

USI Technical Report Series in Informatics

Optimistic Aborts for Geo-distributed Transactions

Theo Jepsen, Leandro Pacheco de Sousa, Huynh Tu Dang,
Fernando Pedone, Robert Soulé

Faculty of Informatics, Università della Svizzera italiana, Switzerland

Abstract

Network latency can have a significant impact on the performance of transactional storage systems, particularly in wide area or geo-distributed deployments. To reduce latency, systems typically rely on a cache to service read-requests closer to the client. However, caches are not effective for write-heavy workloads, which have to be processed by the storage system in order to maintain serializability.

This paper presents a new technique, called *optimistic abort*, which reduces network latency for high-contention, write-heavy workloads by identifying transactions that will abort as early as possible, and aborting them before they reach the store. We have implemented *optimistic abort* in a system called Gotthard, which leverages recent advances in network data plane programmability to execute transaction processing logic directly in network devices. Gotthard examines network traffic to observe and log transaction requests. If Gotthard suspects that a transaction is likely to be aborted at the store, it aborts the transaction early by re-writing the packet header, and routing the packets back to the client. Gotthard significantly reduces the overall latency and improves the throughput for high-contention workloads.

Report Info

Published

Number

USI-INF-TR-2016-05

Institution

Faculty of Informatics
Università della Svizzera italiana
Lugano, Switzerland

Online Access

www.inf.usi.ch/techreports

Introduction

Many distributed applications produce write-intensive workloads that access geographically-distributed storage systems. Examples of such applications span a range of domains, from popular web services offered by companies such as Google, Microsoft, and Amazon [12, 5, 15, 11, 19], to sensor-network based monitoring for military purposes [57] or environmental tracking [13, 56, 60].

For these applications, network latency can have a significant impact on performance. The standard approach for reducing network latency is to deploy proxy servers that service read operations closer to the client with cached data [21, 1, 40]. Unfortunately, although caching works well for read-heavy workloads [33], it provides little benefit for write-heavy workloads, since write operations must be routed directly to the store to ensure serializability (§2).

This paper presents a new technique, called *optimistic abort*, which reduces network latency for high-contention, write-heavy workloads. The key idea is to identify transactions that will abort as early as possible, and abort them before they reach the store. Optimistic abort provides a complimentary technique to caching; rather than moving data closer to the client, it *aborts transactions closer to the client*.

The logic of optimistic abort could be realized in various ways. For example, similar to a read cache, a proxy server could be used to identify and abort transactions. In this paper, we explore a highly efficient, switch-based implementation that leverages recent advances in data plane programmability [53, 9, 7]. Specifically, our system, is implemented in the P4 language [7], allowing it to use the emerging ecosystem of compilers and tools to run on reconfigurable network hardware [58, 44, 61, 46].

With the resulting system, named Gotthard, network devices naturally serve as message aggregation points for multiple clients. Clients issue transaction requests using a custom network protocol header. Gotthard switches examine network traffic to observe and log transaction requests. If Gotthard suspects that the store is likely to abort a transaction, the Gotthard switch proactively aborts the transaction by re-writing the packet header, and forwarding the packet back to the client. Since Gotthard aborts the transaction before the packet travels the network to the storage system, it can significantly reduce the overall latency for processing the request. Moreover, Gotthard’s optimistic strategy reduces commit time, improving overall system throughput.

We have implemented a prototype of Gotthard, and evaluated it both in emulation using Mininet [37], and on Amazon’s EC2 cloud infrastructure. Our experiments include both a set of microbenchmarks that explore the parameter space for different operational conditions, and an implementation of TPC-C [16] to emulate a real-world inspired workload. Our evaluation shows that Gotthard significantly reduces the latency for transactions, and under high-contention workloads, and increases transaction throughput.

As a highlight of some of results, our evaluation shows that when deployed on regions in the western and eastern United States in EC2, Gotthard almost doubles the throughput for TPC-C Payment transactions.

Overall, this paper makes the following contributions:

- It presents a novel algorithm for improving the performance of transaction systems by optimistically aborting transactions before they reach the store.
- It describes an implementation of the optimistic abort technique that builds on emerging technological trends in data plane programming languages.
- It explores the parameter space for network-based transaction processing, and demonstrates significant performance improvements for realistic workloads.

The rest of this paper is organized as follows. We first motivate Gotthard with an example experiment (§2). We then present the design of the Gotthard system in detail (§3). Next, we describe the implementation (§4) and present a thorough evaluation (§5). Finally, we discuss related work (§6), and conclude (§7).

Background and Motivation

Before detailing the design of Gotthard, we briefly describe opportunities due to technological trends, the system model, and present an experiment that motivates moving transaction processing logic into network devices.

Opportunity

Recently, the landscape for network computing hardware has changed dramatically. Several devices are on the horizon that offer flexible hardware with customizable packet processing pipelines, including Protocol Independent Switch Architecture (PISA) chips from Barefoot networks [8], FlexPipe from Intel [27], NFP-6xxx from Netronome [39], and Xpliant from Cavium [62].

Importantly, this hardware is not simply more powerful; it is also more programmable. Several vendors and consortia have proposed domain specific languages that target these devices [53, 9, 7]. These languages allow users to customize their network hardware with new applications and services. Particularly relevant to Gotthard is the P4 language [7]. P4 provides high-level abstractions that are tailored directly to the needs of network forwarding devices: packets are processed by a sequence of tables; tables match header fields, and perform actions that forward, drop, or modify packets. Moreover, P4 allows for stateful operations that can read and write to memory cells called registers.

As a result of these changes, networks now offer a programmable substrate that was not previously available to system designers. Consequently, there is an opportunity to re-visit the traditional separation of application-layer and transport layer logic advocated by the end-to-end principle [50]. In the case of Gotthard, we leverage the programmable network substrate to improve performance for write-heavy workloads.

System Model

We consider a geographically distributed system composed of *client* and *server* processes. Processes communicate through message passing and do not have access to a shared memory. The system is asynchronous: there is no bound on messages delays and on relative process speeds. Processes are subject to crash failures and do not behave maliciously (e.g., no Byzantine failures).

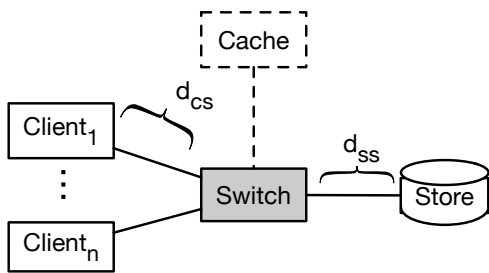


Figure 1: A typical cache would run on a server connected to a switch. For our experiments, we have implemented an idealized cache deployment that runs directly in the switch itself, thus reducing network latency.

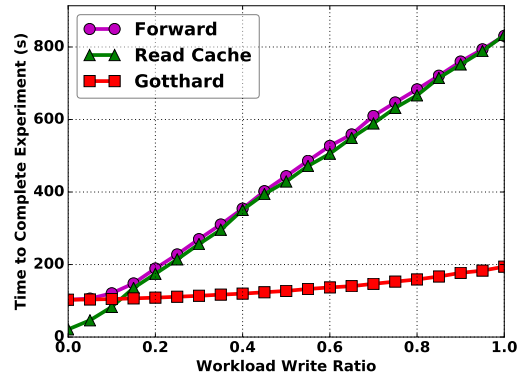


Figure 2: Time to complete 1000 transactions as the workload becomes more write-heavy. The read cache is ineffective as the percentage of write transactions exceeds 15%.

The store contains a set $D = \{x_1, x_2, \dots\}$ of data items. Each data item x is a tuple $\langle k, v \rangle$, where k is a key and v a value. We assume that the store exposes an interface with two operations: $read(k)$ returns the value of a given k , and $write(k, v)$ sets the value of key k to value v . We refer to those transactions that contain only read operations as *read transactions*. Transactions that contain at least one write operation are called *write transactions*.

We assume that clients execute transactions locally and then submit the transaction to the store to be committed. When executing a transaction, the client may read values from their own local cache. Write operations are buffered until commit time.

The isolation property that the system provides is *one-copy serializability*: every concurrent execution of committed transactions is equivalent to a serial execution involving the same transactions [6].

To ensure consistency, the store implements optimistic concurrency control (OCC) [31]. All *read transactions* are served directly by the store. To commit a *write transaction*, the client submits its buffered writes together with all values that it has read. The store only commits a transaction if all values in the submitted transaction are still current. As a mechanism for implementing this check, the system uses a $compare(k, v)$ operation, which returns true if the current value of k is v , and false otherwise. Note that the compare operation is not exposed to the users, but is simply used by the system to implement optimistic concurrency control. In the event of an abort, the server returns updated values, allowing the client to immediately re-execute the transaction.

Motivation

Figure 1 shows a typical deployment for a distributed transactional storage system. Clients submit transactions to a store by sending messages through the network.

The illustration shows a single switch in the network to simplify the presentation, but in practice, the network topology can be large and quite complex. If the distance between the client and the store is large, then the latency incurred by transmitting packets over the network can significantly impact the system performance. To reduce this overhead, systems typically deploy a cache [21, 1]. The cache stores copies of some subset of the data at the store. Client messages are routed to the cache. If the cache has a copy of the requested data, then it can respond to the clients request. Otherwise, it forwards messages to the store.

As long as the cache is closer to the client than the store, then caching is an effective technique for workloads that contain mostly read requests. However, caching does little to help improve performance for write-heavy workloads, since write operations need to be forwarded to the store in order to maintain strong consistency.

To demonstrate this behavior, we performed a simple experiment in which we measured the total time to complete 1,000 transactions for increasingly write-heavy workloads. In the experiment, a client submits two types of transactions to the store: one with a single read operation, and one with both a read and write operation. We varied the proportion of the two transaction types, to create workload scenarios ranging from more read-intensive to more write-intensive.

All the transactions passed through a switch that operated in one of three different modes of execution. In the

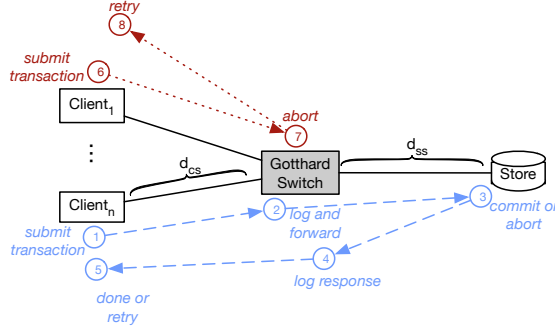


Figure 3: Overview of Gotthard deployment.

first, the switch acted traditionally, and simply forwarded requests to the store. In the second mode, the switch was modified using a data plane P4 program to behave like a cache. In Figure 1, the proxy cache is illustrated with a dashed-line to indicate that in our experiment, the cache operations are actually executed by the switch. This is therefore an “ideal” cache, which has no network latency between the switch and the proxy server hosting the cache. In the third configuration, which will be explained in Section 3, the switch executed Gotthard logic.

The experiment was run in emulation using the hardware described in Section 5. Figure 2 shows the results. We see that as the percentage of write requests increases, the cache becomes less effective. In fact, when the percentage of write requests is above 15%, we see almost no benefit to performance for using the read cache. This is because the client can read stale values from the cache which then cause write transactions to abort at the store.

As a preview of our results, Figure 2 also includes measurements for Gotthard which show how it improves performance for write-intensive workloads. The experiment demonstrates that Gotthard significantly reduces execution time with respect to simple forwarding and caching. When the workload is only at only 25% writes, Gotthard already reduces the completion time by half (i.e., a 2x improvement).

System Design

It is widely recognized that the performance of optimistic concurrency control protocols are heavily dependent on workload contention. As contention increases, so does the number of aborted transactions, causing the overall performance of the system to drop [23]. Gotthard is designed to address this problem.

In some respects, one can view Gotthard as providing complimentary behavior to read caches. For example, SwitchKV [33] forwards all transactions that contain only write operations to the store, and possibly services read operations at the cache. In contrast, Gotthard improves performance by focusing on write-intensive workloads.

However, although the approaches are complimentary, combining the two techniques is not necessarily straightforward. Because all read operations in Gotthard are serviced by the store, Gotthard enforces single-copy serializability (a proof-sketch of the correctness is in Section 3.7). Servicing read operations at a proxy server, or even a switch, would require additional mechanisms to ensure that the cache entries are valid.

Overview

Figure 3 shows a basic overview of Gotthard. From a high-level, clients submit transaction requests to the store. The requests pass through a Gotthard switch, which either forwards the request on to the store, or aborts the transaction, and responds to the client directly. In the figure, the two cases are distinguished by color and line type.

The blue, dashed-line shows the forwarding case. When the client submits the request (1), the switch examines the transaction and logs the operations in its local cache (2). It then forwards the transaction message to the store, which can commit or abort the transaction (3). The store responds to the client with the outcome of the execution. The switch logs the result of the execution (4), and forwards the response to the client. If the client learns that the transaction was aborted, it can re-try. Otherwise, the transaction is complete (5).

The red, dotted-line shows the abort case. As before, the client submits a request (6), and the switch examines the transaction message. When logging the operations, if the switch sees that a transaction is likely to abort based

Algorithm 1: Store Transaction Processing

```
1 Upon receiving TXN(compares, reads, writes):
2   corrections  $\leftarrow \emptyset$ 
3   foreach (key, value)  $\in$  compares do
4     if value  $\neq$  store[key] then
5       | corrections  $\leftarrow$  corrections  $\cup$  (key, store[key])
6   if corrections  $\neq \emptyset$  then respond Abort(corrections)
7   foreach (key, value)  $\in$  writes do
8     | store[key]  $\leftarrow$  value
9   response  $\leftarrow$  writes
10  foreach key  $\in$  reads do
11    | response  $\leftarrow$  response  $\cup$  (key, store[key])
12  respond OK(response)
```

on some previously seen transaction, the switch will preemptively abort the request (7), and send a response to the client. The client can then re-submit the transaction (8).

For a transaction that would have aborted at the store, the intuitive advantage of the Gotthard approach is clear: the message avoids traveling the distance from the switch to the store, d_{ss} , twice. However, within this basic framework, there are a number of subtle design decisions that impact the performance. We discuss these in more detail below.

Note that in the following discussion, we refer to the state that Gotthard uses to log transactions as a cache. However, to be clear, Gotthard only uses this state to make decisions about aborting transactions. It does not service read requests from its cache. As we discuss in the next section, the values that Gotthard keeps in its cache may be stale or incorrect.

In most of the following text, we refer to a basic deployment with a single switch to make exposition simpler. In Section 3.6, we discuss more complex topologies.

Data Store

Algorithm 1 shows the logic executed by the store. A transaction request message (TXN) contains three possibly empty lists of operations: *compares*, *reads*, and *writes*. The store first iterates over the compare operations to check for stale values (i.e., the value of the compare operation does not match the current value in the store). If any comparison operation fails, the store aborts the transaction (line 6). As part of the abort response, the store includes a list of correct values (*corrections*), with the updated values for comparisons that caused the transaction to fail. Otherwise, the store applies any write operations the transaction may include to update the values of the data items (line 8). Then the store responds to the client with all the values that were updated, along with the values that the transaction reads (line 11).

Optimistic Abort

A Gotthard switch aggregates messages from several clients. The switch logs requests and responses in order to determine if a subsequent transaction is likely to abort. Gotthard adopts an aggressive strategy for aborting transactions. It proactively updates its cache with the latest value after the switch has seen a transaction request (step 2 in Figure 3).

We refer to this as an *optimistic abort* strategy. It is optimistic because the switch assumes that any transaction request that is has seen is likely to be committed. As a result, it can make decisions about aborting subsequent transactions sooner. However, this approach may abort transactions that would not have been aborted by the store.

The logic for early optimistic abort is presented in Algorithm 2. The switch has logic for processing both transaction requests (line 1) and their responses (line 15). When the switch receives a transaction from the client, it iterates over the list of compare operations. If any compare operation references a key that is not in the cache (line 4), then the switch cannot reason about the validity of the transaction, so it forwards the request to the store (line 5). If the compare operation references a key that is in the cache, then the switch compares the value in the packet, with that in the store. If the values differ, the current value in the cache is added to a per-transaction set of

Algorithm 2: Gotthard Switch Logic

```
1 Upon receiving TXN(compares, reads, writes) from client:
2   corrections  $\leftarrow \emptyset$ 
3   foreach (key, value)  $\in$  compares do
4     if key  $\notin$  cache then
5       | forward to store
6     else if value  $\neq$  cache[key] then
7       | corrections  $\leftarrow$  corrections  $\cup$  (key, cache[key])
8   if corrections  $\neq \emptyset$  then
9     | respond Abort(corrections)
10  else if reads  $\cup$  writes  $\neq \emptyset$  then
11    foreach (key, value)  $\in$  writes do
12      | cache[key]  $\leftarrow$  value
13    forward to store
14  else respond OK()
15 Upon receiving Abort(corrections) from store:
16  foreach (key, value)  $\in$  corrections do
17    | cache[key]  $\leftarrow$  value
18  forward to client
```

corrections (line 7). After processing the compares, if there is at least one correction (i.e. there was a comparison that failed), the switch immediately sends an abort response to the client with the list of corrections (line 9).

If no comparisons failed, then the switch checks if the request contains *read* or *write* operations. If there are write operations, then it updates its cache with the new values (line 12). Then, it forwards the transaction to the store for processing. (line 13). Otherwise, if the transaction only included compares, the switch returns a successful response immediately to the client (line 14).

An abort message from the store (ABORT) contains a non-empty list of the corrections. The corrections contain the updated values that caused the transaction to fail. When the switch receives an abort message (line 15), it updates its cache with the correct values (line 17). If the switch did not update its cache on aborts, it would incorrectly abort subsequent transactions.

We note that although a Gotthard switch needs to maintain state in its local cache, the size of the cache does not need to be too large to be effective. It is sufficient to reserve enough space to keep for “hot” data items. The amount of space available for the cache will depend on the target platform for deployment. If the size is restricted, Gotthard could use a least-recently used cache eviction policy to make space available for new items.

Version Numbers

In a typical transactional storage system, data items would include a version number. In other words, each data item x would be a tuple $\langle k, v, ts \rangle$, where k is a key, v its value, and ts its version. The store would use the version numbers to determine if a transaction should be aborted due to a stale value. Notably, with this design, the store is responsible for assigning version number to data items.

However, with the optimistic abort strategy, the switch must update its local cache of data items before the store can assign a version number. Therefore, Gotthard cannot use version numbers to check for stale values. Instead, both the switch and the store compare with the current value to see if the data item for a specific key has changed.

Gotthard Messages

All transactions and transaction responses are encoded in a custom Gotthard packet header. Gotthard headers are encapsulated inside UDP transport-layer protocol headers. Prior work has shown that UDP is suitable for high-performance key-store systems [40]. If the set of operations in a transaction exceed the maximum transmission unit (MTU), Gotthard transactions may be split across multiple packets.

The header is shown in P4 syntax in Figure 5. The transaction header, `gotthard_hdr_t`, contains eight fields:

- `msg_type` indicates if the message is a request to or a response from the store.
- `from_switch` is a flag that indicates if an abort was performed by the switch or the store. This is currently used in our experiments for statistical purposes.
- `txn_id` is a unique identifier for each transaction.
- `frag_count` is the total number of fragments, and is used if a set of operations exceeds the packet maximum transmission unit (MTU).
- `frag_seq` identifies which fragment out of the total this packet contains.
- `status` indicates a commit or abort.
- `op_cnt` is the number of operations in the transaction.

The transaction header is followed by an array of operations. Each operation is defined in P4 as a separate header, `gotthard_op_t`. The `op_type` indicates the type of operation (compare, read, or write) and the `key` and `value` are the operation operands. For read operations, the value operand is unused. For compare operations, the value operand contains the value to compare.

An alternative implementation could have used different headers for transaction requests and responses. We chose this implementation because it simplified the P4 logic for aborting a message. The switch simply needs to overwrite the compare fields with the correct values, and change the `msg_type`. In contrast, removing and then appending a new header would have been an expensive operation in the software implementation of the P4 switch.

Expected Deployment

As our evaluation in Section 5 will show, Gotthard is most effective when the Gotthard switch is (i) deployed close to clients, and (ii) when the workload exhibits a high-degree of write contention.

An ideal deployment for Gotthard is one in which client locality correlates with write contention. For example, in Figure 4, we assume that clients $1 \dots n$ frequently access the same data items, and clients $(n + 1) \dots m$ access a different set of data items. Two switches, (Switch 1 and 2), are deployed close to their respective clients, and examine and potentially abort writes from the different sets of data items.

Such partitioning of data and clients occurs naturally in a number of applications. For example, if the clients are sensors, they are likely to read data from overlapping sets of items. In the TPC-C benchmarks that we used in our evaluation, data is naturally partitioned by clients accessing different warehouses that are geographically distributed.

We should note that if clients connected to different switches update data at the store, the correctness of Gotthard is not violated. Clients and switches will learn of new values after an abort message from the store. For example, if client_1 writes a value v_1 for key k_1 , then switch_1 will record v_1 in its cache. However, if client_{n+1} had previously written a value v'_1 for key k_1 , the request from client_1 will pass through switch_1 , but will be aborted by the store. The client and switch_1 will learn of the new value v'_1 in the abort response from the store. They will both then update their local caches, and the client can re-submit the transaction with the latest value.

Correctness

The storage (Algorithm 1) ensures serializable executions by means of optimistic concurrency control. The serialization order is defined by the order transactions arrive at the store. A transaction only commits if all the reads it performed during execution are still up to date at the time the transaction is received by the store.

The correctness of the switch logic (Algorithm 2) follows from the fact that (a) the switch does not commit any transactions, although it may abort transactions, and (b) the switch forwards non-aborted transactions to the store without changing their operations.

Fault Tolerance

It is often desirable for storage systems to provide fault tolerance. Fault-tolerance is an orthogonal problem to those that Gotthard addresses. However, the system could address failures by replicating the store using any of the standard consensus protocols [32, 42, 43, 14], and adding rules to the switch to route traffic to a replica in the event of failure (e.g., similar to OpenFlow's fast fail-over).

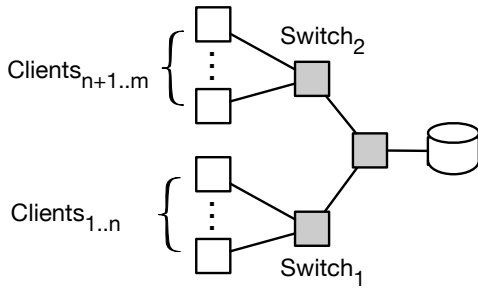


Figure 4: Gotthard is most effective when clients that connect to the same switch exhibit a high-degree of write contention.

```

1 header_type gotthard_hdr_t {
2   fields {
3     bit<1> msg_type; // REQ or RES
4     bit<1> from_switch;
5     bit<32> txn_id;
6     bit<8> frag_seq;
7     bit<8> frag_cnt;
8     bit<8> status;
9     bit<8> op_cnt;
10    // followed by op_cnt gotthard_op_t
        headers
11  }
12 }
13
14 header_type gotthard_op_t {
15   fields {
16     bit<8> op_type;
17     bit<32> key;
18     bit<1024> value;
19   }
20 }

```

Figure 5: Gotthard packet header in P4 syntax.

Implementation

We have implemented a prototype of Gotthard. The switch logic is written in the P4 programming language [7], version 1.10 [48] (i.e., P4-14), although porting to P4-16 should be straightforward. We used P4.org’s reference compiler [47] targeting the Behavioral Model switch [45]. The client and store code is written as standalone Python programs, using Python 2.7.

To be sent over Ethernet links, the size of the Gotthard message cannot exceed the Ethernet MTU of 1,500 bytes, otherwise the message will be fragmented at the network level. If a transaction has more than 10 operations, the message would exceed the MTU, so we fragment the Gotthard message at the application level using the header fields described in Section 3.5. When the Gotthard switch receives a fragmented message, it does execute any of Algorithm 2, but just forwards the message. This is because it is not possible to reason about the validity of a transaction without seeing it in its entirety.

All code, including the switch, client, store, simulation and experiment harness is publicly available with an open-source license.¹

Evaluation

In this section, we describe two sets of experiments that evaluate Gotthard. The first set of experiments are microbenchmarks that explore how Gotthard impacts the performance for executing transactions on a key-value store under various operational conditions, including relative distance of the switch to clients and the store, workload contention, and data locality on a multiple switch setup. In the second set of experiments, we explore the effectiveness of Gotthard for improving the performance for real-world inspired workloads, by using the TPC-C [16] benchmark with the parameters adjusted to increase contention. We ran our experiments using both a network emulator for a controllable environment, and on Amazon’s EC2 [2] for a more realistic setting using virtual machines in the cloud. Overall, the results demonstrate that Gotthard significantly reduces latency and improves throughput for workloads with high contention.

Experimental Setup

For all emulator-based experiments, we used Mininet [37] to create a virtual network. All experiments ran on a server machine with 12 cores (dual-socket Intel Xeon E5-2603 CPUs @ 1.6GHz) and 16GB of 1600MHz DDR4 memory. The operating system was Ubuntu 14.04.

We ran all cloud based experiments using Amazon’s EC2 [2], deploying on EC2 instances of type `t2.medium`. The clients, switch, and store were placed in individual instances, located in geographically distributed regions. Specifically, the clients, switch and store were placed in regions `us-west-2`, `us-west-1` and

¹<https://github.com/usi-systems/gotthard>

| | VA | EU | OR | JP |
|----|-------|---------|---------|---------|
| CA | 70-79 | 143-148 | 20-22 | 102-113 |
| VA | - | 70-76 | 66-85 | 140-153 |
| EU | - | - | 126-131 | 209-218 |
| OR | - | - | - | 89-100 |

Table 1: RTT between AWS regions (milliseconds)

| Parameter | Default |
|-----------|---------|
| RTT | 100 ms |
| delta | 0.2 |
| #clients | 8 |
| %writes | 0.2 |
| locality | 1.0 |

Table 2: Default parameter values in microbenchmarks.

us-east-1, respectively. We observed RTT latencies of 20-25ms between the clients and the switch, and 60-82ms between the switch and the store.

In both the emulator and EC2 experiments, we used the P4 Behavioral Model switch [45] to run the Gotthard P4 program. The same Python code was used in both deployments for the client and store.

In all experiments, we compare the performance of three different approaches:

- *Forwarding.* In this deployment, the switch simply forwards transaction messages to the intended destination, without executing any logic. This approach represents a baseline deployment.
- *Read Cache.* This deployment represents an idealized on-path look-through read cache deployment in which the switch itself caches values, and services read requests if the data is in the cache.
- *Gotthard.* The switch that executes the transaction logic described in this paper.

Except for the multi-switch microbenchmark, we used a topology in which all messages to and from the store are routed through a single switch, as shown in Figure 1.

Calibration

Gotthard is designed to reduce the overhead due to network latency for transaction processing. Network latency is a particular concern for wide area and geo-distributed deployments. In order to calibrate the expected network latencies in our simulations, we first measure round-trip time (RTT) latencies one might expect for a global deployment on EC2.

EC2 is hosted in multiple global regions. Each region corresponds to a geographic location and is divided into multiple availability zones. We started a `m3.large` VM instance in each availability zone, and ran the ping client to measure RTT to each other running instance. Measurements were collected for 2 minutes. Table 1 reports the minimum and maximum average for each region. The intra-region latencies were between 0.8 and 1.2 milliseconds, except for the JP region, where it was 2.5 milliseconds.

These measurements identify the range of network latencies one can expect in geo-distributed deployments. In all microbenchmark experiments in emulation, we used a fixed RTT of 100 milliseconds (i.e., 50 milliseconds one way). This number is well within the range of latencies one would expect on EC2 and is close the latencies observed in our actual EC2 experiments.

Microbenchmarks

To understand the effect of specific operational conditions on the performance of Gotthard, we designed a set of microbenchmarks, each focusing on one the following parameters:

- *Delta ratio.* The delta ratio captures the relative distance of the switch to clients and the store. Using the variables in Figure 3, the delta ratio is defined as $d_{cs}/d_{cs} + d_{ss}$. So, a smaller ratio means that the switch is closer to the client. Note that $d_{cs} + d_{ss} = \frac{RTT}{2} = 50\text{ms}$.
- *Percentage of writes.* Each client in the benchmarks issues a mix of *read* and *write* transactions. Percentage of writes dictates how many of the total number of transactions are writes.
- *Number of clients.* The number of concurrent clients submitting transactions.

- **Contention.** To characterize contention, we model the popularity of keys in the store using a Zipf distribution. Contention is increased by increasing the Zipf exponent; when the exponent is 0, all keys are equally likely to be accessed by clients.
- **Locality on a multi-switch deployment.** In a multi-switch setup such as the one described in Figure 4, it is possible that the optimistic execution of Gotthard results in false positives; transactions built upon optimistic state are aborted at the store. How likely that is to happen will depend on the workload. *Locality* describes the probability of clients accessing keys “owned” by another switch.

In the microbenchmarks, each transaction accesses a single key, representing a counter. A *read* transaction reads the current counter value and a *write* transaction atomically increments a counter.

For each experiment, unless otherwise specified, we used the default parameter settings in Table 2 and the topology in Figure 3. We deliberately chose a lower value of 0.2 %writes to better understand the effect of other parameters on performance, but as Figure 7a shows, Gotthard’s relative performance improves with higher %writes. For each data point, we ran the experiment for three minutes and report throughput in transactions committed per second, and the average latency of committed transactions.

Delta ratio. In this experiment, we quantify the effect of the relative distance of the switch to clients and the store: the *delta ratio*. As the ratio increases, the switch gets further from the client and closer to the store. Figure 6a reports overall throughput as the delta ratio changes from 0 to 1. As expected, as the switch get closer to the client (i.e., delta approaches 0), the overall throughput for Gotthard improves because transactions can be aborted and restarted earlier. Even though the performance of the Read Cache also improves slightly with a smaller delta ratio, its overall throughput is limited by write transaction throughput, since conflicting writes must reach the store to abort. Figure 6b reports average latency. Gotthard’s latency in this scenario is always lower than the other approaches, and is dependent on the latency between client and switch. Under high contention, write latency dominates the overall latency for Read Cache, due to costly aborts.

Percentage of writes. Figures 7a and 7b show throughput and latency, respectively, as %writes changes. These figures clearly demonstrate the effect of write operations, which limit the overall performance of the system. With %writes at 0, Read Cache can always serve requests from the switch, but as the %writes grows, throughput becomes limited due to the high cost of aborting writes. Gotthard’s throughput, on the other hand, degrades slowly, since requests are aborted at the switch and clients can optimistically retry transactions sooner. At 0.2 %writes, Gotthard’s throughput is already 1.5x that of Read Cache, reaching 3.3x that of Read Cache at around 0.5 %writes.

Number of clients. Figure 8a reports throughput, and demonstrates how Gotthard allows a higher rate of transactions under contention, past the saturation point of the other two approaches, reaching more than 4x times the throughput of Read Cache. Using Gotthard, transactions can abort early and optimistically be retried with updated values, *before the previous conflicting transaction commits*. If the optimistic retries are not aborted at the store, more contending transactions can be executed concurrently. Figure 8b shows that Gotthard’s latency remains stable under increasing load, until it starts reaching saturation at around 24 clients. Figure 9 shows the CDF for latencies when load is fixed at 8 clients. The “fat” tail observed for Read Cache and Forward are due to write transactions that repeatedly abort due to contention.

Contention. This experiment characterizes the effect of contention on the performance of Gotthard. Clients submit transactions that can access one of 10 keys, and the popularity of each key is dictated by a Zipf distribution. Figure 10a shows how increasing the Zipf exponent affects throughput. With the Zipf exponent at 0, all keys are accessed with the same probability, and contention increases as the exponent increases. With low contention, Read Cache performs comparatively better than Gotthard. However, as contention increases, the number of aborts grows, limiting the performance of write transactions for Read Cache. Gotthard, on the other hand, by optimistically aborting and retrying transactions closer to the client, is virtually unaffected by the increase in contention.

In Figure 12, we try to characterize the relative performance of Gotthard against the Read Cache approach, under varying %writes and Zipf exponent. The heatmap shows the throughput of Gotthard normalized to that of the Read Cache. Gotthard’s relative throughput is higher the closer we get to the upper-right corner of the

heatmap. With low contention (Zipf exponent at 0), Gotthard has better throughput when %writes is higher than 0.4. As contention increases, Gotthard has better throughput even when %writes get as low as 0.15.

Locality. For the locality experiment, we change the topology of our setup to that illustrated by Figure 4. In this setup, we have two groups of 8 clients, with each group connected to a different switch. The third switch, placed in between the first two switches and the store, is virtually collocated with the store; there is no added latency between the middle switch and the store. We divide a total of 10 keys into two sets of 5 keys, each set being “local” to one of the group of clients. To control locality, we vary the number of clients on each group that accesses “remote” keys, that is, those of the opposite group. Thus, when locality is 1, every client accesses keys in the local set. At the other extreme, when locality is 0.5, half of the clients in each group access non-local keys. Furthermore, keys inside a set are chosen with probability given by a Zipf distribution of exponent 3.

In Figure 11 we can see that, as expected, locality has a marked effect on Gotthard’s performance. In this scenario, Gotthard does worse than the other two approaches when locality is lower than around 0.8. Still, as locality increases, so does Gotthard’s performance. Gotthard should be suitable for more sophisticated deployments as long as the expected workload exhibits some degree of locality.

Summary. Overall, the experiments show that Gotthard consistently improves both throughput and average latency in workloads exhibiting contention. Specifically, *the comparative improvement in performance provided by Gotthard increases as contention goes up.* In more elaborate setups involving multiple switches, we expect Gotthard to still provide performance benefits, as long as clients exhibit some degree of locality in their access patterns.

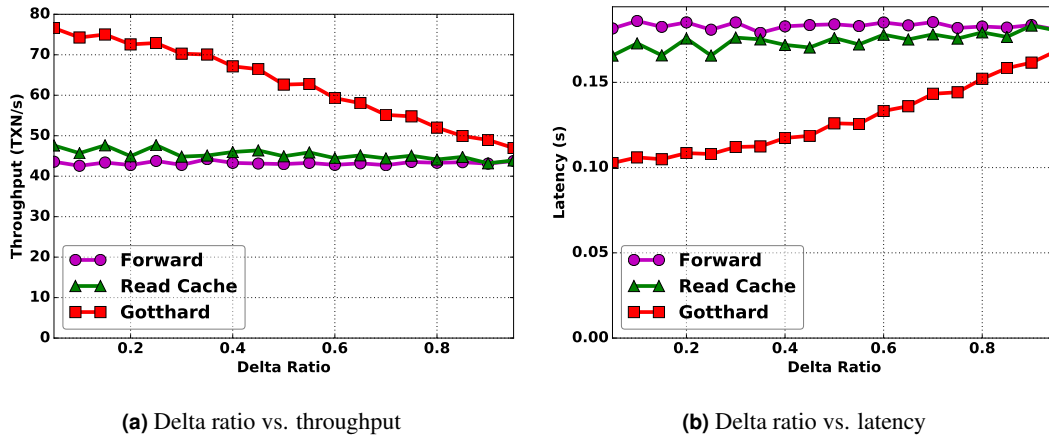


Figure 6: Performance as delta ratio increases.

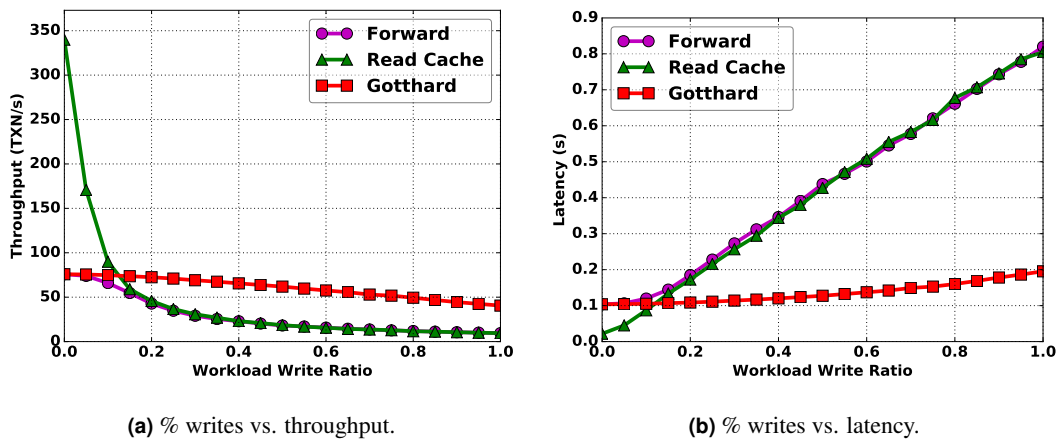
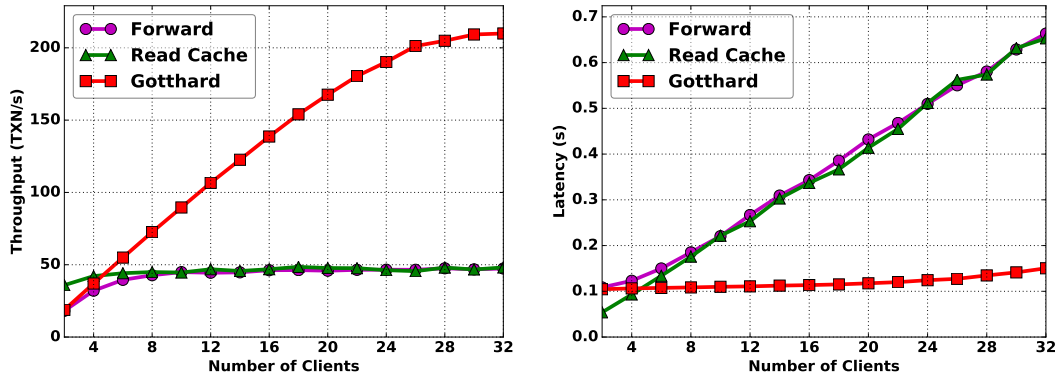


Figure 7: Performance as %writes changes.



(a) #clients vs. throughput.

(b) #clients vs. latency.

Figure 8: Performance as load (#clients) increases.

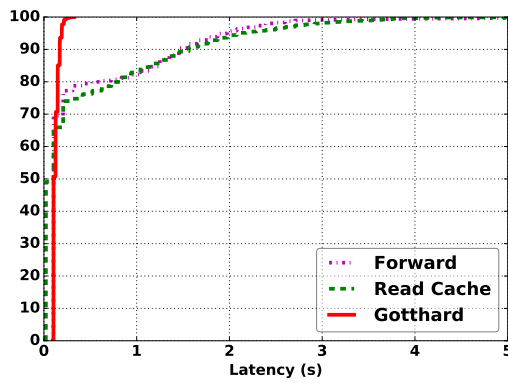
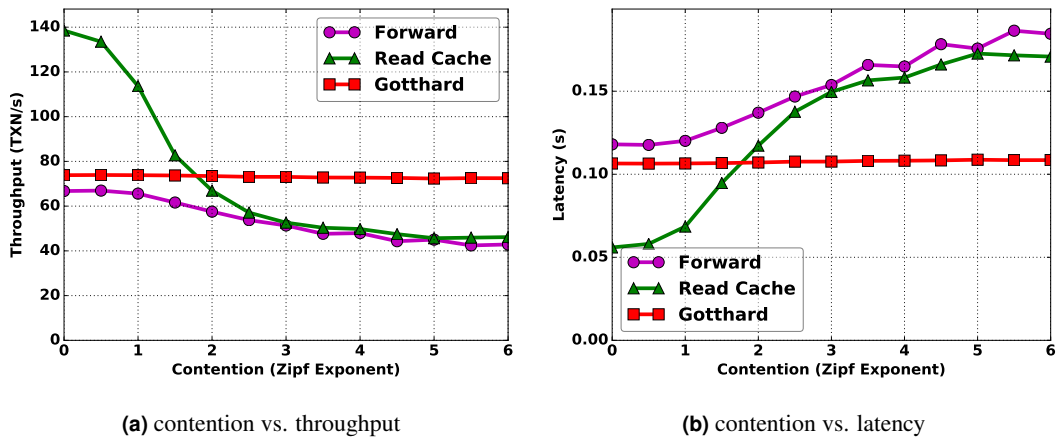


Figure 9: Latency CDF for fixed delta, %writes, clients.



(a) contention vs. throughput

(b) contention vs. latency

Figure 10: Performance with increasing contention.

TPC-C

TPC-C [16] is a widely-used benchmark for transaction processing systems. It models an online transaction processing (OLTP) workload for a fictional wholesale parts supplier that maintains a number of warehouses for different sales districts. Clients issue one of five different transactions against the store to: enter and deliver orders, record payments, check an order status, and monitor inventory at warehouses.

The TPC-C specification states that the benchmark should be run with a set of specific parameter values: 1 warehouse, 10 districts, 3000 customers and 100,000 items. For our evaluation, we differed from the standard parameters to use the following settings: 1 warehouse, 2 districts, 10 customers and 50 items. We altered these

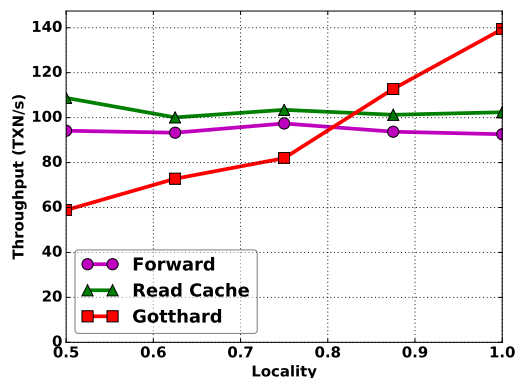


Figure 11: Effect of workload locality on throughput. (Multiple switch topology.)

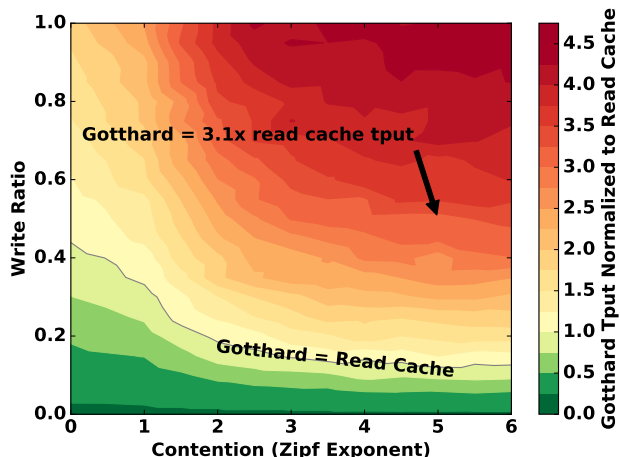


Figure 12: Gotthard’s throughput relative to Read Cache as contention and %writes change. Gotthard does better the closer it is to red (upper right corner).

parameters to increase the contention on the data store.

As a TPC-C driver, we modified an open source Python implementation from Pavlo et al. that was originally used to issue transactions against MongoDB.² Since the store does not have indexes, we modified the transaction executor to only perform exact selects. To represent the TPC-C database in our store, we map each record to a key-value pair. The client driver keeps a copy of all records it has read or written, essentially mirroring the store with a local cache. This eliminates the unnecessary latency of issuing read requests for records that have not changed. Consistency is guaranteed by validating the transaction before committing, ensuring the read values (possibly from the local cache) were up-to-date.

Complete benchmark. Figure 13 shows a summary of the throughput for the three switch approaches. The bars on the left show the average of all transactions, weighted by their frequency. The remaining bars show the throughput for individual transactions. The results on Mininet (Figure 13a) and on EC2 (Figure 13b) are similar. Both show that Gotthard improves throughput for Payment and Order Status. It is less effective for the other three transactions.

The reason that Gotthard is less effective for Delivery, New Order, and Stock Level is because when the client receives an abort response for one of these transactions, it must issue multiple read transactions before they can resubmit. These read transactions are due to the fact that records in these transactions have dependencies on records that are accessed by other transactions.

However, we see that on average, Gotthard improves the overall performance. In the case of the two transactions that can be re-submitted without additional reads, Gotthard significantly improves the throughput compared to the Read Cache (around 2x for Payment, and 1.1x for Order status). Below, we focus on these two transactions in detail.

Payment transaction. Figures 14a and 14e show the throughput and latency for the Payment transaction as we vary the delta ratio. When the switch is closer to the clients (lower delta ratio), Gotthard has lower latency and higher throughput than both the forwarding and Read Cache switches. As the switch is moved further from the clients (higher delta ratio), Gotthard’s performance is more similar to the other switches’.

Figures 14b and 14f show the throughput and latency for the Payment transaction as we vary the number of clients. Gotthard has higher throughput and lower latency as the number of clients increases. The other switches saturate at 8 clients, while Gotthard continues to scale.

Order Status transaction. Figures 14c and 14g show the throughput and latency for the Order Status transaction as we vary the delta ratio. Gotthard has higher throughput than the Read Cache until the switch is halfway

²<https://github.com/apavlo/py-tpcc>

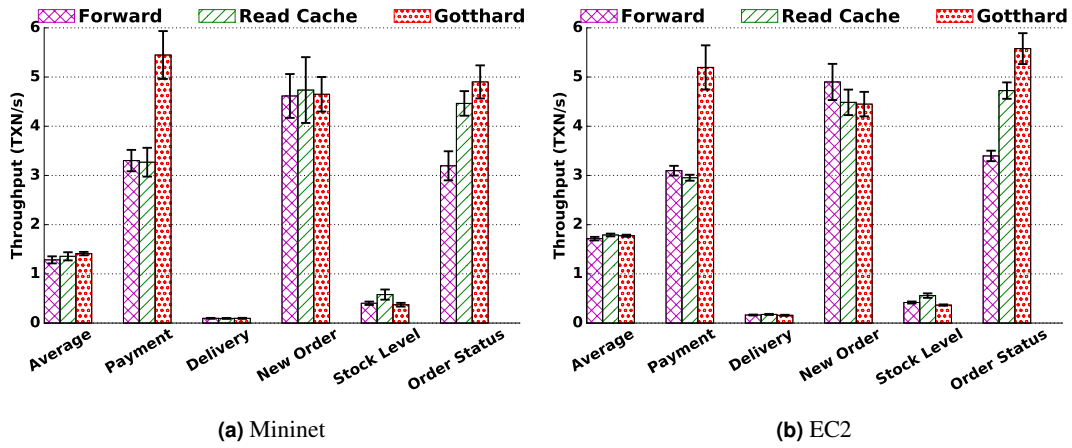


Figure 13: Summary of TPC-C transaction throughput. Gotthard improves the throughput of the Payment and Order Status transactions, which contain many write operations.

between the clients and the store (delta ratio 0.5). Gotthard always has lower latency than the others.

Figures 14d and 14h show the throughput and latency as the number of clients is increased. The other switches saturate at 4 clients, while Gotthard maintains higher throughput until 10 clients.

TPC-C Summary. Overall, the experiments show that Gotthard has higher throughput than the Read Cache for high contention transactions. Gotthard has the greatest benefit for the TPC-C workload when the switch is closer to the client, and still provides better performance when the switch is close to the store. As the number of clients increases, and thus the contention, Gotthard reaches a higher saturation point after the other switches.

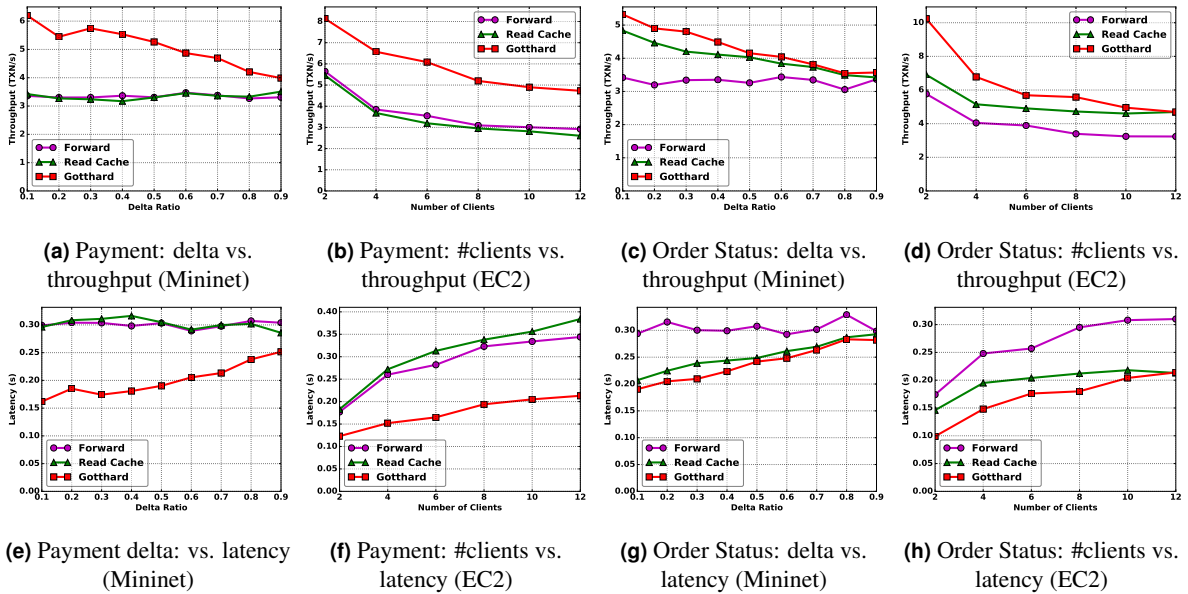


Figure 14: TPC-C Payment and Order Status transactions

Related Work

Proxies and caches The idea of using a proxy to extend distributed services is a well-established idea [52] that has been widely adopted [10, 3, 28, 29]. Proxies are often used to scale services by caching copies of data closer to clients, such as with content distribution networks (CDNs) [41, 21, 59]. CDNs typically are used for static content, although there are examples of proxies used for dynamic content [22].

Prior work has also explored the possibility of leveraging the network to route requests dynamically to proxies to service requests [59]. Notably, recent work on SwitchKV [33] uses OpenFlow-enabled switches to dynamically route read requests to proxy caches.

Gotthard differs from this work in that it is not a cache, *per se*. It keeps copies of transaction requests, but it does not service client read requests. Rather, it uses copies of previous requests to make informed decisions about when to abort transactions early, with the goal of reducing latency for write-heavy workloads.

Geo-distributed databases. Many recent works deal with geo-distribution in the context of transactional databases, most with a focus on replication and partitioning. Works such as [5, 15, 30, 38, 51] aim at providing strong consistency (i.e., serializability) over wide-area networks. To improve latency and availability, many works also propose weaker consistency criteria such as Parallel Snapshot Isolation [55], causal consistency [34, 35] and RAMP transactions [4]. PLANET [49] exposes transaction state to the application, enabling speculative processing and faster revocation using the *guess and apologies* paradigm [24]. Gotthard’s approach is complementary to the aforementioned research, and combining our approach with a full-fledged database solution is part of our future work.

Data plane programming languages. Gotthard is written in P4 [7], although there are several other projects have proposed domain-specific languages for data plane programming. Notable examples including Huawei’s POF [53] and Xilinx’s PX [9]. We chose to focus on P4 because there is a growing community of active users, and it is relatively more mature than the other choices. However, the ideas for implementing Gotthard should generalize to other languages.

Application logic in the network. Several recent projects investigate leveraging network programmability for improved application performance. One thread of research has focused on improving application performance through traffic management. Examples of such systems include PANE [20], EyeQ [26], and Merlin [54] which all use resource scheduling to improve job performance. NetAgg [36] leverages user-defined *combiner* functions to reduce network congestion. Another thread of research has focused on moving application logic into network devices. Dang et al. [18] proposed the idea of moving consensus logic in to network devices. Paxos Made Switch-y [17] describes an implementation of Paxos in P4. István et al. [25] implement Zookeeper’s atomic broadcast on an FPGA.

Conclusion

The advent of flexible hardware and expressive dataplane programming languages will have a profound impact on networks. One possible use of this emerging technology is to move logic traditionally associated with the application layer into the network itself.

With Gotthard, we have leveraged this technology to move transaction processing logic into network devices. Gotthard uses a custom packet header and programmable switches to identify doomed transactions, and abort them in the network as early as possible.

For write-intensive, high-contention workloads, significantly reduce the overall latency for processing transactions. Furthermore, Gotthard reduces load on the store, increasing system throughput.

Overall, Gotthard compliments prior techniques for reducing network latency, such as using a cache to service read requests. Moreover, Gotthard provides a novel application of data plane programming languages that advances the state-of-the-art in this emerging area of research.

References

- [1] Akamai. <https://www.akamai.com/>.
- [2] Amazon EC2. <aws.amazon.com/ec2>.
- [3] T. E. Anderson, M. D. Dahlin, J. M. Neefe, D. A. Patterson, D. S. Roselli, and R. Y. Wang. Serverless Network File Systems. *ACM Transactions on Computer Systems (TOCS)*, 14(1):41–79, Feb. 1996.
- [4] P. Bailis, A. Fekete, J. M. Hellerstein, A. Ghodsi, and I. Stoica. Scalable atomic visibility with ramp transactions. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 27–38. ACM, 2014.

- [5] J. Baker, C. Bond, J. C. Corbett, J. Furman, A. Khorlin, J. Larson, J.-M. Leon, Y. Li, A. Lloyd, and V. Yushprakh. Megastore: Providing scalable, highly available storage for interactive services. In *Conference on Innovative Data Systems Research (CIDR)*, pages 223–234, Jan. 2011.
- [6] P. A. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley Longman Publishing Co., Inc., 1987.
- [7] P. Bosshart, D. Daly, G. Gibb, M. Izzard, N. McKeown, J. Rexford, C. Schlesinger, D. Talayco, A. Vahdat, G. Varghese, and D. Walker. P4: Programming Protocol-Independent Packet Processors. *SIGCOMM Computer Communication Review (CCR)*, 44:87–95, July 2014.
- [8] P. Bosshart, G. Gibb, H.-S. Kim, G. Varghese, N. McKeown, M. Izzard, F. Mujica, and M. Horowitz. Forwarding Metamorphosis: Fast Programmable Match-Action Processing in Hardware for SDN. In *SIGCOMM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM)*, pages 99–110, Aug. 2013.
- [9] G. Brebner and W. Jiang. High-Speed Packet Processing using Reconfigurable Computing. *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 34:8–18, Jan. 2014.
- [10] E. A. Brewer, R. H. Katz, E. Amir, H. Balakrishnan, Y. Chawathe, A. Fox, S. D. Gribble, T. Hodes, G. Nguyen, V. N. Padmanabhan, M. Stemm, S. Seshan, and T. Henderson. A Network Architecture for Heterogeneous Mobile Computing. Technical report, University of California at Berkeley, 1998.
- [11] B. Calder, J. Wang, A. Ogus, N. Nilakantan, A. Skjolsvold, S. McKelvie, Y. Xu, S. Srivastav, J. Wu, H. Simitci, et al. Windows Azure Storage: a highly available cloud storage service with strong consistency. In *ACM Symposium on Operating Systems Principles (SOSP)*, pages 143–157, Oct. 2011.
- [12] D. G. Campbell, G. Kakivaya, and N. Ellis. Extreme Scale with Full SQL Language Support in Microsoft SQL Azure. In *ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 1021–1024, June 2010.
- [13] R. Cardell-Oliver. ROPE: A Reactive, Opportunistic Protocol for Environment Monitoring Sensor Networks. In *IEEE Workshop on Embedded Networked Sensors (EmNets)*, pages 63–70, May 2005.
- [14] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM (JACM)*, 43:225–267, Mar. 1996.
- [15] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, W. Hsieh, S. Kanthak, E. Kogan, H. Li, A. Lloyd, S. Melnik, D. Mwaura, D. Nagle, S. Quinlan, R. Rao, L. Rolig, Y. Saito, M. Szymaniak, C. Taylor, R. Wang, and D. Woodford. Spanner: Google’s Globally-Distributed Database. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 251–264, Oct. 2012.
- [16] T. Council. TPC-C Benchmark, revision 5.11, 2010.
- [17] H. T. Dang, M. Canini, F. Pedone, and R. Soulé. Paxos Made Switch-y. *SIGCOMM Computer Communication Review (CCR)*, 44:87–95, Apr. 2016.
- [18] H. T. Dang, D. Sciascia, M. Canini, F. Pedone, and R. Soulé. NetPaxos: Consensus at Network Speed. In *ACM SIGCOMM Symposium on SDN Research (SOSR)*, pages 59–73, June 2015.
- [19] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vosshall, and W. Vogels. Dynamo: Amazon’s Highly Available Key-value Store. In *ACM Symposium on Operating Systems Principles (SOSP)*, Oct. 2007.
- [20] A. Ferguson, A. Guha, C. Liang, R. Fonseca, and S. Krishnamurthi. Participatory Networking: An API for Application Control of SDNs. In *SIGCOMM Conference on Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM)*, pages 327–338, Aug. 2013.
- [21] M. J. Freedman, E. Freudenthal, and D. Mazières. Democratizing content publication with coral. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 18–18, Mar. 2004.
- [22] R. Grimm, G. Lichtman, N. Michalakakis, A. Elliston, A. Kravetz, J. Miller, and S. Raza. Na Kika: Secure Service Execution and Composition in an Open Edge-Side Computing Network. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 169–182, San Jose, California, May 2006.
- [23] R. E. Gruber. Optimism vs. locking: A study of concurrency control for client-server object-oriented databases. Technical report, Massachusetts Institute of Technology, 1997.
- [24] P. Helland and D. Campbell. Building on quicksand. In *Conference on Innovative Data Systems Research (CIDR)*, Jan. 2009.
- [25] Z. István, D. Sidler, G. Alonso, and M. Vukolić. Consensus in a Box: Inexpensive Coordination in Hardware. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 103–115, Mar. 2016.
- [26] V. Jeyakumar, M. Alizadeh, D. Mazières, B. Prabhakar, A. Greenberg, and C. Kim. EyeQ: Practical network performance isolation at the edge. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 297–312, Apr. 2013.

- [27] L. Jose, L. Yan, G. Varghese, and N. McKeown. Compiling Packet Programs to Reconfigurable Switches. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 103–115, May 2015.
- [28] A. D. Joseph, A. F. deLospinasse, J. A. Tauber, D. K. Gifford, and M. F. Kaashoek. Rover: A Toolkit for Mobile Information Access. In *ACM Symposium on Operating Systems Principles (SOSP)*, pages 156–171, Dec. 1995.
- [29] B. Knutsson, H. Lu, J. Mogul, and B. Hopkins. Architecture and Performance of Server-Directed Transcoding. *ACM Transactions on Internet Technology (TOIT)*, 3(4):392–424, Nov. 2003.
- [30] T. Kraska, G. Pang, M. J. Franklin, S. Madden, and A. Fekete. Mdcc: Multi-data center consistency. In *European Conference on Computer Systems*, pages 113–126. ACM, 2013.
- [31] H. T. Kung and J. T. Robinson. On optimistic methods for concurrency control. *ACM Transactions on Database Systems (TODS)*, 6(2):213–226, June 1981.
- [32] L. Lamport. The Part-Time Parliament. *ACM Transactions on Computer Systems (TOCS)*, 16:133–169, May 1998.
- [33] X. Li, R. Sethi, M. Kaminsky, D. G. Andersen, and M. J. Freedman. Be fast, cheap and in control with switchkv. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 31–44, Mar. 2016.
- [34] W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen. Don’t settle for eventual: scalable causal consistency for wide-area storage with cops. In *ACM Symposium on Operating Systems Principles (SOSP)*, pages 401–416. ACM, 2011.
- [35] W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen. Stronger semantics for low-latency geo-replicated storage. In *Presented as part of the 10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, pages 313–328, 2013.
- [36] L. Mai, L. Rupperecht, A. Alim, P. Costa, M. Migliavacca, P. Pietzuch, and A. L. Wolf. NetAgg: Using Middleboxes for Application-Specific On-Path Aggregation in Data Centres. In *ACM International Conference on Emerging Networking Experiments and Technologies (CoNEXT)*, pages 249–262, Dec. 2014.
- [37] Mininet. <http://mininet.org>.
- [38] F. Nawab, V. Arora, D. Agrawal, and A. El Abbadi. Minimizing commit latency of transactions in geo-replicated data stores. In *ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 1279–1294. ACM, 2015.
- [39] Netronome. NFP-6xxx - A 22nm High-Performance Network Flow Processor for 200Gb/s Software Defined Networking, 2013. Talk at HotChips by Gavin Stark. http://www.hotchips.org/wp-content/uploads/hc_archives/hc25/HC25.60-Networking-epub/HC25.27.620-22nm-Flow-Proc-Stark-Netronome.pdf.
- [40] R. Nishtala, H. Fugal, S. Grimm, M. Kwiatkowski, H. Lee, H. C. Li, R. McElroy, M. Paleczny, D. Peek, P. Saab, et al. Scaling memcache at facebook. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, pages 385–398, Apr. 2013.
- [41] M. Nottingham and X. Liu. Edge Architecture Specification, 2001. http://www.esi.org/architecture_spec_1-0.html.
- [42] B. Oki and B. Liskov. Viewstamped Replication: A General Primary-Copy Method to Support Highly-Available Distributed Systems. In *ACM Symposium on Principles of Distributed Computing (PODC)*, pages 8–17, Aug. 1988.
- [43] D. Ongaro and J. Ousterhout. In Search of an Understandable Consensus Algorithm. In *USENIX Annual Technical Conference (ATC)*, pages 305–320, June 2014.
- [44] Open-NFP. <http://open-nfp.org/>.
- [45] P4 Behavioral Model. <https://github.com/p4lang>.
- [46] P4@ELTE. <http://p4.elte.hu/>.
- [47] P4.org. <http://p4.org>.
- [48] The P4 Language Specification Version 1.1.0. http://p4.org/wp-content/uploads/2016/03/p4_v1.1.pdf.
- [49] G. Pang, T. Kraska, M. J. Franklin, and A. Fekete. Planet: Making progress with commit processing in unpredictable environments. In *ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 3–14, 2014.
- [50] J. H. Saltzer, D. P. Reed, and D. D. Clark. End-to-end Arguments in System Design. *ACM Transactions on Computer Systems (TOCS)*, 2:277–288, Nov. 1984.
- [51] D. Sciascia and F. Pedone. Geo-Replicated Storage with Scalable Deferred Update Replication. In *IEEE International Conference on Dependable Systems and Networks (DSN)*, pages 1–12, June 2013.
- [52] M. Shapiro. Structure and Encapsulation in Distributed Systems: the Proxy Principle. In *6th IEEE International Conference on Distributed Computing Systems (ICDCS)*, pages 198–204, May 1986.
- [53] H. Song. Protocol-oblivious Forwarding: Unleash the Power of SDN Through a Future-proof Forwarding Plane. In *Workshop on Hot Topics in Software Defined Networking (HotSDN)*, pages 127–132, Aug. 2013.

- [54] R. Soulé, S. Basu, P. Jalili Marandi, F. Pedone, R. Kleinberg, E. Gün Sirer, and N. Foster. Merlin: A Language for Provisioning Network Resources. In *ACM International Conference on Emerging Networking Experiments and Technologies (CoNEXT)*, pages 213–226, Dec. 2014.
- [55] Y. Sovran, R. Power, M. K. Aguilera, and J. Li. Transactional storage for geo-replicated systems. In *ACM Symposium on Operating Systems Principles (SOSP)*, pages 385–400. ACM, 2011.
- [56] G. Tolle, J. Polastre, R. Szewczyk, D. Culler, N. Turner, K. Tu, S. Burgess, T. Dawson, P. Buonadonna, D. Gay, and W. Hong. A microscope in the redwoods. In *ACM International Conference on Embedded Networked Sensor Systems (SenSys)*, pages 51–63, Nov. 2005.
- [57] M. P. Đurišić, Z. Tafa, G. Dimić, and V. Milutinović. A survey of military applications of wireless sensor networks. In *2012 Mediterranean Conference on Embedded Computing (MECO)*, pages 196–199, June 2012.
- [58] H. Wang, K. S. Lee, V. Shrivastav, and H. Weatherspoon. P4FPGA: Towards an Open Source P4 Backend for FPGA. In *The 2nd P4 Workshop*, Nov. 2015.
- [59] L. Wang, V. Pai, and L. Peterson. The Effectiveness of Request Redirection on CDN Robustness. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 345–360, Dec. 2002.
- [60] G. Werner-Allen, K. Lorincz, J. Johnson, J. Lees, and M. Welsh. Fidelity and Yield in a Volcano Monitoring Sensor Network. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pages 381–396, Nov. 2006.
- [61] Xilinx SDNet Development Environment. www.xilinx.com/sdnet.
- [62] XPliant Ethernet Switch Product Family. www.cavium.com/XPliant-Ethernet-Switch-Product-Family.html.