# Securing Passive Replication Through Verification

**_Bruno Vavala_**[1,2], Nuno Neves[1], Peter Steenkiste[2]

[1]University of Lisbon (Portugal)

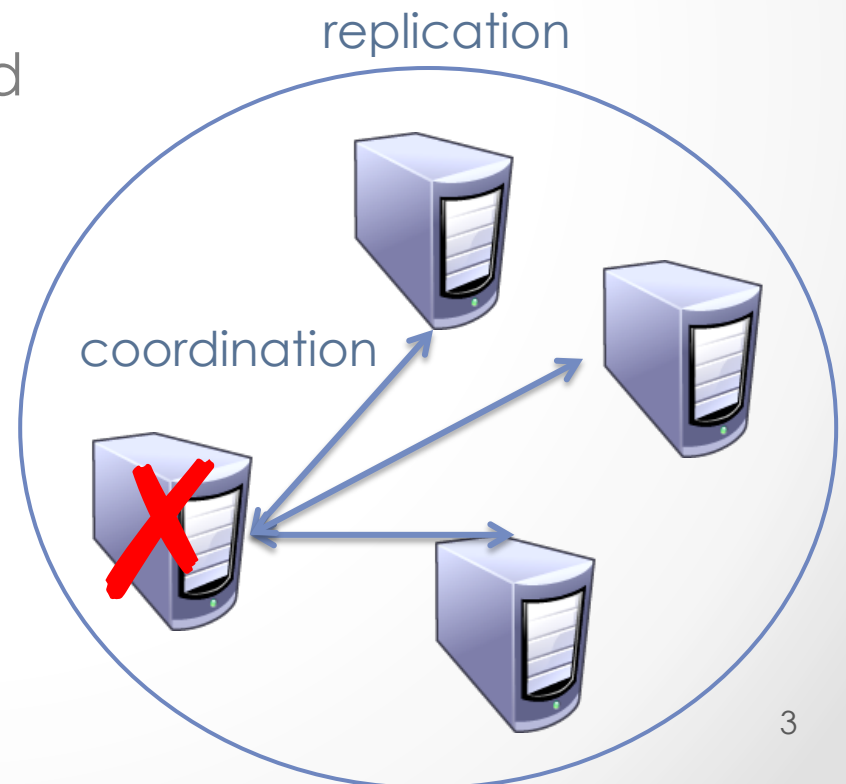[2]Carnegie Mellon University (U.S.)

# Outline

- **Motivation and background**

- Goals

- Architecture Design & System Operations

- Evaluation

- Takeaways

# Fault-Tolerance

- Service continuity has to be ensured in case of failure

- Components have to be replicated

- Replicas must be coordinated

replication

coordination

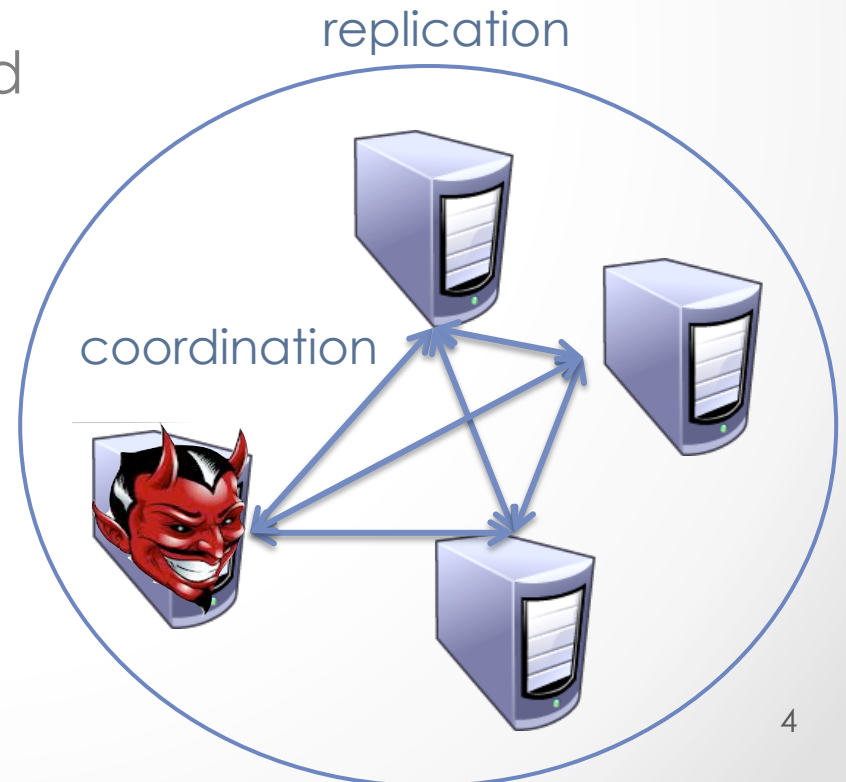# Fault-Tolerance

- Service continuity has to be ensured in case of failure

- Components have to be replicated

- Replicas must be coordinated

- Arbitrary failures require
  **+**replicas
  **+**coordination

replication

coordination

# Replication

2 main design choices

Active Replication
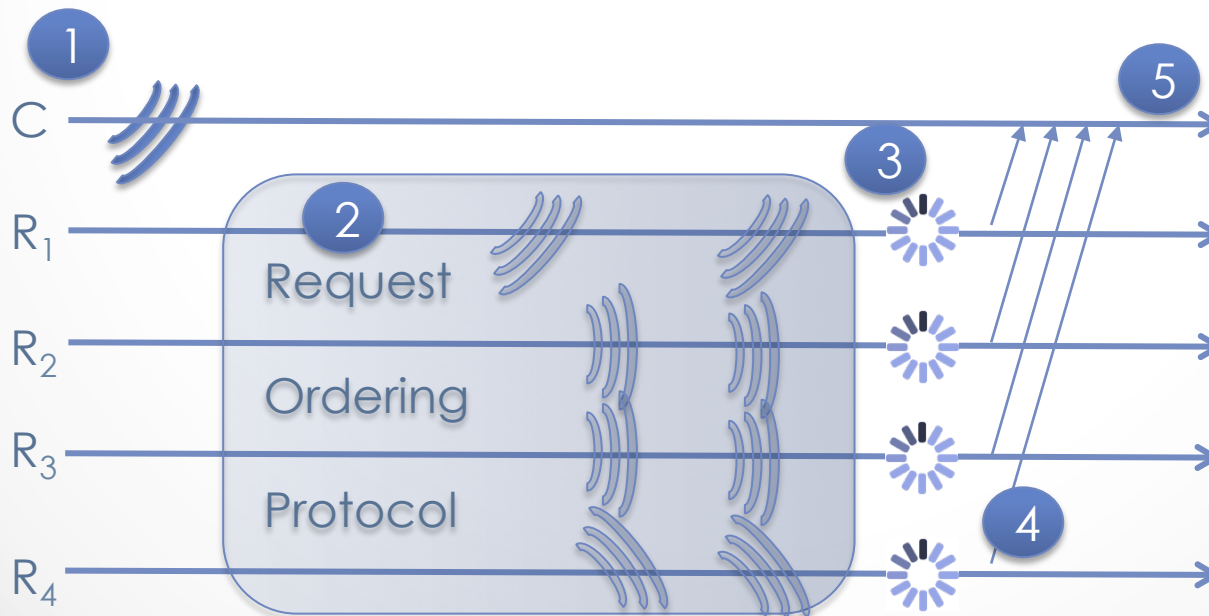(State Machine Replication)

vs.

Passive Replication

# Active Replication (AR)

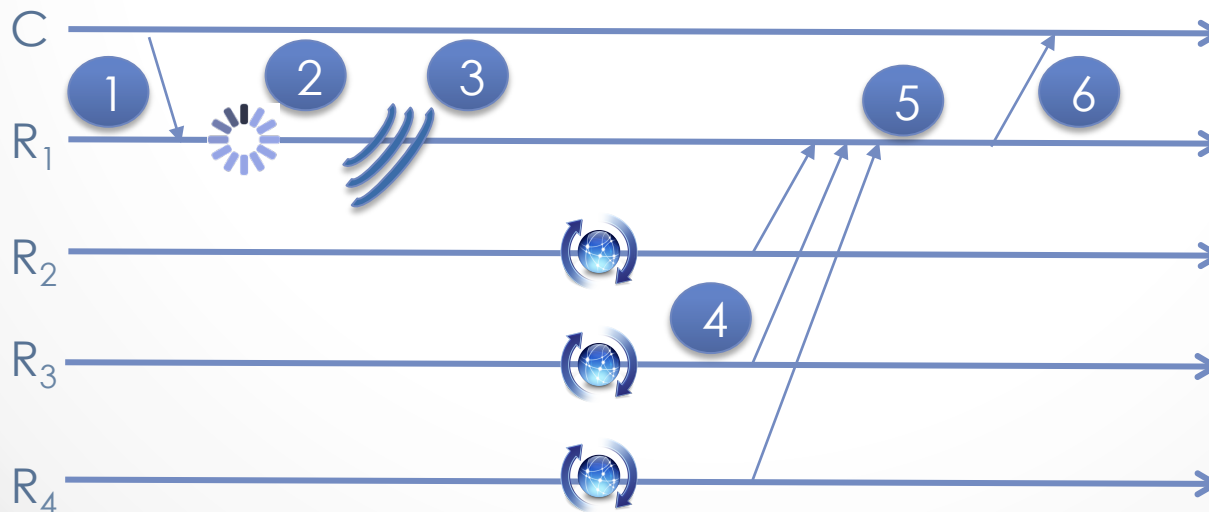State Machine approach:

1. System receives the requests

2. Requests are ordered ("many" messages)

3. Enough replicas execute them

4. Each replica returns an answer

5. Answers are voted

# Passive Replication (PR)

1. Primary receives the requests
2. Requests are executed
3. State updates are broadcast

4. Backups apply updates and return ACK
5. Primary votes on ACKs
6. Primary replies to client



7

# Current BFT Solutions

## AR

- **PBFT (OSDI'99)**
  Seminal practical SMR work
-
-
-
-
-
- **BFT-SMaRt (DSN'14)**
  High performance

…and
many

.

.

many
others!

## PR

∅

# Why no PR solutions?

# Why no PR solutions?

system

client

AR

$R_1$ $\dashrightarrow$

$R_2$ $\dashrightarrow$

Voter

correct answer ✔

$R_3$ $\dashrightarrow$

$R_4$ $\dashrightarrow$

- Enough redundancy to extract correct answer

# Why no PR solutions?

system                                                                    client



AR

- $R_1$
- $R_2$
- $R_3$
- $R_4$

Voter

correct answer ✔

PR

- $R_1$
- $R_2$
- $R_3$

correct ?

?

- Challenge: how to verify the result efficiently?
- Trivial *inefficient* solution: re-execute the service

# Pros & Cons

|  | AR | PR |
|---|---|---|
| **Byzantine FT** | ✔ | ✗ |
| **Replicas** | $2f+1$ | $2f+1$ |
| **Re-Computations** | $O(n)$ | $O(1)$ |
| **Message size** | \|request\| + \|input\| | \|reply\| + \|update\| |
| **Non-determinism** | ✗ | ✔ |

*"While some consensus algorithms, such as Paxos […] have started to find their way into those systems, their uses are limited mostly to the maintenance of the global configuration information in the system, not for the actual data replication."* – L. Lamport et al.

# Outline

- Motivation and background
- **Goals**
- Architecture Design & System Operations
- Evaluation
- Takeaways

# Goals

**Fault-tolerant** & **resource-efficient** & **simple**
replicated architecture for unmodified services

**Challenges**

- Protect the service results from malicious failures

- Efficient verification of the results

- Ensure that state updates are correctly propagated

- Ensure that client gets correct and consistent results

14

# Outline

- Motivation and background
- Goals
- **Architecture Design & System Operations**
- Evaluation
- Takeaways

# V-PR

## Verified Passive Replication

# Best of Both Worlds

| | AR | PR | | V-PR |
|---|---|---|---|---|
| **Byzantine FT** | ✔ | ✘ | | ✔ |
| **Replicas** (w/ trust assumptions) | 2f+1 | 2f+1 | | 2f+1 |
| **Executions** | O(n) | O(1) | | O(1) |
| **Message size** | \|request\| + \|input\| | \|reply\| + \|update\| | | \|reply\| + \|update\| |
| **Non-determinism** | ✘ | ✔ | | ✔ |

# TCC Overview

- Trusted Computing Component
  - It performs actua
  - It provides trusted
  - It has internal regi                                    code

- Primitives
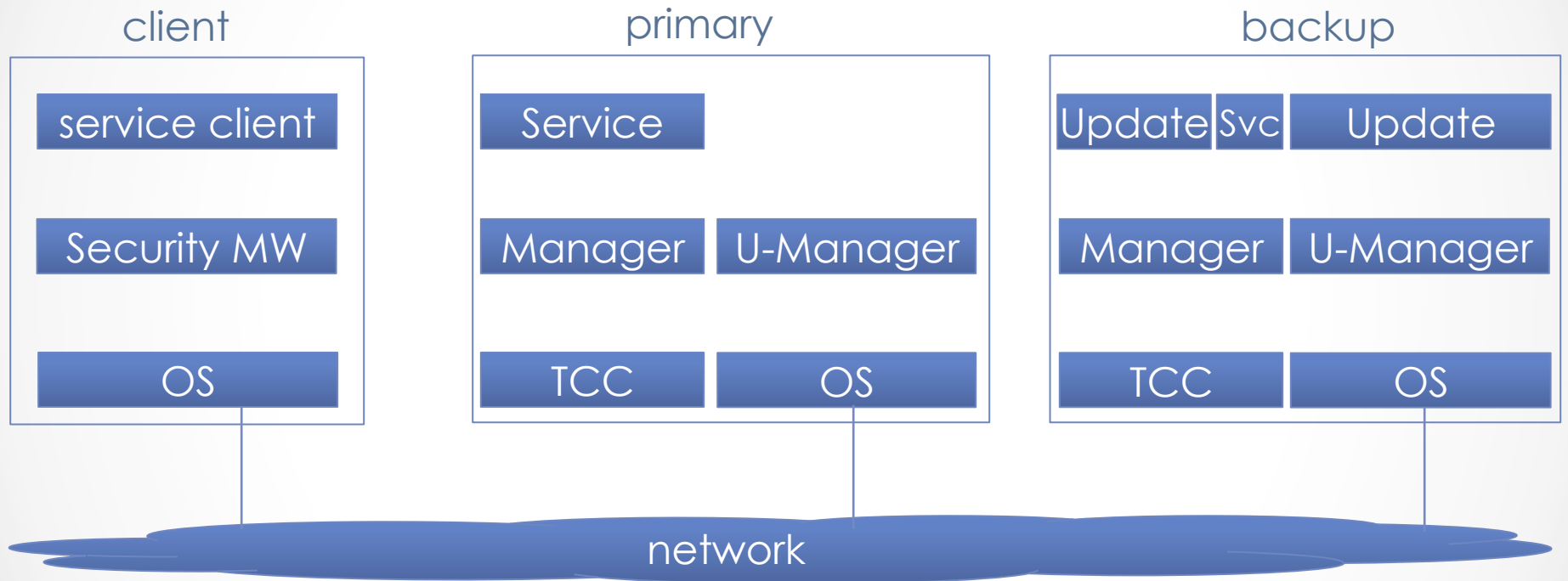  - **put**(data, ID)/**get**                                          ernal
    storage. Only the sa      can store and retrieve data
  - **execute**(code, input). TCC-backed isolated execution of arbitrary code.
    Running code is identified for ID-based operations
  - **attest**(). TCC signature that could carry information on running code and results
  - **create/get/incr_counter**(ID, name). Access controlled Trusted counters. Only ID
    can read or modify them
  - **verify**(). Check validity of attestation, through manufacturer certificate

No different assumptions with respect to previous works, just a more powerful TCC!
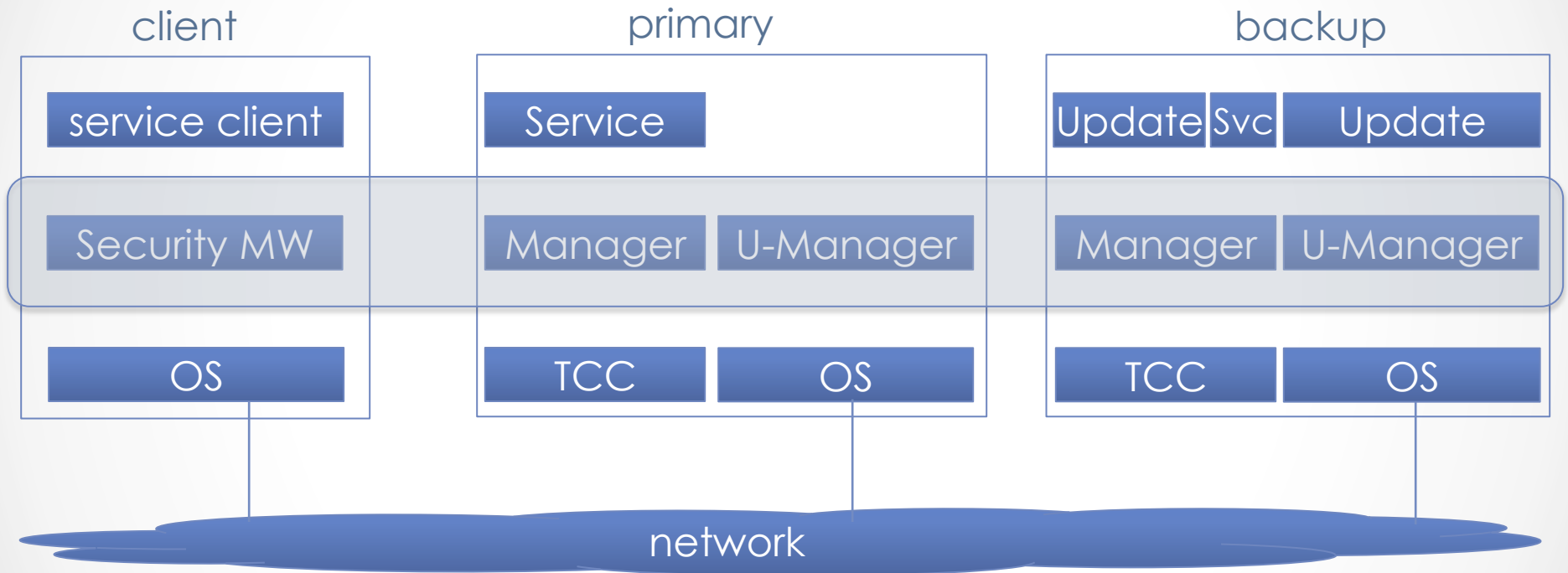
# Model

- TCC is crash-only

  Rest of the system can fail arbitrarily (Byzantine)

- TCC only usable through primitives

- Correct Majority of replicas

- Asynchronous model for safety, partially synchronous oth.

- Model does not consider:

  o Denial of Service attacks

  o Physical tampering (at least not to the TCC hardware)

  o Service vulnerabilities
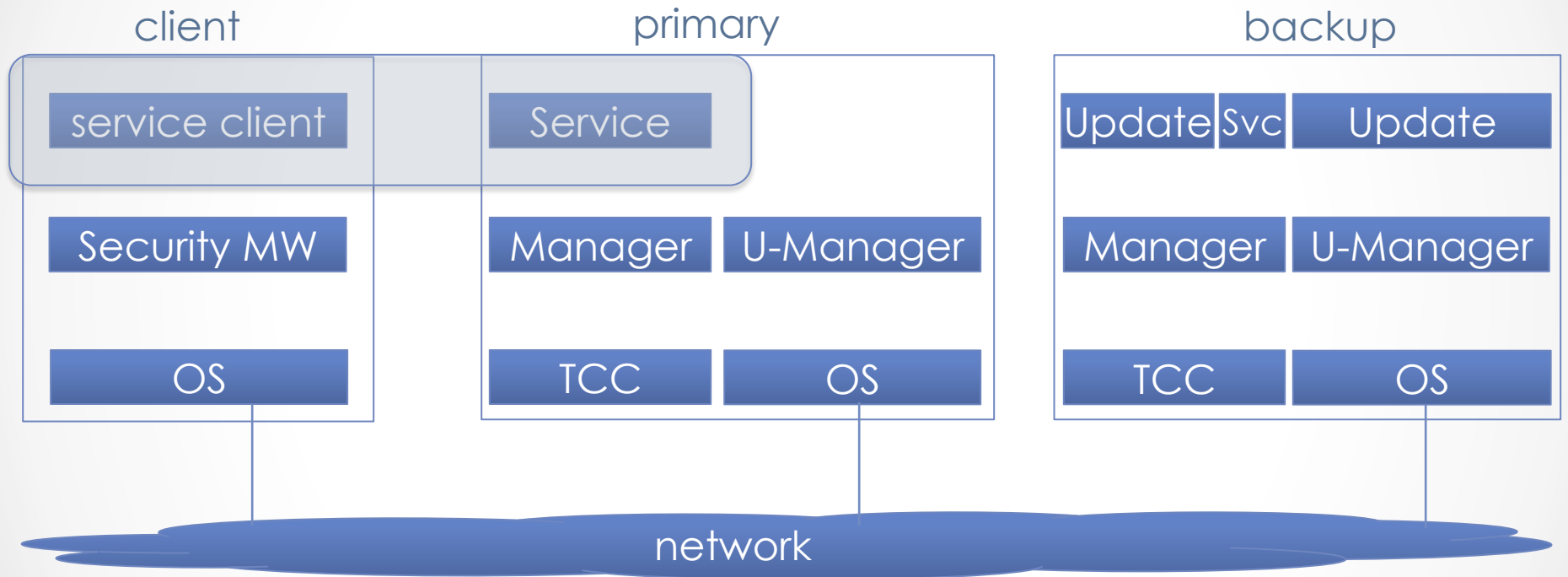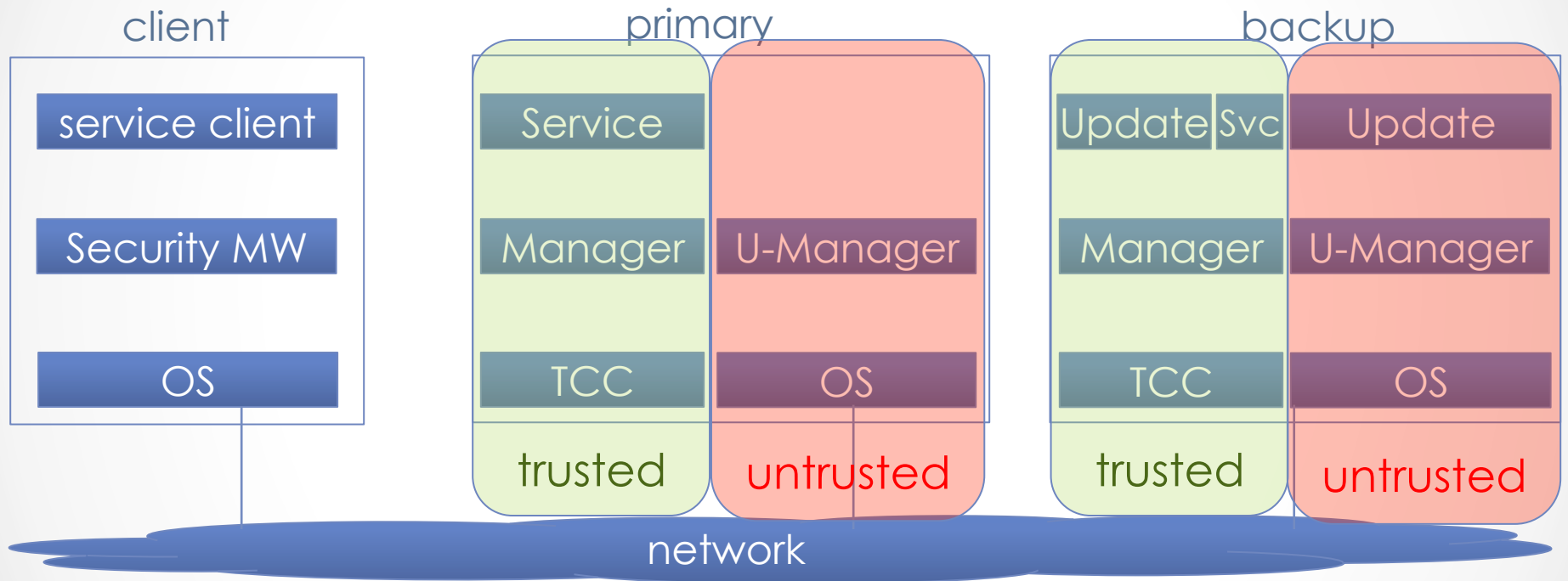
# V-PR Architecture

**client**

- service client
- Security MW
- OS

**primary**

- Service
- Manager | U-Manager
- TCC | OS

**backup**

- Update | Svc | Update
- Manager | U-Manager
- TCC | OS

network

# V-PR Architecture

client                  primary                  backup

| client | primary | backup |
|---|---|---|
| service client | Service | Update \| Svc   Update |
| Security MW | Manager   U-Manager | Manager   U-Manager |
| OS | TCC   OS | TCC   OS |

network

- Core components:  SMW,  Manager,  U-Manager
- Update service only applies state updates

21

# V-PR Architecture

client         primary         backup

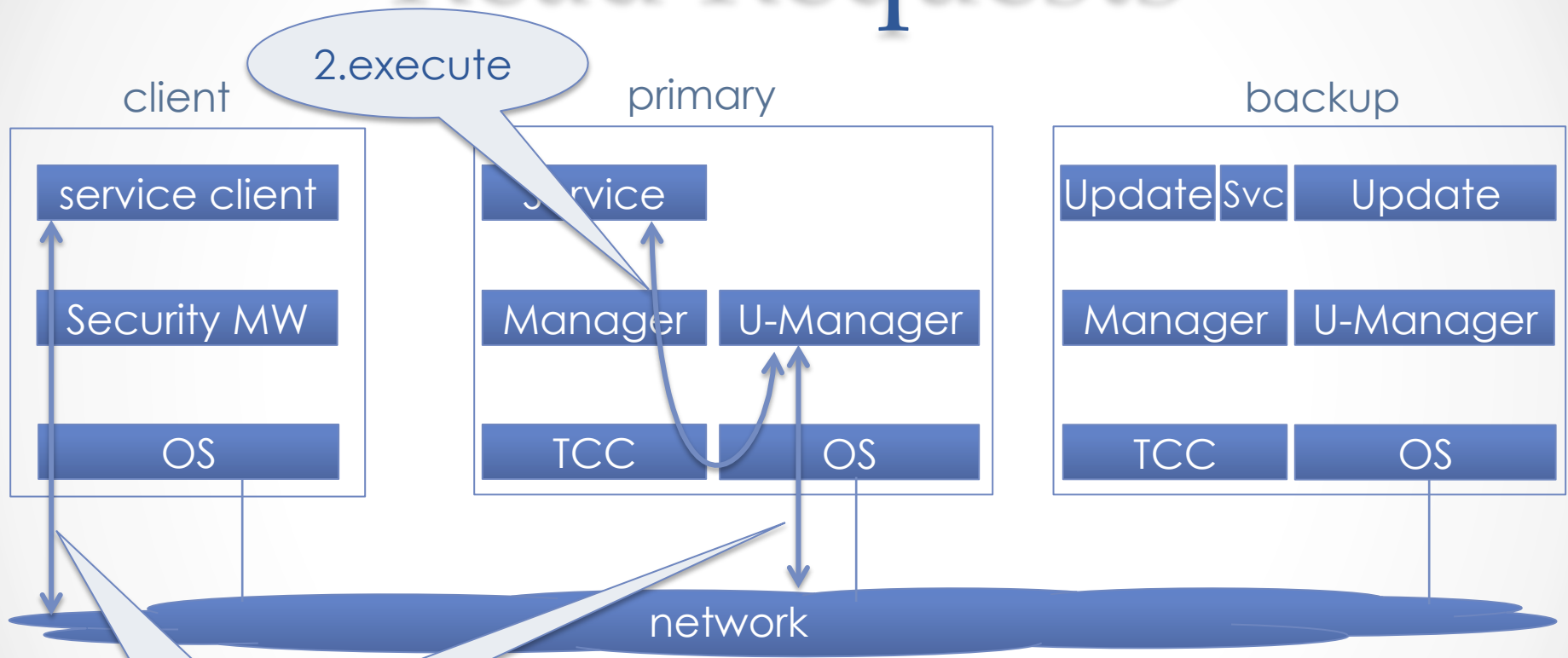| client | primary | backup |
|--------|---------|--------|
| service client | Service | Update Svc Update |
| Security MW | Manager U-Manager | Manager U-Manager |
| OS | TCC OS | TCC OS |

network

- Service Client and Service are not modified
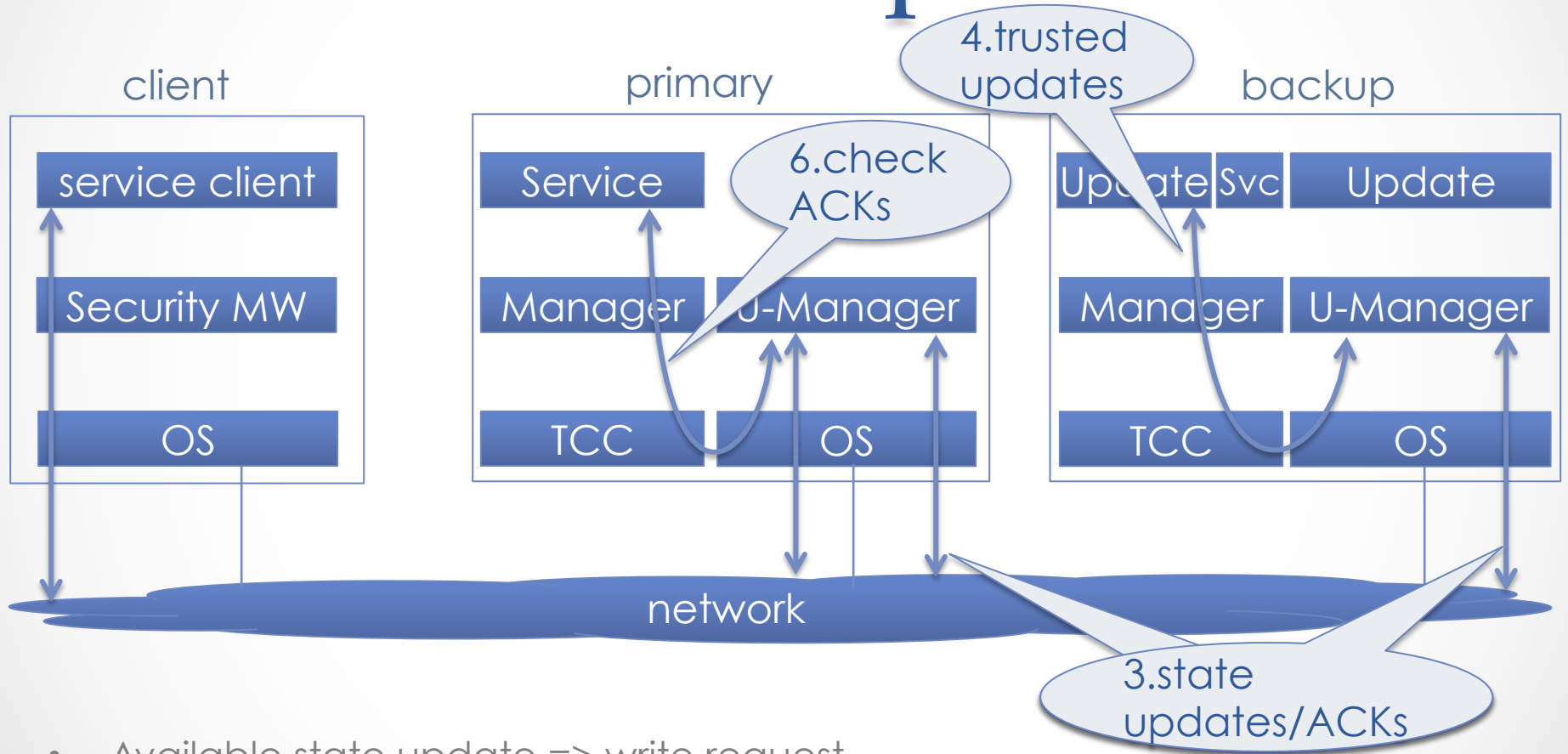- Important effort to make V-PR service oblivious

22

# V-PR Architecture



- Dual failure model (crash+Byzantine)
- Two execution environments with different Trust assumptions
- Entry point: *execute*(Manager) to call TCC service

23

# Read Requests



2.execute

client       primary       backup

service client | Service | Update | Svc | Update

Security MW | Manager | U-Manager | Manager | U-Manager

OS | TCC | OS | TCC | OS

network

1.client request/reply

- Client SMW can verify primary's execution and establish a session key with the Manager
- No state updates => read request
- 2 messages

24

# Write Requests



- Available state update => write request
- 4 steps (of message passing) overall

# Outline

- Motivation and background

- Goals

- Architecture Design & System Operations

- **Evaluation**

- Takeaways

# Evaluation

# Implementation

- Message passing with ZeroMQ

- TCC with XMHF-TrustVisor (S&P'10, S&P'13)

- Full SQLite database engine
  - VPR-ed SQLite
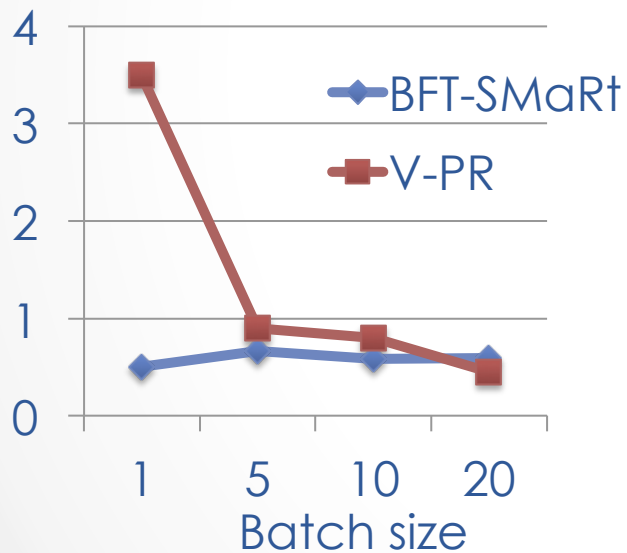
- OS-free implementation
  - very small TCB

- Against recent AR schemes:
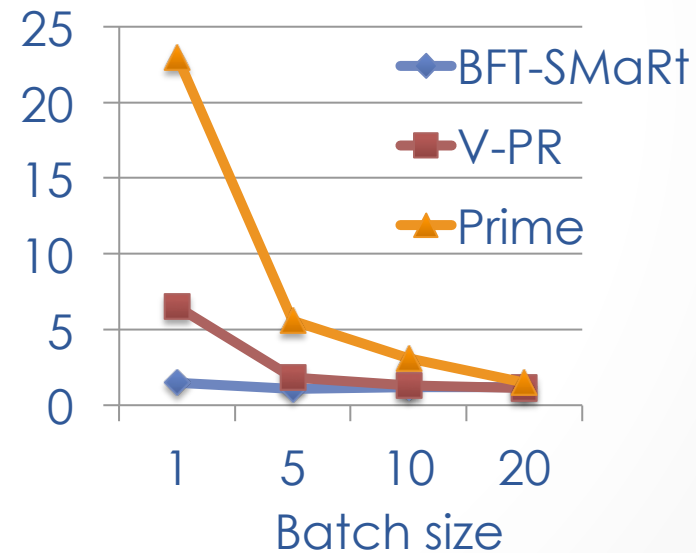  - BFT-SMaRt        (IEEE DSN'14)
  - Prime        (IEEE TDSC'11)

trusted environment

| Service |
|---|
| Manager |
| TrustVisor |
| XMHF |
| Hardware |

**TCC**

# Performance

- Overhead comparison among BFT-SMaRt, Prime and V-PR
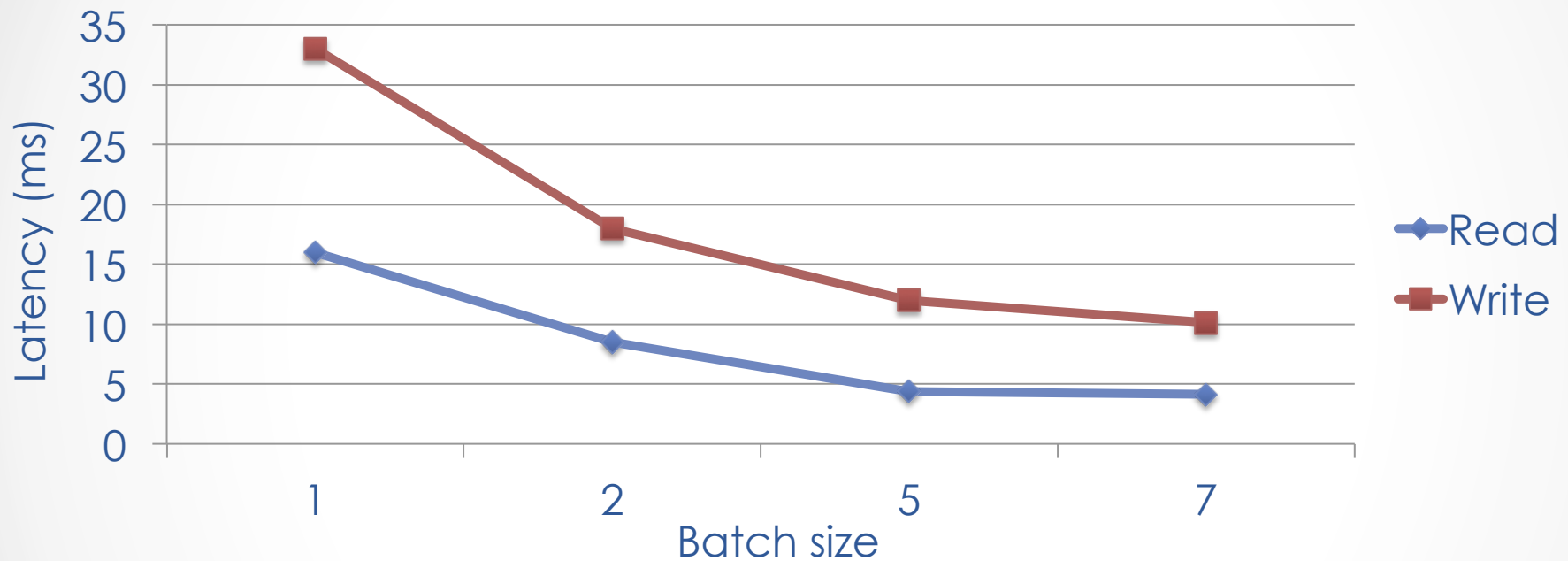
**Read-latency (ms)**



**Write-latency (ms)**



29

# VPR-ed SQLite



- Realistic trusted executions are the bottleneck
  - o 2 TCC execution at the primary (for write requests)
  - o in pessimistic runs, 1 more TCC execution at backups

# Outline

- Motivation and background

- Goals

- Architecture Design & System Operations
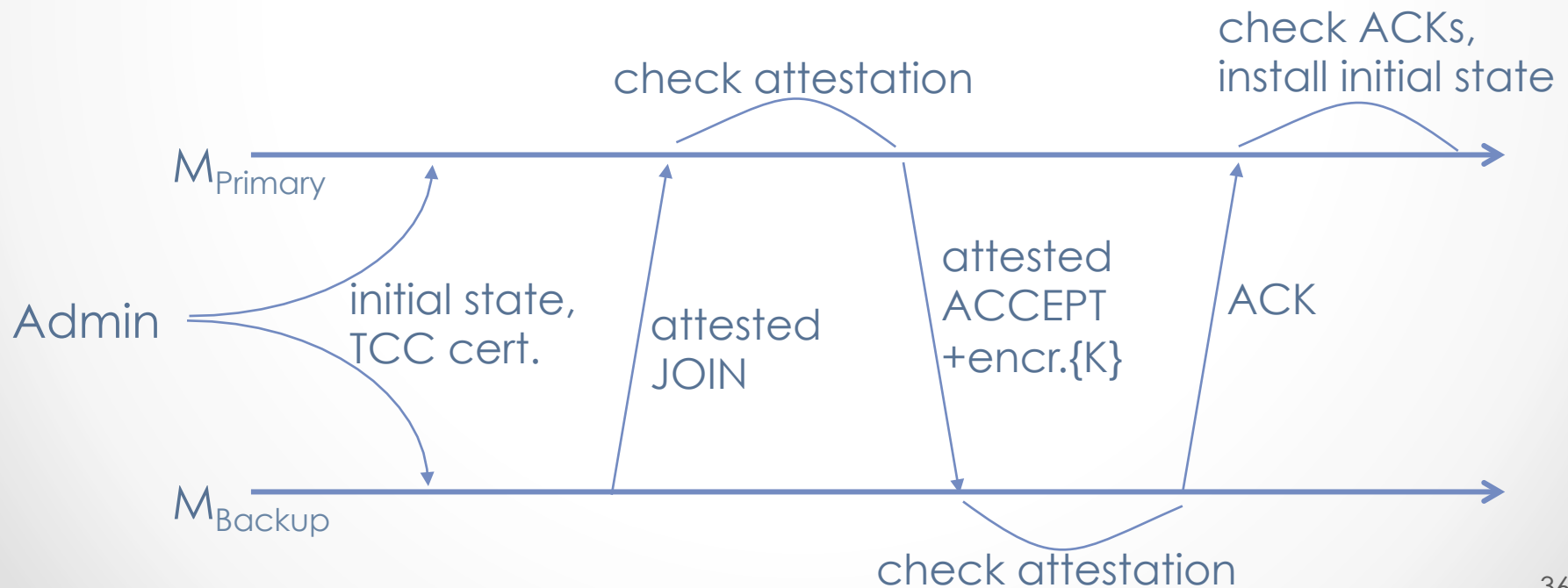
- Evaluation

- **Takeaways**

# Takeaways

- Easy to design fault-tolerant protocols
  using hardware-based security
  - o V-PR is the first fully-passive replication scheme that tolerates Byzantine failures

- No additional assumptions (compared to previous literature)

- Linear factor reduction in executing replicas
  - o Non-determinism supported by design

- Main limitation is the current technology
  - o …but it's making progress, check out Intel SGX

# Thanks.

# System Initialization

- Need to form a secure group
  - If other replicas participate, they could be later shutdown (state loss)
- Share a unique key K (use TCC secure storage for confidentiality)
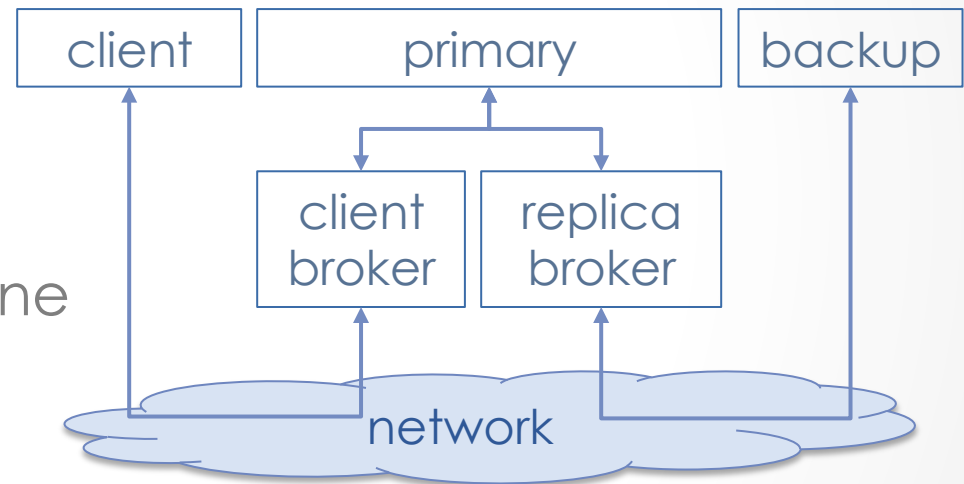- Start from same initial state



check ACKs,
install initial state

check attestation

$M_{Primary}$

Admin

initial state,
TCC cert.

attested
JOIN

attested
ACCEPT
+encr.{K}

ACK

$M_{Backup}$

check attestation

# Primary Change

- Primary identified through local view counter

  o Each replica answer to only one specific primary

- Detect primary's failure through timeouts (partial synchrony)

  o Start primary change protocol, but always answer to primary's updates

  o Exchange messages to increment view counter

  o Eventually, no progress => new primary

- Extreme cases

  o Multiple primaries: safe, because only one can make progress

  o Only one view increment:

    - replica wait for others to change primary

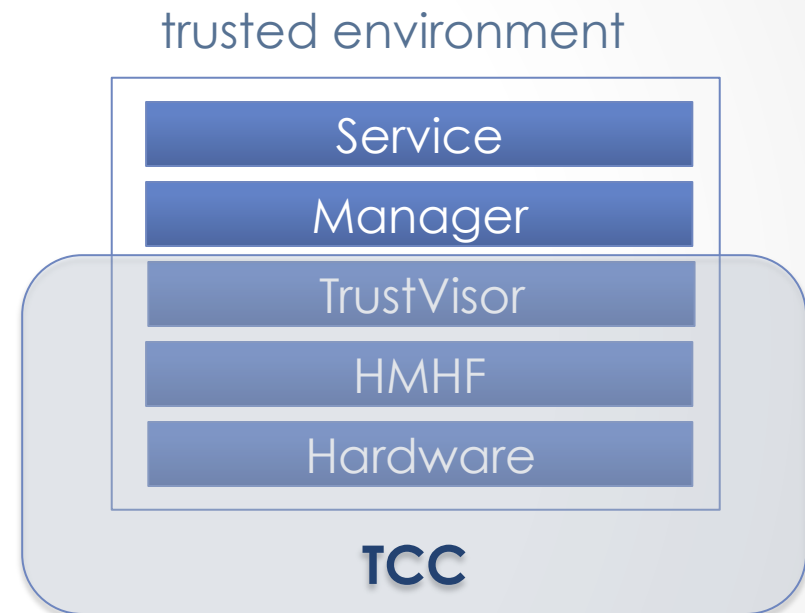    - replica can make progress through consecutive updates anyway

# Implementation

- Message passing w/ high performance library ZeroMQ

- TCC with XMHF (S&P'13) and TrustVisor(S&P'10)

- Full SQLite database engine
  - VPR-ed SQLite

| client | primary | backup |
|--------|---------|--------|

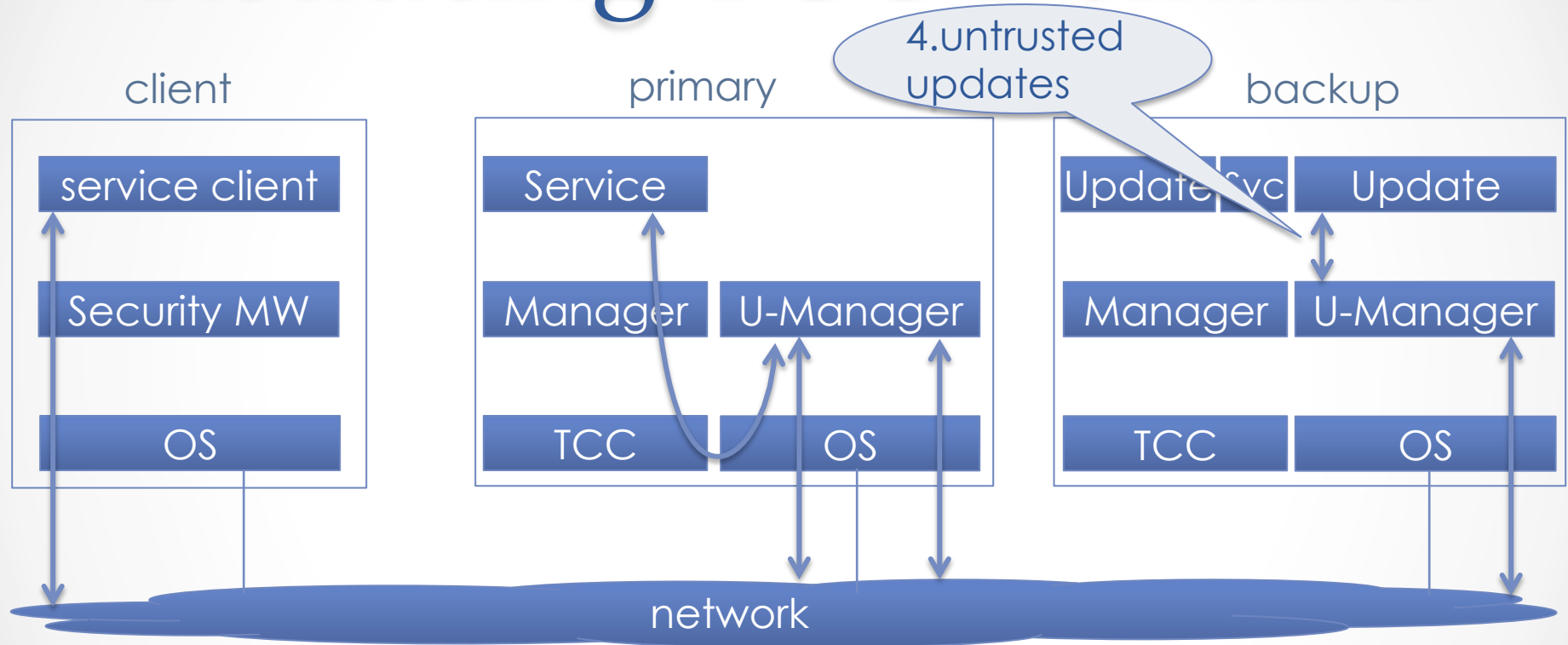| client broker | replica broker |
|---------------|----------------|

network

# Implementation

- Message passing w/ high performance library ZeroMQ

- TCC with XMHF (S&P'13) and TrustVisor(S&P'10)

- Full SQLite database engine
  - VPR-ed SQLite

trusted environment

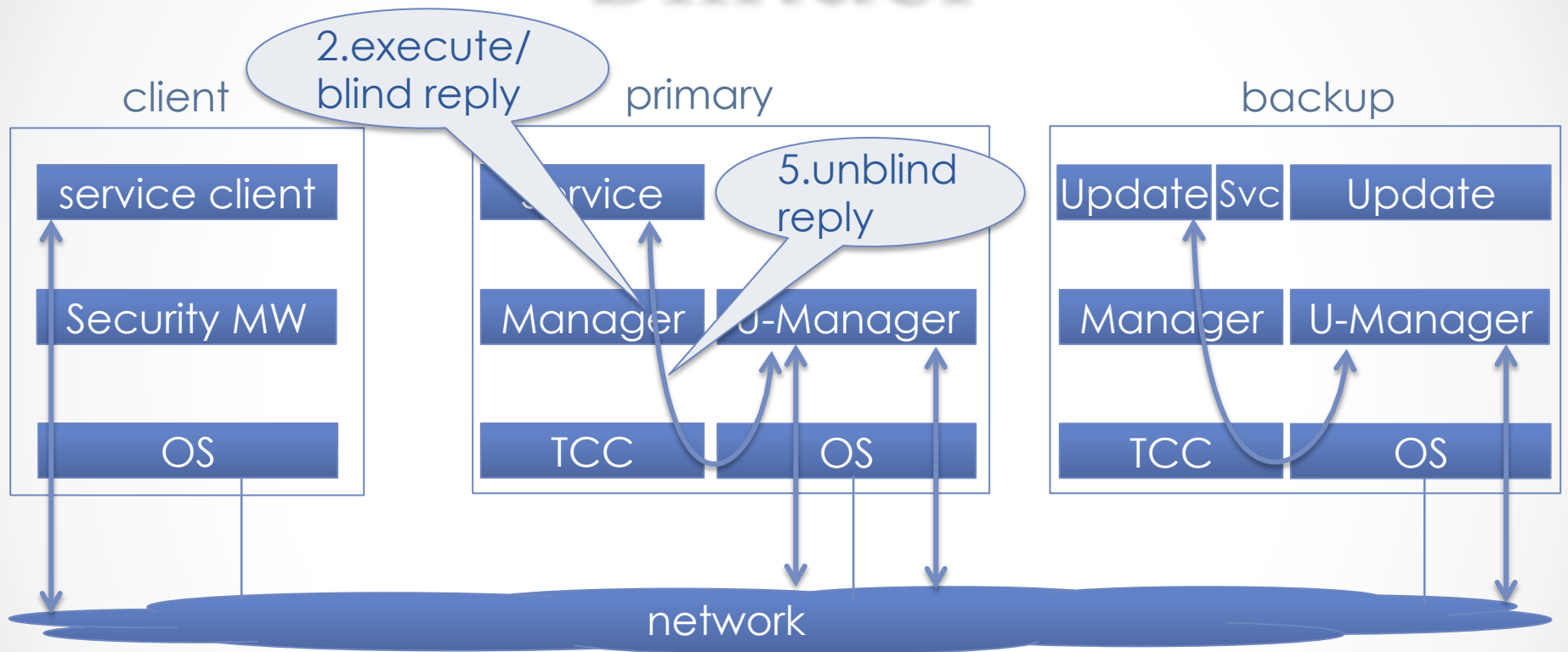| Service |
| Manager |
| TrustVisor |
| HMHF |
| Hardware |

**TCC**

- Some addressed challenges:
  - Extending the hypervisor to provide dynamic resource management and trusted counters
  - Running the service in an untrusted environment (no OS support, no access to devices, like disk): created custom APIs (memory allocation, debugging, etc.), custom filesystem (as a module, so no modification to SQLite)

# Reducing TCC Demand

client             primary           backup

4.untrusted updates

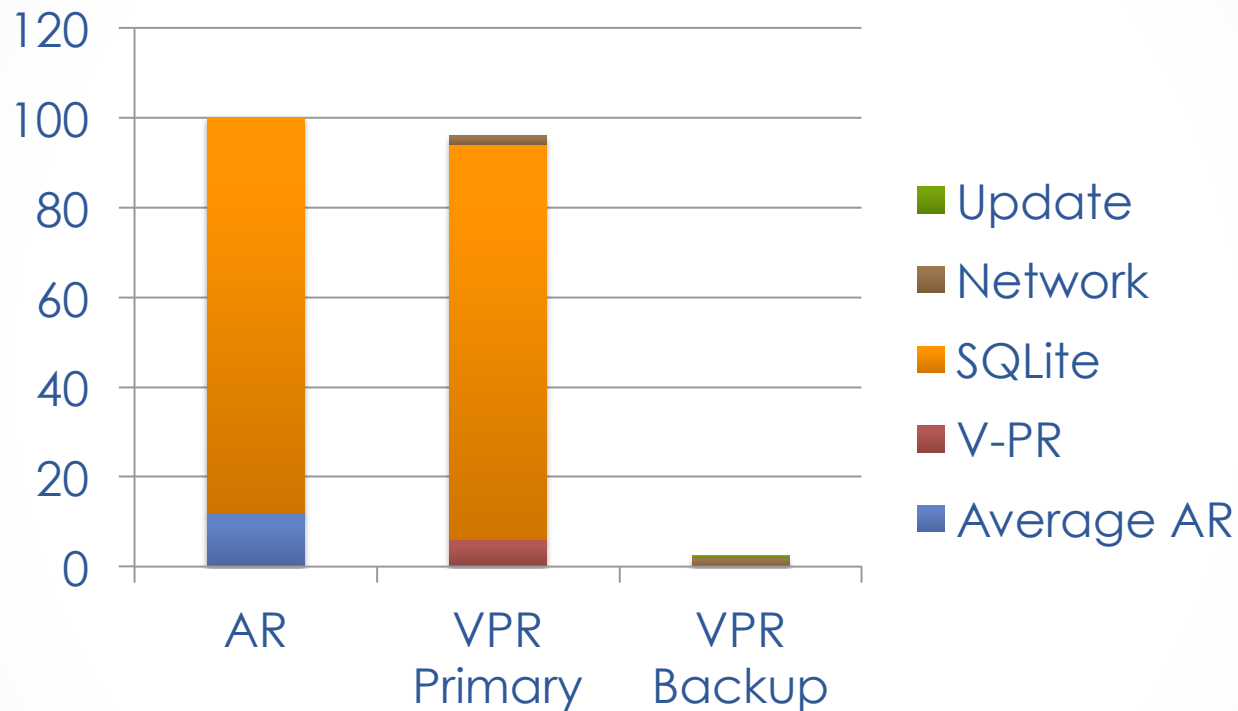| client | primary | backup |
|---|---|---|
| service client | Service | Update Svc   Update |
| Security MW | Manager   U-Manager | Manager   U-Manager |
| OS | TCC   OS | TCC   OS |

network

- Speculative update: validate it and send ACK
- No TCC execution => 1 active TCC and rest are passive
- Backup ACKs required: 2f+1
  (yes, all of them, so at least a correct one always available)

# Blinder



- Reply's authenticator is blinded during update
- U-Manager cannot send it back to client and break consistency
- Reply is unblinded after ACKs are validated

# Code size



Legend:
- Update
- Network
- SQLite
- V-PR
- Average AR

Categories: AR, VPR Primary, VPR Backup

- Actively used code in fault-free scenario
  - KSLoC=thousand lines of source code

- VPR Backup's code is independent from the implemented service
  - Measurement of service code is not included

42