

# Large Scale Distributed Systems

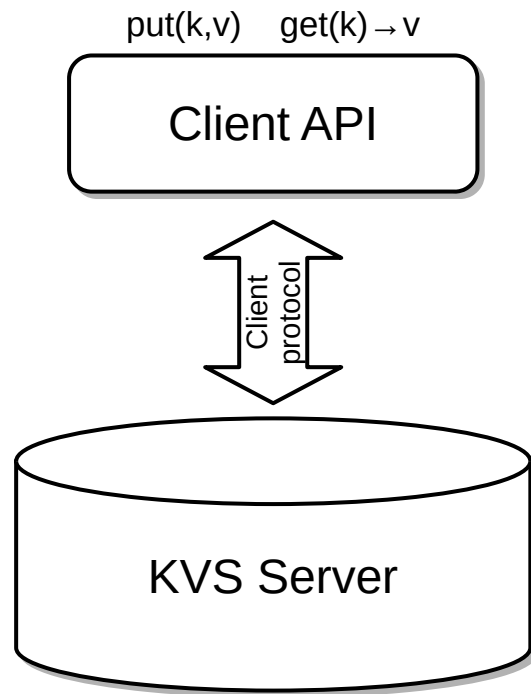
José Orlando Pereira

Departamento de Informática  
Universidade do Minho




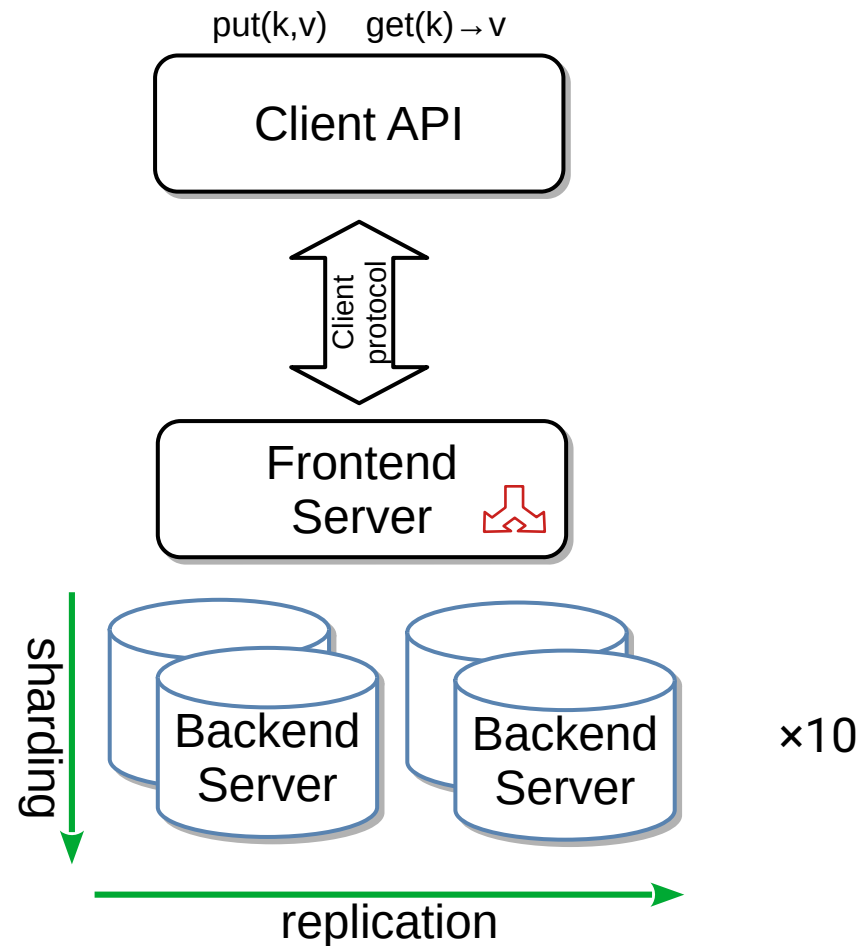
# Key-Value Store

- Simple data model:
  - $\text{Map}\langle K, V \rangle$
- Simple interface:
  - $\text{get}(k) \rightarrow v$
  - $\text{put}(k, v)$
- Avoids session state in server:
  - No multi-item transactions
  - No long lived operations



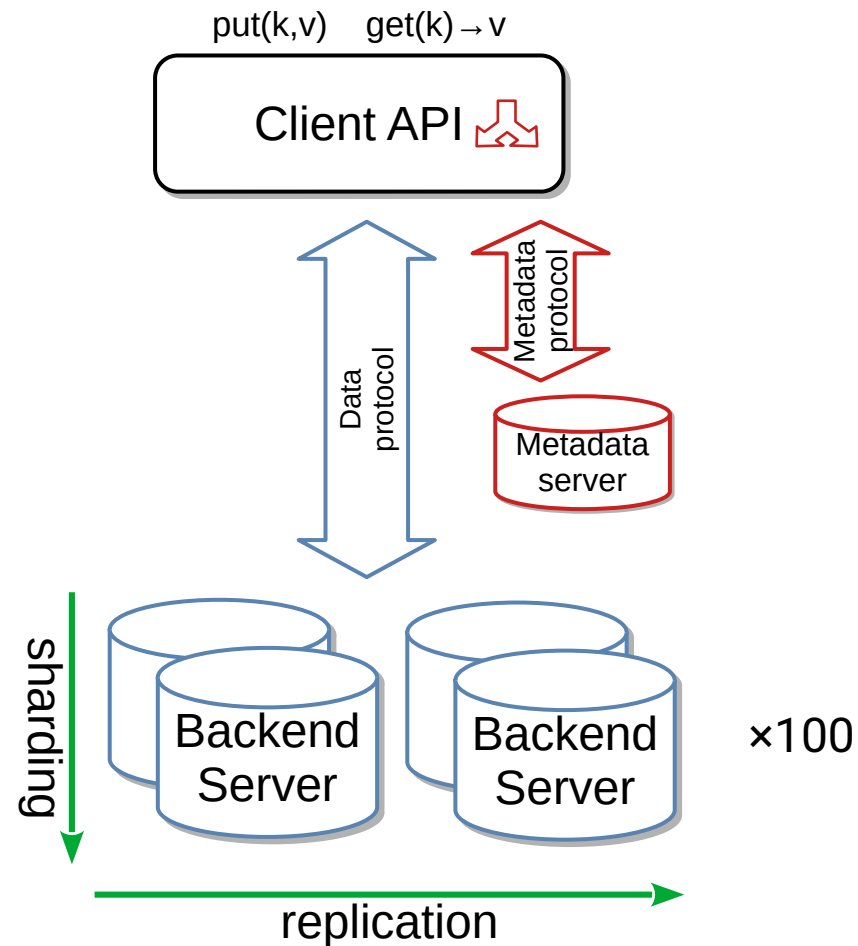
# Centralized

- Replication for availability
- Sharding for scale-out
- Centralized architecture:
  - Frontend server with metadata store
  - Backend servers with data stores
- Frontend is a bottleneck (latency and bandwidth)
- Example:  **mongoDB**





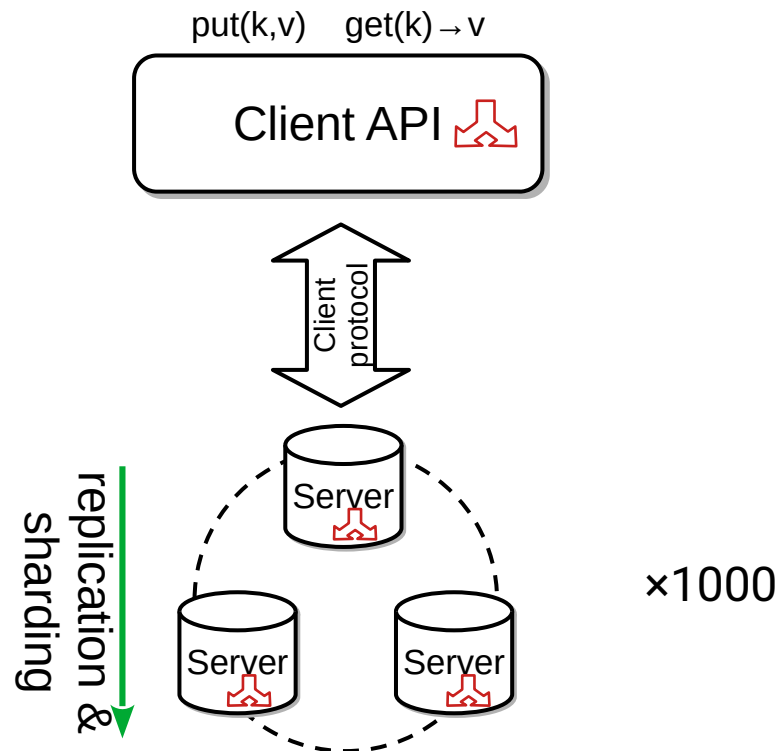
# Decentralized data

- Separate data and metadata servers and protocols
  - Client-side caching of metadata
- Avoids data bottleneck
- Examples:



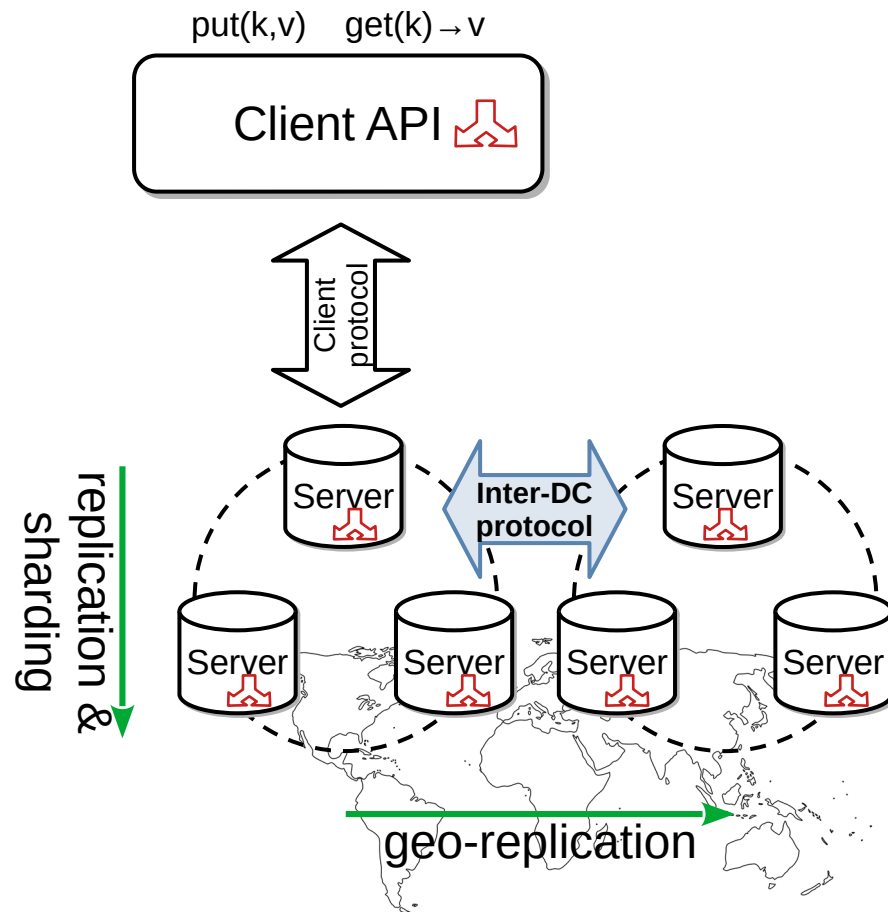
# Fully decentralized

- Consistent hashing
- Epidemic dissemination
- Examples:
  - Dynamo  
(not DynamoDB!) 
  -  Apache CASSANDRA

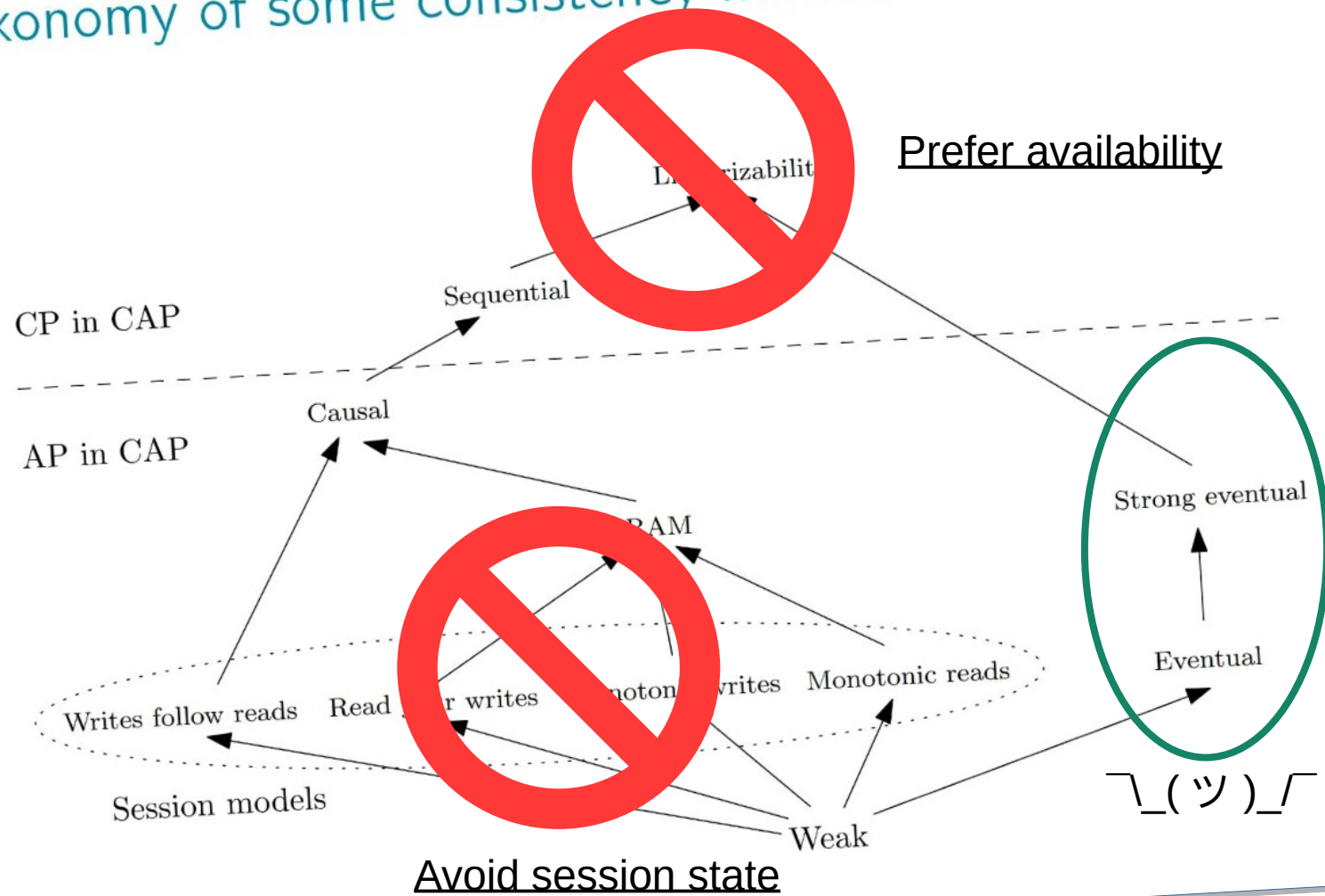


# Geo-replication

- Multi-data center replication
- Main challenge: How to shield clients from Inter-DC protocol?
  - Latency
  - Availability

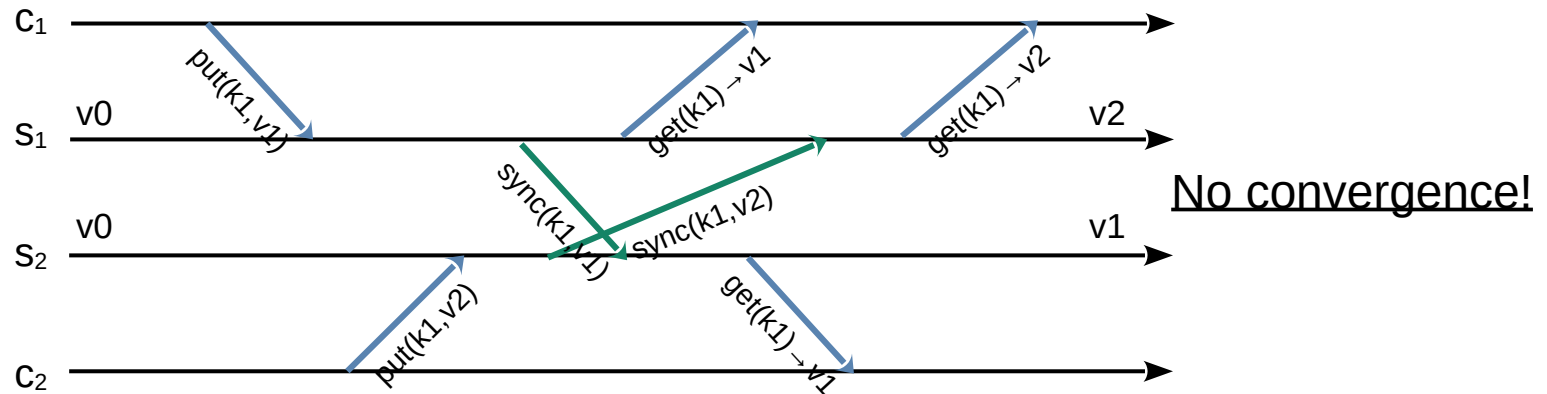


# A taxonomy of some consistency models



# Local writes

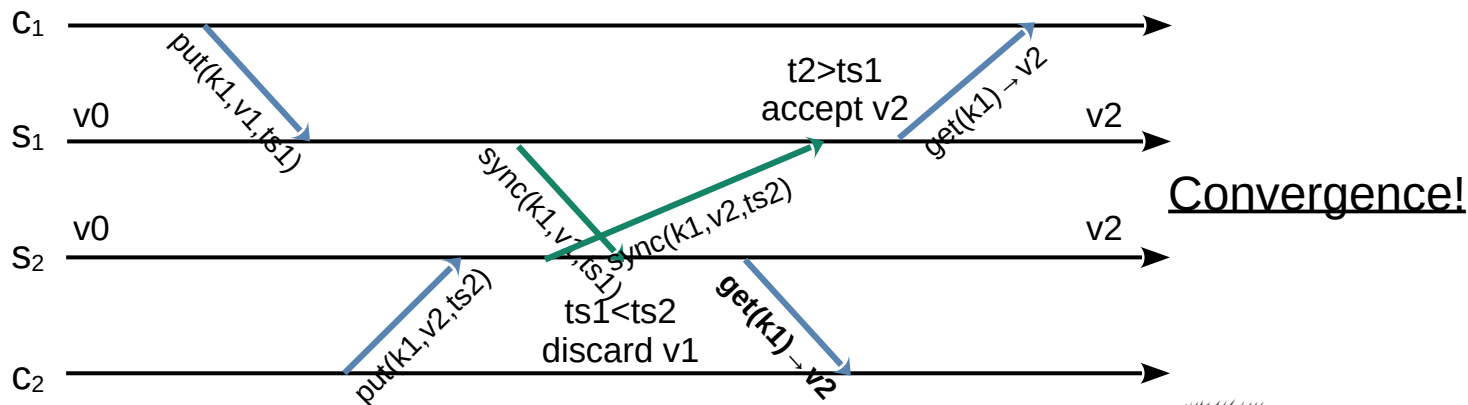
- Writes are issued to one or a few servers in the local data center
- Writes are asynchronously propagated to other data centers





# Convergence

- Simple reconciliation rule: Last Writer Wins (LWW)
  - Attach a timestamp to each data item
  - On conflict, keep the item with the latest timestamp



# Convergence

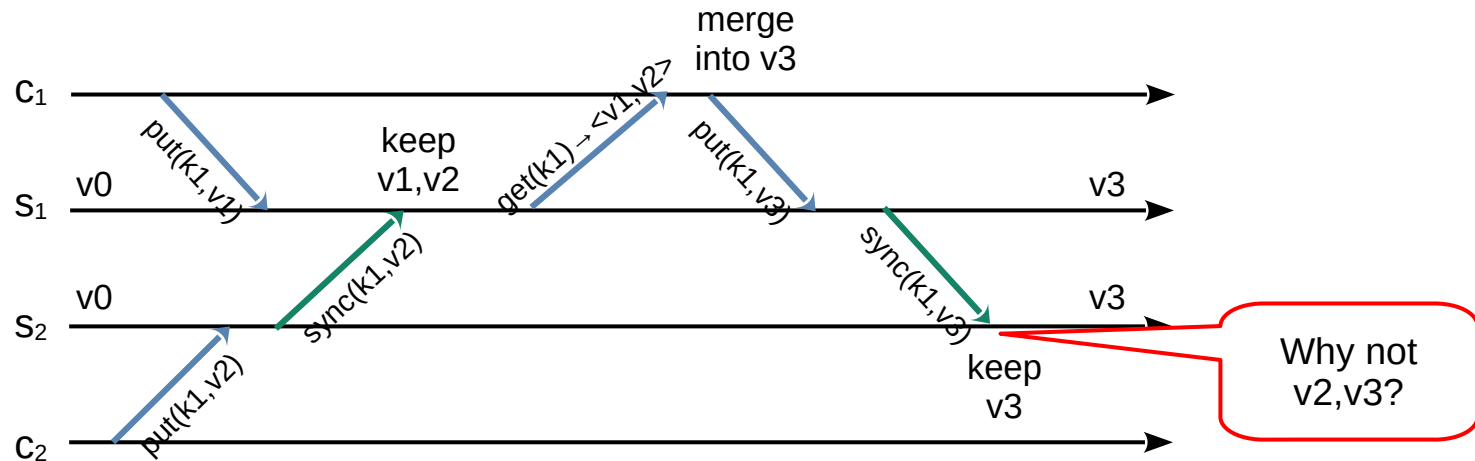
- Generating the timestamp:
  - Physical: when the client is a Web application and the state partitioned by user
  - Logical?
- LWW Register is actually a State-based CRDT... and KVS servers could support a variety of CRDTs
  - Limits the range of possible data structures and operations
  - Pushes complexity / policy / computation into the server
  - Not worth it if divergence is rare (client affinity + network stability)

# Convergence

- Idea: Delegate reconciliation to the client application
  - Simple LWW (Cassandra)
  - Other State-based CRDTs (Cure)
  - Custom code (Bayou)
- Servers cannot callback into the client application
- How to achieve this?

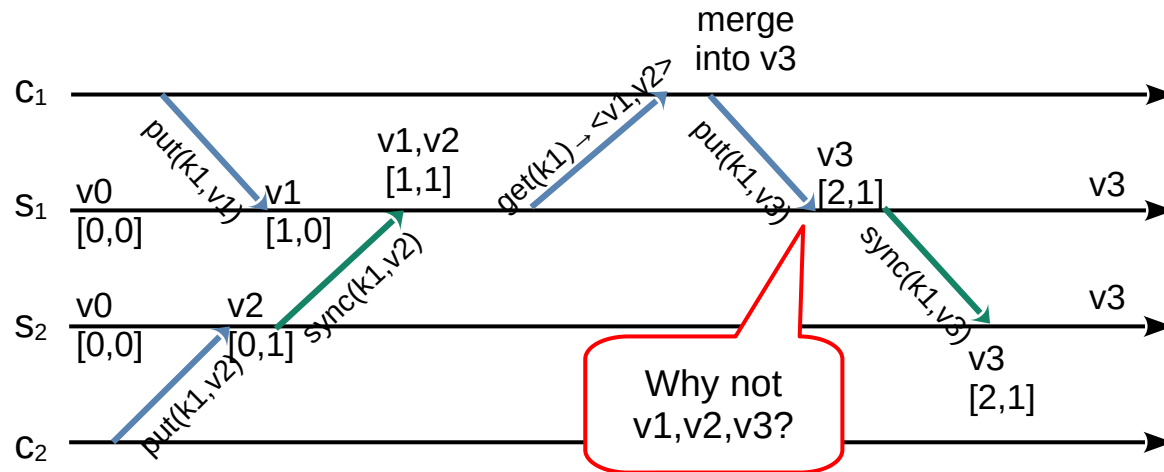
# Convergence

- Keep list of conflicting values in each key and return them to the application:  $\text{get}(k) \rightarrow \langle v_1, \dots, v_n \rangle$
- Let the application merge them and replace the list with a single value



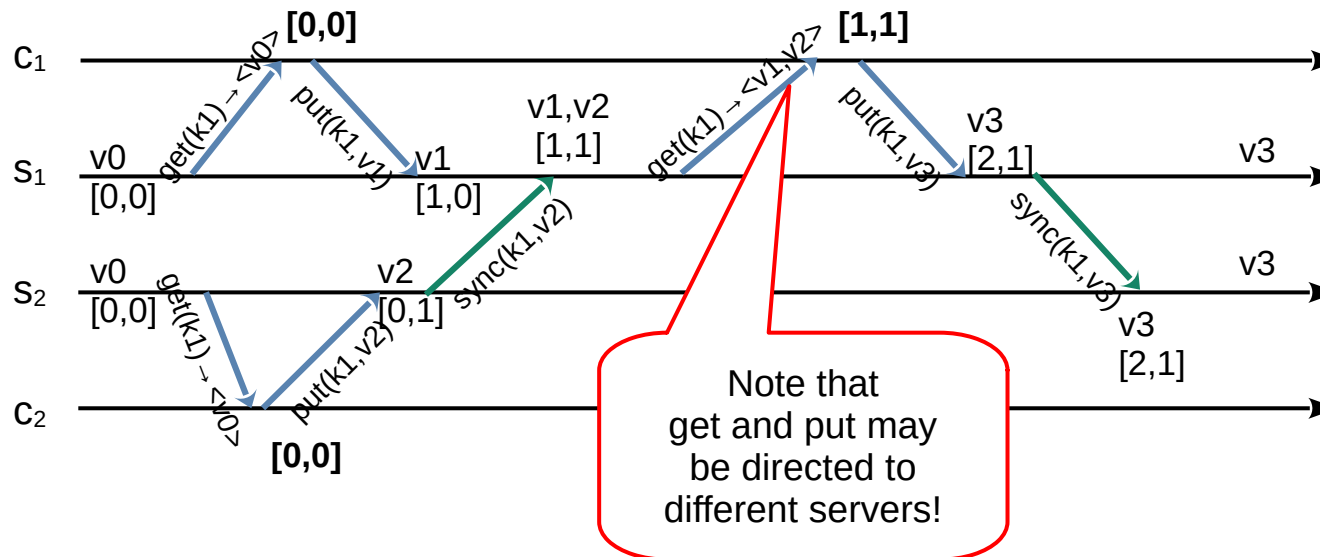
# Ordering propagation requests

- Keep a vector timestamp with each data item
  - An element for each server
- When merging concurrent timestamps, keep distinct values

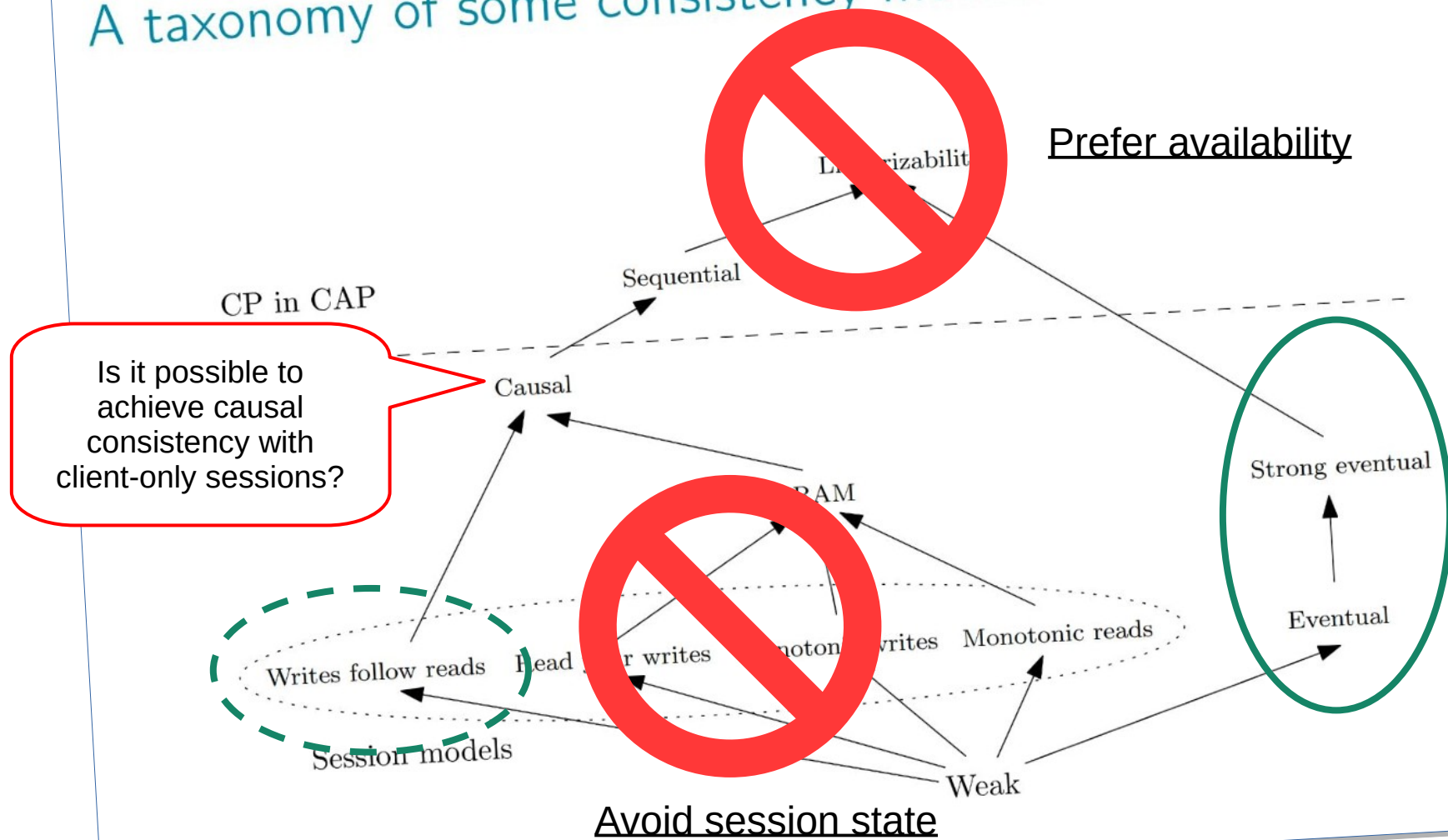


# Ordering write requests

- Client requests are anonymous: How to tell what the client has read before?
  - Force the client to read before writing
  - Make it keep track of context

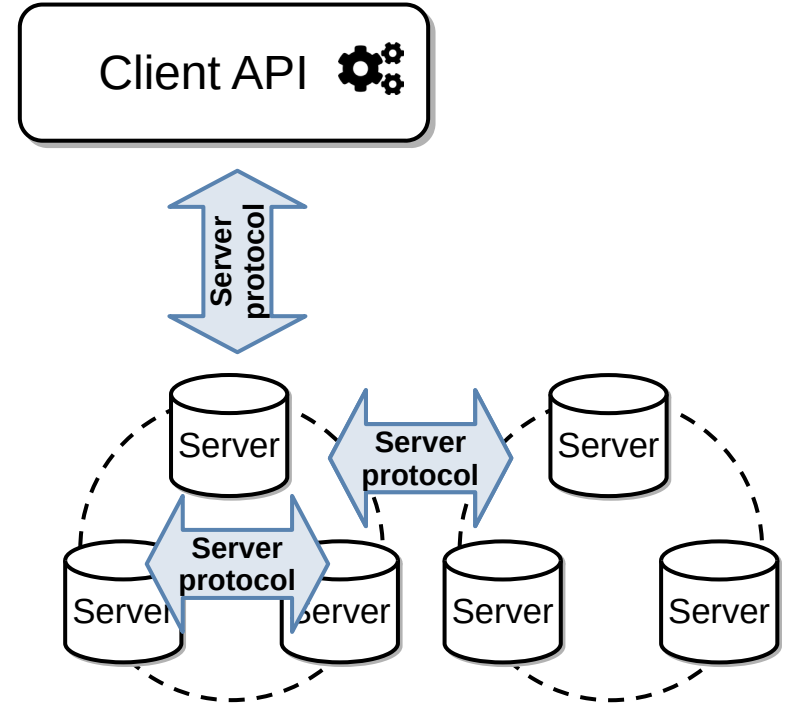


# A taxonomy of some consistency models



# Causal KVS (COPS)

- Client API:
  - `createCtx()` → `ctx`
  - `get(k,ctx)` → `v`
  - `put(k,v,ctx)`
  - `deleteCtx(ctx)`
- Servers keep a scalar timestamp for each (k,v) item (including server id in lower bits)
- Server protocol:
  - `get(k)` → (v,vers)
  - `put_after(k,v,deps,vers=0)` → `vers`
  - `dep_check(k,vers)`



get operation is trivial



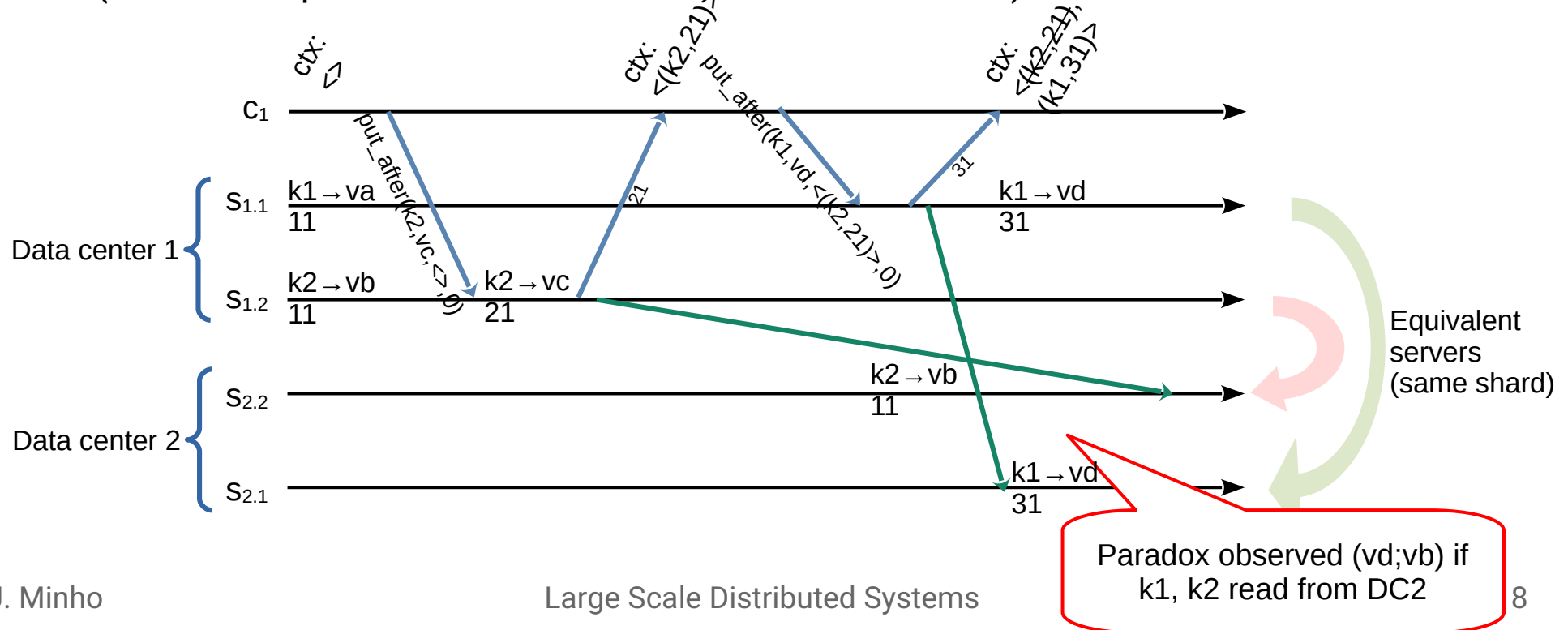
# Managing client context

- Client context needs to keep nearest dependencies
- Simple (approximate) strategy:
  - After a  $\text{put}(k, \dots) \rightarrow \text{vers}$  operation:
    - Clear all dependencies
    - Insert  $(k, \text{vers})$
  - After a  $\text{get}(k, \dots) \rightarrow (v, \text{vers})$  operation:
    - If  $(k, \dots)$  in context, update to  $(k, \text{vers})$
    - If not, insert  $(k, \text{vers})$

(It is approximate because it does not recognize dependencies between separately read keys and keeps them both)

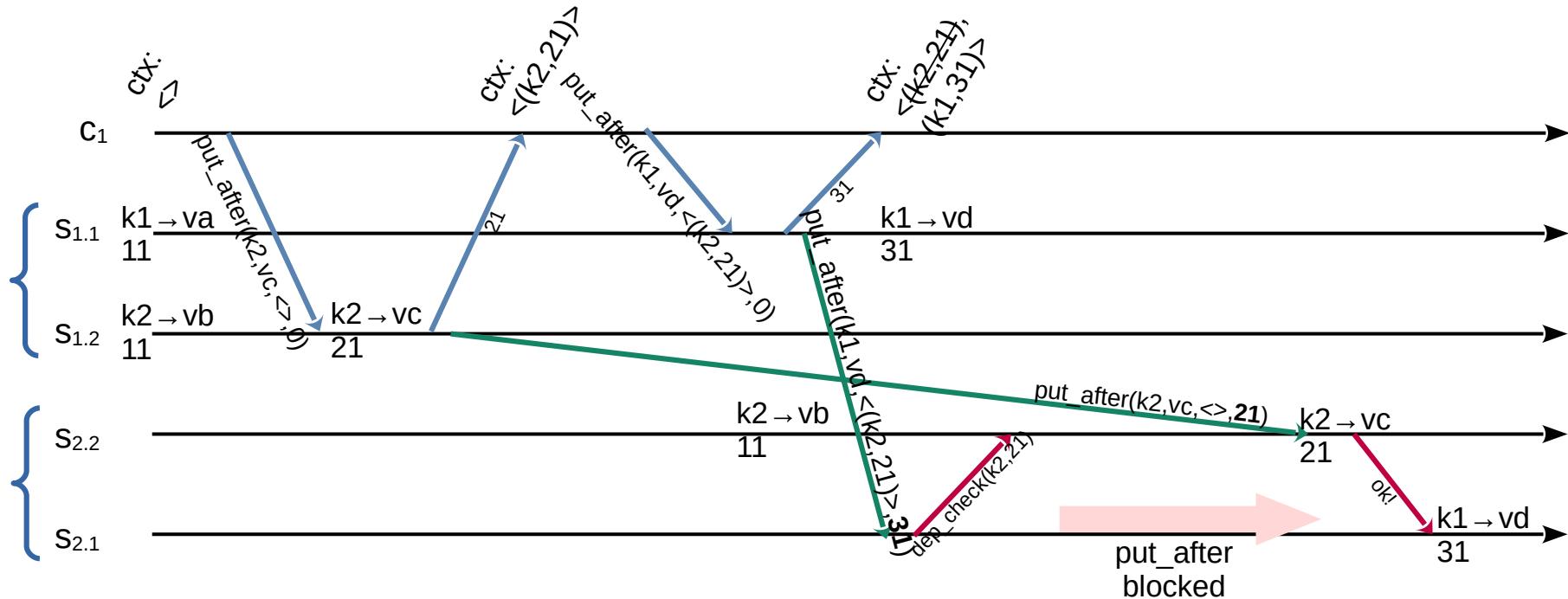
# Write operation

- Synchronous write to server in local data center:
  - append dependencies in local context to request
  - wait for write to commit before returning local clock
 (it is now impossible to read an older k in the same DC)




# Write propagation

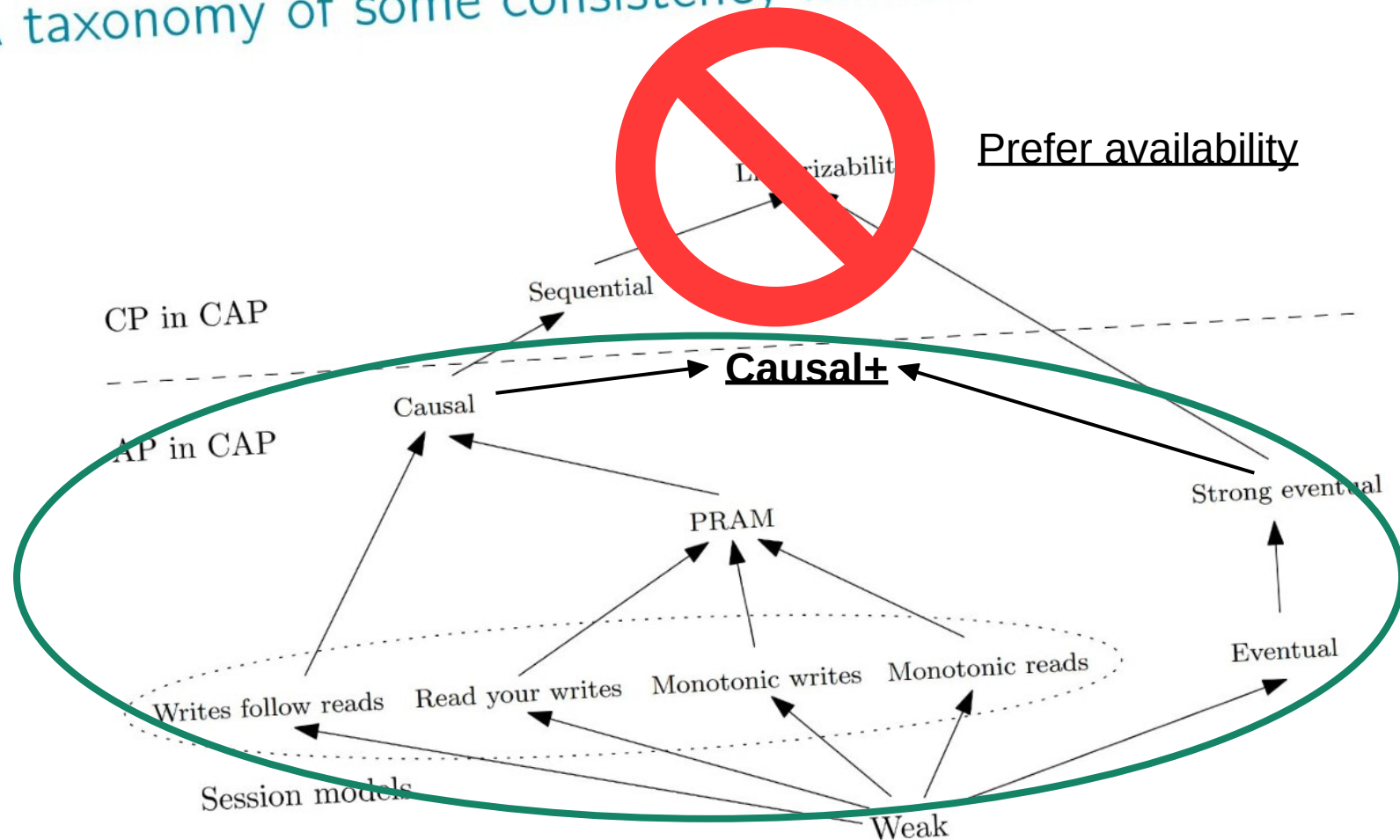
- Propagation is not ordered, as it is potentially from different data centers
- A value is committed remotely only after waiting for dependencies with the dep\_check operation of corresponding servers



# Consistency and convergence

- WFR: reads kept in local context and used in write
- RYW: 
- MR: } synchronous write and client affinity
- MW: }
- Strong Eventual:
  - Logical timestamps including global server mean that there are no ties between concurrent updates to same item
  - Last Writer Wins when propagating to remote DC

## A taxonomy of some consistency models



## No (server-side) session state

# References

- G. DeCandia et al., “**Dynamo: Amazon’s highly available key-value store**,” Oper. Syst. Rev., vol. 41, no. 6, pp. 205–220, Oct. 2007  
<https://doi.org/10.1145/1323293.1294281>
- W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen, “**Don’t settle for eventual: Scalable causal consistency for wide-area storage with COPS**,” in Proc. 23rd ACM Symposium on Operating Systems Principles, 2011  
<https://www.cs.cmu.edu/~dga/papers/cops-sosp2011.pdf>