

Barista: A Distributed, Synchronously Replicated, Fault Tolerant, Relational Data Store

Anant Bhardwaj
MIT
anantb@csail.mit.edu

Manasi Vartak
MIT
mvartak@mit.edu

Rebecca Taft
MIT
rytaft@mit.edu

David Goehring
MIT
dggoeh1@mit.edu

ABSTRACT

Barista is a *distributed, synchronously replicated, fault tolerant, relational data store*. It runs as a middleware service over database instances to provide an abstraction for a distributed relational data store. The data is replicated across a set of Paxos state machines to provide fault-tolerance and strong consistency. The real-time replication provides high availability; clients automatically failover between replicas. Barista supports SQL for data management. Client applications can use Barista with the same SQL code they used before.

1. INTRODUCTION

Barista is a middleware layer written over PostgreSQL [2] (Postgres). It replicates the data over multiple Postgres instances to provide fault tolerance. All writes are propagated synchronously using Paxos [7]. Barista provides strong consistency. Clients automatically failover between replicas. Barista exposes SQL for client applications.

In the demo, we would showcase the following contributions:

A relational data store with the following guarantees:

1. **ACID:** Barista provides ACID guarantees. It can be used as a database backend for any JDBC-compliant application.
2. **Fault Tolerance and Recovery:** Barista can tolerate $\lfloor (b-1)/2 \rfloor$ faults and recover from server crashes by catching up the state of the crashing servers. It can also tolerate $\lfloor (b-1)/2 \rfloor$ disk failures and recover by bringing the new disks up-to-date.

Cross-Language support by providing Thrift bindings for the APIs: Barista APIs can be used in all the popular languages. We have provided sample client code for Go, C++, Java, Python, and JavaScript.

State Safety with ZooKeeper: All the running states of the system are stored in Apache ZooKeeper. ZooKeeper's atomic Read() and Write() APIs provide a safe mechanism for writing/reading system states. The states are necessary for recovering from various

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivs 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/3.0/>. Obtain permission prior to any use beyond those covered by the license. Contact copyright holder by emailing info@vldb.org. Articles from this volume were invited to present their results at the 41st International Conference on Very Large Data Bases, August 31st - September 4th 2015, Kohala Coast, Hawaii.

Proceedings of the VLDB Endowment, Vol. X, No. Y
Copyright 2015 VLDB Endowment XXXX-XXXX/XX/XX.

crashes/failures. Also, ZooKeeper provides an efficient way to perform log purging.

Evaluation with the TPC-C Benchmark: We implemented the industry standard TPC-C benchmark to evaluate the throughput and latency of client operations in a real-time scenario with our system. We describe this in detail in §4.

Performance Optimizations: We made the following optimizations to our paxos-based protocol by implementing a version of Multi-Paxos [9]:

1. Avoid 2 round-trips per agreement by having a server issue Prepare messages ahead of time
2. Avoid dueling leaders under high client load by using a designated leader.

A Comprehensive Test Suite: To test Barista in a range of scenarios including concurrent client requests, network partitions, duplicate requests and machine failure, we wrote a set of automated test cases. The tests run on multiple Postgres instances on the same machine. The tests instantiate a number of Barista servers and clients and perform various combinations of operations including opening connections, executing arbitrary SQL queries and closing connections. We test that the queries are applied in the same sequence across all replicas. Our test cases cover varying numbers of DB connections, different SQL queries and failure scenarios. We also test three recovery scenarios: (1) where a server restarts with disk intact, (2) server restarts without disk and (3) server restarts in the middle of a paxos agreement.

2. IMPLEMENTATION

Barista runs as a middleware service above each database instance. It intercepts client requests to the database. These requests include opening and closing a database connection, beginning and ending a transaction, and executing SQL queries.

2.1 Paxos Agreement

Once a client request arrives, the service initiates a paxos agreement between the replicas to reach consensus on the slot-id in the paxos log for the client request. This ensures that all replicas agree on a single ordering for the client requests. The paxos state and logs are stored in Apache ZooKeeper so that in case of a failure, a recovering machine or a replacement machine can reconstruct its state.

As opposed to gets and puts, Barista supports the core operations required to interact with a database, namely: opening and closing connections to the database; starting, committing and rolling back transactions; and executing queries. Each of these operations is tracked in the paxos log. Since the presence of replication must be

transparent to clients (and therefore they must not need to connect to different instances separately), we track the opening and closing of connections in the log in addition to query operations.

2.2 Enforcing the ordering on Postgres

PostgreSQL is a non-deterministic database. This can cause transactions to appear in a different order than the order agreed upon by the replicas during paxos-agreement. This is because Postgres is multi-threaded and transactions can run as different threads. The threads might get scheduled in any order and thus the commit order can be different from the order in which the transactions were submitted. To work around this, we do not allow more than one pending request on any replica. Before submitting a new query to the db, Barista makes sure the last one has returned. This is done by synchronizing database calls with a lock. This affects the performance in terms of number of client requests per second, but makes consistency and recovery easier.

2.3 State Safety

We use Apache ZooKeeper for Paxos state safety. Apache ZooKeeper provides a distributed configuration service for maintaining configuration information for large distributed systems.

For each paxos instance there is a node in the ZooKeeper with the path `/barista/paxos/machine_name/seq_num = Paxo N_P, N_A, V_A, Decided`. The paxos instances update the state by calling `Set(path, value)` API. ZooKeeper's `Write()` and `Read()` APIs are atomic which guarantees consistency. The following is a bit of sample code that updates the state of a paxos instance:

```
px.path = "/paxos/" + px.Format(px.peers[px.me])
...
...

if args.N_A >= paxo.N_P {
    paxo.N_P = args.N_A
    paxo.N_A = args.N_A
    paxo.V_A = args.N_A
    reply.Status = OK
} else {
    reply.Status = REJECT
}

if px.use_zookeeper {
    px.Write(
        px.path + "/store/" + strconv.Itoa(args.Seq), paxo
    )
} else {
    px.store[args.Seq] = paxo
}
...
...
```

We also do write-ahead logging of database operations. The log is stored as a path `(/barista/sqlpaxos/machine_name/seq_num)` in Apache ZooKeeper. We call this log `sqlpaxos` log. The write-ahead logging allows us to run paxos agreement for multiple client queries in parallel.

We maintain an application pointer (AP) to track the last transaction that has been applied to the database for each replica. *The update to AP and the actual transaction must be atomic* – thus this state variable can't be in a log file outside the database. We maintain a table called `sqlpaxoslog (lastseqnum int)` that tracks the Application Pointer (AP) for that Postgres instance. We change each client transaction by adding the AP update as part of the transaction itself to achieve atomicity. The AP is used to recover the database state as discussed in the recovery section.

2.4 Log Purging/Garbage Collection:

When `paxos.Done()` from other peers updates `paxos.Min()` – all paxos instances in ZooKeeper with `seq_num < paxos.Min()` are purged. This is done by removing all `/barista/paxos/machine_name/seq_num` nodes if the `{seq_num} < paxos.Min()`.

We garbage-collect `sqlpaxos` logs by deleting `/barista/sqlpaxos/machine_name/seq_num` nodes if the `{seq_num} < Global AP`. The Global AP is the sequence number below which all the operations have been applied on all the replicas.

The purging allows us to keep the ZooKeeper logs small.

2.5 Recovery

Recovery from crash & restart (no disk failure): The paxos state is recovered by reading the saved states from the ZooKeeper. The ZooKeeper Get/Set APIs are atomic which ensures that the entire state can be recovered consistently. Below is a snippet that shows how we recover Paxos state for a particular slot:

```
px.path = "/paxos/" + px.Format(px.peers[px.me])
...
...
if px.use_zookeeper {
    paxo, ok = px.Read(
        px.path + "/store/" + strconv.Itoa(args.Seq))
} else {
    paxo, ok = px.store[args.Seq]
}
...
...
```

The `sqlpaxos` write-ahead logs don't have any state other than the application pointer (AP). The AP is updated as the part of the transaction itself in the database and thus can be recovered by reading the `sqlpaxoslog (lastseqnum int)` table in the database. Paxos fills holes in its log to ensure that everything after the AP can be retrieved as part of the paxos protocol.

Recovery from disk failure: This recovery is slow and requires some manual steps. To recover from a complete disk wipe out, we provide a script that copies the database data files from a *healthy_machine*. The recovery requires that the *healthy_machine* is not serving any request during the recovery because if it serves a new request it will change its state during the recovery and would lead to inconsistent data transfer.

3. CLIENT APIS

Barista APIs are exposed as a Thrift IDL file (*barista.thrift*). Thrift [5] is a framework for scalable cross-language services development. It combines a software stack with a code generation engine to build RPC services that work efficiently and seamlessly between C++, Java, Go, Python, Ruby, JavaScript, and various other languages. A Thrift IDL file is processed by the Thrift code generator to produce code for the various target languages to support the defined data types and services in the IDL file. Although we are implementing Barista in Go on the server side, clients can be implemented in any language. We have provided sample client code in Go, C++, Java, Python, and JavaScript.

3.1 Barista Client APIs

Below is the list of data types and methods available to client applications through *barista.thrift*:

```
/* Barista constants */

// version info
const double VERSION = 0.1

/* Database Connection */

// connection parameters
struct ConnectionParams {
    1: optional string client_id,
    2: optional string seq_id,
    3: optional string user,
    4: optional string password,
    5: optional string database
}

// connection info -- must be passed in every execute_sql call
```

```

struct Connection {
    1: optional string client_id,
    2: optional string seq_id,
    3: optional string user,
    4: optional string database
}

/* ResultSet */

// A tuple
struct Tuple {
    1: optional list <binary> cells
}

// A result set (list of tuples)
struct ResultSet {
    1: required bool status,
    2: Connection con,
    3: optional i32 row_count,
    4: optional list <Tuple> tuples,
    5: optional list <string> field_names,
    6: optional list <string> field_types
}

/* Barista Exceptions */

// Database Exception
exception DBException {
    1: optional i32 errorCode,
    2: optional string message,
    3: optional string details
}

/* Barista RPC APIs */

service Barista {
    double get_version()

    Connection open_connection (1: ConnectionParams con_params)
        throws (1: DBException ex)

    ResultSet execute_sql (1: Connection con, 2: string query,
        3: list <binary> query_params) throws (1: DBException ex)

    ResultSet execute_sql_txn (1: Connection con, 2: string query,
        3: list <binary> query_params) throws (1: DBException ex)

    void begin_txn (1: Connection con)
        throws (1: DBException ex)

    void commit_txn (1: Connection con)
        throws (1: DBException ex)

    void rollback_txn (1: Connection con)
        throws (1: DBException ex)

    void close_connection (1: Connection con)
        throws (1: DBException ex)
}

```

3.2 A Sample Python Client

Barista RPC stubs for Python can be generated as follows:

```
thrift --gen py barista.thrift
```

Once RPC stubs are generated, a python program can call Barista APIs. Below is a sample client code snippet in Python:

```

transport = TSocket.TSocket('localhost', 9000)
transport = TTransport.TBufferedTransport(transport)
protocol = TBinaryProtocol.TBinaryProtocol(transport)
client = Barista.Client(protocol)

transport.open()

con_params = ConnectionParams(
    user="postgres", password="postgres", database="postgres",
    client_id="1234567890", seq_id="1")

con = client.open_connection(con_params)
res = client.execute_sql(con=con,
    query="SELECT 6.824 as id, 'Distributed Systems' as name",
    query_params=None)

print "\t".join(res.field_names)
for tuple in res.tuples:
    print "\t".join(tuple.cells)

client.close_connection(con)
transport.close()

```

4. EVALUATION

In order to test the performance of our system, we implemented the TPC-C benchmark [3], an industry standard for comparing the performance of OLTP database systems. TPC-C simulates the operation of a wholesale parts supplier, in which a population of terminal operators executes a set of transactions against a database.

Transaction	%	Response Time (s)		Total	Rollbacks	%
		Average	90th %			
Delivery	3.99	0.137	0.218	61	0	0.00
New Order	43.82	0.160	0.508	670	6	0.90
Order Status	2.88	0.218	0.903	44	0	0.00
Payment	40.35	0.172	0.606	617	0	0.00
Stock Level	2.81	0.220	0.901	43	0	0.00

Table 1: Results from running the benchmark for five minutes on a single, non-replicated instance of PostgreSQL

Transaction	%	Response Time (s)		Total	Rollbacks	%
		Average	90th %			
Delivery	3.48	0.766	2.407	50	0	0.00
New Order	45.06	0.733	2.325	648	5	0.77
Order Status	3.34	0.568	2.191	48	0	0.00
Payment	37.90	0.698	2.332	545	0	0.00
Stock Level	4.10	0.851	2.528	59	0	0.00

Table 2: Results from running the benchmark for five minutes on Barista with standard Paxos

These transactions include monitoring the stock level of a warehouse, creating a new order for a customer, accepting payment from a customer, making a delivery to a set of customers, and checking the status of an order. Clearly, the intent of this benchmark is to simulate a realistic real-time OLTP system.

The TPC-C specification is quite complex, so rather than implementing it ourselves, we modified an existing open-source implementation for PostgreSQL [1]. The open-source implementation uses C code to communicate directly with the database, so in order to test our replicated, fault-tolerant framework, we needed to modify this code to communicate via the Barista client instead. Fortunately, Thrift allows us to create clients in several different popular programming languages, so we created a C++ Barista client and linked it with the modified C benchmark code. In this way, we were able to compare the performance of our fault-tolerant, replicated version of PostgreSQL with vanilla PostgreSQL. The results of running TPC-C on Barista (with standard Paxos and Multi-Paxos variants) and vanilla PostgreSQL are shown in tables 1 to 3.

As expected, vanilla PostgreSQL performs better than Barista since it does not require paxos agreement for each transaction. However, the performance of Barista is still impressive given that each transaction can take 5-8 paxos agreement instances, and each paxos agreement instance can require up to 20 round-trip messages between replicas. Despite this overhead, Barista only increases latency by about 4-5X and barely decreases throughput at all. Furthermore, the average round-trip time between our virtual machines is 0.9 ms, so given the number of messages passing between machines, an added latency of a few hundred ms is reasonable. In addition, Barista is extremely fault tolerant while vanilla PostgreSQL relies on a single point of failure.

Transaction	%	Response Time (s)		Total	Rollbacks	%
		Average	90th %			
Delivery	3.96	0.806	2.482	57	0	0.00
New Order	44.69	0.948	2.508	644	4	0.62
Order Status	2.78	0.756	2.069	40	0	0.00
Payment	39.69	0.719	1.997	572	0	0.00
Stock Level	2.91	0.762	2.240	42	0	0.00

Table 3: Results from running the benchmark for five minutes on Barista with Multi-Paxos

Theoretically, Multi-Paxos should perform faster than traditional Paxos due to the enhancements described (primarily the reduced contention and number of messages needed for agreement) which produce smaller latencies. Generally, for an N replica system the number of messages required for agreement is $2(N-1)$, $N-1$ Accept and Learn messages to other servers from the leader. This is much smaller than 1) the $3(N-1)$ messages required in the case of a single normal Paxos proposer ($N-1$ Prepare, Accept, Learn messages) and 2) a possibly unbounded number messages in the case of dueling proposers. In our 5 replica system this translates to 8 RTTs instead of 12 for the best case. Thus, since the total number of round trips for agreement is minimized, the total end to end latency for a particular operation is also minimized. However, the results indicate that performance is comparable to normal paxos despite the performance improvements. This could be a result of the increased overhead of pinging the leader and updating replicas over a ping. However, we are still investigating why the performance of Multi-Paxos is not as high as it should be.

5. RELATED WORK

Synchronous replication schemes for PostgreSQL have been explored in the literature. The main challenge is ensuring that transactions are committed in the same order on all replicas. In our work we assume that there are no distributed transactions. One method to accomplish consistent ordering is 2 phase commit [4] across all replicas – however, this requires all replicas to be alive and communicating at all times, thus nullifying any benefit from Paxos. Moreover, 2 phase commit significantly increases latency. Systems like H-store [6] adopt a concurrency protocol that is defined by having one thread of execution per data partition and thus eliminating locking in single-partition transactions. In a synchronous replication scheme, the equivalent would be to only execute one transaction at a time on each replica and start the next one only after the current transaction finishes. Another technique that has been explored on unreplicated systems is speculative execution [8] where the system runs subsequent transactions even before the current one has finished. It is likely that some of these transactions may have to be redone based on the results of the current transaction. In our scenario, we want to avoid having to undo transactions. We can get the benefit of speculative execution instead by batching non-conflicting transactions. Some systems also adopt the technique of

using a pre-processor that determines the serial order of all transactions and transactions must request locks in the same order. The Paxos log can help with creating a universal order of transactions.

6. SUPPLEMENTARY MATERIALS

- **Source Code:** <https://github.com/abhardwaj/barista>
- **Benchmark Code:** <http://people.csail.mit.edu/anantb/files/barista/6.824/>

7. ADDITIONAL AUTHORS

8. REFERENCES

- [1] Database Test Suite: DBT-2 OLTP Benchmark. <http://sourceforge.net/apps/mediawiki/osdldbt/>.
- [2] PostgreSQL. <http://www.postgresql.org>.
- [3] TPC Transaction Processing Performance Council: TPC-C OLTP Benchmark. www.tpc.org/tpcc/.
- [4] Two-phase commit protocol. http://en.wikipedia.org/wiki/Two-phase_commit_protocol.
- [5] A. Agarwal, M. Slee, and M. Kwiatkowski. Thrift: Scalable Cross-Language Services Implementation. Technical report, Facebook, 2007.
- [6] R. Kallman, H. Kimura, J. Natkins, A. Pavlo, A. Rasin, S. Zdonik, E. P. C. Jones, S. Madden, M. Stonebraker, Y. Zhang, J. Hugg, and D. J. Abadi. H-store: A High-performance, Distributed Main Memory Transaction Processing System. *Proc. VLDB Endow.*, 1(2):1496–1499, Aug. 2008.
- [7] L. Lamport. The Part-time Parliament. *ACM Trans. Comput. Syst.*, 16(2):133–169, May 1998.
- [8] B. W. Lampson. Lazy and Speculative Execution in Computer Systems. In *Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming, ICFP '08*, pages 1–2, New York, NY, USA, 2008. ACM.
- [9] R. D. Prisco, B. W. Lampson, and N. A. Lynch. Revisiting the Paxos Algorithm. In *Proceedings of the 11th International Workshop on Distributed Algorithms, WDAG '97*, pages 111–125, London, UK, UK, 1997. Springer-Verlag.