# CRDT transactions in a scalable way* with SwiftCloud
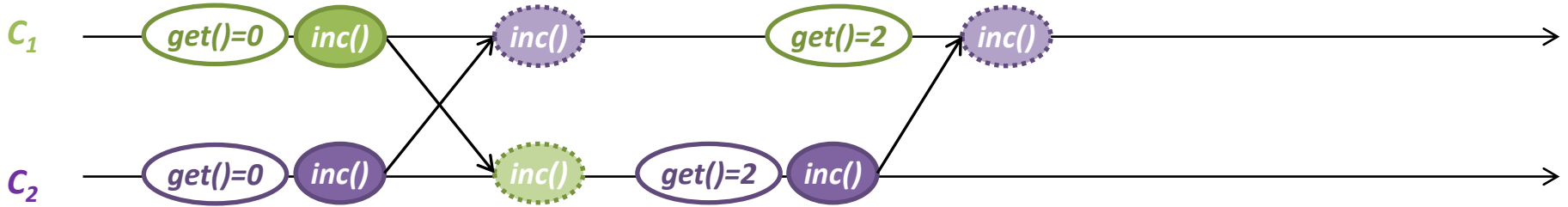
| | |
|---|---|
| Marek Zawirski | UPMC-LIP6 & INRIA |
| Annette Bieniusa | U. Kaiserslautern |
| Valter Balegas | UNL |
| Nuno Preguiça | UNL |
| Sérgio Duarte | UNL |
| Marc Shapiro | INRIA & UPMC-LIP6 |
| Carlos Baquero | U. Minho |

# Object model and transactional guarantees
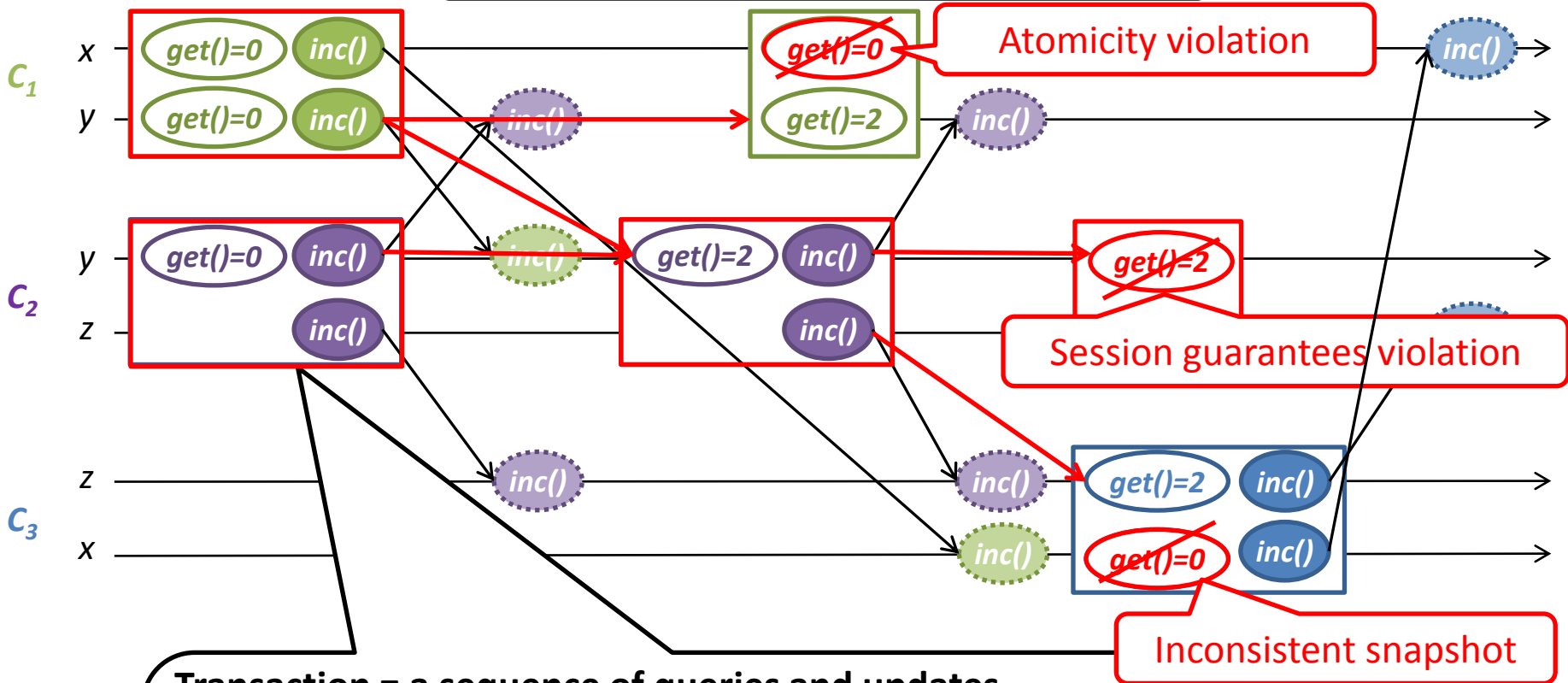
Example: a replicated counter object (CRDT)



$C_1$

get()=0  inc()  inc()  get()=2  inc()

$C_2$

get()=0  inc()  inc()  get()=2  inc()

**A simplified model of object replica**

inc() ----> inc() ----> inc()

- Log of updates, a linear extension of "happened-before" order
- Object version = initial state + application of a valid subsequence
- Suboptimal, later extended with pruning

# Object model and transactional guarantees

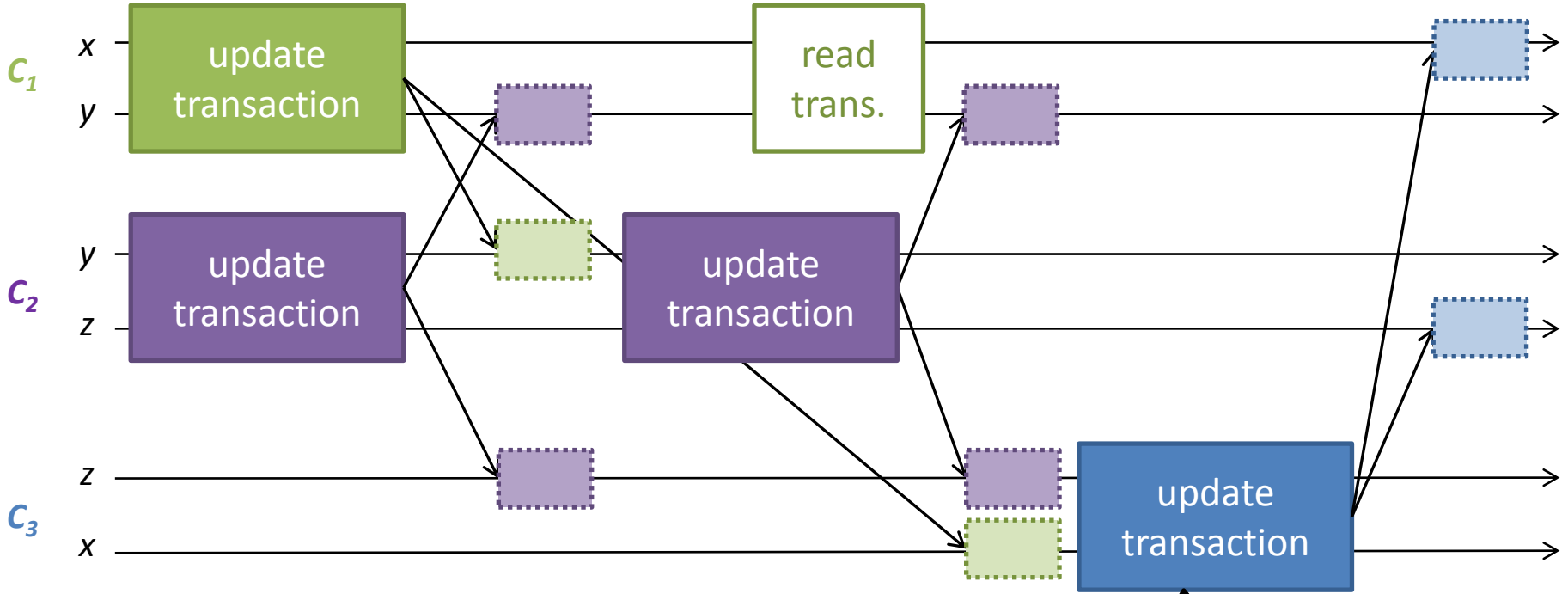Example: a database of replicated counters



**Atomicity violation**

**Session guarantees violation**

**Inconsistent snapshot**

**Transaction = a sequence of queries and updates**
- *Atomicity* => all-or-none updates visible
- *Consients snapshot for reads* => consistent cut based on causality
  - Causality extended across objects through each transaction
- *Session guarantees* => growing snapshot including prior updates
- *No aborts* => asynchronous implementation possible
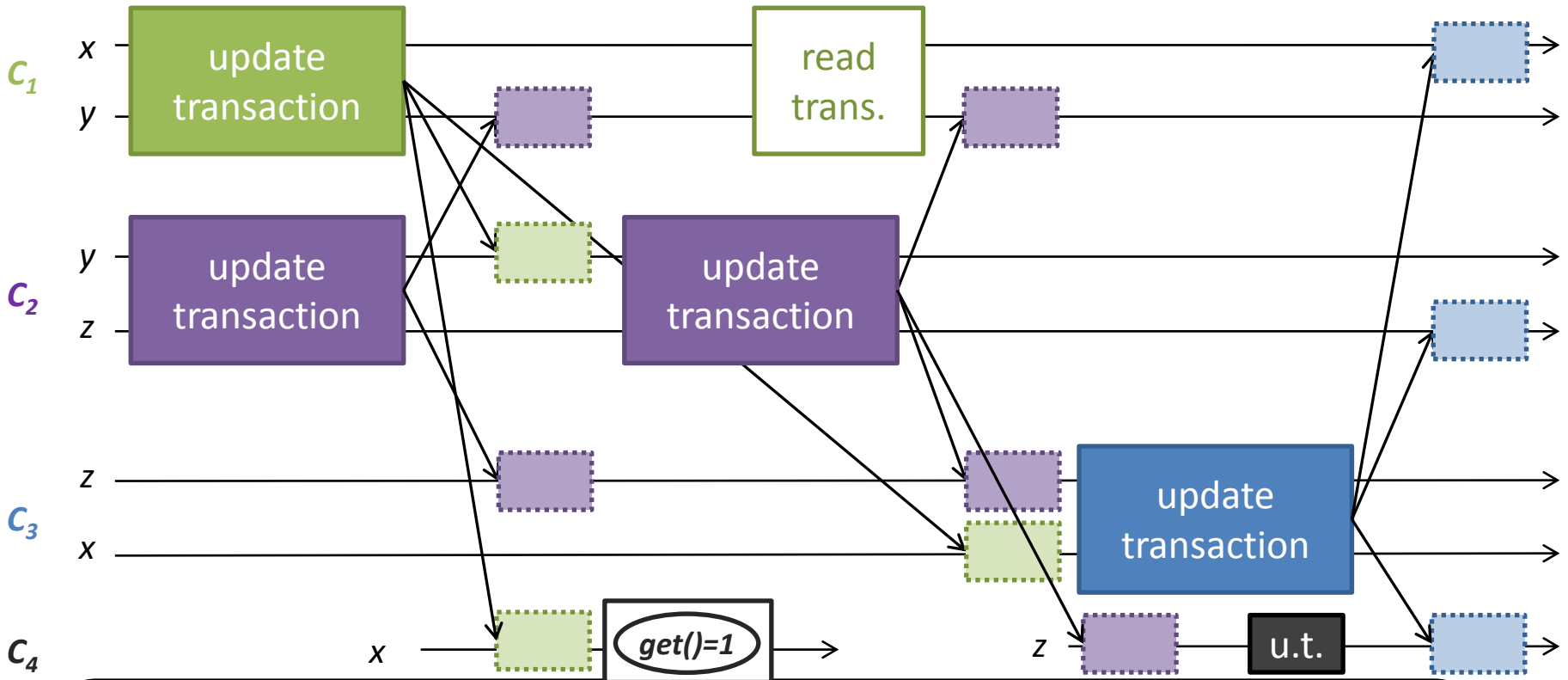
# Object model and transactional guarantees



Example: a database of replicated objects

**Simplified update transaction record**
- Read set = write set
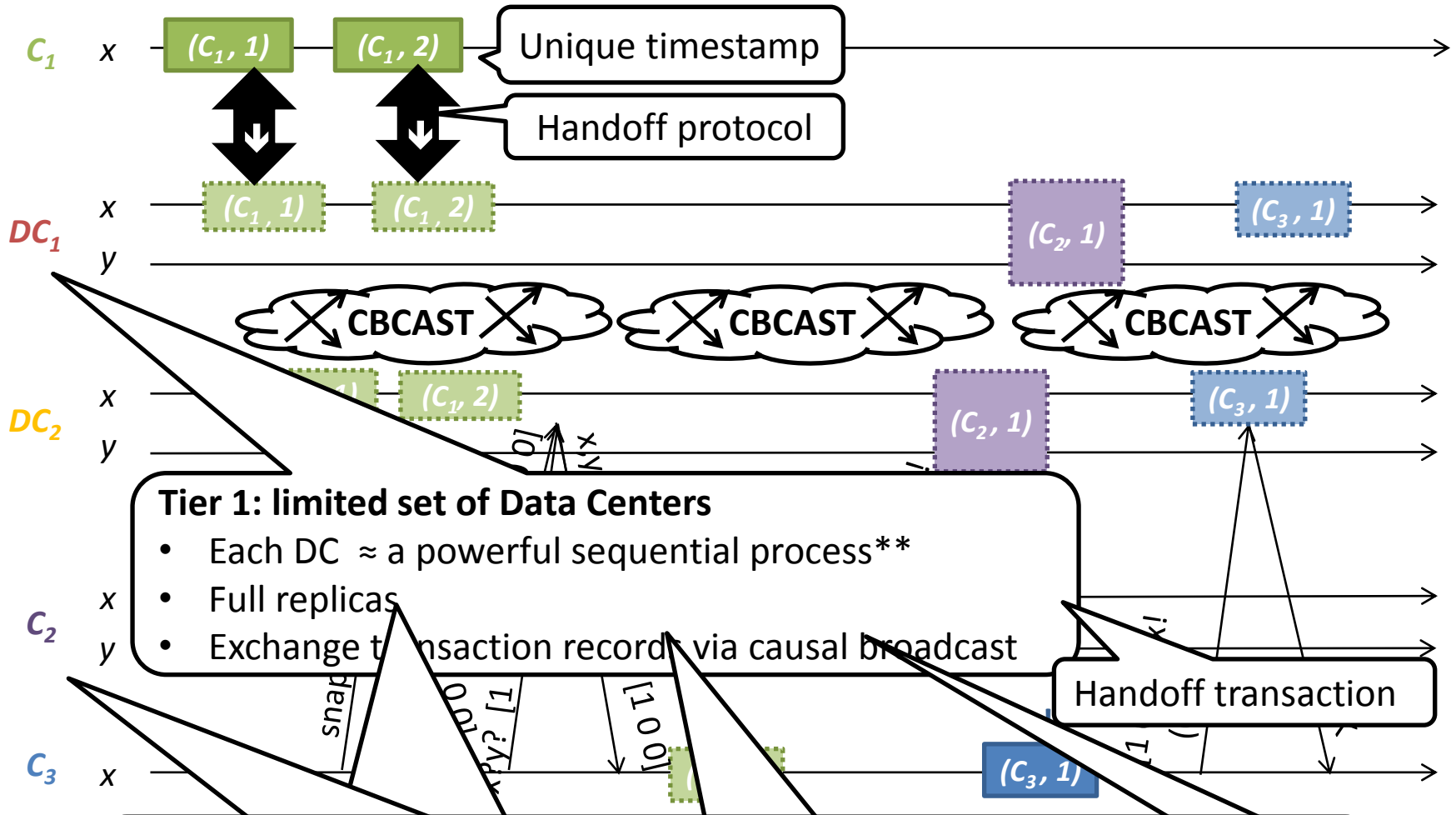
4

# What makes the problem hard? Workaround!



**Versioning and tracking causality is difficult under our assumptions:**
- Partial replication => direct dependency (causality) tracking hard
- Dynamic replica set => locating objects and updates dissemination hard
- High client churn => very big vector clocks

**Workaround:** use Data Centers as full reference replicas processing transactions!

# The 2-tier architecture with (naïve) handoff



$C_1$  x  $(C_1, 1)$  $(C_1, 2)$

Unique timestamp

Handoff protocol

$DC_1$  x  $(C_1, 1)$  $(C_1, 2)$  $(C_2, 1)$  $(C_3, 1)$
y

CBCAST  CBCAST  CBCAST

$DC_2$  x  $(C_1, 2)$  $(C_2, 1)$  $(C_3, 1)$
y

**Tier 1: limited set of Data Centers**
- Each DC ≈ a powerful sequential process**
- Full replicas
- Exchange transaction records via causal broadcast

$C_2$  x
y

Handoff transaction

$C_3$  x  $(C_3, 1)$

New tra...
- Defin...
- Effici...

✓ Dissemination and snapshot problems solved: DC-CBCAST is affordable!
(records delivery order is *stronger* than causality)
— Version Vectors (e.g. snapVV) unacceptably large: O(|Clients|)

...ction
...ates)

# Handoff with DC-assigned alias timestamp



During handoff DC assigns another *alias timestamp* to more efficiently refer to a set of transaction records

✓ DC-assigned aliases allow us to use managable vectors: O(|DCs|)
— But… is handoff asynchronous? Is it fault-tolerant and wait-free?

# Problem 1: making handoff asynchronous



$C_2$ awaits handoff ack+alias, does it prevent executing transaction concurrently?
✓ No! Include previous local transactions in snapshot by referring to their client timestamp rather than alias (session guarantees), e.g. snapVV($C_2$, 2) = [1 0 $C_2$=1]

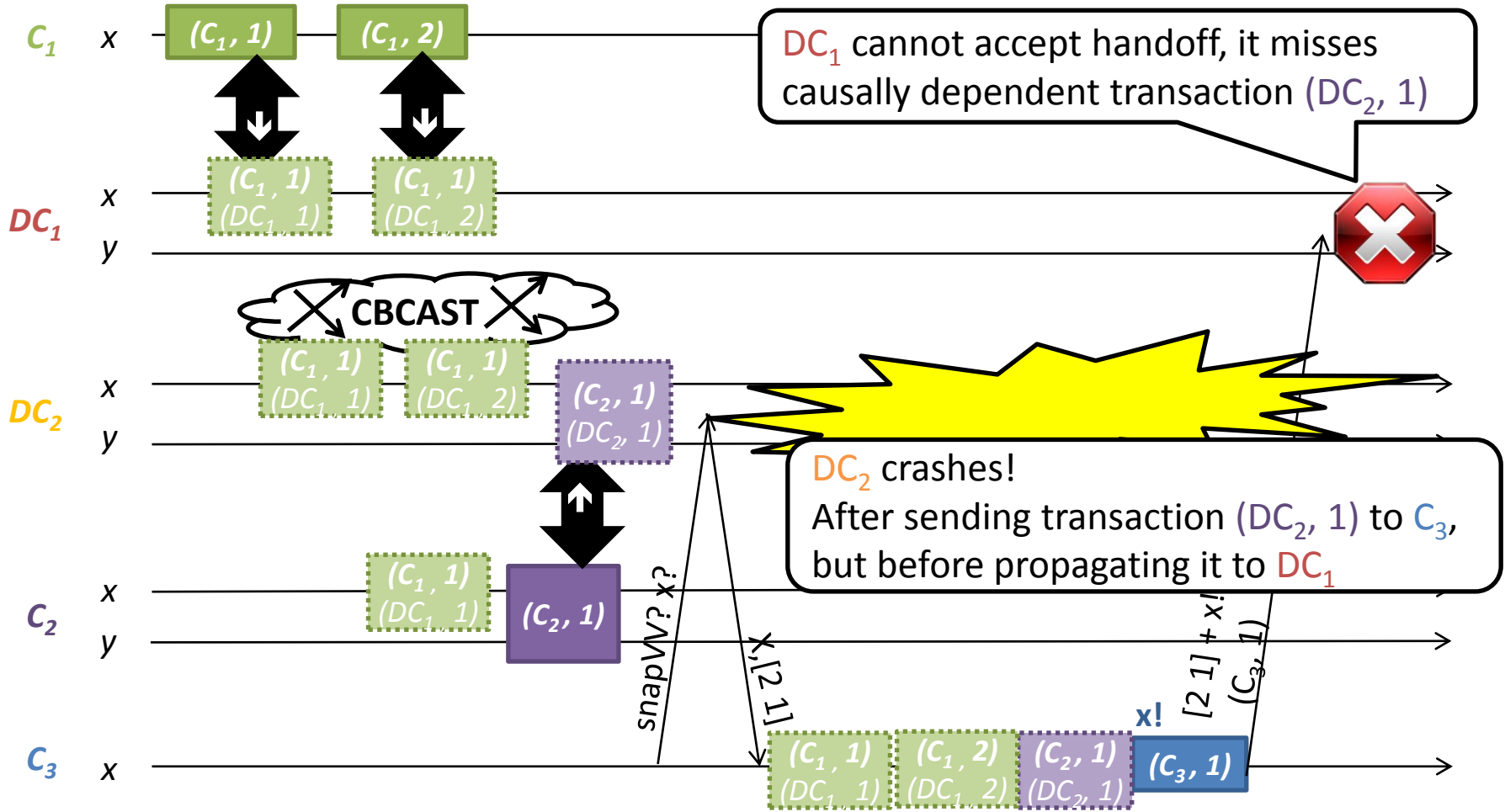# Problem 2: reading partially disseminated trans.

$C_1$   x   $(C_1, 1)$   $(C_1, 2)$

DC$_1$ cannot accept handoff, it misses causally dependent transaction $(DC_2, 1)$

$DC_1$   x   $(C_1, 1)$ $(DC_1, 1)$   $(C_1, 1)$ $(DC_1, 2)$   y

CBCAST

$DC_2$   x   $(C_1, 1)$ $(DC_1, 1)$   $(C_1, 1)$ $(DC_1, 2)$   $(C_2, 1)$ $(DC_2, 1)$   y

$DC_2$ crashes!
After sending transaction $(DC_2, 1)$ to $C_3$, but before propagating it to $DC_1$

$C_2$   x   $(C_1, 1)$ $(DC_1, 1)$   $(C_2, 1)$   y

snapVV? x?

x, [2 1]

[2 1] + x!

$(C_3, 1)$

x!

$C_3$   x   $(C_1, 1)$ $(DC_1, 1)$   $(C_1, 2)$ $(DC_1, 2)$   $(C_2, 1)$ $(DC_2, 1)$   $(C_3, 1)$

Liveness issue: updates of client $C_3$ are invisible to other clients until $DC_2$ recovers
- Side-effect of depending on *partially disseminated* transaction
- Client $C_3$ cannot recover missing transaction, it does not replicate y!

# Solution: DC offers only stable transactions



- ✓ Keep track of *stable transactions*, e.g. disseminated to majority of DCs
- ✓ Offer only stable transactions to the client (modulo his own transactions)
- — Delays visibility of recent transactions

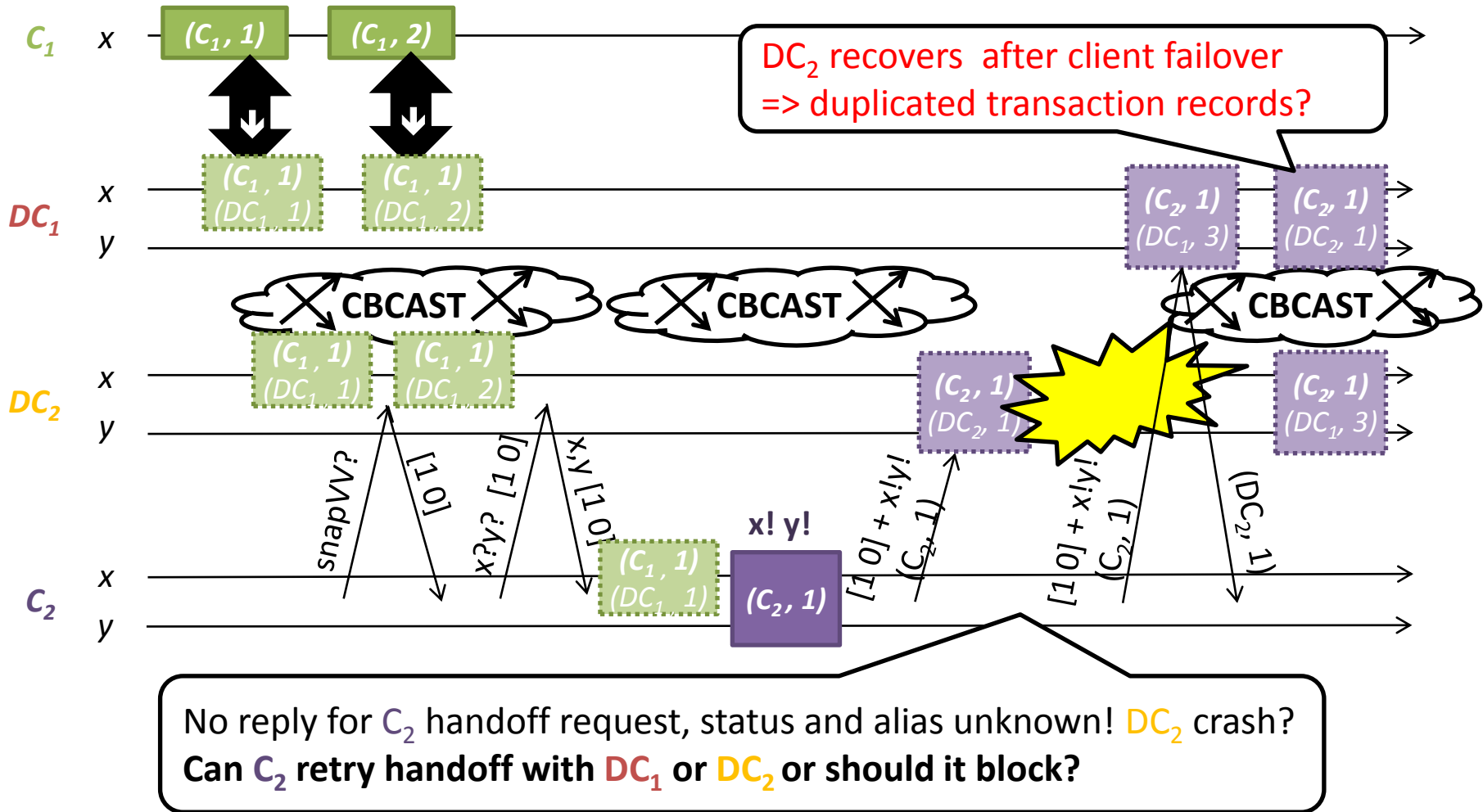# Problem 3: retrying handoff request



$C_1$  x  $(C_1, 1)$  $(C_1, 2)$

DC$_2$ recovers after client failover
=> duplicated transaction records?

$DC_1$  x  $(C_1, 1)$  $(C_1, 1)$  $(C_2, 1)$  $(C_2, 1)$
       y  $(DC_1, 1)$  $(DC_1, 2)$  $(DC_1, 3)$  $(DC_2, 1)$

CBCAST  CBCAST  CBCAST

$DC_2$  x  $(C_1, 1)$  $(C_1, 1)$  $(C_2, 1)$  $(C_2, 1)$
       y  $(DC_1, 1)$  $(DC_1, 2)$  $(DC_2, 1)$  $(DC_1, 3)$

snapVV?  [1 0]  x?y? [1 0]  x,y [1 0]  x! y!  [1 0] + x!y! $(C_2, 1)$  [1 0] + x!y! $(C_2, 1)$  $(DC_2, 1)$

$C_2$  x  $(C_1, 1)$  $(C_2, 1)$
      y  $(DC_1, 1)$

No reply for $C_2$ handoff request, status and alias unknown! DC$_2$ crash?
**Can $C_2$ retry handoff with DC$_1$ or DC$_2$ or should it block?**

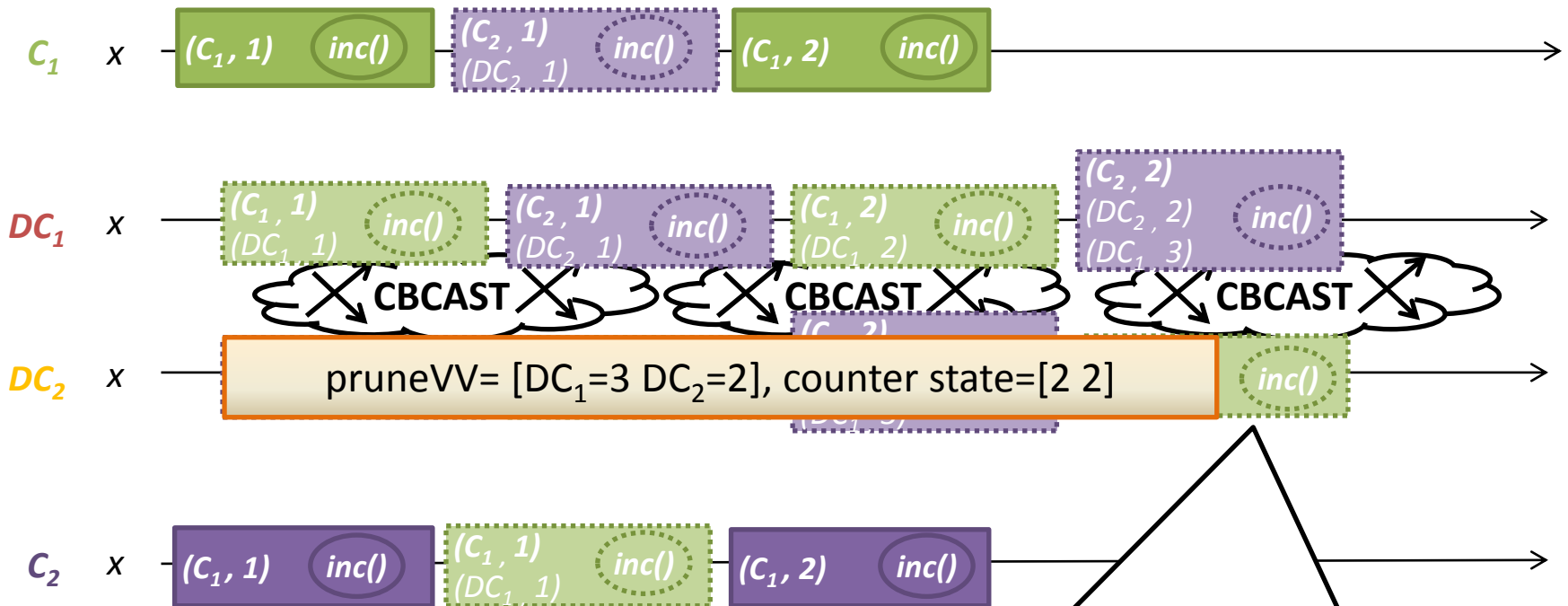# Solution: make handoff idempotent



✓ Merge transaction records with identical client timestamp (allow multiple aliases)

✓ This makes handoff idempotent, allows retries or failover

— Handoff idempotence more difficult with object pruning

12

# Problem 4: pruning object updates log safely



$C_1$  x  $(C_1, 1)$ inc()  $(C_2, 1)$ $(DC_2, 1)$ inc()  $(C_1, 2)$ inc()

$DC_1$  x  $(C_1, 1)$ $(DC_1, 1)$ inc()  $(C_2, 1)$ $(DC_2, 1)$ inc()  $(C_1, 2)$ $(DC_1, 2)$ inc()  $(C_2, 2)$ $(DC_2, 2)$ $(DC_1, 3)$ inc()

CBCAST    CBCAST    CBCAST

$DC_2$  x  pruneVV= [$DC_1$=3 $DC_2$=2], counter state=[2 2]    inc()

$C_2$  x  $(C_1, 1)$ inc()  $(C_1, 1)$ $(DC_1, 1)$ inc()  $(C_1, 2)$ inc()
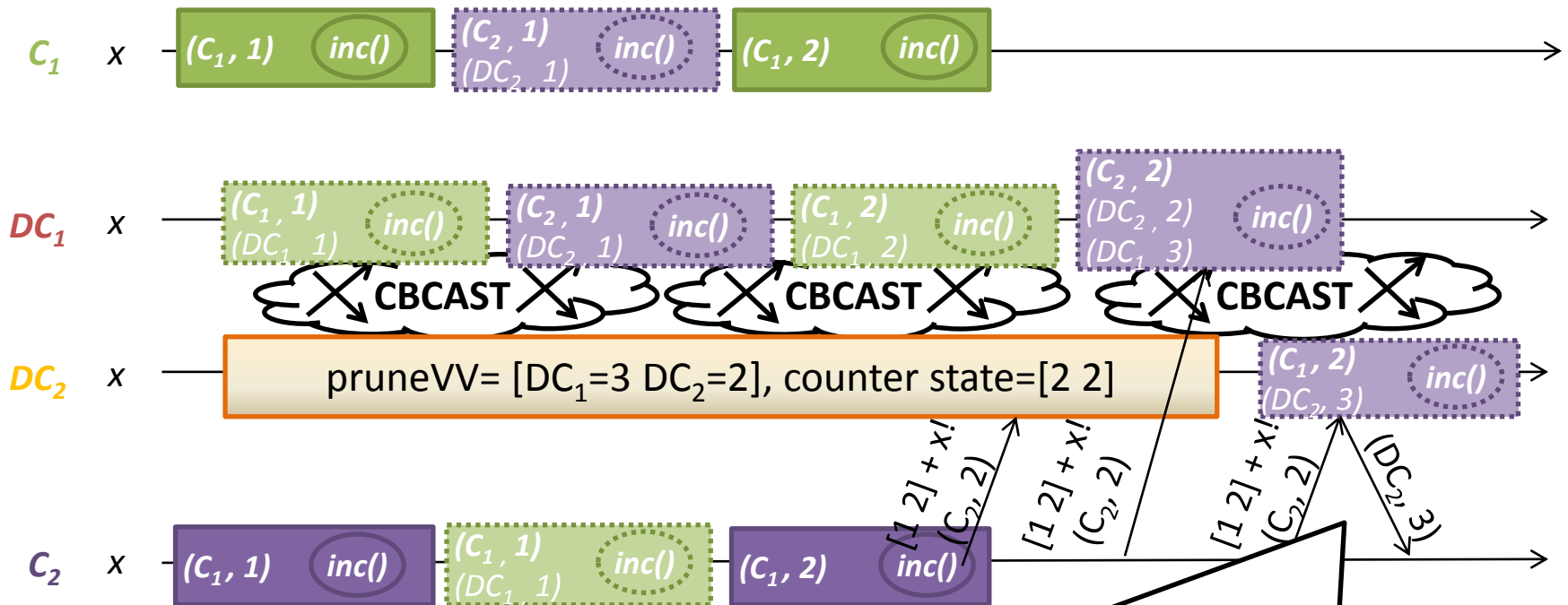
— Log-based object representation very costly:  O(|updates|)

**Replace old log of operation with state-based CRDT and pruneVV**
- For more efficient representation utilize timestamp aliases*
- Update representation must be uniform across DCs
  => pruning when set of aliases is stable (if many, select one)
- Current implementation: every DC prunes independently
  => client may need to merge states from different DCs

13

# Pruning and handoff idempotence



$C_1$  x

$(C_1, 1)$  inc()   $(C_2, 1)$ $(DC_2, 1)$  inc()   $(C_1, 2)$  inc()

$DC_1$  x

$(C_1, 1)$ $(DC_1, 1)$  inc()   $(C_2, 1)$ $(DC_2, 1)$  inc()   $(C_1, 2)$ $(DC_1, 2)$  inc()   $(C_2, 2)$ $(DC_2, 2)$ $(DC_1, 3)$  inc()

CBCAST   CBCAST   CBCAST

$DC_2$  x

pruneVV= [$DC_1$=3 $DC_2$=2], counter state=[2 2]

$(C_1, 2)$ $(DC_2, 3)$  inc()

$C_2$  x

$(C_1, 1)$  inc()   $(C_1, 1)$ $(DC_1, 1)$  inc()   $(C_1, 2)$  inc()

[1 2] + x! $(C_2, 2)$   [1 2] + x! $(C_2, 2)$   [1 2] + x! $(C_2, 2)$   $(DC_2, 3)$

— Retrying handoff of pruned transaction => violated idempotence!
✓ Maintain single clientVV per DC – "idempotence guard"
✓ Never transmitted, only single entry

14

# Lessons learnt

- Implementing CRDT transactions ≈ implementing a huge "database" semi-lattice
  - Difference w.r.t. ordinary object: (dynamic) fragmentation*
  - Use different techniques inside and across objects

- Causality tracking is difficult at scale, both inside/across CRDTs
  - Limit communication topology, here: 2-tier architecture
  - Use handoff protocol with timestamp aliasing and

- Making handoff live and correct despite Tier 1 failures
  - Reading stable versions helps failover
  - Timestamp aliasing helps too
  - When forced to store a big VV, share it across DB