

Restart-Rollback: a novel fault model for Efficient TEE replication

Baltasar Dinis

Instituto Superior Técnico (ULisboa) / INESC-ID /
MPI-SWS

NDSS – February 28 - March 2, 2023

joint work with

Rodrigo Rodrigues (INESC-ID),

Peter Druschel (MPI-SWS),



TÉCNICO
LISBOA



Trusted execution environments (TEEs)



TÉCNICO
LISBOA



Teechain: A Secure Payment Network with Asynchronous Blockchain Access

Joshua Lind

Oded Naor

Ittay Eyal

ROTE: Rollback Protection for Trusted Execution

Sinisa Matetic
ETH Zurich

Mansoor Ahmed
ETH Zurich

Kari Kostiainen
ETH Zurich

Aritra Dhar
ETH Zurich

Avocado: A Secure In-Memory Distributed Storage System

Maurice Bailleu¹, Dimitra Giantsidi¹

Vasilis Gavrielatos¹, Do Le Quoc², Vijay Nagarajan¹, Pramod Bhatotia^{1,3}

¹University of Edinburgh ²Huawei Research ³TU Munich

Robust P2P Primitives Using SGX Enclaves

Yaoqi Jia¹

Shruti Tople²

Tarik Moataz³

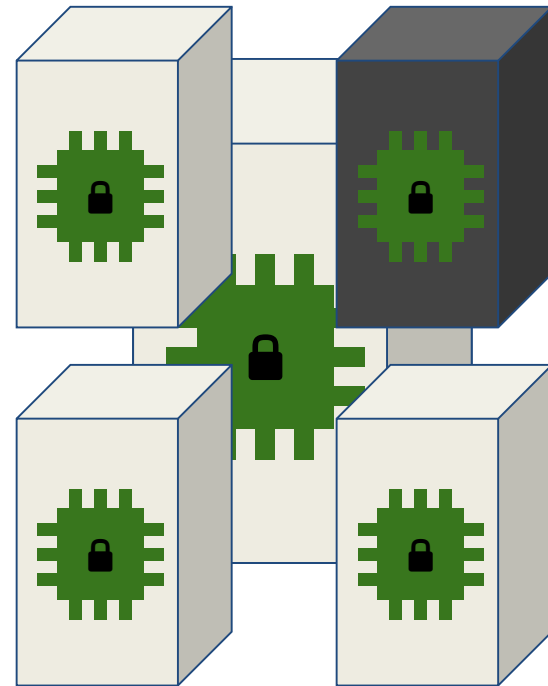
Deli Gong¹

Prateek Saxena⁴

Zhenkai Liang⁴

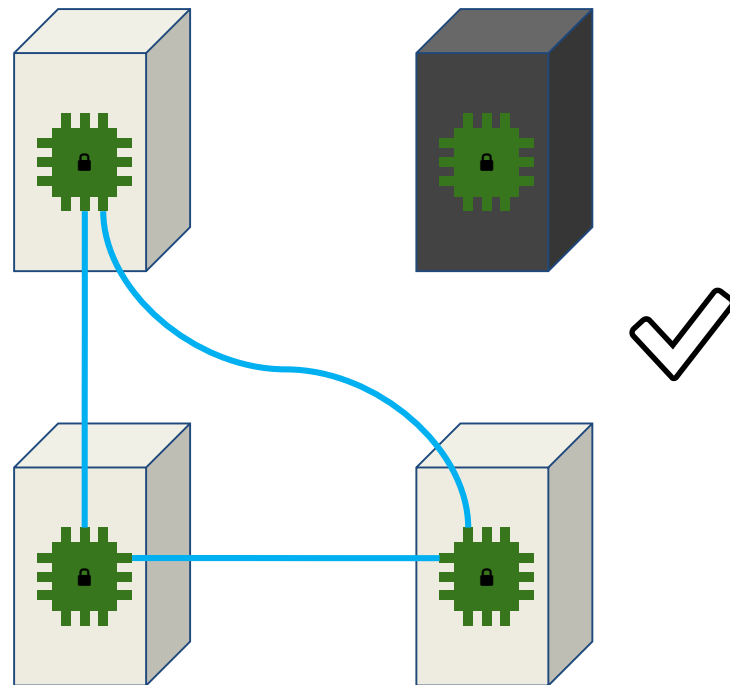
¹ACM Member [jiayaoqijia, gnnnnng]@gmail.com ²Microsoft Research t-shtopl@microsoft.com

³Aroki Systems tarik@aroki.com ⁴National University of Singapore [prateeks, liangzk]@comp.nus.edu.sg



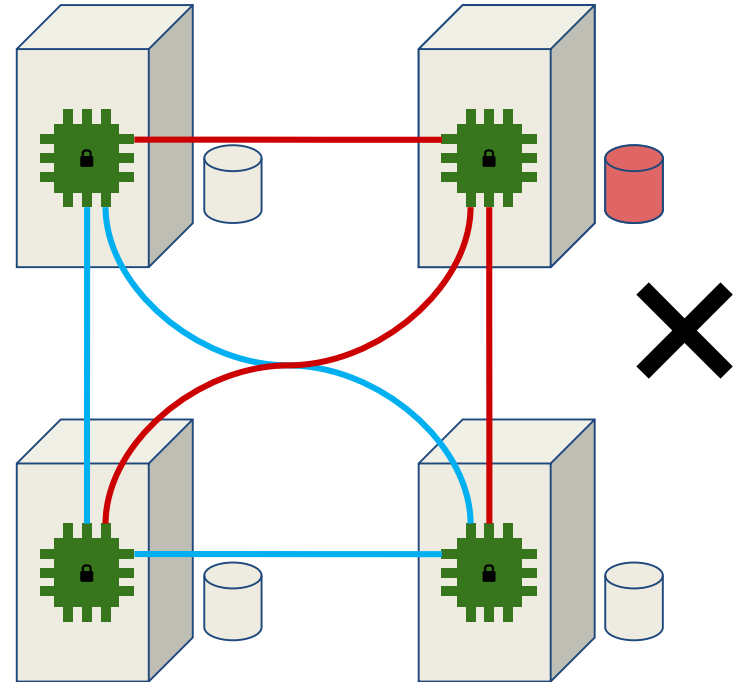
What is the right model for TEE replication?

- Crash Fault Tolerance (CFT) assumes
 - Correct computation
 - Replicas may fail silently



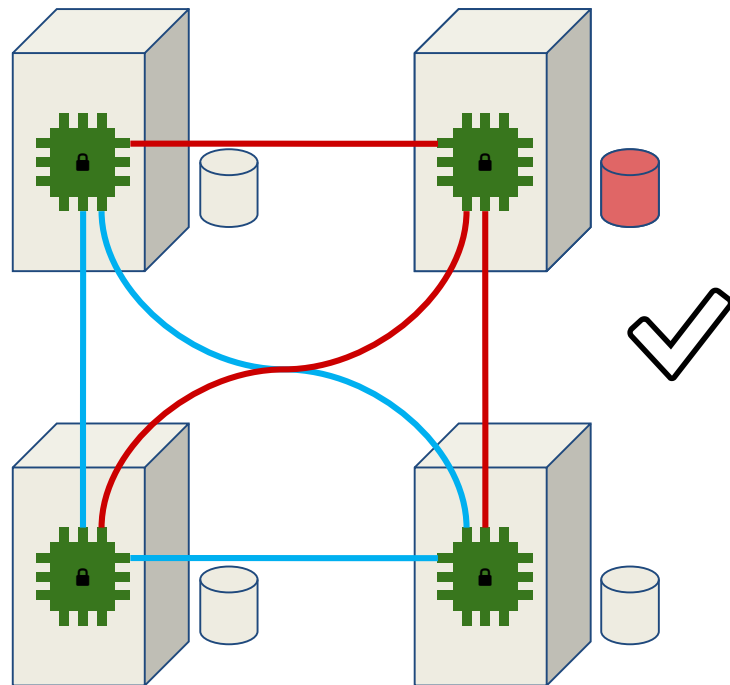
What is the right model for TEE replication?

- Persistent storage complicates things, since it is not under direct control of the TEE
- An adversary may replace the persistent storage with an older version: a **rollback attack**



What is the right model for TEE replication?

- Byzantine Fault tolerance (BFT) tolerates arbitrary faults
 - By extension, rollbacks
- BFT mechanisms are **expensive**



What is the right model for TEE replication?




TÉCNICO
LISBOA



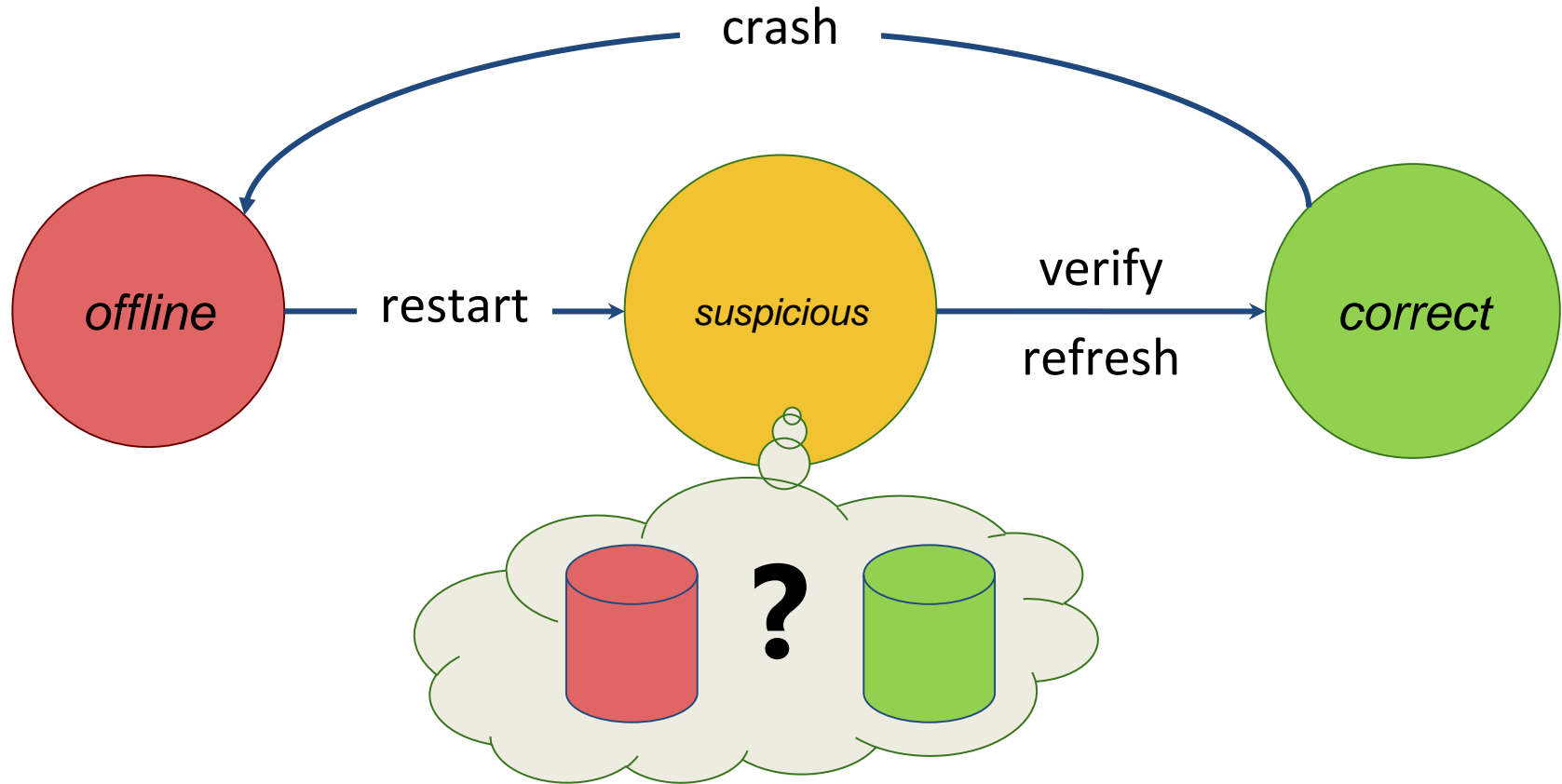
- CFT is sufficient, *unless* the TEEs have external persistent state
 - If so, the system is vulnerable to **rollback attacks**
- For this reason, BFT is necessary
 - More expensive and intuitively sounds **pessimistic**, given TEE guarantees
- Can we achieve the best of both worlds?

Yes, if the fault model accurately reflects TEE behavior

Talk outline

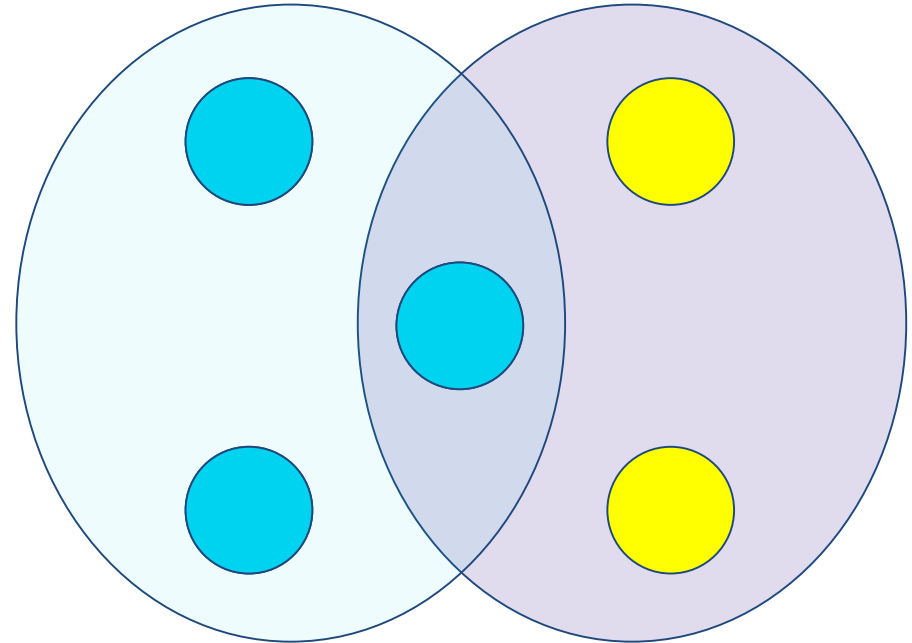
- Motivation
- The Restart-Rollback fault model 
- Adapting replication protocols for the new model
- TEEMS: a metadata service for TEE-grade cloud storage

Restart-rollback model

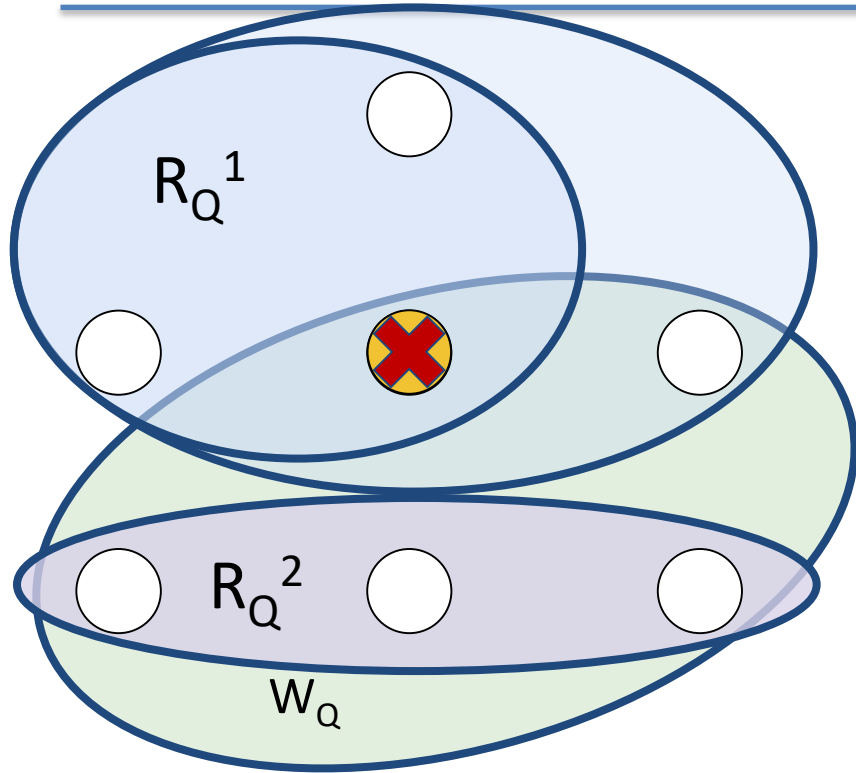


We have a fault model, now what?

- A major primitive for replication protocols are **quorums**:
 - A subset of replicas required to execute an operation
 - Guarantee correctness because pairs of quorums **intersect**
- We need a **quorum system** in the RR model




RR Quorum System Example



- Read and write quorums intersect in at least one non-rolled back replica
- They have different sizes (i.e, they are **asymmetric**)
- Two read quorums may not intersect
- Read quorums scale up when needed (they are **dynamic**)

Talk outline

- Motivation
- The Restart-Rollback fault model
- Adapting replication protocols for the new model 
- TEEMS: a metadata service for TEE-grade cloud storage

Adapting CFT replication protocols



TÉCNICO
LISBOA



1. After a restart, flag replies with **suspicion**
2. Adjust quorum sizes
 - a. When reading system state, use R_Q
 - b. When writing system state, use W_Q
3. Deal with split brain
 1. Because read quorums may not intersect
 2. This will be protocol specific

Distributed register

- Read/write operations
- CFT: ABD
- BFT: Byzantine Quorums

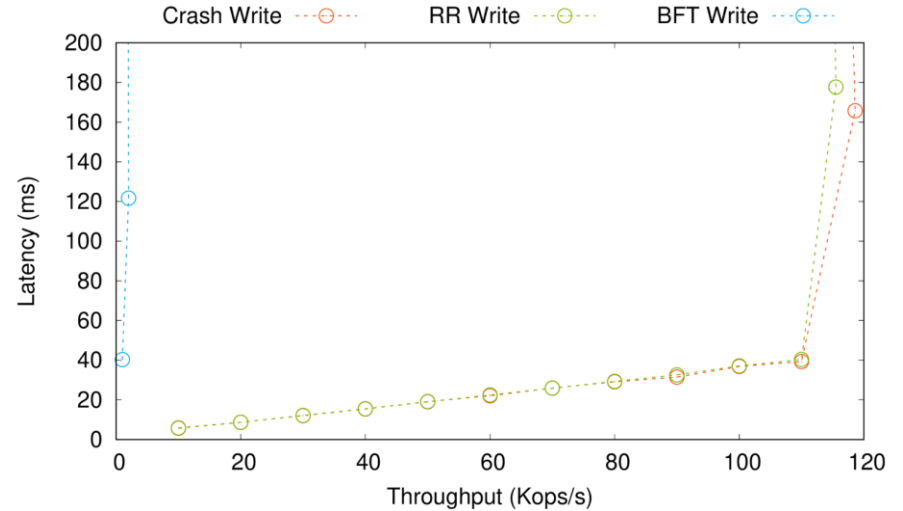
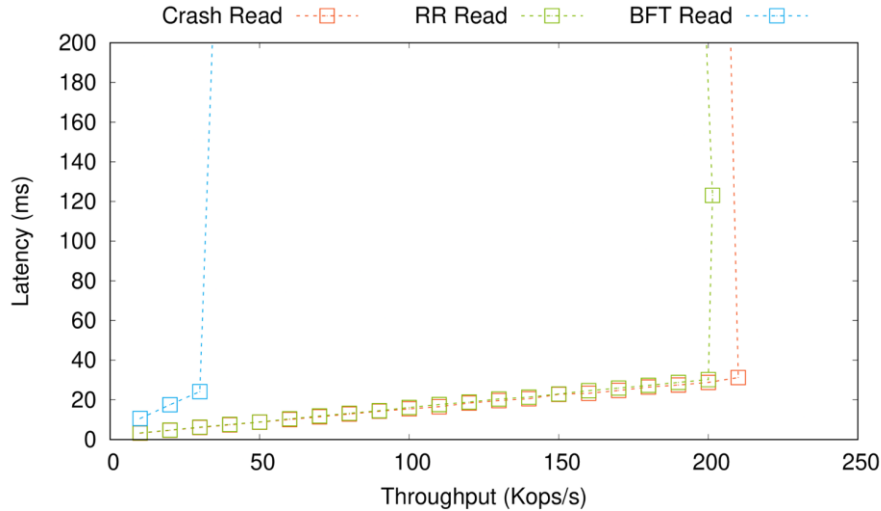
State machine replication

- Any deterministic service
- CFT: Paxos, Raft
- BFT: PBFT, BFT smart, etc.


Evaluation of the read/write protocols



TÉCNICO
LISBOA



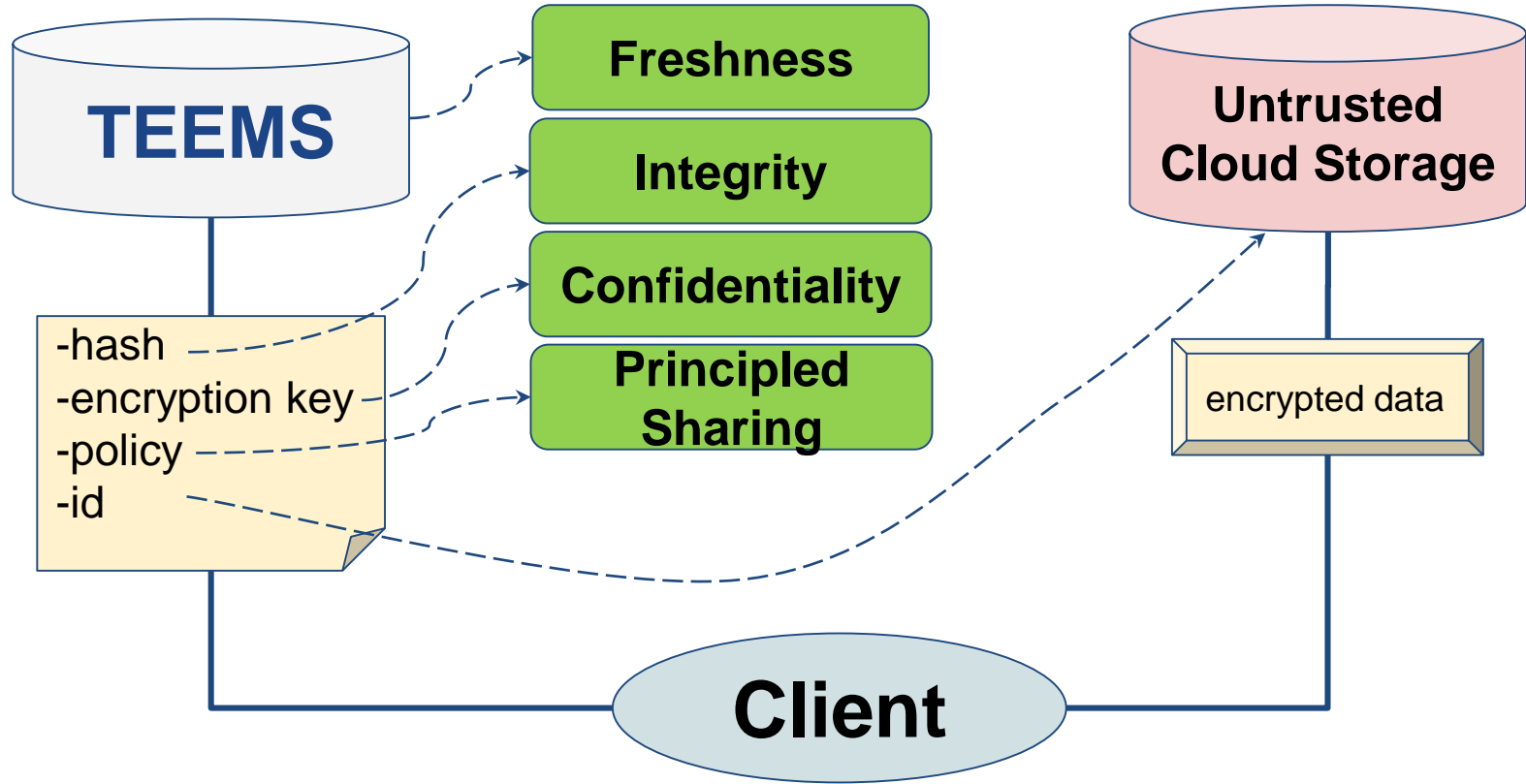
Talk outline

- Motivation
- The Restart-Rollback fault model
- Adapting replication protocols for the new model
- **TEEMS: a metadata service for TEE-grade cloud storage** 

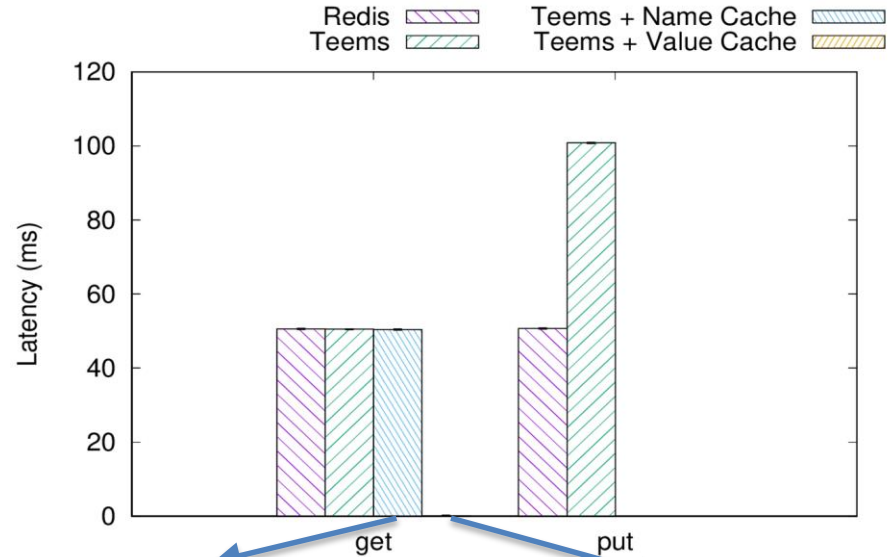
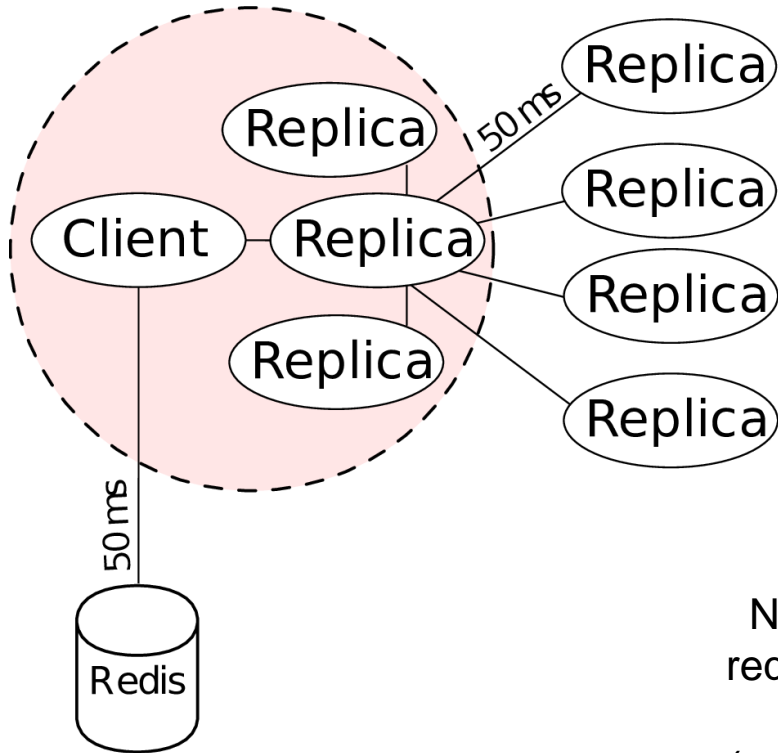
TEEMS (TEE metadata service)

- TEEs are key element to increasing trust in cloud services
 - E.g., Azure confidential computing, Google confidential VMs
- Can we build cloud storage with similar guarantees?
 - Freshness
 - Integrity
 - Confidentiality
 - Principled Sharing
- Strawman solution: Encryption
 - does not ensure **freshness**
 - no mechanism for **principled sharing** among clients

Leverage our new protocols to build **trusted metadata service**



Evaluation of TEEMS-based storage



Name hint still requires fetching the data (expensive step)

Caching the entire value allow for a fast (local) check of the metadata



- Restart-Rollback **precisely** models TEEs with persistent state
 - Ensuring correctness
- With RR, we can intuitively adapt existing CFT protocols, yielding similar performance
- Using RR, we built TEEMS a system that provides the analog to the “confidential VM” paradigm for cloud storage

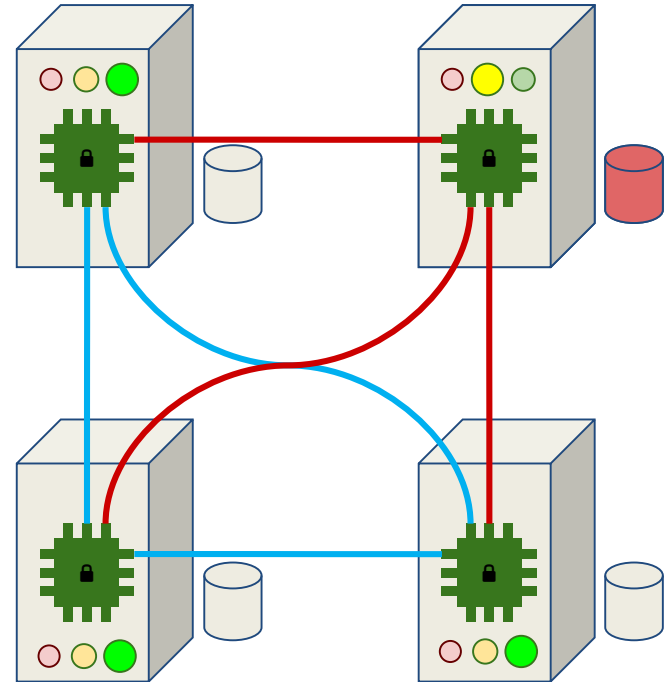
Thank you



1. Pick your favourite CFT protocol (read/write register, replicated state machine)
 - Remember, **computation is correct**
2. Identify quorums where the protocol reads state
 - Replace with R_Q
3. Identify quorums where the protocol writes state
 - Replace with W_Q
4. We're done?

Restart-rollback model

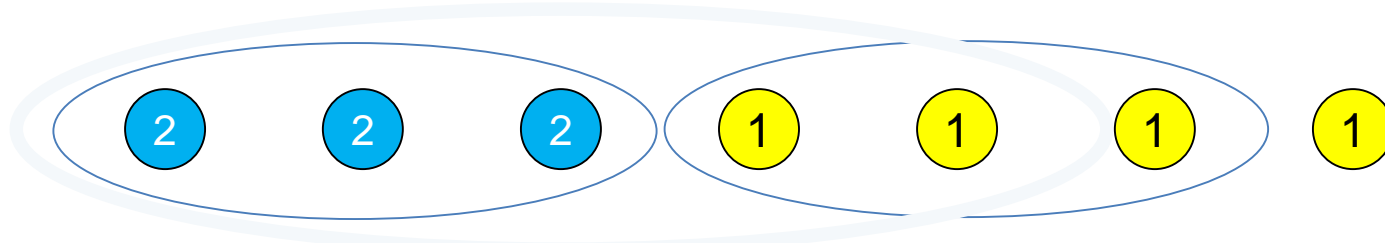
- If more than F replicas are simultaneously in the  state, the system becomes **unavailable**
- If more than M_R replicas have simultaneously been rolled back (i.e, the storage is in the  state) then the system suffers a **rollback**



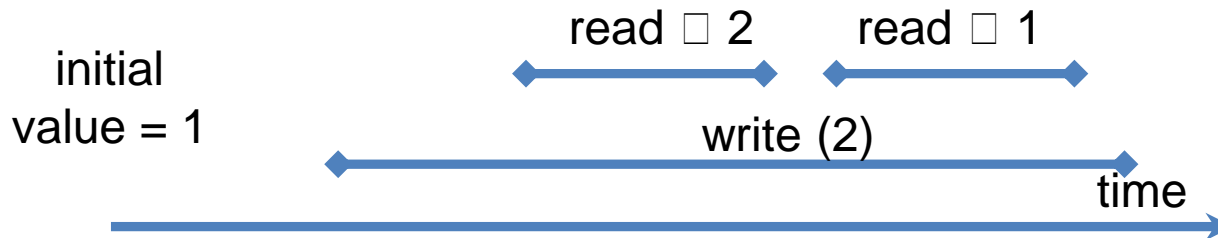
Reading in parallel with ongoing write



- $M_R = 4, F = 2, N=7, W_Q = 5$



- Write replaces version 1 (yellow) with version 2 of data (blue)
- First read sees unanimously version 2
- Second read sees unanimously version 1





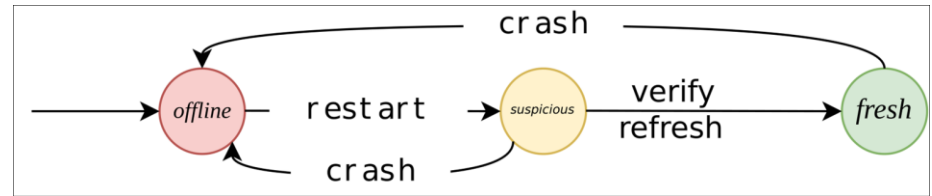
Model	Parameters	N	Quorum Size
Crash [AD ICSE'76]	f	$2f + 1$	$f + 1$
Byzantine [PSL jACM'80]	f	$3f + 1$	$2f + 1$
Hybrid [MP TPDS'91, CLNV SRDS'02, C+ SOSP'09]	u, r	$2u + r + 1$	$u + r + 1$

- Replication in the asynchronous model relies on quorums
- Smaller quorum sizes lead to more efficient operations

Deriving quorum systems

- During each quorum-RPC, count number of replies with suspicion flag: s
 - s is the number of replicas in the quorum that **may** have been rolled back
- Per-RPC maximum number of rolled back nodes: $\min(s, M_R)$

M_R : simultaneous rollbacks
 F : simultaneous crashes



Deriving quorum systems

M_R : simultaneous rollbacks
 F : simultaneous crashes
 s : restarts in this RPC
 $\min(s, M_R)$: maximum number of rolled back nodes in this RPC

- Correctness (safety) intersection requirement becomes:
$$R_Q + W_Q > N + \min(s, M_R)$$
- In conjunction with liveness constraints, we derive the following:

$$N = \max(M_R, F) + F + 1$$

$$W_Q = \max(M_R, F) + 1$$

$$R_Q = F + \min(s, M_R) + 1$$

CFT Distributed Register

Read()

```
Send <Read> to all replicas
Wait for majority of replies
if unanimous:
    return value
else // writeback
    Send <Write, value, hts>
        to replies with ts < hts
    Wait for majority of replies
return value
```

Write(value)

```
Send <Read> to all replicas
Wait for majority of replies
Send <Write, value, hts + 1>
Wait for majority of replies
```

return

to replies with **ts < hts**

Wait for **majority of** replies

return value

Distributed Register in the new model



TÉCNICO
LISBOA



```
Read()
```

```
Send <Read> to all replicas
```

```
Wait for  $R_Q$  replies
```

```
if unanimous:
```

```
    return value
```

```
else // writeback
```

```
    Send <Write, value, hts>
```

```
    to replies with ts < hts
```

```
    Wait for  $W_Q$  replies
```

```
return value
```

```
Write(value)
```

```
Send <Read> to all replicas
```

```
Wait for  $R_Q$  replies
```

```
Send <Write, value, hts + 1>
```

```
Wait for  $W_Q$  replies
```

```
return
```

But: Read quorums may not intersect! 