

Efficient (re)naming in Conflict-free Replicated Data Types (CRDTs)

Matthieu Nicolas, Gérald Oster and Olivier Perrin

March 27, 2018

In order to serve an ever-growing number of users and provide an increasing volume of data, large scale systems such as data stores[3] or collaborative editing tools[4] have to adopt a distributed architecture. However, as stated by the CAP theorem[2], such systems cannot ensure both strong consistency and high availability. As a result, the literature and companies increasingly adopt the optimistic replication model known as eventual consistency to replicate data among nodes. This consistency model allows replicas to temporarily diverge to be able to ensure high availability, even in case of network partition. Each node owns a copy of the data and can edit it, before propagating the updates to others. A conflict resolution mechanism is however required to handle updates generated in parallel by different replicas.

An approach which gains in popularity since a few years proposes to define Conflict-free Replicated Data Types (CRDTs)[6]. These data structures behave as traditional ones, like the *Set* or the *Sequence* data structures, but are designed for a distributed usage. Their specification ensures that concurrent updates are resolved deterministically, without any kind of agreement, and that replicas eventually converge after observing all updates, thus achieving *Strong Eventual Consistency*[7].

In [6], Shapiro et al present two designs of CRDTs : *State-based* CRDTs and *Operations-based* CRDTs.

State-based CRDTs define data structures with monotonically increasing states with idempotent and commutative merge functions. This allows one replica to share its local updates by broadcasting its state to others. Upon the reception of the state of another replica, a node is able to update its own state by merging them, regardless of its concurrent updates. Thanks to the properties of the defined state and merge function, states can be missed or delivered multiple times : as long as the most recent state of each replica is successfully broadcast to others once, each node will converge. Thus, no assumptions are made on the network layer. However, this is achieved by broadcasting the whole state repeatedly, which may be inefficient according to the size of the data structure.

Operations-based CRDTs define data structures with a set of operations to perform updates and a partial order between these operations, usually the causal order[5]. In addition, operations have to be designed such that concurrent

operations commute. This allows to propagate local updates by broadcasting corresponding operations to other replicas. Operations are delivered according to the defined partial order. Upon delivery of an operation, a replica updates its state by applying it. In comparison to *State-based* CRDTs, this solution achieves better performances, especially regarding the bandwidth consumption. Nevertheless, it requires the network layer to keep track of the defined partial order, which may be a complex and costly task.

To achieve convergence, *State-based* and *Operations-based* CRDTs proposed in the literature mostly rely on identifiers to reference updated elements. To generate such element identifiers, nodes often use their own identifiers as well as logical clocks. Thus, according to how are generated node identifiers, the size of element identifiers usually increases with the number of nodes. Furthermore, element identifiers have to comply to additional constraints according to the CRDT, for example forming a dense set[1]. In this case, element identifiers' size also increases according to the number of elements contained in the data structure. Therefore, the size of element identifiers is usually not bounded.

Since the size of identifiers is not bounded, the size of metadata attached to each element increases over time. It exceeds more and more the size of data itself. This impedes the adoption of CRDTs since nodes have to broadcast and store metadata, causing the application's performances and efficiency to decrease over time.

This PhD aims to address this issue by 1. proposing more efficient specifications of identifiers according to their set of constraints, 2. proposing mechanisms to rename identifiers to reduce their size.

References

- [1] Luc André, Stéphane Martin, Gérald Oster, and Claudia-Lavinia Ignat. Supporting adaptable granularity of changes for massive-scale collaborative editing. In *International Conference on Collaborative Computing: Networking, Applications and Worksharing - CollaborateCom 2013*, pages 50–59, Austin, TX, USA, October 2013. IEEE Computer Society.
- [2] Eric Brewer. Towards Robust Distributed Systems, 2000.
- [3] The SyncFree Consortium. AntidoteDB: A planet-scale, available, transactional database with strong semantics.
- [4] Matthieu Nicolas, Victorien Elvinger, Gérald Oster, Claudia-Lavinia Ignat, and François Charoy. MUTE: A Peer-to-Peer Web-based Real-time Collaborative Editor. In *ECSCW 2017 - 15th European Conference on Computer-Supported Cooperative Work*, volume 1 of *Proceedings of 15th European Conference on Computer-Supported Cooperative Work - Panels, Posters and Demos*, pages 1–4, Sheffield, United Kingdom, August 2017. EUSSET.

- [5] Ravi Prakash, Michel Raynal, and Mukesh Singhal. An adaptive causal ordering algorithm suited to mobile computing environments. *J. Parallel Distrib. Comput.*, 41(2):190–204, March 1997.
- [6] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. A comprehensive study of Convergent and Commutative Replicated Data Types. Research Report RR-7506, Inria – Centre Paris-Rocquencourt, January 2011.
- [7] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. Conflict-free Replicated Data Types. In *International Symposium on Stabilization, Safety, and Security of Distributed Systems - SSS 2011*, pages 386–400, Grenoble, France, October 2011. Springer.