# Co-designing *reliability* and *performance* for datacenter memory

**Adarsh Patil**

**Doctoral Examination**

**30th May 2023**

**Advisor:**
  Vijay Nagarajan (UoE)
**Co-examiners:**
  Vilas Sridharan (AMD)
  Antonio Barbalace (UoE)

THE UNIVERSITY *of* EDINBURGH
**informatics**

# About me: My journey so far…

**(2012–14) Datacenter Infra Team – Virtualization & Linux Engineering**
*solutions architect:* platform benchmarking, performance analysis

**(2014–17) Masters by Research – HAShCache [TACO '18], TLB reach [arXiv]**
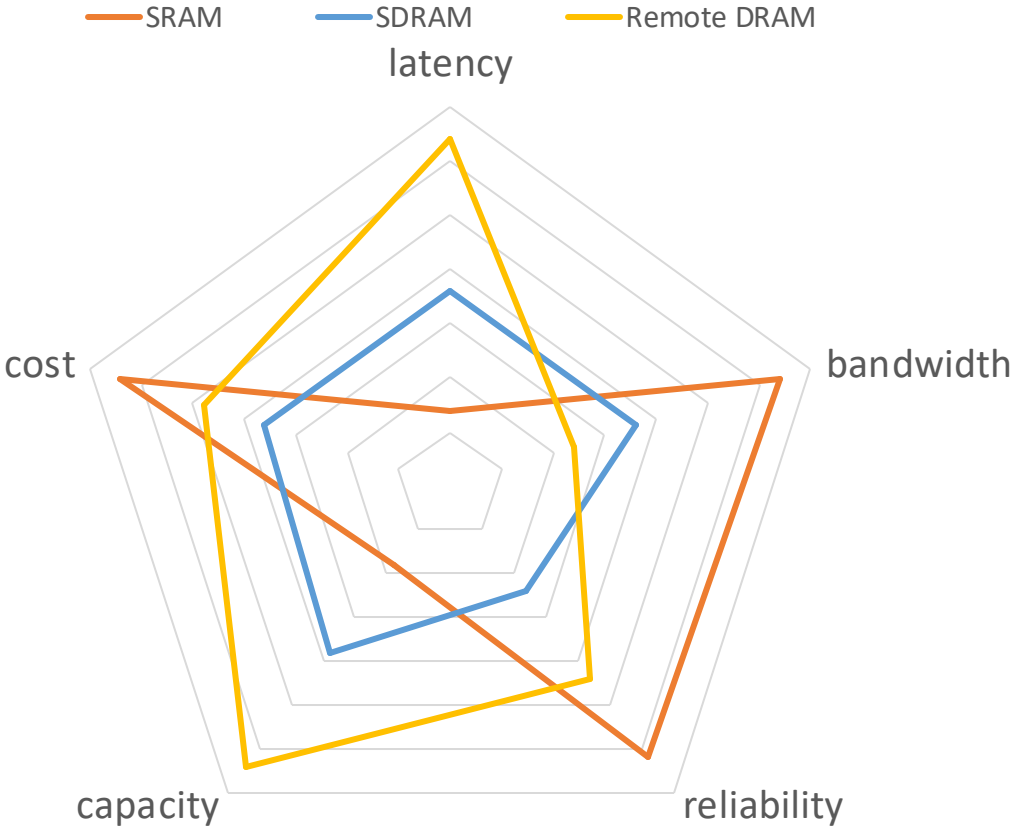*memory architecture:* DRAM cache, heterogeneous SoCs, virtual mem [Advisor: Prof. R. Govindarajan]

**(2017–19) Research Scientist – HPC ecosystem and applications team**
*application understanding:* s/w optimization, h/w architecture for next gen
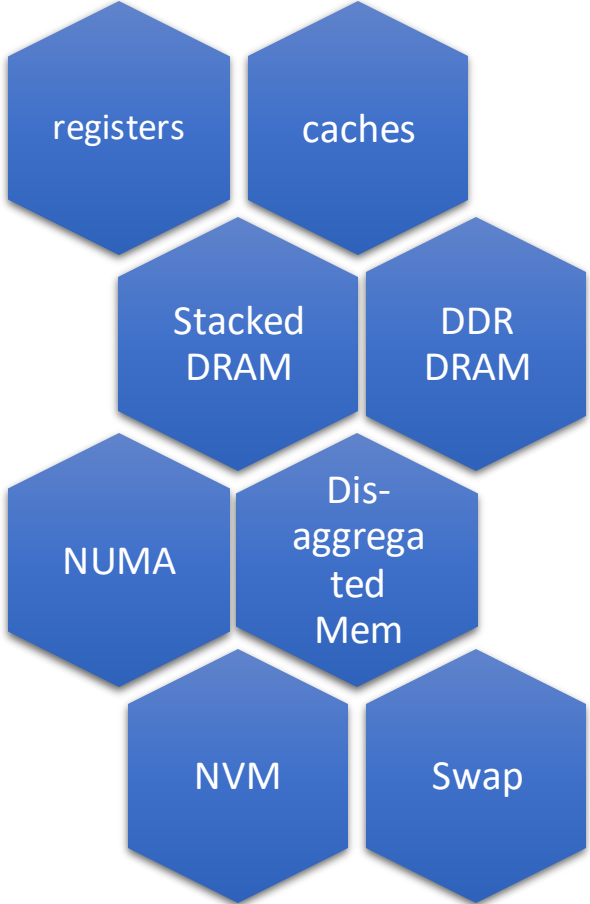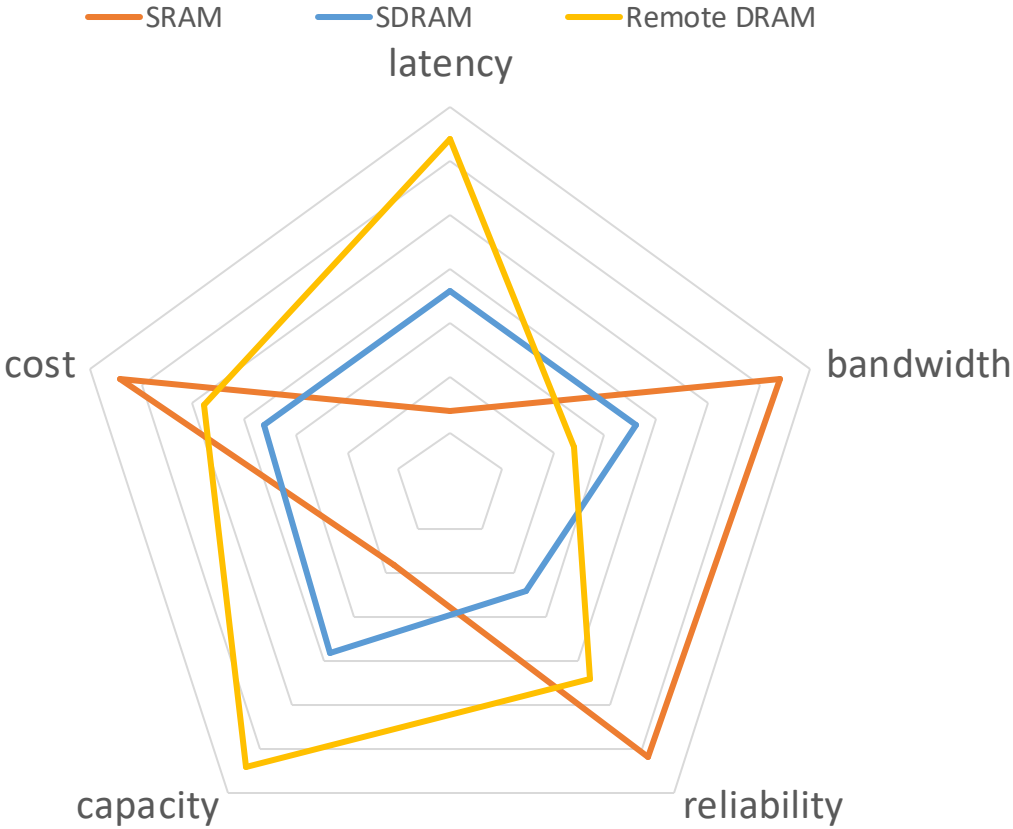
**(2019–now) – Co-designing reliability and performance for the datacenter**
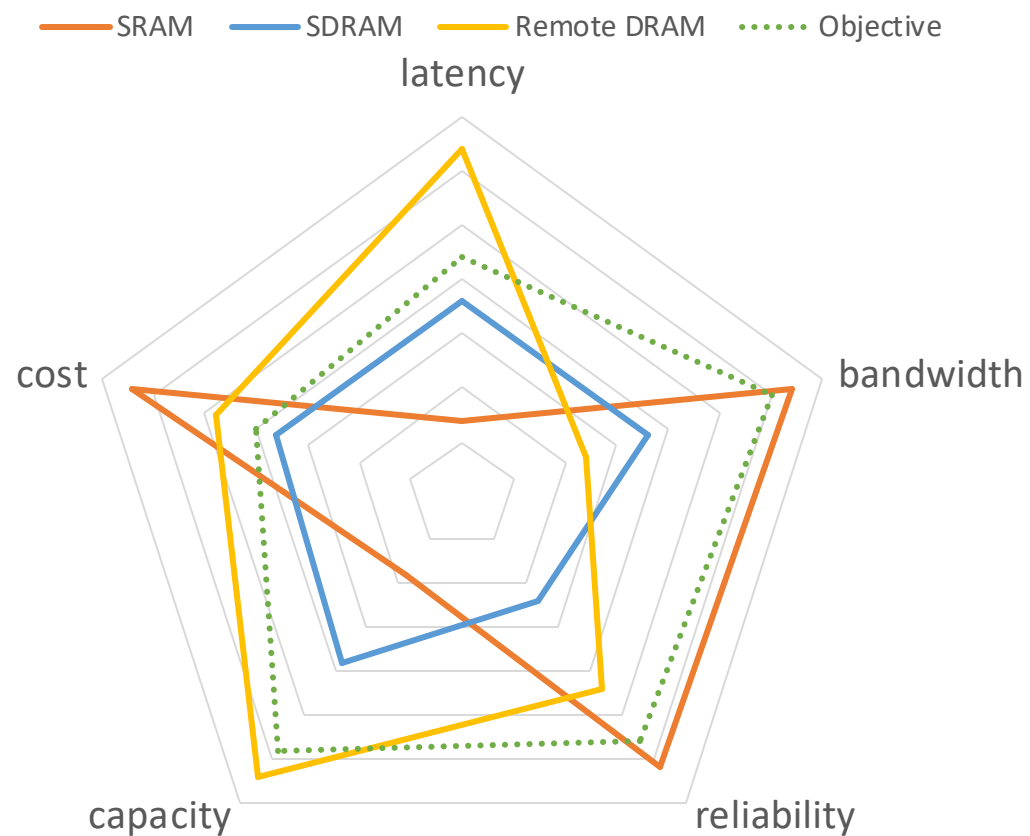*holistic approach:* integrating hardware + application
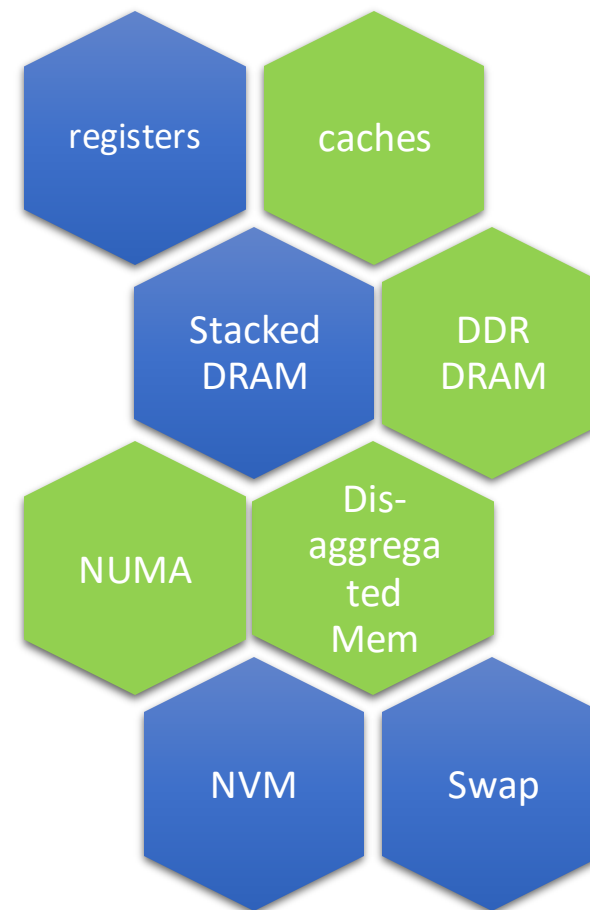
# Memory – a perpetual conundrum!
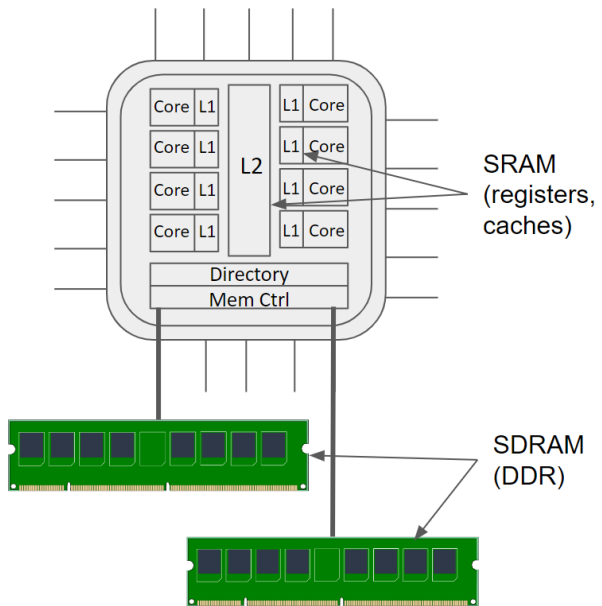
# Memory – a perpetual conundrum!

# Thesis objective

# Datacenter memory



SRAM
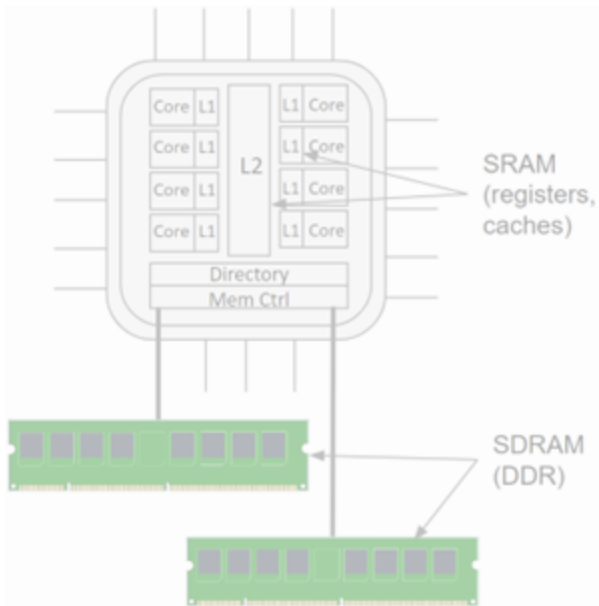(registers,
caches)

SDRAM
(DDR)

# Datacenter memory



SRAM (registers, caches)

SDRAM (DDR)

NUMA

Socket 0

Socket 1

cache coherent

interconnect

Remote DRAM for Socket 1

Remote DRAM for Socket 0

# Datacenter memory



NUMA

Disaggregated memory

# Main memory is comprised of DRAMs



Shared main memory

SRAM (registers, caches)

SDRAM (DDR)

NUMA

Socket 0

Socket 1

cache coherent interconnect

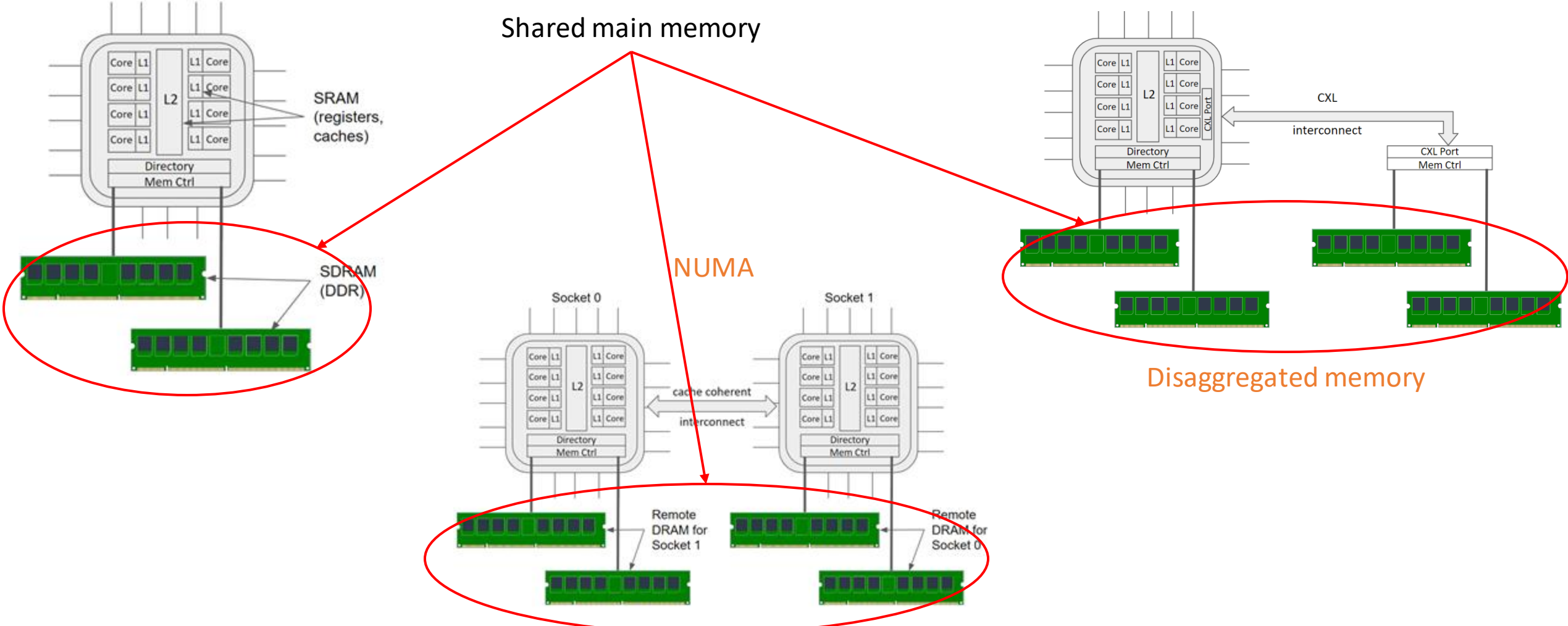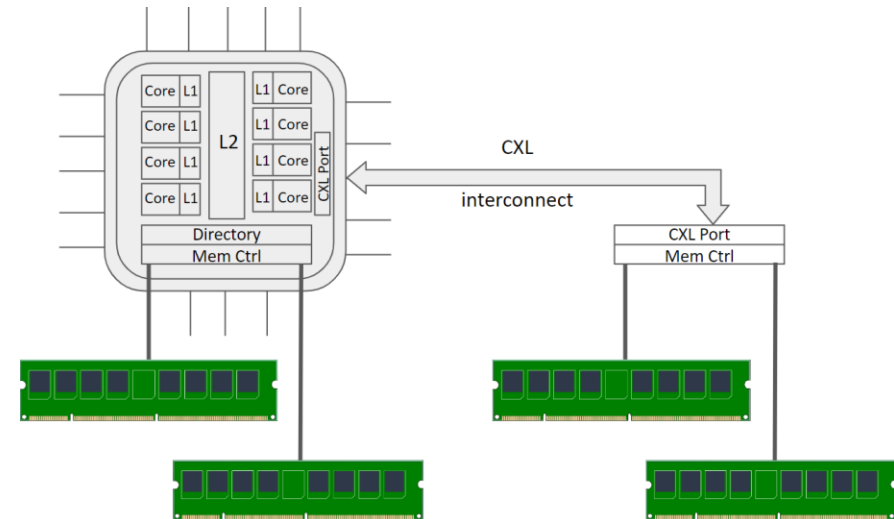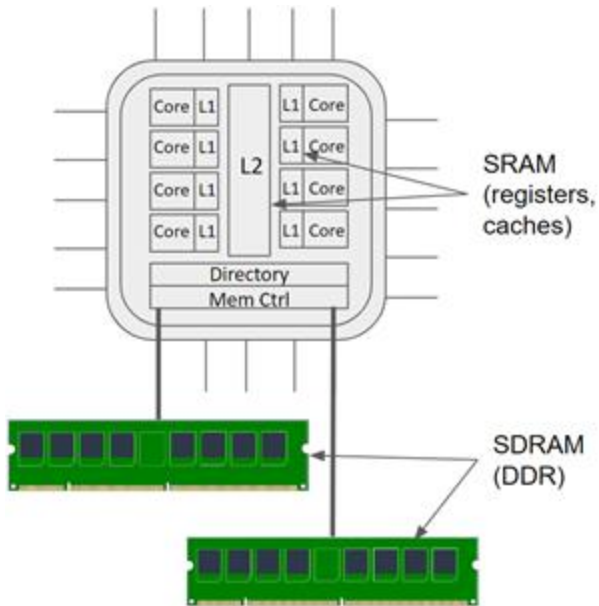Remote DRAM for Socket 1

Remote DRAM for Socket 0

CXL interconnect

CXL Port Mem Ctrl

Disaggregated memory

# Pervasive coherence protocols



SRAM (registers, caches)

SDRAM (DDR)

NUMA

Disaggregated memory

CXL interconnect

# Pervasive coherence protocols



SRAM (registers, caches)

SDRAM (DDR)

**Intra-processor coherence**

NUMA

cache coherent interconnect

Remote DRAM for Socket 1

Remote DRAM for Socket 0

Socket 0

Socket 1

CXL interconnect

Disaggregated memory

# Pervasive coherence protocols

Inter-processor coherence

NUMA

Disaggregated memory

Intra-processor coherence
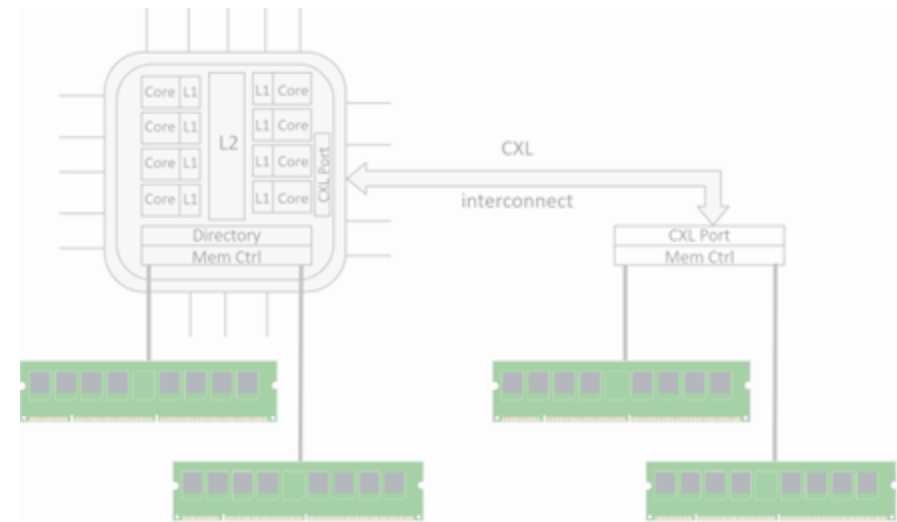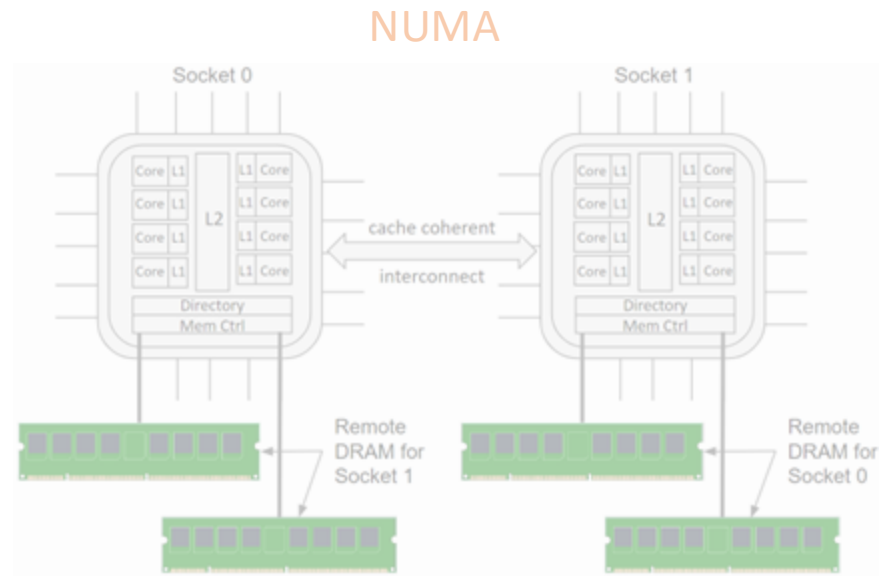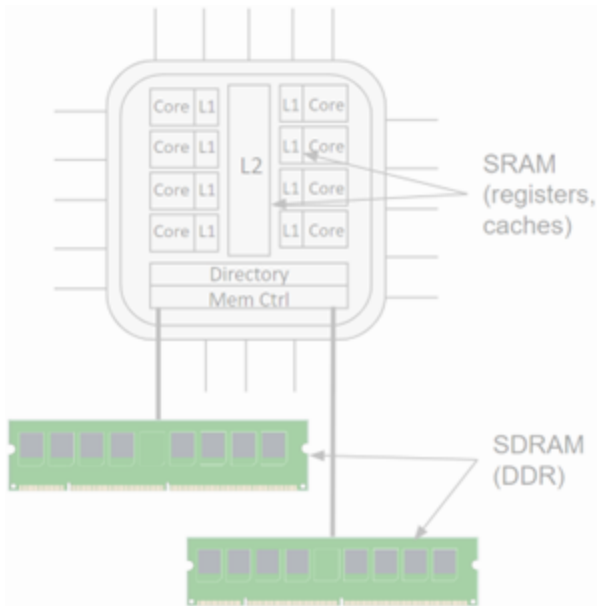
# Pervasive coherence protocols



Intra-processor coherence

Inter-processor coherence

NUMA

Disaggregated memory

Inter-server coherence

# Thesis insights and contributions

**Dvé Memory Replication**

Employ coherence protocols to improve reliability and performance of DRAM memory **[ISCA '21]**

# Thesis insights and contributions

**Dvé**
**Memory Replication**

Employ coherence protocols to improve reliability and performance of DRAM memory **[ISCA '21]**

**Āpta**
**Reliable CXL.mem**

Harden the coherence protocols against common modes of failures **[DSN '23]**

# Thesis insights and contributions



**Improving DRAM Reliability and Performance**

1

# The problem: Increasing DRAM Faults

# Progression of Reliability Mechanisms



Memory
Controller

Channel

DIMM

Rank

Chip

Bank

DRAM
Array

Cell

Row/column sparing

SECDED ECC, in-DRAM ECC

# Progression of Reliability Mechanisms

Memory Controller

Channel — IBM RAIM, Intel Memory Mirroring

DIMM

Rank

Chip — Chipkill ECC / Multi-tier ECC

Bank

DRAM Array — Row/column sparing

Cell — SECDED ECC, in-DRAM ECC

# Progression of Reliability Mechanisms

Memory Controller

Channel — IBM RAIM, Intel Memory Mirroring

DIMM

Rank

Chip — Chipkill ECC / Multi-tier ECC

Bank

DRAM Array — Row/column sparing

Cell — SECDED ECC, in-DRAM ECC

**Performance overheads**

# Replication for Reliability

💡 Dvé insights
- ❑ Full data replica (not ECC code)
- ❑ Keep Replicas as far apart and disjoint as possible
- ❑ Tolerate errors arising from anywhere in the memory path

# Replication for Reliability

Coherent interconnect

X, Y, Z

A, B, C

Dvé insights
- ❑ Full data replica (not ECC code)
- ❑ Keep Replicas as far apart and disjoint as possible
- ❑ Tolerate errors arising from anywhere in the memory path

Coherent interconnect

X, Y, Z

A, B, C

Replicas

A, B, C

X, Y, Z

Dvé insights
- Full data replica (not ECC code)
- Keep Replicas as far apart and disjoint as possible
- Tolerate errors arising from anywhere in the memory path

# Replication for Reliability



Dvé insights
- ❑ Full data replica (not ECC code)
- ❑ Keep Replicas as far apart and disjoint as possible
- ❑ Tolerate errors arising from anywhere in the memory path

For Detection
- ❑ Existing ECC, CRC, Parity
- ❑ Strong detection-only code
- ❑ Other diagnostic capabilities

For Correction
- ❑ Rely on replica

# Replication for Reliability

Coherent interconnect

X, Y, Z

A, B, C

Replicas

A, B, C

X, Y, Z

Memory Controller

Channel

DIMM

Rank

Chip

Bank

Array

Cell

Dvé

Dvé insights
- ❑ Full data replica (not ECC code)
- ❑ Keep Replicas as far apart and disjoint as possible
- ❑ Tolerate errors arising from anywhere in the memory path

For Detection
- ❑ Existing ECC, CRC, Parity
- ❑ Strong detection-only code
- ❑ Other diagnostic capabilities

For Correction
- ❑ Rely on replica

# Coherent Replication for Performance

Dvé insights
- ❑ Use replica to improve performance

# Coherent Replication



A

Coherent interconnect

X, Y, Z

A, B, C

Dvé insights
- ❑ Use replica to improve performance

# Coherent Replication

Coherent

interconnect

**A**

X, Y, Z

A, B, C

Replicas

A, B, C

X, Y, Z

Dvé insights
- ❑ Use replica to improve performance
- ❑ Route memory requests to nearest replica

# Coherent Replication



write(A)

Coherent interconnect

A, B, C

X, Y, Z

A, B, C

X, Y, Z

Replicas

Dvé insights
- ❑ Use replica to improve performance
- ❑ Route memory requests to nearest replica
- ❑ Ensure safe access to replica

# Coherent Replication

Coherent

interconnect

A

Dvé insights
- ❑ Use replica to improve performance
- ❑ Route memory requests to nearest replica
- ❑ Ensure safe access to replica

X, Y, Z

A, B, C

A, B, C

X, Y, Z

Replicas

# Coherent Replication



Coherent interconnect

Replica Dir Ctrl

X, Y, Z

A, B, C

Replicas

A, B, C

X, Y, Z

Dvé insights
- ❑ Use replica to improve performance
- ❑ Route memory requests to nearest replica
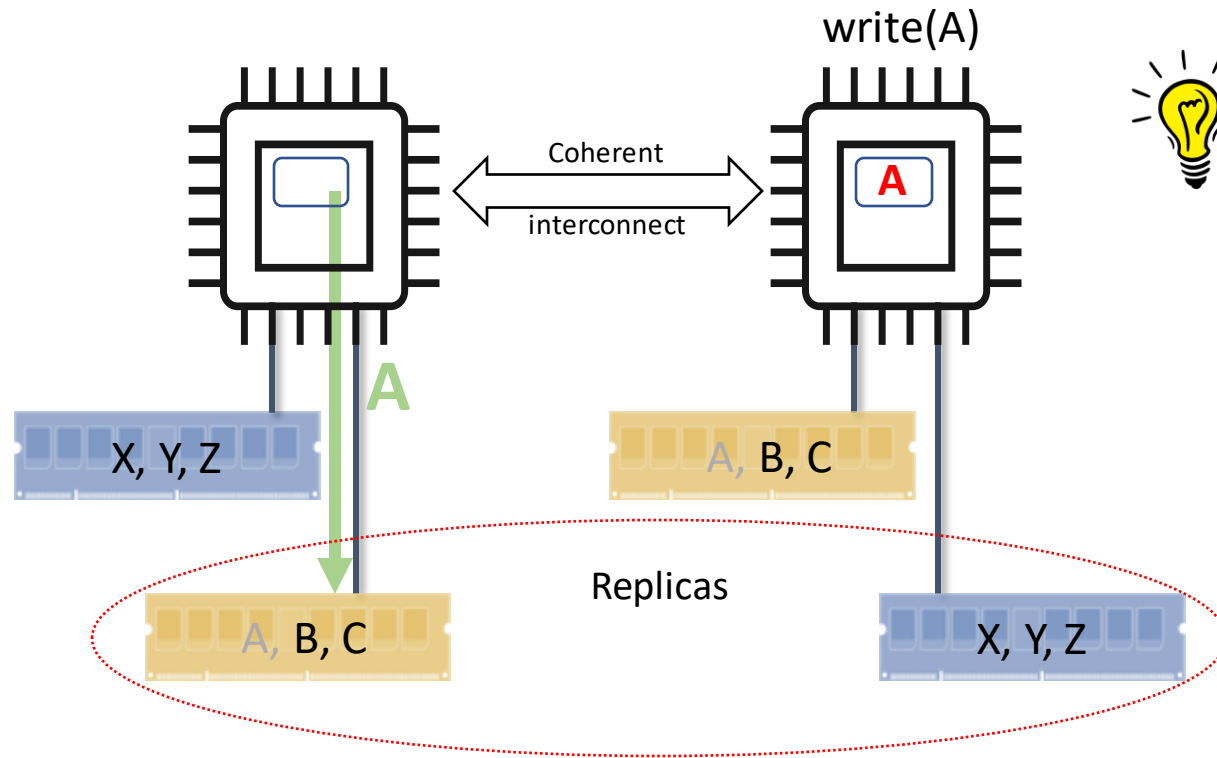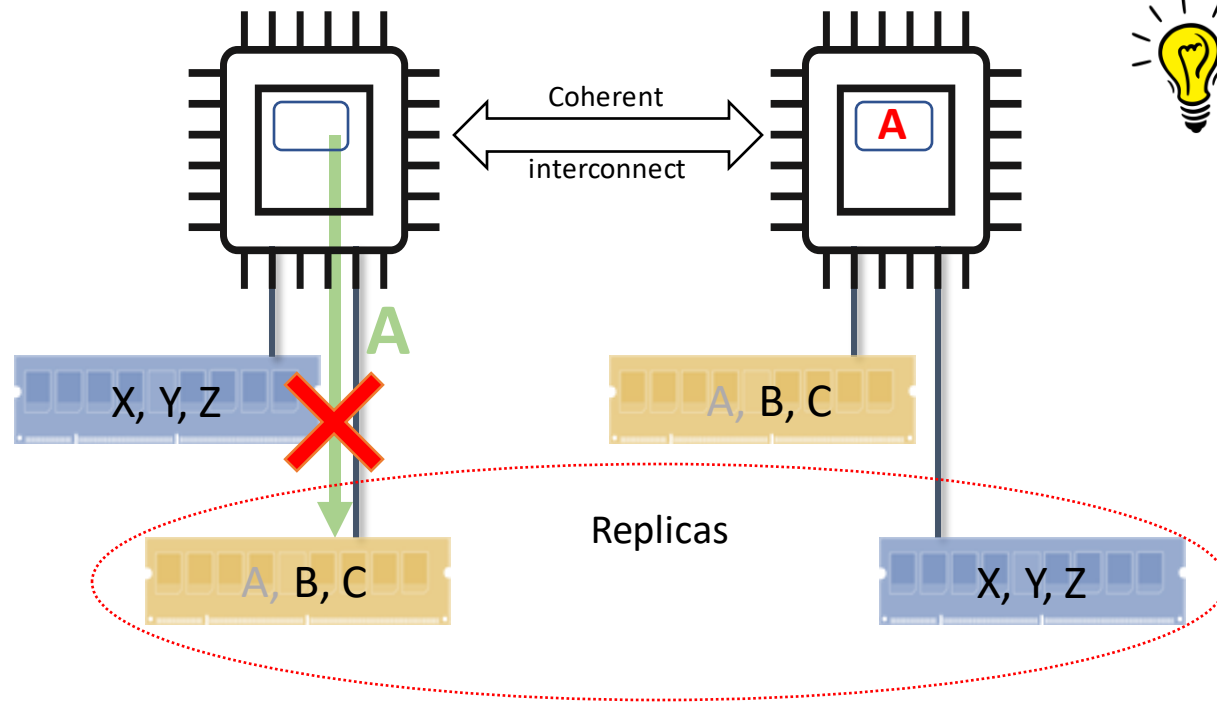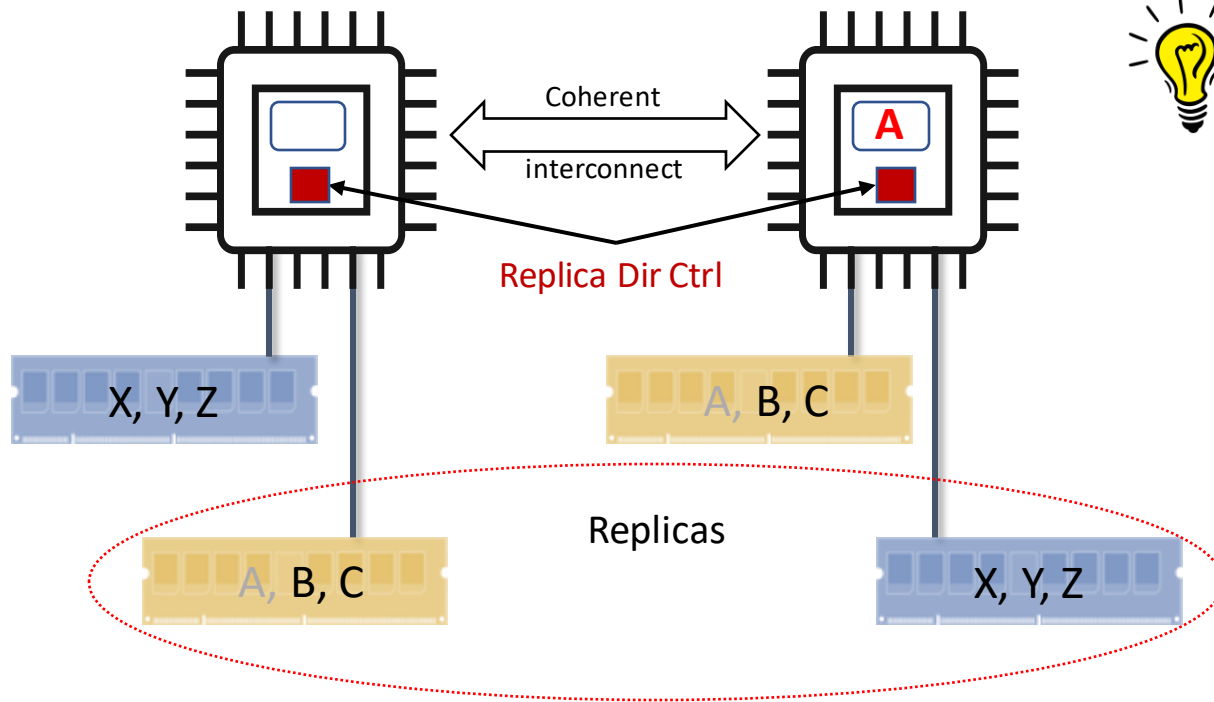- ❑ Ensure safe access to replica

# Coherent Replication



Dvé insights
- Use replica to improve performance
- Route memory requests to nearest replica
- Ensure safe access to replica

Coherent Replication
- Builds on existing cache coherence protocols
- maintain the replicas in sync (for reliability)
- provide coherent access to both replicas during fault-free operation (for performance)

# Coherent Replication



Coherent interconnect

Replica Dir Ctrl

X, Y, Z

A, B, C

Replicas

A, B, C

X, Y, Z

Dvé insights
- Use replica to improve performance
- Route memory requests to nearest replica
- Ensure safe access to replica

Coherent Replication
- Builds on existing cache coherence protocols
- maintain the replicas in sync (for reliability)
- provide coherent access to both replicas during fault-free operation (for performance)

Mechanisms
- Allow-based
- Deny-based

# Capacity overheads?



Reliability
Performance

Capacity

# Capacity overheads?

💡 Dvé insights
- ❑ Utilize idle memory

**Reliability**
**Performance**

**Capacity**

Skewed memory utilization
- ❑ 50% of the memory is idle in 90% of the servers
- ❑ Provisioning for peak

# Capacity overheads?

**Reliability Performance**

**Capacity**

Dvé insights
- ☐ Utilize idle memory
- ☑ Overheads applicable only as and when demanded by the application

Skewed memory utilization
- ☐ 50% of the memory is idle 90% of the servers
- ☐ Provisioning for peak

Interface to allocate high-reliability memory
- ☐ Hardware-software co-design
- ☐ OS support

# Capacity overheads?

Capacity

Reliability
Performance

Dvé insights
- ❏ Utilize idle memory
- ❏ Overheads applicable only as and when demanded by the application

Skewed memory utilization
- ❏ 50% of the memory is idle 90% of the time
- ❏ Provisioning for peak

Interface to allocate high-reliability memory
- ❏ Hardware-software co-design
- ❏ OS support

Flexible trade-off between capacity and reliability

# Summary

**Dvé Memory Replication**

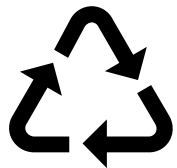### Replication for Reliability

Lowers DUE by

4x over Chipkill
172x over IBM RAIM
11% over Intel Memory Mirroring

### Coherent Replication for Performance

Improves performance by

5% - 117% over baseline NUMA
3% - 107% over an improved
Intel mirroring scheme

### On-demand Replication

hardware-software co-design
using OS/compiler support

### Artifacts available

https://github.com/adarshpatil/dve
https://adar.sh/dve

**Fault-tolerant disaggregated memory for accelerating FaaS**
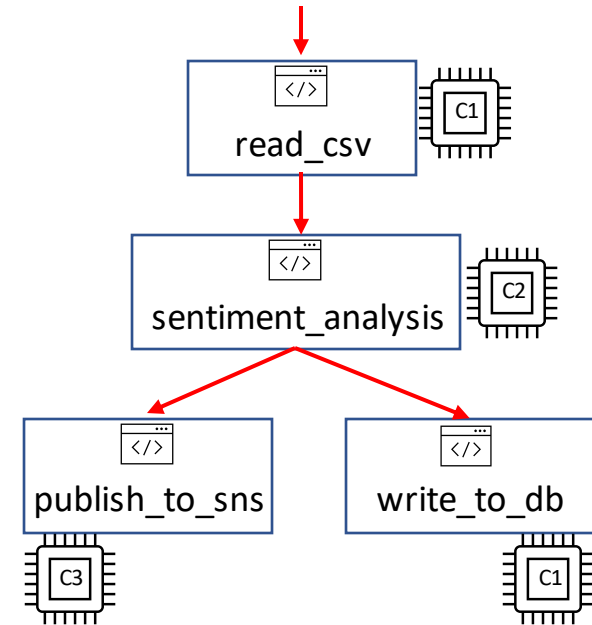
**2**

- State machine workflow of *stateless* functions

# FaaS applications

- State machine workflow of *stateless* functions

- Cloud provider dynamically orchestrates and schedules functions on a fleet of compute servers

- State machine workflow of *stateless* functions

- Cloud provider dynamically orchestrates and schedules functions on a fleet of compute servers

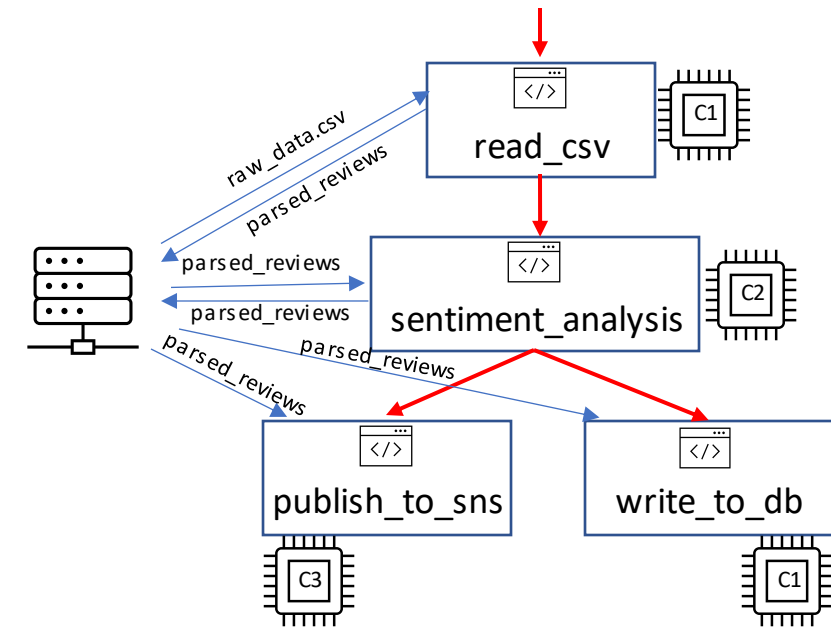- **State maintained externally as objects in a remote data store**

- State machine workflow of *stateless* functions

- Cloud provider dynamically orchestrates and schedules functions on a fleet of compute servers

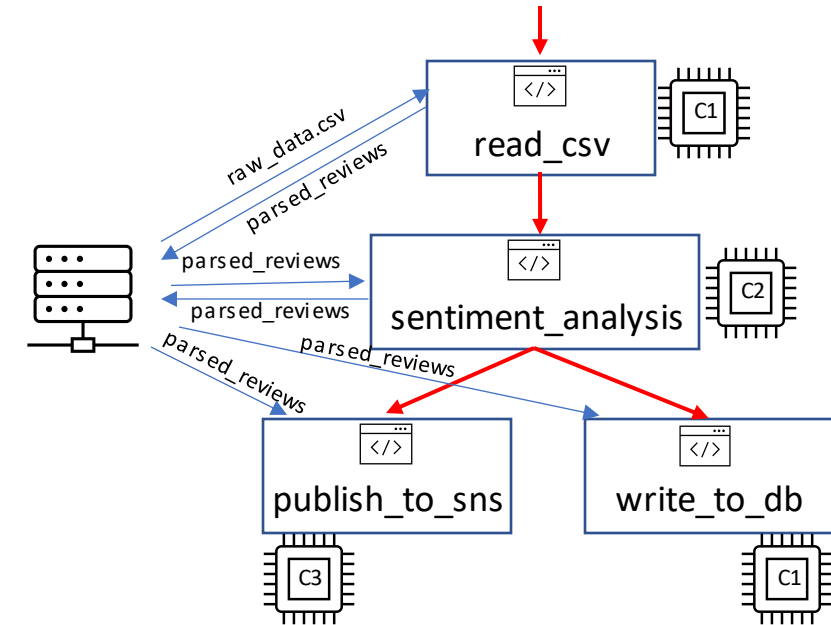- **State maintained externally as objects in a remote data store**
  - × Splitting state-compute adds communication overheads
    - How much?

# Quantifying communication overheads

- Functions from FunctionBench and SeBS benchmark suites

- Compute - Optimized with Intel OneAPI, run on 16-core Skylake CPUs

- Communication - Amazon S3 object store (median of 100 executions)



96% of execution time is spent in retrieving data from S3

- State machine workflow of *stateless* functions

- Cloud provider dynamically orchestrates and schedules functions on a fleet of compute servers

- State maintained externally as objects in a remote data store

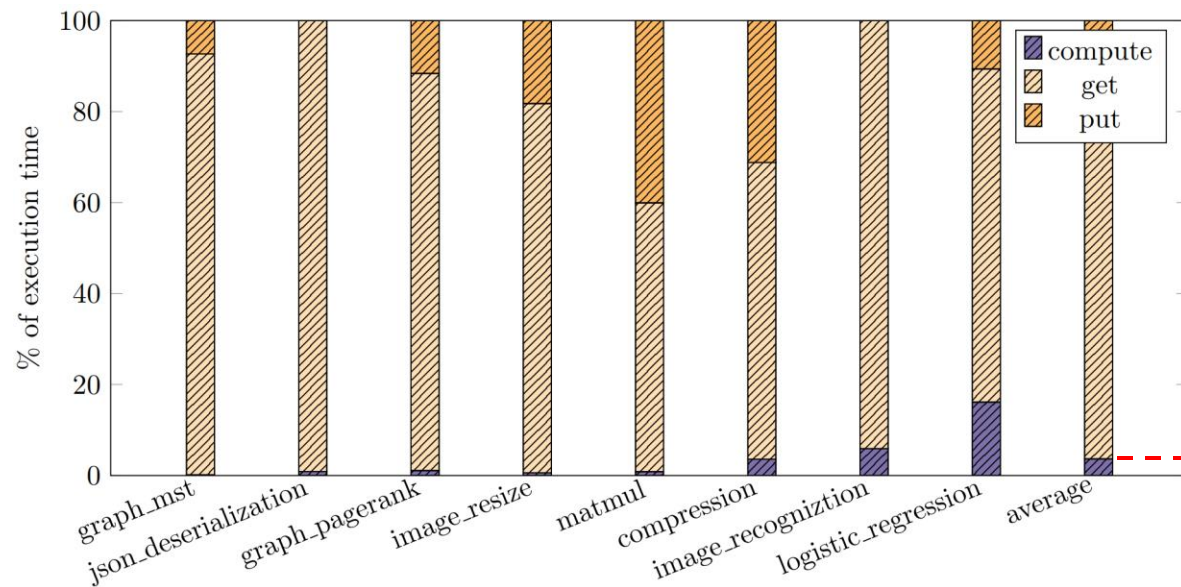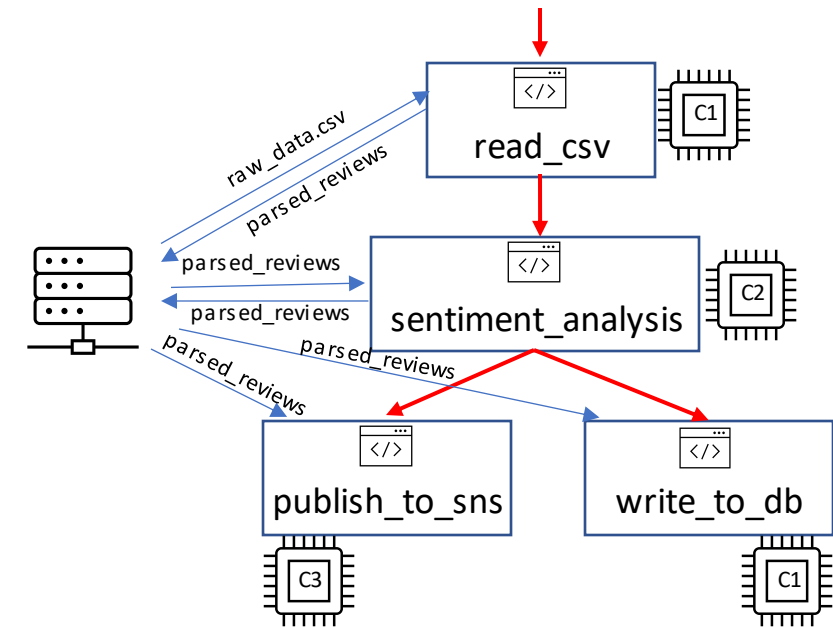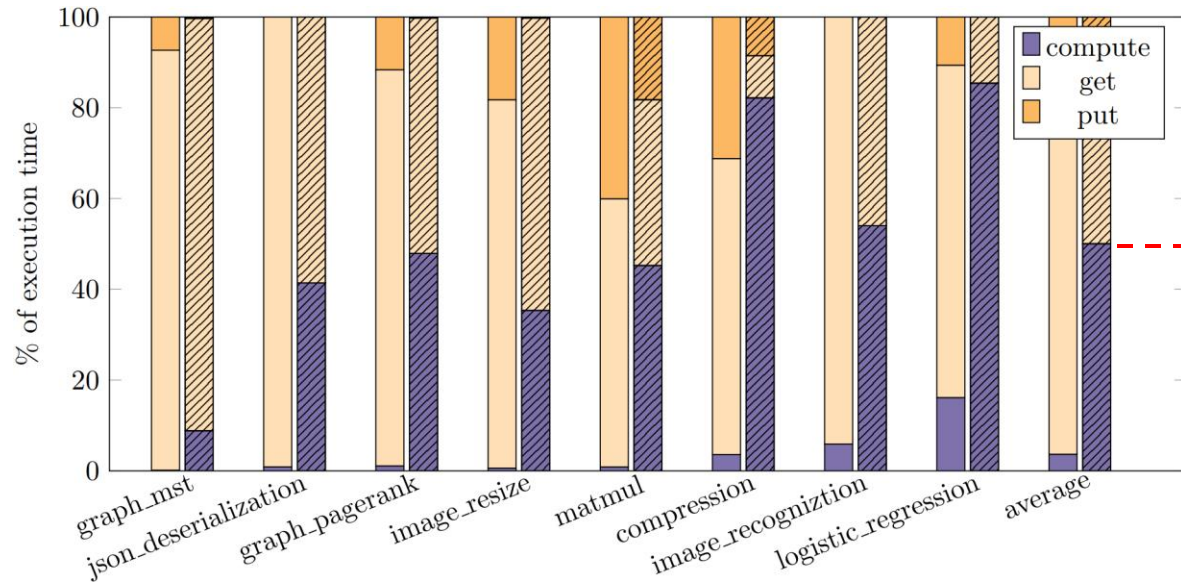  × Splitting state-compute adds communication overheads

× Communication overheads severely limit performance

Can we do better?

- High-performance in-memory object store

- One-sided RDMA verbs to read/write objects

- Infiniband network (Mellanox ConnectX-3 NIC on PCIe-gen3 x16)



51% of execution time is spent in retrieving data from object store

# The problem: Data communication



"The two most expensive operations in terms of cost were the orchestration workflow and when data passed between distributed components."

# Performance potential for Āpta

Object-granular CXL disaggregated memory

# Performance potential for Āpta
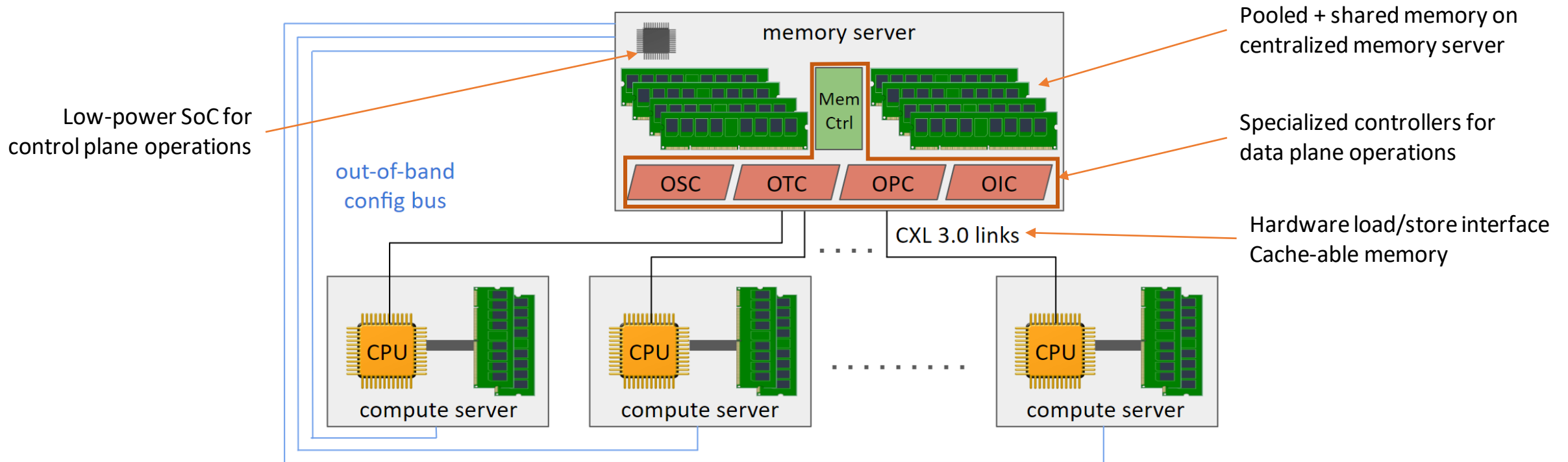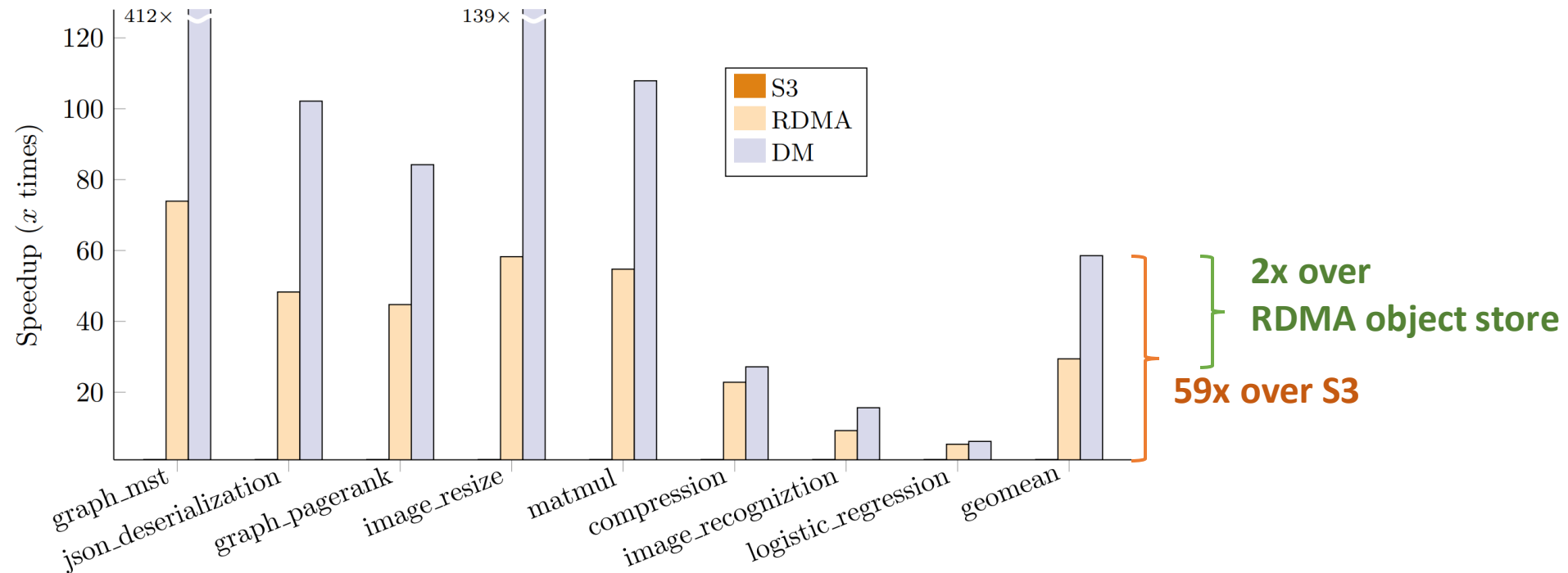
## Object-granular CXL disaggregated memory



Pooled + shared memory on centralized memory server

Low-power SoC for control plane operations

Specialized controllers for data plane operations

Hardware load/store interface Cache-able memory

Fig: Āpta system schematic

# Performance potential for Āpta

With OpenCAPI-like access latency / bandwidth for DM[†]
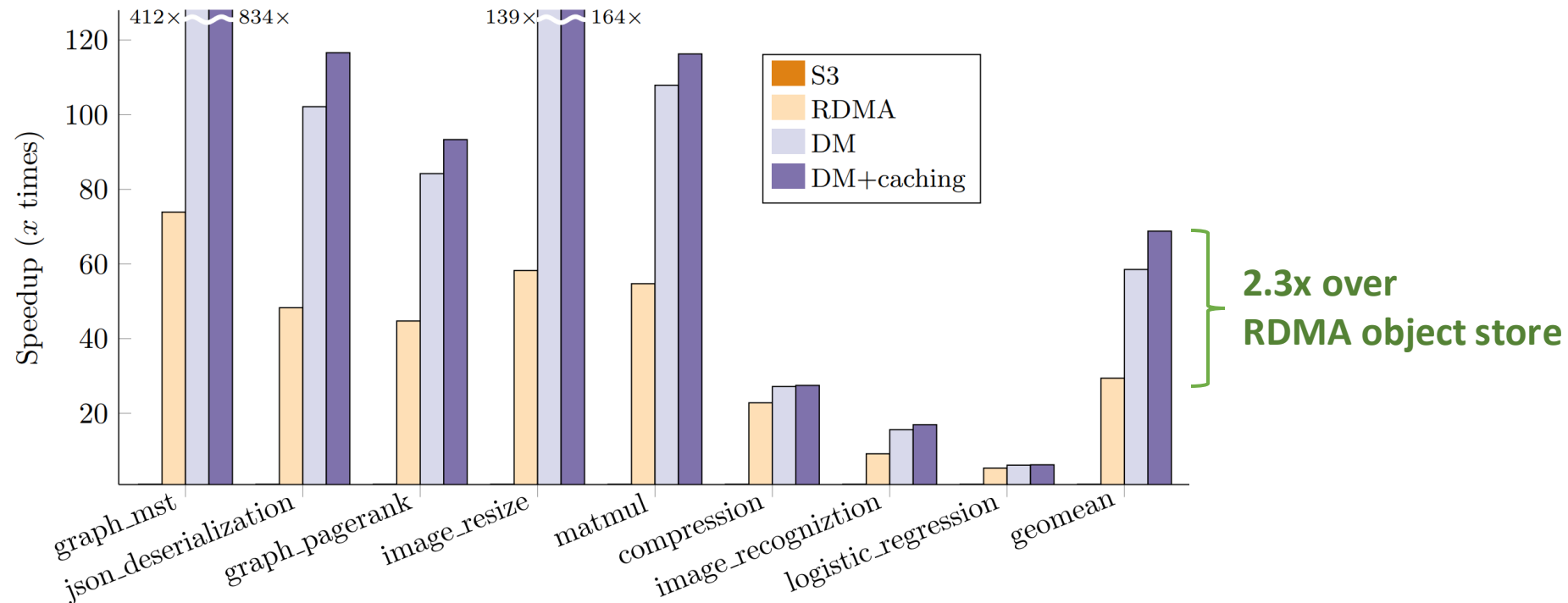
2x over RDMA object store

59x over S3

13% communication overheads (Recall 51% for RDMA-based object store)
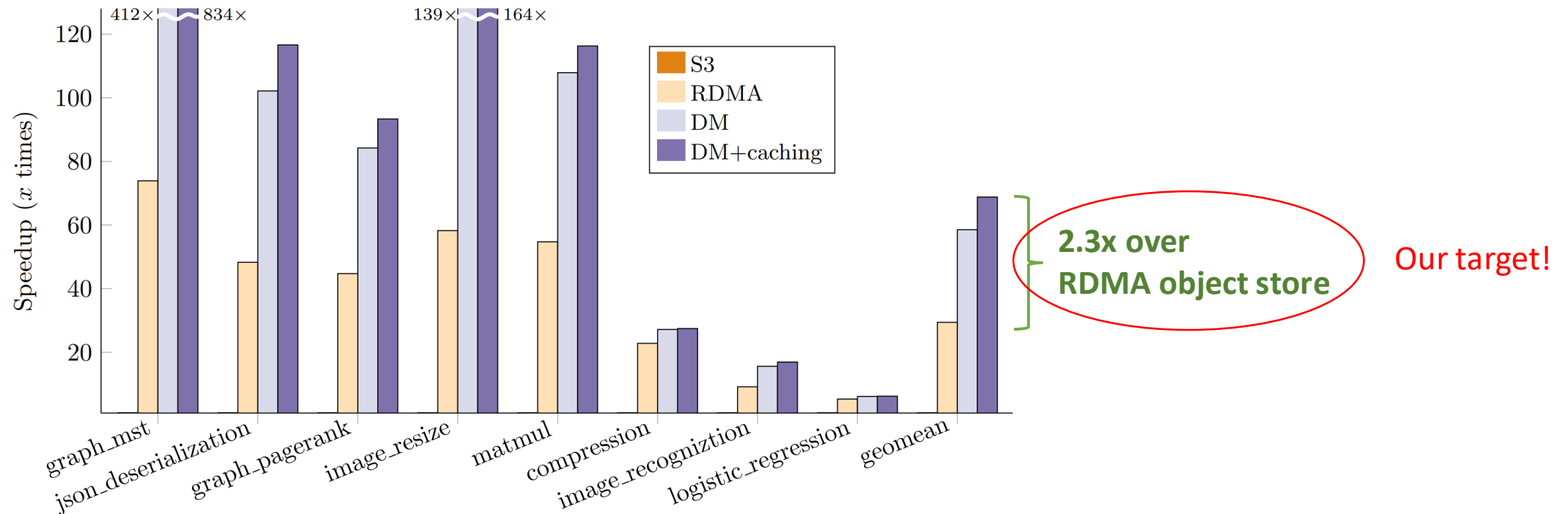
[†] ThymesisFlow [MICRO 20]

# Performance potential for Āpta

Object caching at compute server

# Performance potential for Āpta
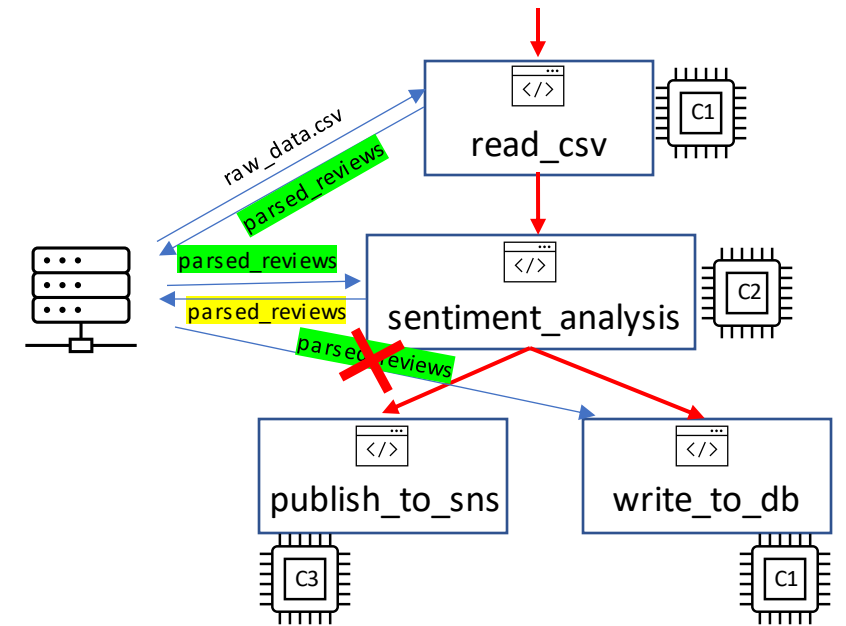
## Object caching at compute server

- Enforcing strong consistency in presence of caching

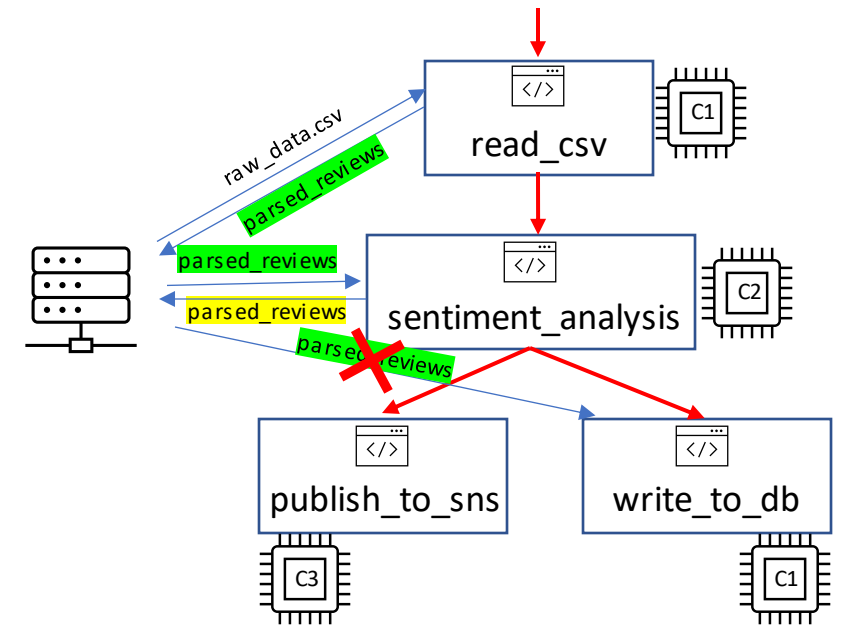- Enforcing strong consistency in presence of caching

- Enforcing strong consistency in presence of caching

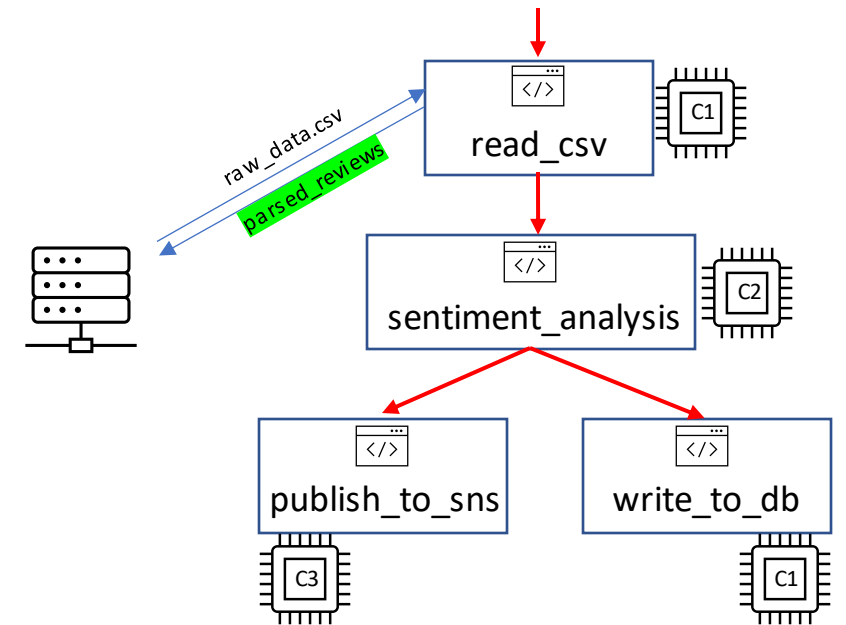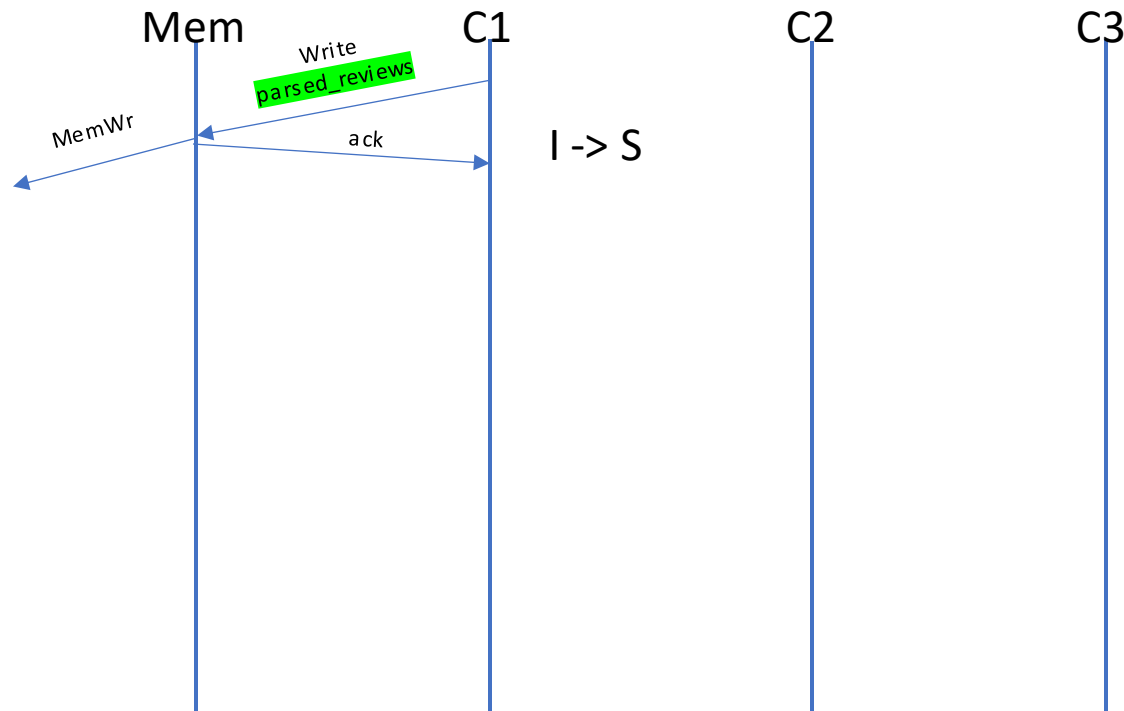    CXL 3.0 inter-node coherence protocol
    Enforces SWMR invariant

# The CXL.mem coherence

- Enforcing strong consistency in presence of caching
  - CXL 3.0 Inter-node coherence protocol
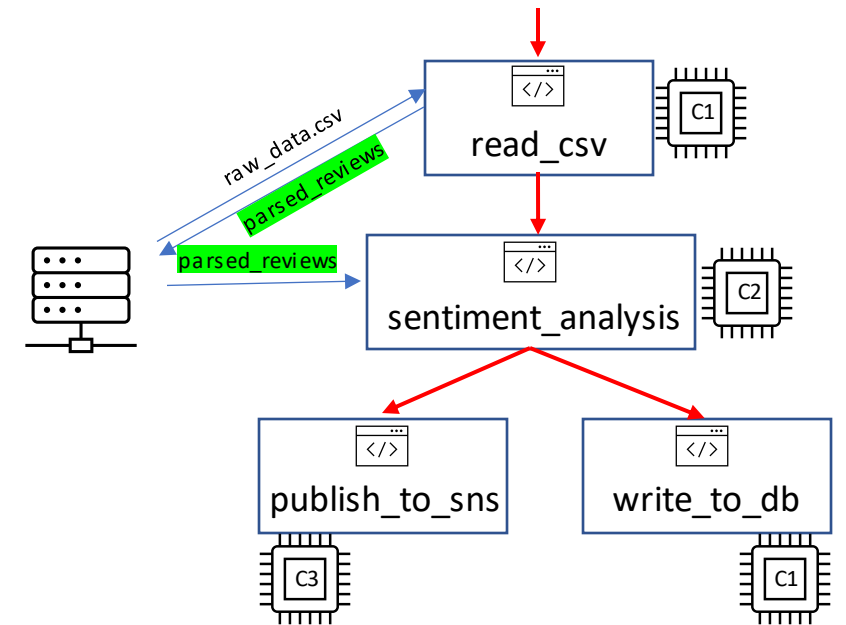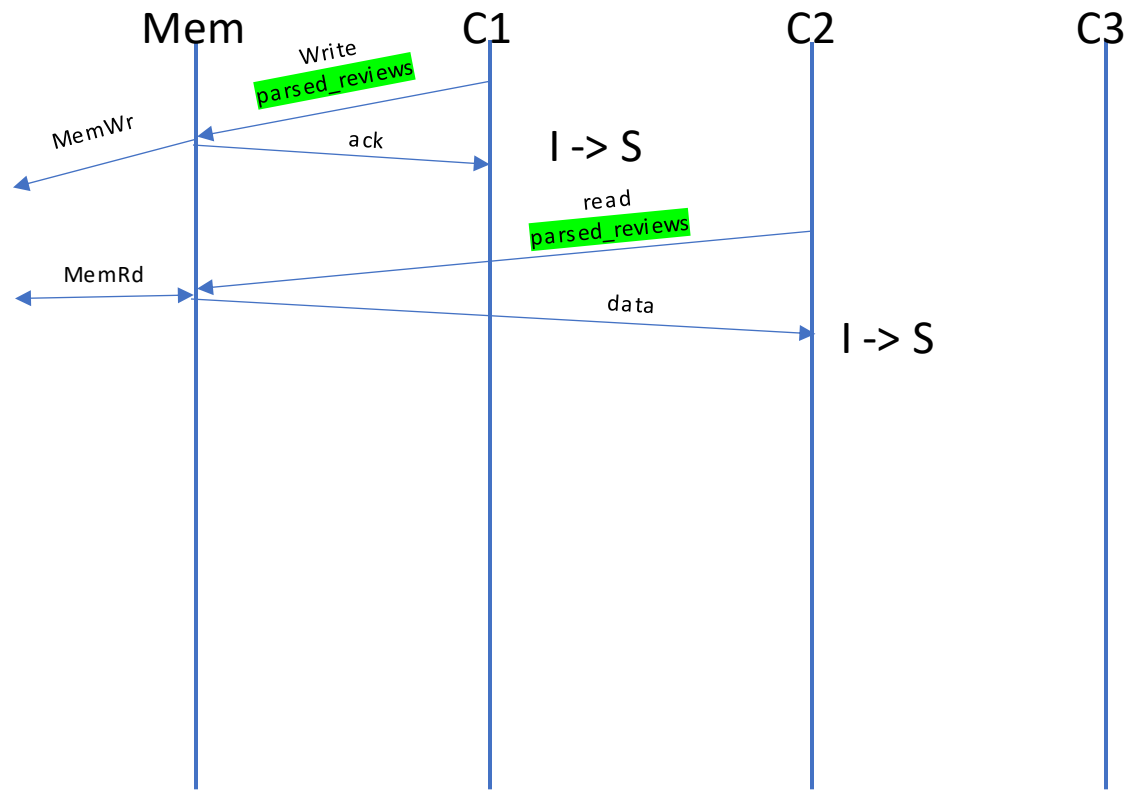  - Enforces SWMR invariant

- Enforcing strong consistency in presence of caching
  Inter-node coherence protocol
  CXL 3.0 protocols enforce SWMR invariant

Mem    C1        C2        C3

Write
parsed_reviews

MemWr          ack          I -> S

read
parsed_reviews

MemRd          data          I -> S

raw_data.csv

parsed_reviews

read_csv    C1

parsed_reviews

parsed_reviews

sentiment_analysis    C2

publish_to_sns        write_to_db

C3                    C1

- Enforcing strong consistency in presence of caching
  - Inter-node coherence protocol
  - CXL 3.0 enforce SWMR invariant

- Enforcing strong consistency in presence of caching
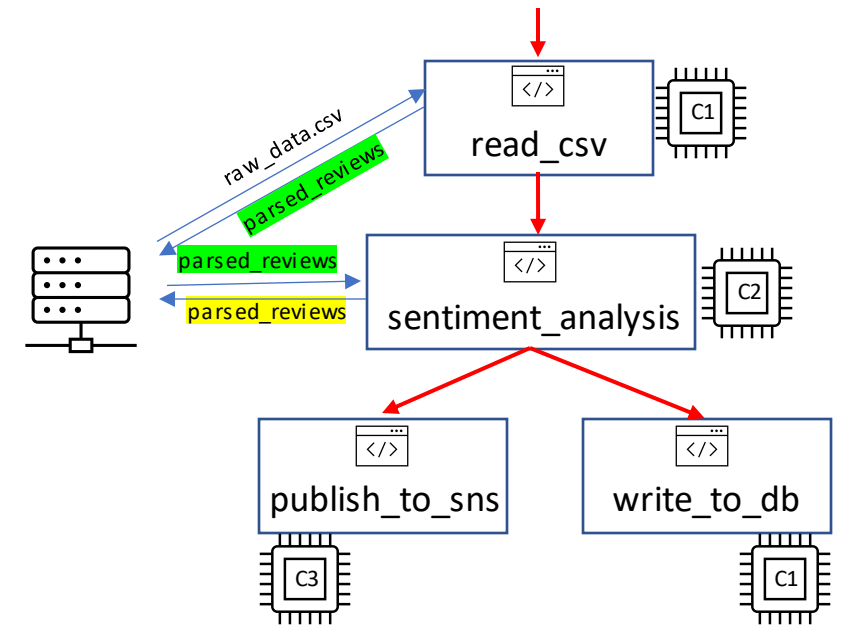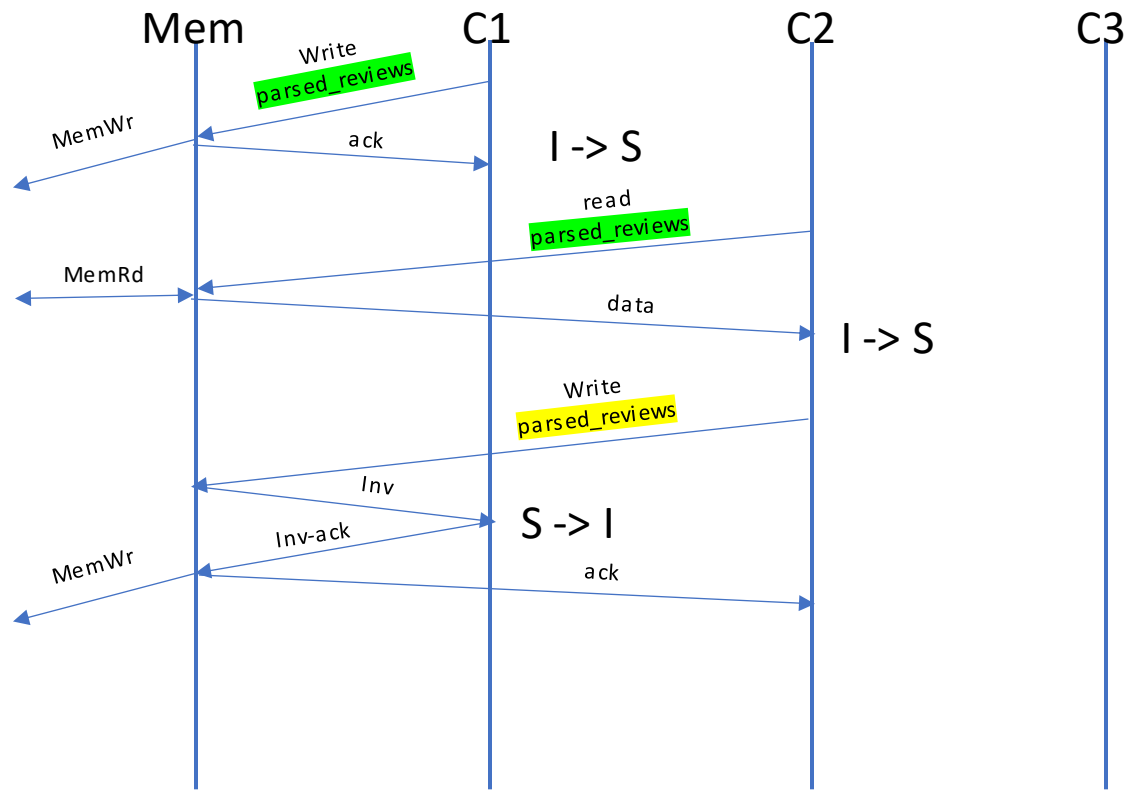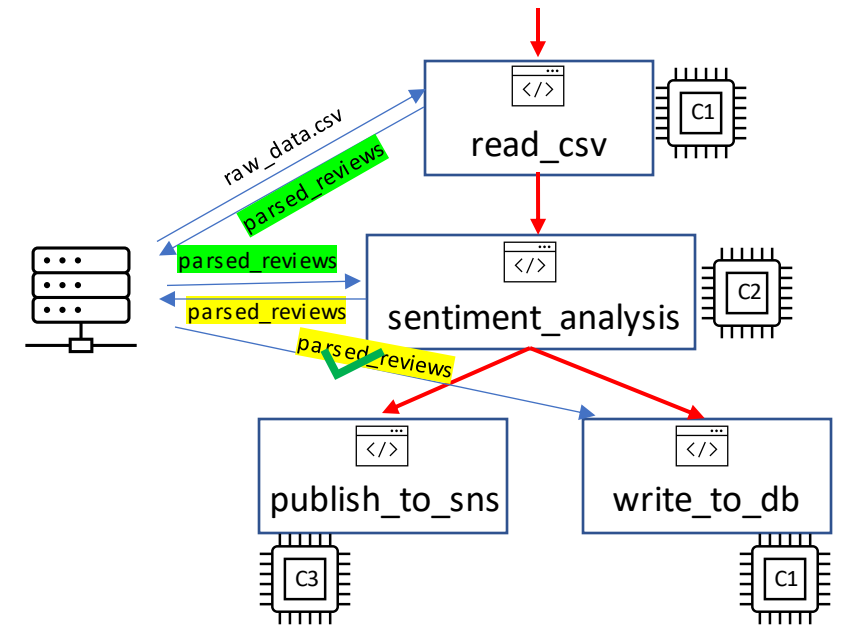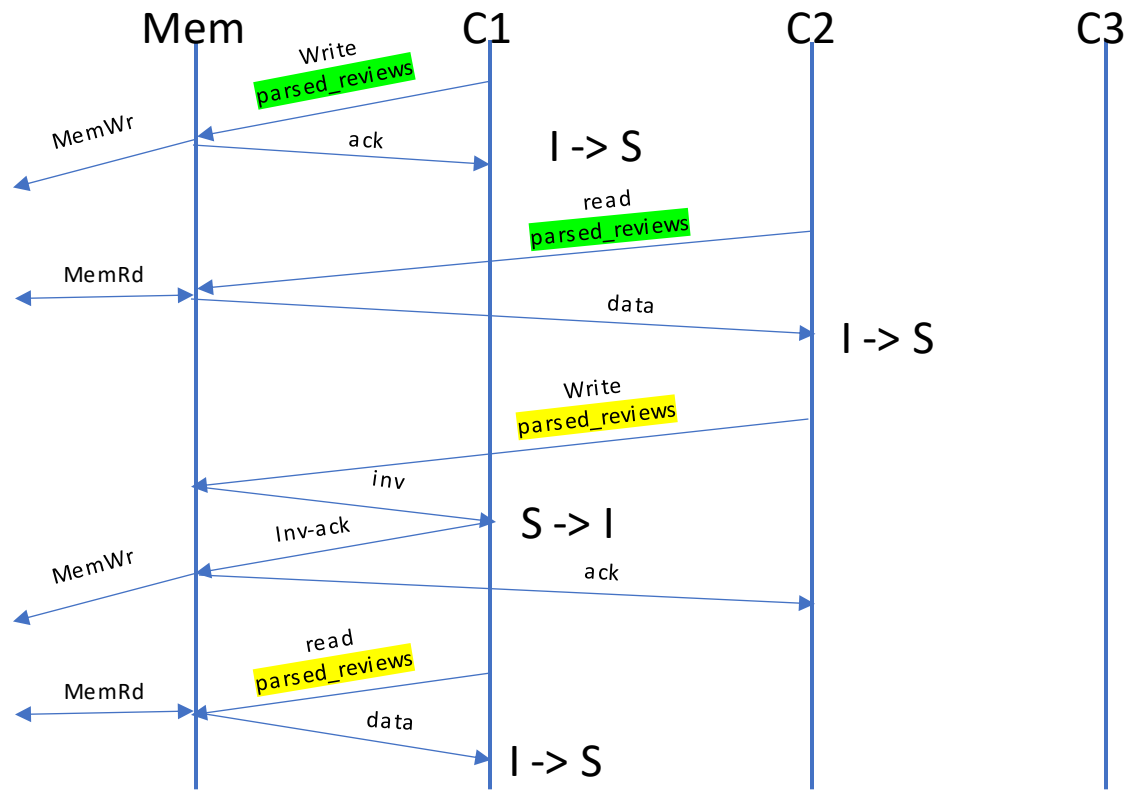  - Inter-node coherence protocol
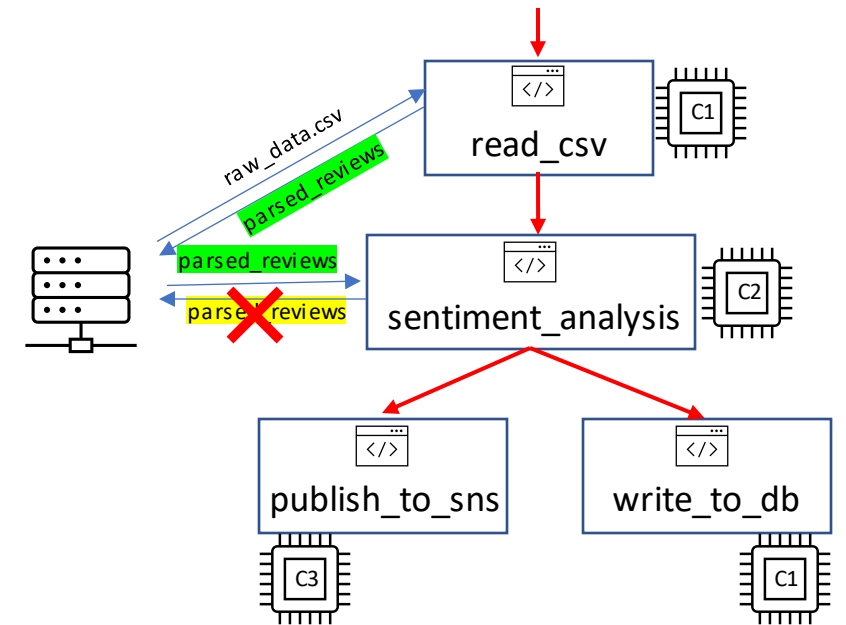  - CXL 3.0 enforce SWMR invariant

# The need for fault-tolerant coherence

- The fault tolerance problem
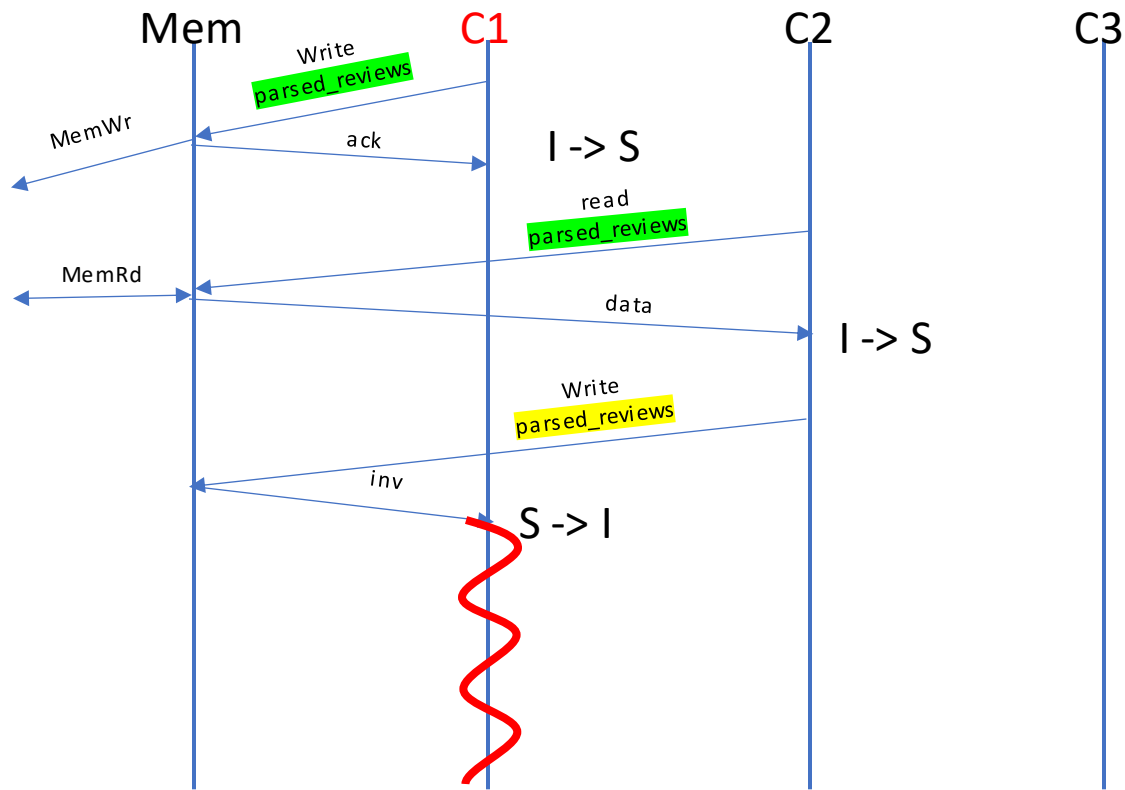  Compute server failures

- The fault tolerance problem

Compute server failures

# The need for fault-tolerant coherence

- The fault tolerance problem

  Compute server failures – blocking

- The fault tolerance problem

  Compute server failures – blocking
  Network congestions – high tail latency

# Key Problem

- Invalidations are in the critical path of a write

# Āpta: Fault-tolerant Coherence Protocol

i. Lazy invalidation policy

Write is acknowledged immediately

Invalidation messages are sent asynchronously and tracked

ii. Coherence-aware function scheduling

Never schedules function invocations on servers with pending invalidation-acknowlegements

Lazy linearizability (Lazy invalidation policy + Coherence-aware scheduling)

# Āpta: Fault-tolerant Coherence Protocol

Lazy linearizability (Lazy invalidation policy + Coherence-aware scheduling)

Ensures compute server fault-tolerant operation

# Āpta: Fault-tolerant Coherence Protocol

Lazy linearizability (Lazy invalidation policy + Coherence-aware scheduling)

Ensures compute server fault-tolerant operation

Provides line-rate coherence

# Āpta Summary

**Accelerating function-as-a-service**

Improves performance by
40% - 142% over RDMA
21% – 90% over RDMA + caching
15% - 42% over un-cached CXL

**Fault-tolerant coherence protocol**

Protocol verified in Murφ model checker
32% lower standard deviation of exec time

**Object-granular Disaggregated Memory**

CXL-based shared memory IPC
Bulk cache-line loads
Transaction atomic durability

**Artifacts available**

https://github.com/adarshpatil/apta
https://adar.sh/apta

Āpta
Reliable CXL.mem

# Summary: Thesis contributions

**Dvé**
**Memory Replication**

- Unique design point in the reliability design space
- Explored a novel extrapolation of two-tier approach
- Introduced flexible / on-demand reliability

**Āpta**
**Reliable CXL.mem**

- Showcases a use case for shared disaggregated memory
- Proposes a lightweight fault tolerance solution
- Consistency & availability via fault-tolerant coherence

# Summary: Retrospective contemplation

| Critical analysis |
|---|
| *Software complexity*: OS, scheduler |
| *Problems of scale*: throughput, co-location |
| *Performance corner cases*: worst-case scenarios |

| Lessons learnt |
|---|
| Mental model of correctness during development |
| Think and reason from first principles |

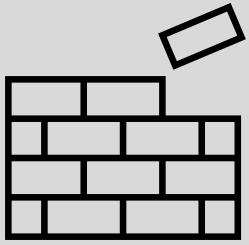| Takeaways |
|---|
| *Robust* reliability is key for next gen memory<br>• technology agnostic, demand reliability (DDR, LPDDR, GDDR)<br>• hardware disaggregated memory (new fault models) |
| Application *driven* architecture<br>• Hardware fault-tolerance must match application evolution<br>• Good understanding of application characteristics |
| *Revisit* design decisions in-step with advances in technology<br>• shared memory systems today are more closely resembling traditional distributed systems |
| End-to-end argument to system design [Saltzer, 1984] |
| Tame complexity through modularization |

# Future Research Directions

- Value-added disaggregated memory

  **Reliability, Availability**, Security, Compression….

- Redesigning distributed datacenter co-ordination services for modern hardware

  **Kubernetes (scheduler)**, Chubby (locks), Kafka (configuration)….

- Efficient shared disaggregated memory

  Heterogenous compute, consistency-directed coherence mechanisms

# Future Research Directions

- Value-added disaggregated memory
    **Reliability, Availability**, Security, Compression….

- Redesigning distributed datacenter co-ordination services for modern hardware
    **Kubernetes (scheduler)**, Chubby (locks), Kafka (configuration)….

- Efficient shared disaggregated memory
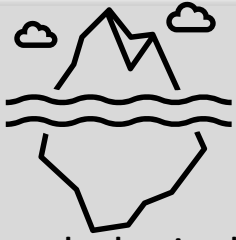    Heterogenous compute, consistency-directed coherence mechanisms

What's next?  #OpenToWork – Industry Research positions
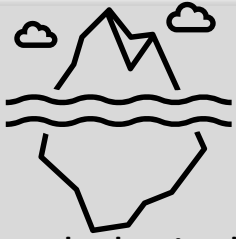Current: Post-doc at University of Edinburgh

Enjoyed… the journey
- Going with a hunch, high-level problem– usually after discussions with Vijay
- Reading related work critically
- Designing experiments to demonstrate the problem (motivation)
- Inception of a workable solution
- Proof of viability – pen/paper, creative descriptions
- Designing experiments to demonstrate the solution
- Refining the idea & solution
- Putting it all together: Presenting, writing the problem, solution, pros/cons

Disliked…
- Convincing reviewers
- The journey of solitude – the imposter syndrome, the large gaps

# What I…

Would do differently (with benefit of hindsight)….

- Better evaluation techniques
  - e.g., Learn HDL, try FPGA-based prototype
- Better paper positioning for maximum impact
  - e.g., Apta is an intersection of KV store + FaaS performance + CXL protocol

- Time-management: better context switching between projects
- Dealing with rejection (still not mastered this)
- Prioritize mental wellbeing: self-reassurance, self-belief, avoid comparing