# Large Scale Distributed Systems
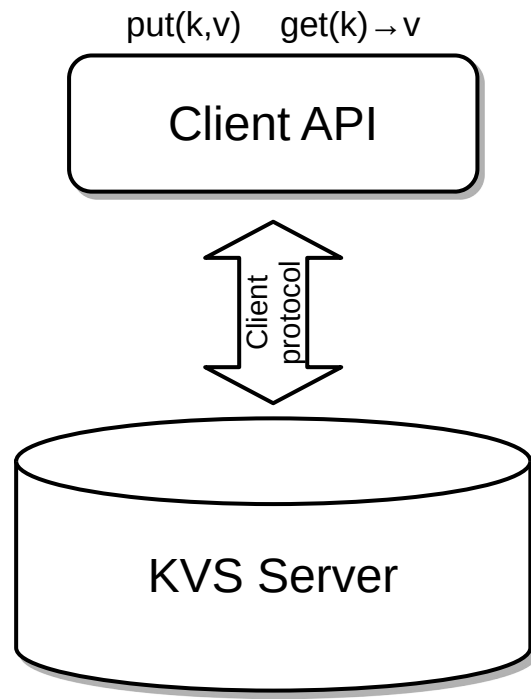
José Orlando Pereira

Departamento de Informática
Universidade do Minho
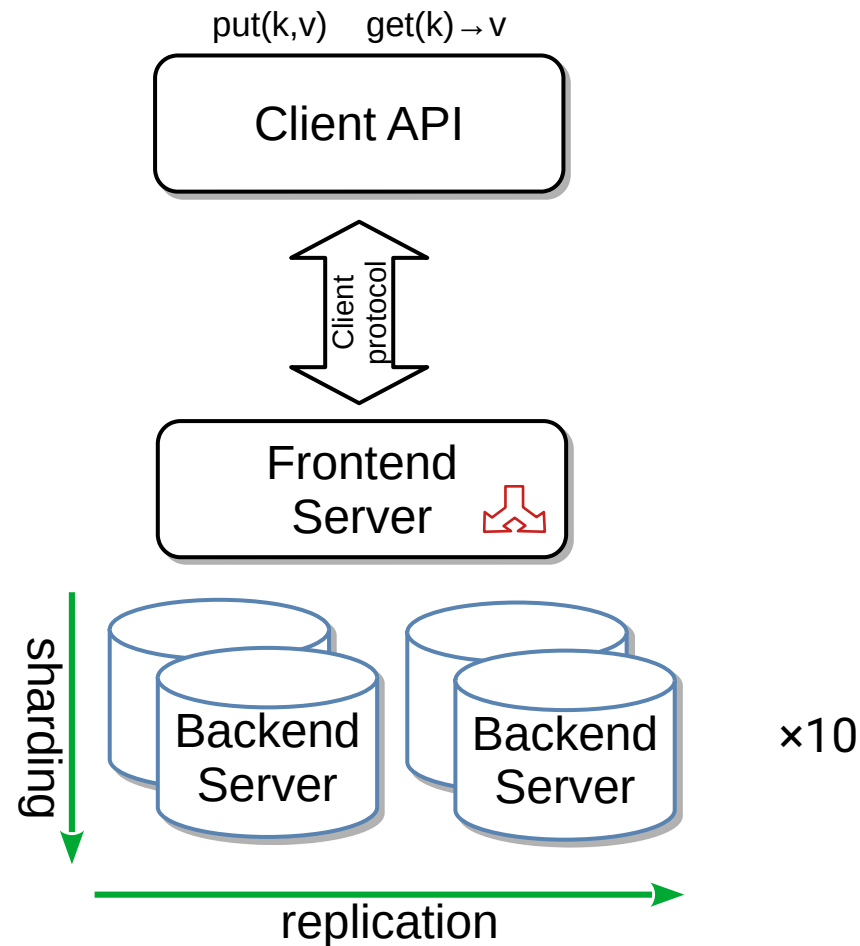
# Key-Value Store

- Simple data model:
  - Map<K,V>
- Simple interface:
  - get(k) → v
  - put(k,v)
- <u>Avoids session state</u> in server:
  - No multi-item transactions
  - No long lived operations

put(k,v)     get(k)→v

Client API

Client protocol
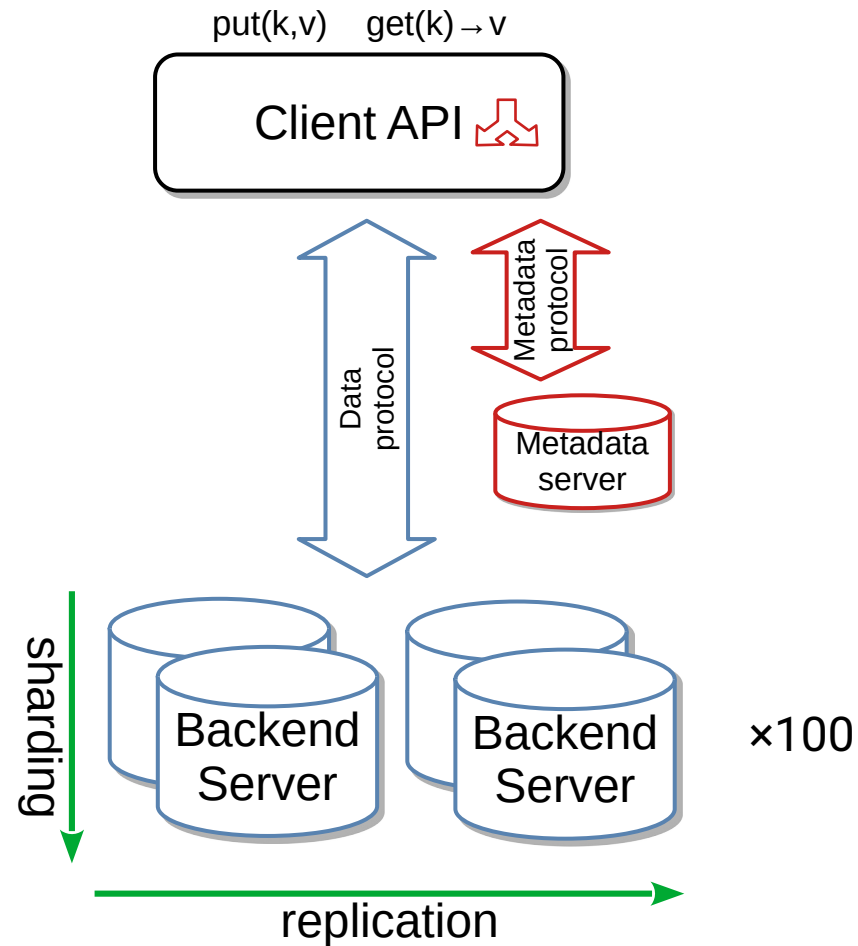
KVS Server

# Centralized

- <u>Replication</u> for availability

- <u>Sharding</u> for scale-out

- Centralized architecture:

  - Frontend server with metadata store

  - Backend servers with data stores

- Frontend is a bottleneck (latency and bandwidth)

- Example: mongoDB



put(k,v)    get(k)→v

Client API

Client protocol

Frontend Server

sharding

Backend Server    Backend Server    ×10

replication

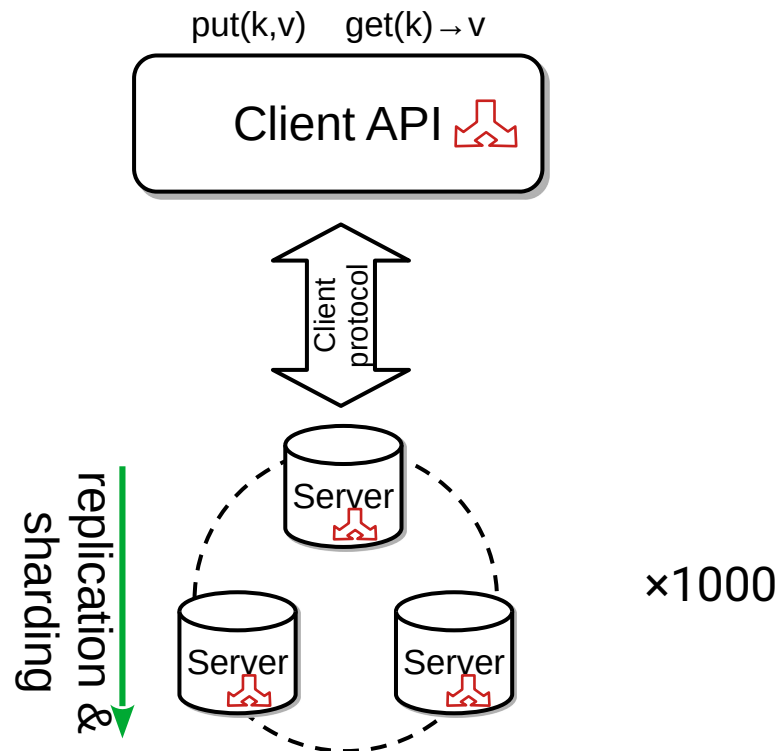# Decentralized data

- <u>Separate data and metadata</u> servers and protocols
  - Client-side caching of metadata

- Avoids data bottleneck

- Examples:
  - Cloud Bigtable
  - APACHE HBASE

put(k,v)    get(k) → v

Client API

Data protocol

Metadata protocol

Metadata server

sharding

Backend Server    Backend Server    ×100

replication

# Fully decentralized

- <u>Consistent hashing</u>

- <u>Epidemic dissemination</u>

- Examples:
    - Dynamo
      (not DynamoDB!) **aws**
    - Apache CASSANDRA



put(k,v)   get(k)→v

Client API

Client protocol

Server

replication & sharding

Server   Server

×1000

# Geo-replication

- Multi-data center replication

- <u>Main challenge</u>: How to shield clients from Inter-DC protocol?
  - Latency
  - Availability

# A taxonomy of some consistency models

Prefer availability

Linearizability

Sequential

CP in CAP

Causal

AP in CAP

PRAM

Strong eventual

Eventual

Writes follow reads    Read your writes    Monotonic writes    Monotonic reads
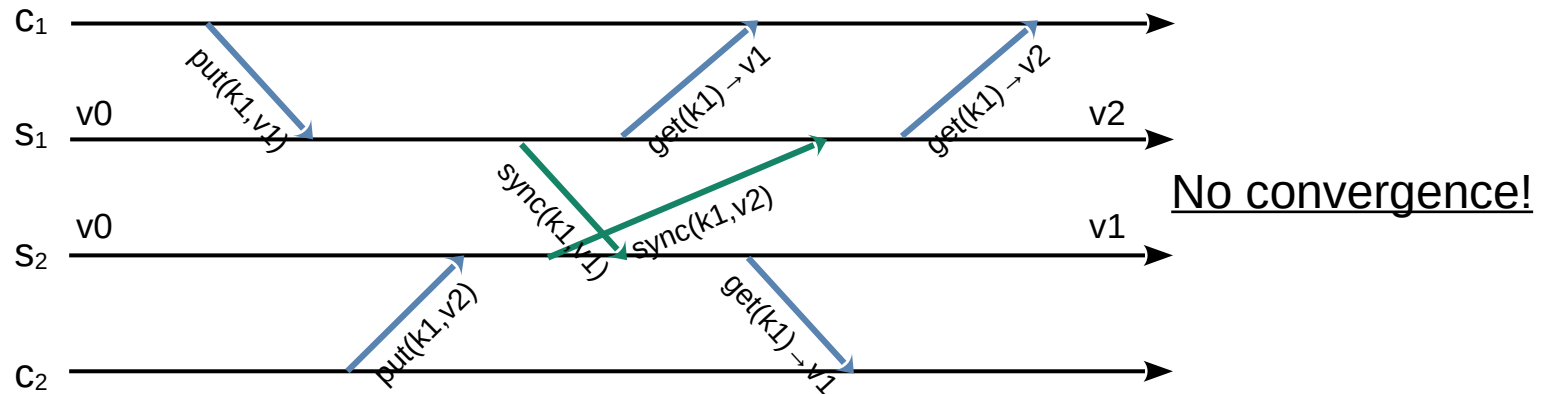
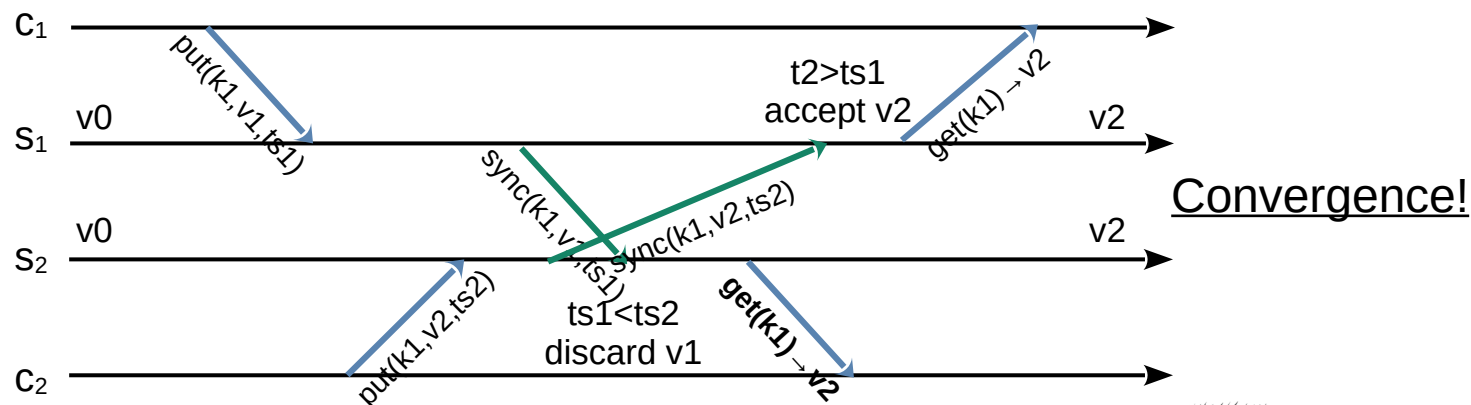Session models

Weak

Avoid session state

¯\_(ツ)_/¯

# 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:
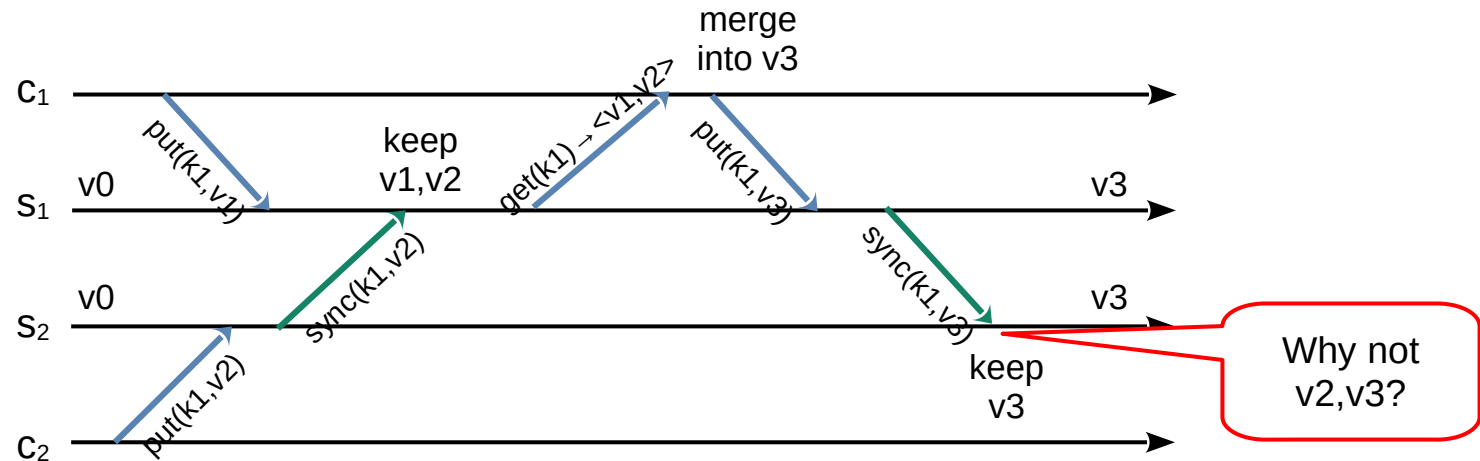  - Physical: when the client is a Web application and the state partitioned by user
  - Logical?

- LWW Register is actually a <u>State-based CRDT</u>… 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 <u>client application</u>
  - Simple LWW (Cassandra)
  - Other State-based CRDTs (Cure)
  - Custom code (Bayou)

- Servers cannot callback into the client application
- How to achieve this?

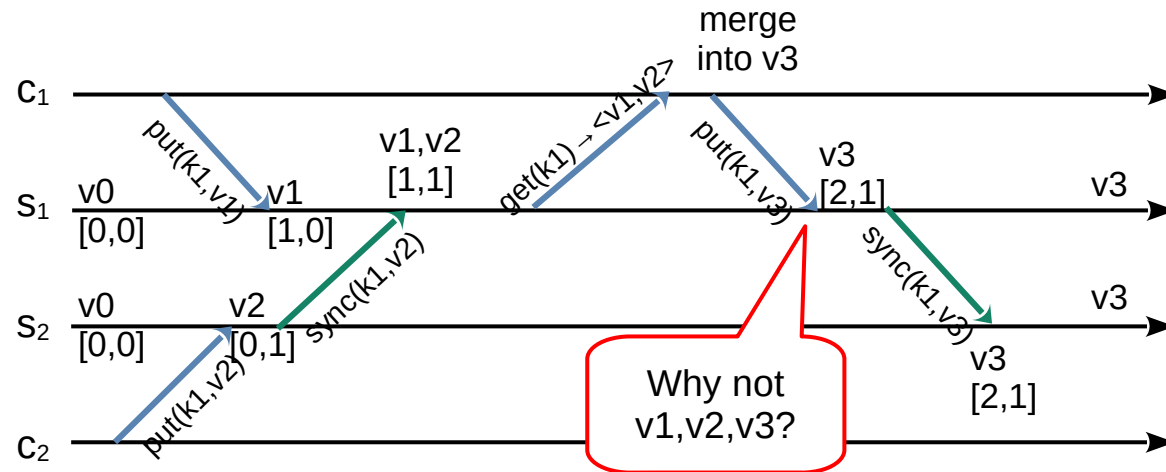# Convergence

- Keep list of conflicting values in each key and return them to the application: $get(k) \rightarrow <v1,...vn>$

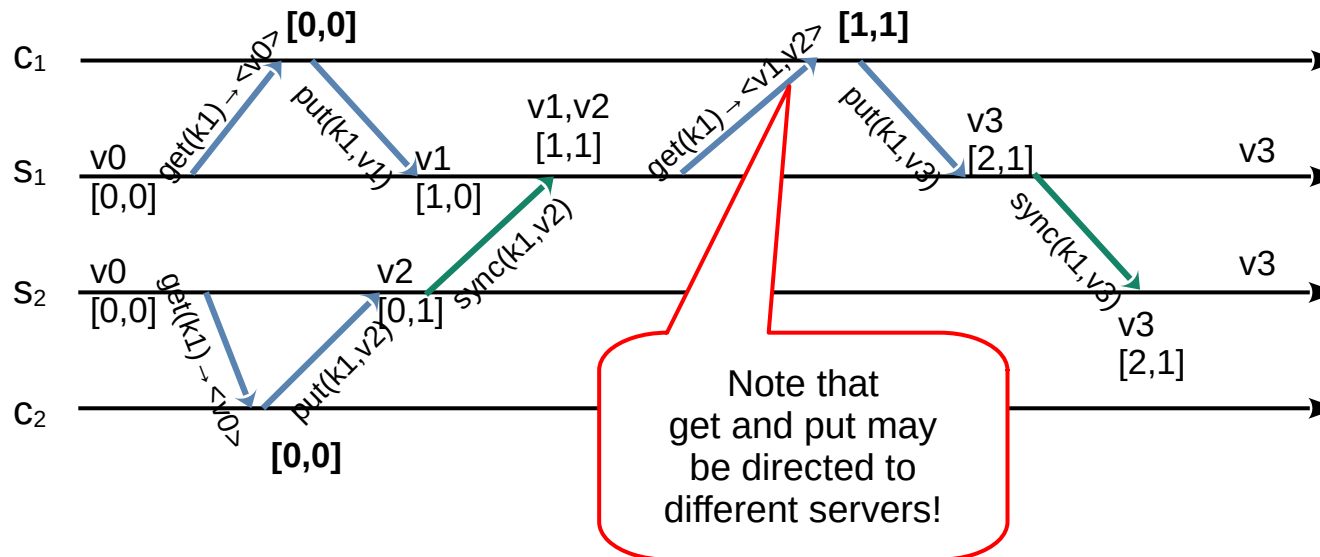- Let the application merge them and replace the list with a single value

# 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

<u>Prefer availability</u>

Linearizability

Sequential

CP in CAP

Is it possible to achieve causal consistency with client-only sessions?

Causal

Strong eventual

PRAM

Eventual

Writes follow reads    Read your writes    Monotonic writes    Monotonic reads

Session models

Weak

<u>Avoid session state</u>

# 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)



Client API

Server protocol

Server

Server protocol

Server

Server protocol

Server

Server

Server

Server

get operation is trivial

# Managing client context

- Client context needs to keep nearest dependencies

- Simple (approximate) strategy:
    - After a put(k,...)→vers operation:
        - Clear all dependencies
        - Insert (k,vers)
    - After a get(k,...)→ (v,vers) operation:
        - If (k,...) in context, update to (k,vers)
        - If not, insert (k,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)



Data center 1

Data center 2

ctx: ⟨⟩

ctx: ⟨(k2,21)⟩

ctx: ⟨(k2,21), (k1,31)⟩

$c_1$

$s_{1.1}$   k1→va   11

put_after(k2,vc,⟨⟩,0)

$s_{1.2}$   k2→vb   11   k2→vc   21

21

put_after(k1,vd,⟨(k2,21)⟩,0)

31

k1→vd   31

$s_{2.2}$   k2→vb   11

$s_{2.1}$   k1→vd   31

Equivalent servers (same shard)

Paradox observed (vd;vb) if k1, k2 read from DC2

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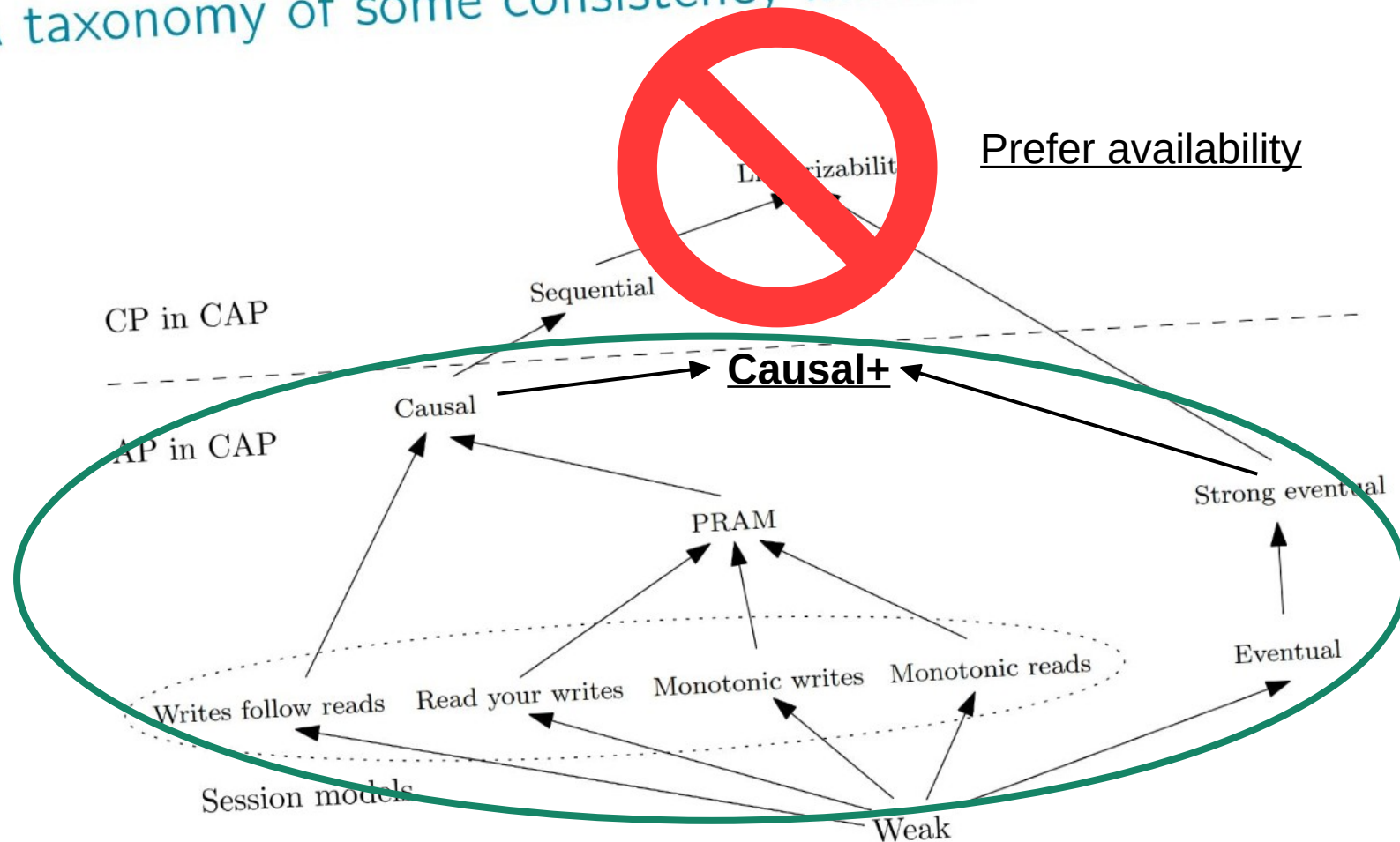
- Propagation is not ordered, as it is potentially from different data centers
- A value is committed remotely only after waiting for dependencies with the dep_check operation of corresponding servers



Large Scale Distributed Systems

# 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

  – <u>Last Writer Wins</u> when propagating to remote DC

# A taxonomy of some consistency models

Prefer availability

CP in CAP

Sequential

Linearizabilit

**Causal+**

Causal

AP in CAP

Strong eventual

PRAM

Writes follow reads    Read your writes    Monotonic writes    Monotonic reads    Eventual

Session models

Weak

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