

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

Matthieu Nicolas  
matthieu.nicolas@inria.fr

Inria, F-54600, Université de Lorraine, LORIA, F-54506, CNRS, F-54506

## Abstract

TODO: Write abstract

In order to serve an ever-growing number of users and provide an increasing volume of data, large scale systems such as data stores or collaborative editing tools have to adopt a distributed architecture. However, as stated by the CAP theorem, such systems cannot ensure both strong consistency and high availability in case of network partitions. As a result, literature and companies increasingly adopt the optimistic replication model paired with the eventual consistency model 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 updates to others. A conflict resolution mechanism is however required to handle updates generated concurrently by different replicas.

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

To achieve convergence, CRDTs proposed in the literature mostly rely on unique identifiers to reference updated elements. To generate such element identifiers, nodes often use their own identifier as well as a logical clock. Thus, regarding to how node identifiers are generated, 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 in case of a *Sequence* data structure. In this case, element identifiers' size also increases according to the updates performed on the data structure. Therefore, the size of element identifiers is usually not bounded.

Since the size of identifiers attached to each element is not bounded, the overhead of the replicated data structure increases over time. Since nodes have to store and broadcast the identifiers, the application's performances and efficiency decrease over time. This impedes the adoption of concerned CRDTs.

We present an approach to address the issue of the growth of identifiers for a family of CRDTs suffering particularly from this issue : *Sequence* CRDTs.

The *Sequence* data structure represents a numbered suite of values. Two operations are usually defined to update the sequence : *insert(index, element)* adds the *element* at the given *index* whereas *delete(index)* removes the element at the given *index*.

This data structure is often used in applications and, over the years, several CRDTs were proposed to represent replicable sequences. To ensure convergence, these data structures attach an identifier to each inserted element. These identifiers are used to achieve transaction-less and commutative updates by uniquely identifying each element and ordering them relatively to each others.

The downside of this approach is the increasing size of the identifiers. Since the identifiers are used to order the elements, they have to form a dense set so that nodes are always able to insert a new element between two others. However, two identifiers of the same size can be contiguous. When inserting a new element between two such identifiers, we have no other choice than to increase the size of the generated

identifier to be able to generate one respecting the intended order. This leads to a specification of identifiers which does not bound their size.

To address the issue of evergrowing identifiers, several approaches were presented. In LSEQ, authors propose several strategies to generate identifiers and design a mechanism to switch between these strategies to limit the speed of the identifiers growth. In Core-Nebula, authors propose a renaming mechanism to reduce the size of currently used identifiers, based on a consensus algorithm. To prevent a system-wide consensus, they divide the system into two tiers : the *Core*, a small set of controlled and stable nodes, and the *Nebula*, an uncontrolled set of nodes. Both can perform updates but only the members *Core* participate in the consensus leading to a rename, while a catch-up mechanism is provided for nodes from the *Nebula* to transform the updates they performed concurrently to a renaming.

TODO: Ajouter une phrase sur le fait qu'effectuer un consensus est coûteux, surtout dans un système très dynamique, pour motiver la conception d'un mécanisme de renommage entièrement distribué

In our approach, we propose a fully distributed renaming mechanism, allowing any node to perform a renaming without any kind of coordination with others.

We define a new operation : *rename*. This operation generates a new, equivalent state of the current sequence but with arbitrary, shorter identifiers. A map, called *renamingMap*, is also generated to link each identifier from the former sequence to the corresponding identifier in the new one. By propagating this map, other nodes are able to apply the same renaming to their own local copy.

However, being a distributed system, nodes can perform updates and broadcast them while another node is performing concurrently the renaming process.

Upon the reception of these concurrent updates, the node which triggered the renaming process cannot apply them to its copy. Indeed, the identifiers having changed, performing the updates would not insert the elements at the correct positions nor delete the correct elements.

To address this issue, we designed rewriting rules for concurrent operations to a *rename*. By using the original identifier and the *renamingMap*, we are either able to find the new corresponding identifier, if it has been renamed, either able to deterministically generate a new identifier preserving the existing order, if it has been concurrently inserted.

However, one *rename* operation can be concurrent to another one. To solve this scenario, we give priority to one operation over the others. Nodes which observed and applied one of the "losing" *renaming* operations have to undo its effect before applying the "winning" one. Thus we designed additional rewriting rules which reverse the effect of a previously applied *rename* while also generating new identifiers preserving the order for identifiers generated after the renaming. To determine which *rename* operation to proceed with, we define and use a total order between *rename* operations.

Based on this approach, we propose a new version of the LogootSplit algorithm implementing this renaming mechanism. We also introduce an optimisation to reduce the bandwidth consumption of the *rename* operation by compressing the *renamingMap*.

The designed renaming mechanism has been implemented in MUTE, the collaborative text editor developed by our research team. It is now in the process of being tested.

However, testing does not ensure correctness. We are thus taking steps to provide a formal proof of the algorithm, either using a proof assistant either using a automatic theorem prover.

TODO: Add introduction to explain that we need to assess the performance of the mechanism

In parallel, we are designing a benchmarker for collaborative editing tools. The goal is be able to replay a collaborative editing session, from its logs, using different conflicts resolution mechanisms to compare their respective performances. We intend to mesure common metrics like the memory and CPU usage, bandwidth consumption but also more specific metrics like the time to propagate and integrate an update and the global time required to converge.

TODO: Write conclusion