

CMCS 2024

---

# Coinductive Reasoning about CRDT Emulation

---

**Nathan E. Liittschwager,**  
Stelios Tsampas,  
Jonathan Castello,  
Lindsey Kuper

# Background: State Machine Replication

Imagine a *service* offered to a *client* in the form of a black-box with I/O behavior and internal state

The client inputs are *requests* in the form of *commands*  $a \in A$  which update the internal state.

Client observations  $b \in B$  are accessible by a *query* method call

Implement the service as a *state machine* (coalgebra):

$$(u, q) : X \rightarrow X^A \times B$$

# Background: State Machine Replication

To achieve *fault tolerance* the state machine  $(u, q) : X \rightarrow X^A \times B$  is *replicated* on  $n$  servers or nodes:  $(u_1, q_1), \dots, (u_n, q_n) : X \rightarrow X^A \times B$  with same initial state  $x \in X$

Client input seen as a totally ordered sequence  $(a_1, \dots, a_k) \in A^*$

Each replica  $(x \in X, (u_i, q_i))$  must compute the same sequence  $(a_1, \dots, a_k) \in A^*$ , thus obtain the same state (*linearizability*)

**Fault tolerance achieved:** if one replica goes down, another replica can take its role

# Background: State Machine Replication

To achieve *fault tolerance* the state machine  $(u, q) : X \rightarrow X^A \times B$  is replicated on  $n$  servers or nodes:  $(u_1, q_1), \dots, (u_n, q_n) : X \rightarrow X^A \times B$  with same initial state  $x \in X$

Client input seen as a totally ordered sequence  $(a_1, \dots, a_k) \in A^*$

Each replica  $(x \in X, (u_i, q_i))$  must compute the same sequence  $(a_1, \dots, a_k) \in A^*$ , thus obtain the same state (*linearizability*)

**Fault tolerance achieved:** if one replica goes down, another replica can take its role

# Background: State Machine Replication

To achieve *fault tolerance* the state machine  $(u, q) : X \rightarrow X^A \times B$  is replicated on  $n$  servers or nodes:  $(u_1, q_1), \dots, (u_n, q_n) : X \rightarrow X^A \times B$  with same initial state  $x \in X$

Client input seen as a totally ordered sequence  $(a_1, \dots, a_k) \in A^*$

Non-trivial!

Each replica  $(x \in X, (u_i, q_i))$  must compute the same sequence  $(a_1, \dots, a_k) \in A^*$ , thus obtain the same state (*linearizability*)

**Fault tolerance achieved:** if one replica goes down, another replica can take its role

# Background: CRDTs

Totally ordering commands  $(a_1, \dots, a_k) \in A^*$  non-trivial in asynchronous distributed systems

Conflict-free Replicated Data Types (CRDTs) solve this problem by using data structures which don't require total order

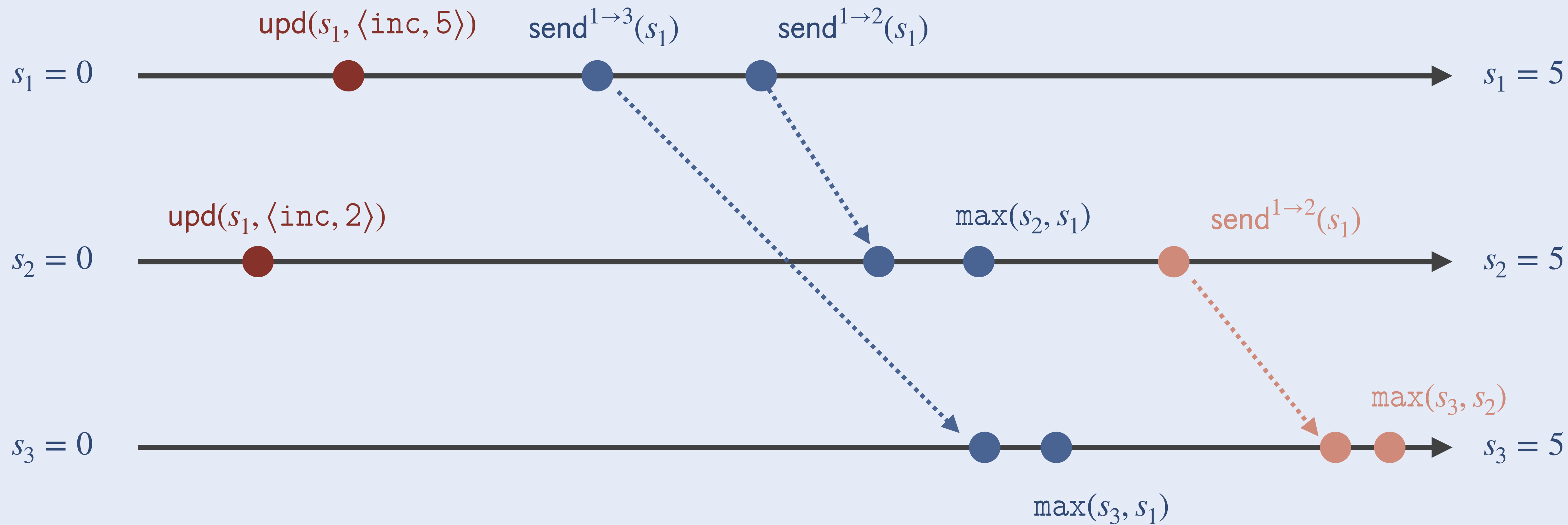
Two major flavors of CRDT: “State-based” and “Operation-based”

Both achieve *strong convergence* - if they know about the same set of messages, they have the same state.

# Background: State-Based CRDTs

State-based CRDTs use a join-semilattice  $(S, \sqcup)$  as the state-space

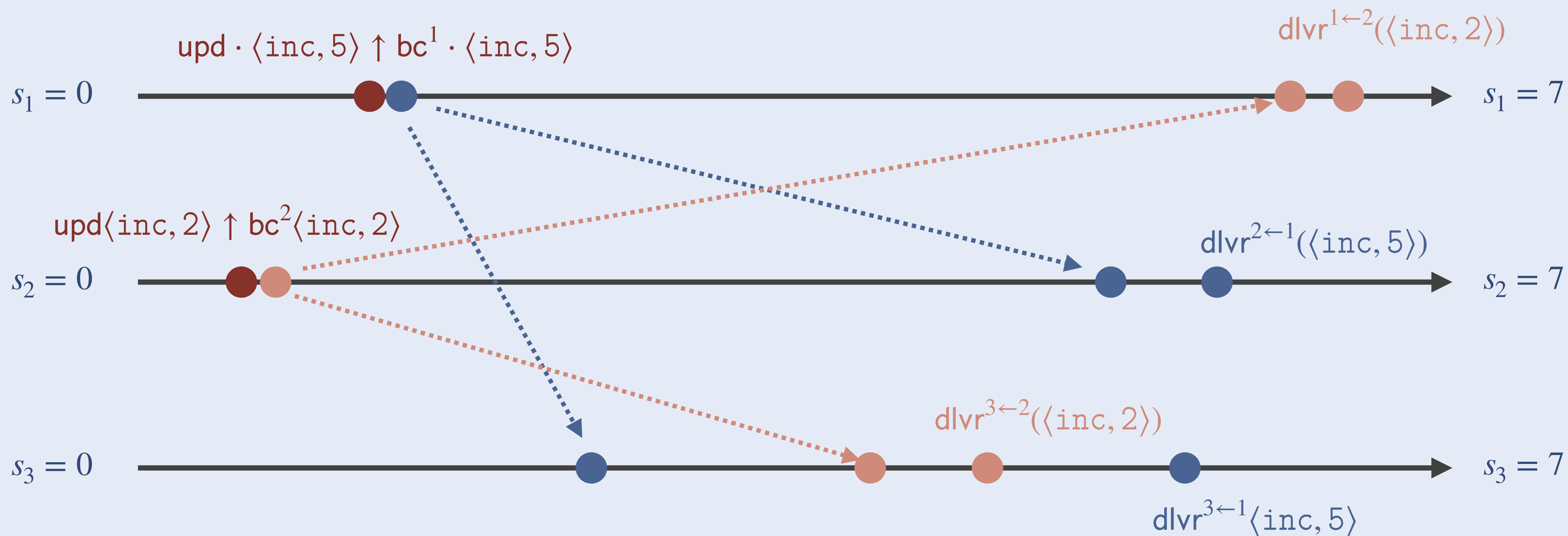
Updates must be *inflationary*:  $s \sqcup \text{upd}(s, a) = \text{upd}(s, a)$  for all  $s \in S, a \in A$



# Background: Operation-based CRDTs

Operation-based CRDTs require messages (operations) be a *partial order*  $(M, <)$

Each replica executes the same set of messages in a way consistent with  $<$ ,  
communication is by *broadcast*





# Background: CRDT Emulation and Equivalence

State-based and op-based CRDTs are often considered to be *equivalent*.

The reasoning is that they “*emulate*” each other: there are a pair of maps  $\mathcal{F}, \mathcal{G}$  to translate between the two types [Shapiro et al. 2011]

## Definition (Emulation Maps)

State-based CRDT:  $((S, \sqcup), s_0, u, q)$       Op-based CRDT:  $(S, s_0, M, u, \mathfrak{t}, e, q)$

$$((S, \sqcup), s_0, u, q) \xrightarrow{\mathcal{F}} (S, s_0, S, u, u, \sqcup, q)$$

$$(S, s_0, M, u, \mathfrak{t}, e, q) \xrightarrow{\mathcal{G}} ((\mathcal{P}_{fin}(M), \cup), \emptyset, u', q')$$

effect-msg  
prepare-msg

where  $u'(H, a) = H \cup \{\mathfrak{t}(\llbracket H \rrbracket, a)\}$ , and  $q'(H) = q(\llbracket H \rrbracket)$

and  $\llbracket \cdot \rrbracket : \mathcal{P}_{fin}(M) \rightarrow S$  is appropriate map to translate sets of messages into state

# Background: CRDT Emulation and Equivalence

Maps  $\mathcal{F}$  and  $\mathcal{G}$  are intuitively, potentially correct

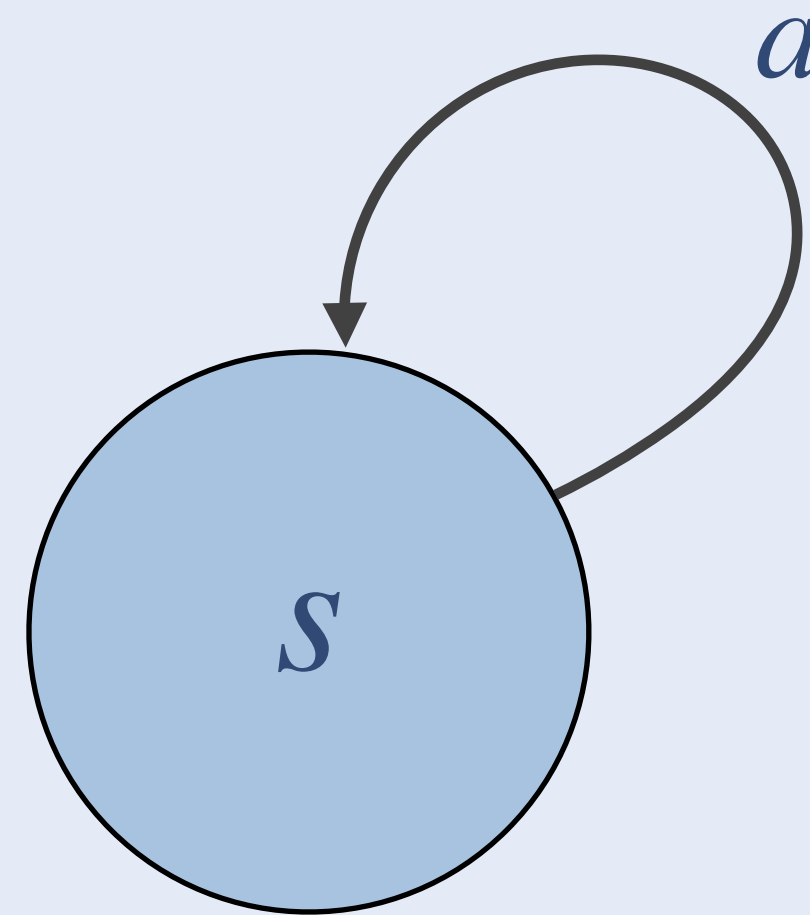
But, emulation is not rigorously defined: no formal requirements on behavior, only informal arguments about strong eventual convergence

What if we defined  $\mathcal{F}$  to map each replica to a trivial state machine?

$$S = \{s\}$$

$$\forall a \in A . u(s)(a) = s$$

$$q(s) = \top$$



Convergence!

# Strong Bisimulation?

What if we required the original CRDT and the  $\mathcal{F}$ -emulator (or  $\mathcal{G}$ -emulator) to exhibit a bisimulation?

For example, if  $\mathcal{G}$  were a coalgebra homomorphism - Very strong notion of “equivalence”!

But is there even such a bisimulation?

# Strong Bisimulation?

What if we required the original CRDT and the  $\mathcal{F}$ -emulator (or  $\mathcal{G}$ -emulator) to exhibit a bisimulation?

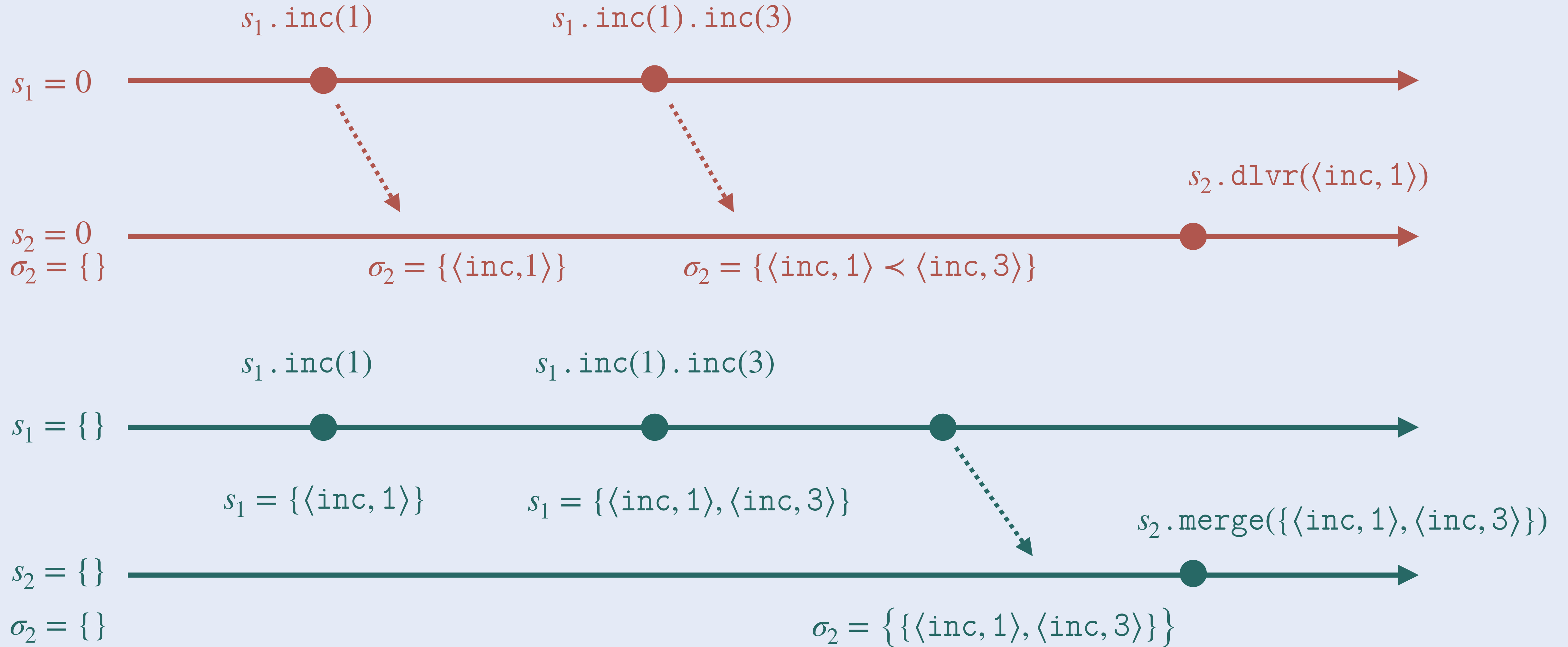
For example, if  $\mathcal{G}$  were a coalgebra homomorphism - Very strong notion of “equivalence”!

But is there even such a bisimulation? NO.

The semantics of Op-based CRDTs treat  $\text{upd}(a) \uparrow \text{bc}^i(m)$  events as *atomic* (uninterruptible)

But there is **no** such requirement for  $\text{upd}(a)$  and  $\text{send}^{i \rightarrow j}(s)$  events on state-based CRDTs

# Bisimulation Game



# Bisimulation Game

Let  $\llbracket \cdot \rrbracket : \mathcal{P}_{fin}(M) \rightarrow S$  be a map which interprets sets of messages to states

Define the query map  $q_{op} = id_S$

Define the query map  $q_{st} = id_S \circ \llbracket \cdot \rrbracket$

$$s_1 . inc(1) . inc(3) \implies q(s_1) = 4$$

$$s_2 . dlvr(\langle inc, 1 \rangle) \implies q(s_2) = 1$$

$$s_1 . inc(1) . inc(3) \implies q'(s_1) = \llbracket \{ \langle inc, 1 \rangle, \langle inc, 3 \rangle \} \rrbracket = 4$$

$$s_2 . merge(\{ \langle inc, 1 \rangle, \langle inc, 3 \rangle \}) \implies q'(s_2) = \llbracket \{ \langle inc, 1 \rangle, \langle inc, 3 \rangle \} \rrbracket = 4$$

# Background: CRDT Emulation and Equivalence

Despite this, the notion of emulation is “load bearing” in the CRDT literature

*“...our techniques [on op-based CRDTs] naturally extends to state-based CRDTs since they can be emulated by an op-based model...”* - [Nagar et al., 2019]

*“... [our work on synthesis of state-based CRDTs]... can always be translated to op-based CRDTs if necessary...”* - [Laddad et al., 2022]

Our contributions: We close this gap by showing that  $\mathcal{F}$  and  $\mathcal{G}$  induce a pair of *weak simulations* between the original CRDT and its “emulator”

# Coinductive Reasoning about CRDT Emulation

## Definition (Weak Simulation)

Let  $(X, (h, \text{obs}_1))$  and  $(Y, (g, \text{obs}_2))$  be coalgebras of endofunctor  $\mathcal{P}(-) \times B$  and let  $g^* : Y \rightarrow \mathcal{P}(Y)$  be the *reflexive, transitive closure* of  $g$ .

A *weak simulation* of  $(X, h)$  and  $(Y, g)$  is a relation  $R \subseteq X \times Y$  s.t.

$\forall (x, y) \in X \times Y$ . if  $(x, y) \in R$ , then

1.  $\text{obs}_1(x) = b \implies \text{obs}_2(y) = b$
2.  $x' \in h(x) \implies \exists y' \in g^*(y) \wedge (x', y') \in R$



# Coinductive Reasoning about CRDT Emulation

## Definition (Op-based CRDT Systems)

$$\frac{\alpha \notin \text{update} \quad x_j \longrightarrow_{\text{op}} x'_j \uparrow (a, m)}{\langle \alpha, (x_i)_{i \in n} \rangle \rightsquigarrow_{\text{op}} \langle \text{upd}^j(a, m), (x_i)_{i \in n} [x_j \leftarrow x'_j] \rangle} \quad [\text{OpUpdate}]$$

$$\frac{(x'_i)_{i \in n} = \text{bcast}_m^j(x_i)_{i \in n}}{\langle \text{upd}^j(a, m), (x_i)_{i \in n} \rangle \rightsquigarrow_{\text{op}} \langle \text{bc}^j(m), (x'_i)_{i \in n} \rangle} \quad [\text{OpBroadcast}]$$

$$\frac{\alpha \notin \text{update} \quad x_j \longrightarrow_{\text{op}} x'_j \text{ via deliver } m}{\langle \alpha, (x_i)_{i \in n} \rangle \rightsquigarrow_{\text{op}} \langle \text{dlvr}^j(m), (x_i)_{i \in n} [x_j \leftarrow x'_j] \rangle} \quad [\text{OpDeliver}]$$

(Events:  $\alpha := \top \mid \text{upd}^i(a, m) \mid \text{bc}^i(m) \mid \text{dlvr}^i(m)$ )

# Coinductive Reasoning about CRDT Emulation

## Definition (State-based CRDT Systems)

$$x_j = (s_j, \sigma_j) \quad x_j \longrightarrow_{\text{st}} x'_j \quad (u(s_j, a), \sigma_j) = x'_j$$

$$\langle \alpha, (x_i)_{i \in n} \rangle \rightsquigarrow_{\text{st}} \langle \text{upd}^j(a), (x_i)_{i \in n} [x_j \leftarrow x'_j] \rangle$$

[StUpdate]

$$x_i = (s_i, \sigma_i) \quad x_j = (s_j, \sigma_j) \quad x'_j = (s_j, \sigma_j \cup \{s_i\})$$

$$\langle \alpha, (x_i)_{i \in n} \rangle \rightsquigarrow_{\text{st}} \langle \text{send}^{i \rightarrow j}(s_i), (x_i)_{i \in n} [x_j \leftarrow x'_j] \rangle$$

[StSend]

$$x_j = (s_j, \sigma_j) \quad x_j \longrightarrow_{\text{st}} x'_j \quad (s_j \sqcup s, \sigma_j \setminus \{s\}) = x'_j$$

$$\langle \alpha, (x_i)_{i \in n} \rangle \rightsquigarrow_{\text{st}} \langle \text{dlvr}^j(s), (x_i)_{i \in n} [x_j \leftarrow x'_j] \rangle$$

[StDeliver]

(Events:  $\alpha := \top \mid \text{upd}^i(a) \mid \text{send}^{i \rightarrow j}(s) \mid \text{dlvr}^i(s)$ )

# Weak Simulation (State-based $\rightarrow$ Op-based)

## Theorem (Weak Simulation)

Let  $(\rightsquigarrow_{\text{st}}, q_{\text{st}})$  be the state-based CRDT system for  $c = ((S, \sqcup), s_0, u, q)$  and  $(\rightsquigarrow_{\text{op}}, q_{\text{op}})$  the op-based *emulator* CRDT system for  $\mathcal{F}(c)$ . There are a pair of weak simulations  $Q_1$  and  $Q_2$  such that,

1.  $Q_1$  is a weak simulation for  $(\rightsquigarrow_{\text{op}}, q_{\text{op}})$  and  $(\rightsquigarrow_{\text{st}}, q_{\text{st}})$
2.  $Q_2$  is a weak simulation for  $(\rightsquigarrow_{\text{st}}, q_{\text{st}})$  and  $(\rightsquigarrow_{\text{op}}, q_{\text{op}})$

# Weak Simulation (Op-based $\rightarrow$ State-based)

## Theorem (Weak Simulation)

Let  $(\rightsquigarrow_{\text{op}}, q_{\text{op}})$  be the op-based CRDT system for  $o = (S, s_0, M, u, t, e, q)$  and  $(\rightsquigarrow_{\text{st}}, q_{\text{st}})$  the state-based *emulator* CRDT system for  $\mathcal{G}(o)$ . There are a pair of weak simulations  $R_1$  and  $R_2$  such that,

1.  $R_1$  is a weak simulation for  $(\rightsquigarrow_{\text{op}}, q_{\text{op}})$  and  $(\rightsquigarrow_{\text{st}}, q_{\text{st}})$
2.  $R_2$  is a weak simulation for  $(\rightsquigarrow_{\text{st}}, q_{\text{st}})$  and  $(\rightsquigarrow_{\text{op}}, q_{\text{op}})$

## Future Work

- It would be interesting to try to capture this notion of *emulation* (translation + simulation) in higher generality. Perhaps in different categories (e.g., Kleisli), perhaps with different systems.
- Our model is based on more classical distributed systems theory (“vectors” of transition systems) but is nonetheless coalgebraic. A general coalgebraic treatment of distributed systems theory would be interesting - the challenge is reconciling the different notions of “simulation” under coalgebraic lens
- What are the interesting properties preserved by a translation + simulation?

# References

1. Shadaj Laddad et al. “Katara: Synthesizing CRDTs with Verified Lifting”. In: Proc. ACM Program. Lang. 6.OOPSLA2 (2022). doi: 10.1145/3563336. url: <https://doi.org/10.1145/3563336>.
2. Marc Shapiro et al. “Conflict-Free Replicated Data Types”. In: Stabilization, Safety, and Security of Distributed Systems. Ed. by Xavier Défago, Franck Petit, and Vincent Villain. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 386–400. isbn: 978-3-642-24550-3.
3. Kartik Nagar and Suresh Jagannathan. “Automated Parameterized Verification of CRDTs”. In: Computer Aided Verification. Ed. by Isil Dillig and Serdar Tasiran. Cham: Springer International Publishing, 2019, pp. 459–477. isbn: 978-3-030-25543-5.

Questions?