

# Entwicklung eines Datenbanksynchronisations-Mikroframeworks auf Basis konvergenter, replizierter Datentypen

Sebastian Götte und Matti Möll

27.10.2014

## 1 History of CRDTs and the current state of research

Data synchronisation has been an important area of research for several decades. Recently, discourse has been heightened due to increasing parallelism of contemporary system architectures. Traditionally, such systems have either relied on a central authority<sup>1</sup> or been largely limited to read-only affairs<sup>2</sup>.

Well-defined automatic merge strategies to reconcile diverging realities in different network locations have been in development for some time now, especially under the aegis of distributed version control systems (DVCSs) such as git, Mercurial and Darcs but have not yet seen widespread deployment.

One possible solution to this problem is what we call Convergent Replicated Data Types (CRDTs). CRDTs have been used by computer scientists for at least 40 years. The concept was only formalized and its name coined five years ago in [4]. Since then, there has been some further research into the formal properties of CRDTs<sup>3</sup> and first real-world implementations of the concept using the term *CRDT* have appeared<sup>4</sup>.

In 2007, Amazon.com published a paper detailing the architecture of their *dynamo* key-value store. From a technical standpoint, Dynamo's replication system is behaving like a CRDT set implementation.

We want to give an overview of the history and existing implementations of CRDTs followed by an overview of some concepts that may be used in conjunction with CRDTs to provide useful higher-level behavior and end with some exemplary possible real-world use cases.

## 2 Convergent replication

CRDTs have first been defined in [4]. [4] are using the term *CRDT* for *Collision-free Replicated Data Type* and are distinguishing between *CvRDTs* (*Convergent Replicated Data Types*) which they are also calling *State-based CRDTs*, and *CmRDTs* (*Commutative Replicated Data Types*) which they are also calling *Operation-based* or *Op-based CRDTs*. Since both concepts are shown

---

<sup>1</sup>See every RDBMS in existence, OpenLDAP and ActiveDirectory, RSS and apt among others

<sup>2</sup>See e.g. BitTorrent and other classical file sharing systems

<sup>3</sup>See [5]

<sup>4</sup>See e.g. [6] resp. [8] and [7] and [2]

to be exactly equivalent, we are using the acronym *CRDT* for either approach and settled on the long form *Convergent Replicated Data Type* since it is describing both very well and is sufficiently handy.

## 2.1 Operation-based CRDTs

An op-based CRDT is a data type whose value is defined solely by a set of commutative operations applied to a common base state. For illustration, consider an up/down counter as a simple example. The initial state would be 0 and the operations would be to add or subtract a number. The current counter value is defined as the sum of all add/subtract operations performed on this instance so far. Since addition is commutative, the order of these operations does not matter. In a distributed system, additions and subtractions on the same counter can be performed simultaneously on multiple nodes, and the resulting conflict can be resolved by each node telling each other node all operations that have been performed on its instance.

## 2.2 State-based CRDTs

A state-based CRDT is a data type which is associated with a partial ordering on its value space. A simple example for a state-based CRDT is an add-only set. The associated partial ordering is the subset relation. The merge operation is the set union.

State- and Operation-based CRDTs are equivalent in that each can be used to emulate the other. Most existing implementations use a state-based storage, where some do incorporate some op-based-like behavior for more efficient sync<sup>5</sup>.

## 3 Possible directions for research

## 4 Possible use cases

## Literatur

- [1] Scott Aaronson. Eigenmorality. <http://www.scottaaronson.com/blog/?p=1820&re=1>, 2014. [Online; accessed Oct 27 2014].
- [2] Peter Bourgon, Tomás Senart, Björn Rabenstein, and Johan Uhle. Roshi: a crdt system for timestamped events. <https://developers.soundcloud.com/blog/roshi-a-crdt-system-for-timestamped-events>, 2014. [Online; accessed Oct 27 2014].
- [3] DeCandia, Hastorun, Jampani, Kakulapati, Lakshman, Pilchin, Sivasubramanian, Voshall, and Vogels. Dynamo: Amazon’s highly available key-value store. 2007.
- [4] Mihai Letia, Nuno Preguiça, and Marc Shapiro. Crdts: Consistency without concurrency control. 2009.
- [5] Marc Shapiro, Nuno Preguica, Carlos Baquero, and Marek Zawirski. A comprehensive study of convergent and commutative replicated data types. 2011.

---

<sup>5</sup>e.g. Amazon Dynamo is computing a delta using a Merkle tree and then only transmitting the subset of changed entities, which is equivalent to the set of operations since last launch.

- [6] Basho Technologies. Riak. <http://basho.com/riak/>, 2009. [Online; accessed Oct 27 2014].
- [7] Basho Technologies. Distributed data types in riak 2.0. <http://basho.com/distributed-data-types-riak-2-0/>, 2014. [Online; accessed Oct 27 2014].
- [8] Basho Technologies. Riak 2.0 announcement. <http://basho.com/riak-2-0-new-capabilities-new-use-cases-available-for-download/>, 2014. [Online; accessed Oct 27 2014].