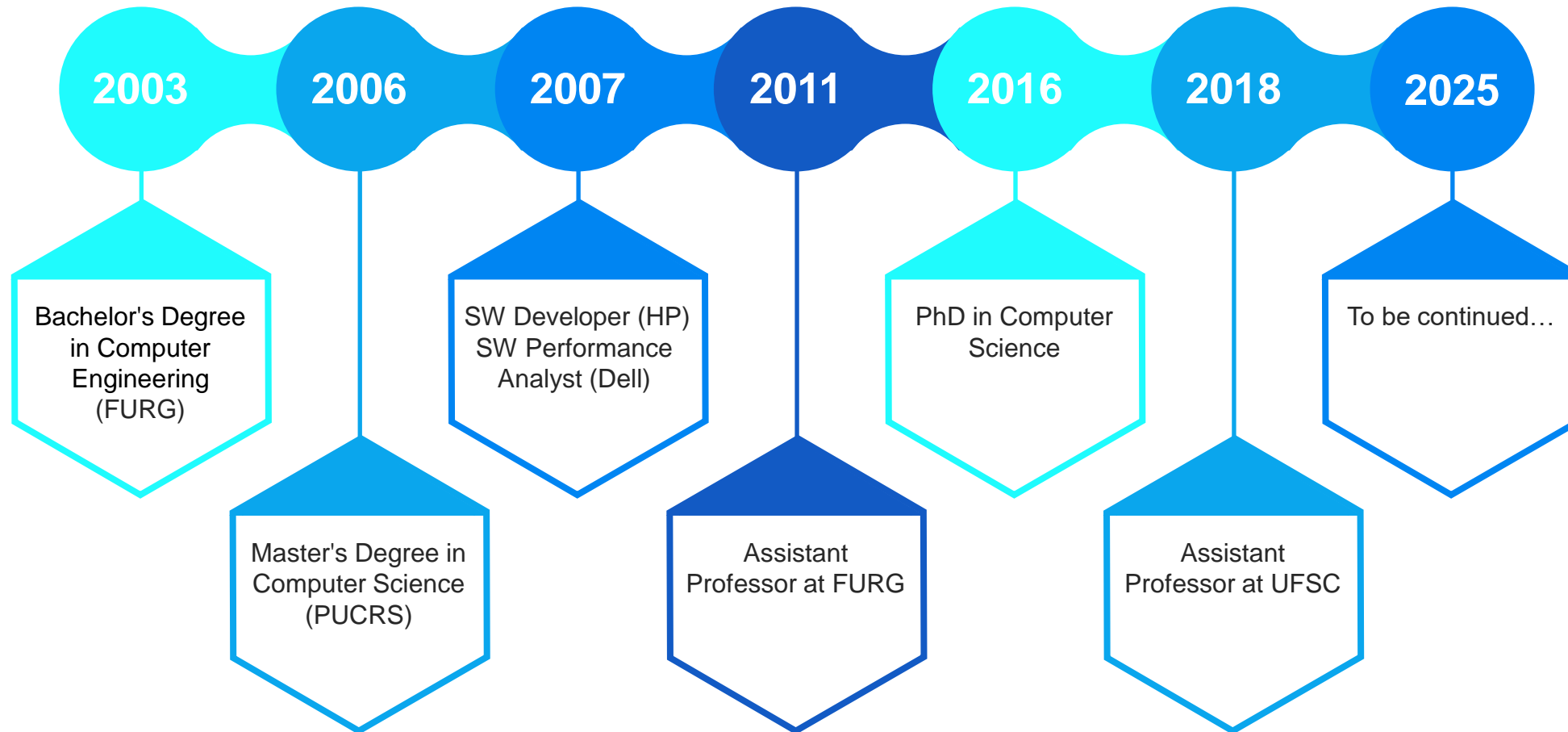# Composing State Machine Replicas

Odorico Machado Mendizabal
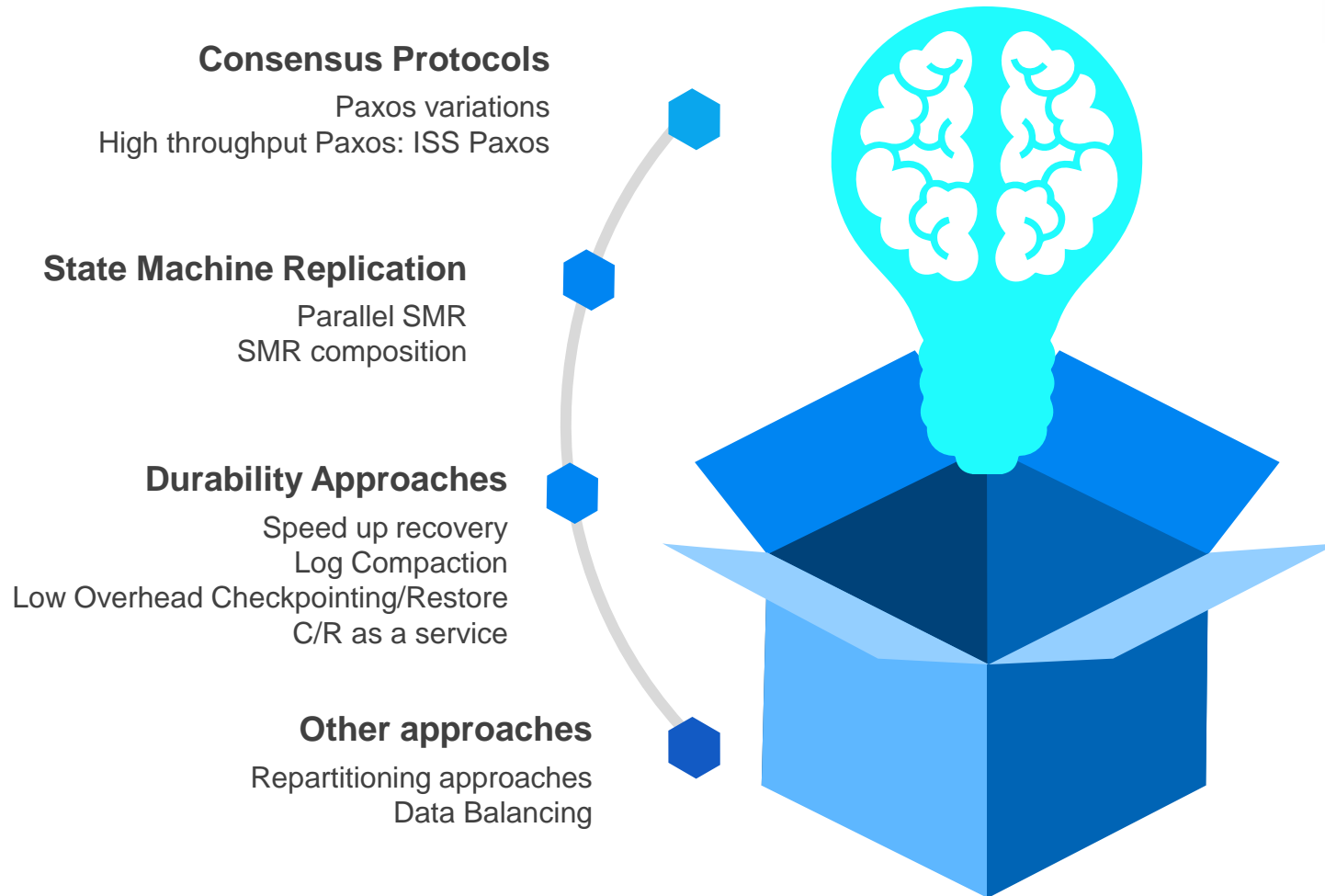
Workshop Suíça–Brasil: Um Olhar Atual sobre Sistemas Distribuídos:
Da Pesquisa à Aplicação no Mundo Real

15 e 16 de abril de 2025

# Breaking the Ice



**2003** — Bachelor's Degree in Computer Engineering (FURG)

**2006** — Master's Degree in Computer Science (PUCRS)

**2007** — SW Developer (HP) SW Performance Analyst (Dell)

**2011** — Assistant Professor at FURG

**2016** — PhD in Computer Science

**2018** — Assistant Professor at UFSC

**2025** — To be continued…

# Research Interest and Recent Results



**Consensus Protocols**

Paxos variations
High throughput Paxos: ISS Paxos

**State Machine Replication**

Parallel SMR
SMR composition

**Durability Approaches**

Speed up recovery
Log Compaction
Low Overhead Checkpointing/Restore
C/R as a service

**Other approaches**

Repartitioning approaches
Data Balancing

# State Machine Replication
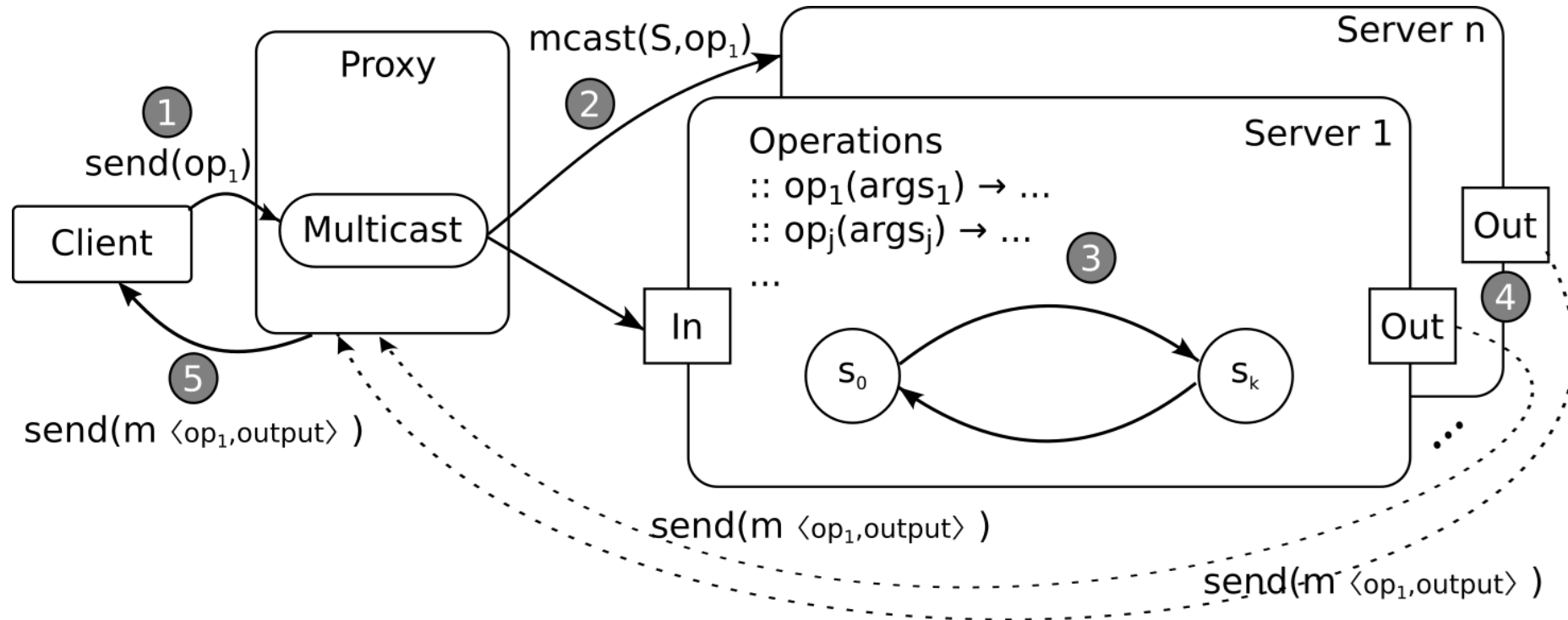
A set of servers behaves as replicated state machines
- Replicas start in the same initial state
- Service operations are deterministic

Clients issue commands to every replica through a consensus or atomic broadcast protocol:
- correct replicas receive every command
- if a replica processes a command $c_1$ before $c_2$, then no replica process $c_2$ before $c_1$

# State Machine Replication

# State Machine Replication – Applications

**Google Chubby**

Lock service, used to help coordination in distributed environment using locking semantics
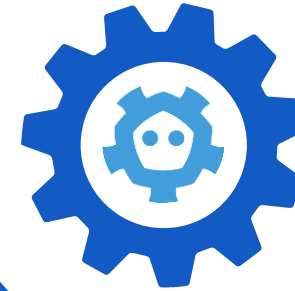
**Key-value store services**

Ex. etcd, used to replicate metadata on Kubernetes cluster managers

**Apache Zookeeper**

Open source server for highly reliable distributed coordination

**Google File System**

Distributed file system that provides efficient and reliable data access

**Online services**

# This talk

> Over the past decades, SMR has gained popularity, leading to extensive research that has enhanced its resilience, performance, and scalability. However, one aspect not yet addressed in SMR is service composition

## Composing State Machine Replication

Caroline Martins Alves [ Universidade Federal de Santa Catarina | caroline.martins@posgrad.ufsc.br ]
Matheus Antonio de Souza [ Universidade Federal de Santa Catarina | matheus.souza.m.a.s@grad.ufsc.br ]
Thais Bardini Idalino [ Universidade Federal de Santa Catarina | thais.bardini@ufsc.br ]
Odorico Machado Mendizabal [ Universidade Federal de Santa Cata-
rina | odorico.mendizabal@ufsc.br ]

Department of Computer Science and Statistics , Universidade Federal de Santa Catarina, R. Delfino Conti, s/n,
Trindade, Florianópolis, SC, 88040-900, Brazil.

**Abstract** High availability is a fundamental requirement in large-scale distributed systems, where replication strate-
gies are central in keeping applications operational despite a bounded number of failures. State Machine Replication
(SMR) is one of the most widely adopted approaches for implementing highly available, fault-tolerant services, as
increases uptime while ensuring strong consistency. In recent years, research on SMR
variations tailored to enhance resilience, performance, and scalability
perspective by introducing Composi...
...composit...

# SMR Research

A (very) quick view of approaches over the years

# SMR – other approaches

**Byzantine faults**

Adding protocols capable of tolerating arbitrary faults

**Recovery and Reconfiguration**

Aimed to improve resilience

**Parallel SMR**

Independent requests executed in parallel

# Composing State Machine Replication (CSMR)

Build services by combining separate instances of SMR

- Microservice-based systems

- Flexible and loosely coupled solutions

- Development in a modular way

SMR 1

SMR 2
.
.
.

SMR n

# CSMR – Selected use cases

### SMR 1 – Lock service

```
boolean acquire (string key)
boolean release (string key)
```

### SMR 2 – key-value store

```
string get(string key)
void put(string key, string value)
```
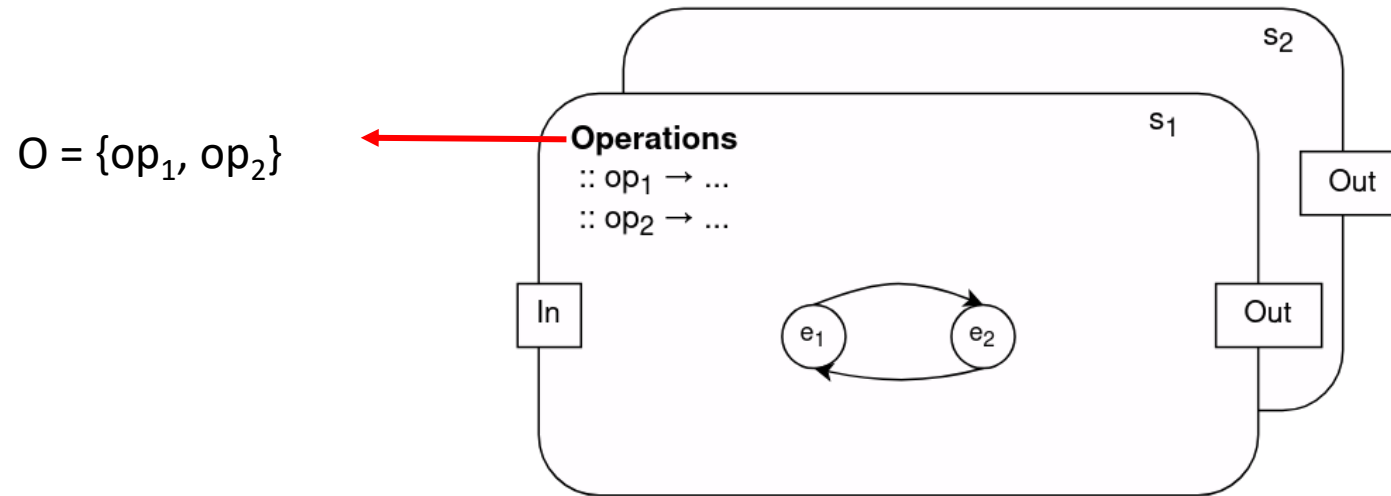
### SMR 3 – Logging service

```
void append(Object entry)
Object[] retrieve(int first, int last)
boolean truncate(int index)
```
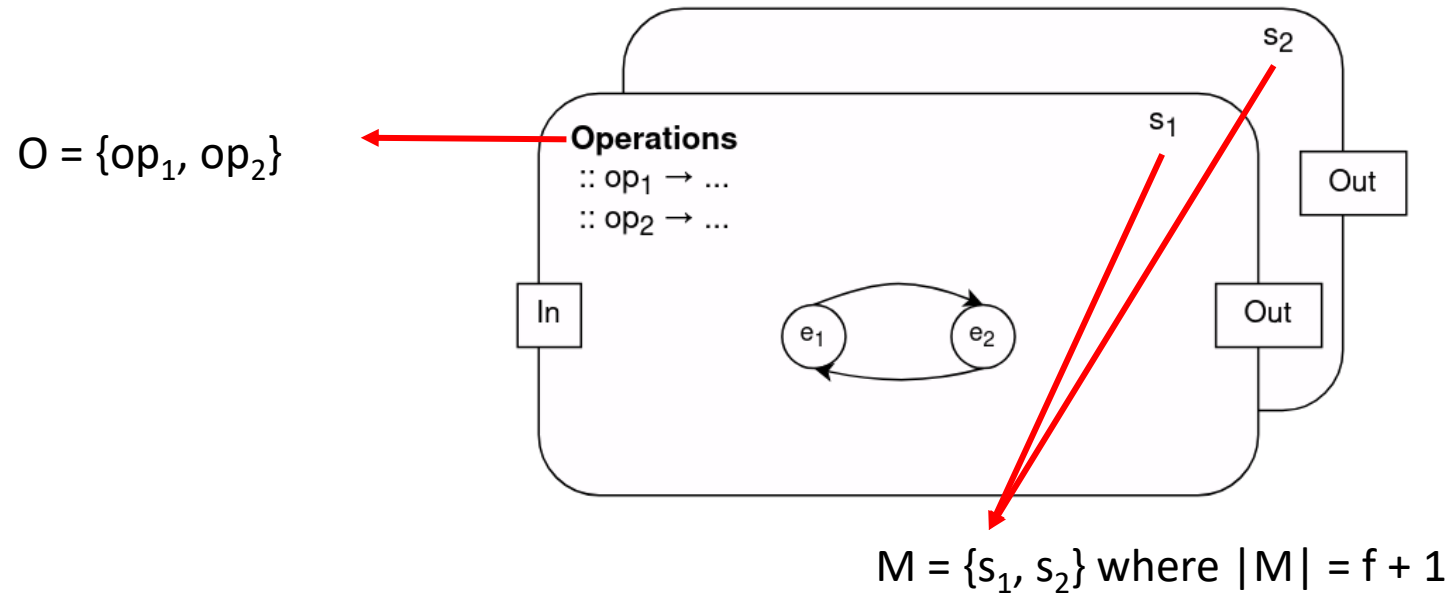
# SMR Formalization: Definitions

- Replicated service

- Operations (with arguments)
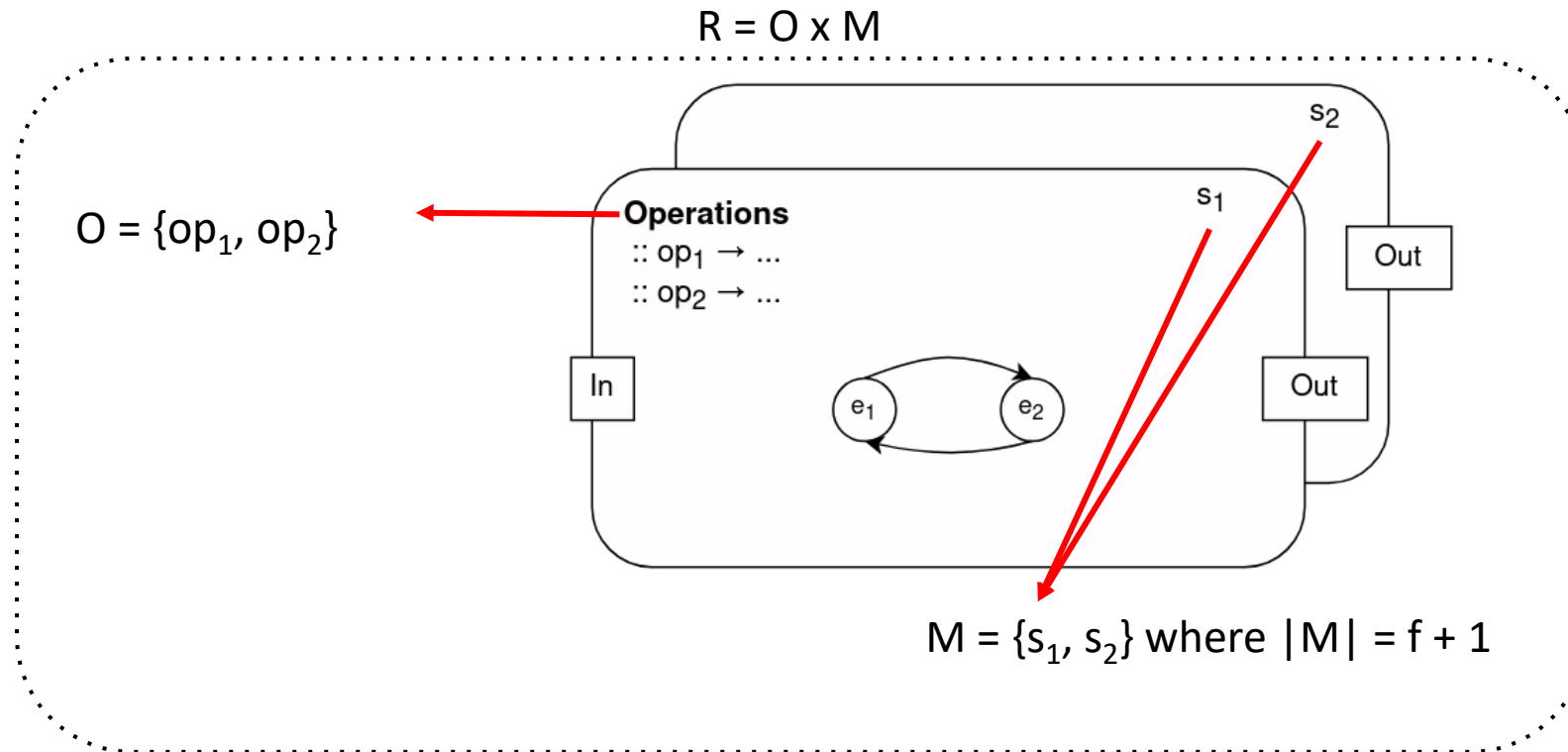
- Execution

- Output

# SMR Formalization: Definitions

$O = \{op_1, op_2\}$

# SMR Formalization: Definitions



$O = \{op_1, op_2\}$

$M = \{s_1, s_2\}$ where $|M| = f + 1$

# SMR Formalization: Definitions

$$R = O \times M$$

$$O = \{op_1, op_2\}$$

**Operations**
:: $op_1 \rightarrow \ldots$
:: $op_2 \rightarrow \ldots$

In

$e_1$ $e_2$

$s_2$

$s_1$

Out

Out

$$M = \{s_1, s_2\} \text{ where } |M| = f + 1$$

# SMR Formalization: Definitions

- **Definition 1** - Replicated service
  - **Example: SMR 2 - Key-value store**

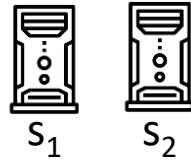**SMR 2 - Key-value store**



O = {get, put}

# SMR Formalization: Definitions

- **Definition 1** - Replicated service
  - **Example: SMR 2 - Key-value store**

**SMR 2 - Key-value store**

$O$ = {get, put}

**Replicas**

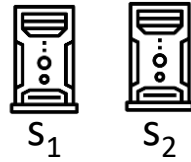$s_1$ $s_2$

$M$ = {$s_1$, $s_2$}

# SMR Formalization: Definitions

- **Definition 1** - Replicated service
  - **Example: SMR 2 - Key-value store**

**SMR 2 - Key-value store**



$O = \{get, put\}$

**Replicas**



$s_1 \quad s_2$
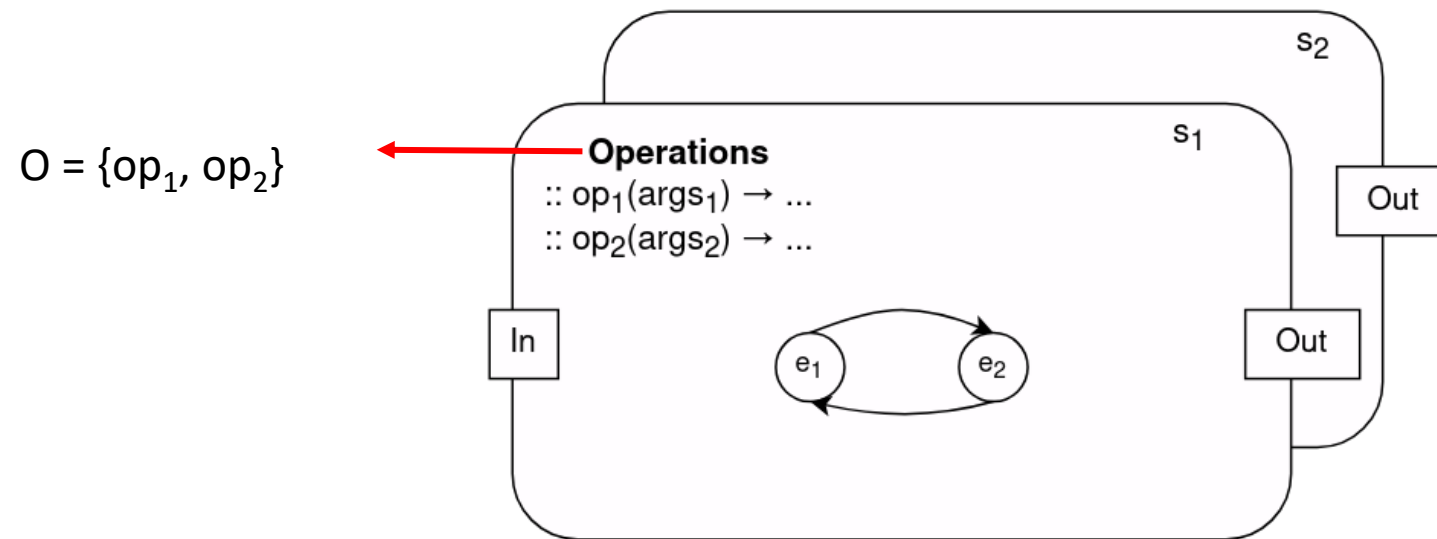
$M = \{s_1, s_2\}$

**Replicated Service**
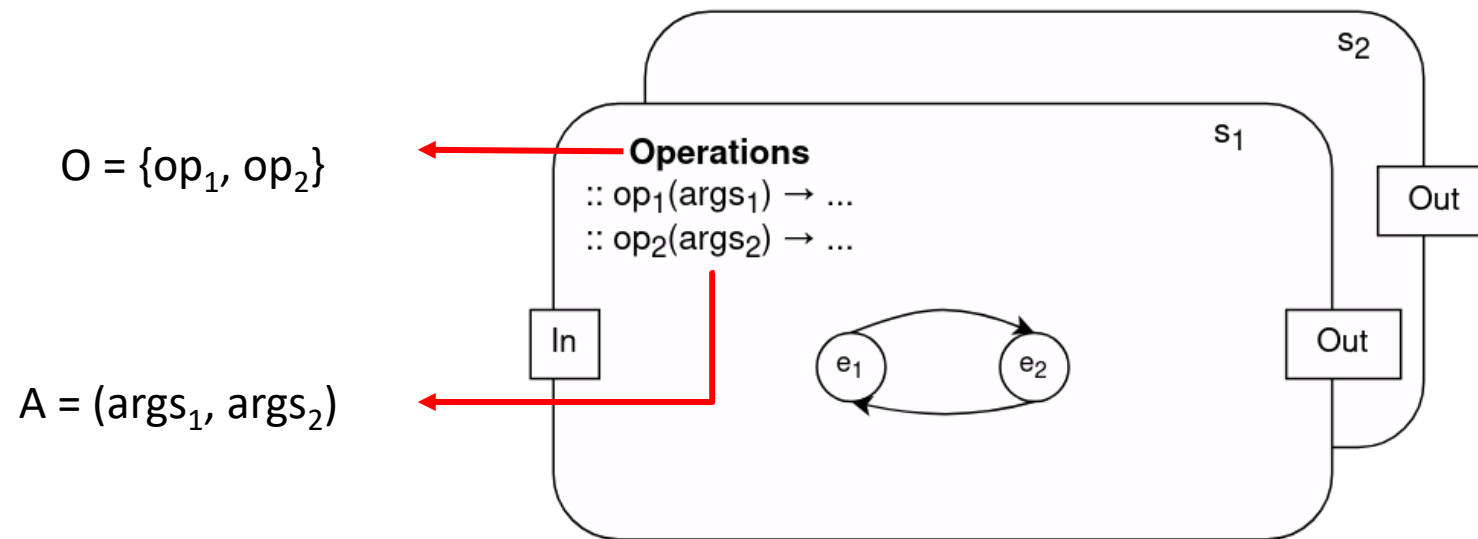
$R = \{(get, s_1), (get, s_2), (put, s_1), (put, s_2)\}$

# SMR Formalization: Definitions

- **Definition 2** - Operations with arguments

$O = \{op_1, op_2\}$
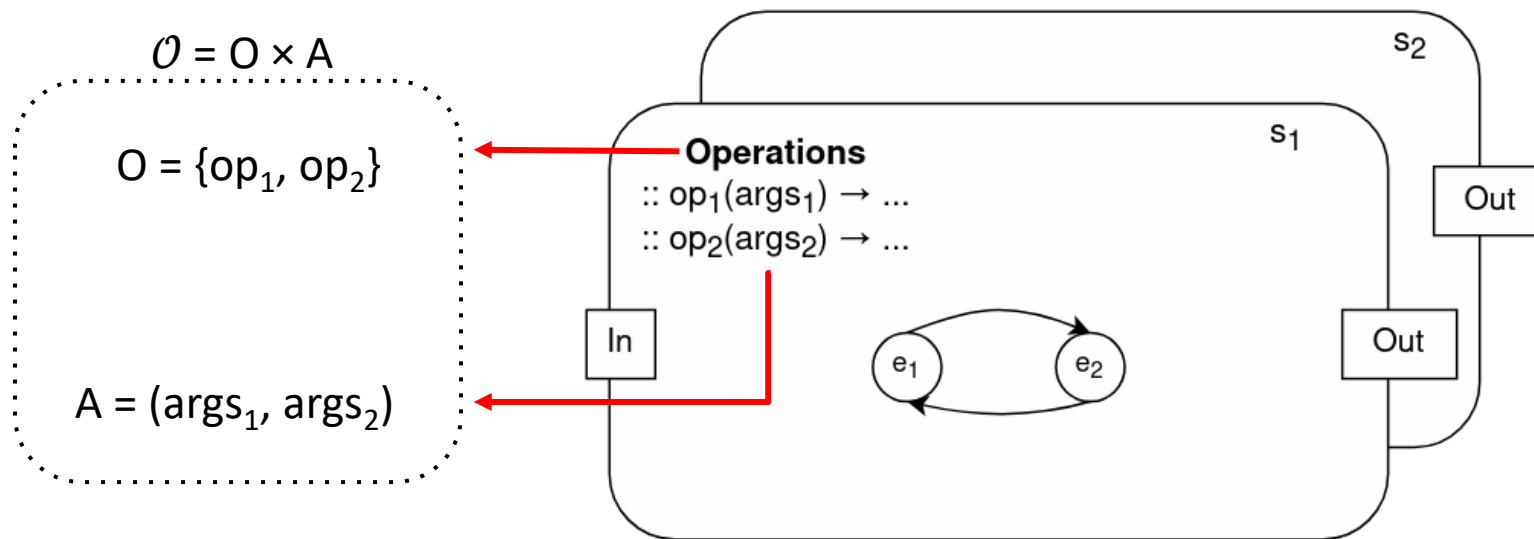
# SMR Formalization: Definitions

- **Definition 2** - Operations with arguments



$O = \{op_1, op_2\}$

$A = (args_1, args_2)$

# SMR Formalization: Definitions

- **Definition 2** - Operations with arguments

$$\mathcal{O} = O \times A$$

$$O = \{op_1, op_2\}$$

$$A = (args_1, args_2)$$

# SMR Formalization: Definitions

- **Definition 2** - Operations with arguments
  - **Example: SMR 2 - Key-value store**

**SMR 2 - Key-value store**

O = {get, put}

# SMR Formalization: Definitions

- **Definition 2** - Operations with arguments
  - **Example: SMR 2 - Key-value store**

**SMR 2 - Key-value store**



$O = \{get, put\}$

**Arguments**

$A = \{a, b, ..., z\}$

A* set of all strings over A

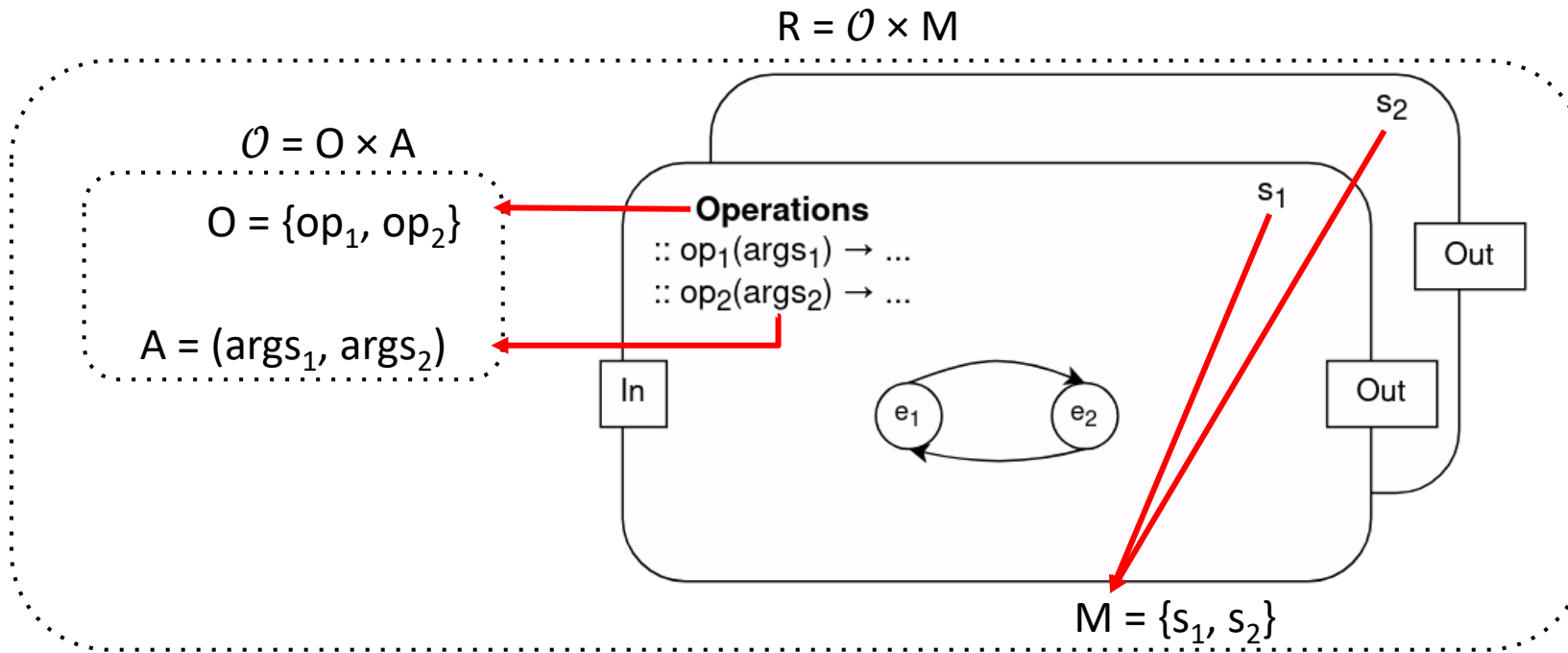**Operations with arguments**

$\mathcal{O}_{get} = \{get\} \times A*$

(get, "firstkey"); (get, "secondkey")

# SMR Formalization: Definitions

- **Definition 3** - Replicated service with arguments



$$R = \mathcal{O} \times M$$

$$\mathcal{O} = O \times A$$

$$O = \{op_1, op_2\}$$

$$A = (args_1, args_2)$$

**Operations**
:: $op_1(args_1) \rightarrow \ldots$
:: $op_2(args_2) \rightarrow \ldots$

In

Out

Out

$e_1$   $e_2$

$s_1$

$s_2$

$$M = \{s_1, s_2\}$$

# SMR Formalization: Definitions

- **Definition 4** – Execution
- **Definition 5** – Output

# SMR Formalization: Example

**SMR 2 - Key-value store**



$O = \{get, put\}$

**Arguments**

$A = \{a, b, ..., z\}$

$A*$ set of all strings over $A$

**Operations with arguments**

$\mathcal{O}_{get} = \{get\} \times A*$

(get, "firstkey"); (get, "secondkey")

**Execution and output**

$E \subseteq \mathcal{O} \times U$

$E \subseteq (\{get\} \times A*) \times A*$

(get, "firstkey") $\times$ ("firstvalue")

# Composing SMR

A new perspective

# Composing State Machine Replication (CSMR)

**Literature**
Composition can be a powerful strategy

**Existing SMR**
Add new features combining different SMRs

**Extending a service**
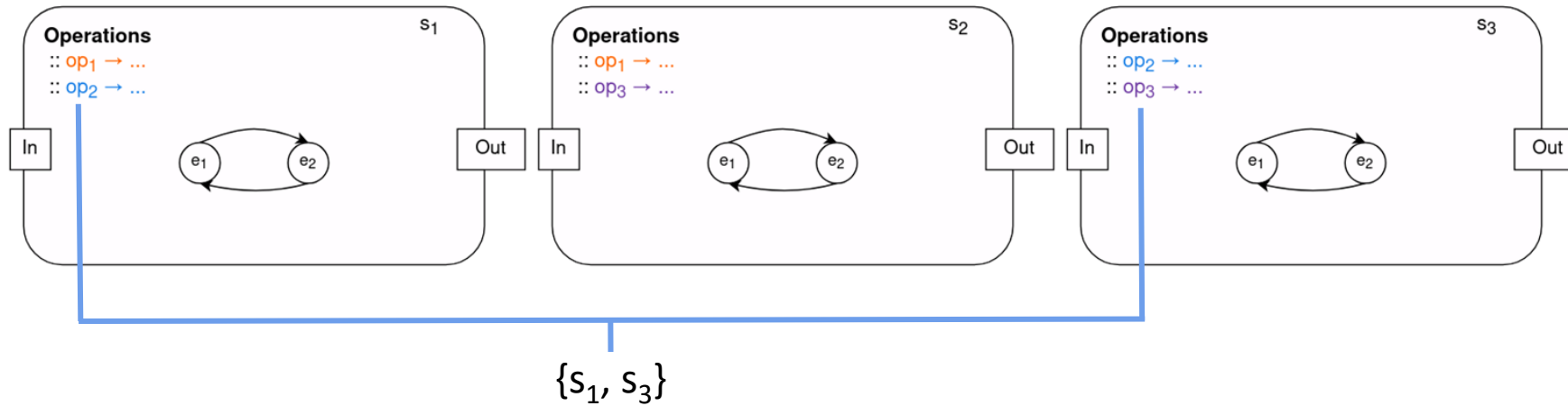Incorporating additional functionalities

$\neq$ **Different operations**
Now there is no longer any requirement for all replicas to perform the same operations
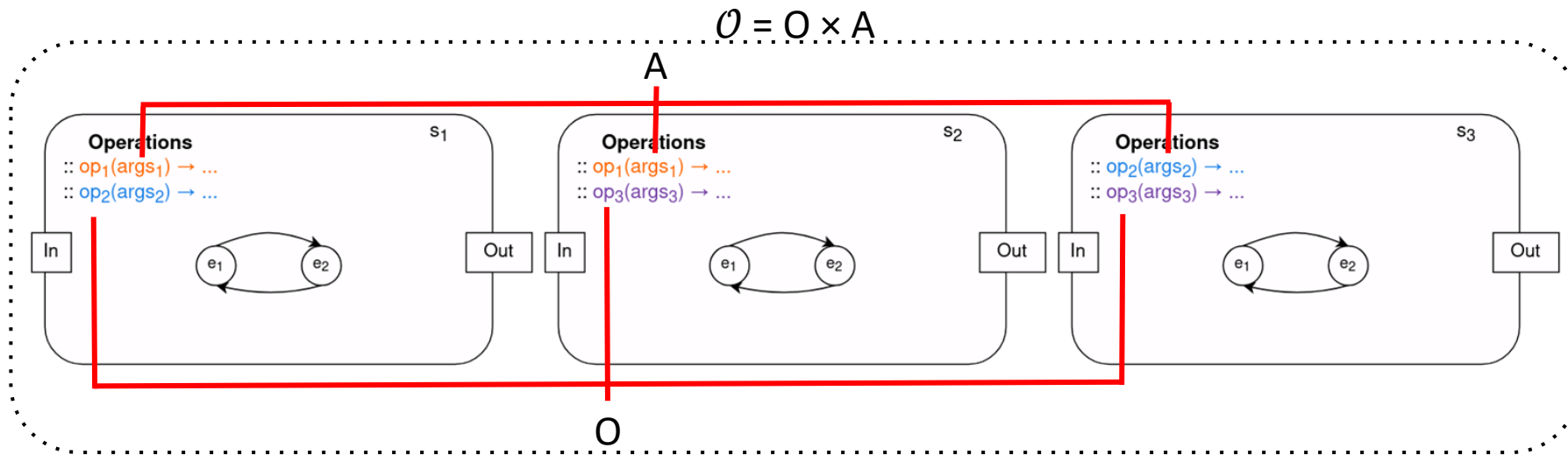
# CSMR: Definitions

- **Definition 6** - Replication set
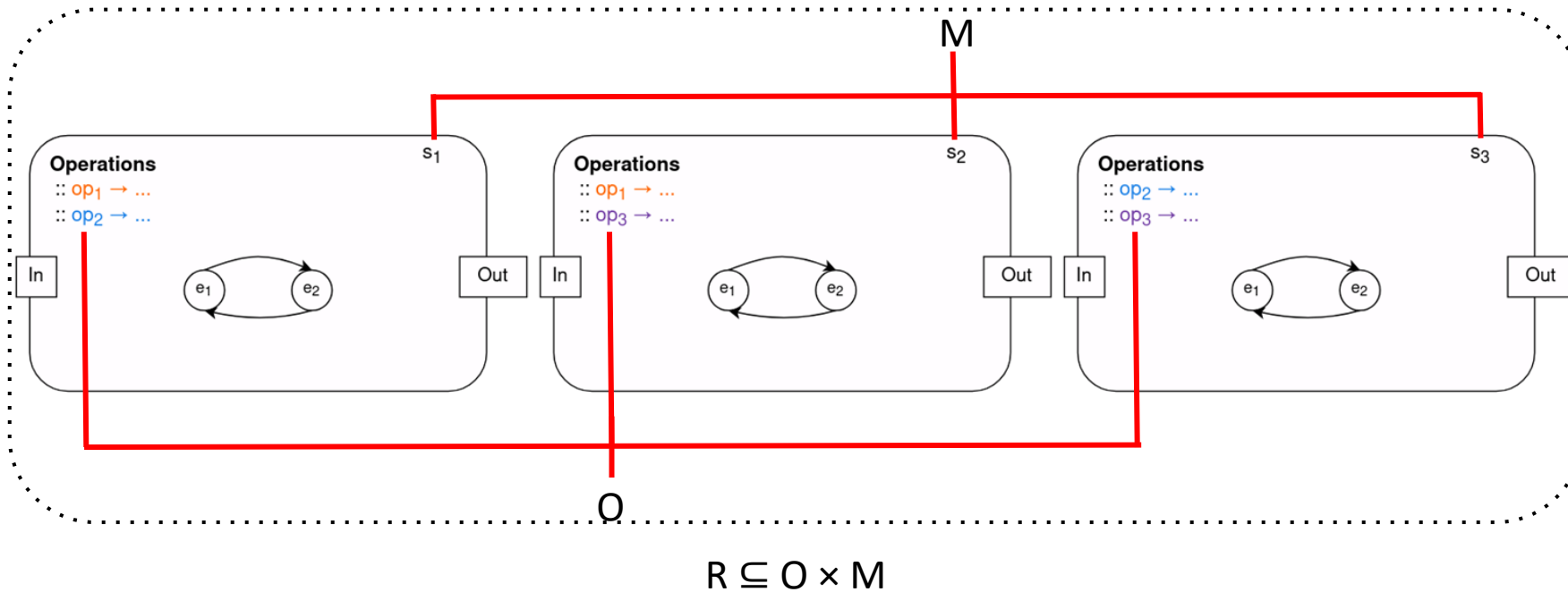
$$R(x) = \{s_i \in M : (x, s_i) \in R\}$$



$\{s_1, s_3\}$

# CSMR: Definitions

- **Definition 7** - Replication set with arguments



$$\mathcal{O} = O \times A$$

- $R(x) = \{s_i \in M : ((x, a), s_i) \in R, \text{ with } a \in A\}$

# CSMR: Definitions
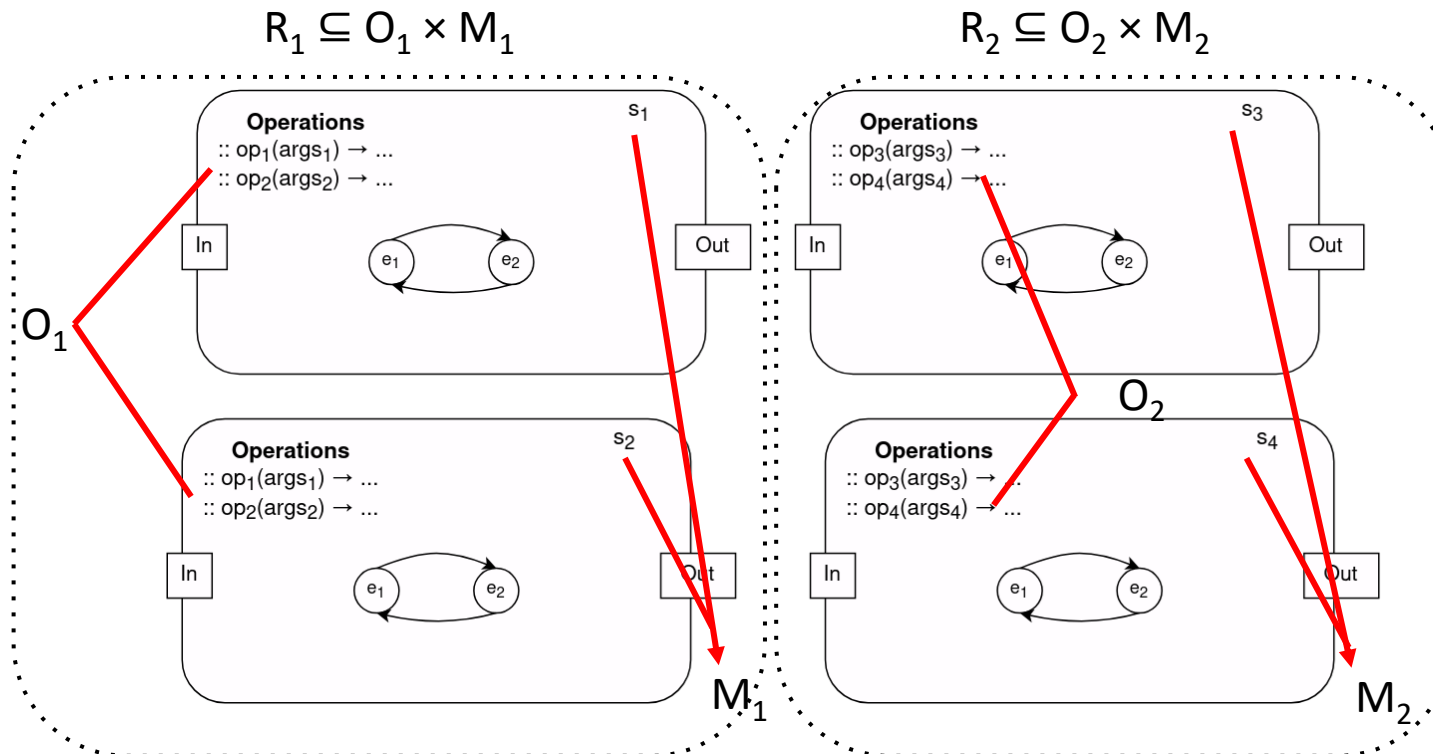
- **Definition 8** - Composable replicated service



$$R \subseteq O \times M$$

# CSMR: Definitions
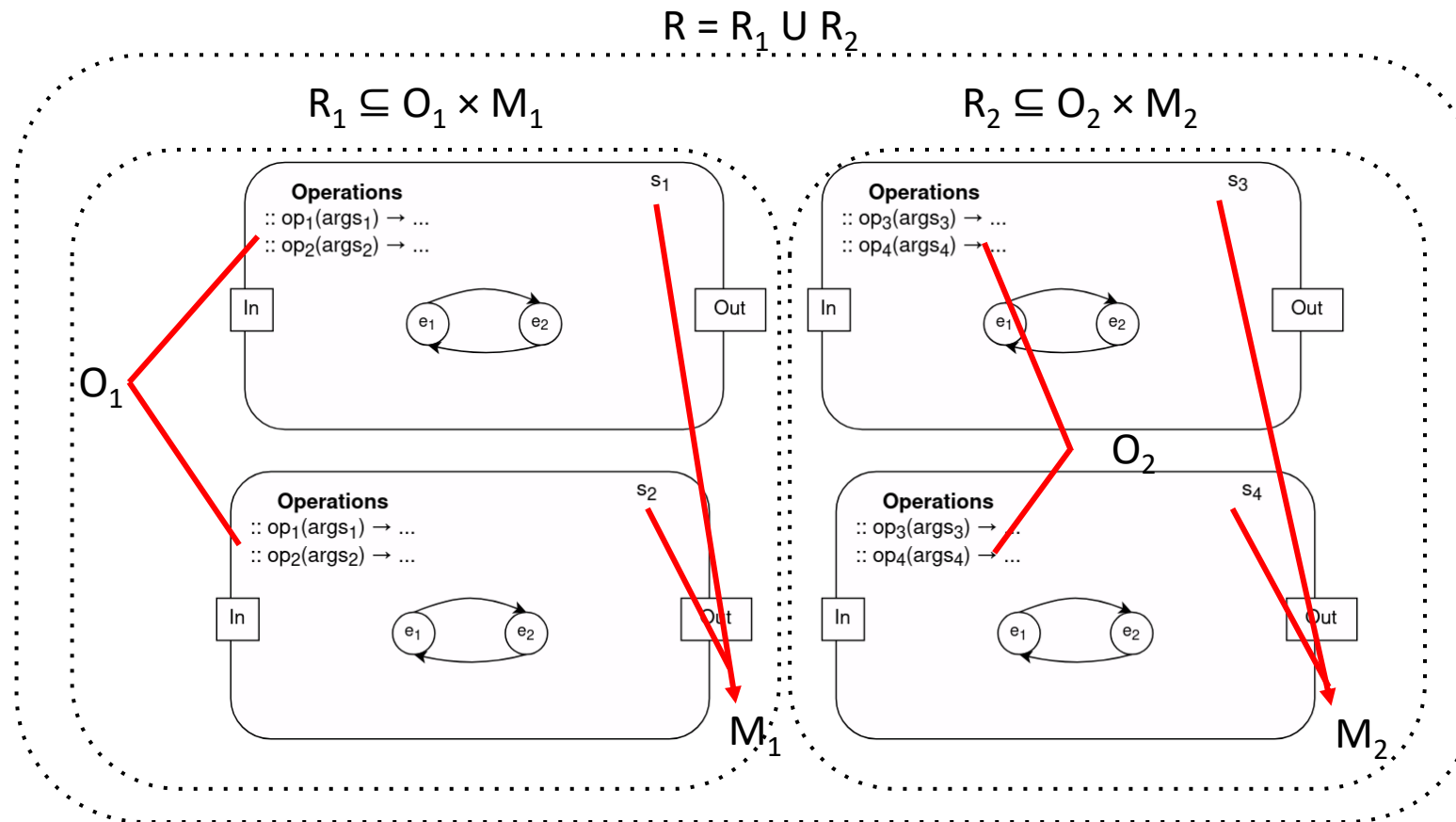
- **Definition 9** - Composition

# CSMR: Definitions

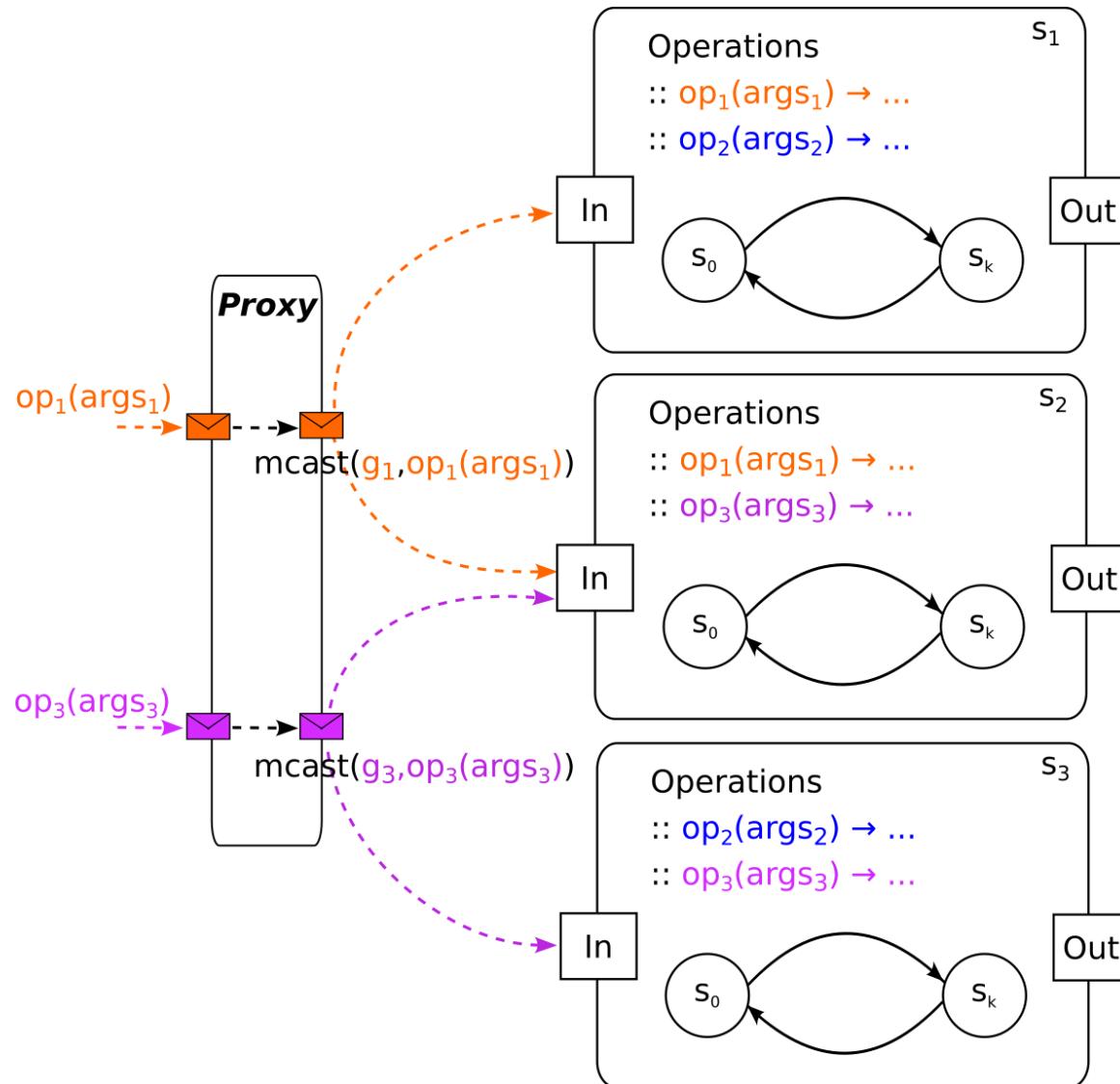- **Definition 9** - Composition

# CSMR: Definitions

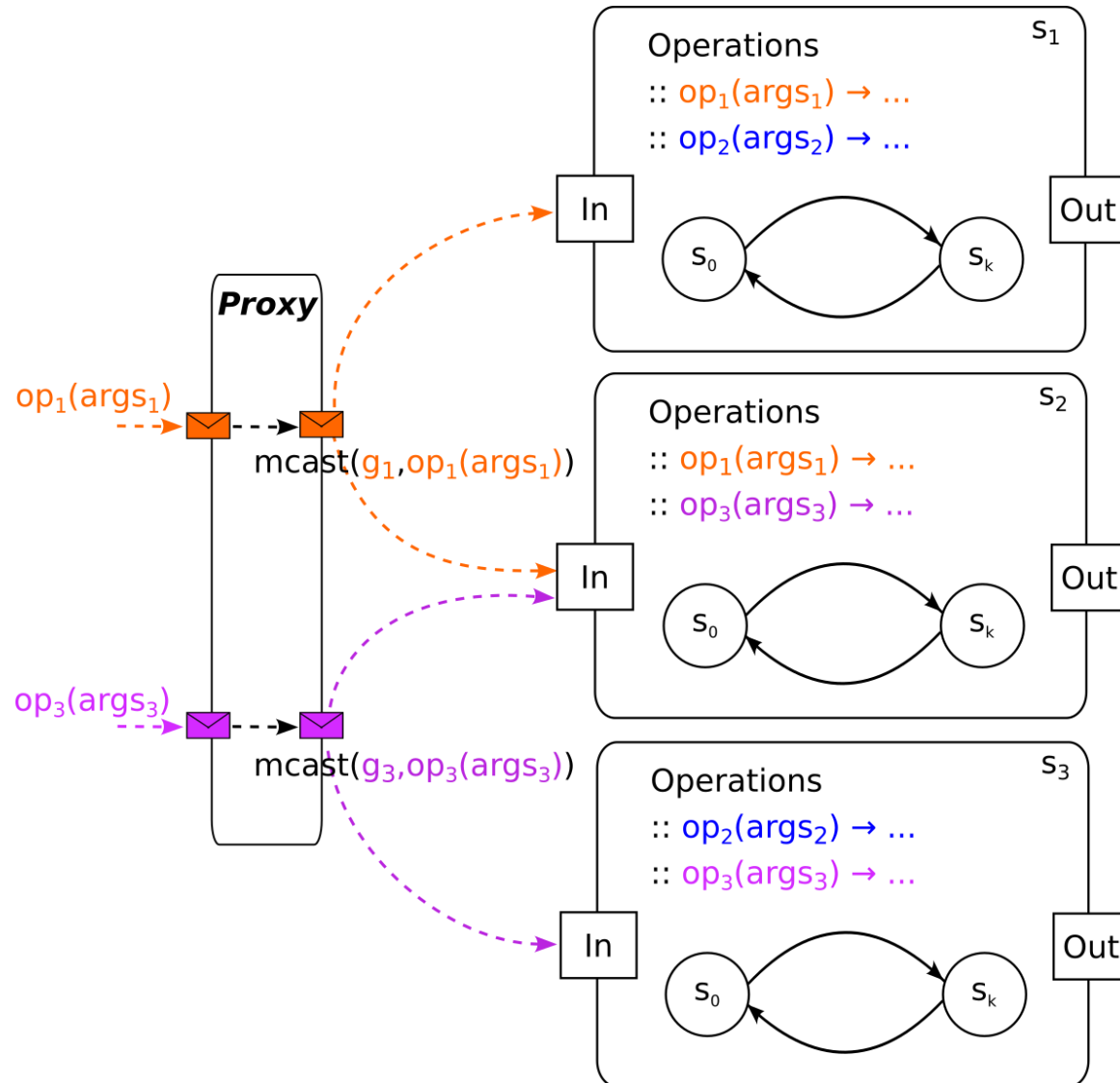- **Definition 9** - Composition

# Composing strategies

- Adding SMR operations

- Extending SMR operations' execution

- Argument partition

# Composing strategies – Adding operations

# Composing strategies – Adding operations

# Composing strategies – Extending operations execution
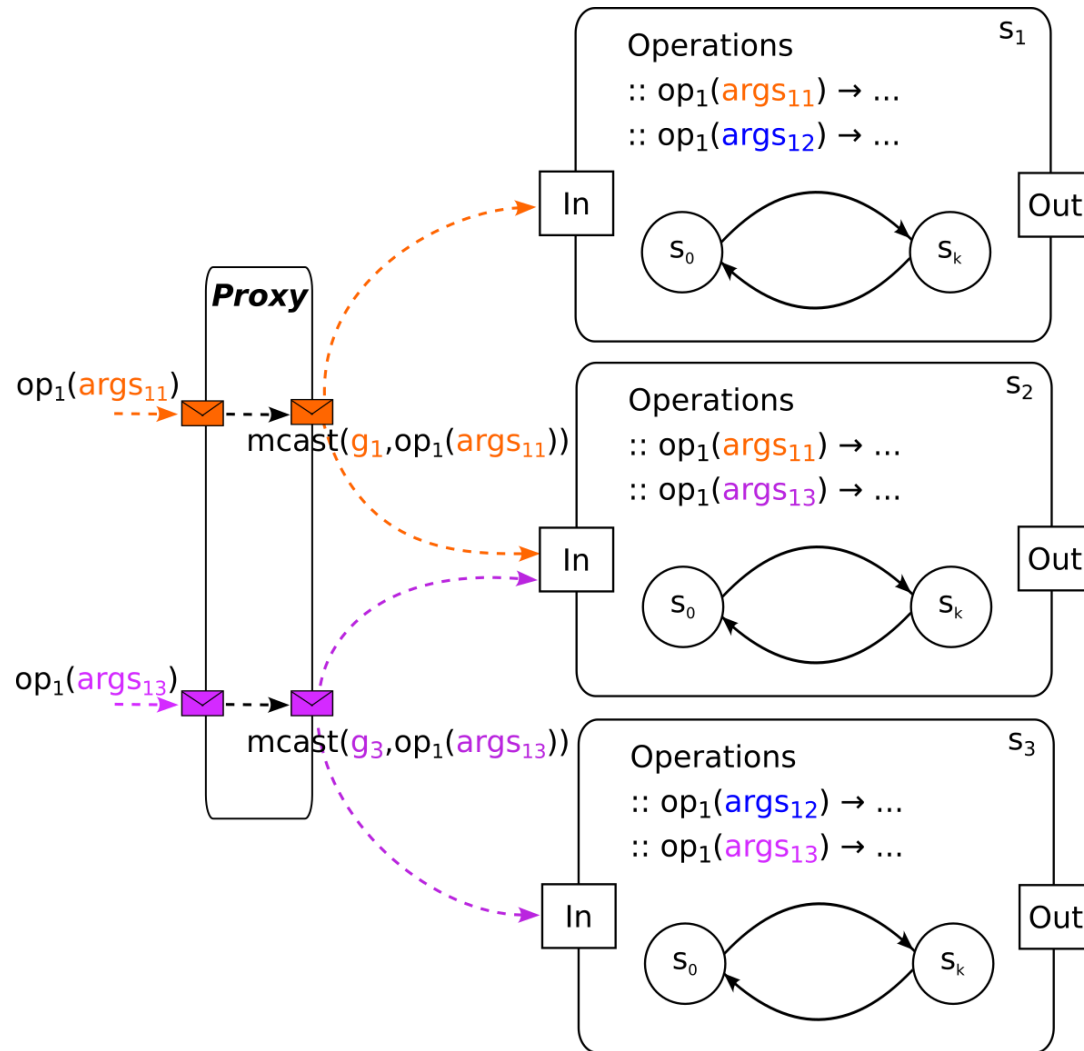
# Composing strategies – Extending operations execution

- Example: Key-value store with logging

# Composing strategies – Argument partition



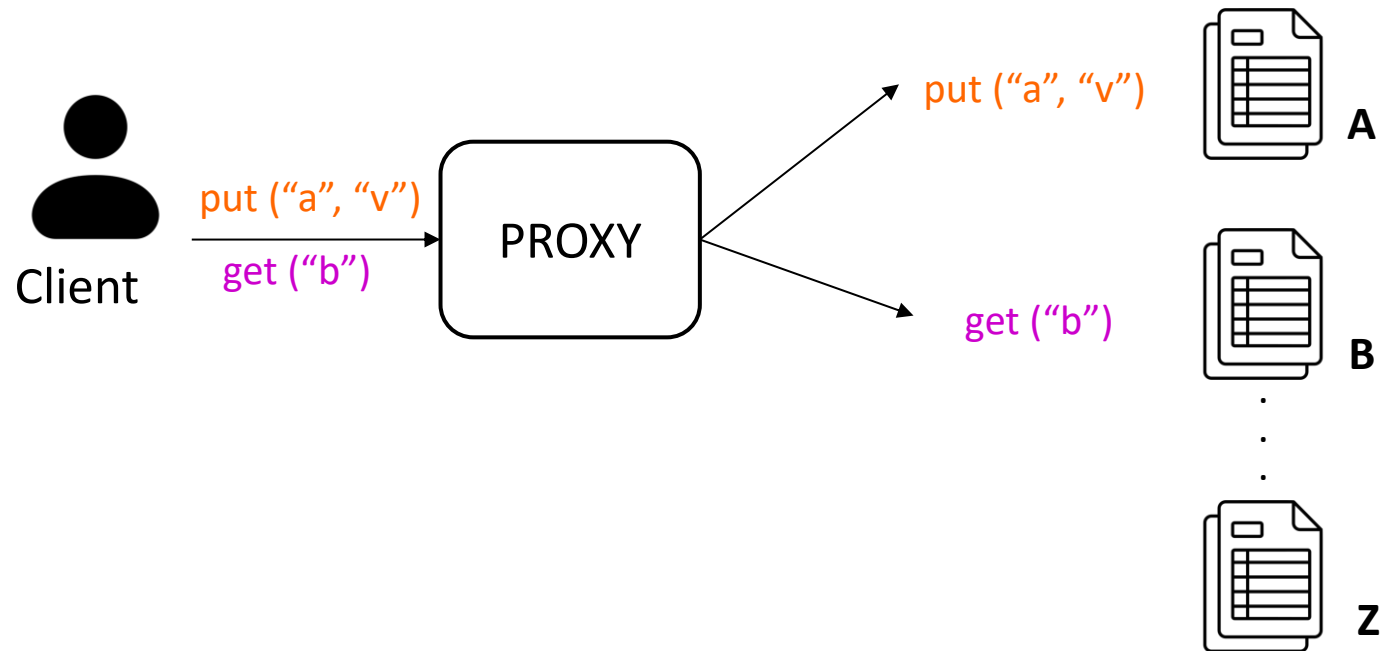SMR 2 – Key-value store

$O = \{get, put\}, A = \{A_a, A_a \times A\}$

$+$

SMR 2 – key-value store

$O = \{get, put\}, A = \{A_b, A_b \times A\}$

# Composing strategies – Argument partition

# Final remarks

● CSMR enables the construction of more complex applications

  ○ Still, preserving fault tolerance

● Modular approach encourages development of loosely coupled architectures

  ○ Fits well to microservices and cloud applications

● Some of the composing strategies resemble previous contributions in the literature

  ○ Bezerra et al. Scalable state-machine replication (DSN, 2014)

  ○ Xavier et al. Scalable and decoupled logging for state machine replication (SBRC, 2021)

# Future work

- Define an RPC API for client requests invocation

- Define a declarative configuration for CSMR (YAML file)

- Implement the Proxy

  ○ Many challenges

- Propose new use cases

  ● Do you have any ideas? ☺

# THANK YOU

Composing State Machine Replicas

## Contacts:

**Odorico Machado Mendizabal**
odorico.mendizabal@ufsc.br
https://www.inf.ufsc.br/~odorico.mendizabal/

**Caroline Martins Alves**
caroline.martins@posgrad.ufsc.br

**Thaís Bardini Idalino**
thais.bardini@ufsc.br
https://thaisidalino.github.io/