminates the interest and sends back the data packet. Then the data are cached at node 02 and 12. 3) Node 13 sends its interest along the path {13, 12, 02, 01, 00}. Then the data is replied by node 12, since it has cached the data. 4) Node 33 sends its interest along the path {33, 32, 02, 01, 00}, and node 02 returns the data packet.

Note that the execution order of the four steps is not important. They can be executed in an arbitrary order, and still achieve the same result. The reason is that, in the end, all the steps together cover the same Steiner tree by traversing every link of the tree exactly once.

### B. Datacast congestion control algorithm

Datacast congestion control algorithm works for a single Steiner tree. It is one of the most important part of Datacast to realize its design goal, i.e., to achieve scalability and high bandwidth efficiency. Since Datacast turns hard group states into soft-state based packet caching, it is natural to require that the cache size in intermediate devices for each group is as small as possible (so as to support more groups), and the rates of receivers are synchronized (so as to improve bandwidth efficiency). If the rates of receivers are synchronized, only one copy of each packet is delivered in a Steiner tree. When

receivers have different receiving bandwidth, we expect all the rates of receivers are synchronized to the receiving rate of the slowest receiver.

A synchronized scheme may suffer from significant throughput degradation if a receiver in the group has a small receiving rate. In this case, we may either kick out the very slow receivers, or split the data delivery group into multiple ones. These topics are our future work.

Datacast uses the classical AIMD for congestion control. This is not new. What is new in Datacast is how congestion is detected. Datacast uses *duplicate interests* as congestion signals. A duplicate interest is an interest requiring the same data which has been asked before. The source receives a duplicate interest in the following two cases: 1) The network is congested, so some packets are dropped. Then the receiver will retransmit the interest, which serves as a duplicate interest. 2) Receivers are out of sync. When slow receivers cannot keep up with the fast ones, their interests will not be served by the cache of the intermediate devices. The interests will finally be sent to the data source, which serves as duplicate interests. In these cases, the source needs to slow down its sending rate. On the other hand, if there is no congestion and the rates of receivers are well synchronized, there will be no duplicate interests, and the source should increase its sending rate.

After congestion is detected, the rate adjustment becomes easy: when the source receives a duplicate interest, it decreases its sending rate by half; when no duplicate interest is received in a time interval  $T$ , the source increases the sending rate by  $\delta$ . Datacast congestion control is therefore rate-based. The source maintains and controls a sending rate  $r^2$ . Note that the sending rate of the duplicate data packet is not constrained by the congestion control, since the corresponding duplicate interest packets are from the slowest receiver, and the receiving rate of the slowest receiver should not be further reduced.

At the receivers' side, each receiver is given a fixed number of credit,  $w$ , which means that one receiver can send at most  $w$  interests into the network. When a receiver sends out an interest, the credit is decremented by one. When it receives a data packet, its credit is incremented by one. In Datacast, the guideline for setting  $w$  is to saturate the pipe. In a DCN with 1Gbps link, when the RTT is 200us (which is a typical network latency in a data center environment),  $w = 16$  can saturate the link. To achieve reliability, the receiver retransmits an interest if the data packet does not come back after a timeout. The timeout is calculated in the same way as TCP.

To summarize, Datacast congest control algorithm works as follows.

$$r = \begin{cases} \frac{r}{2} & \text{when a duplicate interest is received.} \\ r + \delta & \text{when there is no duplicate interest in } T. \end{cases}$$

As we can see, Datacast congestion control algorithm is simple. The source does not need to know which receiver is the slowest one, and what is the available bandwidth of that slowest receiver. In Section V-D, we will show analytically that Datacast uses small cache sizes and results in few duplicate data transmissions.

<sup>2</sup>To be exact, this is the rate of the source's token bucket.

### C. Cache management

To prevent cache interferences among different transmission trees, we use a *per-tree based* cache replacement algorithm. Each device uses a per Datacast tree based cache with size  $C$ . This is possible due to the following reasons: 1) A Datacast tree can be uniquely identified by a global unique tree transmission id (assigned by Master). 2) The cache size needed by each tree is small (as we will show in the next subsection).

In each tree, we find that the most popular data packets are the new ones, since new data packets will always be accessed by other receivers in the future. To keep new data packets in caches and erase old data packets, Datacast chooses First In First Out (FIFO) as its per-tree cache replacement policy. To prevent unpopular data packets from being put into caches, Datacast does not cache duplicate data packets.

Note that although this is a per-tree strategy, it is a scalable solution. The reasons are: 1) Compared with IP multicast, we do not need any protocol (e.g., IGMP) to maintain Datacast's per-tree states. Switches just use local decisions to manage its cache. 2) Datacast can work efficiently with small caches, e.g., 125KB, and large memory is expected for future network devices, e.g., 16GB memory for a switch. If it uses 4GB as Datacast cache, a network device can support up to 32k ( $\approx \frac{4GB}{125KB}$ ) simultaneous trees.

### D. Properties of Datacast congestion control algorithm

In this subsection, we study the following questions: 1) What is the condition for Datacast to work at the full rate (i.e., the receiving rate of the slowest receiver)? 2) When Datacast works at the full rate, how many duplicate data will be sent from the data source? We define the *duplicate data ratio* as the ratio of the duplicate data sent by the source to all the new data sent. To answer these questions, we have built a fluid model and derived the following theorems<sup>3</sup>. (Details are presented in Appendix.)

*Theorem 1:* Datacast works at the full rate, i.e., the rate of the slowest receiver,  $R$ , if the cache size,  $C$ , satisfies

$$C > \frac{R^2 T}{2\delta} - (w \cdot MTU - R \cdot RTT_m) \quad (1)$$

where  $RTT_m$  is the slowest receiver's minimum round trip time (the pingback RTTs).

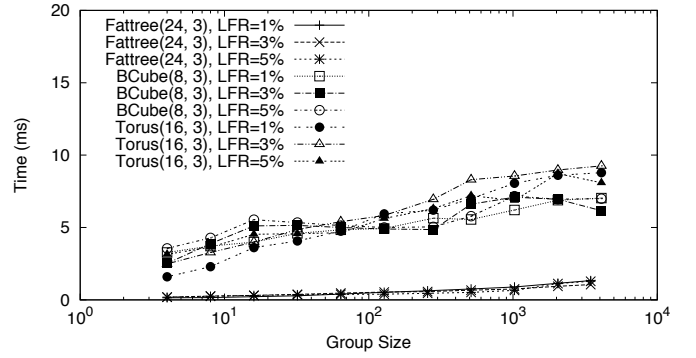
*Theorem 2:* When Datacast works at the full rate, the duplicate data ratio of Datacast is lower than or equal to

$$\frac{\frac{\delta}{T}}{\frac{\delta}{T} + \frac{R}{2MTU + RTT_m}}$$

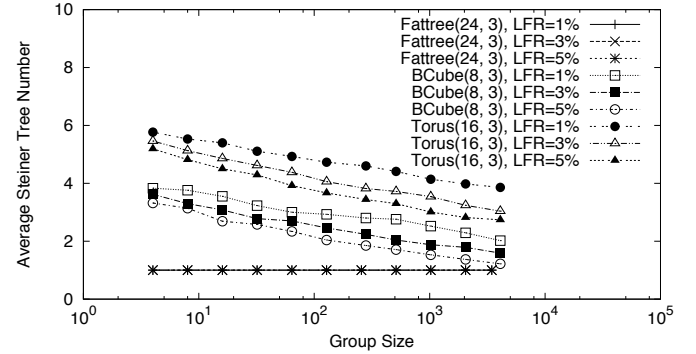
the equal sign is true when  $RTT_m = 0$ .

Theorem 1 tells us Datacast works at the full rate when the cache size is greater than  $\frac{R^2 T}{2\delta} - (w \cdot MTU - R \cdot RTT_m)$ . For example, when  $\delta = 5\text{Mbps}$ ,  $T = 1\text{ms}$ ,  $R = 100\text{Mbps}$  and the credit number is just enough to saturate the pipe (i.e.,  $w \cdot MTU = R \cdot RTT_m$ ), Datacast works at the full rate when the cache size is larger than 125KB. Theorem 2 reveals the bandwidth efficiency of Datacast. In the above example,

<sup>3</sup>These results nicely fall back to the ones in our previous work [10] when latencies are ignored.



(a) The running times.



(b) The numbers of Steiner trees.

Fig. 5. Performance of our multiple Steiner trees algorithm.

the duplicate data ratio is 1.19% when RTT is ignorable. Theorem 1 and 2 tell us that Datacast achieves the goal of high bandwidth efficiency, and also meets the requirement of using small cache size in the intermediate devices.

## VI. SIMULATION

### A. Evaluation of the multiple Steiner trees algorithm

To study the performance of the multiple Steiner trees algorithm, we use a Dell PowerEdge R610 server, which has two E5520 Intel Xeon 2.26GHz CPUs and 32GB RAM. We study our algorithm under three topologies, Fattree(24, 3), BCube(8, 3) and Torus(16, 3). The BCube and Torus contain 4096 servers, while the Fattree contains 3456 servers. For each simulation, we randomly generate link failures. The link failure rates (LFR) include 1%, 3% and 5%. We ignore the cases when the network is not connected.

1) *Running time:* Figure 5(a) shows the running times of our algorithm. From the results, we can see that our algorithm can finish all of the tree calculations within 10ms. We compared our algorithm with the generic algorithm which first calculates the spanning trees using Po's algorithm [25], then prunes them to get Steiner trees. The time complexity of the generic algorithm is dominated by the spanning tree calculation. The times needed for calculating spanning trees for Fattree(24, 3), BCube(8, 3) and Torus(16, 3) are 1, 39 and 42 seconds respectively. This algorithm therefore cannot be used in Datacast.

2) *Steiner tree number:* Figure 5(b) shows the numbers of Steiner trees constructed by our algorithm. For BCube and Torus, the numbers of Steiner trees decrease as the group size

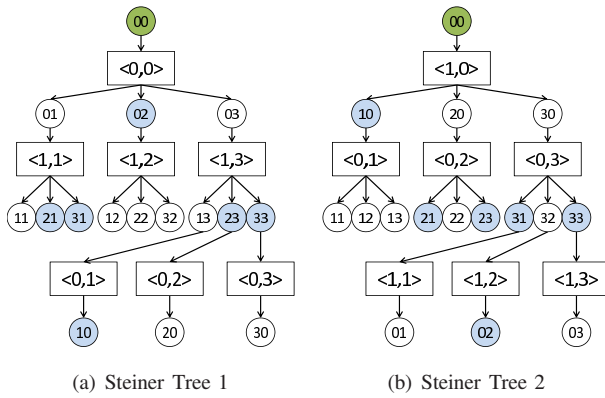


Fig. 6. The simulation and experiment setup.

and the link failure rate increase. This is expected, since a large group would experience more link failures, and more link failures will break more trees. Though Fattree has only one Steiner tree, our algorithm helps on failure recovery when the original tree is broken by link failures.

To check whether our algorithm can create enough Steiner trees, we have derived an upper bound of the Steiner tree number, which is the minimum value of the out-degree of the source and the in-degrees of all the receivers. The Steiner tree numbers produced by our algorithm are only 0.8% less than the bounds on average.

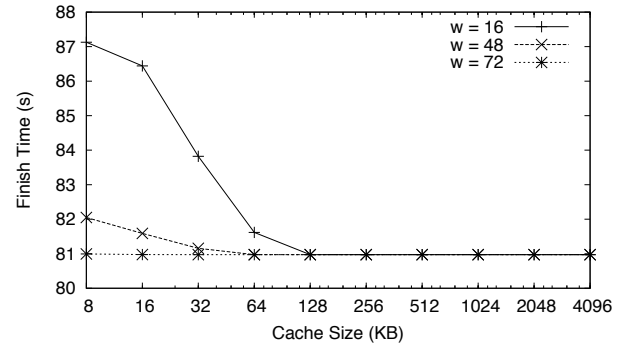
3) *Steiner tree depths*: Our algorithm also guarantees small tree depths. For example, when the link failure rate is 1%, the average Steiner tree depths for BCube, Torus and Fattree, are 9.99, 24.31 and 6.00, respectively.

### B. Micro benchmarks for Datacast congestion control algorithm

We have built Datacast in NS3. In this subsection, we use micro benchmarks to study Datacast congestion control algorithm in a BCube(4, 1). We use a single multicast tree shown in Figure 6(a). The green node, 00, is the source, while the blue ones, 02, 10, 21, 23, 31 and 33, are the receivers.  $\delta = 5\text{Mbps}$ ,  $T = 1\text{ms}$  and  $\text{MTU} = 1.5\text{KB}$ . The link rates are  $1\text{Gbps}$ , and the propagation delays are  $5\mu\text{s}$ . We slow down the link from switch  $\langle 0,0 \rangle$  to node 02 to  $100\text{Mbps}$  to make node 02 the slowest receiver. The queue size for each link is 100 packets. The headers of the interest and data packets are both 16 bytes. The initial rate of the source is  $500\text{Mbps}$ .

1) *Full Rate Cache Requirement*: We first verify Theorem 1. We vary the cache sizes from  $8\text{KB}$  to  $4096\text{KB}$ . Given the credit numbers are 16, 48 and 72 packets, the bounds derived by Theorem 1 are  $102\text{KB}$ ,  $54\text{KB}$  and  $18\text{KB}$  respectively. The simulation results are shown in Figure 7(a). The results suggest that Datacast works at the full rate when the cache size is larger than the bound. Its throughput,  $98.799\text{Mbps}$ , is very close to the optimal results, which is  $98.933\text{Mbps}$  ( $= 100\text{Mbps} \times \frac{1500-16}{1500}$ ). The results also suggest that Datacast experiences graceful throughput degradation when there is not enough cache.

2) *Duplicate Data Ratio*: To verify Theorem 2, we vary the rate increase,  $\delta$ , from  $0.10\text{Mbps}$  to  $102.40\text{Mbps}$ . In an empty network (no traffic), the round trip time is ignorable, so the



(a) Datacast's finish times under different cache sizes.

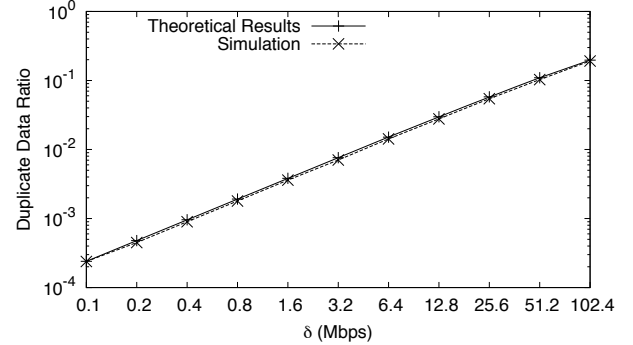
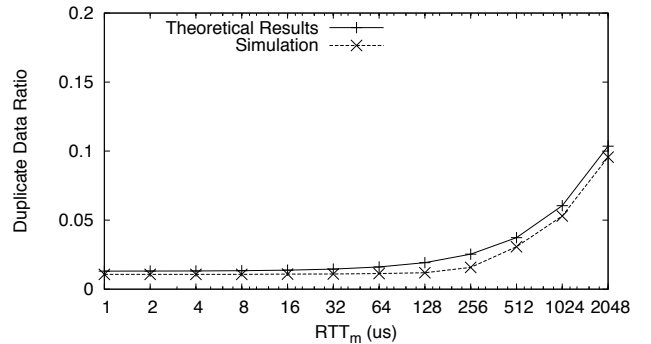
(b) Duplicate data ratio vs.  $\delta$ (c) Duplicate data ratio vs.  $RTT_m$ 

Fig. 7. The finish times and duplicate data ratios of Datacast.

duplicate data ratio is  $\frac{\delta}{T} / (\frac{\delta}{T} + \frac{R^2}{2\text{MTU}})$ . From the results shown in Figure 7(b), we can see that the duplicate data ratio derived from our model is consistent with the simulation results.

We also study the duplicate data ratio under the congestion case. We add a queueing delay at the slow link, which varies from  $1\mu\text{s}$  to  $2\text{ms}$ . The results are shown in Figure 7(c), which suggest that Theorem 2 captures the trend of the increase of duplicate data ratios as the latency grows. From the results, we can also see that even if congestion happens, the duplicate data ratio is still lower than 0.1.

3) *Performance under packet losses*: To see whether Datacast is resilient to packet losses, we randomly drop data packets at the link from switch  $\langle 0,0 \rangle$  to node 02. The cache sizes are set to  $128\text{KB}$ . When the packet loss rate is  $1.02\%$ , the finish time only increases by  $2.76\%$  and the duplicate ratio is  $1.23\%$ .

4) *Fairness*: In this simulation, we set all the links back to  $1\text{Gbps}$ . To study intra-protocol (inter-protocol) fairness, we use the Datacast group to compete with nine other Datacast

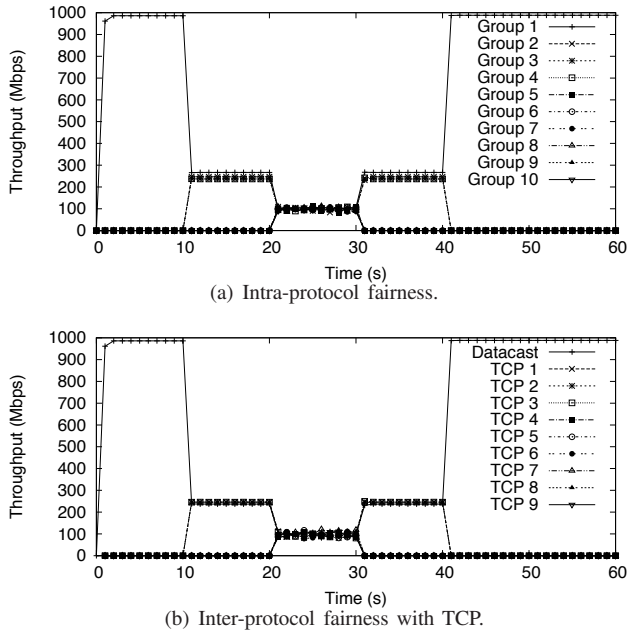


Fig. 8. Intra-protocol and inter-protocol fairness.

groups (TCP flows). Three of them start at 10s, end at 40s. Six of them start at 20s, end at 30s. Figure 8 shows the results. We can see that Datacast achieves good intra-protocol (inter-protocol) fairness.

Datacast achieves good inter-protocol fairness with TCP, since their additive increase parts are at the same magnitude. In this simulation, we measure that the RTT of TCP is about 1ms when there are nine TCP flows and one Datacast group. TCP increases its rate at the speed of 12Mbps ( $= \frac{MTU}{RTT}$ ) per RTT (1ms), while Datacast increases its rate at the speed of 5Mbps per millisecond. Therefore, Datacast and TCP achieve good inter-protocol fairness.

5) *Cache replacement algorithms*: We study the performance of Datacast with three different cache management policies, Least Recently Used (LRU), Least Frequently Used (LFU) and First In First Out (FIFO). The cache miss ratios for LRU, LFU and FIFO are 3.90%, 1.63% and 1.12%, respectively. FIFO achieves the minimum duplicate data ratio of them, since it always stores new data packets in the cache, which will be used in the future.

### C. Performance comparison

BitTorrent was originally designed for P2P file sharing in the Internet. Since a data center is a collaborative environment and the network topology can be known in advance, we use techniques similar to Cornet [11] to improve the original BitTorrent. Cornet improvements include: a server does not immediately leave the system after it receives all the content; no SHA1 calculation per block; use large block size (4MB). Cornet suggests using large block size (4MB). Our simulations demonstrate that smaller block size results in better performance. We choose 108KB as the block size in the simulations. We call the Cornet optimized version BT-Cornet. Similar to Cornet, we also consider the topology awareness. Since we have rich topological information, we design the following

neighbor selection algorithm: a server selects 10 peers (when the group size is less than 10, all the members are peers). It sorts the group members via the distance. It prefers peers with a small distance, but guarantees that at least one member (if it exists) is selected as its peer at each distance range. Similar to Cornet, tit-for-tat and choke-unchoke are disabled. We call the optimized version BT-Optimized.

We use two metrics for the comparison. The first metric is the network stress, which is the sum of all the bytes transmitted on all the links. The second is the finish time.

In all the simulations, the source sends 500MB data. Figure 9 shows the performance of Datacast, BT-Cornet, BT-Optimized under different group sizes for three different topologies, Fattree(24, 3), BCube(8, 3) and Torus(16, 3). The group size varies from 8 to 1024. Our results clearly demonstrate that Datacast is better than BT-Cornet and BT-Optimized in terms of the network stress and the finish time. On BCube and Torus, Datacast is much faster since each server has multiple 1Gbps ports. In all the simulations, the network stresses of BT-Optimized are 1.2-3.5X than Datacast, and Datacast is 1.1-3.7X faster than BT-Optimized.

We also note that in our simulations, when the topology is Fattree, the finish time with BT-Cornet is smaller than with BT-Optimized. This is because with BT-Optimized, we prefer peers that are close with each other. This preference may result in small cliques which may not be fully connected. BCube does not have such an issue because its structure does not have hierarchy.

In the experiments, Datacast's finish times are quite close to the ideal cases. There is one Steiner tree in Fattree(24, 3), and there are four Steiner trees in BCube(8, 3), and six in Torus(16, 3). Therefore the ideal finish times are 4s, 1s and 0.67s for Fattree(24, 3), BCube(8, 3) and Torus(16, 3), respectively. The finish times of Datacast are 0.67% larger than the ideal cases on average. Datacast is also efficient. The average link stress of Datacast is only 1.002, which means that each packet only traverses each Steiner tree link 1.002 times on average.

## VII. IMPLEMENTATION

### A. ServerSwitch based implementation

We have implemented Datacast using the design shown in Figure 2. Fabric Manager, Master, data source and receivers are all implemented as user-mode applications. Each node in the data center runs a Datacast daemon, which is responsible for forwarding and receiving signalling messages. When Datacast is trying to start a group for data transmission, it first starts a Master process. The Master process calculates multiple Steiner trees, and then sends signalling messages to the group members. The daemons on these nodes will start the data source process and the receiver processes. Then the transmission starts.

To cache data packets in intermediate nodes, we use the ServerSwitch platform [8]. ServerSwitch is composed of an ASIC switching chip and a commodity server. The switching chip is connected to the server CPU and memory using PCI-E. ServerSwitch's switching chip is programmable. It uses a TCAM table to define operations for specific types of packets. To implement data packet caching in switches, we use User



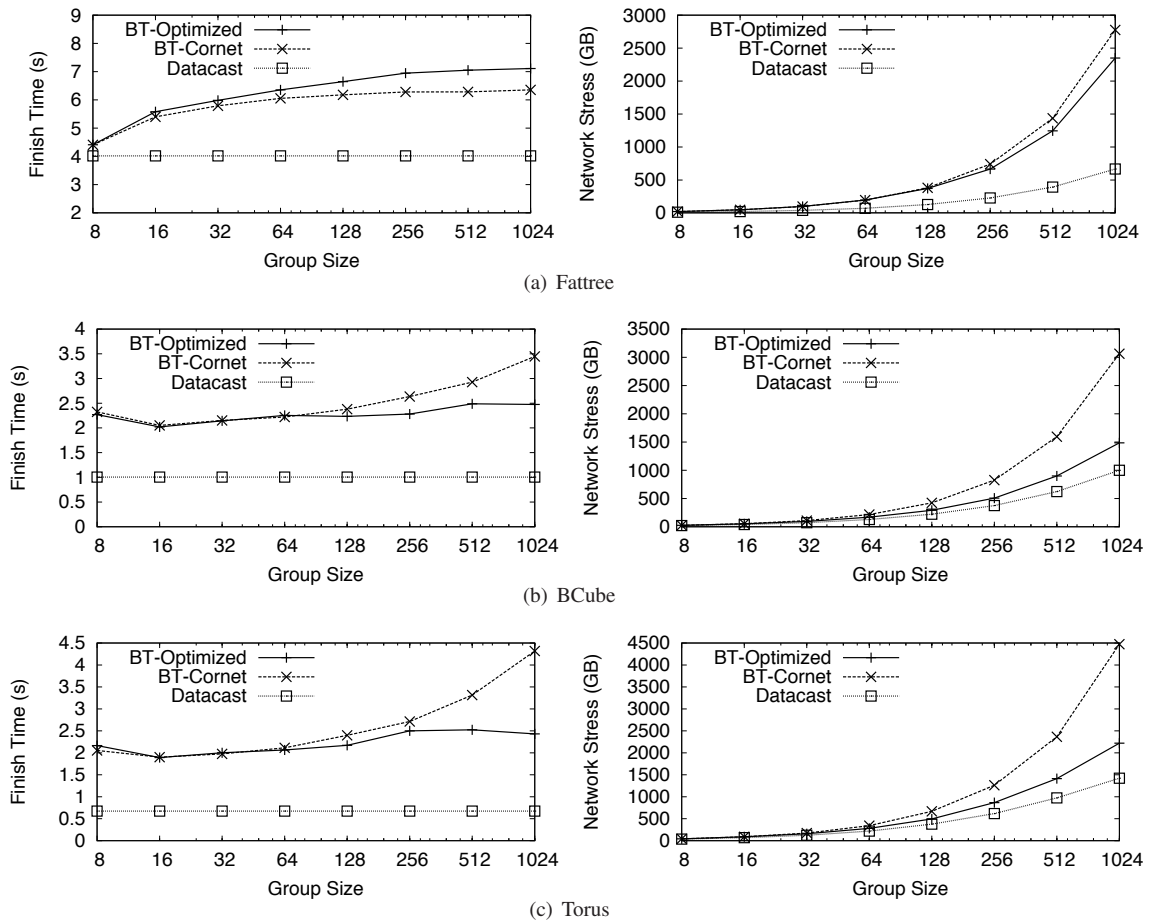


Fig. 9. Performance comparison of Datacast and BitTorrent.

Defined Lookup Keys (UDLK) to forward data packets to the Datacast kernel mode driver at branch nodes. The driver is used to do the in-network data packet caching. At non-branch nodes, the data packets are directly forwarded by hardware.

### B. Evaluation

In this subsection, we use our real testbed implementation to evaluate Datacast. We use a BCube(4, 1) with 1Gbps links for our study.

1) *Efficiency study*: We study Datacast's performance when different cache sizes are set for branching nodes. We use a single Steiner tree shown in Figure 6(a) and slow down the link from switch  $\langle 0,0 \rangle$  to node 02 to 100Mbps. We let  $\delta = 5\text{Mbps}$ ,  $T = 1\text{ms}$  and  $w = 8$ . The minimum round trip time is about 350 $\mu\text{s}$ . Based on Theorem 1, Datacast works at the full rate when the cache size is larger than 120KB. When we use 64KB (or 32KB) cache, the average throughput is 96.774Mbps (or 88.757Mbps), which is still acceptable due to the graceful throughput degradation of Datacast. When the cache size is 128KB, the average throughput is 98.684Mbps, and the duplicate data ratio is 1.45%, which is lower than the theoretical bound derived by Theorem 2, 2.87%.

2) *Performance comparison*: We compare the performance of Datacast with BitTorrent (we use  $\mu\text{torrent}$ ). In this experiment, we use both Datacast and BitTorrent to transfer 4GB

TABLE I  
PERFORMANCE COMPARISON OF DATACAST AND BITTORRENT.

	Finish Time (s)	Link Stress
Datacast	16.9	1.01
BitTorrent	41-52	1.39

data. The cache size on each branch node is 512KB. For Datacast,  $\delta = 125\text{Mbps}$  and  $T = 1\text{ms}$ .

Datacast finishes the transmission within 16.9s. The source achieves 1.89Gbps throughput on average, which is close to the 2Gbps capacity of the two 1Gbps Steiner trees. The link stress of Datacast is 1.01. This means that Datacast achieves high bandwidth efficiency, since each packet only traverses each Steiner tree link 1.01 times on average. We compare Datacast with BitTorrent. Using BitTorrent, the receivers finish the downloading in 41-52s, and the link stress is 1.39. BitTorrent is 2.75 times slower than Datacast on average, while its link stress is 1.38 times larger.

3) *Failure handling*: To study the failure handling of Datacast, we manually tear down the slow link. Our Fabric Manager detects the link failure in 483ms, and then notifies all the Masters. The Master uses the signalling protocol proposed in Section III to deliver the signalling messages to all the receivers in 2.592ms. (As a comparison, using TCP to send the signalling messages to receivers in parallel takes 20.122ms.) Then the transmission continues.

## VIII. DISCUSSION

In this paper, we focus on Datacast for RGDD communication within a data center (intra-DC). We also study whether the Datacast protocol can be extended for inter data center (inter-DC) RGDD communication. The biggest challenge here is that the network latency for inter data center communication can be large, which will result in high duplicate data ratio. For example, our measurements show that the network latency between data centers located in east coast and west coast of the US is around 71ms. If we use the configuration in our simulation (i.e.,  $\delta = 5\text{Mbps}$ ,  $T = 1\text{ms}$  and  $R = 100\text{Mbps}$ ), the bound of duplicate data ratio will be as high as 78.3% based on Theorem 2.

In order to address this issue, we can first select representative nodes in each data center and use existing high speed TCP variants (e.g., CUBIC [18]) to deliver data from the source to these nodes, and then start Datacast to do RGDD within each data center. The detailed design and evaluation of this inter-DC approach will be our future work.

## IX. RELATED WORK

RGDD is an important traffic pattern, which has been studied for decades. Existing solutions can be classified into two categories.

**Reliable IP multicast.** The design space of reliable IP multicast has been nicely described in [12]. IP multicast has scalability issues for maintaining a large number of group states in the network. Adding reliability to IP multicast is also hard due to the ACK implosion problem [13].

We compare Datacast with two representative reliable multicast systems: pgm congestion control (pgmcc) [22] and Active Reliable Multicast (ARM) [26]. Pgmcc needs to explicitly track the slowest receiver for congestion control, and the congestion control protocol needs to be run between the sender and the slowest receiver. Datacast does not need to track which receiver is the slowest. This is because Datacast uses the duplicate interest packets as congestion signals, hence congestion control becomes the local action of the sender. ARM uses the active network concept and network devices also cache packet, but the cached packets are used only for re-transmission. Hence most likely the cached packets will not be used even once. Furthermore, re-transmitted packets are broadcasted along the whole sub-tree in ARM, whereas they are delivered only to the needed receivers in Datacast.

**End-host based overlay system.** End-host based overlay system overcomes the scalability issue by transmitting data among peers. No group states are needed in network devices, and reliability is easily achieved by directly using TCP. It is widely used in the Internet. However, end-host based overlay systems suffer from low bandwidth efficiency. For example, the worst-case link stress of SplitStream can be tens [3], and the average and worst-case link stresses of End System Multicast (ESM) [19] are 1.9 and 9, respectively.

Recently, in the work of Orchestra [11], Cornet is proposed, which is an optimized version of BitTorrent for DCNs. Different from the distributed manner of Cornet, Datacast is a centralized approach. Due to the fact that a data center network is built and managed by a single organization,

centralized designs become possible (e.g., software-defined networking [17]). Due to its centralized nature, Datacast is able to utilize multiple Steiner trees for data delivery, and achieve minimum finish time. Since the routing path from a receiver to data source is predetermined, high cache utilization is achieved. Furthermore, as we have demonstrated in the paper, the intermediate device only needs to maintain small cache per Steiner tree. All these benefits are hard, if not totally impossible, to be achieved by distributed approaches like Cornet.

## X. CONCLUSION

In this paper, we have presented the design, analysis, implementation and evaluation of Datacast for RGDD in data centers. Datacast first calculates multiple edge-disjoint Steiner trees with low time complexity, and then distributes data among them. In each Steiner tree, by leveraging in-network packet caching, Datacast uses a simple, but effective congestion control algorithm to achieve scalability and high bandwidth efficiency.

By building a fluid model, we show analytically that the congestion control algorithm uses small cache size for each group (e.g., 125KB), and results in few duplicate data transmissions (e.g., 1.19%). Our analytical results are verified by both simulations and experiments. We have implemented Datacast using the ServerSwitch platform. When we use Datacast to transmit 4GB data in our 1Gbps BCube(4, 1) testbed with two edge-disjoint Steiner trees, the link stress is only 1.01 and the finish time is 16.9s, which is close to the 16s lower bound.

## APPENDIX

To build the model, we first analyze under what condition a duplicate interest is received at the data source. Figure 10 shows a scenario with three caching switches between the source and the slowest receiver. We assume that these switches are “shared with” (i.e., also connected to) a number of fast receivers. From the figure, we can see that the caches that are farther to the slowest receiver will store newer data (shown in the shadow areas), while the ones that are closer to the slowest receiver will store older data. The reason is that data packets are propagated from the source to the slowest receiver. However, the interest is sent from the slowest receiver to the source. If the last shared switch (i.e., switch 3) does not have the corresponding data, others will not have it either. The last shared switch is therefore very critical to cache misses. We defined it as *the critical caching node*. When the critical caching node cannot serve an interest, the interest will be sent to the source as a duplicate interest. The critical caching node does not change over time for a given transmission tree, since it is determined by the structure of the transmission tree and the positions of slow and fast receivers, i.e., the last shared caching node of the slowest receiver and fast receivers.

After understanding when a duplicate interest is received at the source, we build a fluid model to analyze the performance of Datacast, based on the following assumptions: 1) The (desired<sup>4</sup>) rate of the slowest receiver,  $R$ , does not change

<sup>4</sup>Here “desired” means that the rate of the slowest receiver is not constrained by the sending rate of the data source.

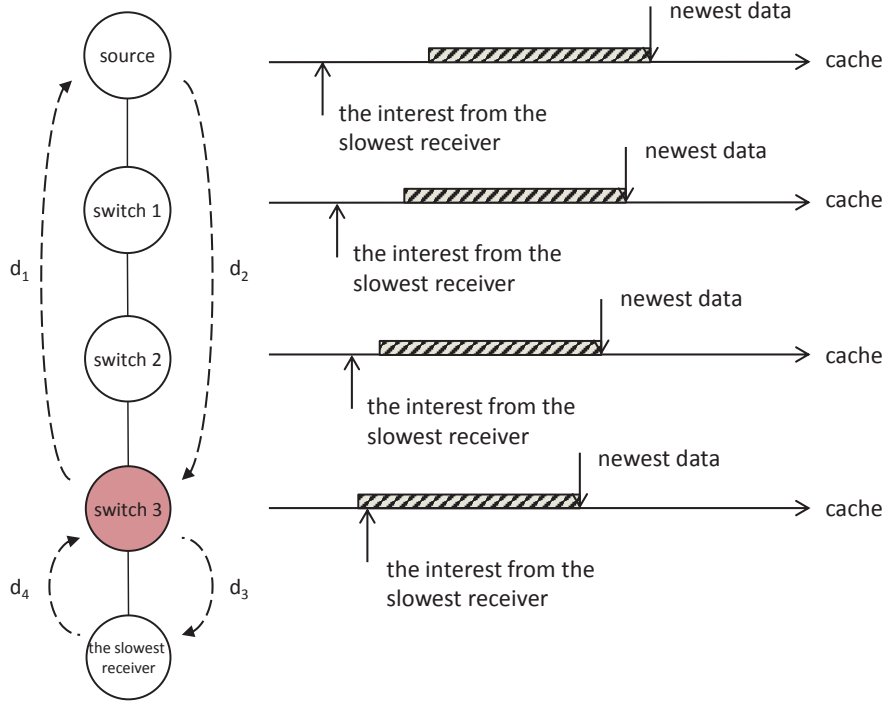


Fig. 10. The critical caching node.  $d_1$ ,  $d_2$ ,  $d_3$  and  $d_4$  are the latencies.

TABLE II  
NOTATIONS USED IN THE FLUID MODEL.

Notation	Meaning
$t$	The current time.
$x_s(t), x_r(t)$	The data sequence positions of the data source and the slowest receiver.
$R$	The rate of the slowest receivers.
$C$	The size of the cache (the content store).
$MTU$	The size of a full Datacast data packet.
$\delta, T$	The two parameters of Datacast congestion control, which are proposed in Section V-B.
$t_a$	The start time of state 0.
$t_b$	The end time of state 0, and the start time of state 1.
$t_c$	The end time of state 1.
$\Delta x(t)$	$x_s(t - d_1 - d_2) - x_r(t - d_1 + d_3)$

over time. 2) The credit number  $w$  is large enough to saturate the pipe. 3) The queue is large enough so that there is no packet drop due to the buffer overflow of a queue. Table II shows the notations that are used in the analysis. Our fluid model can be described by the following equations:

$$x_s''(t) = (1 - p(t)) \frac{\delta}{T} - p(t) \frac{x_s'(t)}{2} \frac{x_r'(t - d_1 - d_4)}{MTU} \quad (2)$$

$$x_r'(t) = \begin{cases} R & \text{if } x_r(t) < x_s(t - d_2 - d_3) \\ \max\{R, x_s'(t - d_2 - d_3)\} & \text{if } x_r(t) = x_s(t - d_2 - d_3) \end{cases} \quad (3)$$

$$p(t) = \mathbb{1}_{\{x_s(t - d_1 - d_2) - x_r(t - d_1 + d_3) > C + w \cdot MTU - (d_3 + d_4)R\}} \quad (4)$$

In this model, Equation (3) captures the slowest receiver's

(actual) rate. At time  $t$ , the slowest receiver wants data  $x_r(t)$ , and the newest data it can get from the data source is  $x_s(t - d_2 - d_3)$ . When  $x_r(t) < x_s(t - d_2 - d_3)$ , it means that there are packets in the queues between the source and the slowest receiver, so the slowest receiver's rate is  $R$ . When  $x_r(t) = x_s(t - d_2 - d_3)$ , the queues between the source and the slowest receiver are empty, so the slowest receiver is constrained by both the source's rate at time  $t - d_2 - d_3$  and  $R$ . Equation (4) is an indicator function.  $p(t) = 1$  when the data source receives a duplicate interest, otherwise  $p(t) = 0$ . If the data source receives a duplicate interest at time  $t$ , the interest will not be served by the critical caching node at time  $t - d_1$ . When the slowest receiver is retrieving data from the critical caching node, the data in the queues between the critical caching node and the slowest receiver are  $w \cdot MTU - (d_3 + d_4)R$ . At time  $t - d_1$ , the interest from the slowest receiver is retrieving  $x_r(t - d_1 + d_3) + w \cdot MTU - (d_3 + d_4)R$  from the critical caching node, while the newest data is  $x_s(t - d_1 - d_2)$ . So if the distance between them is larger than  $C$ ,  $p(t) = 1$ . Otherwise,  $p(t) = 0$ . Equation (2) models the rate control at the data source.  $\frac{\delta}{T}$  captures a constant rate increase  $\delta$  in every time period  $T$  if there is no duplicate interest. The second term is the rate decrease when duplicate interests are received (i.e.,  $p(t) = 1$ ). When  $p(t) = 1$ , the data source receives one duplicate interest from the slowest receiver in every time period  $\frac{MTU}{x_r'(t - d_1 - d_4)}$ , and decreases its sending rate by half. The decreasing rate therefore is  $\frac{x_s'(t)}{2} / \frac{MTU}{x_r'(t - d_1 - d_4)} = \frac{x_s'(t)}{2} \frac{x_r'(t - d_1 - d_4)}{MTU}$ .

We say the system is in **state 0** when  $p(t) = 0$ , in **state 1** when  $p(t) = 1$ . It is easy to see that the system will oscillate between the two states, since  $x_s''(t) > 0$  in state 0, and  $x_s''(t) < 0$  in state 1. We call it a **cycle** from the start of state

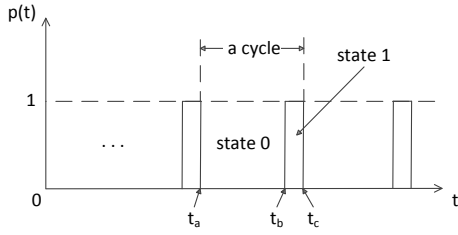


Fig. 11. An illustration of the state changes in Datacast.

0 to the end of state 1. Figure 11 gives us an illustration of the state changes in Datacast.

### Proof of Theorem 1:

*Proof:* We first prove that if Inequality (1) is true, the rate of the slowest receiver is not reduced, i.e.,  $x'_r(t) = R$ . To prove that, we first prove  $\Delta x(t) > 0$ . It is easy to see it holds in state 1, since  $\Delta x(t) > C + w \cdot MTU - (d_3 + d_4)R$  in state 1, and  $w$  is enough to saturate the pipe, i.e.,  $w \cdot MTU \geq R T T_m \cdot R$ , where  $R T T_m = d_1 + d_2 + d_3 + d_4$ . Next, we prove that it is also true in state 0.

In state 0, when  $t \in (t_a, t_a + d_1 + d_2]$ , we have

$$\begin{aligned} \Delta x(t) &> x_s(t_a - d_1 - d_2) - x_r(t_a - d_1 + d_3) - (t - t_a)R \\ &\geq C + (w \cdot MTU - (d_3 + d_4)R) - (d_1 + d_2)R \\ &= C + (w \cdot MTU - R T T_m \cdot R) > \frac{R^2 T}{2\delta} \end{aligned}$$

When  $t \in (t_a + d_1 + d_2, t_b)$ , we have

$$\begin{aligned} \Delta x(t) &= \int_{t_a + d_1 + d_2}^t \Delta x'(t) dt + \Delta x(t_a + d_1 + d_2) \\ &> \int_{t_a}^{t - d_1 - d_2} (x'_s(t) - R) dt + \frac{T}{2\delta} R^2 \\ &= \frac{\delta}{2T} (t - d_1 - d_2 - t_a)^2 \\ &\quad + (x'_s(t_a) - R)(t - d_1 - d_2 - t_a) + \frac{T}{2\delta} R^2 \\ &\geq -\frac{T}{2\delta} (R - x'_s(t_a))^2 + \frac{T}{2\delta} R^2 > 0 \end{aligned}$$

So  $\Delta x(t) > 0$  is also true in state 0. Putting  $\Delta x(t) > 0$  into (3), we get  $x'_r(t) = R$ , which means that the slowest receiver's rate is not slowed down. Actually, it can be further proved that the average sending rate of the data source will converge to  $R$  (which is omitted due to the space limitation), i.e., Datacast works at the full rate when Inequality (1) is satisfied. ■

Theorem 1 provides a sufficient condition to guarantee  $x'_r(t) = R$ . When  $C$  is not large enough,  $x'_r(t)$  can possibly be constrained by  $x'_s(t - d_2 - d_3)$  in state 0. However,  $x'_s(t - d_2 - d_3)$  will grow at a constant speed,  $\frac{\delta}{T}$ .  $x_s(t - d_2 - d_3)$  will soon be greater than  $x_r(t)$ , which means that the slowest receiver's rate is back to  $R$ . Even when  $C$  is not large enough, the system will experience graceful performance degradation instead of abrupt performance changes, as we have observed in the simulations and experiments.

### Proof of Theorem 2:

*Proof:* The duplicate ratio can be calculated as  $\frac{(t_c - t_b)R}{x_s(t_c) - x_s(t_a)}$ .  $(t_c - t_b)R$  is the amount of duplicate data that the slowest receiver requested in state 1, while  $x_s(t_c) - x_s(t_a)$  is the amount of new data sent from the source in the whole cycle. On entering the stable state, in each cycle, the data source and the slowest receiver move forward by the same distance, i.e.,  $x_s(t_c) - x_s(t_a) = x_r(t_c) - x_r(t_a)$ . Since  $x'_r(t) = R$ ,  $x_r(t_c) - x_r(t_a) = (t_c - t_a)R$ . The duplicate data ratio can be simplified as  $\frac{t_c - t_b}{t_c - t_a}$ . To calculate it, we first derive the links between the two states. At time  $t_a$ ,  $t_b$  and  $t_c$ , we have

$$\begin{aligned} x'_s(t_a) &= x'_s(t_c) \\ x'_s(t_b) &= x'_s(t_a) + \frac{\delta}{T}(t_b - t_a) \\ x'_s(t_c) &= x'_s(t_b) e^{-\frac{R}{2MTU}(t_c - t_b)} \end{aligned}$$

At time  $t_b$  and  $t_c$ , we have  $\Delta x(t_b) = \Delta x(t_c) = C + w \cdot MTU - (d_3 + d_4)R$ . So we have  $x_s(t_c - d_1 - d_2) - x_s(t_b - d_1 - d_2) = x_r(t_c - d_1 + d_3) - x_r(t_b - d_1 + d_3)$ . The right item is  $(t_c - t_b)R$ , since  $x'_r(t) = R$ . The left item can be divided into two parts,  $x_s(t_c - d_1 - d_2) - x_s(t_b)$  and  $x_s(t_b) - x_s(t_b - d_1 - d_2)$ . We calculate them separately, and then we get

$$\begin{aligned} (t_c - t_b)R &= x'_s(t_b)(d_1 + d_2) - \frac{\delta}{2T}(d_1 + d_2)^2 \\ &\quad + \frac{2MTU}{R} x'_s(t_b) (1 - e^{-\frac{R}{2MTU}(t_c - t_b - d_1 - d_2)}) \end{aligned} \quad (5)$$

From Equation (5), we can derive:

$$\begin{aligned} &(t_c - t_b)R \\ &\leq \frac{2MTU}{R} x'_s(t_a) (e^{\frac{R}{2MTU}(t_c - t_b)} - e^{\frac{R}{2MTU}(d_1 + d_2)}) \\ &\quad + x'_s(t_b)(d_1 + d_2) \\ &\leq \frac{2MTU}{R} x'_s(t_a) (e^{\frac{R}{2MTU}(t_c - t_b)} - 1) - \frac{R}{2MTU}(d_1 + d_2) \\ &\quad + x'_s(t_b)(d_1 + d_2) \\ &= (x'_s(t_b) - x'_s(t_a))(d_1 + d_2) + \frac{2MTU}{R} (x'_s(t_b) - x'_s(t_a)) \\ &= \frac{\delta}{T} \left( \frac{2MTU}{R} + d_1 + d_2 \right) (t_b - t_a) \end{aligned} \quad (6)$$

From (6), we can finally derive the bound of the duplicate data ratio

$$\frac{t_c - t_b}{t_c - t_a} \leq \frac{\frac{\delta}{T}}{\frac{\delta}{T} + \frac{R}{2MTU + d_1 + d_2}} \leq \frac{\frac{\delta}{T}}{\frac{\delta}{T} + \frac{R}{2MTU + R T T_m}}$$

the equal sign is true when  $R T T_m = 0$ . ■

### REFERENCES

- [1] M. Al-Fares, A. Loukissas, and A. Vahdat. A Scalable, Commodity Data Center Network Architecture. In *SIGCOMM*, 2008.
- [2] Ashok Anand, Archit Gupta, Aditya Akella, Srinivasan Seshan, and Scott Shenker. Packet Caches on Routers: The Implications of Universal Redundant Traffic Elimination. In *SIGCOMM*, 2008.
- [3] Miguel Castro, Peter Druschel, Anne-Marie Kermarrec, Animesh Nandi, Antony Rowstron, and Atul Singh. SplitStream: High-Bandwidth Multicast in Cooperative Environments. In *SOSP*, 2003.

- [4] Bram Cohen. Incentives Build Robustness in BitTorrent. In *Workshop on Economics of Peer-to-Peer Systems*, 2003.
- [5] J. Dean and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In *OSDI*, 2004.
- [6] J. Edmonds. Edge-disjoint branchings. In R. Rustin, editor, *Combinatorial Algorithms*, pages 91–96. Algorithmics Press, New York, 1972.
- [7] C. Guo et al. BCube: A High Performance, Server-centric Network Architecture for Modular Data Centers. In *SIGCOMM*, 2009.
- [8] Guohan Lu et al. ServerSwitch: A Programmable and High Performance Platform for Data Center Networks. In *NSDI*, 2011.
- [9] Hussam Abu-Libdeh et al. Symbiotic Routing in Future Data Centers. In *SIGCOMM*, 2010.
- [10] J. Cao et al. Datacast: A Scalable and Efficient Group Data Delivery Service for Data Centers. In *CoNEXT*, 2012.
- [11] M. Chowdhury et al. Managing Data Transfers in Computer Clusters with Orchestra. In *SIGCOMM*, 2011.
- [12] M. Handley et al. The Reliable Multicast Design Space for Bulk Data Transfer, Aug 2000. RFC2887.
- [13] Sally Floyd et al. A Reliable Multicast Framework for Light-weight Sessions and Application Level Framing. *IEEE Trans. Netw.*, Dec 1997.
- [14] Van Jacobson et al. Networking Named Content. In *CoNEXT*, 2009.
- [15] S. Ghemawat, H. Gobioff, and S. Leung. The Google File System. In *SOSP*, 2003.
- [16] R. L. Graham and L. R. Foulds. Unlikelihood That Minimal Phylogenies for a Realistic Biological Study Can Be Constructed in Reasonable Computational Time. *Mathematical Bioscience*, 1982.
- [17] K. Greene. Special reports 10 emerging technologies 2009. MIT Technology Review, 2009. <http://www.technologyreview.com/biotech/22120/>.
- [18] Sangtae Ha, Injong Rhee, and Lisong Xu. Cubic: a new tcp-friendly high-speed tcp variant. *SIGOPS Oper. Syst. Rev.*, 42(5):64–74, July 2008.
- [19] Yang hua Chu, Sanjay G. Rao, Srinivasan Seshan, and Hui Zhang. A Case for End System Multicast. *IEEE J. Sel. Areas Commun.*, Oct 2002.
- [20] M. Isard, M. Budi, and Y. Yu. Dryad: Distributed Data-Parallel Programs from Sequential Building Blocks. In *EuroSys*, 2007.
- [21] Force10 networks. Force10 s7000. [www.force10networks.com](http://www.force10networks.com).
- [22] Luigi Rizzo. pgmcc: a TCP-friendly Single Rate Multicast Congestion Control Scheme. In *SIGCOMM*, 2000.
- [23] Shyue-Ming Tang, Jinn-Shyong Yang, Yue-Li Wang, and Jou-Ming Chang. Independent Spanning Trees on Multidimensional Torus Networks. *IEEE Trans. Computers*, Jan 2010.
- [24] David L. Tennenhouse and David J. Wetherall. Towards an Active Network Architecture. *SIGCOMM CCR*, Apr 1996.
- [25] Po Tong and E. L. Lawler. A Fast Algorithm for Finding Edge-disjoint Branchings. *Information Processing Letters*, Aug 1983.
- [26] Li wei H. Lehman, Stephen J. Garland, and David L. Tennenhouse. Active Reliable Multicast. In *INFOCOM*, 1998.
- [27] J. Widmer and M. Handley. TCP-Friendly Multicast Congestion Control (TFMCC): Protocol Specification, August 2006. RFC 4654.



**Chuanxiong Guo** is a Principal Development Lead in the Windows Azure Group of Microsoft. Before that, he was a Senior Researcher in the Wireless and Networking Group of Microsoft Research Asia (MSRA). He received his Ph.D. degree from the Institute of Communications Engineering in Nanjing China. His research interests include networked systems design and analysis, network security, data center networking, networking support for operating systems. He is currently working on data center networking (DCN) and Cloud Computing.

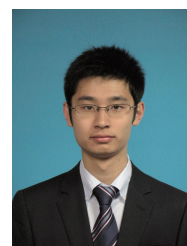


**Guohan Lu** received the B.S. degree in mechanical engineering, M.S. and PhD degrees in electronic engineering, both from Tsinghua University, China. He is currently an Associate Researcher in Microsoft Research Asia. His research interests are on network measurement and monitoring, network security and data center networks.



journals in the areas of

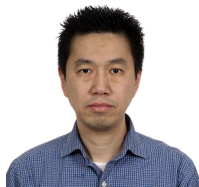
**Yongqiang Xiong** is now with Wireless and Networking Group at Microsoft Research Asia as a researcher. Dr. Xiong received his B.S., M.S., and Ph.D degrees from Tsinghua University, Beijing, China in 1996, 1998 and 2001, respectively, all in computer science. His research interests include data center and peer-to-peer networking, routing protocols for both MANETs and overlay networks, and network security. He has published over 40 papers, and served as TPC member or reviewers for the international key conferences and leading journals in the areas of wireless and networking. Dr. Xiong is member of IEEE.



**Yixin Zheng** received his BS degree from Tsinghua University, China, in 2012. He is currently an MS candidate in the Electronic Engineering Department at Tsinghua University. His research interests are in networking systems and data mining applications, with a focus on communication protocols and real-time data mining service in sensor networks.

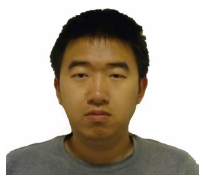


**Jiaxin Cao** received the bachelor degree and Ph.D. degree from University of Science and Technology of China in 2008 and 2013, respectively. During his Ph.D. program, he worked as a research intern in the W&N group of Microsoft Research Asia. His major research interests are data center networking and software defined networking. He is a Research Software Development Engineer in Microsoft now.



**Yongguang Zhang** is a Principal Researcher at Microsoft Research Asia, where he leads the Wireless & Networking research group. He received his Ph.D. in computer science from Purdue University in 1994. From 1994 to 2006 he was a research scientist at HRL Labs (Malibu, California) where he led various research efforts in internetworking techniques, system developments, and security mechanisms for satellite networks, ad-hoc networks, and 3G wireless systems, including as a co-PI in a DARPA Next Generation Internet project and as technical leads

in five other DARPA-funded wireless network research projects. From 2001 to 2003, he was also an adjunct assistant professor of Computer Science at the University of Texas at Austin. Yongguang Zhang's current interests include mobile systems and wireless networking. He has published over 50 technical papers and one book, including top conferences and journals in his fields (Sigcomm, NSDI, MobiCom, MobiSys, ToN, etc.). He recently won a string of Best Paper Awards (NSDI'09, CoNEXT'10, and NSDI'11) as well as 5 Best Demo Awards in a roll (MobiSys'07, SenSys'07, MobiSys'08, NSDI'09, and SIGCOMM'10). He is an Associate Editor for IEEE transactions on Mobile Computing, was a guest editor in an ACM MONET Journal, and has organized and chaired/co-chaired several international conferences, workshops, and an IETF working group. He was a General Co-Chair for ACM MobiCom'09.



**Yibo Zhu** is a second year PhD student in Department of Computer Science, University of California, Santa Barbara. He is working at Sand Lab co-advised by Prof. Ben Y. Zhao and Prof. Heather Zheng. Yibo's research interests include data center and wireless networks. He co-authored several papers published in top networking conferences such as ACM SIGCOMM'12, WWW'12 and CoNEXT'12. Yibo worked as an intern in Microsoft Research, Redmond in 2013 and Microsoft Research, Asia in 2011.



**Chen Chen** is a second-year Ph.D student at University of Pennsylvania. His research interest lies at clouding, software-defined network(SDN), security and formal verification. His current work involves virtualization in data center network(DCN) and formal verification on secure routing protocols.



**Ye Tian** received the bachelor's degree in electronic engineering and the master's degree in computer science from the University of Science and Technology of China (USTC), in July 2001 and 2004, respectively. He received the PhD degree from the Department of Computer Science and Engineering, The Chinese University of Hong Kong in December 2007. He is an associate professor at the School of Computer Science and Technology, USTC. He joined USTC in August 2008. His research interests include Internet and network measurement, peer-to-peer networks, online social networks, and multimedia networks. He is a member of the IEEE, and a senior member of the China Computer Federation (CCF). He is currently serving as an associate editor for Springer Frontiers of Computer Science.

peer networks, online social networks, and multimedia networks. He is a member of the IEEE, and a senior member of the China Computer Federation (CCF). He is currently serving as an associate editor for Springer Frontiers of Computer Science.