

CAP Theorem Impact in Reliable Data Processing

Pavel Krobot, Faculty of Information Technology - Brno University of Technology -
xkrobo01@stud.fit.vutbr.cz

Extended abstract

This paper deals with the subject of reliable distributed processing of IP flow data. Presented thoughts are based on Brewer's CAP theorem and its consequences. This theorem is still one of the most important findings for distributed databases. The need for processing data in a distributed manner emerged from a constant growth of datasets, which have to be analyzed. Existing ACID (Atomicity Consistency Isolation Durability) databases have almost reached their limits and they are not feasible for reliable processing of huge datasets in real time, which is what we want in processing of IP flow data. Studying new principles is inevitable. One of such is BASE paradigm (Basically Available Soft-state Eventually consistent). As a result of BASE consequence and implication analysis, Brewer came up with CAP theorem. This theorem deals with three main properties of distributed data processing systems. Consistency (C) is equivalent to having a single up-to-date copy of data. High availability (A) provides access to the data at any time. Partition tolerance (P) states that the distributed system can operate as usual when a network partition occurs. CAP theorem itself claims that any networked system with shared data can have at most two of these three desirable properties. Despite the fact that this theorem was formally proven by Gilbert and Lynch, some of its ideas could be misleading.

In its formulation CAP theorem presents all three properties as equal. This is not accurate in practical use. Firstly, consistency and

availability can be measured in a spectrum, while partition tolerance is rather binary. Definition of the partition tolerance can vary, but at the end we can only say if the distributed system supports the partition tolerance or not. Secondly, even if partitions arise rarely they can never be forfeited absolutely. As a result, in the real system we have to decide between consistency and availability. On the other hand as, this decision does not have to be binary and we could choose a tradeoff between both, consistency and availability, according to the system needs.

Analyzing this tradeoff in search for maximal consistency and availability in distributed IP flow data processing system is the main focus of this paper. One possibility is to consider partitions as rare and have system with both, consistency and availability along with partition detection. When the partition is detected, system switches into another mode of operation with limited functionality. Some operations, like data updates have to be executed carefully or postponed until system is recovered from partition. Second choice is to have system prepared for occurrence of partitions all the time. It is necessary to have some tradeoff between consistency and availability in this scenario. CAP theorem as such ignores latency, which is in the real system very important aspect. For example, it has great impact on partition detection in the first proposed solution. Both examined solutions have slightly different characteristics and require fulfillment of different conditions, which is identified and evaluated in this study.

References

- [1] SALOMÉ, Simon. Report to Brewer's CAP Theorem [online]. 2012 [cit. 2015-10-11]. Dostupné z: <https://fenix.tecnico.ulisboa.pt/downloadFile/845043405442708/10.e-CAP-3.pdf>
- [2] BREWER, Eric. *CAP twelve years later: How the "rules" have changed*. 45(2): 23-29. DOI: 10.1109/MC.2012.37. ISSN 0018-9162. Dostupné také z: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6133253>
- [3] HALE, Coda. You Can't Sacrifice Partition Tolerance. <http://codahale.com> [online]. 2010 [cit. 2015-10-11]. Dostupné z: <http://codahale.com/you-cant-sacrifice-partition-tolerance/>