

Design and Verification of Byzantine Fault Tolerant CRDTs

Liangrun Da

School of Computation, Information, and Technology
Technical University of Munich

November 15, 2024

- ▶ Wikipedia is currently a centralized system and needs funding to maintain the server.
- ▶ A decentralized alternative can be run by volunteers without any central server, hence no funding is needed.
- ▶ However, decentralized collaboration systems face a major problem: **no consistency guarantee**.

What are Conflict-free Replicated Data Types?

RGA as an Example



Hllo

p

Hllo

r

Hllo

q

What are Conflict-free Replicated Data Types?

RGA as an Example

Hello

p

⟨Insert “e” after “H” ⟩

Hllo

r

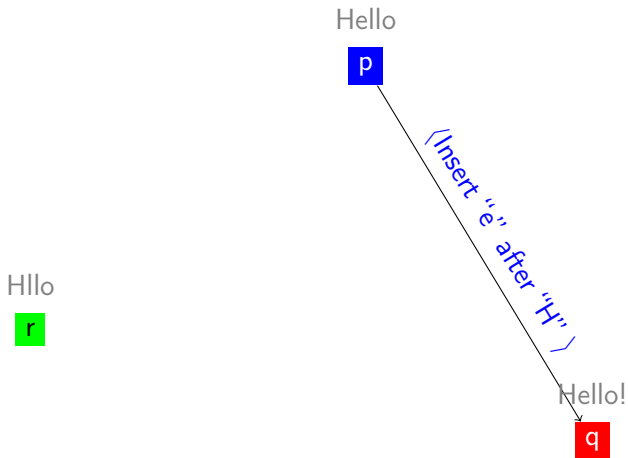
Hllo!

q

⟨Insert “!” after “o” ⟩

What are Conflict-free Replicated Data Types?

RGA as an Example



What are Conflict-free Replicated Data Types?

RGA as an Example

r

Hello!

p

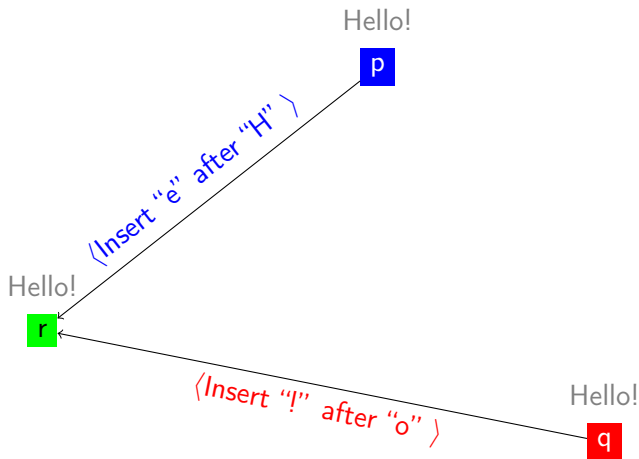
<Insert "!" after "o">

Hello!

q

What are Conflict-free Replicated Data Types?

RGA as an Example



- ▶ Eventual Delivery: an update delivered at some correct node is eventually delivered at all correct nodes
- ▶ Convergence: correct replicas that have delivered the same set of updates have equivalent state.
 - ▶ Operations are delivered in causal order.
 - ▶ Concurrent operations commute.
- ▶ Termination: All method executions terminate.

Termination is usually easy to guarantee, so we will focus on the other two.

¹M Shapiro et al. "Conflict-free replicated data types". In: *Stabilization, Safety, and Security of Distributed Systems: 13th International Symposium, SSS 2011, Grenoble, France, October 10-12, 2011. Proceedings 13*. Springer. 2011, pp. 386–400.

- ▶ CRDTs are widely believed to be the solution for replication in peer-to-peer systems because it doesn't require a central server.

- ▶ CRDTs are widely believed to be the solution for replication in peer-to-peer systems because it doesn't require a central server.
- ▶ Most CRDTs algorithms assume participating peers strictly follows the protocol, i.e. they are not Byzantine fault tolerant (BFT).

- ▶ CRDTs are widely believed to be the solution for replication in peer-to-peer systems because it doesn't require a central server.
- ▶ Most CRDTs algorithms assume participating peers strictly follows the protocol, i.e. they are not Byzantine fault tolerant (BFT).
- ▶ However, open peer-to-peer systems allow anyone to join or leave, hence the assumption is not safe.

- ▶ BFT Consensus requires a total order broadcast, which is unnecessary overhead for CRDTs.

- ▶ BFT Consensus requires a total order broadcast, which is unnecessary overhead for CRDTs.
- ▶ BFT Consensus provably requires the Byzantine peers to be less than $1/3$ of the total peers, where as BFT CRDTs should be able to tolerate any number of Byzantine peers.

- ▶ BFT Consensus requires a total order broadcast, which is unnecessary overhead for CRDTs.
- ▶ BFT Consensus provably requires the Byzantine peers to be less than $1/3$ of the total peers, where as BFT CRDTs should be able to tolerate any number of Byzantine peers.
- ▶ BFT Consensus and BFT CRDTs are fundamentally different.

Eventual Delivery

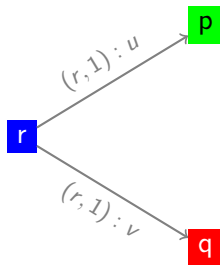
Version vectors are not safe



Eventual Delivery

Version vectors are not safe

- ▶ Peer A sends two different updates with the same version vector.



Eventual Delivery

Version vectors are not safe

- ▶ Peer A sends two different updates with the same version vector.
- ▶ Peer B and C exchange their version vectors and believe that their updates are the same.

r

(r:1, p:0, q:0)

{u}

p



exchange

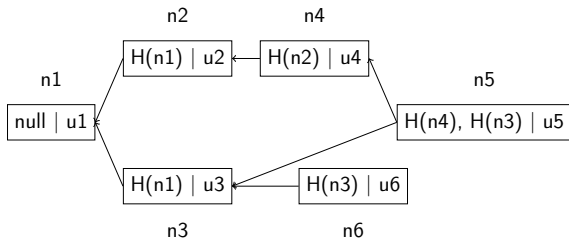
q

(r:1, p:0, q:0)

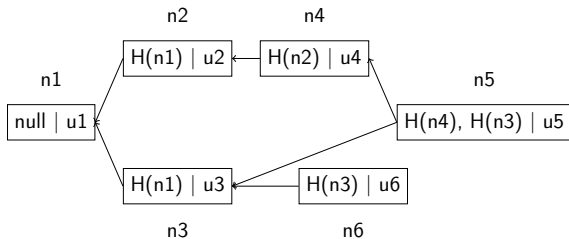
{v}

- ▶ If a peer periodically sends updates it doesn't receive from another peer, and vice versa, the eventual delivery is guaranteed.
- ▶ However it is highly inefficient.

- ▶ A hash DAG is a directed acyclic graph where each node has a value (i.e. an update in our context) and a set of hashes, which resolve to its predecessors. (like Git)



- ▶ A hash DAG is a directed acyclic graph where each node has a value (i.e. an update in our context) and a set of hashes, which resolve to its predecessors. (like Git)
- ▶ The ancestor relation between the operations in the hash DAG reflects the causal dependency between the operations.



Theorem

The heads of a hash DAG is the set of nodes with no successors. If two peers have the same set of heads, their hash graphs, including CRDT operations contained in the hash graphs must be same.

- ▶ Traditionally, a CRDT operation's validity is checked before it is broadcast to the peers. However, a malicious peer can skip the validity check to cause problems.

Hllo

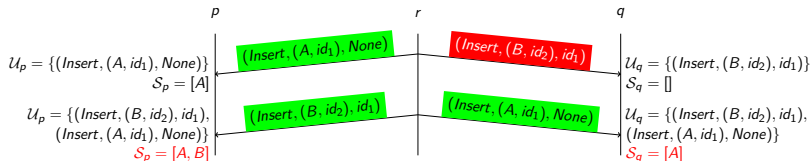
p

⟨Insert "e" after "X" ⟩ **invalid!**

- ▶ Traditionally, a CRDT operation's validity is checked before it is broadcast to the peers. However, a malicious peer can skip the validity check to cause problems.
- ▶ Instead of relying on validity at the sender side, we shift the validity check responsibility to receiver side.

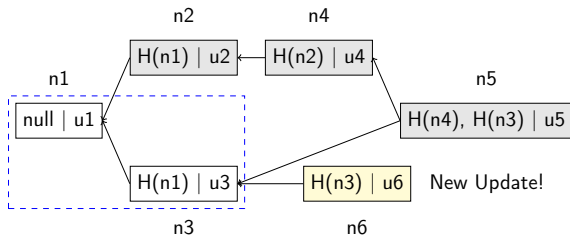
However, naively moving the validity check to the receiver side could cause inconsistent states.

$(Insert, (A, id_1), id_2)$ means insert A with ID id_1 after the element with ID id_2 . If id_2 is $None$, then A is inserted as the first element.



- ▶ If we can ensure that all peers make the same validity decision, then we can avoid such problems.

- ▶ If we can ensure that all peers make the same validity decision, then we can avoid such problems.
- ▶ A way to achieve this is to have each peer only validate operations based on the ancestors of the operation in the hash DAG. This makes the validity check deterministic since all peers share the same ancestors of the operation.



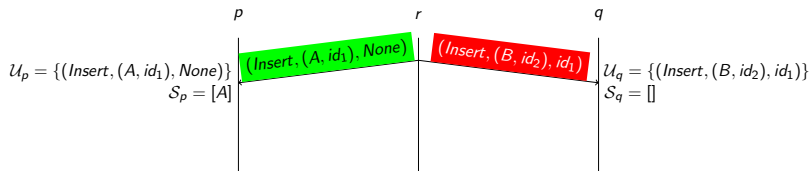
Theorem

If all correct peers follow the protocol, and each operation's validity is checked only against its ancestors in the hash DAG, then a valid operation on one correct peer must also be a valid operation on all other correct peers.

Convergence

Use ancestor-only checks

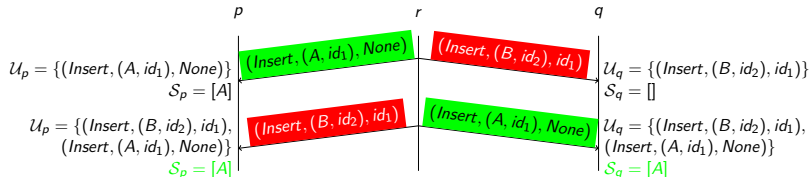
Even though p has element with ID id_1 , it still reject the update because the element with ID id_1 is not an ancestor of $(Insert, (B, id_2), id_1)$.



Convergence

Use ancestor-only checks

Even though p has element with ID id_1 , it still reject the update because the element with ID id_1 is not an ancestor of $(Insert, (B, id_2), id_1)$.

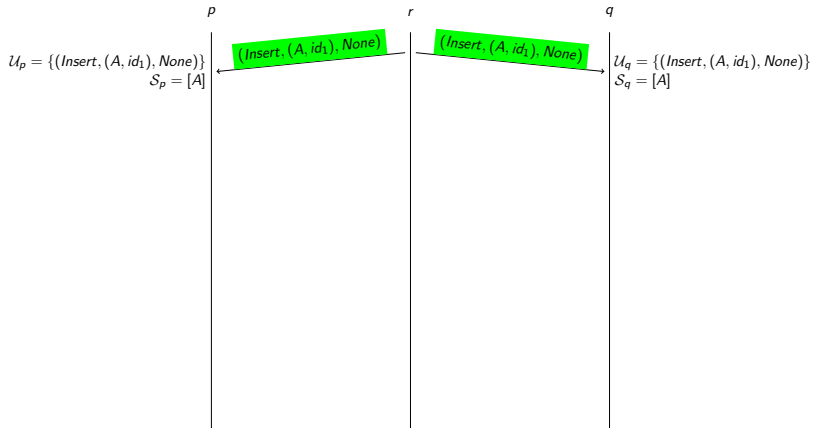


Many CRDT algorithms require unique IDs². It is often achieved by (unique peer-id, counter) traditionally, but a malicious peer can simply use the same (unique peer-id, counter) for different operations, resulting in divergence.

²Including ORSet, RGA, Logoot, Treedoc, etc.

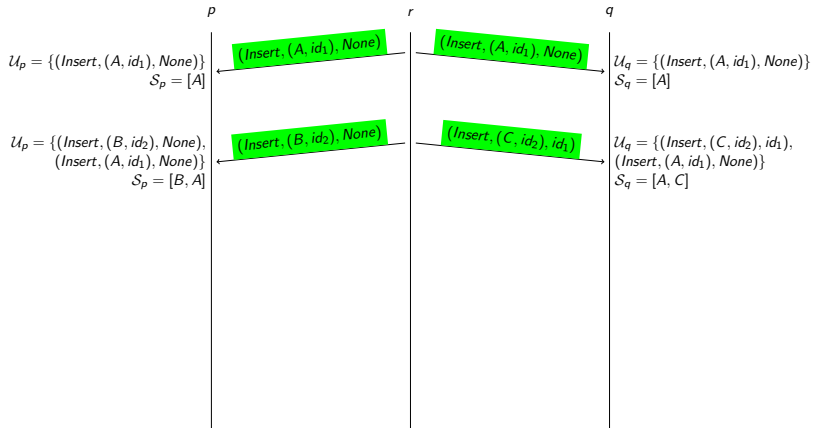
Convergence

RGA as an example



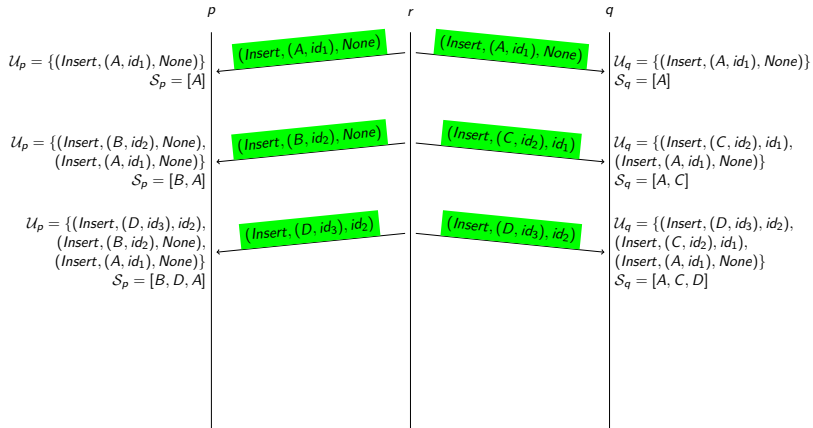
Convergence

RGA as an example



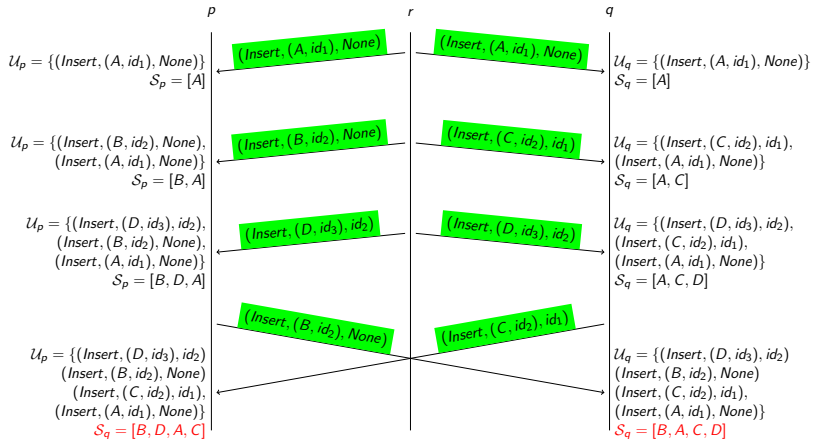
Convergence

RGA as an example



Convergence

RGA as an example



- ▶ Instead of relying on the sender to provide an unique ID, we let each peer generate the ID by using the hash of the node that contains the operation.
- ▶ Since the hash function is collision-resistant, and the nodes in the hash DAG are unique, the IDs are unique.

- ▶ How to prove the theorems we stated in the previous sections in Byzantine environment?
- ▶ How to ensure we have ruled out all possible vulnerabilities for a particular CRDT?

- ▶ We model a peer-to-peer system consisting of indefinitely many peers in Isabelle/HOL, a formal proof assistant.

- ▶ We model a peer-to-peer system consisting of indefinitely many peers in Isabelle/HOL, a formal proof assistant.
- ▶ The system only assumes a collision-resistant hash function, and we place no other assumptions on the system.

- ▶ We model a peer-to-peer system consisting of indefinitely many peers in Isabelle/HOL, a formal proof assistant.
- ▶ The system only assumes a collision-resistant hash function, and we place no other assumptions on the system.
- ▶ We proved the theorems that we claimed previously in Isabelle/HOL.

- ▶ We model a peer-to-peer system consisting of indefinitely many peers in Isabelle/HOL, a formal proof assistant.
- ▶ The system only assumes a collision-resistant hash function, and we place no other assumptions on the system.
- ▶ We proved the theorems that we claimed previously in Isabelle/HOL.
- ▶ Using those theorems, one can verify the correctness of the resulting BFT CRDTs by proving some simple properties.³

³L Da and M Kleppmann. *A Framework for Designing and Verifying Byzantine Fault Tolerant CRDTs*. Version 0.0.1. <https://github.com/LiangrunDa/bft-crdt-isabelle>. Oct. 2024.

To prove a BFT CRDT correct, one needs to prove the following three properties:

- ▶ Concurrent operations commute.
- ▶ The validity check only depends on the ancestors of the operation.
- ▶ It never fails on a valid operation.

Then the following theorems are automatically proved:

theorem *sec-convergence*:

assumes $\langle \text{heads } (\text{graph } i) = \text{heads } (\text{graph } j) \rangle$

shows $\langle \text{apply-operations } (\text{delivered-nodes } i) = \text{apply-operations } (\text{delivered-nodes } j) \rangle$

theorem *sec-progress*: $\langle \text{apply-operations } (\text{delivered-nodes } i) \neq \text{None} \rangle$

- ▶ We applied our method to two well-known CRDTs, *RGA* and *ORSet*, and obtained two BFT CRDTs.
- ▶ We proved the correctness of the two BFT CRDTs formally using only 244 and 522 LoC respectively.

- ▶ We analyzed the possible vulnerabilities of traditional CRDTs under Strong Eventual Consistency model.

- ▶ We analyzed the possible vulnerabilities of traditional CRDTs under Strong Eventual Consistency model.
- ▶ We proposed several approaches to prevent the vulnerabilities.

- ▶ We analyzed the possible vulnerabilities of traditional CRDTs under Strong Eventual Consistency model.
- ▶ We proposed several approaches to prevent the vulnerabilities.
- ▶ We formalized the system in Isabelle/HOL and proved the correctness of the proposed approaches.

- ▶ We analyzed the possible vulnerabilities of traditional CRDTs under Strong Eventual Consistency model.
- ▶ We proposed several approaches to prevent the vulnerabilities.
- ▶ We formalized the system in Isabelle/HOL and proved the correctness of the proposed approaches.
- ▶ Our framework along with the theorems can be used to verify the correctness of a BFT CRDT by proving some simple properties. We evaluated our framework on ORSet and RGA.

locale *bft-strong-eventual-consistency* = *peers-with-arbitrary-history* +

assumes *sem-check-only-ancestors-relevant*:

$\langle (\text{ancestor-nodes-of } n) \subseteq \text{fset } G \implies \text{is-sem-valid-set}$

$(\text{ancestor-nodes-of } n) \ n \longleftrightarrow \text{is-sem-valid } G \ n \rangle$

assumes *concurrent-ops-commute*: $\langle \text{hb.concurrent-ops-commute}$
(delivered-nodes i) \rangle

assumes *step-never-fails*: $\langle \text{apply-history} ([], \{||\}) \ ns = (dn, G) \implies$
no-failure dn \implies

$\text{check-and-apply} (dn, G) (hs, v) = (dn', G') \implies \text{no-failure } dn' \rangle$