

# HP: Hybrid Paxos for WANs\*

Dan Dobre, Matthias Majuntke, Marco Serafini and Neeraj Suri  
{dan,majuntke,marco,suri}@cs.tu-darmstadt.de

## Abstract

*Implementing a fault-tolerant state machine boils down to reaching consensus on a sequence of commands. In wide area networks (WANs), where network delays are typically large and unpredictable, choosing the best consensus protocol is difficult. During normal operation, Classic Paxos (CP) requires three message delays, whereas Fast Paxos (FP) requires only two. However, when collisions occur, due to interfering commands issued concurrently, FP requires four extra message delays. In addition, FP uses larger quorums than CP. Therefore, CP can outperform FP in many situations.*

*We present Hybrid Paxos (HP), a consensus protocol that combines the features of FP and CP. HP implements generalized consensus, where collisions are caused only by interfering commands. In the absence of collisions HP requires two message delays, and only one extra message delay otherwise. Our evaluation shows that when collisions are rare, the latency of HP reaches the theoretical minimum. When collisions are frequent, HP behaves like CP.*

## 1 Introduction

WAN replication offers protection against catastrophic failures of a single site and can be used to enhance the resilience of critical services. Implementing a deterministic service in a fault-tolerant manner boils down to reaching consensus on a sequence of system commands.

In the standard state-machine approach, a sequence of instances of a consensus protocol are used to choose the sequence of client commands, where the  $i$ th instance chooses the  $i$ th command. In this paper we consider generalized consensus [18], where a single instance of consensus is used to choose an increasing history of commands. A history is an equivalence class of command sequences, where two command sequences are equivalent iff executing them produces the same state and output. The underlying observation is that often commands *commute*, so it does not matter in which order they are executed.

The consensus problem is stated in terms of *proposers* that propose commands, *acceptors* that choose an increasing command history and *learners* that learn what history has been currently chosen. In a client/server system, clients might play the roles of proposer and learner and servers might play the role(s) of acceptor (and learner). A leader is elected among the acceptors to coordinate their actions.

In a WAN environment, where network delays are large and unpredictable, the *latency* of a consensus protocol matters. Latency is defined as the number of message delays between when a client proposes a command and when that command is learned by a learner.

Consensus protocols which attain the optimal latency [20] are the well known Classic Paxos (CP) [17] and the more recent Fast Paxos (FP) [18]. Their message patterns are illustrated in Figure 1. In normal operation, CP requires three message delays. The communication pattern during normal operation is Client  $\rightarrow$  Leader  $\rightarrow$  Acceptors  $\rightarrow$  Learners. FP saves one message delay by having the clients send their proposals directly to the acceptors, bypassing the leader. This works fine if the acceptors receive the same sequence of interfering commands. However, when commands are proposed concurrently, commands may be accepted in interfering orders, resulting in no command being chosen. In order to guarantee progress, FP then runs a collision recovery procedure, which adds four message delays. Thus, if collisions are frequent, FP has a significantly higher latency and a lower throughput than CP.

We found that even in the absence of collisions, depending on the the layout of clients and servers, CP can outperform FP (for many clients). This comes from the fact that in order to be fast, FP needs larger quorums than CP, called *fast quorums* [20].

When clients have direct access to a local replica, the recently developed consensus protocol Mencius [21] has been shown to outperform CP. However often, clients and servers are not co-located. When clients are using a remote service replicated for disaster tolerance, none of the mentioned protocols has the final say.

---

\*Research funded in part by Inco-Trust, Microsoft Research and DFG GRK 1362 (TUD GKM).

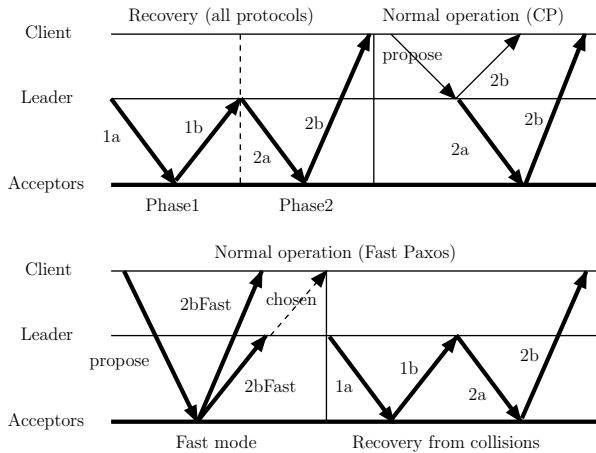


Figure 1. Paxos message patterns

**Paper Contributions** In this paper we present a generalized consensus protocol called Hybrid Paxos (HP). HP is essentially CP with an additional “fast mode” that enables fast learning in the absence of collisions. In presence of collisions, HP requires three message delays as CP does. These latencies are optimal [20] and they are attained using a linear number of messages and the optimal number of  $2f + 1$  servers, where  $f$  is the bound on crash-failures. Compared to Mencius, HP uses weaker synchrony assumptions, resulting in higher availability in WANs.

We show for the first time that generalized consensus can be used in practice to build efficient replicated services. The key to efficiency is that fast learning must not impact the bottleneck, which in CP is the leader. Additional messages in HP are exchanged only between clients (which are both proposers and learners) and acceptors. Thus, HP exploits the relative underutilization of the acceptors and offers a better latency up to 70% of the peak throughput of CP.

In addition, fast learning is enabled only if spare capacity is available. This is done by adaptively switching it on and off based on the load. Our evaluation using Emulab [26] shows that the latency of HP reaches the theoretical minimum. In the presence of collisions and with increasing load, HP behaves like CP.

## Motivation

**There is no clear winner with Fast and Classic Paxos** We argue that the quorum size matters by showing that even in the absence of collisions, CP can outperform FP.

We have sampled delays among Planetlab [3] nodes and used them to simulate the normal operation (best case) of CP and FP in four different WAN settings (Table 1). The client distribution is as follows: 56% are located in the US, 38% in Europe, and 6% in Asia. All topologies use 11 servers. FP requires a fast quorum (9 servers), while CP only requires a (majority) quorum. The simulation results in

Figure 2 suggest that: (a) in many settings, some clients are better off using FP and others prefer CP and (b) the crash of a single server can turn a setting beneficial for FP into one beneficial only for CP. Thus, there exist practical settings where neither of the two protocols always outperforms the other.

Table 1. WAN server layout (11 servers)

Topology	Europe	World	CLUS-5(4)
Leader Site	Hungary	Japan	Switzerland
Backups	Europe	Global	Europe
Clusters	No	No	1 with 5 (4) nodes, 3 with 2 nodes
Quorums	6 servers for CP, 9 servers for FP		

In the *Europe* setting, servers are located at 11 different sites in Europe. For most clients, the distance between them and the servers is close to uniform. Thus, the FP pattern leads to good results: 28% of the clients observe that FP is at least 10% faster than CP, and 10% of the clients even observe a 20% improvement. However, 12% of the clients find that CP is 10% better, as they do not have good connectivity with three additional servers required by FP. This supports observation (a).

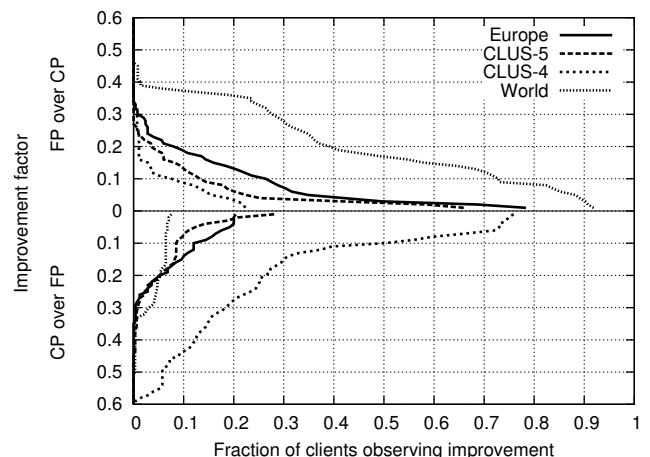


Figure 2. Improvement factor of FP over CP and vice versa

The *World* setting models a world-wide setting in which acceptors are spread over the US, Europe and Asia, and the leader is located in Asia. In this scenario the advantage of the FP pattern is even clearer: 75% of the clients observe a 10% improvement over CP, and for 20% of the clients the improvement is  $> 35\%$ . The reason is the additional message delay to reach the leader, which is large for most of the clients. However, there are some clients which prefer CP: 6% find it to be 20% more efficient than FP, supporting

observation (a). This is essentially the fraction of Asian clients which can quickly reach the leader.

The *CLUS* topology considers the case when servers are clustered at four different sites in Europe, providing cheap disaster tolerance. The only distinction between the *CLUS-5* and *CLUS-4* is that in the latter, one node in the largest cluster is crashed (Table 1). Before the crash, a fast quorum (9 servers) can be reached by contacting three sites. Note that for 13% of the clients, FP outperforms CP by at least 10%. Two sites are sufficient for CP to sample a quorum (6 servers). However, after the crash FP accesses *all* sites. This results in a shift of the performance profile, with CP dominating FP for 50% of the clients, supporting observation (b).

## 2 System Model and Definitions

We consider a distributed system consisting of  $n$  servers and any number of clients. For simplicity, we assume reliable FIFO channels. Further we consider a crash-stop model in which clients and servers fail only by crashing and nonfaulty servers never crash.<sup>1</sup> We assume that at most a minority of servers fails, which is necessary to solve consensus [7]. The system is *asynchronous*, with no bounds on message delay or processor speed. However, each server has access to a failure detector  $\Omega$ , that eventually outputs at all servers the same nonfaulty server [6].

**Mathematical Preliminaries** In the generalized consensus problem [18] a single consensus instance is performed on a monotonically increasing and partially ordered set of commands, called *command history*. A command history, or history for brevity, is defined as an equivalence class of command sequences. Two command sequences are equivalent iff they can be transformed into one another by permuting its elements such that the order of every pair of interfering commands is preserved. Two commands are *interfering* if it matters in which order they are executed. Else, they are called *commutable*.

Histories are constructed by appending a command sequence  $\sigma$  to the initially empty history  $\perp$  using the special append operator  $\bullet$ . The resulting history is  $\perp \bullet \sigma$ . Histories  $\perp \bullet \sigma$  and  $\perp \bullet \tau$  are equal iff  $\sigma$  and  $\tau$  are equivalent sequences.

The prefix relation  $\sqsubseteq$  is defined as a partial order on the set of histories. For two histories  $h$  and  $h'$ ,  $h \sqsubseteq h'$  iff there is a command sequence  $\sigma$  such that  $h \bullet \sigma = h'$ . We say that  $h$  is a prefix of  $h'$  (or equivalently that  $h'$  is an extension of  $h$ ). A history  $h$  is isomorphic to a directed graph  $G(h)$  whose nodes are the commands. There is an edge between

<sup>1</sup>The algorithm can be easily extended to a model in which crashed nodes may recover [1] and links are fair-lossy [2]. This however, lies outside the scope of the paper.

any two interfering commands  $c_i$  and  $c_j$  from  $c_i$  to  $c_j$  in  $G(h)$  iff  $i < j$  in  $h$ . For two histories  $h$  and  $h'$ , it holds that  $h \sqsubseteq h'$  iff the graph  $G(h)$  is a prefix of the graph  $G(h')$ .  $G(h) = G(h')$  iff  $h = h'$ .

A lower bound of a set  $H$  of histories is a history that is a prefix of every element in  $H$ . The greatest lower bound (*glb*) of  $H$  is a lower bound of  $H$  that is an extension of every lower bound of  $H$ . We write the *glb* of  $H$  as  $\sqcap H$  and we let  $h \sqcap h'$  equal  $\sqcap\{h, h'\}$  for any two histories  $h$  and  $h'$ . The least upper bound (*lub*) is defined in the analogous manner. We write *lub* of  $H$  as  $\sqcup H$  and we let  $h \sqcup h'$  equal  $\sqcup\{h, h'\}$ . Intuitively, the *glb* (resp. *lub*) of a set of histories is the largest common prefix (resp. the smallest common extension).

We define two histories  $h$  and  $h'$  to be *compatible* iff they have a common upper bound, i.e., there is some history  $g$  with  $h \sqsubseteq g$  and  $h' \sqsubseteq g$ . A set of histories  $H$  is compatible iff every pair of histories in  $H$  are compatible.

**Consensus Properties** Generalized consensus ensures that if any two command histories are learned, then they are compatible (*Consistency*). To rule out choosing a default value, it must hold that if history  $h$  is chosen then there exists a proposed command sequence  $\sigma$ , such that  $h = \perp \bullet \sigma$  (*Nontriviality*). Finally, if a learner process learns history  $h$ , then  $h$  was chosen (*Conservatism*). Liveness requires that if command  $c$  is proposed, then eventually some history containing  $c$  is learned (*Progress*).

## 3 Generalized Consensus and Paxos

In this section we review Paxos and describe it as consensus on a growing command history [18]. When needed, we differentiate between FP and CP. This description serves as a basis for HP which is introduced in Section 4.

In the client/server system that we consider, clients are both proposers and learners. The servers are acceptors and cooperate to choose a single command history. Acceptors query the  $\Omega$  failure detector that elects a leader among them. Safety is guaranteed even if no leader or multiple leaders are elected, but a unique leader is required to ensure progress. Paxos operates on a set of round numbers. The round numbers are partitioned among the potential leaders such that each leader has its disjoint set of round numbers.

As mentioned in the introduction, Paxos assumes predefined sets of acceptors called *quorums*. The requirement for CP is that any two quorums intersect. FP requires larger quorums, called *fast quorums*, and the requirement is that the intersection of any fast quorum  $FQ$  and any quorum is larger than  $n - |FQ|$ .

Following the Paxos protocol description [17], we divide the acceptor and leader actions in Phase 1 and Phase 2 actions. Phase 1 actions are executed when a new round is

started (e.g. after a leader crash). Phase 2 actions (1) complete the choosing all the histories that failed to be chosen in an earlier round and (2) they are repeatedly executed during normal operation.

We now describe the algorithms' actions below (see also Fig. 1 as an illustration). Note that the focus lies on consensus, and therefore the execution of commands is omitted from the description.

### Phase1: Start a new round

**(1a)** Leader  $l$  picks a new round number  $r$  from its set of round numbers and sends a  $\langle \text{"1a"}, r \rangle$  message to all acceptors.

**(1b)** When acceptor  $a$  receives a  $\langle \text{"1a"}, r \rangle$  message from leader  $l$ , if it has not yet received a message with a higher round number, then it replies with a  $\langle \text{"1b"}, r, \dots \rangle^2$  message. We say that  $a$  has moved to round  $r$  and considers  $r$  as its current round from now on. Moreover,  $a$  stops accepting proposals from clients.

If  $a$  has already received a message with round number  $r' > r$  then it sends a message to the leader, indicating that it is ignoring the  $\langle \text{"1a"}, r \rangle$  message. (Upon receiving such a message the leader performs step (1a) with a round number  $> r'$  if it still believes to be the leader.)

### Phase2: Complete earlier rounds

**(2a)** If leader  $l$  has received  $\langle \text{"1b"}, r, \dots \rangle$  messages from a quorum of acceptors, then it sends a  $\langle \text{"2a"}, r, h \rangle$  message to the acceptors where  $h$  is the history that has been determined from the received "1b" messages. Further, the leader adopts  $h$  as the history currently chosen. The rule of picking  $h$  depends on the type of protocol and is described later (see Sections 3.1 and 4.2.1).

**(2b)** If acceptor  $a$  receives a  $\langle \text{"2a"}, r, h \rangle$  message in its current round  $r$  (i.e. it has not yet received any message with a higher round number), it stores  $h$  as the accepted history and sends a  $\langle \text{"2b"}, r, h \rangle$  message to every learner. Next,  $a$  starts accepting proposals from clients.

**(Learn)** If a learner receives identical  $\langle \text{"2b"}, r, h \bullet c \rangle$  messages from a quorum, then it learns that history  $h \bullet c$  is chosen.

<sup>2</sup>Note that "... " will be replaced by protocol specific information later.

### Normal operation CP

**(Propose  $c$ )** Client  $cl$  sends a  $\langle \text{"propose"}, c \rangle$  message to the leader.

**(2aCP)** When the leader receives a  $\langle \text{"propose"}, c \rangle$  message from client  $cl$ , it appends command  $c$  to  $h$  and sends a  $\langle \text{"2a"}, r, h \bullet c \rangle$  message to the acceptors.

**(2bCP)** If acceptor  $a$  receives a  $\langle \text{"2a"}, r, h \bullet c \rangle$  message from the leader in its current round  $r$  (i.e. it has not yet received any message with a higher round number), then it accepts  $h \bullet c$  and sends a  $\langle \text{"2b"}, r, h \bullet c \rangle$  message to all learners. Learning is done as described in the (Learn) step.

### Normal operation FP

**(ProposeFP  $c$ )** Client  $cl$  sends  $\langle \text{"propose"}, c \rangle$  messages to the acceptors.

**(2bFast)** If acceptor  $a$  receives a  $\langle \text{"propose"}, c \rangle$  message from client  $cl$ , then  $a$  appends  $c$  to its command history  $h$  and sends  $\langle \text{"2bFast"}, r, h \bullet c \rangle$  messages to the learners and to the leader.

**(Collision Handling)** If the leader receives identical  $\langle \text{"2bFast"}, r, h \bullet c \rangle$  messages from a fast quorum, it indicates to the learners that  $h \bullet c$  is chosen by sending  $\langle \text{"chosen"}, r, h \bullet c \rangle$  messages to the learners. Else, the leader initiates collision recovery, which entails starting a new round (Phase 1) and recovering from earlier rounds (Phase 2).

**(Fast Learn)** If a learner receives identical  $\langle \text{"2bFast"}, r, h \bullet c \rangle$  messages from a fast quorum, or equivalently a  $\langle \text{"chosen"}, r, h \bullet c \rangle$  message then it learns that  $h \bullet c$  is chosen. "Slow" learning is done as in (Learn).

## 3.1 The rule of picking a history

We now explain the core of the Paxos protocol and why it satisfies Consistency. For this purpose, we now informally describe the rule of picking a history based on the  $\langle \text{"1b"}, r, \dots \rangle$  messages received by the leader in step (2a). A formal and complete treatment appears in an earlier work by Lamport [18].

**Invariant** Paxos maintains the following invariant for safety: if a history  $h$  is chosen in round  $r$  and a history  $h'$  is chosen in a higher numbered round, then  $h \sqsubseteq h'$ . Intuitively, Consistency follows from this invariant and the fact that once a quorum of acceptors has joined a higher numbered round, no history can be chosen in previous rounds anymore.

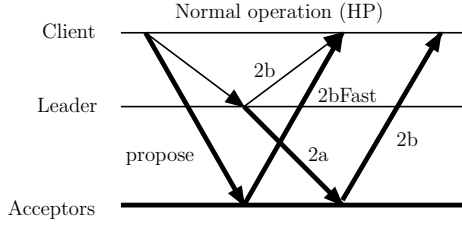


Figure 3. HP message pattern

**Pick classic** In CP, if a history  $h$  has been chosen then learning implies that a quorum of acceptors has accepted  $h$ . In step (1b), each acceptor reports the history it has accepted. By the quorum intersection property, at least one acceptor reports  $h$ . The picking rule is to select the *lub* among the reported histories. It is not difficult to see that  $h \sqsubseteq \text{lub}$ .

**Pick fast** In FP, if a history  $h$  has been chosen, then fast learning implies that a fast quorum  $FQ$  has accepted  $h$ . In step (1b), each acceptor reports the history it has accepted. Let  $Q$  be the set of all reported histories collected by the leader in step (2a). By the intersection property of  $FQ$  with a quorum,  $Q$  contains at least  $n - |FQ| + 1$  (possibly incompatible) extensions of  $h$  and at most  $n - |FQ|$  histories which are not extensions of  $h$ . Hence, there is a majority subset  $M \subseteq Q$  containing the extensions of  $h$ . The goal is now to find a history which is an extension of  $h$  using this knowledge. First, the *glb* is computed for every majority subset  $M \subseteq Q$ . As all majorities intersect, the *glbs* are pairwise compatible. Next, the leader picks the *lub* of these *glbs*. Note that one of the *glbs* is an extension of  $h$ , and therefore  $h \sqsubseteq \text{lub}$ .

## 4 The HP Protocol

As mentioned in the introduction, HP is essentially CP with an integrated “fast mode”, which allows fast learning in the absence of collisions. Therefore, the progress property of HP is inherited from CP. Hence, throughout the section, we will focus on safety.

The roles played by clients and servers and their interaction are the same as in Paxos (see Section 3). Phase 1 and 2 do not change and so they are as depicted in Figure 1. Figure 3 illustrates the message pattern of HP during normal operation.

### 4.1 Overview

We now briefly summarize the main differences between HP and Paxos.

First, fast learning is refined to accommodate that learning and fast learning are done in parallel, such that a learner can learn the quickest outcome. A naïve composition would

fail, as two incompatible histories could be learned in the (Fast Learn) and (Learn) steps. We prevent this problem by replacing fast learning with *hybrid learning*. The idea of hybrid learning is that a learner waits for a fast quorum of identical “*2bFast*” and “*2b*” messages, of which at least one is of type “*2b*”. Note that hybrid learning is fast because the leader is an acceptor (see Figure 3).

Secondly, the rule of picking a history in step (2a) is extended accordingly. In HP, each acceptor keeps two separate histories, a *classic* history updated by “*2a*” messages and a *fast* history updated by client proposals. Both histories are reported to the leader in step (1b). The leader applies the picking rules as described in Section 3.1 to each type separately. Resulting are two histories  $h$  and  $fh$ , where  $fh$  is determined by the fast histories. The final history is determined as the *lub* of  $h$  and the largest prefix of  $fh$  which is compatible with  $h$ .

### 4.2 The Protocol

We now describe the actions of the HP protocol during normal operation. The focus lies on highlighting the difference to CP. A complete description in pseudocode and proofs can be found in the full paper [9].

Phase 1 and 2 actions as well as actions (2aCP), (2bCP) and (Learn) are actions of the HP protocol. Since they do not change, they are not listed below.

#### Normal Operation

**(ProposeHP  $c$ )** Client  $cl$  executes the actions (Propose  $c$ ) and (ProposeFP  $c$ ).

**(2bFastHP)** If acceptor  $a$  receives a  $\langle \text{“propose”}, c \rangle$  message from client  $cl$ , then  $a$  appends  $c$  to the local fast history  $fh$  and sends  $\langle \text{“2bFast”}, r, fh \bullet c \rangle$  messages to the learners. (Note that the difference to action (2bFast) is that  $c$  is appended to the fast history, and that no “*2bFast*” messages are sent to the leader).

**(Hybrid Learn)** If a learner receives identical  $\langle \text{“2bFast”}, r, fh \bullet c \rangle$  messages from a fast quorum and one  $\langle \text{“2b”}, r, h \bullet c \rangle$  message and  $h = fh$  then it learns that  $h \bullet c$  is chosen. “Slow” learning is done as in (Learn).

#### 4.2.1 The rule of picking a history

We now explain the rule of picking a history in HP just as we did for Paxos. We will widely reuse the steps in Section 3.1 and refer to them when needed.

In HP, the  $\langle \text{“1b”}, r, \dots \rangle$  messages report two separate, accepted histories, the (classic) history and the fast history. The leader uses the reported (classic) histories to pick a history  $h$  as described in **Pick classic**. Next, the leader uses

the reported fast histories to pick a history  $fh$  as described in **Pick fast**. Note that each of the histories satisfies the invariant. History  $h$  is an extension of any history learned in the (Learn) step and  $fh$  is an extension of any fast history learned in the (Hybrid Learn) step.

**Pick hybrid** We now explain the key difference between HP and Paxos. To be safe, ideally we would pick the  $lub$  of  $h$  and  $fh$  and initialize the acceptors with  $lub$  in phase (2a). However, if  $h$  and  $fh$  are incompatible, then their  $lub$  is undefined. Therefore, the idea is (1) to determine the largest prefix  $pfh$  of  $fh$  which is compatible with  $h$  and (2) to pick the  $lub$  of  $pfh$  and  $h$ . This would be safe only if we can guarantee that any history  $lh$  learned in step (Hybrid Learn) is a prefix of  $pfh$ .

We will now argue that this is the case. We know that  $lh$  is a prefix of  $fh$ . By the choice of  $pfh$ , all prefixes of  $fh$  which are compatible with  $h$  are also prefixes of  $pfh$ . Thus, it suffices to show that  $lh$  is compatible with  $h$ . Hybrid learning implies that some acceptor has accepted  $lh$  as (classic) history. Thus,  $lh$  is a prefix of some history sent by the leader in a “2a” message. Clearly, this holds for  $h$  too. Any two histories sent by the leader (of the same round) are prefixes of each other. So, if  $max$  is the largest of the two histories, then  $max$  is a common extension of  $lh$  and  $h$ . Thus,  $lh$  and  $h$  are compatible.

#### 4.2.2 Implementation Considerations

Now that we have argued about the correctness (safety) of HP, in this section we describe how HP can be tweaked to be practical. We have identified a set of optimizations and listed them below.

**O1** The leader does not have to send the entire history  $h \bullet c$  in step (2aCP), it suffices to send  $c$ . When an acceptor receives  $c$  in a “2a” message, the FIFO property implies that it has already received  $h$ .

If the state machine is implemented at the servers, then there is no reason to send the entire history to the clients. All a client needs to learn is (a) the execution result of its last issued command and (b) that the history producing the result is chosen.

**O2** The solution to (a) is to have the servers speculatively execute commands. Specifically, when the leader receives a proposal from a client, it immediately executes the command and includes the result in the “2b” message it sends back to the client. Speculation at the leader avoids rollbacks and history replays during normal operation.

**O3** In order to attain (b) without sending the history, we replace the histories in the “2b” and “2bFast” messages with *history digests*. Two history digests are equal iff the corresponding histories are equal. Thus, in the (Learn) and (Hybrid Learn) step, clients check history digests for equality.

Intuitively, a history digest function takes as arguments a history  $h$  and a command  $c$  contained in that history. It then computes the smallest prefix of  $h$  containing  $c$  and returns the digest thereof. We refer the reader to the full paper [9] for an incremental digest implementation based on hashing.

**O4** If the classic and fast histories diverge during normal operation, the protocol as described above prevents hybrid learning. A simple solution would be to periodically start a new round. However, this imposes a considerable overhead. Therefore, the idea is to have each acceptor locally *align* the fast history  $fh$  to the classic history  $h$  as follows. Periodically,  $fh$  is replaced with the  $lub$  of  $h$  and the largest prefix of  $fh$  that is compatible with  $h$ . We know from Section 4.2.1 that this is safe.

**O5** We have optimized HP to adapt to a changing workload. Specifically, HP uses a double threshold ( $T_{high}, T_{low}$ ) such that when the load increases above  $T_{high}$  (resp. decreases below  $T_{low}$ ) the fast mode is switched off (resp. switched on). Changes in the load are monitored by the leader, who is counting proposals per time unit. The leader simply tells the clients (in “2b” messages) to stop (respectively start) sending commands to the acceptors.

### 4.3 Discussion

In this section, we provide a comparison of HP and FP in terms of their performance during normal operation. We argue that the cost of collision recovery in FP outweighs the gain obtained from fast learning. The original FP paper [19] says: “If collisions are very rare, then starting a new round might be best. If collisions are too frequent, then classic Paxos might be better than Fast Paxos.”

The latency of HP and FP equals two message delays in the absence of collisions. In the presence of collisions, HP requires three message delays and FP requires *six*. Hence, if FP is recovering from collisions only 25% of the time, then there is no (average) gain from fast learning.

The message complexity of HP is  $4n$  messages per request. FP requires  $3n + 1$  messages in the absence of collisions and  $(6 + l)n$  messages otherwise, where  $l$  is the number of learners. If we consider that  $l = n \geq 3$ , and FP is recovering from collisions only 12% of the time, then FP has a higher (average) message complexity. Even without collisions, in FP the leader collects “2bFast” messages and checks if collisions occurred, thus becoming the bottleneck.

Our experiments with HP have revealed that with increasing load, the collision rate is growing faster than the server capacity utilization rate. For instance, we have observed that the servers are still underutilized when the rate of hybrid learning drops under 50% (with 99% commutable commands). In this situation, FP would spend  $> 50\%$  of the time recovering from collisions, thus performing poorly compared with HP.

## 5 Evaluation

This section explores the performance characteristics of HP and compares it with existing approaches. As argued in Section 4.3 above, we expect HP to outperform FP in most situations, and therefore we omit a direct comparison. We substantiate our claim by showing that HP’s latency often attains the theoretical minimum. We compare HP with CP and show that it performs significantly better under low to medium load *and* equally well under high load. Where appropriate, we also compare HP with Mencius [21].

### 5.1 Experimental Settings

We have implemented a simple banking system in which multiple clients share a bank account. Clients can deposit or withdraw money. The state consists of the balance of the shared account, and clients can issue *withdraw* or *deposit* commands. Executing *withdraw* \$100 subtracts \$100 in a state with at least \$100 and produces \$100 as output. Executing *deposit* \$20 adds \$20 and produces OK as output. Note that any two *deposit* commands are commutable because executing them in either order has the same effect. However, when one of the two operations is a *withdraw*, the order matters.

A scenario is modeled in which clients frequently deposit small amounts of money and less frequently withdraw larger amounts. Where the rate of *withdraw* commands matters, we use “HP- $x$ ” to denote runs of the HP protocol, where on average, one out of  $x$  commands is a *withdraw*. We use “CP3” (respectively “CP4”) to denote the specific CP protocol where a command can be learned by a client after three (respectively four) message delays; CP4 relates to CP without speculation.

We ran experiments in the Emulab testbed [26] and we implemented all protocols in Java using the Neko [25] framework. The protocols are evaluated in a system with five servers ( $f = 2$ ) except for fault-scalability, where the number of servers is scaled up to 21 ( $f = 10$ ).

Client and server nodes are connected by links with a one-way delay of  $20ms$  and a bandwidth of 100 Mbps. The chosen delays are comparable to the “Europe” WAN setting analyzed in section 1. The chosen network bandwidth models modern high-end WAN links such Geant2 [13]. Server nodes are 600 Mhz PCs with 256 MB memory running Fedora 6.

### 5.2 Latency

Figure 4(a) shows the average latency of HP under low and medium load as the rate of *withdraw* operations is varied between 0% and 100%. Note that the withdraw rate corresponds to the probability of collisions. For a *withdraw*

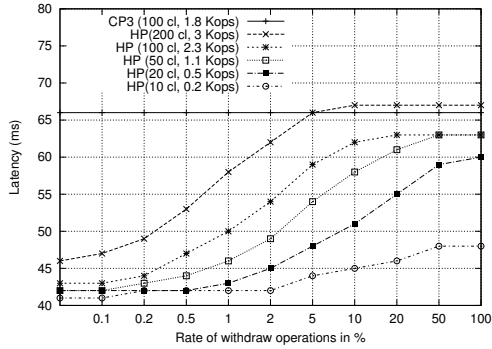
rate of 0.5% and load offered by 100 clients, HP has a 32% lower latency than CP3. This is close to the theoretical minimum. For a *withdraw* rate of 100% and load offered by 10 clients (between 0.1 and 0.2 Kops), HP still features a latency of 20% below the optimum of CP3.

Figure 4(b) compares the latency and throughput of HP-500, CP3, CP4 with and without request batching (of 20 commands) as we vary the offered load. As illustrated, under low load, batching at the leader increases the latency of CP4 and CP3 but not that of HP because most commands ( $> 90\%$ ) are chosen in the fast mode. On the other hand, batching increases peak throughput. In fact, with batching all protocols converge to the same peak throughput. Starting from a throughput of 6 Kops, the curves of HP and CP3 coincide because the fast mode is switched off.

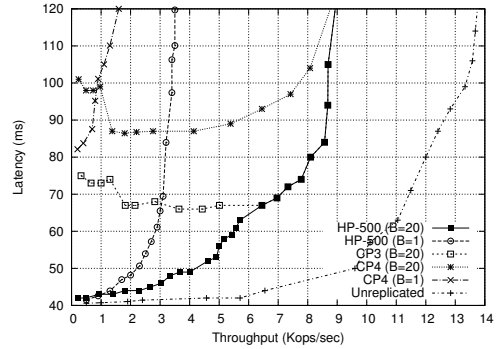
Figures 4(c) and 4(d) illustrate the effectiveness of adaptive switching by means of a dynamic workload. The workload is organized as follows: 50% of the commands are sent under moderate load generated by 100 clients and 50% under high load generated by 1000 clients between  $t = 68$  and  $t = 98$ . Figure 4(c) compares the average latency of HP with that of CP3 and CP4. We have measured the latency of HP with and without adaptive switching. The latter is referred to as nonadaptive. Figure 4(c) clearly shows that during the high load burst, the nonadaptive version of HP performs worst among all protocols. The explanation is the following. Batching offloads the leader and the acceptors become the bottleneck nodes because they process more messages. In contrast, the adaptive version shows a short spike after the load burst starts and a short tail after the burst ends. These can be attributed to conservative thresholds. Overall, the adaptive version of HP features the minimum latency of all protocols. Figure 4(d) compares the cumulative latency distribution of the four protocols under the same workload and confirms that adaptive HP performs best both under moderate and high load.

The experiments presented so far have been conducted in a setting in which the network is always timely. We now add a Pareto distribution to each link using the NetEm [15] utility. The one-way network delay now varies between  $20ms$  and  $60ms$ . Pareto is a heavy tailed distribution, which models the fact that wide-area links are usually timely (e.g. 80% of the time) but can present high latency occasionally.

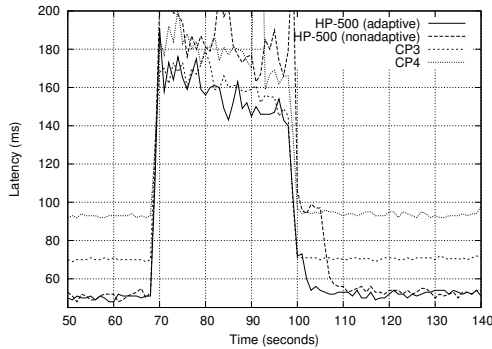
Figure 5(a) compares the latency and throughput of HP-500, CP3, CP4 with batching as we vary the offered load. The trends are the same as in a situation with no network variance. An important point is that all protocols have lower peak throughputs, including the unreplicated system. High variance results in packet reordering and packet retransmission at the transport protocol level (TCP), causing additional load in the bottleneck node. HP outperforms CP3 up to 60% of the peak throughput. Up to a throughput of 1 Kops, HP and the unreplicated system have compa-



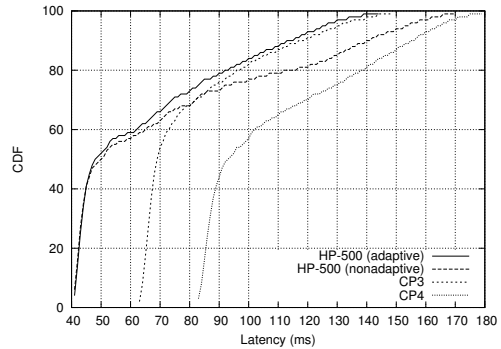
(a) Latency versus *withdraw* rate



(b) Latency versus throughput



(c) Average latency under a changing load ( $B = 20$ )



(d) Latency CDF under a changing load ( $B = 20$ )

Figure 4.

rable latencies. The performance profile of HP is somewhat surprising because with high network variance, the likelihood of collisions increases. E.g., under a link variance of  $40ms$ , if two interfering commands are sent within  $40ms$  from each other, they might be accepted in different orders.

Figure 5(b) supports the above observation showing that under network variance, the latency of HP converges much faster to that of CP3. Nevertheless, under low load and a small fraction of withdraws, HP shows a latency improvement of up to 40% over CP3, which is more than the theoretical maximum latency reduction of  $1/3$ . An explanation could be that the longer it takes to run an instance of a protocol, the more likely it is to depend on a slow link in the critical path. In this particular case, this effect adds to the latency of CP3 and explains the measured latency difference.

Figure 5(c) compares the latency of HP-500, CP3 and Mencius [21] as more servers are added to the system. We are simulating a lightly loaded scenario with 20 clients. With Mencius, a server has to wait for all other servers to skip or to propose a command. For a fair comparison, all commands are commutable and thus Mencius can commit in only one message roundtrip after receiving a reply from all servers, which is optimal. Mencius' dependency on slow links grows as more servers are added and therefore its latency increases. In contrast, the latency of HP and CP3 re-

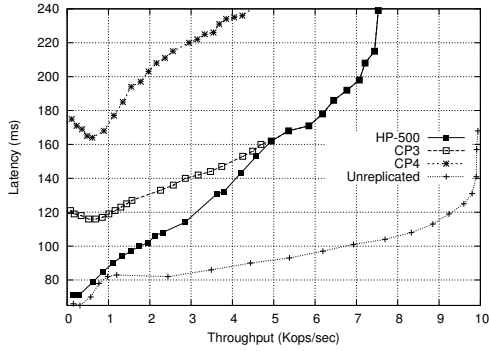
mains roughly constant (CP3's latency even drops) because they wait for the fastest quorum. These results suggest that the latency of CP and HP strongly depends on how large is the fraction of nodes that form a quorum. We observe that the latency oscillates (and even drops) with this fraction.

### 5.3 Throughput

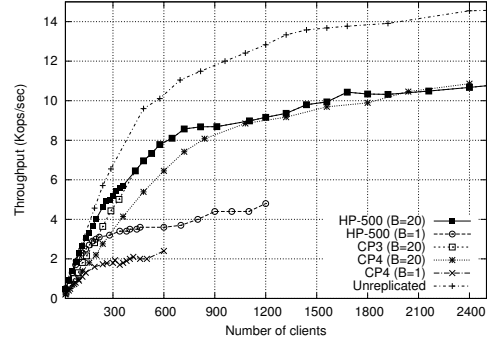
We show now that the lower latency of HP does not come at the cost of lower throughput compared to CP.

Figure 6(a) shows the throughput of HP-500, CP3 and CP4 with and without batching as the number of clients increases. All protocols scale equally well when batching is used; CP4 without batching scales poorly. Figure 6(b) compares the peak throughputs of HP (that equals CP3), CP4 and Mencius as the number of faulty servers tolerated increases. The throughput of Mencius is scaled down from 3GHz machines to ours (600MHz) using a factor of  $1/4$ . The results show that HP outperforms all other protocols except in the case of  $f = 2$  with batching, when its peak throughput is comparable to CP4. The fault scalability of HP is superior to that of CP4 with and without batching. For  $f = 10$  with batching, HP features 73% of the peak throughput for  $f = 2$ . In contrast, CP4's peak throughput drops down to 50%.

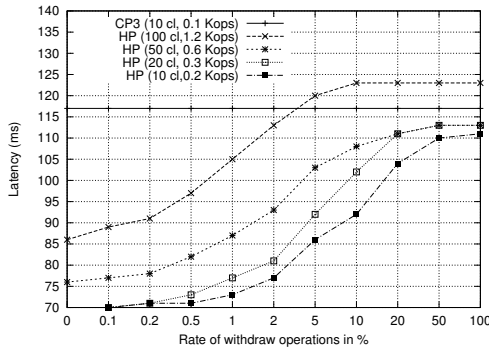




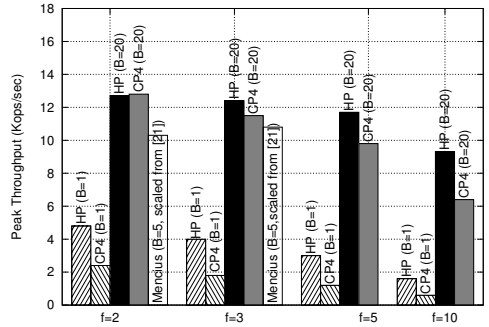
(a) Latency versus throughput ( $B = 20$ )



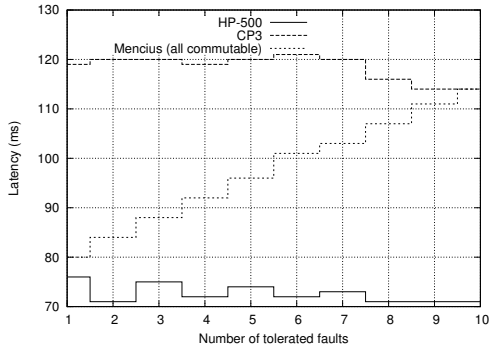
(a) Throughput as the number of clients increases.



(b) Latency vs. withdraw rate



(b) Throughput as  $f$  increases.



(c) Latency as  $f$  increases (20 clients)

**Figure 5. Effect of network variance**

**Figure 6.**

candidate for latency-critical applications because the read phase is done for infinitely many consensus instances together. Lamport's recent work on Fast Paxos [19] is based on the observation that the latency of CP can be further reduced if clients directly propose commands. At the heart of FP lies the idea of one-step consensus [4]. Pedone *et al.* [23, 24] have developed latency-efficient atomic broadcast algorithms based on the execution of a sequence of one-step consensus instances. The mentioned protocols suffer from collisions (which results in degraded latency) when multiple commands are sent at about the same time. In our prior work [10] we have developed consensus protocols tackling this problem and degrading gracefully in the presence of collisions. Guerraoui and Raynal [14] have developed a gracefully degrading consensus protocol that quickly chooses a value accepted by a fixed quorum. Charron-Bost and Schiper [8] present a consensus protocol with the minimum latency of FP and CP but only in failure-free runs.

Pedone and Schiper [22] have introduced the problem of generic broadcast, where one message delay can be saved by delivering messages in different but equivalent order. Recent work by Lamport on generalized consensus [18] borrows ideas from this work and extends FP to solve consensus on a growing set of partially ordered commands. Zielinski [27] proposes a protocol that combines FP and

## 6 Related Work

In the asynchronous model, the possibility of a single crash makes deterministic consensus impossible [12]. It has been shown that the FLP result can be circumvented by additional timing assumptions [11]. Chandra and Toueg [7] have introduced the concept of failure detectors, encapsulating timing assumptions, and  $\Omega$  [6] has been shown to be the weakest failure detector for solving consensus. In his seminal Paxos paper [17], Lamport describes how to build a replicated state machine from consensus. Unlike many other consensus algorithms, Paxos is a perfect

CP into a latency-optimal generic broadcast protocol. The protocol is not resilience optimal and incurs the expense of quadratic messages. Mencius [21] is a state machine replication protocol based on CP. The goal of Mencius is to reduce the load at the leader in order to prevent it from becoming a bottleneck. The main idea is to partition the set of consensus instances among several leaders. To reduce latency, Mencius assumes that each client has a local area connection with some leader. However, this cannot be guaranteed in a general system setting. Further, the system can only make progress if all leaders are correct. Therefore an eventually perfect failure detector [7] is assumed. These are all assumptions we do not make.

An earlier version of this work first introducing the concept of Hybrid Paxos is also discussed in some work by Junqueira *et al.* [16]. They study, by performing network simulations, when CP has a lower latency than FP.

Camargos *et al.*[5] developed a generalized consensus protocol that aims at improving the availability of CP by allowing multiple leaders to coexist. Under collisions, the protocol faces similar problems as FP and requires collision resolution. In contrast HP eliminates collision resolution, improving the latency and throughput of FP in stable runs.

## 7 Conclusion

We have developed Hybrid Paxos, a generalized consensus protocol that features minimal latency and maximum throughput in most situations. The core idea of HP is to add fast learning to CP. HP is to our knowledge the first generalized consensus protocol that attains the optimal latency of two message delays in the absence of collisions and three otherwise. Moreover, it has optimal latency, resilience and number of messages. We have shown that generalized consensus is a practical approach to replication in a WAN. Our experimental results demonstrate that HP can outperform state of the art protocols.

**Acknowledgements** We thank Flavio Junqueira for the inputs and discussions that helped to significantly enhance the paper.

## References

- [1] M. K. Aguilera, W. Chen, and S. Toueg. Failure detection and consensus in the crash-recovery model. In *Proc. of DISC*, pages 231–245, 1998.
- [2] A. Basu, B. Charron-Bost, and S. Toueg. Simulating reliable links with unreliable links in the presence of process crashes. In *Proc. of WDAG*, pages 105–122, 1996.
- [3] A. Bavier, M. Bowman, B. Chun, D. Culler, S. Karlin, S. Muir, L. Peterson, T. Roscoe, T. Spalink, and M. Wawrzoniak. Operating system support for planetary-scale network services. In *Proc. of NSDI*, pages 253–266, 2004. *et al.*
- [4] F. V. Brasileiro, F. Greve, A. Mostefaoui, and M. Raynal. Consensus in one communication step. In *Proc. of PACT*, pages 42–50, 2001.
- [5] L. J. Camargos, R. M. Schmidt, and F. Pedone. Multicoordinated agreement protocols for higher availability. *NCA*, pages 76–84, 2008.
- [6] T. D. Chandra, V. Hadzilacos, and S. Toueg. The Weakest Failure Detector for Solving Consensus. In *Proc. of PODC*, pages 147–158, 1992.
- [7] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *JACM*, (2):225–267, 1996.
- [8] B. Charron-Bost and A. Schiper. Improving fast paxos: being optimistic with no overhead. In *Proc. of PRDC*, 2006.
- [9] D. Dobre, M. Majuntke, M. Serafini, and N. Suri. Hp: Hybrid paxos for wans. 2009. <http://www.deeds.informatik.tu-darmstadt.de/dan/ghp.pdf>.
- [10] D. Dobre and N. Suri. One-step consensus with zero-degradation. In *Proc. of DSN*, pages 137–146, 2006.
- [11] C. Dwork, N. Lynch, and L. Stockmeyer. Consensus in the presence of partial synchrony. *JACM*, (2):288–323, 1988.
- [12] M. J. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of distributed consensus with one faulty process. *JACM*, (2):374–382, 1985.
- [13] Geant2. Pan-european backbone network. Website. <http://www.geant2.net>.
- [14] R. Guerraoui and M. Raynal. The information structure of indulgent consensus. *IEEE Trans. Comput.*, 53(4), 2004.
- [15] S. Hemminger. Network emulation with netem. In *Proc. of LCA*, 2005.
- [16] F. Junqueira, Y. Mao, and K. Marzullo. Classic paxos vs. fast paxos: caveat emptor. In *HotDep workshop*, Berkeley, CA, USA, 2007. USENIX Association.
- [17] L. Lamport. The part-time parliament. *ACM Trans. Computer Systems*, (2):133–169, 1998.
- [18] L. Lamport. Generalized consensus and paxos. In *MSR-TR-2005-33*, 2005.
- [19] L. Lamport. Fast paxos. *Distrib. Comp.*, 19(2), 2006.
- [20] L. Lamport. Lower bounds for asynchronous consensus. *Dist. Comp.*, 19(2), 2006.
- [21] Y. Mao, F. P. Junqueira, and K. Marzullo. Mencius: Building efficient replicated state machine for wans. In *OSDI*, 2008.
- [22] F. Pedone and A. Schiper. Handling message semantics with generic broadcast protocols. *Dist. Comp.*, 15(2), 2002.
- [23] F. Pedone and A. Schiper. Optimistic atomic broadcast: A pragmatic viewpoint. *Journal of Theoretical Computer Science*, 291(1):79–101, 2003.
- [24] F. Pedone, A. Schiper, P. Urbán, and D. Cavin. Solving agreement problems with weak ordering oracles. In *Proc. of EDCC*, pages 44–61, 2002.
- [25] P. Urbán, X. Defago, and A. Schiper. Neko: A single environment to simulate and prototype distributed algorithms. In *Proc. of Information Networking*, pages 503–511, 2001.
- [26] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. An integrated experimental environment for distributed systems and networks. In *Proc. of OSDI*, pages 255–270, 2003.
- [27] P. Zielinski. Optimistic generic broadcast. In *Proc. of DISC*, pages 369–383, 2005.