

CAP in practice: HBase and MongoDB

Thomas Uyttendaele

Abstract

This paper introduces a new, quantified method to test distributed database systems in regard to their behaviour to availability in case of a failing node and consistency of data from the users perspective. This method is tested on HBase and MongoDB, two consistent and partition tolerant systems, with the most important outcomes included.

Keywords:

1. Introduction

New online services and more online users means more data, data and load a single server can't handle. More and more database systems are distributed, for higher availability in case of an unexpected crash but also for horizontal distribution: the data of a single database is spread over different systems. These new services have different requirements; e.g., an update doesn't need to be visible for all users immediately, it can take time, this behaviour is called eventual consistency.

Over the past years, many new systems have been built on this wave of changes, categorized as NoSQL. Some applications contain a high volume of data, have the need for consistent data and being able to work in a highly distributed environment, this are the CP systems in the CAP Theorem. Two examples that offer these guarantees are HBase and MongoDB. They greatly differ in supported queries on their system, but what happens if you only use the bare essentials and compare their behaviour in a distributed environment going from expected shut downs of instances towards crashes and network partitions?

In this article, a comparison of both these systems on these 3 behaviours will be made. In chapter 2 a brief overview of the CAP Theorem is given, chapter 3 discusses HBase and MongoDB on paper. Chapter 4 gives an overview of the used test method and chapter 5 presents the results. The future work is presented in chapter 6, an overview of related work is given in chapter 7 and a conclusion is made in chapter 8.

2. The CAP Theorem[1][2]

The CAP Theorem was introduced by E. Brewer [1] in 2000 and discusses 3 properties of which each network shared-data system can guarantee at most 2: (definitions based on [2])

- (Strict)Consistency: The system acts like there is only single storage
- High availability: The system is available (for updates)

- Partition tolerance: A split in the network let the different partitions still act as a single system.

Designers used this model to explain their design decisions, others used it to compare different systems and sometimes it was misused. As E. Brewer explains 12 years after the launch, the "2 of 3" can be misleading.

One of the reasons is that there exist several types of consistency, partition tolerance and different availability guarantees. The choice for the properties needs to be made several times in the different subsystems and the end solution is not black or white.

At first glance, it also looks like partition tolerance has to be implemented and therefore consistency and availability needs to be given up. In practice, partition splits are only rare and therefore both consistency and availability can be allowed most of the time.

Each of the 3 choices will be discussed in more detail, how their implementation could work, what the influences are and some examples.

2.1. CA

When forfeiting partition tolerance, these systems provide all the time consistent data available to all nodes, except when there are one or more nodes unavailable. In that case, write requests will be not allowed.

These systems can be built around the 2 phase commit and have cache invalidation protocols. Examples of this type are the typical relational databases roll out in clusters.

2.2. CP

A consistent system with partition tolerance will provide all the time the last data, even in the case of network splits. This comes with the loss of the availability of all nodes all the time.

The system can allow operations only on the majority partition. In case multiple splits are present and no partition has a majority of nodes, the whole system can be unavailable. These systems can be built around a master/slave principle where the operations will be directed to the master, the slaves are present to continue operation when the master fails.

In practice, systems like MongoDB, HBase and Redis select CP.

2.3. AP

In a highly available systems with partition tolerance, it is possible to read inconsistent data. As read and write operations are still allowed when there are different partitions, it is possible that the database has other content depending on the used node. When the split is dissolved, a need for manual conflict resolution can be needed. In case a record is adapted in both partitions, the user will need to choose the correct version.

Example systems following AP are Cassandra, Riak and Voldemort.

3. Overview of HBase and MongoDB

In this article, 2 CP systems will be discussed more in detail regarding their choices to forfeit availability and the influence on their behaviour in practice. The systems are HBase and MongoDB, an architectural overview will be given in this section.

3.1. HBase

HBase[3] is an open-sourced, distributed, versioned database designed after Google's BigTable [4]. HBase relies on Zookeeper for the distributed task coordination and the persistent storage can be done on the local hard disk, Hadoop Distributed File System or Amazon S3. In this article is chosen for Hadoop.

HBase nodes exists out of HMaster and HRegionServers, the coordination of the system is done by one HMaster, the handling of data is done by the HRegionServer. To store the data, multiple Hadoop datanode instance should be deployed for data storage, preferable one on each HRegionServer. The data will be replicated to a configurable amount of other nodes, default modus is 3. The data is stored in a table, which is split in one or more regions. A region is leased to a given HRegionServer for a defined time. During this time, only this server will provide the data of the region to the different users. This way the consistency of data can be guaranteed because there is for each record only a single system responsible. Consistency on a single record is provided by a readers/writer lock on a single record for the according queries, this way there is a guarantee to atomicity on a single record, the full procedure is explained by Lars Hofhansl[5].

To be partition tolerant, the partition with the majority of the Zookeeper servers and a HMaster will appoint regions to available HRegionServers, let's call this the data serving partition. In this approach, it is important to place the Zookeeper and HMaster servers in diverse location as otherwise a partition of only management servers will make the whole system unavailable.

In HBase, a node will be able to answer to requests if the node is present in the data serving partition. Only in rare cases that all data copies of Hadoop are stored in unavailable servers, the data will be unreachable. The nodes not in the data serving partition, will be unable to complete any requests.

When a server goes down, he can release the lease in case there is a graceful shut down (the HBase server is notified) and another server can get the lease immediately. In other cases, a new lease can only be given after the decay of the old lease, if the server comes back online in meantime, he will still be responsible.

3.2. MongoDB

MongoDB[6] is an open-sourced, distributed database designed for document storage, this are data entries where the format of each record can be different. According to their website they provide high performance and high availability, but this is incorrect to the given definition in this article.

MongoDB provides data replication and data distribution, the first is done by grouping different MongoDB servers into a ReplicaSet, the second is done by grouping different of these ReplicaSets.

A ReplicaSet exists out of different MongoDB servers whom work as a master/slave configuration. A master is a primary and a slave is called a secondary. The primary is responsible for the write actions, by default a query will succeed once it has a confirmation that the write has been executed on the primary. The read operation will go by default on the primary as well. Both query methods are configurable to give other guarantees, for a write operation there are multiple *write concerns*, it is possible to wait till it has written on hard disc or a number of secondary servers, however all need a primary. For read operation there are multiple *read preferences*, it is possible to read from a secondary or the closest server.

In the default configuration, MongoDB provides a consistency guarantee.

In a ReplicaSet, there is at least half of the ReplicaSet needed for the primary election. As there is always a primary needed, the system has partition tolerance but no high availability, contrary to the statement on the documentation. However, it is possible to read from a secondary but writing is not possible.

The state of the different members of a ReplicaSet is maintained by a heartbeat system: a server is marked as offline if no beat has been received for 10 seconds. In case the primary goes offline, election will be started to re-elect a new one. In other words, a primary has a lease of 10 seconds. This value of 10 seconds is non-configurable.

De data distributions happens by merging replica sets in a cluster. Furthermore, there is the need for access server (as many as you want) and configuration servers (1 or 3). The availability and partition tolerance of the data is the same as in the ReplicaSets as it handled by the ReplicaSets. In case a access server can't reach a primary, another access server will be needed to write data. If a majority of the configuration servers are not reachable, their will be no reconfiguration of the data over the different servers.

3.3. Differences between databases

Both systems provide consistency and partition tolerance and forfeit high availability, but some differences are in their implementation.

First of all, they differ in their handling of partition tolerance, in HBase there are dedicated management servers (HMaster and Zookeeper) to distribute the responsibilities of regions and if the management servers are in a partition with a minority of the data servers, data will be not available. In MongoDB the management of a ReplicaSet is done internally with as result that the write queries will be available in the partition with the majority of the data servers.

Both decisions have their reasons and possible motivations, in HBase it is possible to write every record, as a new region can be created if the old region is unavailable. In MongoDB it is possible that in sharding you can read all the data from multiple secondaries but only write given ranges of records.

In availability, there is a small difference: where in MongoDB it is possible to read from secondaries, this is impossible in HBase.

In section 5, a detailed analyse will be done to the consistency and the behaviour of the systems in the case of network or server failure.

4. Test method

To test the behaviour of database systems towards consistency and availability, the Yahoo Cloud Serving Benchmarking (YCSB [7]) has been extended with event support for availability and reader-writers for consistency.

Each of this tests follows the same steps: calibrating the system, records are preloaded, the test is started and in the end all the records are removed again and they are executed for HBase and MongoDB. During the calibration, a workload is chosen so there is a medium load on the different databases.

4.1. Event support

The implementation of event support is integrated so it supports the execution of UNIX commands at specified moments, the execution time and result code is logged.

Before the test are 300 000 records stored in the database to enable the sharding in all database softwares. In these tests were 3 kinds of tests executed of which each one took 900 second. The first action takes place at 300 seconds, the second at 600 seconds.

- Graceful shut down of a data service and restart of the service
- Hard kill of the data service and a restart of the service
- Blocking of all network traffic from and to a server and the allowance of all network traffic

Each element tests the behaviour of the system under different circumstances, the first two checks what happens in case of a respectively planned and unexpected shut down, the latter check the behaviour in case the network fails.

The tests are executed on each of the data nodes of HBase and MongoDB, caching and buffering on the client side are disabled in both tests.

4.2. Consistency support

To implement consistency support, the YCSB software is extended with an extra workload. A graphical representation of an example workload is shown in figure 1. This workload ex-

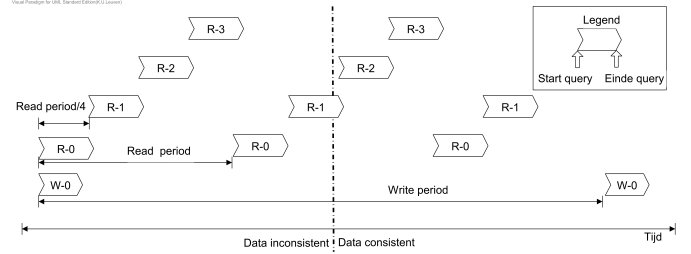


Figure 1: Example of consistency support for one write period with 1 writer and 4 readers

ists out of 1 writer and a user-defined amount of readers. Each writer will inserts or updates a value in the database with at a user-defined pace. The readers will read the record till they read the last written data element, each reader reads in a period defined by the user. The different readers are scheduled uniformly within the reading period.

All this data is logged and gives possibilities to analyse when a record is visible for different users, compared to the time the queries where started or ended.

Before the test, there are 30 000 records stored in the database, each tests takes 500 seconds and results start to be gathered after 30 seconds. In the tests is chosen to write every 0.5s and read with 5 readers every 10ms for MongoDB, for HBase there are 10 readers who read every 30ms.

5. Results

To execute the tests, both systems were deployed on a virtual platform of OpenStack. Each instance has 2 CPUs, 4GB RAM and 50GB disc space. The machines are connected with a gigabit Ethernet and an average ping takes 0.4ms ($\sigma = 0.2$ on 10 000 ping's).

HBase is configured with 5 instances, of which 1 for management (HMaster, Hadoop namenode and Zookeeper) and 4 for data storage (each has a HRegionServer and Hadoop datanode).

MongoDB is configured with 6 instance, grouped by 3 in a ReplicaSet. There was a single configuration server and 3 instances had an access server.

YCSB was deployed on a single instance and used to calibrate the load on both systems. An individual record has 10 fields which each field a size of 100 bytes. The workload existed out of 20% inserts and updates, 40% selects and 20% scans of an uniform spread between 1 and 100. The requests are spread according to Zipfian¹ The basic load for HBase is an average of 600 queries/second spread over 50 threads, for MongoDB there are 15 threads with a total of 200 queries/second.

¹Some record are popular, others rare to be used.

5.1. Availability

When reading and writing in the default setting of MongoDB, both MongoDB and HBase will read and write from a leader of a set of data. Both have a lease period for the leader, which can't be set for MongoDB (10 seconds), but can be configured for HBase (default 180 seconds).

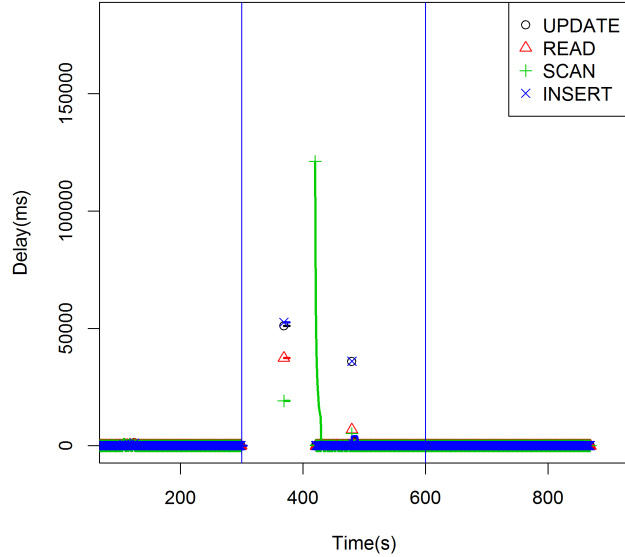


Figure 2: Example of network partition. The requests block for 110 seconds

In case of a stop of a HBase server in this configuration, the queries will halt till the lease for the region has been expired, this can take between 0 seconds and 180, depending on the moment of the action. In case of a graceful stop of the service, the service will release the lease and a faster handover is possible. During the tests it happened that with a hard stop of the service, the queries are impossible on the connections till the server is brought back online. In this scenario it was needed to manually dis- en reconnect to the database. A reason why this happens here is not known.

For MongoDB, there is no difference between the graceful or hard stop of an instance; in case it was a secondary, no influences on the latency will be seen. In case a primary was stopped, the data will be in both cases temporarily unavailable for a few seconds. It seems that in a hard stop, there is still a messaging of the shut down to the other nodes. In case of a network interruption, the results are varied, going from no influence, a temporary unavailability for a few seconds, or no completed queries the whole time the server is down. When in the last case, all connections are reinitialised, the queries can happen, when this is not done the table seems empty.

A short overview of all the reactions is shown in table 1.

5.2. Consistency

In consistency HBase and MongoDB have a different approach, in MongoDB it is possible to configure the read and

	Graceful stop	Hard stop	Network partition
HBase	Few seconds	Dozens of seconds to unlimited	Dozens of seconds
MongoDB	1/3 of the cases, Few seconds	1/3 of the cases, Few seconds	Few seconds to unlimited

Table 1: Availability: Overview of different reaction when stopping an instance

Reader	Start time (ms)	Percentage correct read
1	0 ms	2.6%
2	3 ms	68%
3	6 ms	90%
4	9 ms	93%
5	12 ms	94%
6	15 ms	96%
7	18 ms	96%
8	21 ms	96%
9	24 ms	97%
10	27 ms	97%

Table 2: Consistency: Percentage of the queries started at the given time moment that will read the newest inserted data. The average latency in a random read transaction is around 6ms.

write queries, in HBase there is the choice to use cache at client side. In the tests the cache is disabled but the different options of MongoDB are tested.

HBase: If a record is read while a write query is creating or updating the value, the read query will wait till the completion of the write query and return immediately the correct value. If a read query is send too soon, it will return the old data, but each query returned after the completion of a write query, will return the new data.

MongoDB: A read query can return the old or new data before the write query has been completed. However a readers/writer lock mechanism is used, it seems that after releasing the lock, the write query still has to execute extra steps. When reading from a secondary or the nearest server, there will be monotone read consistency, except if the driver of MongoDB selects another server between the reads, as there is no guarantee that the other servers already have the value. The process of finding the nearest server happens periodically in the background.

From the test results, the writer configuration has no influence on the speed the replication will have, as the read queries get the same consistency window based on the start moment. The difference is the guarantee a user has when a write query has finished.

An overview of the consistency windows is shown in table 2 and 3 for respectively HBase and MongoDB.

6. Future work

The test run on the different systems are a start to have a quantified approach on the consistency and availability of the different systems. However multiple extensions can be made to this test method, next to testing more systems.

	nearest	primary	primary- preferred	secondary
safe	28, 69, 89, 91, 92	80, 98, 98, 99, 99	74, 99, 99, 99, 99	0, 65, 83, 85, 88
normal	24, 68, 87, 89, 92	72, 99, 100, 100, 100	75, 98, 98, 98, 98	0, 69, 85, 89, 92
fsync_safe	28, 73, 87, 90, 90	68, 96, 98, 98, 98	78, 97, 98, 98, 98	0, 66, 80, 85, 86
replicas_safe	24, 74, 87, 88, 91	75, 98, 99, 99, 99	79, 98, 98, 98, 98	1, 67, 84, 87, 89
majority	26, 77, 91, 91, 92	73, 98, 99, 99, 99	77, 99, 99, 99, 100	0, 61, 82, 85, 89

Table 3: Consistency: Percentage of the queries started at 0, 2, 4 and 8ms that will read the newest inserted data with the each row for a write configuration and a column for the read configuration. The average latency on a read transaction is for all read queries around 1ms.

First of all, the strange behaviour of MongoDB and HBase of not allowing connections, could be researched in more detailed.

Secondly, the third element of CAP can be included: Partition tolerance. Right now the connection nodes are chosen when setting up the test, an extension could be to test all connections and see what happens in the case of a partition tolerance. How is the availability of each node, how long is a primary in MongoDB still accepting queries however it is cut off?

Thirdly, the availability and consistency tests could be tested together: when a node is shut down or cut off, is there the loose of data? Especially in MongoDB it could be the case that the data is already read on the primary but not yet replicated to the primary, what happens if the primary is cut off?

At last, the test parameters could be adapted, what happens with a different network infrastructure when the network distance is not equal any more, what happens with a higher and lower load? It could be possible in the end to have a mathematical formula which could predict the average consistency window for example.

7. Related work

The research towards database systems and their consistency guaranties is rare to have a measured approach. In recent paper (February 2014), Golab et al. states that there is only a limited amount of research done towards eventual consistency [8]. They present a new view on consistency, were already some research is done towards active analyse (how long before the data is replicated to all nodes), the amount of passive analyse is limited (what do the users see). Compared to the consistency results from this paper, both are discussed: as well the delays before the data is present everywhere but also on the acting of the specific systems, in example the behaviour of HBase.

Another extension of YCSB called YCSB++[9], provides more logging information on all systems in the first place, but they also test the consistency of HBase in regards of the client caches. The reasoning for this is that it is the standard run configuration of HBase, but the submitting of the client cache towards the system depends not only on the time, but also on the amount of traffic that is being submitted. In the study they compare different cache sizes and this shows already a difference in time. Furthermore, it is possible to disable this caching in case there are records in the need of this strict consistency.

For availability benchmarking, there was no research found on related databases. However, research from 2004 [10] discuss

a way to let a standalone system recover and provide a starting benchmark for it.

8. Conclusion

A new testing method to quantify consistency and availability has been presented in this paper. With this method the effect of a stop of an instance can be measured and also the consistency window for different readers.

Another contribution are the results for two consistent and partition tolerant systems, HBase and MongoDB. These 2 systems are tested with the new method and show their relevance. According to the documentation, MongoDB and HBase both guarantee strict consistency but based on our results their characteristics are different: HBase blocks read queries till the completion of the write queries to guarantee a read will always be the same, MongoDB returns the ready query soon with new or old data before the write operation has been finished.

In the availability of the systems, HBase is in the standard configuration longer unavailable but with a change in the session time out, they can be made more the same. However both have a strange behaviour of not returning any results, in the case of HBase this is for a hard stop, in the case of MongoDB it is for a network partition. This doesn't happen in the other situations and in case of reconnecting to the database, the problems are solved.

- [1] E. A. Brewer, Towards robust distributed systems, in: PODC, 2000, p. 7.
- [2] E. Brewer, Cap twelve years later: How the "rules" have changed, Computer 45 (2) (2012) 23–29.
- [3] Hbase, apache hbase.
URL <https://hbase.apache.org/>
- [4] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, R. E. Gruber, Bigtable: A distributed storage system for structured data, ACM Transactions on Computer Systems (TOCS) 26 (2) (2008) 4.
- [5] L. Hofhansl, Hbase: Acid in hbase (3 2012).
URL <http://hadoop-hbase.blogspot.be/2012/03/acid-in-hbase.html>
- [6] The mongodb 2.6 manuel.
URL <http://docs.mongodb.org/manual/>
- [7] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, R. Sears, Benchmarking cloud serving systems with ycsb, in: Proceedings of the 1st ACM symposium on Cloud computing, ACM, 2010, pp. 143–154.
- [8] W. Golab, M. R. Rahman, A. AuYoung, K. Keeton, X. S. Li, Eventually consistent: not what you were expecting?, Communications of the ACM 57 (3) (2014) 38–44.
- [9] S. Patil, M. Polte, K. Ren, W. Tantisiroj, L. Xiao, J. López, G. Gibson, A. Fuchs, B. Rinaldi, Ycsb++: benchmarking and performance debugging advanced features in scalable table stores, in: Proceedings of the 2nd ACM Symposium on Cloud Computing, ACM, 2011, p. 9.
- [10] J. Mauro, J. Zhu, I. Pramanick, The system recovery benchmark, in: Dependable Computing, 2004. Proceedings. 10th IEEE Pacific Rim International Symposium on, IEEE, 2004, pp. 271–280.