# *Revolutionizing Mobile and Cloud via Coherence*

**Vijay Nagarajan**

Nicolai Oswald

Adarsh Patil

Mahesh Dananjaya

Carr Reece

Theo Olausson

Antonis Katsarakais

Vasilis Gavrielatos

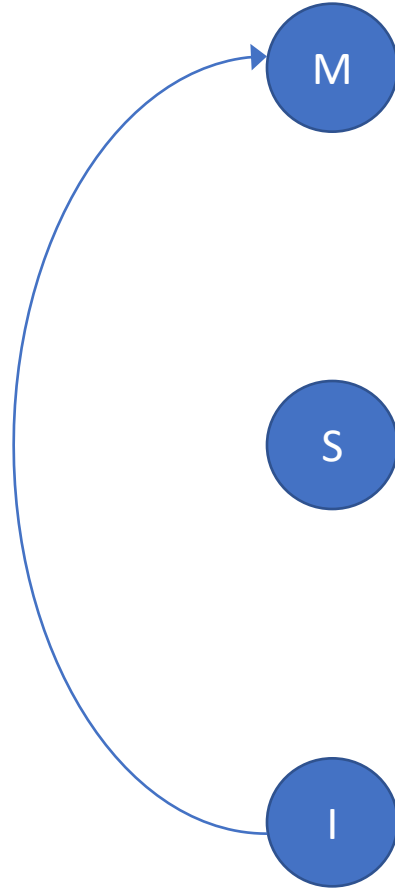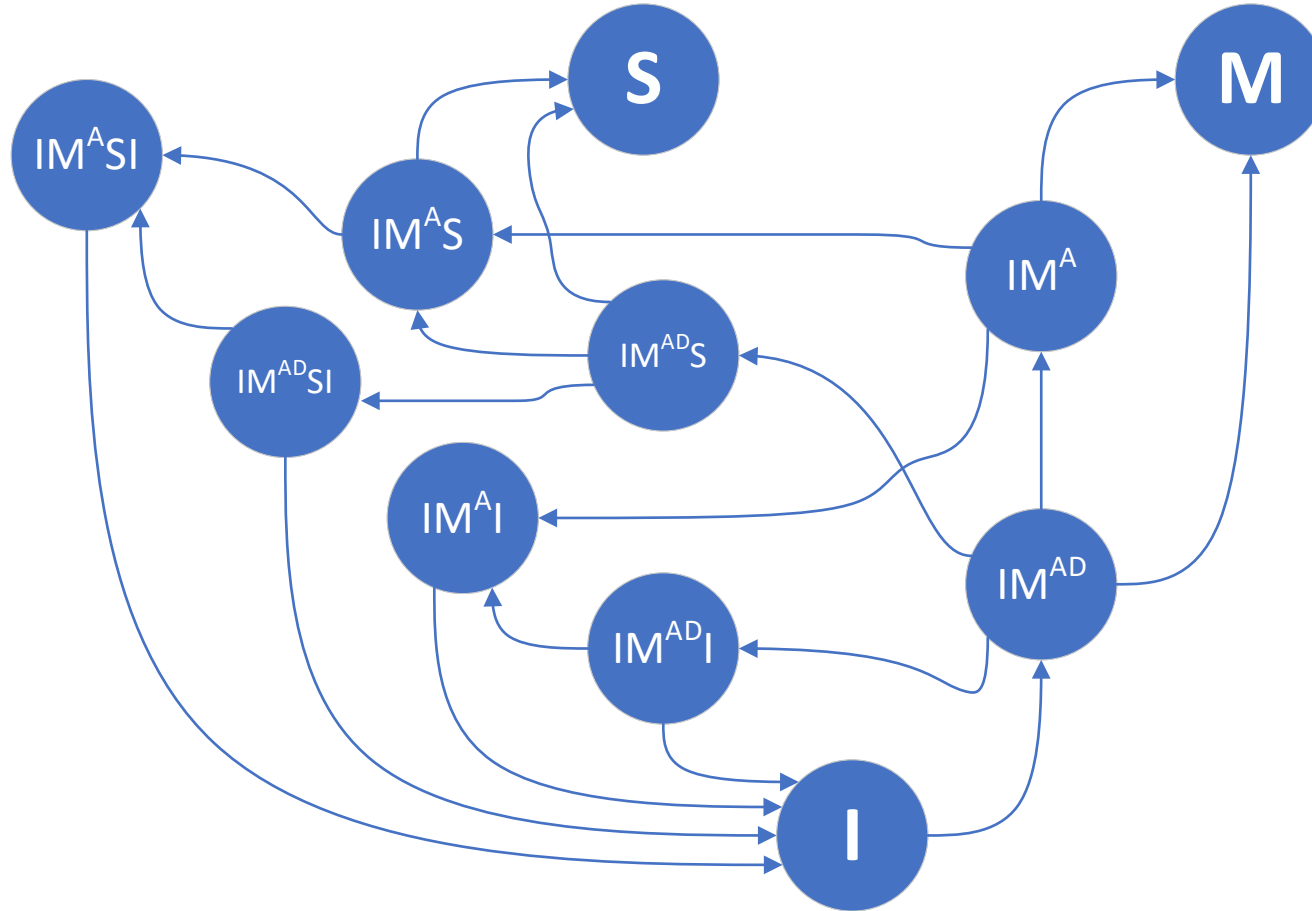Boris Grot

Dan Sorin
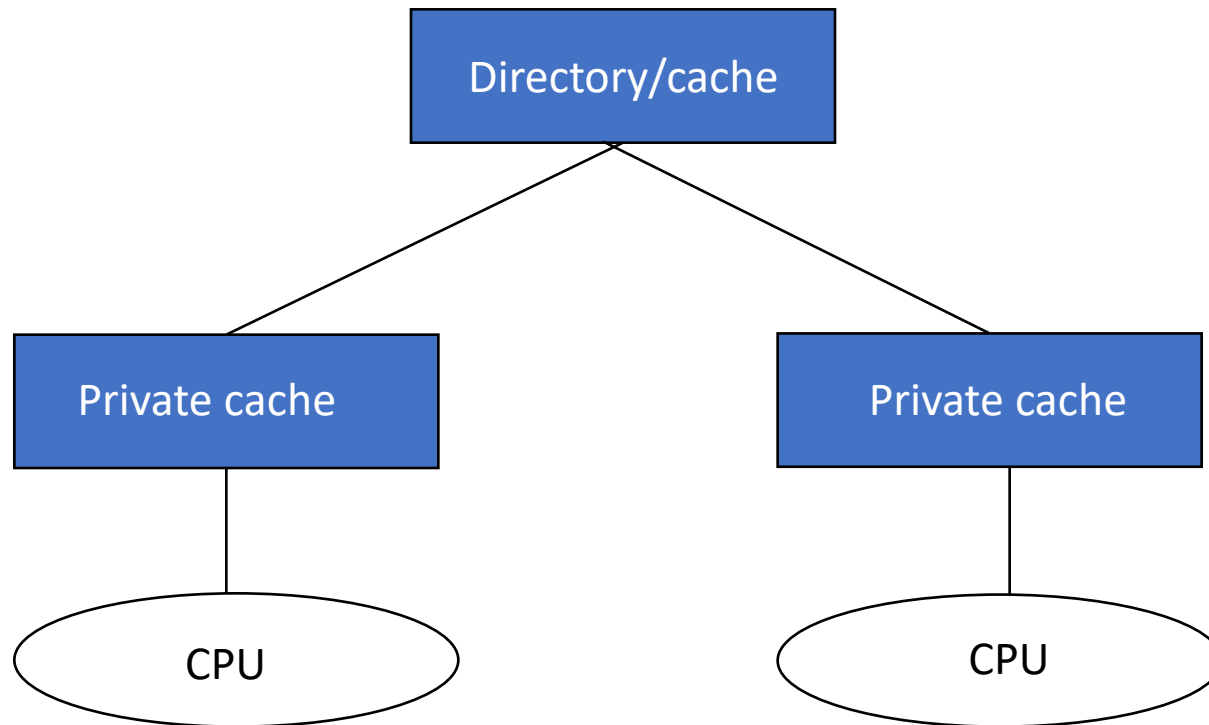
Nikos Nikoleris

Tobias Grosser

**M** Modified

**S** Shared

**I** Invalid

Concurrency!

# Hierarchy!

# Existing approach and its limitations

- Suppose one wants to build a multiprocessor SoC

  - Read ~100-500 page prose document (e.g., CHI).

  - Implement protocol by hand in Verilog

AMBA⋅ 5 CHI
**Architecture Specification**

**arm**
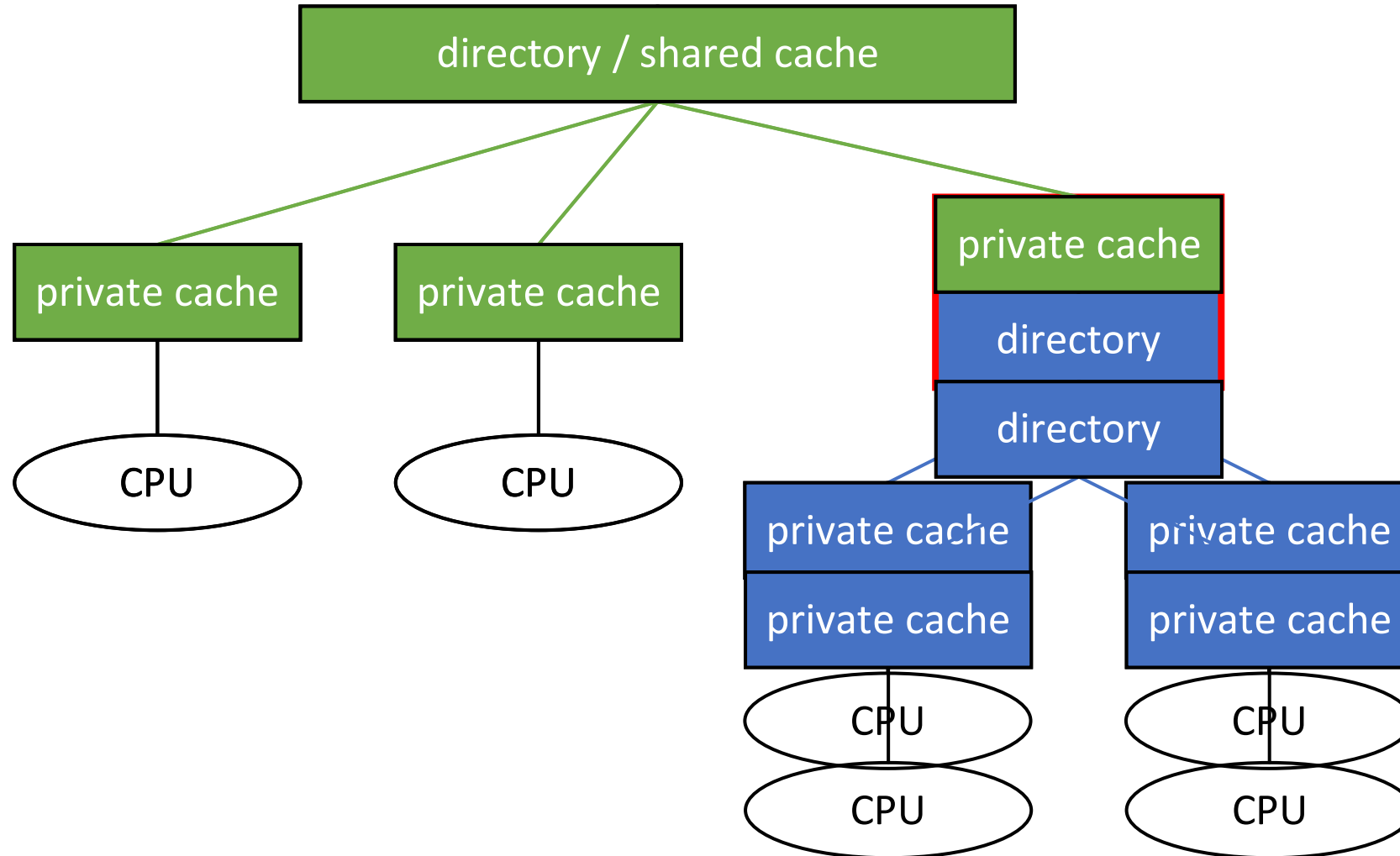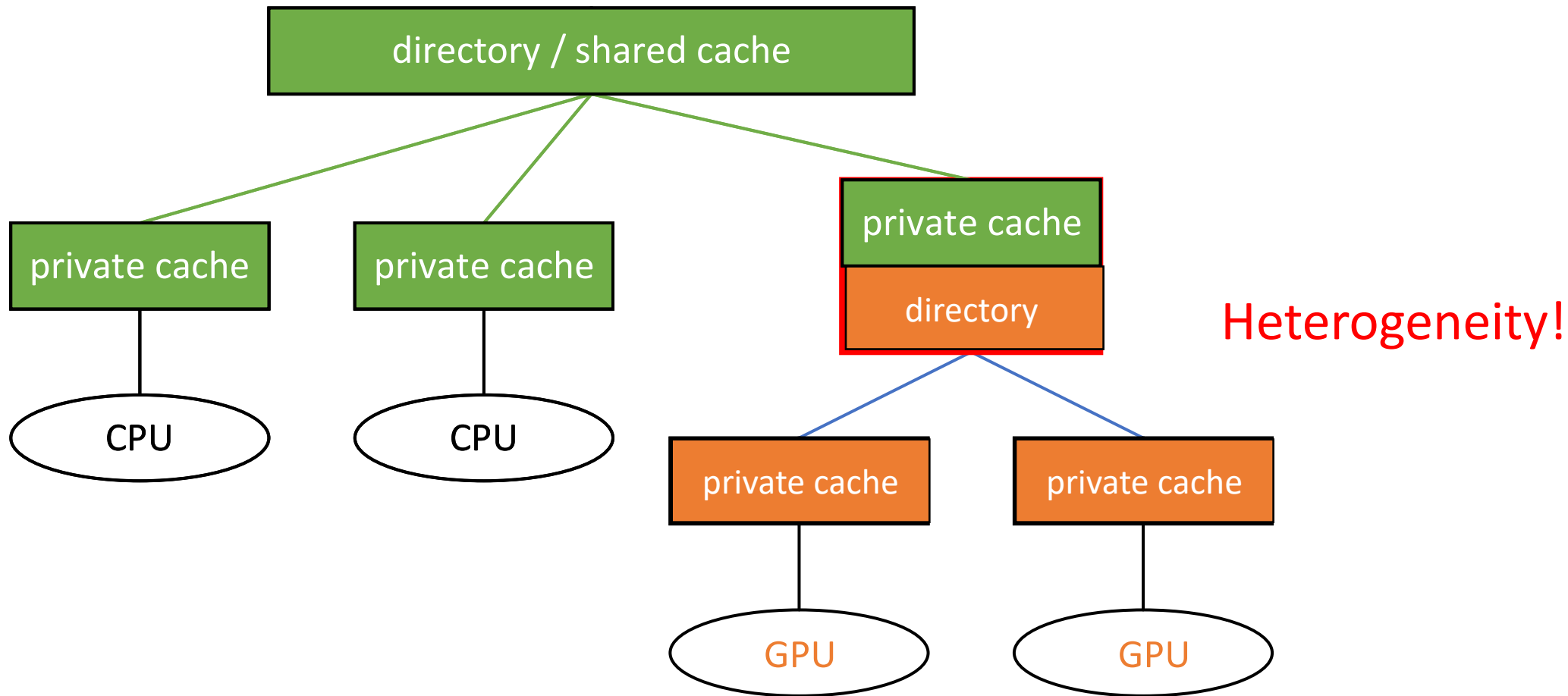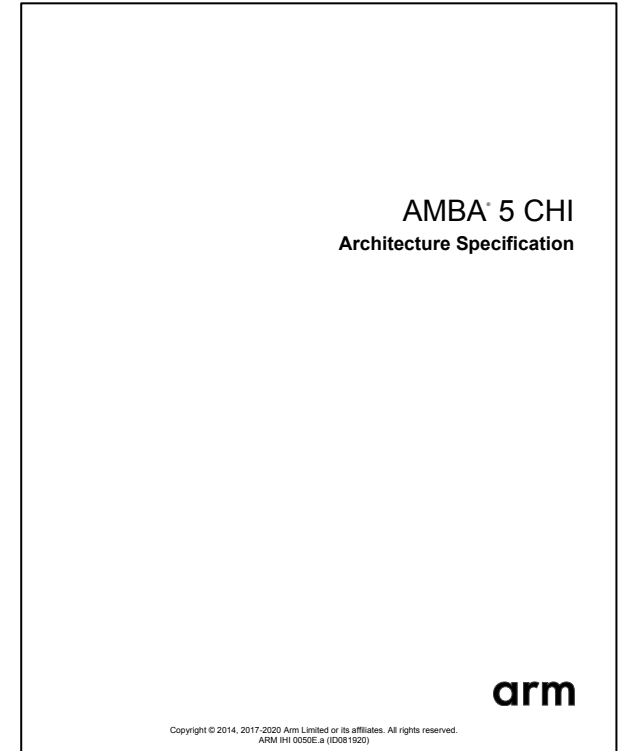
THE UNIVERSITY of EDINBURGH
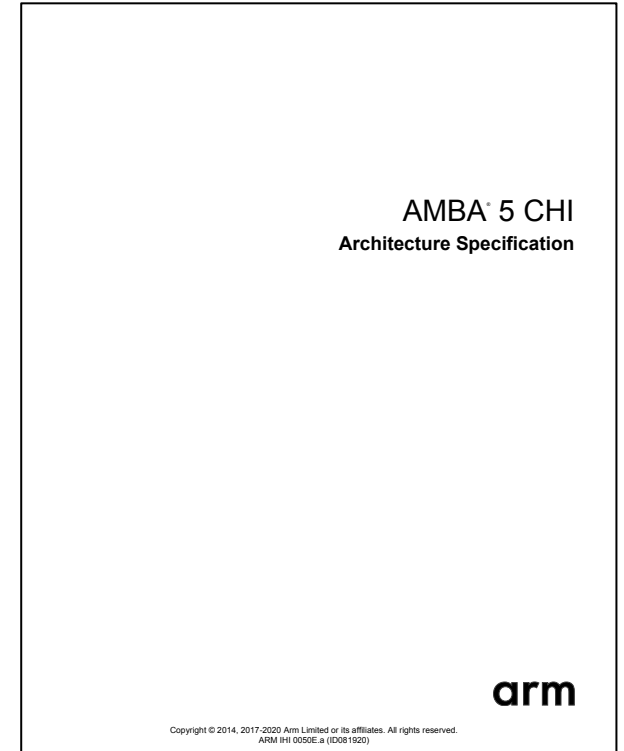**informatics**

# Existing approach and its limitations

- Suppose one wants to build a multiprocessor SoC

  - Read ~100-300 page prose document (Tilelink, CHI).

  - Implement protocol by hand in Verilog

- Limitations
  - Prose = Imprecise
  - Non-exhaustive and conservative
  - Only MOESI

AMBA 5 CHI
**Architecture Specification**

arm

# Our Approach: i/p

# Our Approach: o/p



directory / shared cache

private cache

private cache

private cache
directory

Murphi
Gem5
Verilog

CPU

CPU

private cache

private cache

GPU

GPU

# Outline

- Background and Motivation

- Concurrency: ProtoGen

- Hierarchy: HieraGen

- Heterogeneity: HeteroGen

- Coherence for the cloud

# Cache Coherence


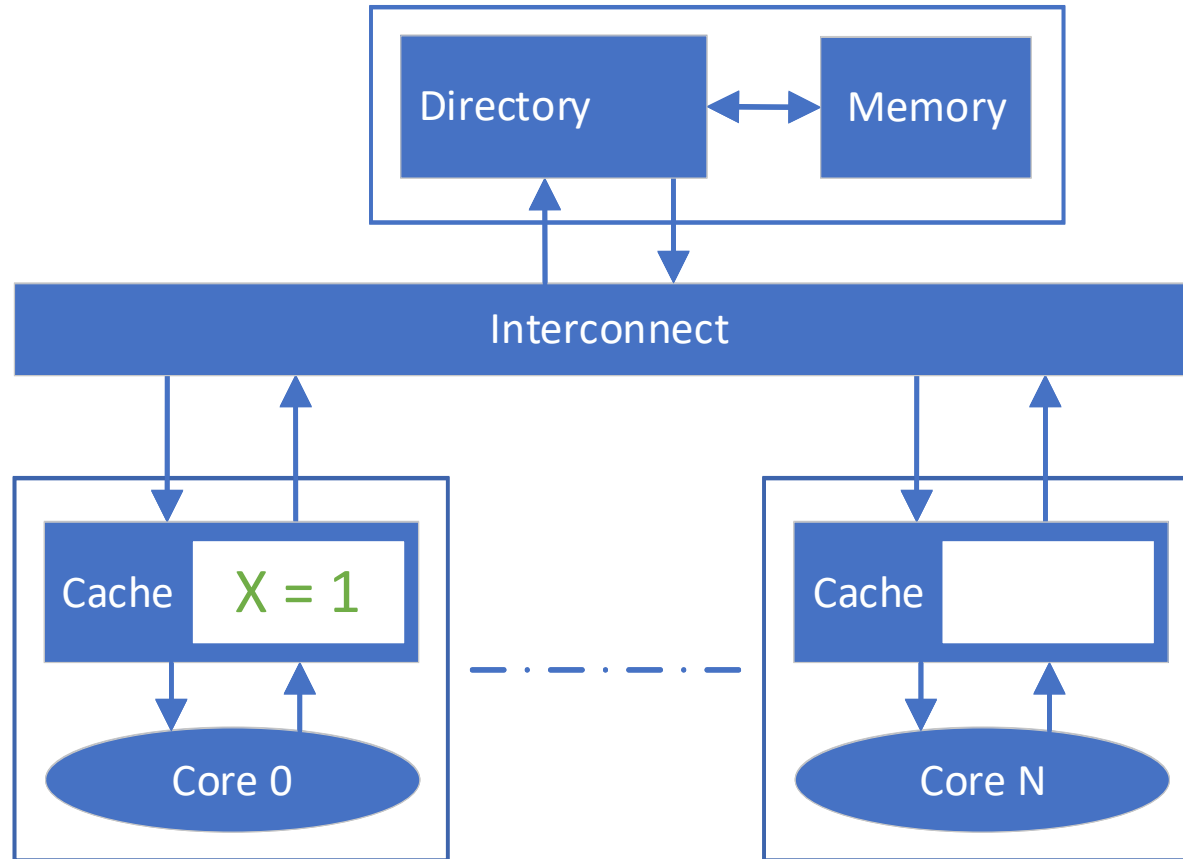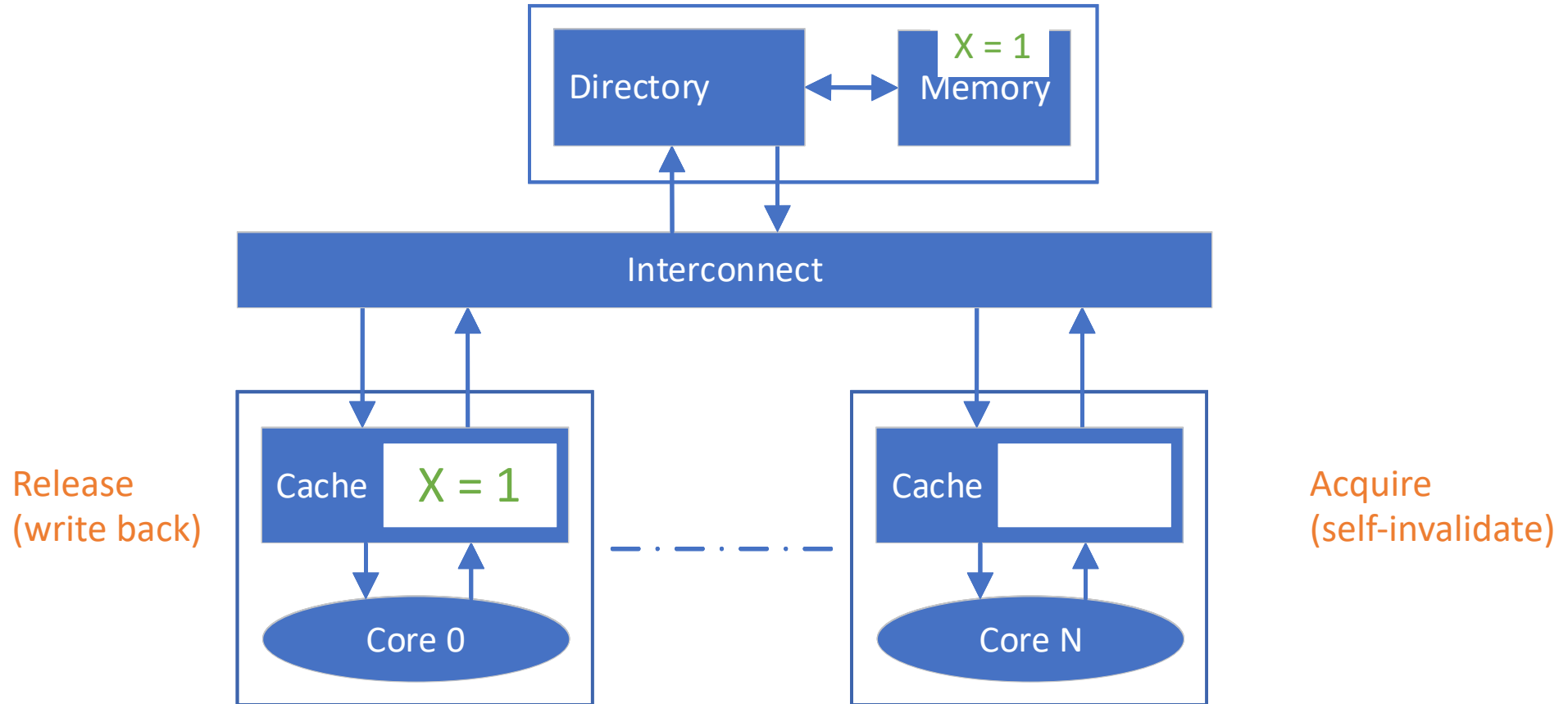
- SWMR: single-writer, multiple-reader invariant

# Consistency-directed Cache Coherence

"Cache coherence protocols are notoriously difficult to design and verif[y] [... Memory Systems, 2004]

"The coherence problem is difficult, because it requires coordinating events across nodes" [IEEE Concurrency 2000]

"... directory-based cach[e] ...

an[d] ...usly

"Sophis[ticated] difficult ...

"... designing and verifying a new hardware coherence protocol is difficult" [Spandex: A Flexible Interface for Efficient Heterogeneous Coherence - ISCA 2018]

...lt to design and implement correctly" [ASPLOS 2017]

"Cache coherence protocols for distributed shared memory multiprocessors are notoriously difficult to design" [ICFS 1996]

# Bugs in the Wild

| No. | Errata Description |
|-----|--------------------|
| 1 | **63 TLB Flush Filter Causes Coherency Problem in Multiprocessor Systems**

**Description**

If the TLB flush filter is enabled in a multiprocessor co... between the page tables in memory and the translations ... the possible use of stale translations even after software ...

**Potential Effect on System**

Unpredictable system failure.

**Suggested Workaround**

In MP systems, disable the TLB flush filter by setting HWCR.FFDIS (bit 6 of MSR 0xC001_0015).

**Fix Planned**

Yes |
| 69 | Multiprocessor Coherency Problem with Hardware Prefetch Mechanism |
| 70 | Microcode Patch Loading in 64-bit Mode Fails To Use EDX |

From AnandTech: "… coherency was broken and manually disabled on the Galaxy S 4. The implications are serious from a power consumption (and performance) standpoint."

THE UNIVERSITY of EDINBURGH
**informatics**

# Why is Coherence Hard?

- Concurrency

- Hierarchy

- Heterogeneity
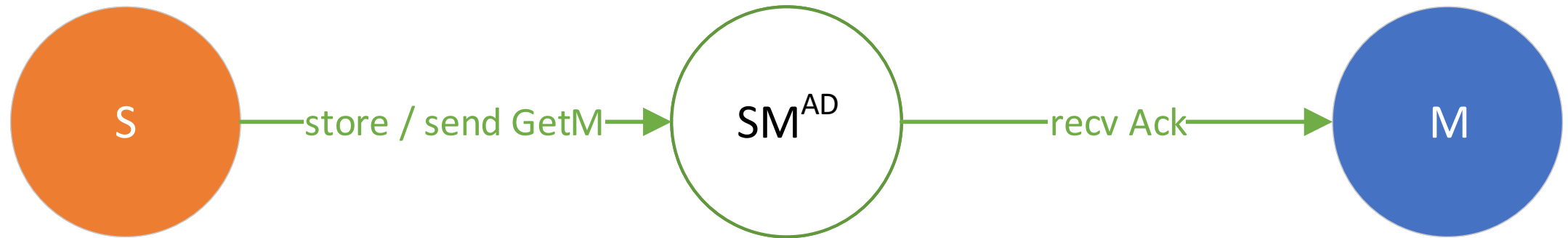
THE UNIVERSITY of EDINBURGH
**informatics**

# Outline

- Background and Motivation

- Concurrency: ProtoGen

- Hierarchy: HieraGen

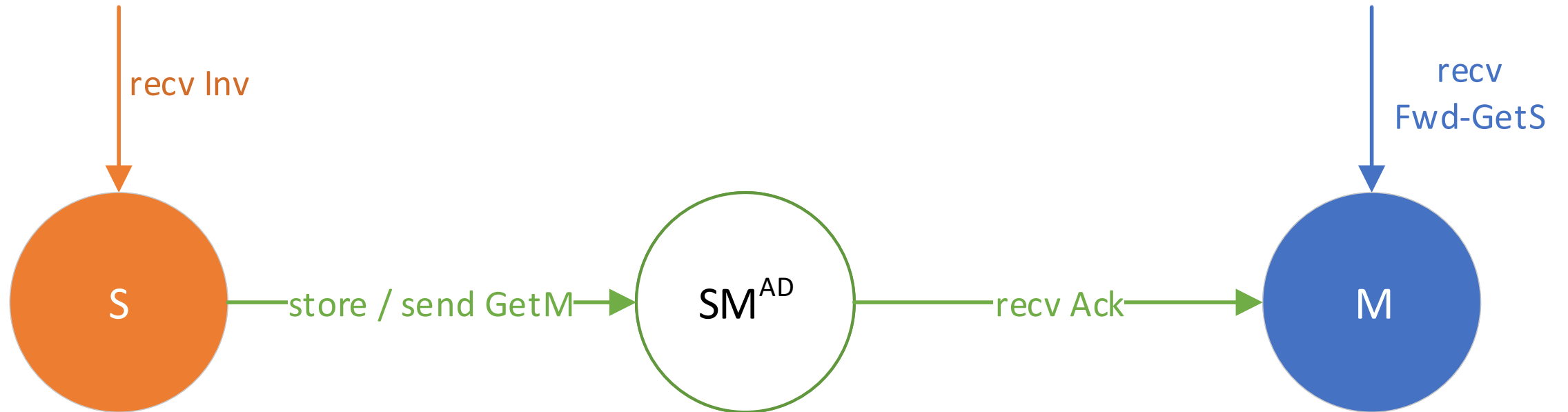- Heterogeneity: HeteroGen

- Coherence for the cloud

physical atomicity
logical atomicity

# Atomic S to M Transition

# Transient States



S → (store / send GetM) → SM$^{AD}$ → (recv Ack) → M
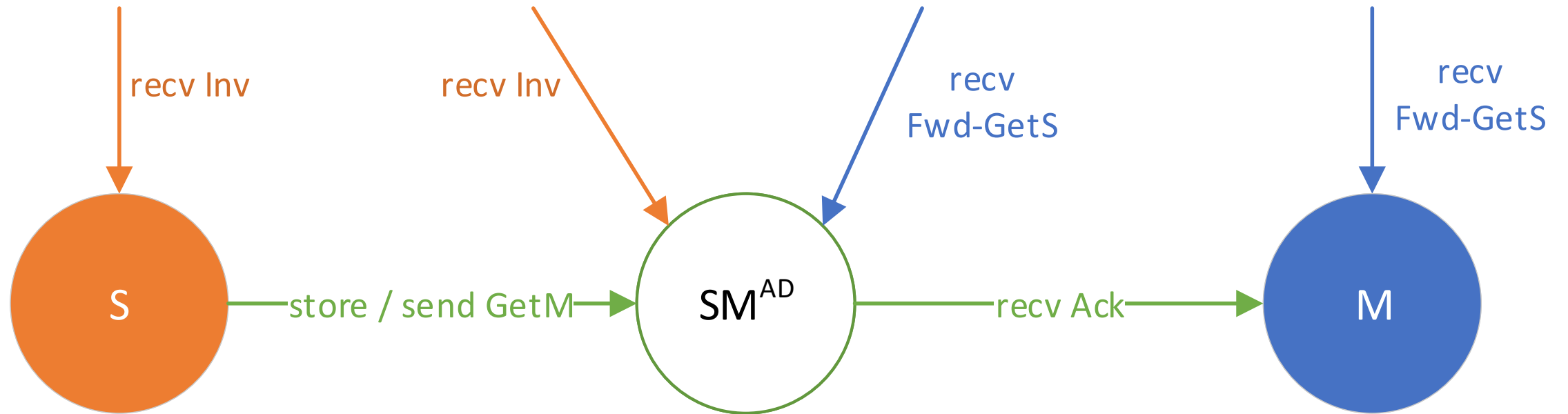
non-atomic transaction

# Concurrent Transactions

# Concurrent Transactions



non-atomic transactions + concurrency = complexity

# To Summarize...

- Stable state protocols assume physically atomic transactions

- Need to support concurrency for performance

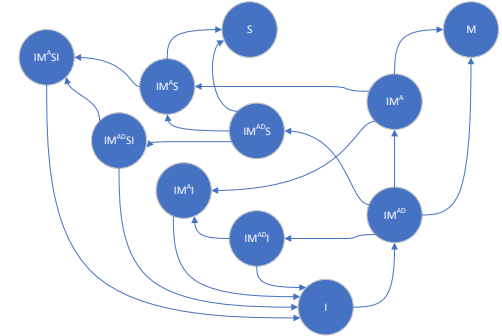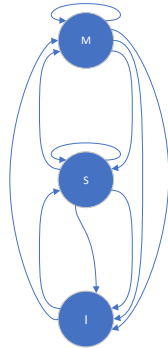- Transient states required to provide logically atomic transactions

# Key realization…

- Stable state protocol is a *sequential* specification

- The final protocol is a *non-blocking concurrent* implementation

- Transient states are *synchronization* operations
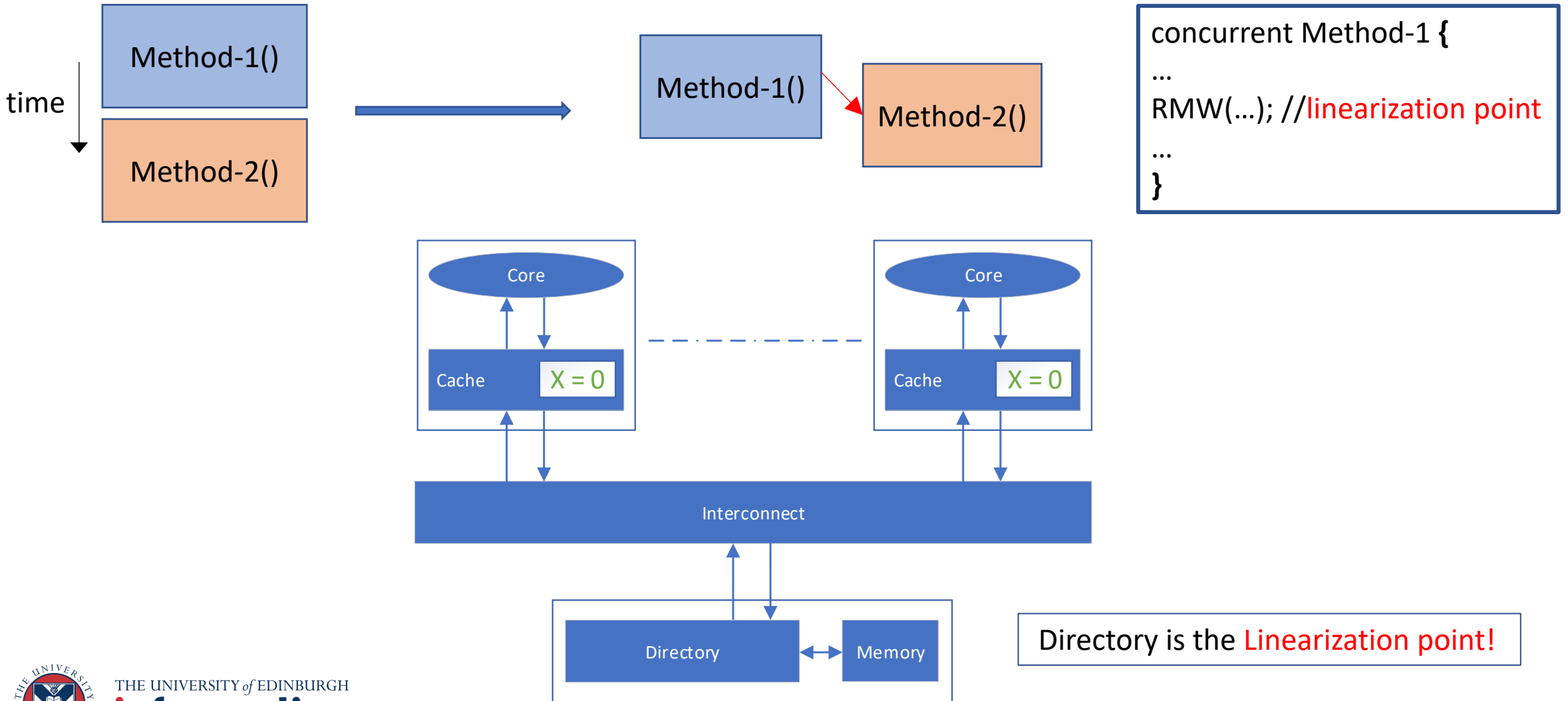
# Insight



Sequential object {
...
...
...
}

Non-blocking concurrent {
...
...
...
}

No wonder cache coherence protocols are Hard!

# Insight
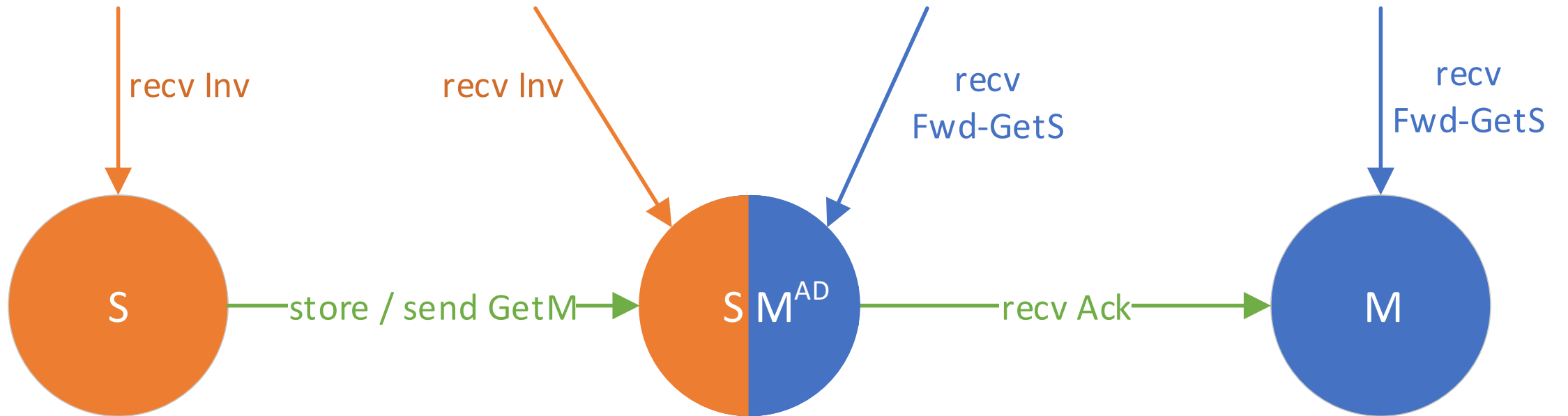
# Demystifying Transient States

How do transient states provide logical atomicity?

- Convey directory serialization order to caches

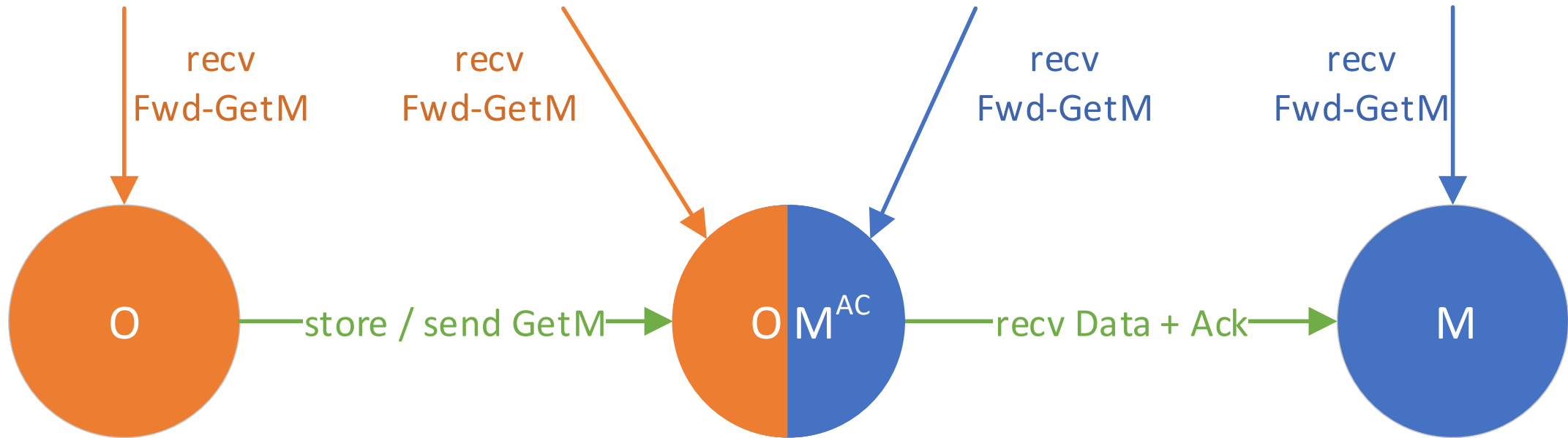- Transient states ensure that caches obey this order

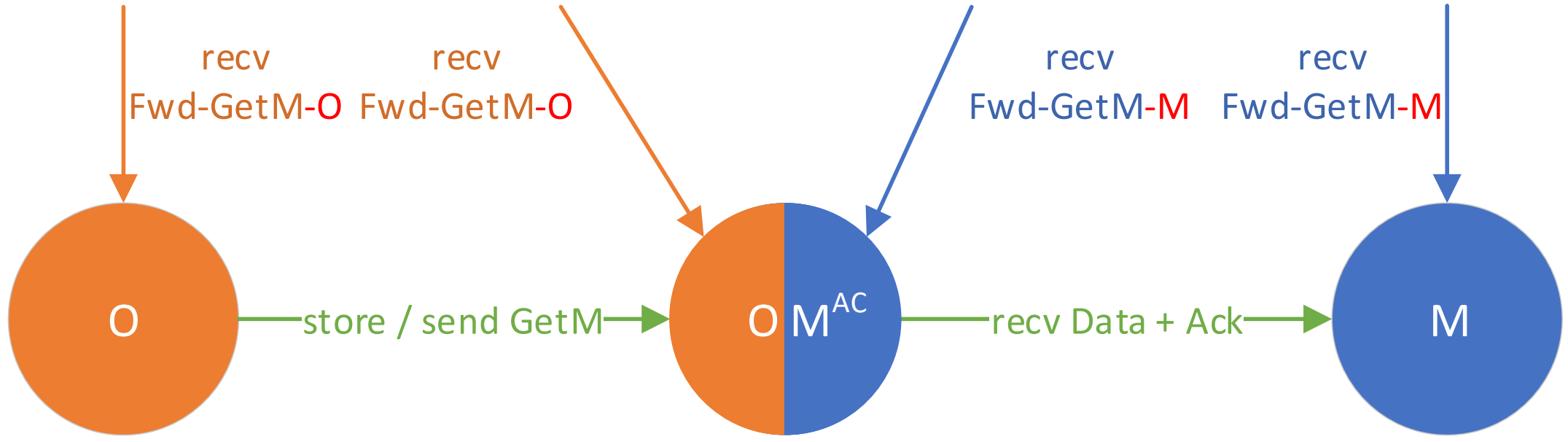ProtoGen  automates by leveraging this insight!

THE UNIVERSITY of EDINBURGH
informatics

# How does cache infer serialization order?

# How to resolve name conflicts?
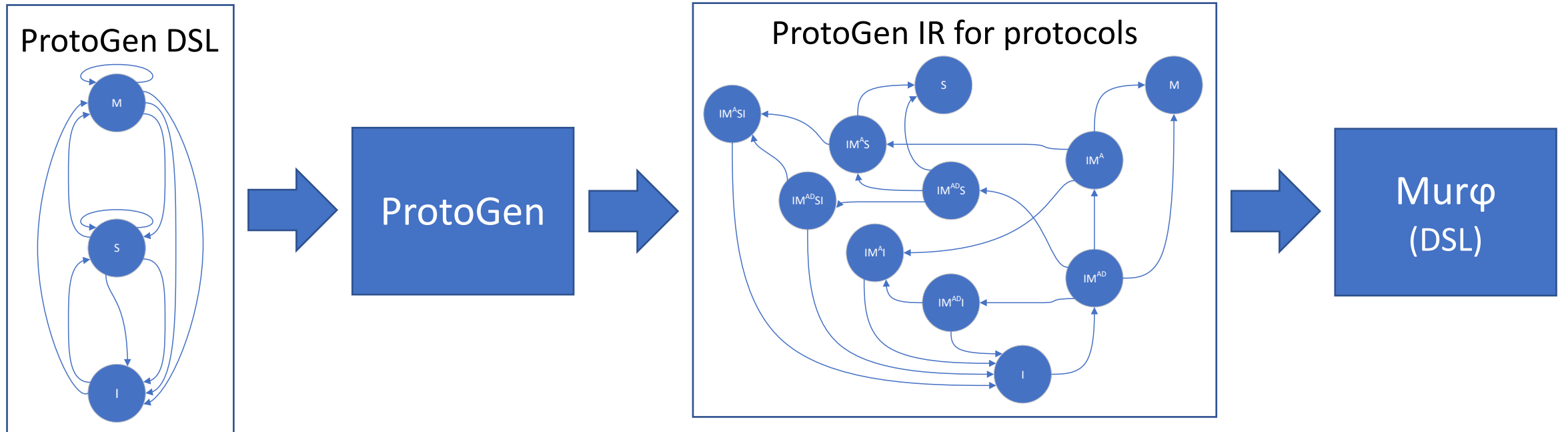
# Rename Messages

# ProtoGen Summary

- **Infer** serialization order from incoming messages

- **Rename** messages in order to achieve this

- **React** like in stable state
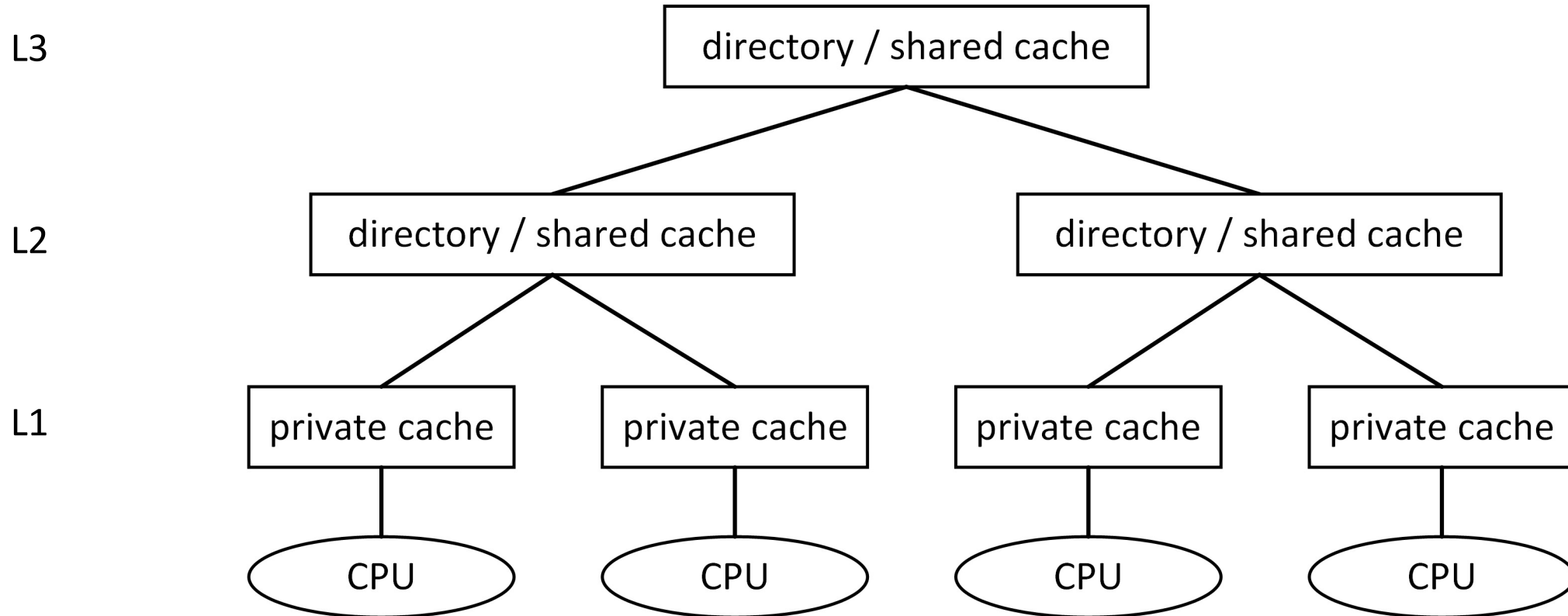
# ProtoGen Tool



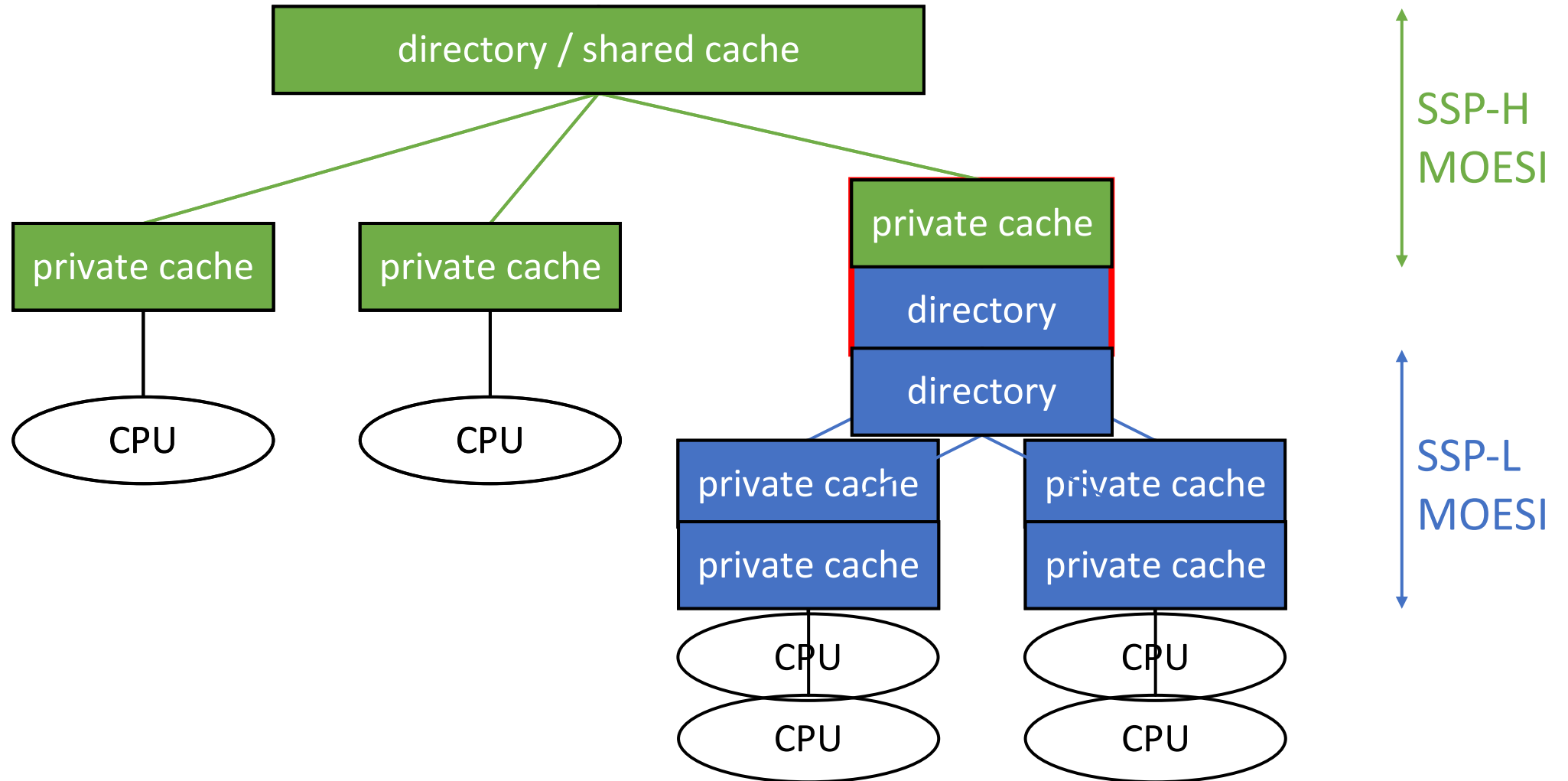ProtoGen as good (or better) than manually generated protocols

# Outline

- Background and Motivation

- Concurrency: ProtoGen

- Hierarchy: HieraGen
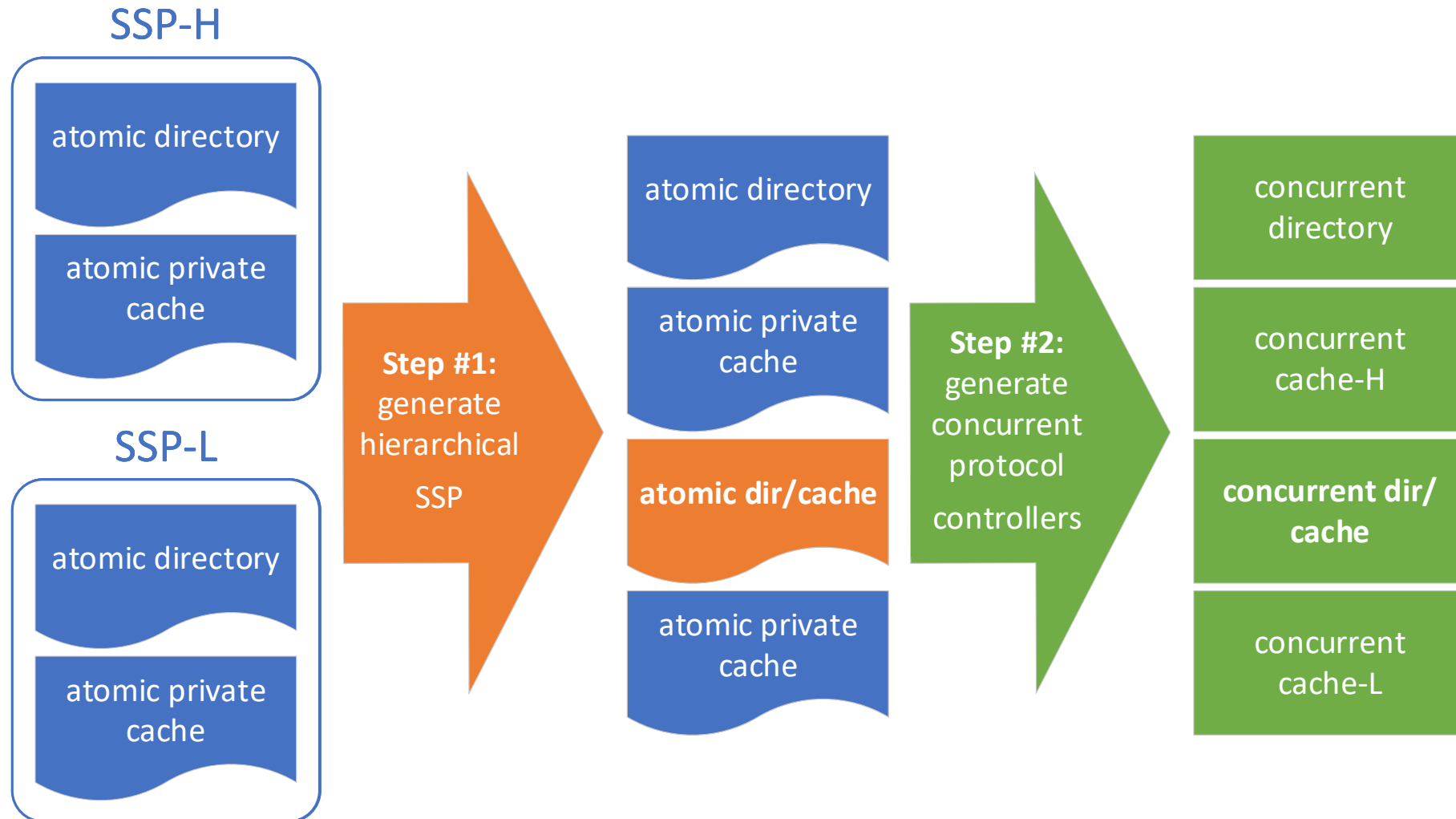
- Heterogeneity: HeteroGen

- Coherence for the cloud

THE UNIVERSITY of EDINBURGH
informatics

# Hierarchical protocols

# The Complexity of Hierarchical protocols

# HieraGen* tool flow

SSP-H

atomic directory

atomic private cache

SSP-L

atomic directory

atomic private cache

**Step #1:** generate hierarchical SSP

atomic directory

atomic private cache

**atomic dir/cache**

atomic private cache

**Step #2:** generate concurrent protocol controllers

concurrent directory

concurrent cache-H

**concurrent dir/ cache**

concurrent cache-L
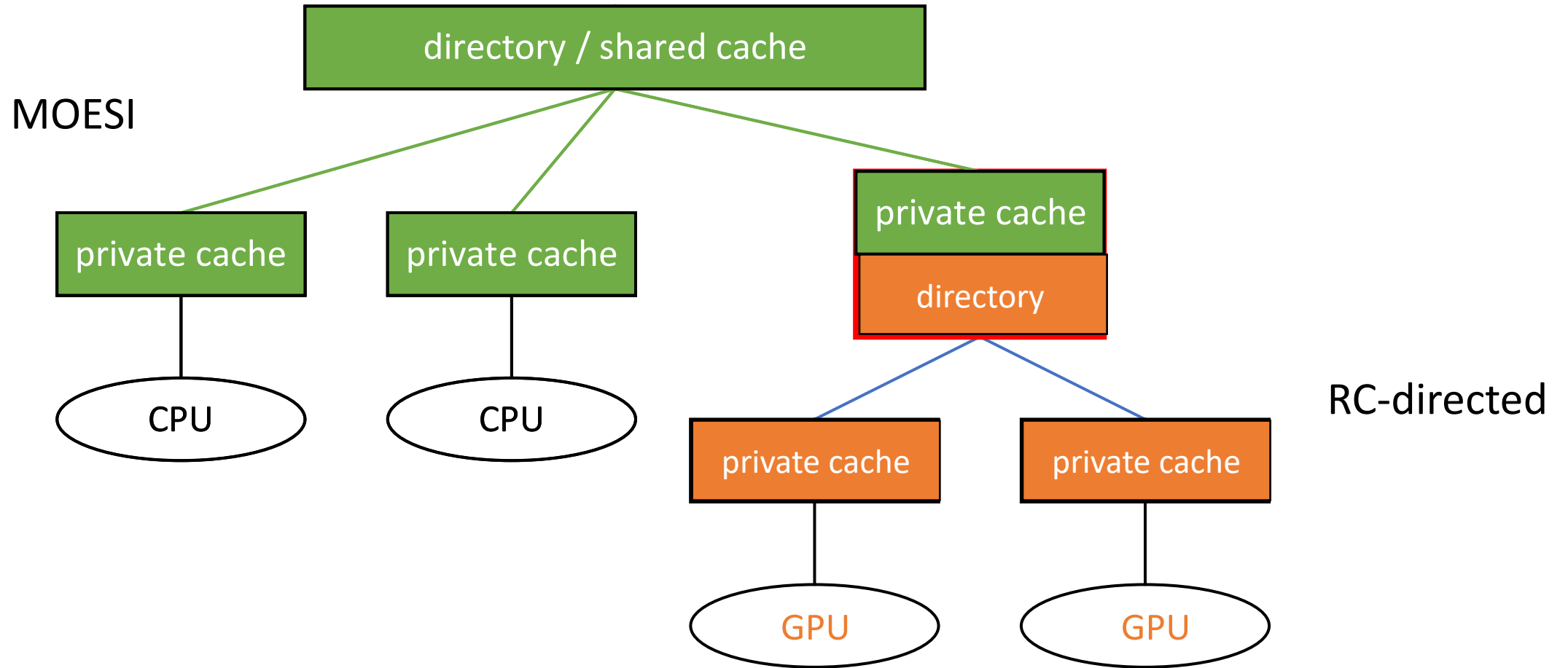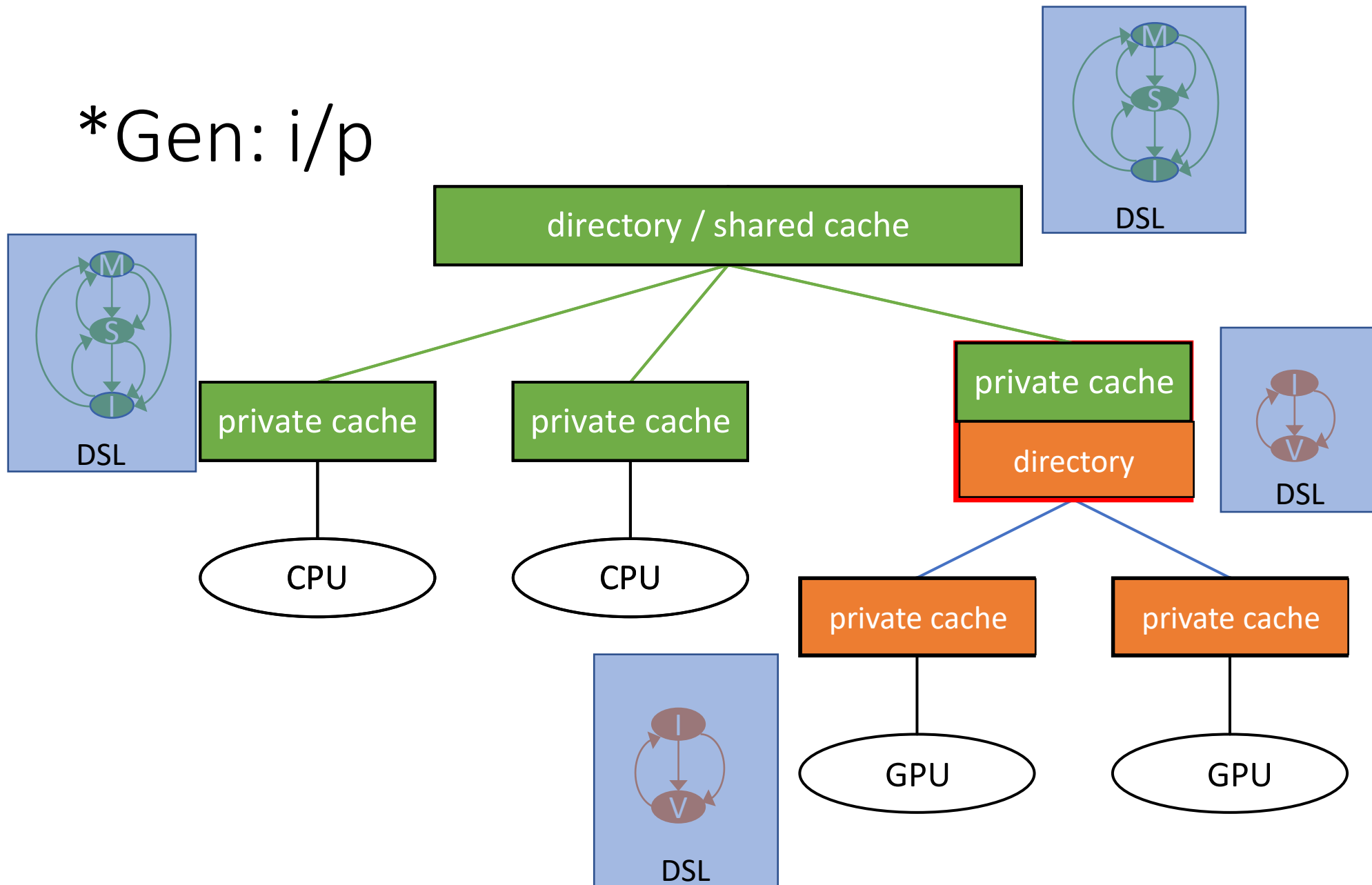
*ISCA'20

# HeteroGen

- How do you stich together two different protocols?
  - HieraGen should work!

- What is the correctness condition?
  - MOESI style protocols SC
  - RC-style protocols…RC
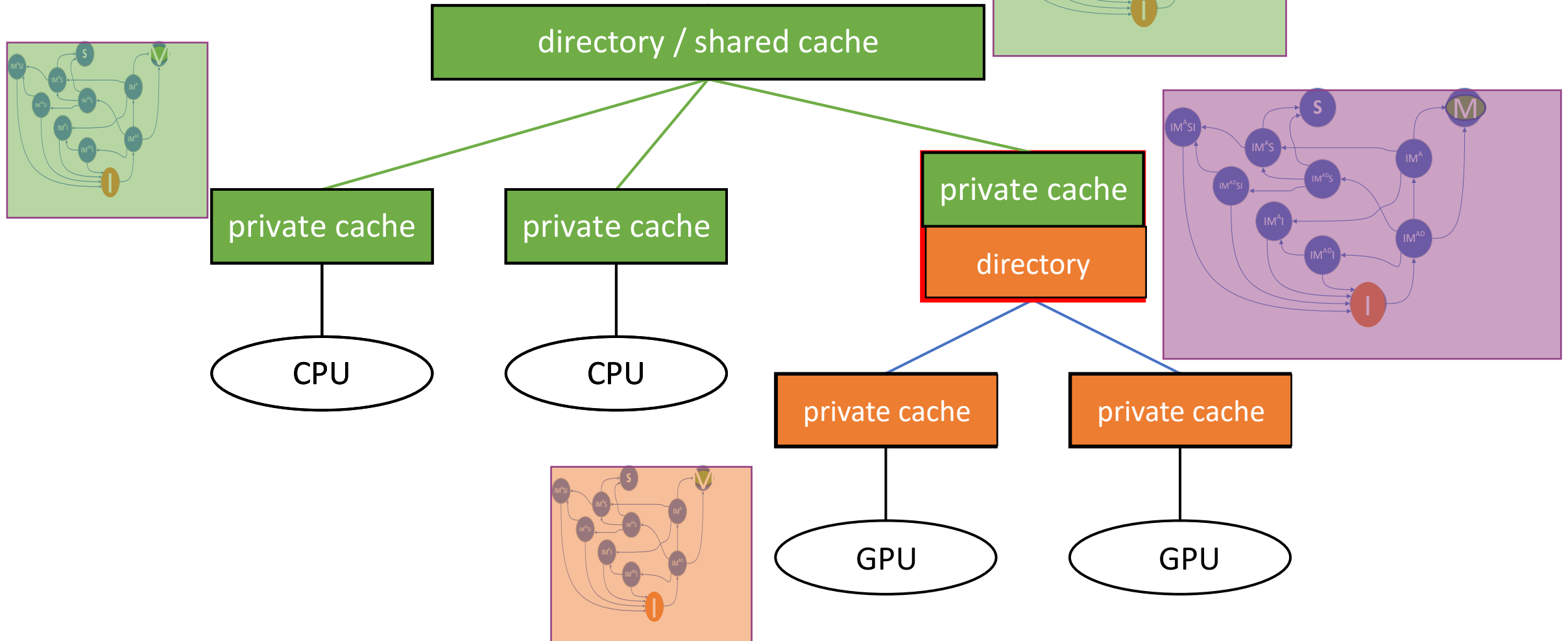  - Compound consistency models!

[under submission]

# Compound Consistency

- Correctness condition for heterogeneous coherence

- Foundation for heterogeneous consistency

- Each cluster can assume its own memory model

*Gen: i/p

# *Gen: o/p

# Outline

- Background and Motivation

- Concurrency: ProtoGen

- Hierarchy: HieraGen

- Heterogeneity: HeteroGen

- **Coherence for the cloud**

THE UNIVERSITY of EDINBURGH
**informatics**

# Datacentre Distributed datastores

In-memory with read/write API

Backbone of online services

Distributed Datastore

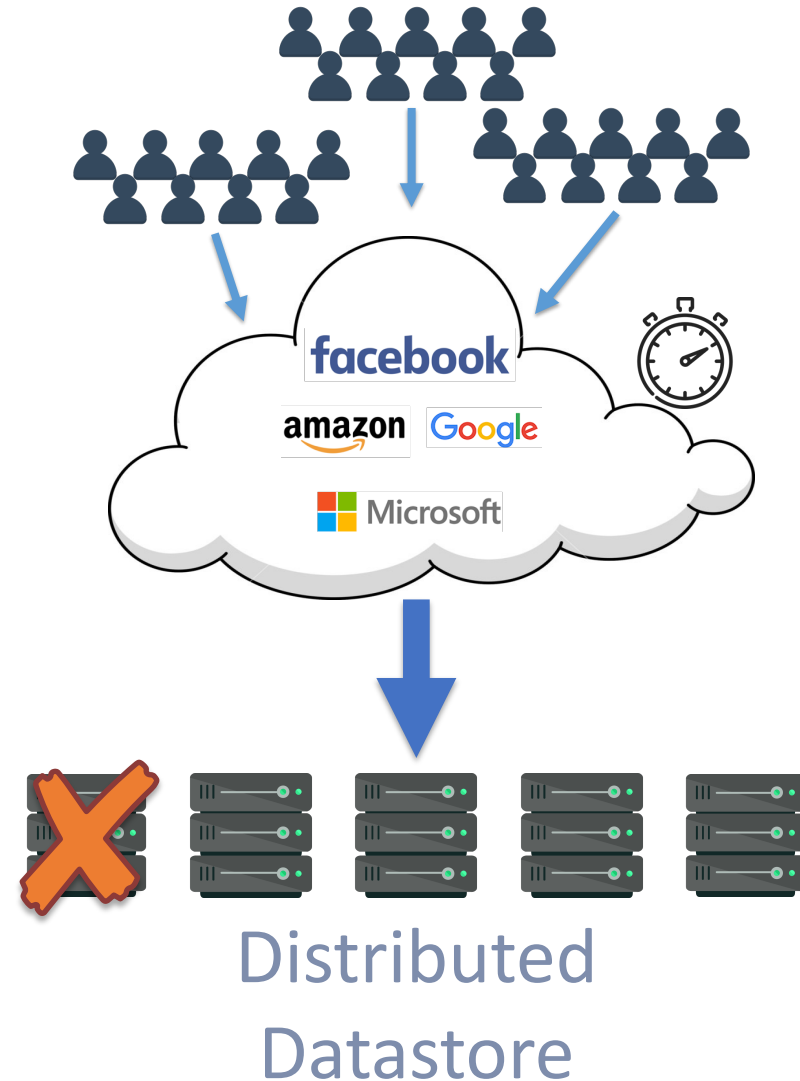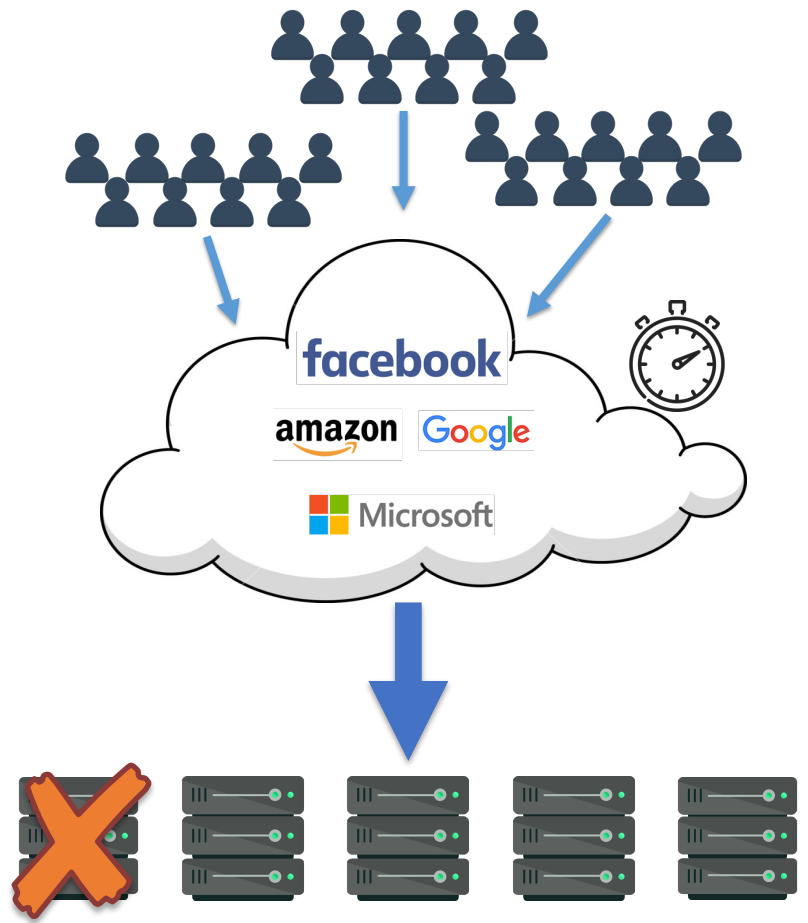# Distributed datastores

In-memory with read/write API

Backbone of online services

Need:

**Consistency (Programmability)**

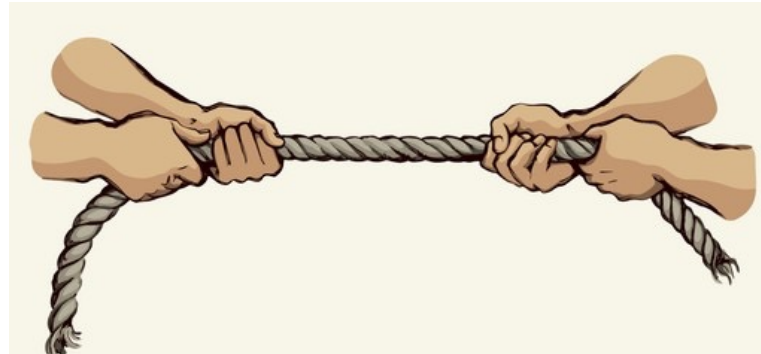**High performance**

**Fault tolerance (Availability)**



Distributed
Datastore

# Distributed datastores

In-memory with read/write API

Backbone of online services

Need:

**Consistency (Programmability)**

**High performance**

**Fault tolerance (Availability)**



Distributed

High-performance Reliable Replication Protocol!

# The problem



Performance
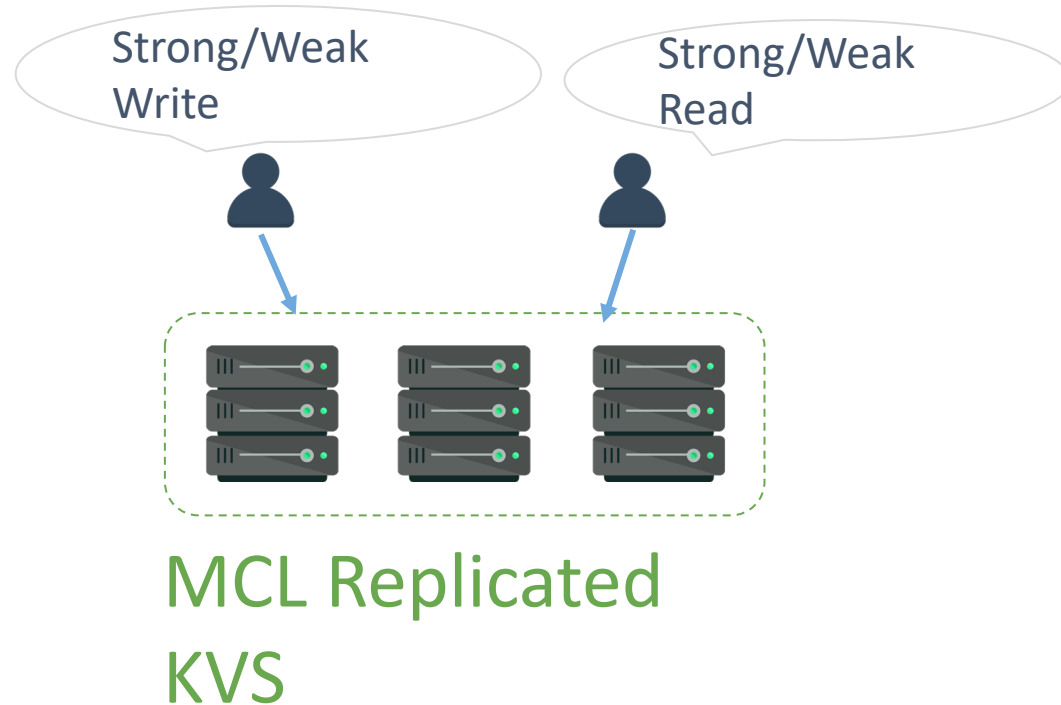
Strong
Consistency

Replication ⇒ Performance vs Consistency

# Existing Solution: Multiple Consistency Levels (MCL)

| | |
|---|---|
| aws | **Amazon DB** |
| G | **App Engine** |
| YAHOO! | **PNUTS** |
| (twitter) | **Manhattan** |
| Microsoft Azure | **Pileus** |

Strong/Weak Write

Strong/Weak Read

MCL Replicated KVS

What about programmability?

THE UNIVERSITY of EDINBURGH
**informatics**

# Datacentre Distributed datastore

# Datacentre Distributed datastore

get()
put()
Tx()

P    P    P    P

**Shared-memory multiprocessor!**

Rare

FAIL

In-memory

Rep    Rep    Rep    Rep

Fast LAN

Datacentre replication = Fast-path coherence + slow-path consensus?

# Coherence-inspired Hermes



*Broadcast-based invalidating*

# Coherence-inspired Hermes



*Broadcast-based invalidating*

Local reads

Fast, concurrent writes

Protocol reliable but blocking

# Shared Memory World

# Can we do the same for KVSes?



**DRF-SC**
Programming
Paradigm

Alice

**Release
Consistency**
DRF-compliant
memory model

Replicated KVS

THE UNIVERSITY of EDINBURGH
informatics
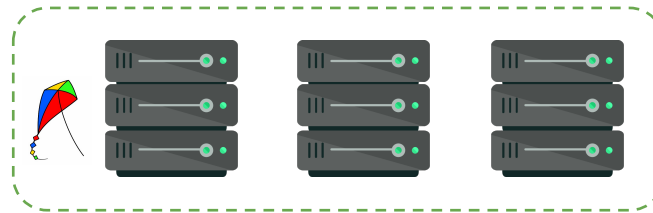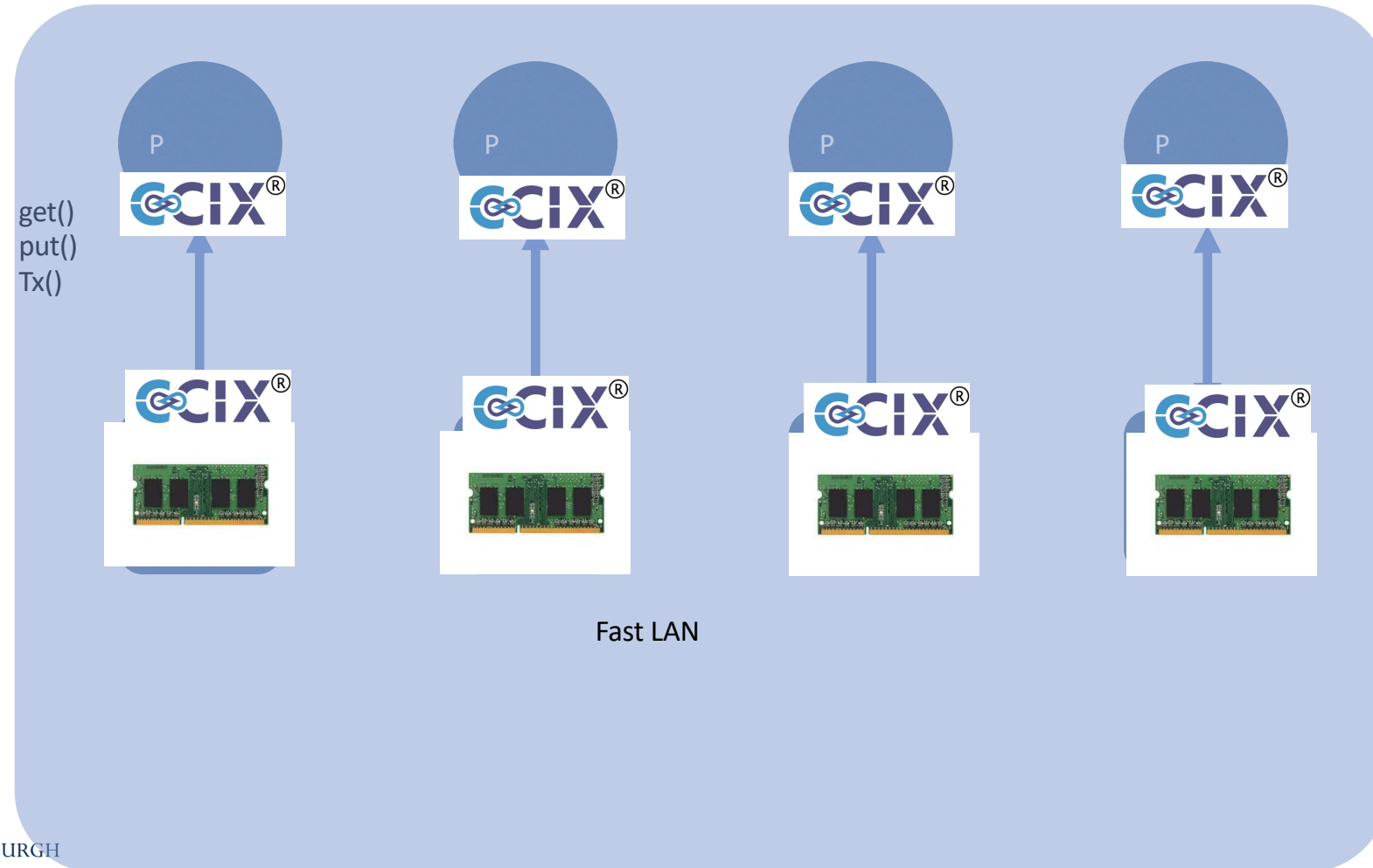
# Kite

A Replicated KVS with
- ➤ Release Consistency
- ➤ High Availability

Kite Replicated KVS

# Rethinking Datacentre Memory

# Revolutionizing Mobile/Cloud via Coherence

- Raise abstraction of coherence protocol design and automate
  - Concurrency
  - Hierarchy
  - Heterogeneity

- Datacentre coherence a great opportunity but needs new family of high-performance fault-tolerant coherence protocols.

THE UNIVERSITY of EDINBURGH
**informatics**