Optimizing All-to-All Data Transmission in WANs

Hao Tan¹⁰ and Wojciech Golab, *Member, IEEE*

Abstract—All-to-all data transmission is a typical data transmission pattern in both consensus protocols and blockchain systems. Developing an optimization scheme that provides high throughput and low latency data transmission can significantly benefit the performance of those systems. This paper investigates the problem of optimizing all-to-all data transmission in a wide area network (WAN) using overlay multicast. We prove that in a hose network model, using shallow tree overlays with height up to two is sufficient for all-to-all data transmission to achieve the optimal throughput allowed by the available network resources. Upon this foundation, we build ShallowForest, a data plane optimization for consensus protocols and blockchain systems. The goal of ShallowForest is to improve consensus protocols' resilience to skewed client load distribution. Experiments with skewed client load across replicas in the Amazon cloud demonstrate that ShallowForest can improve the commit throughput of the EPaxos consensus protocol by up to 100% with up to 60% reduction in commit latency.

Index Terms—Network overlays, multicast, consensus, state machine replication.

I. INTRODUCTION

EING highly available in the presence of machine **B** failures and network partitions is crucial to today's network services. State machine replication (SMR) [1] is a well-established technique to build fault-tolerant distributed systems. By having a group of replicated state machines collectively play the role of a server, the service can continue to operate when some of the machines fail. In SMR, each state machine executes an unbounded sequence of commands that update the current state. To make server state consistent across replicas, all replicated state machines must execute the same sequence of commands. To solve this challenging problem, replicated state machines communicate according to a specific consensus protocol to agree upon on a single sequence of commands to execute. Due to the asynchrony of the system, where messages can be delayed arbitrarily and processes can become arbitrarily slow, the replicas of a replicated state machine cannot always be in exactly the same state.

Traditionally, consensus protocols have been crucial building blocks in modern distributed systems for replicating

The authors are with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON N2L 3G1, Canada (e-mail: hao.tan@uwaterloo.ca; wgolab@uwaterloo.ca).

Digital Object Identifier 10.1109/TNSM.2021.3071025

important data and providing a strict ordering of updates to a small number of machines [2], [3]. Recently, blockchain [4], [5], [6], [7], [8] has become an emerging category of systems that require large scale consensus involving hundreds of nodes across different geographical regions connected by a wide-area network (WAN). Both consensus protocols and blockchain systems require multicasting data to a group of receivers. The following communication pattern dominates the normal operation of consensus protocols: upon receiving client requests, a replica broadcasts a message with commands to all other replicas and commits the request after receiving a certain number of responses. Such a communication pattern can be abstracted as an all-to-all data transmission, where each node in the cluster broadcasts an infinite stream of data to all other participating nodes.

Leader-centric consensus protocols like Paxos [9] and Raft [10] have a stable leader to handle all client requests. Since Internet protocol (IP) multicast is not generally available in a WAN environment, the stable leader in those protocols sends the data directly to all other replicas using multiple unicast transmissions. Assuming each site in the network is associated with an uplink capacity that limits the aggregated throughput of outgoing flows to other sites, this approach would render the leader as the bottleneck. As the number of replicas grows, each transmission will have less share of the available uplink capacity at the leader. Some protocols [11], [12], [13], [14] addressed this issue by handling data transmission using one or more ring overlays to maximize bandwidth utilization. However, a ring overlay is not an ideal option in a WAN environment due to the high latency of WAN links. Other protocols [15], [16], [17] alleviate the single leader bottleneck by distributing the load of data dissemination across all nodes. This strategy works best when the load is spread uniformly across all replicas. However, workloads in the real world can be highly skewed across different geo-areas and continuously changing over time. Prior measurements [18], [19], [20] point out that realworld workloads often exhibit diurnal variation, where the client load peaks around daytime and reaches the bottom at night time. For geo-distributed services deployed across multiple time zones, this type of user activity pattern will result in a skewed client load distribution across different regions. When each replica sends data directly to other replicas, it may yield sub-optimal throughput and lead to a load imbalance across replicas. Therefore, we argue that data dissemination should be handled in a more flexible way for consensus protocols to achieve high commit throughput and low commit latency.

This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 License. For more information, see https://creativecommons.org/licenses/by-nc-nd/4.0/

Manuscript received July 20, 2020; revised December 4, 2020 and March 16, 2021; accepted March 23, 2021. Date of publication April 5, 2021; date of current version September 9, 2021. Authors supported in part by a Ripple Graduate Fellowship, Ripple Faculty Fellowship, and by the Natural Sciences and Engineering Research Council (NSERC) of Canada. The associate editor coordinating the review of this article and approving it for publication was S. Kanhere. (*Corresponding author: Hao Tan.*)

<u> </u>	D
Symbol	Definition
G(V, E)	the network topology
s_i	the client data stream originates at $v_i \in V$
R_i	the incoming rate of s_i
R'_i	the residual rate of s_i to be partitioned
U_i^{\prime}	the reserved uplink capacity of v_i for sending out s_i once
U_i	the residual uplink capacity of v_i after reserving R_i for U
$s_{i,j}$	the j_{th} sub-stream of s_i
$r_{i,j}$	the rate of $s_{i,j}$

TABLE I TABLE OF NOTATIONS

In light of these challenges, we propose ShallowForest, an algorithm that optimizes data transmission for consensus protocols using overlay broadcast. ShallowForest computes data transmission overlays according to client load and available network capacity at each replica to make consensus protocols achieve high throughput and low latency. The central idea behind ShallowForest is inspired by overlay multicast protocols such as SplitStream and Bullet [21], [22], which partition the data stream at the sender and broadcast each stream partition with a potentially different tree overlay. However, the primary difference between ShallowForest and prior overlay multicast protocols is that ShallowForest only uses shallow tree overlays to reduce the network latency subject to the data transmission. This preference over shallow tree overlays is backed by the sufficiency of shallow tree overlays in achieving the optimal throughput of all-to-all data transmission in a hose network model.

This paper is an extended version of our prior work [23]. Besides the contributions mentioned earlier, this paper includes detailed proofs for the theorems in [23]. This paper also presents the result of benchmarking the throughput of data transmission between geo-areas on EC2, which further validates our choice of the hose network model. Last but not least, this paper introduces new theories along with visualizing examples, which extends the sufficiency of using shallow tree overlays from all-to-all data transmission to many-to-many data transmission.

II. PRELIMINARIES

A. Network Model

In this paper, we use the hose network model [24] to represent the WAN. The hose model abstracts the network as a set of sites connected by a core network with unlimited capacity. All sites can send and receive data from each other, bottlenecked only by the edge link capacity between each site and the core network. The network topology is represented as a directed complete graph G(V, E) with n vertices. Each vertex in V represents a geo-distributed site, and each edge in E represents the logical link between two sites. Due to a WAN's heterogeneous bandwidth availability, there are two functions $C_u: V \to \mathbb{R}^+$ and $C_d: V \to \mathbb{R}^+$, which respectively define the uplink and downlink capacity of the edge link between a site and the core network. The network latency between each pair of sites is denoted by $L: E \to \mathbb{R}^+$. The uplink and downlink capacity of a site are shared by all unicast data transmissions associated with that site. For instance,

TABLE II Throughput of One-to-One Data Transmission

		Virginia	Oregon	Tokyo	Ireland	
Virg	irginia –		753 Mbps	702 Mbps	749 Mbps	
Oreg	gon	761 Mbps	-	760 Mbps	705 Mbps	
Toky	/0	715 Mbps	738 Mbps	-	694 Mbps	
Irela	nd	763 Mbps	700 Mbps	687 Mbps	-	
800 700 - 600 - (stop) 500 - 400 - 300 -						
200-					-	
100-					_	
0-	Virginia	a Oreg	on Tc	okyo	reland	

Fig. 1. Aggregated throughput of one-to-many data transmission.

a sender directly multicasting to n receivers at the rate of R will consume nR of the sender's uplink capacity and R of each receiver's downlink capacity.

The hose network model is not only simple but also comply with the measurement we conducted on EC2. We benchmark the throughput of one-to-one as well as one-to-many data transmission on an Amazon EC2 cluster with four geodistributed VMs. We initiate 30 TCP flows from the sender VM to each receiver VM and keeping all flows running for 60 seconds. We repeat the test 5 times at different time of the day and record the average aggregated data transmission throughput. Table II demonstrates the throughput of one-to-one data transfer between each pair of regions. The throughput of one-to-many data transmission is presented in Figure 1, where the x-axis represents the location of the sender VM and each bar represents the aggregated throughput of the one-to-many data transmission at a specific location. Each segment with a particular colour scheme represents the throughput of the data transmission between the sender and a particular receiver. The graph demonstrates that, when there are multiple receivers, the aggregated throughput of data transmission to all receivers is capped around 750 Mbps. This rate is roughly equivalent to the per VM rate limit enforced by the public cloud provider. Due to the congestion control mechanism of TCP, the size of the sending window of a closer receiver grows faster than that of a distant receiver. Similar trends are also observed in many-to-one data transfers where a single receiver receives from multiple senders at the same time.

B. Terminology

Overlay: In a network G(V, E), an overlay O(V, E') is a spanning tree of G rooted at some site $v \in V$. It defines a broadcast transmission with site v as the sender. Each edge $(v_i, v_j) \in E'$ represents the transmission of v's data from v_i to v_j .

Client Data Stream: In a network G(V, E), a client data stream s is an infinite sequence of data bits from clients to be broadcast to all other sites in the network. The rate R_i of a client data stream s_i represents the incoming rate of client data at site $v_i \in V$. For instance, letting r be the number of incoming client requests per second at site v and letting b be the size of each request, the client data rate at site v is rb. We assume that a site's client data stream does not consume its downlink capacity as client requests arrive through the local area network. For all-to-all data transmission, each site $v_i \in V$ is associated with a client data stream s_i that must be received by all other sites.

Partitioning Scheme: Assume for simplicity of analysis that a client data stream can be split at arbitrary fine granularity. A partitioning scheme $P(s_i, k)$ of a client data stream s_i with rate R_i splits s_i into k streams $s_{i,1}, \ldots, s_{i,k}$ with rates $r_{i,1}, \ldots, r_{i,k}$ such that $\sum_{j=1}^{k} r_{i,j} \leq R_i$. Each split of the stream is referred to as a *sub-stream* of s_i .

Aggregated Throughput: In an all-to-all data transmission in a network G(V, E) with *n* sites, each site v_i broadcasts its data to all other sites at the rate R_i without violating the uplink and downlink capacity at any site. Then the aggregated throughput R_{total} of this all-to-all data transmission equals to $\sum_{i=1}^{n} R_i$.

C. Motivating Examples

In a network consisting of n geo-distributed sites v_1, \ldots, v_n , consider the case where each site has equal uplink and downlink capacities equal to B Mbps and the network latency between each pair of sites is L ms. Assume the incoming rate of client data is B Mbps at v_1 and zero elsewhere. v_1 needs to broadcast the data to all other sites.

Example 1: Have v_1 send the data directly to each site, each site can only receive the data at the rate of $\frac{B}{n-1}$ and the latency incurred by each receiver to receive each bit of data equals to L ms.

Example 2: Transfer data on a path joining all sites starting at v_1 such that data transmission only happens between adjacent sites. Using a path overlay, v_1 broadcasts data at the rate of *B* Mbps. However, the communication latency incurred by the last site on the path equals to (n - 1)L ms.

The above examples demonstrate the suboptimality of using a single overlay for data dissemination in terms of either throughput or latency. The example below shows how to achieve the optimal throughput without significantly compromising latency using multiple overlays.

Example 3: Equally partition the incoming data stream into n - 1 streams which are first sent to v_2, \ldots, v_n respectively. Upon receiving the data, each site then broadcasts the data to the remaining n - 2 sites. By using this approach, the network latency incurred by the data transmission becomes 2L ms while the transmission throughput remains *B* Mbps.

Although the above case considers a simplified network which involves only one sender, it raises a fundamental problem: given a network comprising of a set of nodes, with each node having a stream of data to broadcast to all other nodes, how can we maximize the aggregated broadcast





Fig. 2. Two types of base overlays in a cluster of four nodes.

throughput while minimizing the latency for each node's data to reach all other nodes?

III. SHALLOW TREE OVERLAYS SUFFICE FOR All-to-All Data Transmission

For all-to-all data transmission in a network G = (V, E) with *n* sites, each site is associated with a client data stream s_i with the incoming rate R_i that must be received by all other sites. Ideally, the system should broadcast each client data stream at its incoming rate. However, such is not always the case due to the heterogeneity of available bandwidth and non-uniformly distributed client load across all sites. This section characterizes the minimum requirement for a set of client data streams to be broadcast at their incoming rates and proves that using shallow tree overlays with height up to two is sufficient for achieving the optimal broadcast throughput.

Definition 1: Client data streams s_1, \ldots, s_n with rates R_1, \ldots, R_n are said to be **sustainable** if the following four conditions are all met:

1)
$$\forall v_i \in V, R_i \leq C_u(v_i)$$

2)
$$\forall v_i \in V, \sum_{j \neq i} R_j \leq C_d(v_i)$$

3)
$$(n-1)\sum_{i=1}^{n} R_i \leq \sum_{i=1}^{n} C_u(v_i)$$

4) $(n-1)\sum_{i=1}^{n} R_i < \sum_{i=1}^{n} C_d(v_i)$

4)
$$(n-1) \sum_{i=1}^{n} R_i \leq \sum_{i=1}^{n} C_d(v_i)$$

Intuitively, being sustainable is the minimum requirement for a set of client data streams to be broadcast at their incoming rates. Condition (1) ensures that each site has enough uplink capacity to send out its data at least once to other nodes. As each site has to receive from all other peers, condition (2) ensures that the aggregated rate of incoming streams does not exceed a site's downlink capacity. Condition (3) derives from the fact that the client data at each site must be sent at least n - 1 times. Similarly, the fact that client data at each site has to be received n - 1 times lead to condition (4). Note that, condition (4) is the direct result of condition (2) by summing over all possible *i*. If any of the above conditions are violated, the R_{total} will be less than $\sum_{i=1}^{n} R_i$.

Definition 2 (Base Overlay): refers to the following types of overlays: 1-level tree, or 2-level tree with exactly one non-leaf node (excluding the root). Figure 2 demonstrates these two base overlays in a network with four sites.

Definition 3: Let S be a set of sub-streams and \mathbb{O} be the set of all overlays in a network G(V, E). A one-to-one mapping $f : S \to \mathbb{O}$ is said to be **sustainable** if each sub-stream $s \in S$ can be transmitted at its rate using f(s) without violating downlink and uplink capacity constraints at any site.

Algorithm	1:	Computing	Sub-Stream	Rate
¹ Hgoi Iumm	. .	computing	Sub Subuli	ruic

s **Input** : G(V, E) $C_u: V \to \mathbb{R}+$ $R_i \ldots R_n$ **Output:** $r_{i,j}, 1 \leq i, j \leq n$ 1 $r_{i,j} \coloneqq 0, 1 \le i, j \le n; //$ initialize sub-stream rates 2 $U_i \coloneqq C_u(v_i) - R_i, 1 \le i \le n; //$ initialize each residual uplink capacity 3 for i := 1 to n do // go through all source sites $R'_i \coloneqq R_i;$ 4 for $j \coloneqq 1$ to n do // loop through all 5 sub-streams $\begin{array}{l} \text{if } (n-2)R'_i > U_j \text{ then} \\ \big| \quad r_{i,j} \coloneqq \frac{U_j}{n-2}; \end{array}$ 6 7 else 8 $\begin{array}{l} | & r_{i,j} \coloneqq R'_i; \\ U_j \coloneqq U_j - (n-2)r_{i,j}; \\ R'_i \coloneqq R'_i - r_{i,j}; \\ \mathbf{if} \ R'_i = 0 \ \mathbf{then} \end{array}$ 9 10 11 12 break; 13 14 return $r_{1,1} \ldots r_{n,n}$;

Theorem 1: For client data streams s_1, \ldots, s_n with sustainable rates R_1, \ldots, R_n , there exists a partitioning scheme for each client data stream and a sustainable mapping from sub-streams to overlays such that:

- 1) Each sub-stream's overlay is a base overlay.
- 2) The resulting aggregated throughput equals to $\sum_{i=1}^{n} R_i$.

A. Proof of Theorem 1

The general idea for proving Theorem 1 is to construct a partitioning scheme for each client data stream and associate each sub-stream with a base overlay such that the resulting data transmission will not violate downlink and uplink capacity constraints at any site.

1) Constructing Sub-Stream Overlays: Each client stream s_i will be split into *n* sub-streams $s_{i,1}, \ldots, s_{i,n}$ with rates $r_{i,1}, \ldots, r_{i,n}$. The data of the special sub-stream $s_{i,i}$ is sent directly from v_i to all the remaining sites. The data of substream $s_{i,j}$ for $i \neq j$ is sent from v_i to v_j first, and then v_j will broadcast the data to the rest of the sites. All overlays defined previously are base overlays.

2) Computing Sub-Stream Rates: Algorithm 1 computes the rate of each sub-stream defined in the previous section. It does not aim to compute the latency-optimal partitioning scheme, which favors 1-level tree overlays. The purpose of Algorithm 1 is to construct a partitioning scheme and an overlay mapping for proving Theorem 1. Table I summarizes the notations used in Algorithm 1.

For each node v_i , its uplink capacity is divided into two parts: $U'_i = R_i$ and $U_i = C_u(v_i) - R_i$. U'_i represents the reserved uplink capacity for v_i to send out all its data at least once and U_i is the residual uplink capacity such that U_i +

 $U'_i = C_u(v_i)$. The algorithm iterates over all site pairs in $\{(i,j)|1 \leq i \leq n, 1 \leq j \leq n\}$ in lexicographical order to compute sub-stream rates. Using such an order is just for the clarity of the proof and has no impact on the correctness of the output of Algorithm 1.

For the iteration when $r_{i,j}$ is computed, the algorithm greedily allocates as much of U_j as possible to $r_{i,j}$ until either U_j is exhausted or the aggregated sub-stream rate reaches R_i . According to the overlay trees defined in the previous section, sending $s_{i,j}$ consumes $r_{i,j}$ of U'_i and $(n-2)r_{i,j}$ of U_i . This rule also applies to the case i = j, where sending $s_{i,i}$ consumes $r_{i,i}$ of U'_i (sending $s_{i,i}$ from v_i to v_j , $v_j \in G(V) \setminus \{v_i\}$ and $(n-2)r_{i,i}$ of U_i (multicasting $s_{i,i}$ to each site in $G(V) \setminus \{v_i, v_j\}$). As a result, sending $s_{i,i}$ consumes in total $(n-1)r_{i,i}$ of $C_u(v_i)$.

3) Correctness Criteria: The output of Algorithm 1 is a set of sub-stream rates $r_{1,1}, \ldots, r_{n,n}$. A correct output satisfies the following three criteria for all $1 \le i \le n$:

- Valid Partition Constraint: The aggregated rate of all sub-streams of a client data stream is equal to that client data stream's rate, which is equivalent to $\sum_{j=1}^{n} r_{i,j} =$ R_i .
- Uplink Capacity Constraint: The aggregated rate of all sub-streams sent by v_i is less than or equal to $C_u(v_i)$.
- Downlink Capacity Constraint: The aggregated rate of all sub-streams received by v_i is less than or equal to $C_d(v_i).$

4) Correctness of Algorithm 1: Let $U_i[\alpha]$ represent the value of U_i at the start of iteration α of the outer loop.

Proposition 1: For all α such that $1 \leq \alpha \leq n$, if $\sum_{i=1}^{n} U_i[\alpha] \ge (n-2)R_{\alpha}$, then $\sum_{i=1}^{n} r_{\alpha,i} = R_{\alpha}$ at the end of iteration α of outer loop.

Proposition 2: For all α such that $1 \leq \alpha \leq n$, $\sum_{i=1}^{n} U_i[\alpha] \ge (n-2) \sum_{i=\alpha}^{n} R_i.$

Proposition 3: The output of Algorithm 1 satisfies the Valid Partition Constraint.

Proof: This proposition holds as the direct outcome of Proposition 1 and Proposition 2. According to Proposition 2:

$$\sum_{i=1}^{n} U_i[\alpha] \ge (n-2) \sum_{i=\alpha}^{n} R_i \ge (n-2)R_{\alpha}, \forall \alpha : 1 \le \alpha \le n$$
(1)

and Proposition 1:

$$\sum_{i=1}^{n} U_i[\alpha] \ge (n-2)R_\alpha \implies \sum_{i=1}^{n} r_{\alpha,i} = R_\alpha \qquad (2)$$

We have:

$$\sum_{i=1}^{n} r_{\alpha,i} = R_{\alpha}, \forall \alpha : 1 \le \alpha \le n$$
(3)

Equation (3) implies the Valid Partition Constraint. Lemma 1: $U_i[\alpha] \ge 0$ for all i, α such that $1 \le i, \alpha \le n$. *Proof:* This lemma is proved by induction on α .

Base case $\alpha = 1$: By line 2, we have: $U_i[1] = C_u(v_i) - R_i$ for all *i* such that $1 \le i \le n$. According to condition (1) of sustainable rates, $U_i[1] \ge 0$ holds for all *i* such that $1 \le i \le n$. **Induction step:** For an arbitrary number α such that $1 < \alpha \le n$, assume $U_i[\alpha - 1] \ge 0$. Line 6 and line 7 guarantee $r_{\alpha,i} \le \frac{U_i[\alpha - 1]}{n-2}$ and $U_i[\alpha] = U_i[\alpha - 1] - (n-2)r_{\alpha,i}$ according to line 10. As a result, $U_i[\alpha] \ge 0$.

Proposition 4: The output of Algorithm 1 satisfies the Uplink Capacity Constraint.

Proof: From Lemma 1, we have $U_i \ge 0$ throughout the execution of Algorithm 1 for all *i* such that $1 \le i \le n$. According to the overlay defined in the previous Section III-A1, sending $s_{i,j}$ consumes $r_{i,j}$ of U'_i and U'_i is consumed only by sending v_i 's sub-streams. By Proposition 3, $\sum_{j=1}^n r_{i,j} = R_i$ for all *i* such that $1 \le i \le n$. Since we reserve R_i for U'_i , sending all of v_i 's sub-streams will consume exactly the amount of its reserved uplink capacity. Because $C_u(v_i) = U_i + U'_i$, no uplink capacity constraint is violated.

Proposition 5: The output of Algorithm 1 satisfies the Downlink Capacity Constraint.

Proof: Since every sub-stream is broadcast using a tree overlay, each site receives all other site's data exactly once and receives the data at the same rate as the source is sending. According to the condition (2) of sustainable rates, there is also no violation of downlink capacity constraint.

B. Throughput Improvement

Section II-C briefly discussed the benefit of using shallow tree overlays when there is a single sender broadcasting data. In this section, we extend the analysis to general cases that involve multiple senders. According to Theorem 1, any sustainable rates in all-to-all data transmission are achievable using shallow tree overlays. Comparing to directly broadcasting to other sites, the equation (4) captures the throughput improvement achieved by the shallow tree overlay approach for sustainable client data rates R_1, \ldots, R_n :

$$\frac{\sum_{i=1}^{n} R_i}{\sum_{i=1}^{n} \min\left(R_i, \frac{C_u(v_i)}{n-1}\right)} \tag{4}$$

The numerator is the aggregated throughput achieved by the shallow tree overlay approach while the denominator is the aggregated throughput achieved by direct broadcasting. If some site v_i does not have sufficient uplink capacity to broadcast its data directly at the rate R_i , using the shallow tree overlay approach results in higher aggregated throughput. Such a situation arises when there is a mismatch between the distribution of client load and the distribution of available uplink capacity. For instance, the client load at some site is much higher than the client load at other sites, or some site has limited uplink capacity compared to other sites.

C. Extending the Applicability of Shallow Tree Overlays

Theorem 1 implies that using shallow tree overlays with height up to two is sufficient for all-to-all data transmission to achieve the optimal throughput. However, in an arbitrary many-to-many data transmission, a site can multicast its data to a subset of sites in the network. Are shallow tree overlays still sufficient to achieve the optimal throughput for many-tomany data transmission? This section provides an affirmative answer to this question. In the context of many-to-many data transmission, we made the following adjustment to some terms used in previous sections:

Definition 4: In a network G(V, E), a **client data stream** s is an infinite sequence of data bits from clients to be multicast to a set of sites V', where $V' \subseteq V$.

Definition 5: In a network G(V, E), an overlay O(V', E') is a tree rooted at some site $v \in V'$ such that $V' \subseteq V$ and $E' \subseteq E$.

Lemma 2: Let s be a client data stream transmitted using an arbitrary tree overlay T at rate R, there exists a partitioning scheme for s and a sustainable mapping from sub-streams to overlays such that:

- 1) Each sub-stream's overlay is a base overlay consisting of all nodes of *T*.
- 2) The sum of all sub-stream rates equals to R.

Proof: Let T consist of k nodes $\{v_1, v_2, \ldots, v_k\}$ with v_1 being the root and V' be the set of all non-leaf nodes in T. Note that, there can be more than k nodes in the entire network. We prove the theorem by constructing a partitioning scheme and a valid sustainable mapping:

- For each non-leaf node v_i ∈ V', map a sub-stream with rate bi_iR/k-2 to O_i, where O_i is a 2-level tree base overlay rooted at v₁ with v_i being the other internal node, and b_i is the number of subtrees of v_i in T.
- 2) Map a sub-stream with rate $R \sum_{v_i \in V'} \frac{b_i R}{k-2}$ to O_1 , where O_1 is a 1-level tree base overlay rooted at v_1 .

By summing up the sub-stream rates defined above, each non-root node in T receives data at the rate:

$$\sum_{v_i \in V'} \frac{b_i R}{k-2} + \left(R - \sum_{v_i \in V'} \frac{b_i R}{k-2}\right) = R \tag{5}$$

It remains to prove that the mapping is sustainable. As each node receives data at the same rate R, there is no violation of downlink capacity constraint. There is also no violation of uplink capacity constraint at each node in V'. The resulting data transmission consumes the same amount of uplink capacity at each node in V' as the data transmission over T does, which equals to $b_i R$. For v_1 , it sends data to k - 1 nodes at rate $R - \sum_{v_i \in V'} \frac{b_i R}{k-2}$ plus sending data to each node in V' at rate $\frac{b_i R}{k-2}$. Therefore, the uplink capacity of v_1 consumed by data transmission is calculated as follows:

$$(k-1)\left(R - \sum_{v_i \in V'} \frac{b_i R}{k-2}\right) + \sum_{v_i \in V'} \frac{b_i R}{k-2}$$
(6)

Since T is a tree, $\sum_{i=1}^{k} b_i$ equals to k-1 and $\sum_{v_i \in V'} b_i = k-1-b_1$. Therefore, equation (6) can be rewritten as:

=

$$(k-1)R - (k-2)\sum_{v_i \in V'} \frac{b_i R}{k-2}$$
 (7)

$$= (k-1)R - R\sum_{v_i \in V'} b_i \tag{8}$$

$$= (k-1)R - (k-1-b_1)R$$
(9)

$$b_1 R \tag{10}$$



Fig. 3. The visualization of Lemma 2.

The resulting data transmission also consumes the same amount of v_1 's uplink capacity as the original data transmission over T does. As a result, the mapping we construct is sustainable.

Figure 3 demonstrates two examples of transforming the data transmission using an arbitrary tree overlay into data transmission using multiple base overlays. The first example partitions the data stream s transmitted using T into three sub-streams s_1 , s_2 , and s_3 transmitted using O_1 , O_2 and O_3 respectively. Edge weights represent the rate of data stream transmitted using each overlay. s has rate 5, s_1 has rate 1 and both s_2 and S_3 has rate 2. The resulting data transmission using O_1 , O_2 and O_3 sends data to each non-root node of T at rate 5, which is equivalent to the rate of s. Also, the resulting data transmission consumes the same amount of v_1 's, v_2 's and v_3 's uplink capacity (equals to 10) as the original data transmission over T does. Likewise, the data transmission using O_1' , O_2' and O_3' .

Theorem 2: Using shallow tree overlays with height up to two is sufficient for achieving the optimal throughput for many-to-many data transmission.

Proof: This theorem can be proved by contradiction. Assume a data stream s with the rate R must be transmitted by a tree overlay T to achieve the optimal throughput and the height of T is greater than two. According to Lemma 2, s can be partitioned into a set of sub-streams transmitted using base overlays. The resulting data transmission sends data to every node in T (except for the root) at the rate R and consumes the same amount of uplink capacity at each node that sends data in T. Therefore, the original data transmission using T can be replaced by the data transmission using base overlays.

Theorem 2 extends the applicability of shallow tree overlays to many-to-many data transmission. Note that, all-to-all data transmission is a special form of many-to-many data transmission. Therefore, Theorem 2 further validates the correctness of Theorem 1. However, Theorem 2 is not stronger than Theorem 1 since it does not answer what is the optimal achievable throughput given a hose network topology.

IV. SHALLOWFOREST ALGORITHM

ShallowForest is a two-phase algorithm that optimizes allto-all data transmission in a WAN environment for consensus protocols and blockchain systems. We assume that network capacity is the critical performance-limiting resource for such systems in a WAN environment. The primary optimization goal of ShallowForest is to maximize the aggregated data transmission throughput while the secondary goal is to minimize the network latency subject to the data transmission rate. As a result, the first phase computes the maximum achievable data transmission throughput constrained by the network capacity and client load across all sites. In the second phase, ShallowForest computes the optimal way to partition each client data stream and associates each sub-stream with an overlay such that the resulting collection of overlays achieves the optimal throughput obtained from the first phase. As network delay is not negligible in a WAN environment, the second phase also minimizes the aggregated latency weight of the resulting overlays. In the sections below, we describe the ShallowForest algorithm in detail.

A. Throughput-Optimal Broadcast Rate

During the first phase, ShallowForest computes the maximum achievable aggregated broadcast throughput R_{total} . We first demonstrate under what conditions R_{total} becomes achievable in a network with limited resources.

According to Theorem 1, if client data streams already have sustainable rates, $R_{total} = \sum_{i=1}^{n} R_i$. Otherwise, it is impossible to broadcast all client data streams at their incoming rates. In such a case, R_{total} is computed using the following LP formulation:

free variables:
$$R'_i \quad 1 \le i \le n$$
 (11)

maximize:
$$R_{total} = \sum_{i=1}^{N} R'_i$$
 (12)

subject to: $R'_i \leq \min(C_u(v_i), R_i) \quad \forall i : 1 \leq i \leq n$

$$\sum_{j \neq i} R'_j \le C_d(v_i) \quad \forall i : 1 \le i \le n$$
 (14)

(13)

$$(n-1)\sum_{i=1}^{n} R'_{i} \le \sum_{i=1}^{n} C_{u}(v_{i})$$
(15)

In the above LP formulation, variables R'_1, \ldots, R'_n represent some set of sustainable client data rates. Equation (12) defines the optimization objective. Constraints 13-15 ensure that R'_1, \ldots, R'_n meet the first three conditions of being sustainable. After computing R_{total} , ShallowForest proceeds to the next phase.

B. Latency-Optimal Overlays

The goal of the second phase is to compute a partitioning scheme for each client data stream that achieves the aggregated broadcast throughput R_{total} , and construct overlays for all substreams such that the network latency incurred by the data transmission is minimized.

Definition 6: The latency weight: l(O) of an overlay O(V, E) rooted at $v \in V$ is the aggregated network latency incurred by all receivers to receive v's data. Let $P_i \subseteq E$ be the path in O from v to some $v_i \in V$, we have l(O) = $\sum_{v_i \in V} \sum_{e \in P_i} L(e).$

Problem Statement: Given a network G(V, E) with C_u : $V \to \mathbb{R}^+$, C_d : $V \to \mathbb{R}^+$, L : $E \to \mathbb{R}^+$, client data streams $s_1 \ldots s_n$ with rates $R_1 \ldots R_n$, and a target aggregated broadcast throughput R_{total} , find a partitioning scheme $P(s_i, n_i) = \{s_{i,1}, \dots, s_{i,n_i}\}$ and the corresponding overlay of each sub-stream $O_{i,1}, \ldots, O_{i,n_i}$ such that:

- 1) Each sub-stream $s_{i,j}$ can be broadcast at its rate $r_{i,j}$ without violating downlink and uplink capacity constraints at any site.
- 2) $\sum_{i=1}^{n} \sum_{j=1}^{n_i} r_{i,j} = R_{total}.$ 3) $\sum_{i=1}^{n} \sum_{j=1}^{n_i} l(O_{i,j})r_{i,j}$ is minimized.

The term $l(O_{i,j})r_{i,j}$ is the product of sub-stream rate and its overlay's latency weight, which represents the network latency subject to the data transmission over $O_{i,j}$ at rate $r_{i,j}$. Minimizing the sum of this term over all overlays will promote transmitting more data on overlays with low network latency to reduce the average network latency incurred by the entire all-to-all data transmission.

1) Choosing Overlay Candidates: It is impractical to consider all possible types of overlays due to their sheer number. However, overlay candidates used by the second phase have a critical impact on the resulting aggregated latency weight. Selected overlay candidates must not impair achieving the R_{total} computed by the first phase, and are expected to be as shallow as possible to minimize the overhead of network latency. Based on Theorem 1, it is sufficient to only consider base overlays to achieve the optimal aggregated broadcast throughput.

2) LP Formulation: We first set up a partitioning scheme for each client stream, and pair each sub-stream with an overlay. Since there are n overlay candidates for each site, each client data stream s_i will be split into n sub-streams $s_{i,1}, \ldots, s_{i,n}$ with rates $r_{i,1}, \ldots, r_{i,n}$. Those sub-stream rates are the variables to be optimized. We assign an overlay $O_{i,j} = (V, E_{i,j})$ to a sub-stream $s_{i,j}$ such that the data transmission is handled in the following way: (1) the data of $s_{i,i}$ is sent directly from v_i to all the remaining sites; (2) the data of $s_{i,j}$ for $i \neq j$ is sent from v_i to v_j first, and then v_j broadcasts the data to the rest of the sites. The resulting LP formulation is as follows:

free variables:
$$r_{i,j} \quad \forall i, j : 1 \le i, j \le n$$
 (16)

minimize:
$$\sum_{i=1}^{n} \sum_{j=1}^{n} l(O_{i,j}) r_{i,j}$$
 (17)

subject to:
$$U_i \le C_u(v_i) \quad \forall i : 1 \le i \le n$$
 (18)

$$\sum_{j \neq i} r_{j,i} \le C_d(v_i) \quad \forall i : 1 \le i \le n \quad (19)$$

$$\sum_{i=1}^{n} r_{i,j} \le R_i \quad \forall i : 1 \le i \le n$$
(20)

$$\underset{n}{r_{i,j} \ge 0} \quad \forall i,j : 1 \le i,j \le n$$
 (21)

$$\sum_{i=1}\sum_{j=1}r_{i,j} = R_{total} \tag{22}$$

For all *i*, *j* such that $1 \le i, j \le n$:

$$U_{i} = (n-1)r_{i,i} + (n-2)\sum_{j \neq i} r_{j,i} + \sum_{j \neq i} r_{i,j}$$
(23)

Constraint (23) represents the amount of v_i 's uplink capacity consumed by the data transmission with respect to the overlay setup. Constraint (18) characterizes the uplink capacity constraint at a specific site: the aggregated rates of data sent out of a site should be less than or equal to that site's uplink capacity. Constraint (19) characterizes the downlink capacity constraint at a specific site: the aggregated rates of data received by a site should be less than or equal to that site's downlink capacity. Constraint (20) enforces the sum of substream rates being less than or equal to the rate of the original stream. Constraint (21) enforces all sub-stream rates to be nonnegative. Constraint (22) enforces the sum of all sub-stream rates equals to be R_{total} , which is computed in the first phase (see Section IV-A).

C. The Complexity of ShallowForest

In a network with n sites. The LP formulations above involve $O(n^2)$ variables and constraints. Using Vaidya algorithm [25], the optimization problem can be solved with time complexity $O(n^5)$. Note that the actual throughput achieved



Fig. 4. The software architecture of APaxos. Red lines represent protocol messages and blue lines represent client operations.

by the overlays computed in this phase is approximate to the optimal throughput for two reasons: (1) The LP formulations relax the integrality constraint on the rate of each sub-stream; in reality, you cannot split a data stream at a granularity finer than one bit. (2) A software LP solver may introduce rounding error.

V. AMOEBA PAXOS: WORKLOAD-AWARE CONSENSUS

EPaxos [16] is a state-of-the-art decentralized consensus protocol that performs favorably in a WAN environment. However, its workload-agnostic approach to handle data transmission will lead to sub-optimal performance when dealing with skewed load across replicas. To make EPaxos workloadaware, we build Amoeba Paxos (APaxos) on top of the publicly available EPaxos implementation [26] by applying ShallowForest to the data transmission.

A. Overview

Figure 4 depicts the software architecture. There are three major components in APaxos: the ordering plane, the data plane and a centralized controller. The ordering plane receives incoming client requests and orders them using the original EPaxos protocol. Instead of broadcasting messages with client operations directly to other replicas, the ordering plane replaces actual client operations in protocol messages with client operation IDs and offloads the job of broadcasting client operations to a co-located data plane thread.

The data plane broadcasts client operations with specific overlays according to its overlay configuration updated by the controller. The overlay configuration determines how much data to transmit using a specific overlay. The data plane also transmits the client operations received from other replicas based on the overlay information encapsulated in the received data. Besides handling data transmission, the data plane also buffers the received client operations and reassembles them into protocol messages required by the ordering plane.

The controller applies ShallowForest to compute the optimal partitioning scheme and overlays for each site and updates each site's overlay configuration through RPC calls. In the prototype implementation, we hard-code the client data rates and available network resources in the controller.

B. The Ordering Plane

There are three types of messages in EPaxos that enclose client operations: PreAccept, TryPreAccept and PrepareReply. The ordering plane replaces client operations in PreAccept messages with client operation IDs. The resulting message is referred to as PreAcceptLight to distinguish it from the original PreAccept message. The ordering plane only separates client operations from PreAccept messages because broadcasting PreAccept messages consumes the greatest amount of bandwidth in normal operation while the latter two messages are only involved in the EPaxos's recovery process. The ordering plane broadcasts PreAcceptLight messages and handles the protocol messages from other replicas in the same way as an EPaxos replica does.

C. The Data Plane

This section presents salient details of the data plane.

1) Overcoming the Per-Flow Rate Limit: To overcome the per-flow rate limit enforced by public cloud providers, APaxos sets up multiple TCP connections between each pair of replicas located in different areas. For a specific recipient, the data plane picks the TCP connection from the pool in a round-robin fashion and transmits one client operation using a selected TCP connection in a separate thread. The purpose of letting each TCP connection have equal chances to transmit client operations is to make each TCP connection have a similar congestion control window size. With a high incoming rate of client operations, there could be multiple TCP connections concurrently sending client operations to the same recipient.

2) Overlay Configuration: The data plane thread running on site v sends client operations based on a local overlay configuration overlay_config, an array with the same size as the number of overlay candidates. Each entry overlay_config[i] is the amount of data out of a configurable window size w KB that should be broadcast using the *i*th overlay. That is, among w KB of data broadcast by v, overlay_config[i] KB of data should be broadcast using the *i*th overlay. The *i*th entry of the overlay configuration is also referred to as the quota of the *i*th overlay. In our implementation, w is set to 200KB to achieve the best performance.

3) Overlay Information: Another advantage of using only base overlay candidates is minimizing the overhead of overlay information in the data transmitted by the data plane process. To broadcast a client operation γ , the data plane process simply piggybacks relay bit to the original message based on its transmission overlay. The bit is set to 0 for a 1-level tree overlay and 1 otherwise. When a data plane process receives data from other replicas, it checks the piggybacked relay bit. If the bit equals 1, the data plane process will broadcast the message to the remaining replicas with the relay bit set to 0. For overlay management, the data plane process also append the message

Algorithm 2: Overlay Management						
1 Var $view = 0$: view number						
2 Var $i = 0$: configuration number						
3 Var \mathcal{R} : set of client data rates						
4 Var C : set of overlay configurations						
5 Procedure UpdateOverlays()						
6 for $v \in G(V)$ do						
7 $\mathcal{R}[v] \leftarrow \mathbf{RequestClientDataRate}(v)$						
8 if any call fails then						
9 return Abort						
10 $\mathcal{C} \leftarrow \mathbf{ShallowForest}(\mathcal{R})$						
11 $ver \leftarrow (view, i++)$						
12 for $v \in G(V)$ do						
13 ApplyConfig $(v, C[v], ver)$						
14 return OK						
15 Procedure <i>HandleViewChange</i> ($view', A, D$)						
16 if $view' > view$ then						
17 $view \leftarrow view', i \leftarrow 0$						
18 UpdateState $(\mathcal{A}, \mathcal{D})$						
19 UpdateOverlays()						
20 return OK						

with a field *ver*—the version of current overlay configuration, which will be described in Section V-D2

4) Assemble Ordering Plane Messages: Upon receiving a PreAcceptLight message, the data plane assembles a PreAccept message by retrieving all client operations referenced by the PreAcceptLight message and feeds it to the ordering plane. As client operations are transmitted with different overlays, it is possible that some referenced client operations are not present in the cache at the time the PreAcceptLight message arrives. In such a case, the data plane waits for a configurable period of time for the missing client operations to appear in the cache.

D. Controller

In this section, we present the design of the controller. Overlay management is the main responsibility of the controller. It applies ShallowForest to compute the optimal partitioning scheme and overlays for each site in response to changes in workload distribution and participant churn. The controller relies on the following procedures for overlay management:

- RequestClientDataRate (v): The controller invokes this procedure to request v's estimated average client data rate R_v .
- ApplyConfig (v, c, ver): The controller invokes this procedure to push an overlay configuration c with version ver to a site v.
- HandleViewChange (view, A, D): Upon participant churn, the controller call this procedure to update network information. The numerical value view is the view number associates with the current set of active participants. A, D are respectively the set of newly joined nodes and the set of departed nodes.

1) View Change: Consensus protocols and blockchain systems [6], [10], [27], [28] implement view change mechanisms (also referred to as reconfiguration) to handle updates to system configurations such as the set of active participants and current leader. Rather than re-inventing the wheel, the controller leverages the ordering plane to update its local state, which includes the participant set, available bandwidth and latency between participants. When the view changes, the ordering plane notifies the controller by sending a message $\langle view, \mathcal{A}, \mathcal{D} \rangle$, which includes the view number view, a set of newly joined nodes \mathcal{A} , and a set of leaving nodes \mathcal{D} . Upon receiving the message, the controller invokes procedure HandleViewChange presented in Algorithm 2. The controller maintains a copy of the last view number received from the ordering plane and ignores any message with a view number smaller than its local copy. The controller updates the local state through procedure UpdateState by adding information of nodes in \mathcal{A} and removing those of nodes in \mathcal{D} . Each newly joined node in \mathcal{A} is associated with its available bandwidth and network latency to existing participants. The controller then triggers an immediate overlay update via calling the UpdateOverlays procedure. To benefit from the optimized overlay, the ordering plane must wait for the completion of HandleViewChange before processing new client requests.

In APaxos, the ordering plane is EPaxos. EPaxos is not a consensus protocol designed for a dynamic network, where participants may join and leave at a high rate. Similarly to other classic consensus protocols, it operates on a static set of participants. For such protocols, view change is only triggered when the ordering plane detects node failures and replaces failed nodes with live ones.

2) Update Overlays: UpdateOverlays procedure presented in Algorithm 2 updates the overlay configuration of each site in the network. It is invoked either by a configured timeout period δ or upon view change. It collects client data rates through calling **RequestClientDataRate(v)** for each site $v \in G(V)$. After receiving the request, each site v reports its estimated average client data rate R_v computed by the exponential weighted moving average (EWMA) as follows:

$$R_v = (1 - \alpha)R_v + \alpha R'_v$$

Parameter α is a configurable weight and R'_v is the average client data rate at site v within the most recent time window of δ time units. By using EWMA, the estimated client data rate can rapidly converge to recent measurements as the old estimation decay exponentially. A successful invocation of **RequestClientDataRate(v)** returns v's current client data rate. Once the controller collects client data rates from all sites, it applies ShallowForest to calculate the latest partitioning scheme C, and updates each site v's current overlay configuration through calling the **ApplyConfig** procedure. We assume **RequestClientDataRate(v)** eventually terminates if both the controller and the remote site v are alive. If any call to **RequestClientDataRate** fails, the controller will not update overlays until the ordering plane triggers view change.

When **ApplyConfig**(v, C[v], ver) is triggered, the controller pushes overlay configuration C[v] with version ver to v. The

version *ver* is a tuple $\langle view, i \rangle$ representing the *i*-th configuration update within the view numbered by *view*. The data plane process updates its local overlay configuration to C[v] only when the version *ver* is greater than the version of its current local overlay configuration. Due to potential failures and unexpected network conditions, a site's latest overlay configuration can be lost or delayed. In such a scenario, the data plane process can detect outdated overlay configuration by inspecting the version field included in the message broadcast by other sites (described in Section V-C3). If a data plane process detects that its overlay version is outdated, it queries the controller for the latest configuration.

3) Overhead of Overlay Management: As mentioned earlier, the controller invokes procedure **UpdateOverlays** over a configured time interval δ and let $\overline{\delta}$ be the average view change interval. Assume each function invocation incurs a constant network capacity overhead β (such as RPC headers and network packet headers). The invocation of procedure **RequestClientDataRate** consumes β of v's downlink capacity and $\beta + \mu$ of a site v's uplink capacity, where μ is the size of reported client data rate in bytes. Let ω be the total size of **ApplyConfig**'s parameters, **ApplyConfig** consumes $\beta + \omega$ of a site v's downlink capacity and β of v's uplink capacities. After summing up aforementioned terms, overlay management consumes a total U_m of each site's uplink capacity and D_m of each site's downlink capacity. U_m and D_m are computed as follows:

$$U_m = \frac{(2\beta + \mu)(\delta + \overline{\delta})}{\delta\overline{\delta}} \quad D_m = \frac{(2\beta + \omega)(\delta + \overline{\delta})}{\delta\overline{\delta}}$$

ShallowForest targets workloads that fluctuate on timescales of hours or minutes. It is sufficient to set δ to 60 seconds. Assume in a deployment with 10 sites, $\overline{\delta}$ is around 30 seconds. Based on our measurement using gRPC [29], β is around 1 KB, μ is 8 bytes and ω is approximately 100 bytes. After substituting those values, both U_m and D_m are approximately 100 Bps, which are light-weight compared to the network capacity consumed by disseminating client data.

E. Handling Failures

EPaxos does not rely on a controller to determine data transmission overlays. To avoid the single point of failure, APaxos can deploy multiple controllers. When the current controller fails, a backup controller will continue to update each site's overlay configurations.

The other difference between EPaxos and APaxos lies in the way PreAccept messages are sent out. EPaxos assumes message passing is asynchronous between replicas, and introducing a data plane does not break this assumption. As a result, APaxos inherits the *safety* property from EPaxos. For the *liveness* property, EPaxos guarantees the client operation will eventually be committed if there are f + 1 non-faulty replicas and the client retries (possibly with another replica) if it does not receive a response within a timeout period. As a result, the data plane should guarantee all non-faulty replicas finally receive both PreAcceptLight and all client operations it references. Since direct broadcast will guarantee that all non-faulty replicas receive the data, client operations transmitted using a 1-level tree overlay and PreAcceptLight require no additional mechanism to ensure data delivery to all non-faulty replicas. However, if a message is transmitted using a 2-level tree overlay, all leaf nodes will not receive the message when the node in the middle crashes. To preserve the property that the protocol can make progress with f + 1 non-faulty replicas, the middle node sends ACK to the root replica after it completes sending the message to the remaining replicas. If the ACK from the middle node is not received after a timeout period, the root node broadcasts the message directly to other replicas and marks the middle node as a potentially crashed node, which needs to be avoided in future data transmissions.

VI. DISCUSSION

ShallowForest does not aim to directly optimize data transmission for each participant in large-scale systems (e.g., more than tens of nodes). State machine replication and consensus protocols generally do not exhibit improvements in performance through the addition of nodes. As a result, the scale of systems under consideration is naturally limited. Adding hundreds (or even tens) of nodes might render the upper layer consensus protocol the performance bottleneck before the performance becomes restricted by the data transmission. If the number of participants is large, blockchain systems usually limit the scale of consensus participants. For instance, in Algorand [8], a subset of participants called a committee is selected through verifiable random functions to participate in consensus. Except for limiting the scale of participants, organizing participants using a hierarchical topology is another approach adopted by large-scale consensus. Canopus [30] organizes participants into a multi-level leaf-only tree hierarchy, where a node at each level represents a group of nodes at the immediate lower level while the bottom level nodes are actual participants. ShallowForest has the potential to provide overlay optimization for dozens of nodes at the top level. Note that, in such a setting, the churn of logical nodes will not be frequent. It only happens when the network partition cuts off the entire region, which is rare.

ShallowForest relies on the ordering plane to handle participant churn and requires minor modifications to the ordering plane's view change process. As safety properties are defined for a specific view, consensus protocols cannot process new requests until the completion of the view change process. Therefore, participant churn will not happen within a view and it is redundant having a separate mechanism to handle it at the data plane.

VII. EVALUATION

This section presents the evaluation of ShallowForest, particularly the benefit of the ShallowForest optimization in terms of commit throughput and commit latency.

A. Experiment Setup

For the experiment, APaxos is deployed across nine Amazon EC2 regions: Tokyo (TK), Singapore (SG), Sydney (SY),

TABLE III Network Latency (ms) Between Each Pair of Sites Used in the Experiment

	IR	CA	VA	TK	OR	SY	FF	LD	SG
CA	146	-							
VA	74	64	-						
ТК	228	113	166	-					
OR	135	22	79	101	-				
SY	275	152	203	116	143	-			
FF	25	149	90	245	165	288	-		
LD	13	140	80	236	144	274	15	-	
SG	184	178	244	74	165	173	179	173	-

Frankfurt (**FF**), Ireland (**IR**), Oregon (**OR**), Virginia (**VA**), London (**LD**), and California (**CA**). Table III summarizes the network latency measured using ping between each pair of regions. In each region, the experiment uses an m4.xlarge VM instance with four 2.4 GHz Intel Xeon E5-2676v3 processors and 16 GB main memory. The OS version on each VM is Ubuntu 16.04 and the golang version used to compile APaxos is 1.9.4. A client process and an APaxos replica are executed on each VM, as well as a controller process on a single VM chosen at random. The client process sends requests only to its co-located APaxos replica.

The purpose of the experiment is to evaluate the effectiveness of the ShallowForest optimization when the network is the bottleneck. Therefore, the workload consists of only write requests as they involve broadcasting a significant amount of payload data. All requests are committed on the fast-path as each request is associated with a distinct key. For saturating the network resource provisioned to each VM, the request size is set to 4 KB and 20 TCP connections are established between each pair of VMs. Those parameter values are picked by increasing both request size and the number TCP connections until the throughput of APaxos cannot be improved further. In our experiments, we use the CUBIC congestion control algorithm [31], which is the default option on our VM instances. We consider CUBIC the preferred congestion control algorithm for our experiment environment, where the RTT between two replicas is significant. With CUBIC, window growth during congestion avoidance only depends on the real time between two consecutive congestion events. As a result, CUBIC can provide better fairness for flows competing for the same bottleneck, independently from their RTTs.

The client process can be configured to issue requests to an APaxos replica at a specific rate in an open loop. In the experiments, client request rates are enforced to be sustainable. The experiment uses the Zipfian distribution to model the skewed client load across replicas. For the network topology, the experiment uses the average network capacity measured in 1-minute intervals for a total of 30 minutes, and the average latency measured by 3 pings between each pair of replicas. APaxos+SF denotes APaxos optimized using ShallowForest in all results.

B. The Effect of Using Multiple TCP Flows

Besides the ShallowForest optimization, APaxos differs from original EPaxos by using multiple flows between each replica for data transmission. This experiment evaluates the



Fig. 5. Latency vs. throughput for 5 replicas with different flow numbers.

effectiveness of using multiple TCP flows by comparing the performance of original EPaxos and APaxos using a different number of flows between each pair replicas. We compare the performance of three candidates: EPaxos, APaxos using 5 TCP connections and APaxos using 20 TCP connections. For each candidate, the experiment uses five replicas located in IR, CA, VA, TK, and OR. The client loads on all replicas are equivalent.

Figure 5 presents the experimental results where each data point is the average of 4 runs and the error bar represents the standard deviation of 4 runs. Each run lasts for 20 seconds and the VM is warmed up through executing the workload for 20 seconds before each run. The result demonstrates that using a larger number of flows leads to a higher commit throughput. APaxos using 20 TCP connections achieves approximately 4.2X throughput compared to EPaxos, which uses only one TCP connection between each pair of replicas. We also note that the candidate using a smaller number of TCP connections can achieve lower commit latency for low client load in Figure 5. The reason behind it is that the candidate using a smaller number of TCP connections has higher load per connection. According to the CUBIC algorithm, if the congestion window size cwnd is less than the slow-start threshold ssthresh, cwnd is incremented additively per ACK. Therefore, higher load per connection results in faster convergence to optimal cwnd and lower transmission delay.

C. Different Skewness Levels

Workloads in the real world can be highly skewed across different geo-areas, and continuously changing over time. For this experiment, we use five geo-distributed replicas to evaluate the effectiveness of ShallowForest in dealing with skewed client load across replicas. We vary the exponent parameter s of the Zipfian distribution to tune the skewness level of client load distribution. Client load is uniformly distributed when s equals to zero, and increasing s leads to a more skewed client load distribution. Table IV demonstrates the load on each replica at different skewness levels. For each skewness level, we increase the aggregated client request rate and measure the commit throughput as well as the corresponding average commit latency. Figure 6 presents the experimental results where each data point is the average of 4 runs and the error bar represents the standard deviation of 4 runs. Each run lasts for



Fig. 6. Latency vs. throughput for 5 replicas with Zipfian workloads.

TABLE IV LOAD ON REPLICAS UNDER DIFFERENT SKEWNESS LEVELS

s	IR	CA	VA	ТК	OR
0.0	20%	20%	20%	20%	20%
0.5	38%	20%	16%	13%	12%
1.0	62%	15%	9%	6%	5%
1.5	82%	8%	4%	2%	1%

20 seconds and each VM is warmed by executing the workload for 40 seconds before each run.

As shown in Figure 6, the commit latency of APaxos without the ShallowForest optimization grows more rapidly with the increasing commit throughput due to contention for uplink capacity at replicas with high client load. When optimized using ShallowForest, APaxos achieves higher commit throughput with lower commit latency. For instance, when s = 1.5, ShallowForest improves the commit throughput of APaxos by 100% with 60% reduction in commit latency. Figure 6 also shows that ShallowForest only improves APaxos slightly when the client load is moderately skewed. When s = 0.5, ShallowForest improves the commit throughput of APaxos by 10% with 30% reduction in commit latency. As each VM instance is provisioned with similar uplink and downlink capacity, the optimal overlays become increasingly 1-level tree dominated with a more uniformly distributed load. This also explains why ShallowForest brings no performance improvement when s = 0 in this experiment.

D. Effect of Replication Factor on Performance

We also compare the effectiveness of ShallowForest for various replication factors. For this experiment, we measure the commit throughput of APaxos with five, seven and



Replication Factor

Fig. 7. Throughput of different numbers of replicas.

nine geo-distributed replicas. For all replication factors, we set s to 1 and the aggregated data rate of incoming client requests to 750 Mbps. Figure 7 demonstrates the experimental results, where each bar is the average of 4 runs and the error bar represents the standard deviation. ShallowForest improves the commit throughput of APaxos by 43%, 32% and 32% for replication factors 5, 7, and 9. When optimized using ShallowForest, the commit throughput drops faster with the increasing number of replicas. This is due to the higher network latency yielded by 2-level tree overlays in a larger scale deployment, which contains replicas deployed in more distant regions (Sydney and Singapore). However, APaxos optimized using ShallowForest still achieves higher commit throughput for all replication factors resulting from more effective use of the network capacity at replicas with low client load.

VIII. RELATED WORK

Consensus Over WAN Mencius [15] is a variant of Paxos that rotates the leader for each command to distribute the load evenly across replicas. Mencius also addresses the issue of unevenly distributed client load across replicas. It allows a replica with low client load to voluntarily skip its leader term to favour replicas with higher client load. Mencius is not completely leaderless and skipping a leader term cannot help a replica with high client load to utilize network resources at a replica with low client load. E-Paxos [16] further improves scalability and reduces commit latency by removing the role of the leader. Each client is able to send the request to the nearest replica, which is referred to as the command leader. However, unlike APaxos, EPaxos does not consider the availability of network capacity and unbalanced client loads across replicas. Canopus [30] is a network-aware consensus protocol that parallelizes the dissemination of messages according to a leaf only tree (LOT). LOT organizes nodes into several consensus groups based on locality. The main goal of using LOT for data transmission is to minimize the usage of highly contented links. In terms of data transmission, LOT might incur higher network latency in a WAN environment. For *n* nodes, LOT requires $O(\log n)$ transfers for data to reach all nodes, while ShallowForest requires at most two transfers.

Decoupling Data Transmission From Ordering Decoupling data transmission from ordering is a common technique used by many consensus protocols [11], [12], [17] and blockchain systems. To separate the ordering plane from data transmission, S-Paxos [17] associates each batch of client requests with a unique ID and uses Paxos as its ordering plane protocol to order batch IDs. Disseminating client requests is handled by a separate process. The data plane of S-Paxos is leaderless, meaning that a client may contact any replica to broadcast the request to other replicas. Ring Paxos [11], [12] handles data dissemination using one or more logical ring overlays. To multicast messages to a group of receivers, all servers are placed on a logical ring. The sender just sends data once to its immediate successor and all subsequent receivers store-andforward the message until the last receiver receives it. Logical ring overlay minimizes the data replication at the sender to achieve high throughput data transmission and the optimal network utilization. However, when using a ring overlay, the network latency of the data transmission grows proportionally with the number of participants. Both Ring Paxos and S-Paxos handle the data dissemination with a static overlay that does not change adaptively to various client loads across replicas. Some permissionless blockchain systems [4], [5], [32] use gossip protocols [33] for high-throughput data dissemination. ShallowForest does not apply to those systems as it requires the location and identity of each participant to be known a priori.

Application-Level Multicast Application level multicast [21], [22], [34] has been studied extensively for content distribution in P2P networks. Among those systems, ShallowForest is most similar to SplitStream and Bullet network [21], [22]. SplitStream [21] partitions the data to be broadcast into several disjoint sections called stripes and

constructs a separate broadcast tree for each stripe. To receive the complete stream, a node must be presented in every broadcast tree. SplitStream enforces that any two broadcast trees must be interior-node-disjoint, which means every node is an interior node in precisely one tree and a leaf node in all other trees. This property improves the robustness of the system because the failure of a node only affects the delivery of a single stripe. Bullet [22] divides the data into multiple disjoint blocks which are further divided into packet-size objects. For data dissemination, Bullet uses an epoch-based algorithm called RanSub for membership management and overlay construction. Both protocols focus on optimizing data transmission throughput and do not impose any constraints on the height of overlays. ShallowForest only uses shallow tree overlays and optimizes both the throughput and the network latency subject to the data transmission.

Compared to prior works, our contribution is twofold. First of all, we mathematically prove that consensus protocols can achieve the optimal throughput by using overlays with height up to two. This result provides a constant bound on latency for the throughput-latency trade-off in consensus protocols and state machine replication. Second, our work is the first to address the problem of skewed client load distribution from the angle of optimizing data transmission for consensus protocols. Prior works [15], [16] focus on protocol level optimizations such as rotating leadership and leaderless architecture.

IX. CONCLUSION

In this paper, we presented a method of optimizing data transmission in a WAN environment, called ShallowForest, and applied to the widely-cited EPaxos consensus protocol. The key idea of ShallowForest is to partition the data stream and use shallow tree overlays for data transmission. The experimental results demonstrate that ShallowForest can make a consensus protocol more resilient to skewed load by handling the data transmission in a more workload-aware and network-aware manner. In future work, we plan to build a fully autonomous controller that can automatically estimate client load and available network capacity.

REFERENCES

- L. Lamport, "Using time instead of timeout for fault-tolerant distributed systems." ACM Trans. Program. Lang. Syst., vol. 6, no. 2, pp. 254–280, Apr. 1984.
- [2] M. Burrows, "The Chubby lock service for loosely-coupled distributed systems," in *Proc. 7th USENIX Symp. Oper. Syst. Design Implement.*, 2006, pp. 335–350.
- [3] P. Hunt, M. Konar, F. P. Junqueira, and B. Reed, "Zookeeper: Wait-free coordination for Internet-scale systems," in *Proc. USENIX Annu. Tech. Conf.*, 2010, pp. 145–158.
- [4] S. Nakamoto, Bitcoin: A Peer-to-Peer Electronic Cash System. 2008.
 [Online]. Available: http://bitcoin.org/bitcoin.pdf
- [5] V. Buterin, Ethereum: A Next-Generation Smart Contract and Decentralized Application Platform. 2014. [Online]. Available: https://github.com/ethereum/wiki/wiki/White-Paper
- [6] E. Androulaki *et al.*, "Hyperledger fabric: A distributed operating system for permissioned blockchains," in *Proc. 13th ACM EuroSys Conf.*, 2018, pp. 1–15.
- [7] M. Zamani, M. Movahedi, and M. Raykova, "Rapidchain: Scaling blockchain via full sharding," in *Proc. ACM SIGSAC Conf. Comput. Commun. Security*, 2018, pp. 931–948.

- [8] Y. Gilad, R. Hemo, S. Micali, G. Vlachos, and N. Zeldovich, "Algorand: Scaling byzantine agreements for cryptocurrencies," in *Proc. 26th Symp. Oper. Syst. Principles*, 2017, pp. 51–68.
- [9] L. Lamport, "Paxos made simple," ACM SIGACT News, vol. 32, no. 4, pp. 51–58, Dec. 2001.
- [10] D. Ongaro and J. Ousterhout, "In search of an understandable consensus algorithm," in Proc. USENIX Annu. Tech. Conf., 2014, pp. 305–319.
- [11] P. J. Marandi, M. Primi, N. Schiper, and F. Pedone, "Ring Paxos: A high-throughput atomic broadcast protocol," in *Proc. IEEE/IFIP Int. Conf. Depend. Syst. Netw.*, 2010, pp. 527–536.
- [12] P. J. Marandi, M. Primi, and F. Pedone, "Multi-ring Paxos," in Proc. IEEE/IFIP Int. Conf. Depend. Syst. Netw., 2012, pp. 1–12.
- [13] R. van Renesse and F. B. Schneider, "Chain replication for supporting high throughput and availability," in *Proc. 6th USENIX Conf. Oper. Syst. Design Implement.*, 2004, pp. 91–104.
- [14] J. Terrace and M. J. Freedman, "Object storage on CRAQ: Highthroughput chain replication for read-mostly workloads," in *Proc. Conf.* USENIX Annu. Tech. Conf., 2009, p. 11.
- [15] Y. Mao, F. P. Junqueira, and K. A. Marzullo, "Mencius: Building efficient replicated state machines for wans," in *Proc. 8th USENIX Conf. Oper. Syst. Design Implement.*, 2008, pp. 369–384.
- [16] I. Moraru, D. G. Andersen, and M. Kaminsky, "There is more consensus in egalitarian parliaments," in *Proc. 24th ACM Symp. Oper. Syst. Principles*, 2013, pp. 358–372.
- [17] M. Biely, Z. Milosevic, N. Santos, and A. Schiper, "S-Paxos: Offloading the leader for high throughput state machine replication," in *Proc. 31st IEEE Symp. Reliable Distrib. Syst.*, 2012, pp. 111–120.
- [18] K. Yamada, H. Takayasu, T. Ito, and M. Takayasu, "Solvable stochastic dealer models for financial markets," *Phys. Rev. E*, vol. 79, May 2009. Art. no. 051120.
- [19] G. Amvrosiadis, J. W. Park, G. R. Ganger, G. A. Gibson, E. Baseman, and N. DeBardeleben, "On the diversity of cluster workloads and its impact on research results," in *Proc. USENIX Conf. Usenix Annu. Tech. Conf.*, 2018, pp. 533–546.
- [20] Z. Cao, S. Dong, S. Vemuri, and D. H. Du, "Characterizing, modeling, and benchmarking rocksdb key-value workloads at facebook," in *Proc.* 18th USENIX Conf. File Storage Technol., Santa Clara, CA, USA, 2020, pp. 209–223.
- [21] M. Castro, P. Druschel, A.-M. Kermarrec, A. Nandi, A. Rowstron, and A. Singh, "Splitstream: High-bandwidth multicast in cooperative environments," in *Proc. 19th ACM Symp. Oper. Syst. Principles*, 2003, pp. 298–313.
- [22] D. Kostić, A. Rodriguez, J. Albrecht, and A. Vahdat, "Bullet: High bandwidth data dissemination using an overlay mesh," *SIGOPS Oper. Syst. Rev.*, vol. 37, no. 5, pp. 282–297, Oct. 2003.
- [23] H. Tan and W. M. Golab, "Optimizing all-to-all data transmission in WANs," in *Proc. IEEE Int. Conf. Blockchain Cryptocurrency (ICBC)*, 2020, pp. 1–9.
- [24] N. G. Duffield, P. Goyal, A. Greenberg, P. Mishra, K. K. Ramakrishnan, and J. E. van der Merive, "A flexible model for resource management in virtual private networks," *SIGCOMM Comput. Commun. Rev.*, vol. 29, no. 4, pp. 95–108, Aug. 1999.
- [25] P. M. Vaidya, "An algorithm for linear programming which requires $O((m + n)n^2 + (m + n)^{1.5}n)L)$ arithmetic operations," *Math. Program.*, vol. 47, no. 2, pp. 175–201, 1990.

- [26] I. Moraru, D. G. Andersen, and M. Kaminsky. *Epaxos Source Code*. 2014. [Online]. Available: https://github.com/efficient/epaxos
- [27] M. Castro and B. Liskov, "Practical byzantine fault tolerance," in Proc. 3rd USENIX Symp. Oper. Syst. Design Implement., 1999, pp. 173–186.
- [28] B. Liskov and J. Cowling, "Viewstamped replication revisited," MIT Comput. Sci. Artif. Intell. Lab., Cambridge, MA, USA, Tech. Rep. MIT-CSAIL-TR-2012-021, Jul. 2012.
- [29] gRPC Homepage. Accessed: Oct. 5, 2020. [Online]. Available: https://grpc.io/
- [30] S. Rizvi, B. Wong, and S. Keshav, "Canopus: A scalable and massively parallel consensus protocol," in *Proc. 13th Int. Conf. Emerg. Netw. Exp. Technol.*, 2017, pp. 426–438.
- [31] S. Ha, I. Rhee, and L. Xu, "Cubic: A new TCP-friendly high-speed TCP variant," SIGOPS Oper. Syst. Rev., vol. 42, no. 5, pp. 64–74, 2008.
- [32] J. Wang and H. Wang, "Monoxide: Scale out blockchains with asynchronous consensus zones," in *Proc. 16th USENIX Symp. Netw. Syst. Design Implement.*, 2019, pp. 95–112.
- [33] A. Demers et al., "Epidemic algorithms for replicated database maintenance," in Proc. 6th Annu. ACM Symp. Principles Distrib. Comput., 1987, pp. 1–12.
- [34] O. Papaemmanouil, Y. Ahmad, U. Çetintemel, J. Jannotti, and Y. Yildirim, "Extensible optimization in overlay dissemination trees," in *Proc. ACM SIGMOD Int. Conf. Manag. Data*, 2006, pp. 611–622.



Hao Tan received the B.Sc. degree in computer science from the University of Toronto, and the M.Math. degree in computer science from the University of Waterloo, where he is currently pursuing the Ph.D. degree in electrical and computer engineering. His research focuses on building reliable, fast, and smart distributed systems.



Wojciech Golab (Member, IEEE) received the Ph.D. degree in computer science from the University of Toronto in 2010. After a Postdoctoral Fellowship from the University of Calgary, he spent two years as a Research Scientist with Hewlett-Packard Labs, Palo Alto. He then joined the University of Waterloo in 2012, where he is currently an Associate Professor in Electrical and Computer Engineering. His research agenda focuses on algorithmic problems in distributed computing with applications to the design, optimization, and verification of soft-

ware systems for data storage, and analytics. He is a recipient of an Ontario Early Researcher Award, and two Google Faculty Research Awards. His research publications have won two Best Papers Awards, and his doctoral work on shared memory algorithms was distinguished by the ACM Computing Reviews as one of 91 "notable computing items published in 2012."