



# **Don't Settle for Eventual: Scalable Causal Consistency for Wide- Area Storage with COPS**

Authors: Lloyd et al.  
Presenter: Abhishek A. Singh



# Ideas

- ALPS – **A**vailability, Low **L**atency, **P**artition Tolerance and **S**calability
- Causal+,
- COPS GT
- Design of COPS/COPS GT key-value store



# **What is Causal?**



# Causal + X = Causal+

- *X = convergent conflict handling*
- Implemented via **COPS**
  - **C**lusters of **O**rders-**P**reserving **S**ervers



# System properties desired

- Availability
  - Ops. do not fail; do not block indefinitely; no error
- Low latency
  - Ops. complete “quickly” (in milliseconds)
- Partition tolerance
  - Data Store continues to operate under partitions
- High scalability
  - Linear increase in system performance on addition to resources
- Stronger consistency
  - Causal+ consistency. Weaker than linearizability.

# Causal+ consistency

- Potentially causal  $\rightarrow$  (in the paper a squiggly arrow)
  - If  $a$  and  $b$  are two operations in a single thread of execution, then  $a \rightarrow b$
  - If  $a := \text{puts}()$  operation, and  $b := \text{gets}(a)$ , then  $a \rightarrow b$ 
    - $\text{gets}(a)$  returns value put by  $a$
  - If  $a \rightarrow b$  and  $b \rightarrow c$ , then  $a \rightarrow c$

# Define causal+ consistency

- Causal consistency + convergent conflict handling
- Convergent conflict handling
  - all conflicting puts should be handled in the same way across all replicas using some handler function  $h$
  - Problem to be solved: avoid state divergence across the network.
    - Eg. *last-writer-wins*



# Consistency foodchain

**Linearizability** > **Sequential** > Causal+ > Causal > FIFO  
> Per-Key Sequential > Eventual



# Causal+ in COPS

- Remember COPS?
- **C**lusters of **O**rders-**P**reserving **S**ervers
- Defines 2 abstractions:
  - *Versions*
    - If  $x_i \rightarrow y_j$ , then  $i < j$
    - Only return causally later version of a key (progressing property)
  - *Dependencies*
    - If  $x_i \rightarrow y_j$ , then  $x_i$  must be written before  $y_j$
    - $y_j$  **depends on**  $x_i \leftrightarrow \text{put}(x_i) \rightarrow \text{put}(y_j)$



# Scalable causality?

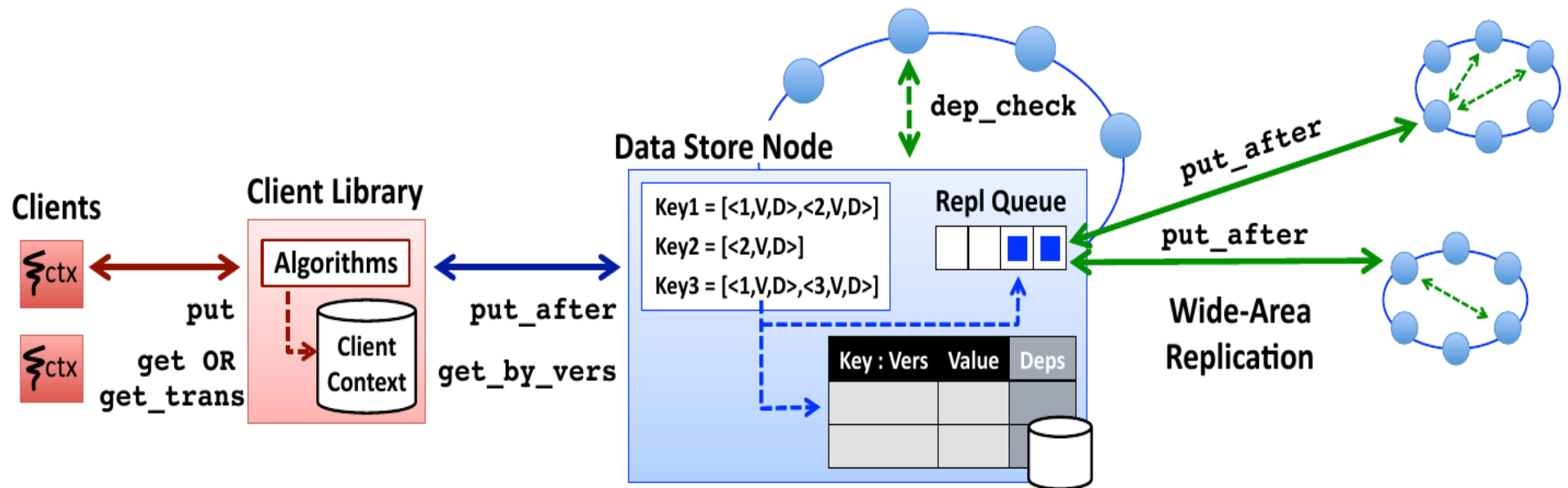
- Dependencies encoded in each key's metadata.
- During key replication each receiving data center performs dependency check before committing the incoming version



# The design of COPS

- Local datacenter clusters implement linearizability
- Remote datacenters: causal replication

# COPS Architecture



# Client library interface

- `ctx id ← createContext()`
- `bool ← deleteContext(ctx id)`
- `bool ← put (key, value, ctx id)`
- `value ← get (key, ctx id) [In COPS]`

OR

- `[values] ← get trans (hkeysi, ctx id) [In COPS-GT]`

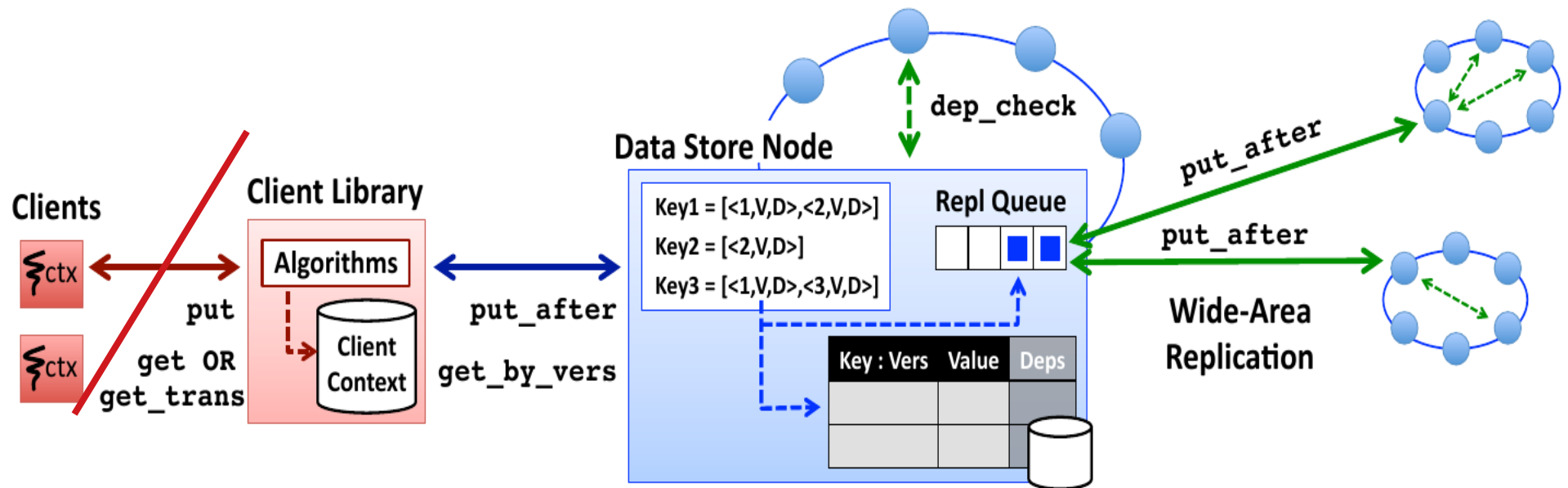
# Read, write and replicate in COPS/COPS-GT

- $[\text{bool}, \text{vers}] \leftarrow \text{put after}(\text{key}, \text{val}, [\text{deps}], \text{nearest}, \text{vers}=\emptyset)$
- 
- $\text{bool} \leftarrow \text{dep check}(\text{key}, \text{version})$
- $[\text{value}, \text{version}, \text{dep}] \leftarrow \text{get by version}(\text{key}, \text{version}=\text{LATEST})$

# GT part of COPS-GT

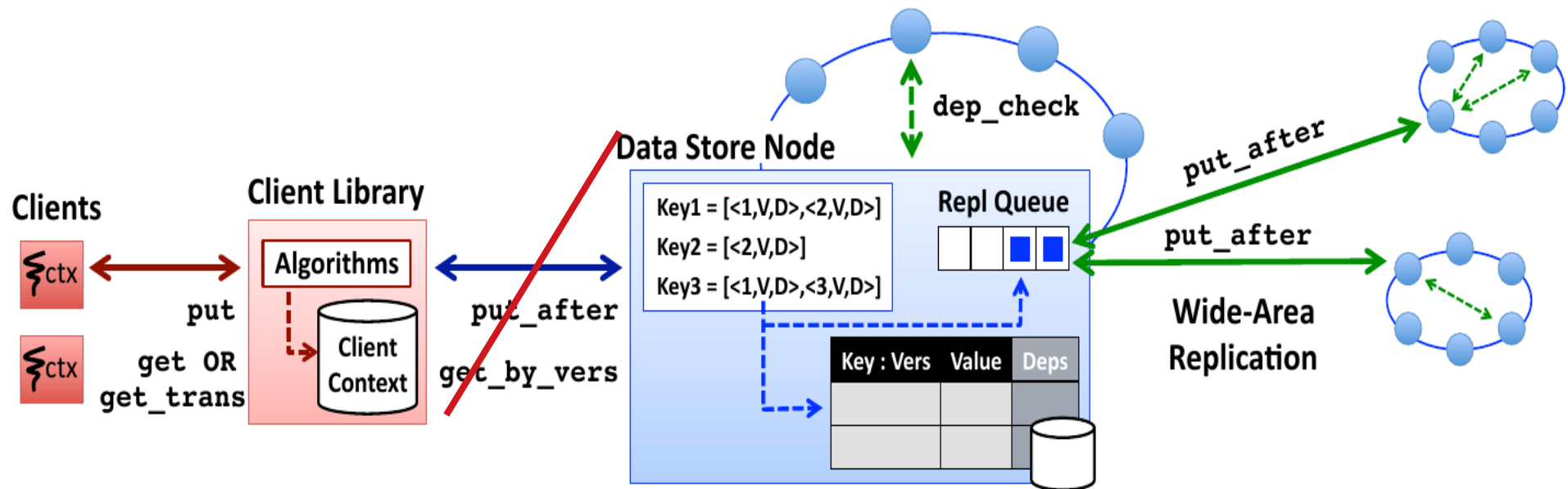
- Get Transactions API
  - $[\text{values}] \leftarrow \text{get\_trans}(\text{keys}, \text{ctx\_id})$
- Why does this have to be exposed to the client and not handled internally?
- Is COPS-GT really “stronger” than COPS?  
How?

# Faults: Scenario 1

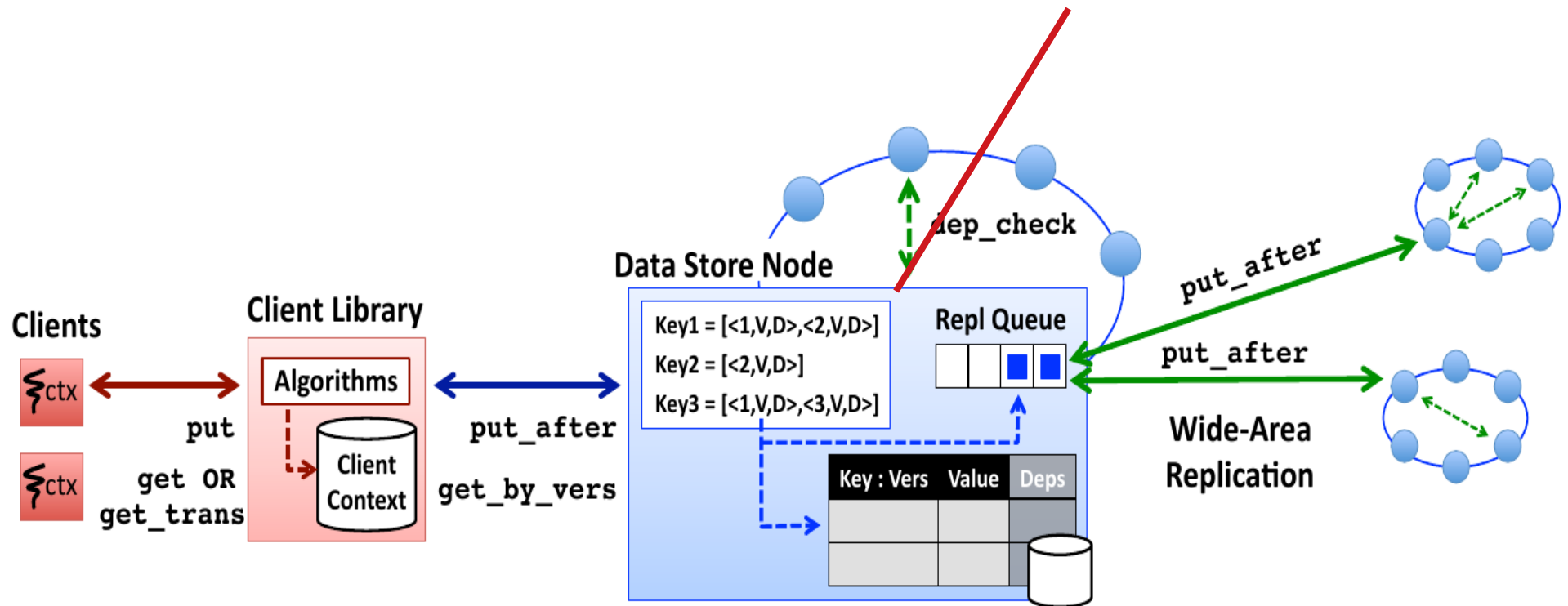




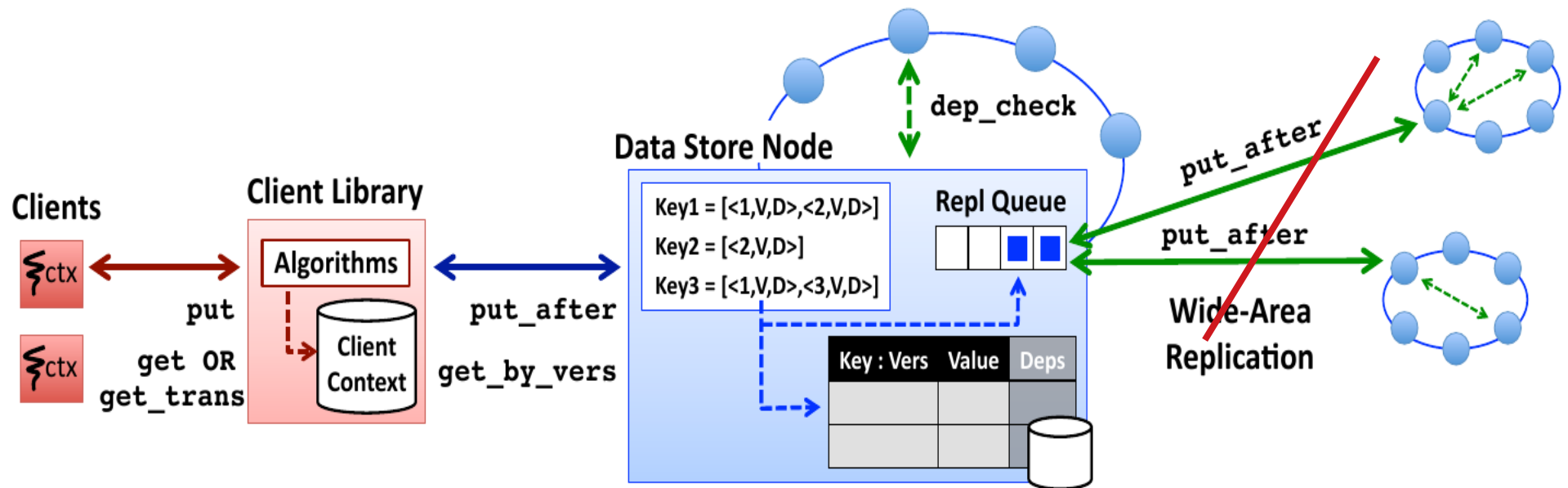
# Faults: Scenario 2



# Faults: Scenario 3



# Faults: Scenario 4



# Faults: All you can fault

