

Dependable Distributed Shared Memory Suitable for Large, Strongly Consistent Objects

Candidate

Andria Trigeorgi

Supervisor

Chryssis Georgiou

PhD Defence

in the Department of Computer Science

University of Cyprus

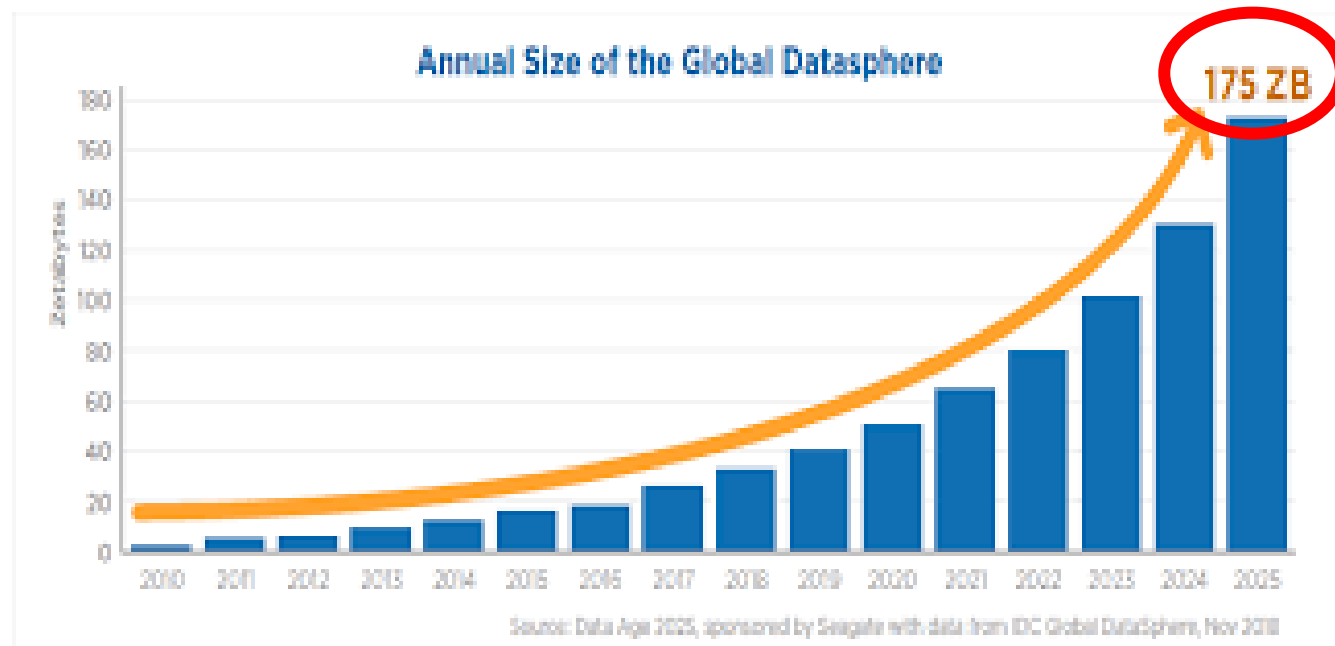


RESEARCH
& INNOVATION
FOUNDATION



University
of Cyprus

Growth of the Digital Universe



“Growth of global data ,
rising from 26 ZB in 2017 to a
projected 175 ZB by 2025.”

- International Data
Corporation (IDC)

“The data in the digital universe
doubles every two years.”

- EMC Digital Universe

- 1 ZB is 10^{21} bytes = the estimated number of stars in the universe!

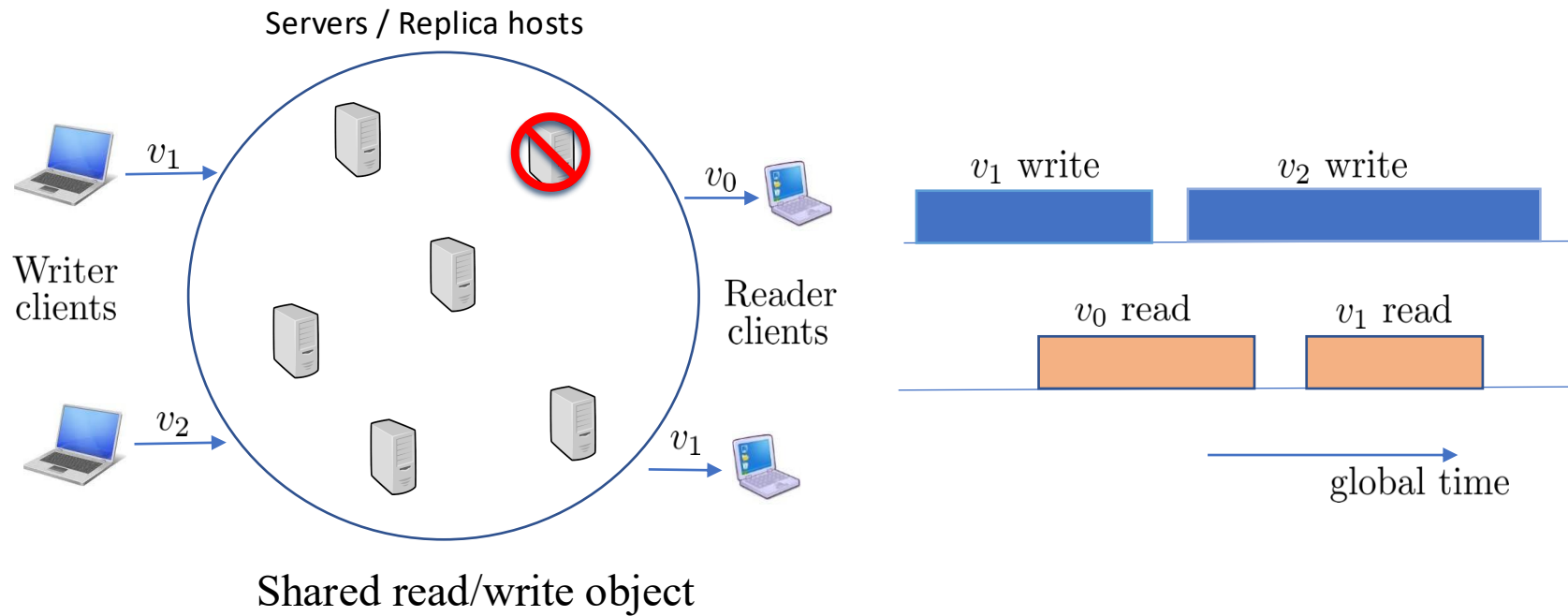
Distributed Storage Systems (DSS)

Commercial Solutions:



- Often provide **Weak Consistency** guarantees (e.g., Dropbox)
- Limited **concurrency** (e.g., HDFS - one writer at a time)
- Often rely on **centralized solutions** to provide strong consistency (e.g., HDFS)
 - Drawback: Performance Bottleneck
- They do not provide rigorously *provable guarantees*

Distributed Shared Memory Emulations (DSMs)



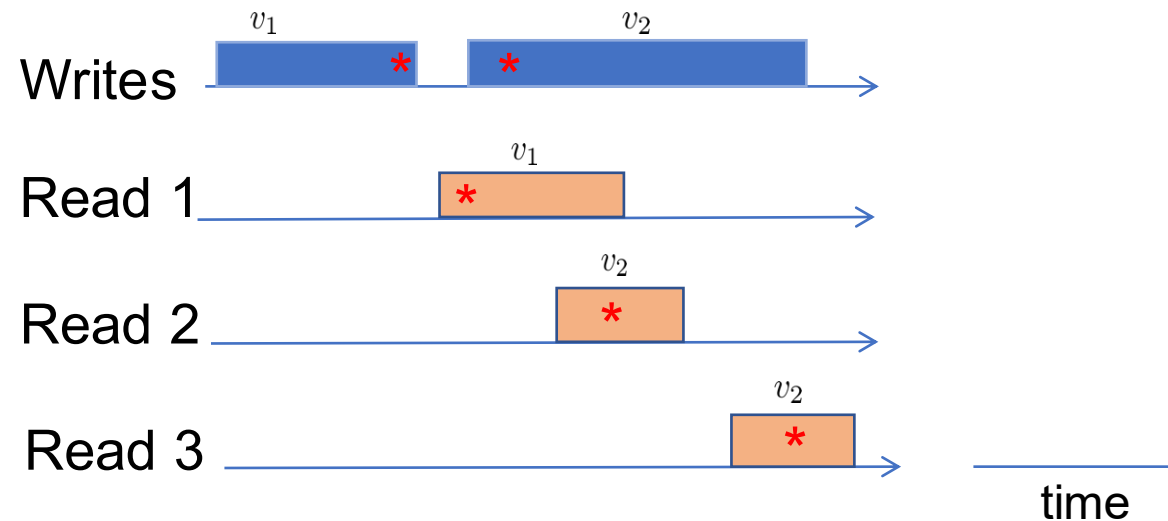
Implementing a **fault-tolerant** shared object in an **asynchronous, message-passing** environment:

- Availability + Survivability => **use redundancy**
- Asynchrony + Redundancy => **concurrent operations**
- Behavior of concurrent operations => **consistency semantics**
 - Safety, Regularity, **Atomicity** (Atomic DSMs) [Lamport 1986]

Atomicity/Linearizability

- Provides the **illusion** that operations happen in a **sequential order**
 - a read returns the **value of the preceding write**
 - a read returns a **value at least as recent as** that returned by **any preceding read**

[Herlihy, Wing 1990]



Seminal Algorithm - ABD

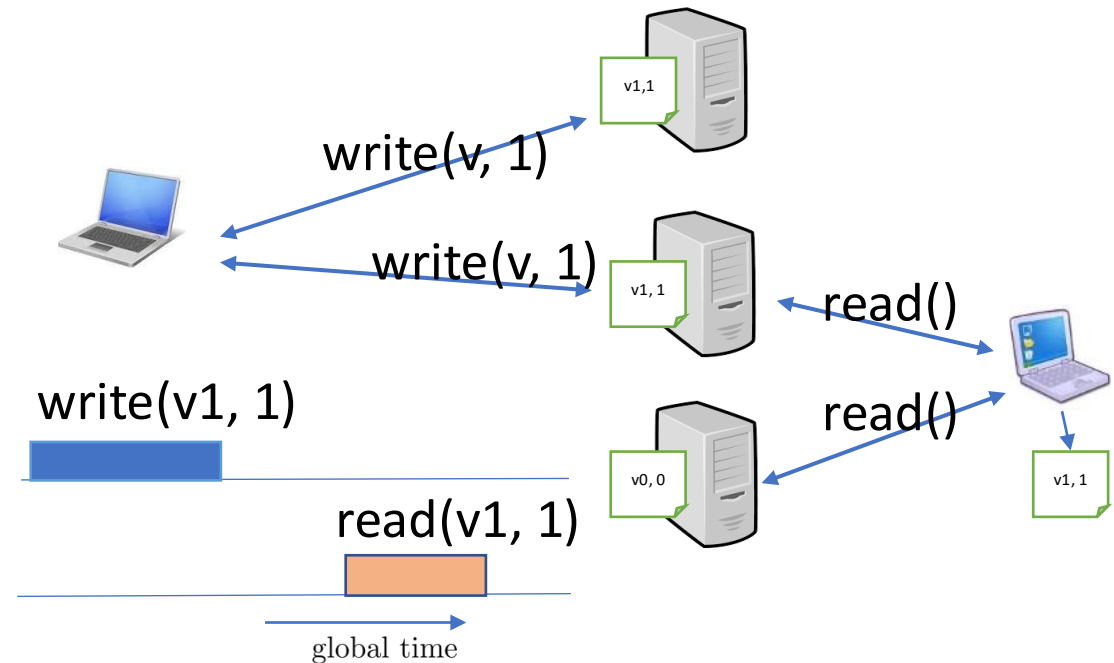
An elegant, intuitive solution that

- uses the power of the **majority**, and
- assigns **logical timestamps** to written values for ordering the operations.

[Attiya, Bar-Noy, Dolev 1995]

Dijkstra Prize 2011

- SWMR atomic registers
- S servers, $f < n / 2$
- 1 writer
- R readers



Extended by Lynch and Schwarzmann

in 1997 for **MWMR**, assigning tags $\langle ts, wid \rangle$ – **MW-ABD**

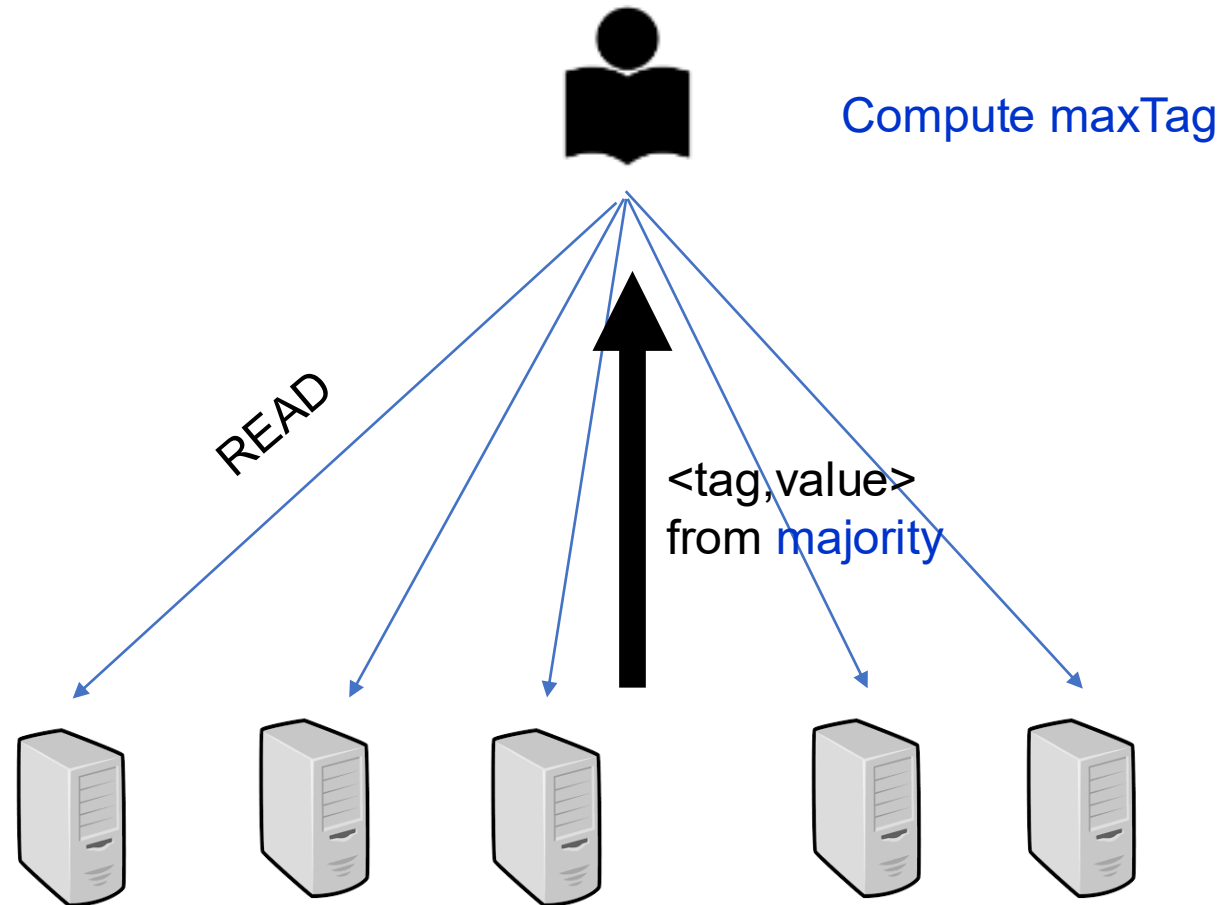
Fault tolerance & Consistency Guarantees have been rigorously proved

H. Attiya, A. Bar-Noy, and D. Dolev, "Sharing Memory Robustly in Message-Passing Systems," Journal of the ACM (JACM), vol. 42, no. 1, pp. 124–142, 1995.

N. Lynch, A. Shvartsman.. "Robust emulation of shared memory using dynamic quorum-acknowledged broadcasts," In Proc. of FTCS pp. 272–281 (1997).

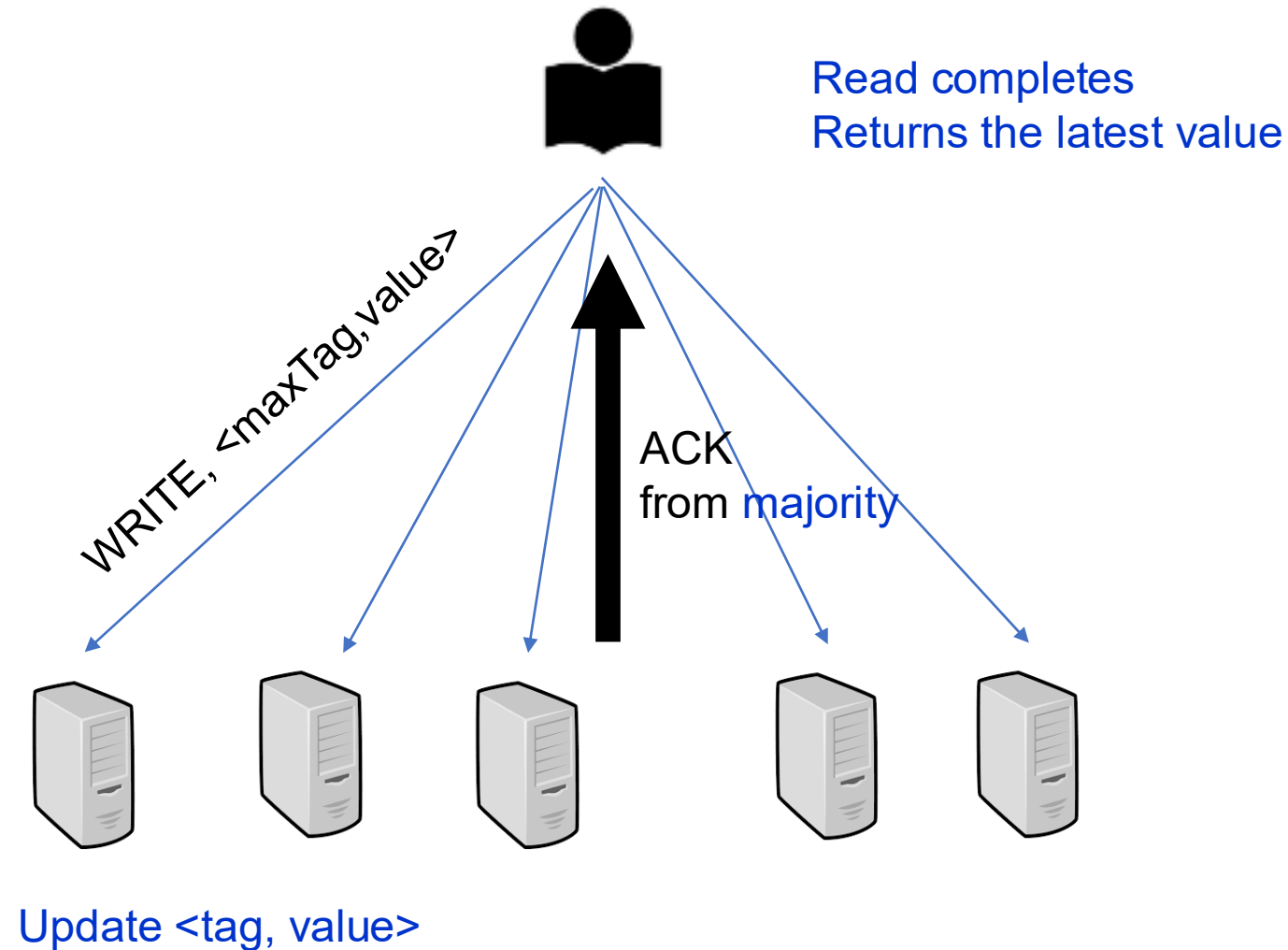
READ Protocol of MW-ABD

Phase1: **Query** - Discover maximum tag and associated value



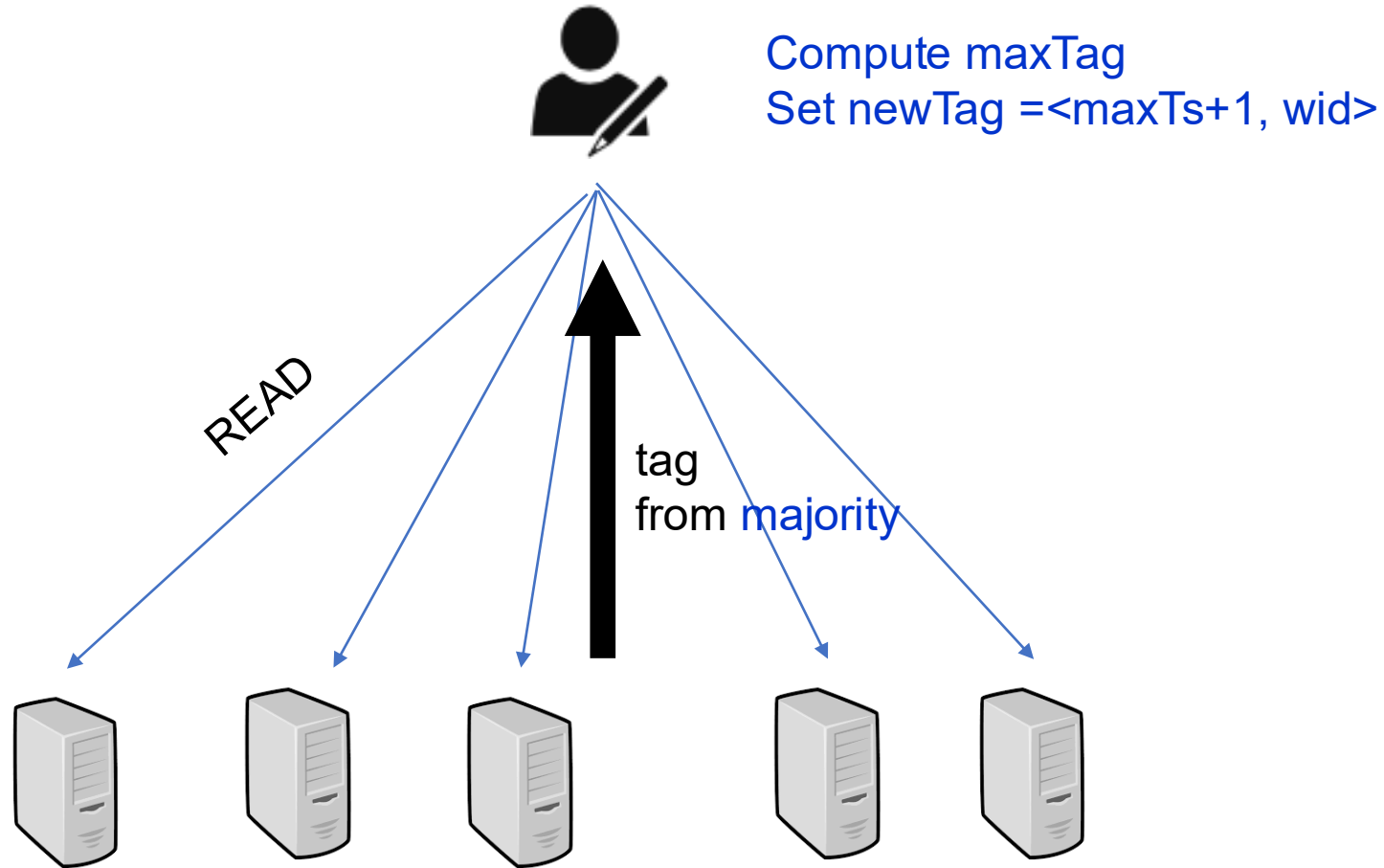
READ Protocol of MW-ABD

Phase 2: **Propagate** $\langle \text{maxTag}, \text{value} \rangle$



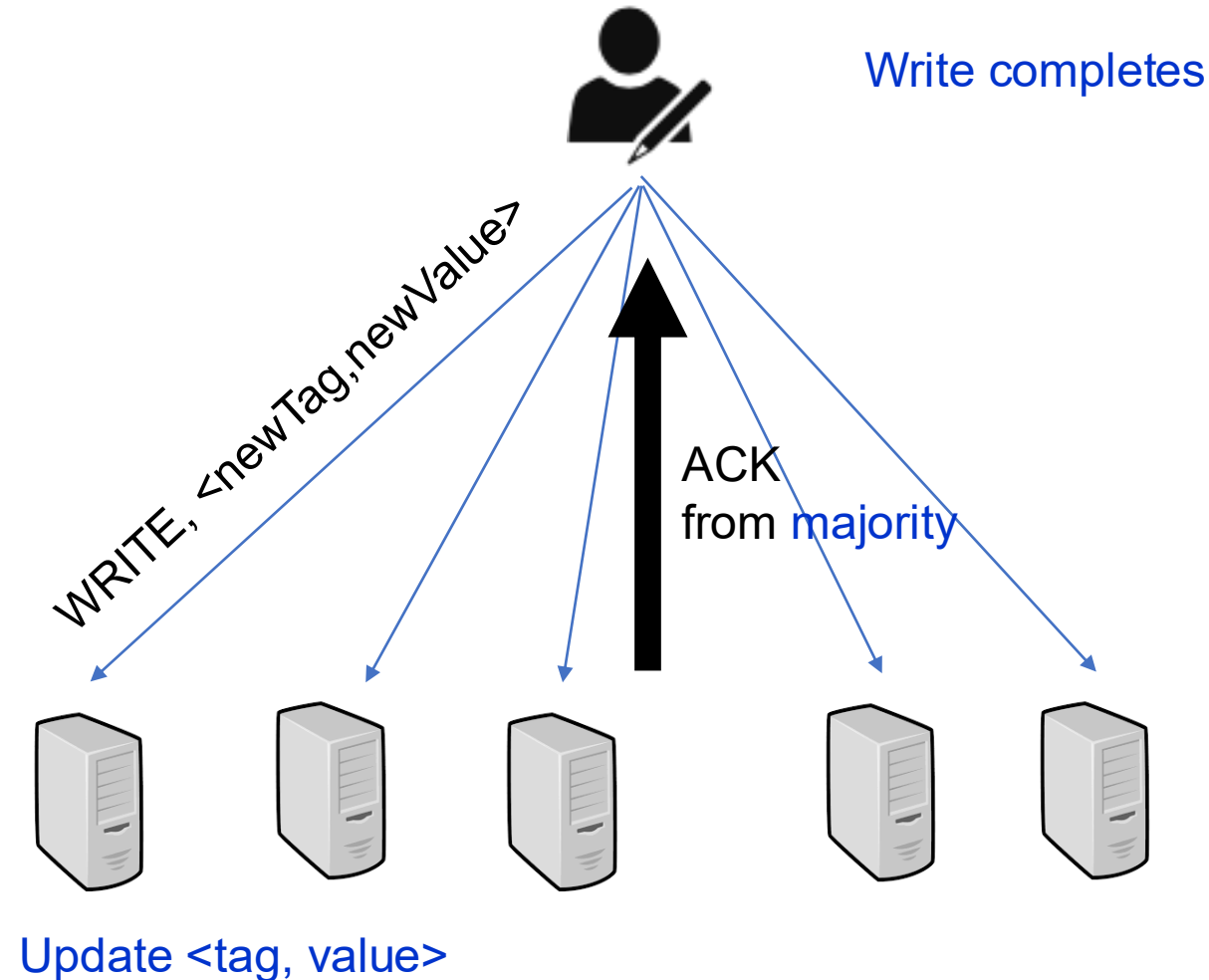
WRITE Protocol of MW-ABD

Phase1: **Query** - Discover maximum tag



WRITE Protocol of MW-ABD

Phase 2: **Write** <newTag, newValue>

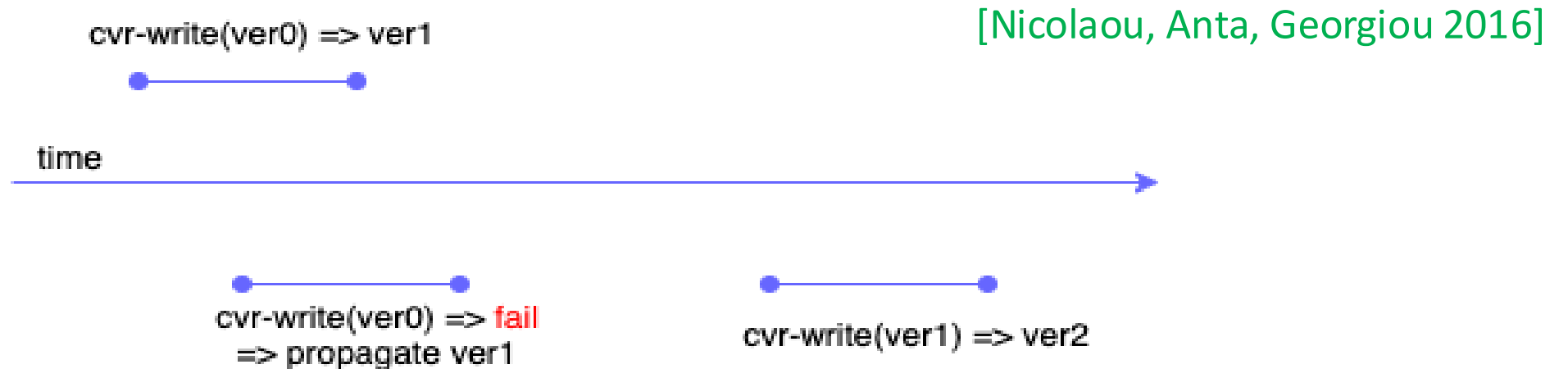


Handling Large Objects

- Many DSMs, for both **static** (fixed servers) and **dynamic** (remove/add servers) systems, were developed over the years.
- These protocols are:
 - efficient for **small objects**
 - expensive solutions that are difficult to implement in an asynchronous, fail prone, message passing environment.
- **Limitation 1**: No enforced dependence between versions of the object [versioning]
- **Limitation 2**: Unable to handle write operations working on different parts of a large object [striping]

Versioning and Coverability

- **Versioned object**: a R/W object where each value written is associated with a version (= tag)
- **Coverability** is defined over a **totally ordered set of versions**, and extends linearizability by guaranteeing that:
 - A write on the object **succeeds only** when the write is associated with the **latest version** of the object. Otherwise, it becomes a read operation (to obtain the latest version)



- **CoABD: Modified** MW-ABD that guarantees coverability

Concurrent Access on Large Objects

- Concurrent write operations may **overwrite** one another
- In some cases this is unavoidable. But what if two changes take place on different parts of a large file?



Ultimate Objective

The development of an **efficient** Distributed Shared Memory that provides **provable atomic consistency guarantees** and **high access concurrency** at large scale under an asynchronous, crash-prone, distributed and even **dynamic** environment.

Methodology

1. Implemented the most efficient Atomic Shared Object Algorithm.
2. Specified a Data Fragmentation Strategy.
3. Designed and implemented the framework of our system, and introduce new consistency guarantees.
4. Implemented ARES to introduce a dynamic solution.
5. Combined our Fragmentation Strategy with ARES.
6. Evaluated our implementation against commercial solutions.
7. Deployed and evaluated the system in network testbeds (Emulab, AWS, Fed4FIRE+)
8. Identified any performance bottlenecks using distributed tracing and optimized them.

System Settings

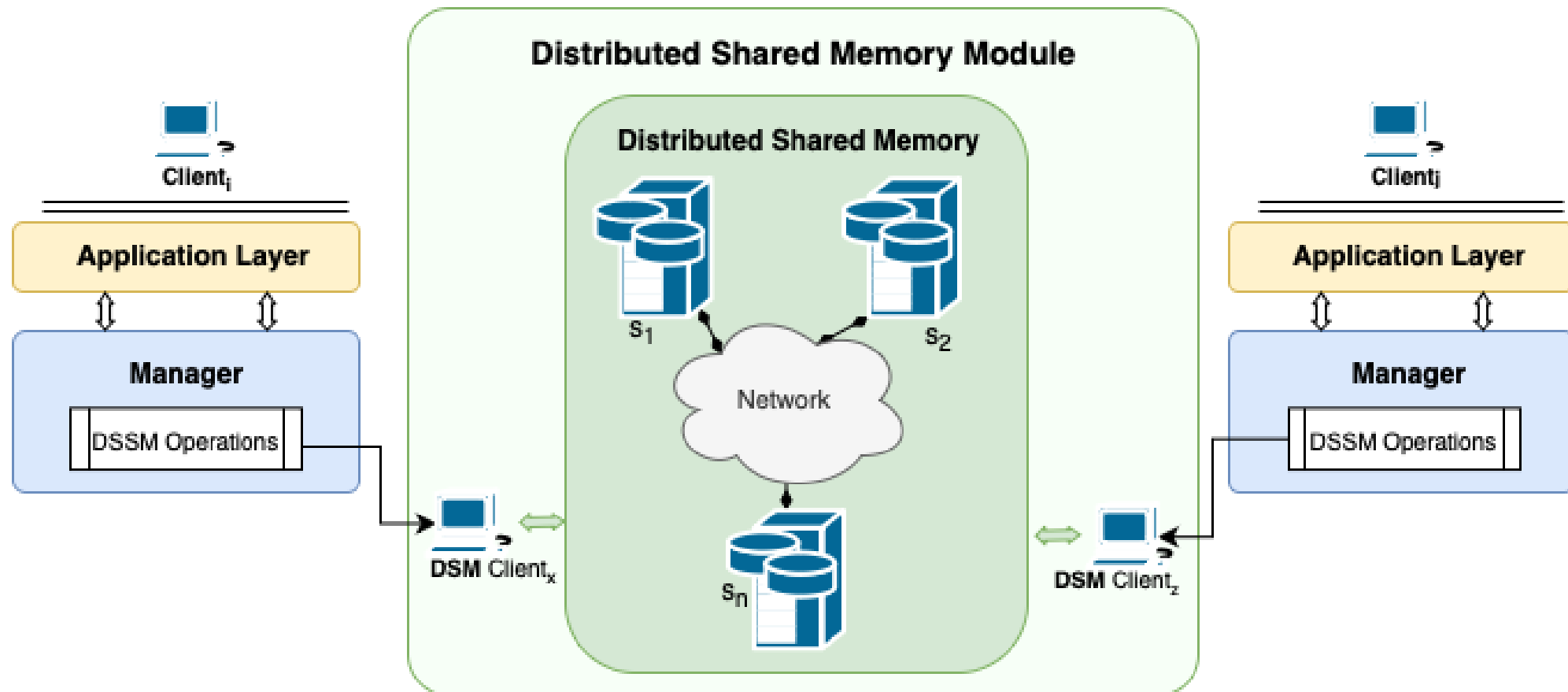
To implement an atomic consistent DSM that:

- supports large shared R/W objects
- with two main distinct sets of processes: a set **C** of clients and a set **S** of servers

(**Configuration** = the set **S** of servers and some additional info)

- Fixed Configuration -> **Static** environment
- Reconfiguration -> **Dynamic** environment
- **C** = a set **W** of writers, a set **R** of readers, and a set **G** of reconfigurers

Basic Architectures



Evaluation Setup



Testbeds



Emulab: a network testbed with tunable and controlled environmental parameters.



AMAZON Web Services (AWS) EC2: a web service that provides scalability and performance.



Fed4FIRE+: a federation of testbeds.



Performance Metric

- Average Operation latency of all clients (Communication + Computation Overhead).
- Update/Write Success Ratio

PART I - Fragmented Objects: Boosting Concurrency of Shared Large Objects

covering Stages 1-3, 7 in Methodology

Approach and Contribution

- Define concurrent objects: (i) the *block* object, and (ii) the *fragmented* object.
- Define the consistency that the fragmented object provides (*fragmented coverable linearizability*).
- Implement CoBFS, a Framework which implements the fragmented objects.
- Performed an experimental evaluation of CoBFS on Emulab.

Fragmented Objects: Boosting Concurrency of Shared Large Objects

Antonio Fernández Anta ¹ Chryssis Georgiou ²
Theophanis Hadjistasi ³ Nicolas Nicolaou ³ Efstathios Stavrakis ³
Andria Trigeorgi ²

¹ IMDEA Networks Institute, Madrid, Spain

² University of Cyprus, Nicosia Cyprus

³ Algolysis Ltd, Limassol, Cyprus

SIROCCO 2021

algolysis
algorithmic solutions



University
of Cyprus

institute
imdea
networks

2016-2020
ReSTART
SEARCH
ΠΡΟΓΡΑΜΜΑ ΕΡΕΥΝΑΣ, ΤΕΧΝΟΛΟΓΙΚΗΣ ΑΝΑΤΙΤΘΗΣ ΚΑΙ ΚΑΙΝΟΤΟΜΙΑΣ



RESEARCH
& INNOVATION
FOUNDATION



European Union
European Regional
Development Fund

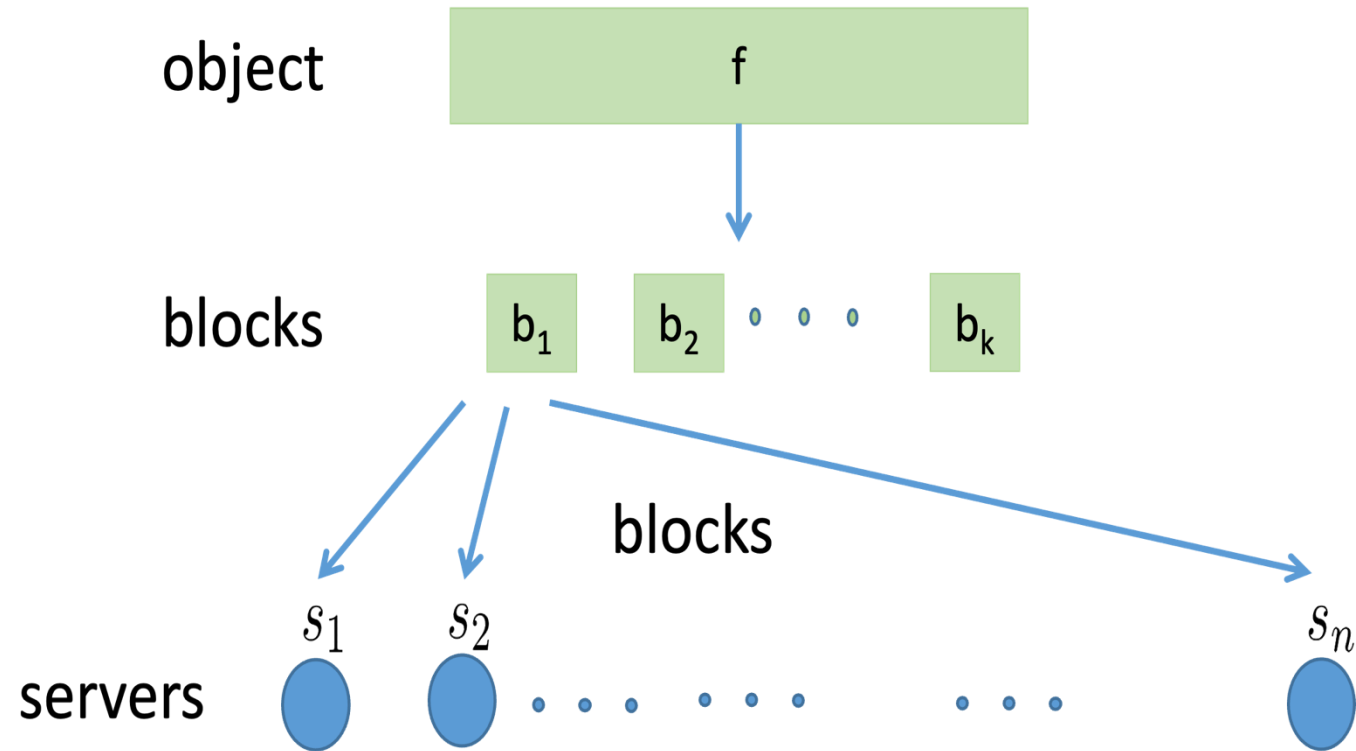


Republic of Cyprus

Solution: Fragmentation

Each object is fragmented into blocks

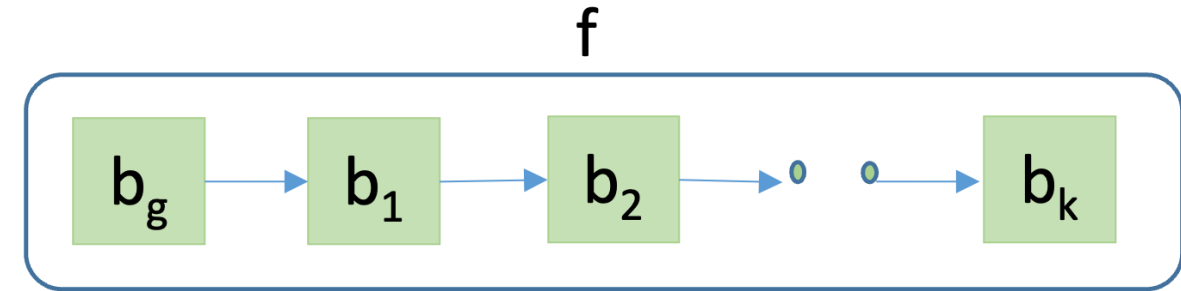
- Allows big amounts to be distributed all over the servers
- Avoids contention for concurrent accesses to different blocks



Solution: Fragmentation

- **Fragmented object**

- Each f is a *list of blocks*
- Each block has the id of its next block
- Each block is linearizable and coverable
- The first block is the genesis block b_{gen}



- **Write Operation write(f)**

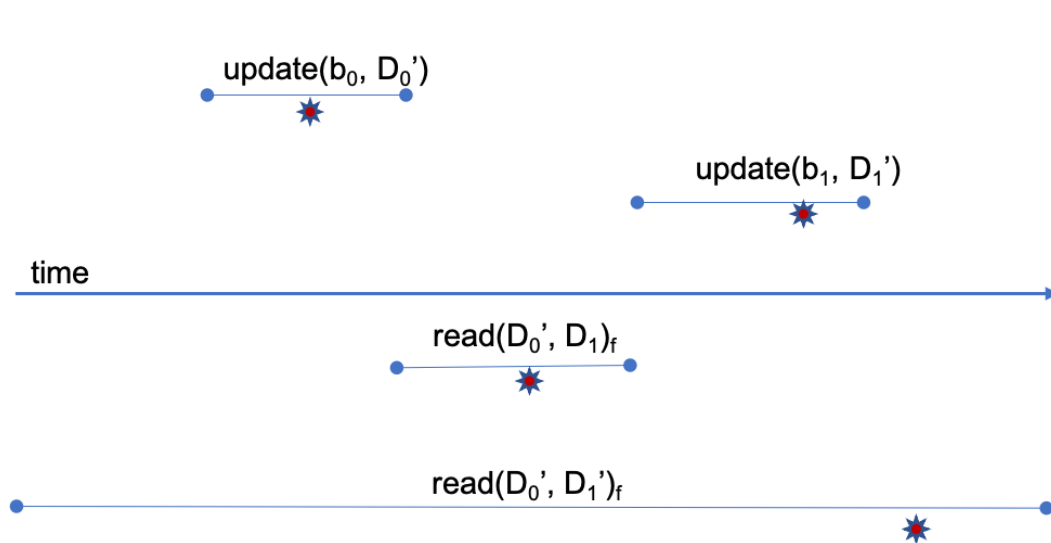
- Propagate **only** modified and new blocks



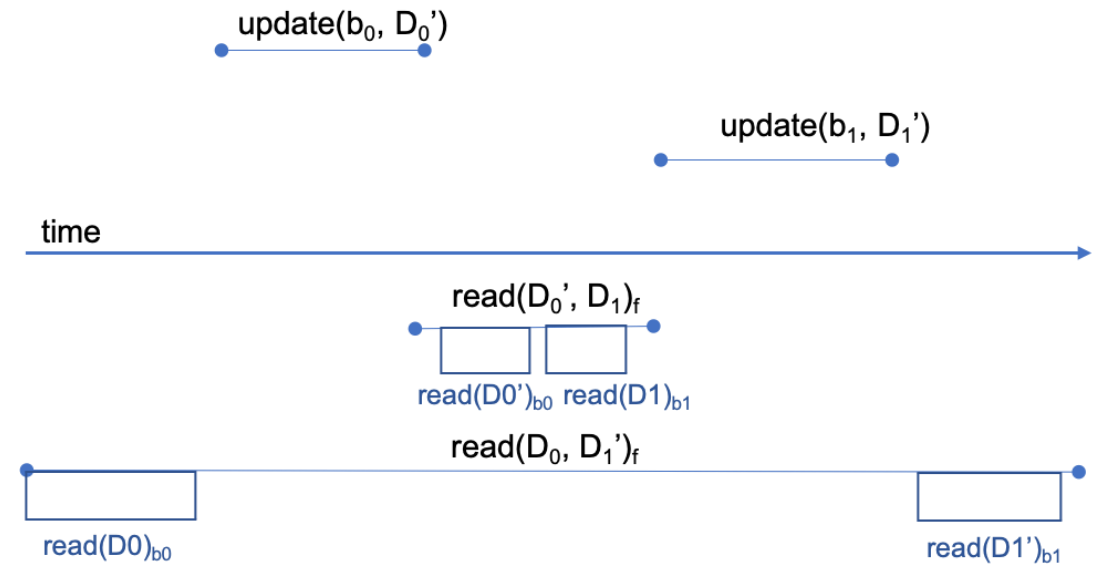
- **Read Operation read(f)**

- Start from b_{gen} and read all the blocks

Fragmented Coverable Linearizability

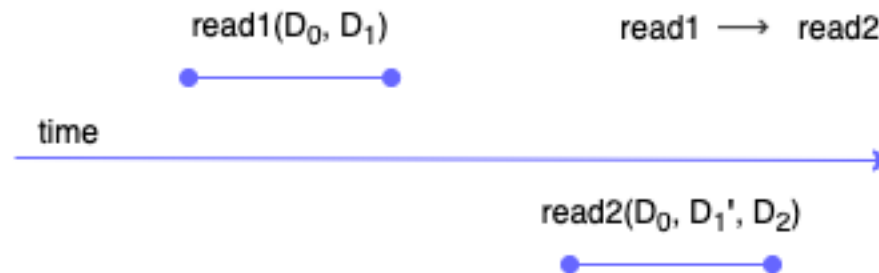


(a) Linearizability on the whole object

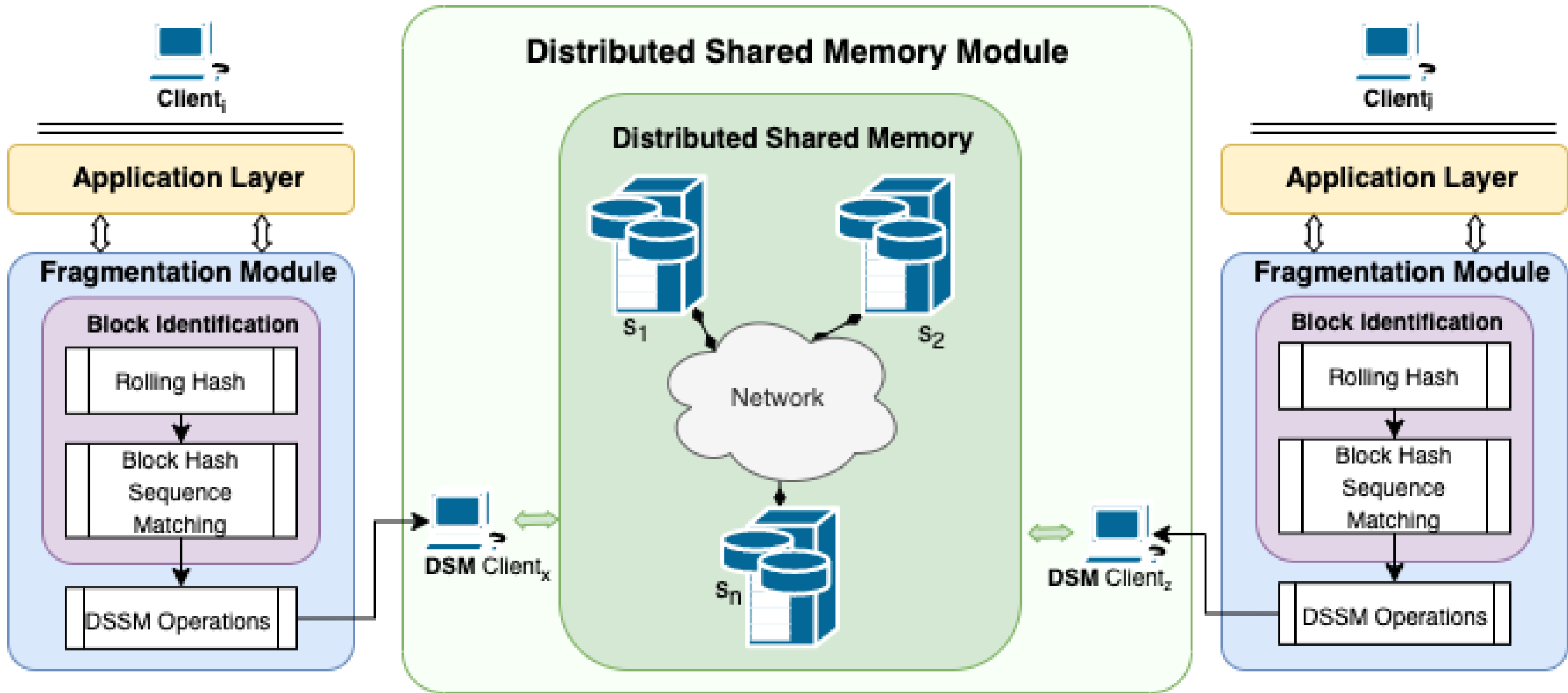


(b) Fragmented Linearizability

Fragmented Coverable Linearizability guarantees that all concurrent operations on different blocks prevail, and only concurrent operations on the same blocks are conflicting.



CoBFS Framework



CoABD-F: Integration of CoABD with CoBFS

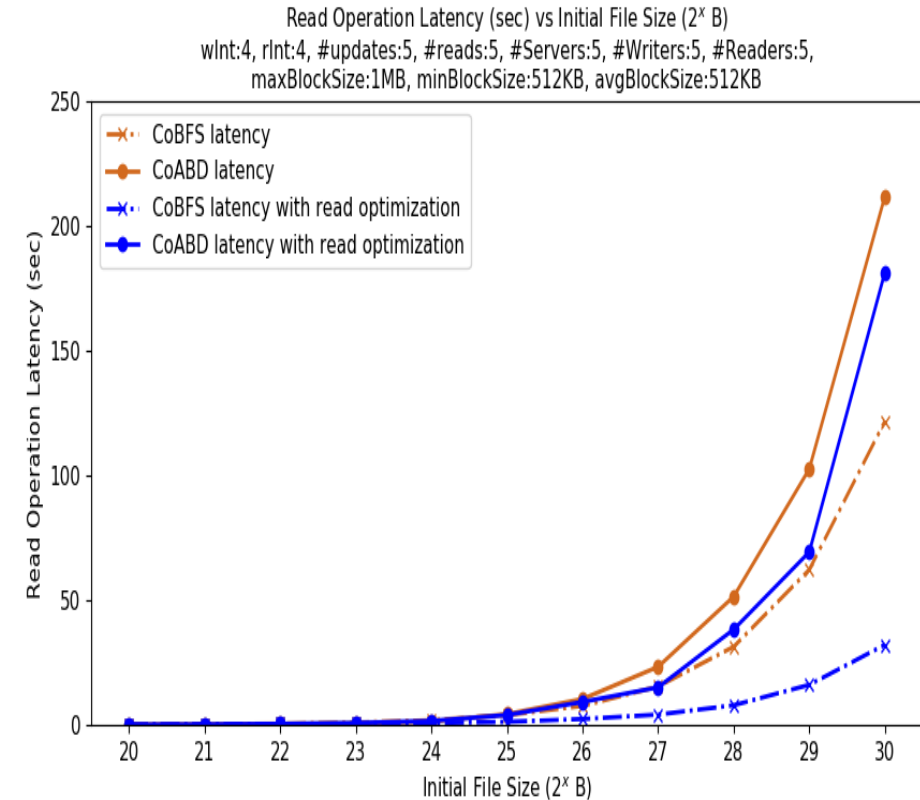
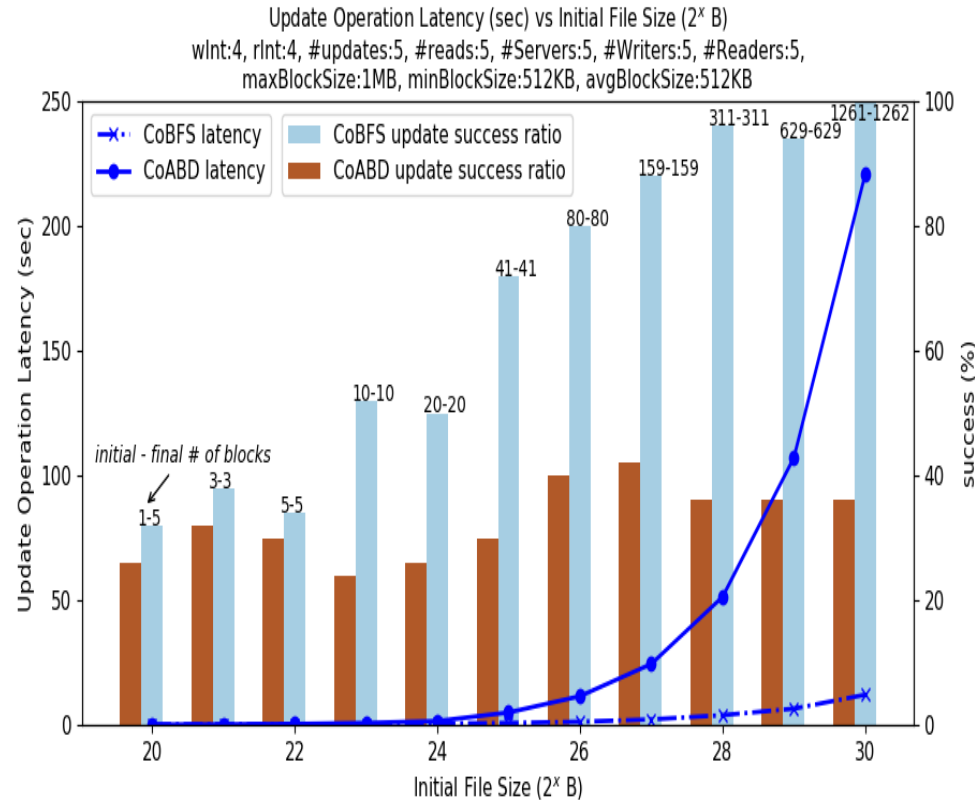
Read Operation

- FM issues `dsmm-read` for each block
- `crv-read` executes the read operation on DSMM
- **CoABD Optimization**: Only blocks that have changed are transferred

Update Operation

- FM divides objects into blocks $\langle D_0, \dots, D_k \rangle$ using **Block Identification (BI)**
- FM updates each block with its corresponding data
 - If $k=0$: `dsmm-write` is used
 - If $k>0$: New blocks are created using `dsmm-create`, modified blocks are written using `dsmm-write`, and block pointers are updated
 - Block Operations: Sequential creation and write operations for chunks of data
- The block operations are executed using `cvr-write` on the DSMM

CoABD VS. CoABD-F (Emulab)



File Size:

- the update latency of CoABD-F remains at extremely low levels, although the file size increases.
- a read optimization decreases significantly the CoABD-F read latency, since it is more probable for a reader to already have the last version of some blocks.

Overview of Results

- CoABD-F has significantly lower write and read operation latency (especially in larger files)
- For the read operation latency of smaller sizes (1MB), suggest that there is room for improvement
- Trade-off between block size, operation latency and write success rate
- ...

Recap and Next Goal

- CoABD-F implements a Robust, Strongly-consistent DSM in an asynchronous message-passing system and supports Versioning, Data Striping, and High Access Concurrency for Large Objects

static



- Next Goal: Obtain such as DSM for Dynamic systems, where servers (replica hosts) change over time, without interrupting the read/write operations or violating data consistency

PART **II** - Implementation and Experimental Evaluation of ARES

covering Stages 4 & 6, 7 in Methodology

Approach and Contribution

- Implement ARES to enable dynamic reconfiguration.
- Performed an experimental evaluation of ARES on Emulab.
- Set up two open-source commercial solutions (Redis and Cassandra) to compare them with ARES.
- Performed experiments in real testbeds (supported by Fed4FIRE+ project), distributed in the European Union (EU) and the USA.

ARES: Adaptive, Reconfigurable, Erasure coded, Atomic Storage

Nicolas Nicolaou¹ Viveck Cadambe²
N. Prakash³ Andria Trigeorgi⁴ Kishori M. Konwar⁵
Muriel Medard⁵ Nancy Lynch⁵

¹ Algolysis Ltd, Limassol, Cyprus

² Pennsylvania State University, US

³ Intel Corp.

⁴ University of Cyprus, Nicosia Cyprus

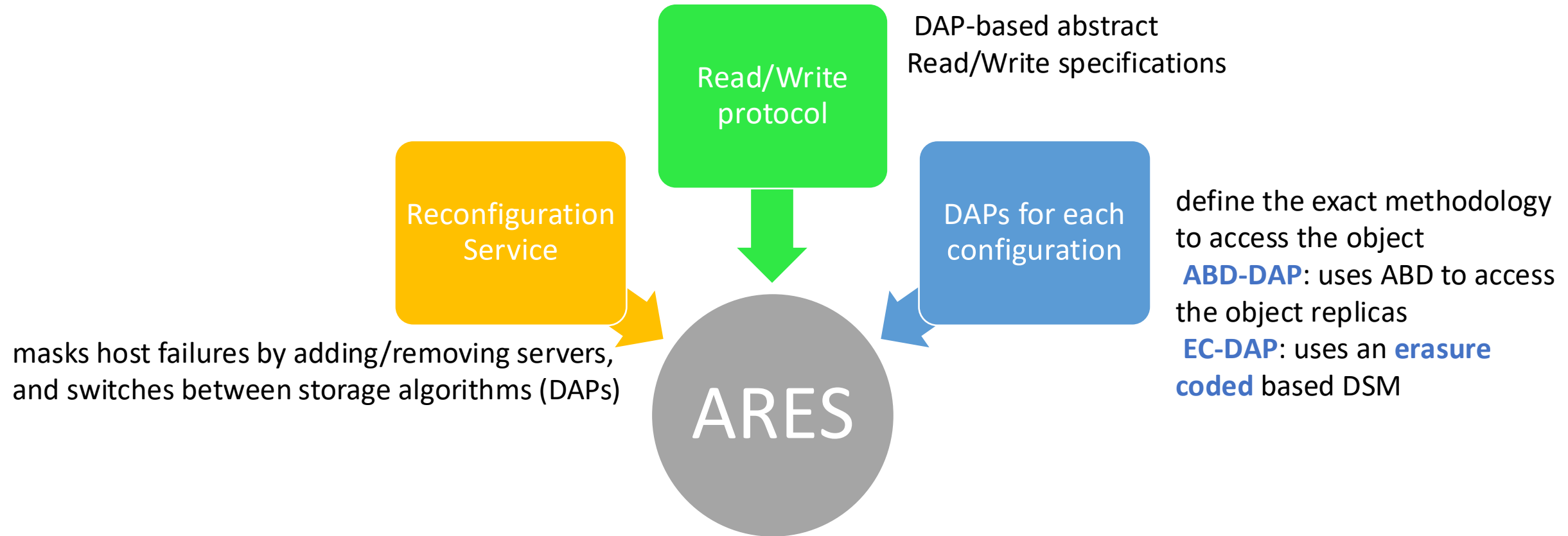
⁵ Massachusetts Institute of Technology, USA

ACM Transactions on Storage



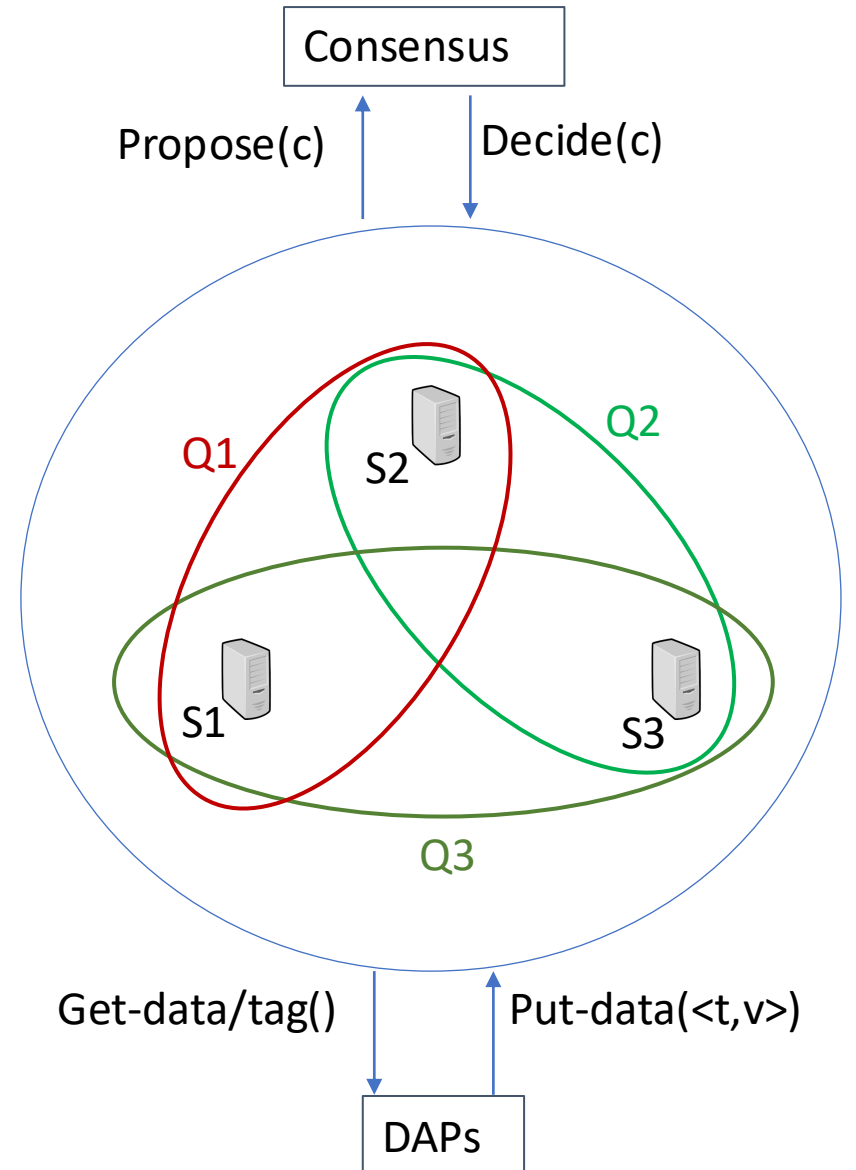
ARES (MWMMR) - Adaptive, Reconfigurable, Erasure Code, Atomic Storage

[Nicolaou et al 2022]

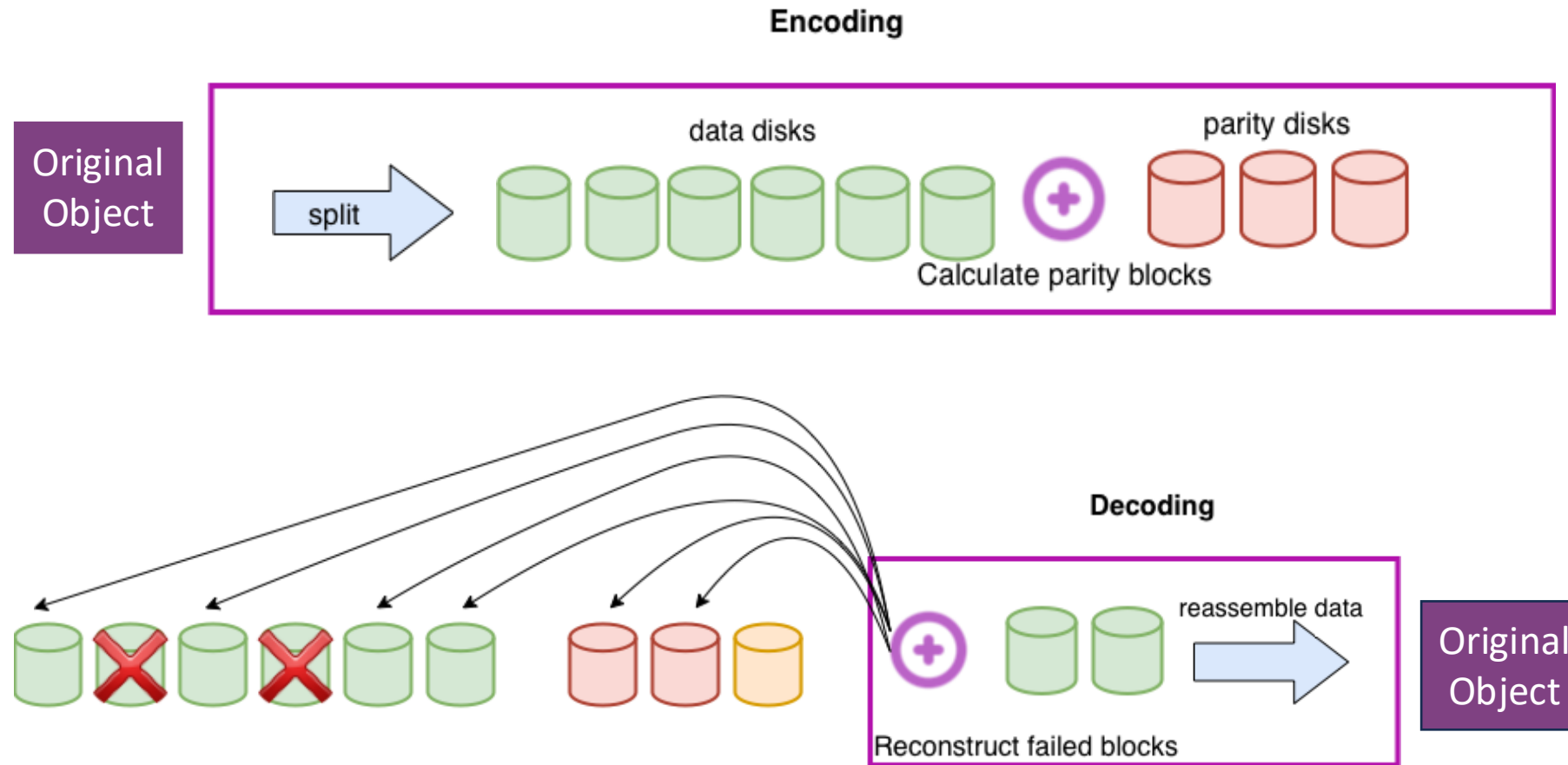


Configurations

- A configuration c is characterized by:
 - A unique identifier
 - A set of servers
 - A quorum set system on servers
 - A consensus instance
 - A DAP implementation
 - D1. $c.get - tag()$: returns a tag $\tau \in T$
 - D2. $c.get - data()$: returns a tag – value pair $(\tau, v) \in T \times V$
 - D3. $c.put - data(<\tau, v>)$: the tag – value pair $(\tau, v) \in T \times V$ as argument
 - ABD-DAP & EC-DAP



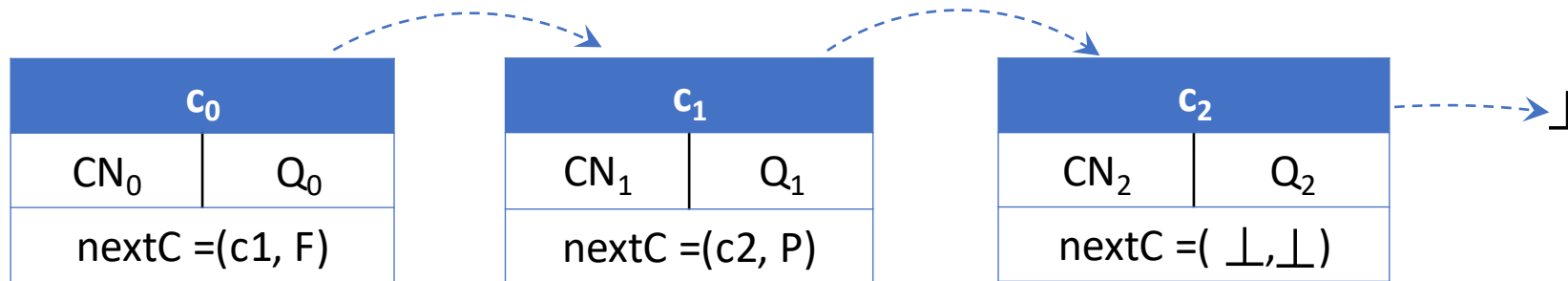
ARES (MWMMR) - Adaptive, Reconfigurable, Erasure Code, Atomic Storage



(n, k)-Reed-Solomon code: n=servers, k=data servers, m=parity servers
BUT reads and writes are still applied on the entire object

Configuration Sequence

- Global configuration sequence G_L
- **nextC**: each server points to the next configuration
 - Same nextC to all servers of a single config c (due to consensus)
- **Flags {P, F}**: pending, finalized
 - Pending: not yet a quorum of servers received msgs
 - Finalized: new configuration propagated to a quorum of servers



Reconfiguration Service

- A reconfig operation performs 2 major steps:
 - 1) Configuration *Sequence Traversal*
 - 2) Configuration *Installation*
 - Transfers the object state from the old to the new configuration

```
6: operation reconfig(c)
    if  $c \neq \perp$  then
8:    $cseq \leftarrow \text{read-config}(cseq)$            attempt get to the latest configuration } (1)
       $cseq \leftarrow \text{add-config}(cseq, c)$        introduce the new configuration
10:   $\text{update-config}(cseq)$                      move the latest value to the new config } (2)
       $cseq \leftarrow \text{finalize-config}(cseq)$     let servers know it is good to be finalized
12: end operation
```

This service guarantees that if $cseq1$ and $cseq2$ are obtained by two clients resp., then either $cseq1$ is a prefix of $cseq2$ or vice versa

Read/Write Operations using DAPs

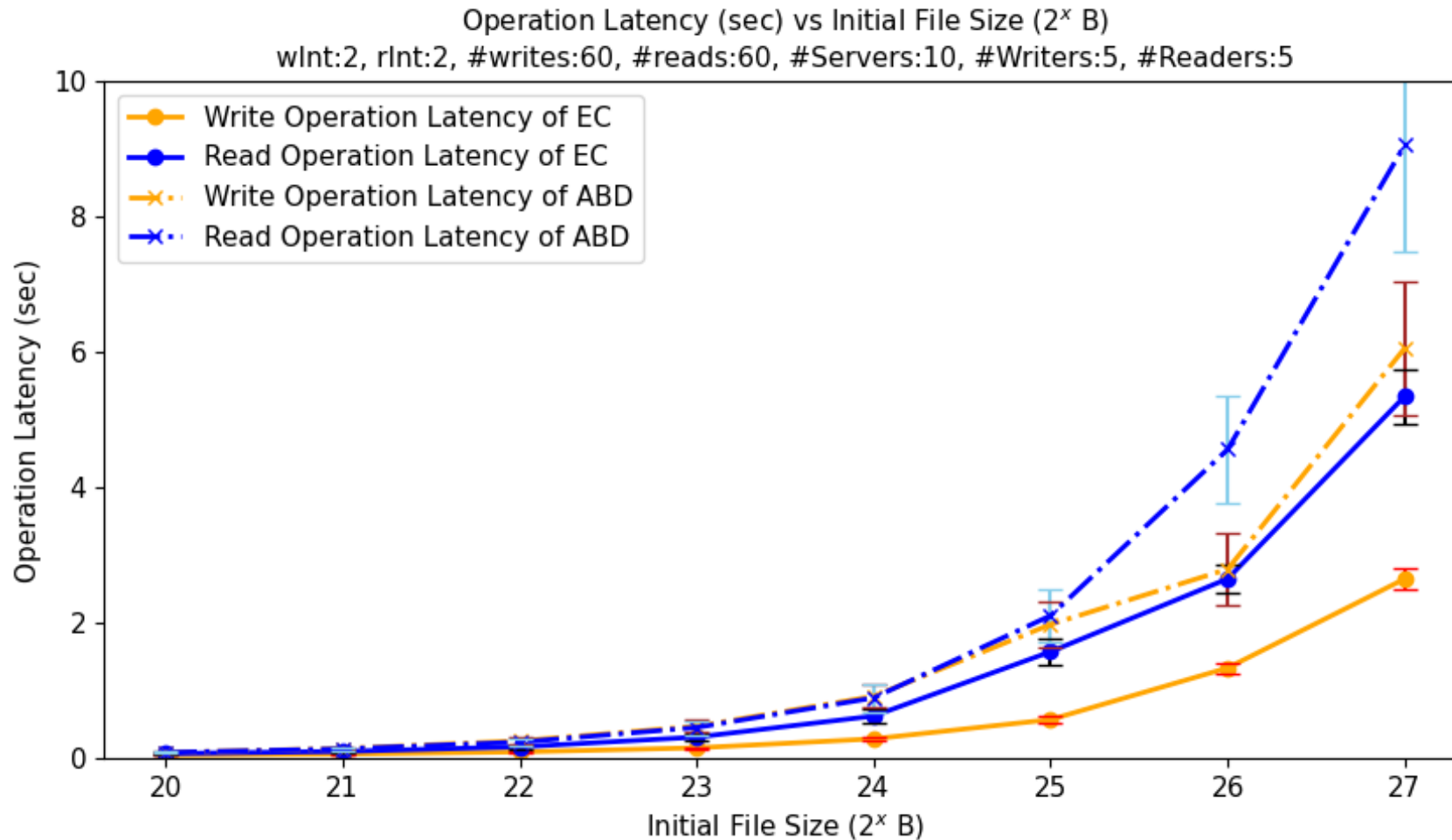
Reader Protocol

- Traverse Config Sequence `cseq`
- Find $\mu = \max(\langle c, F \rangle)$ in `cseq`
- Set $v = \text{last}(\langle c, * \rangle)$ in `cseq`
- Discover for $\mu \leq i \leq v$
`(t,v)=max(cseq[i].get-data())`
- Do
 - `cseq[v].put-data(t,v)`
 - Traverse Sequence `cseq`
- `while(|cseq| > v)`

Writer Protocol(val) (at w_i)

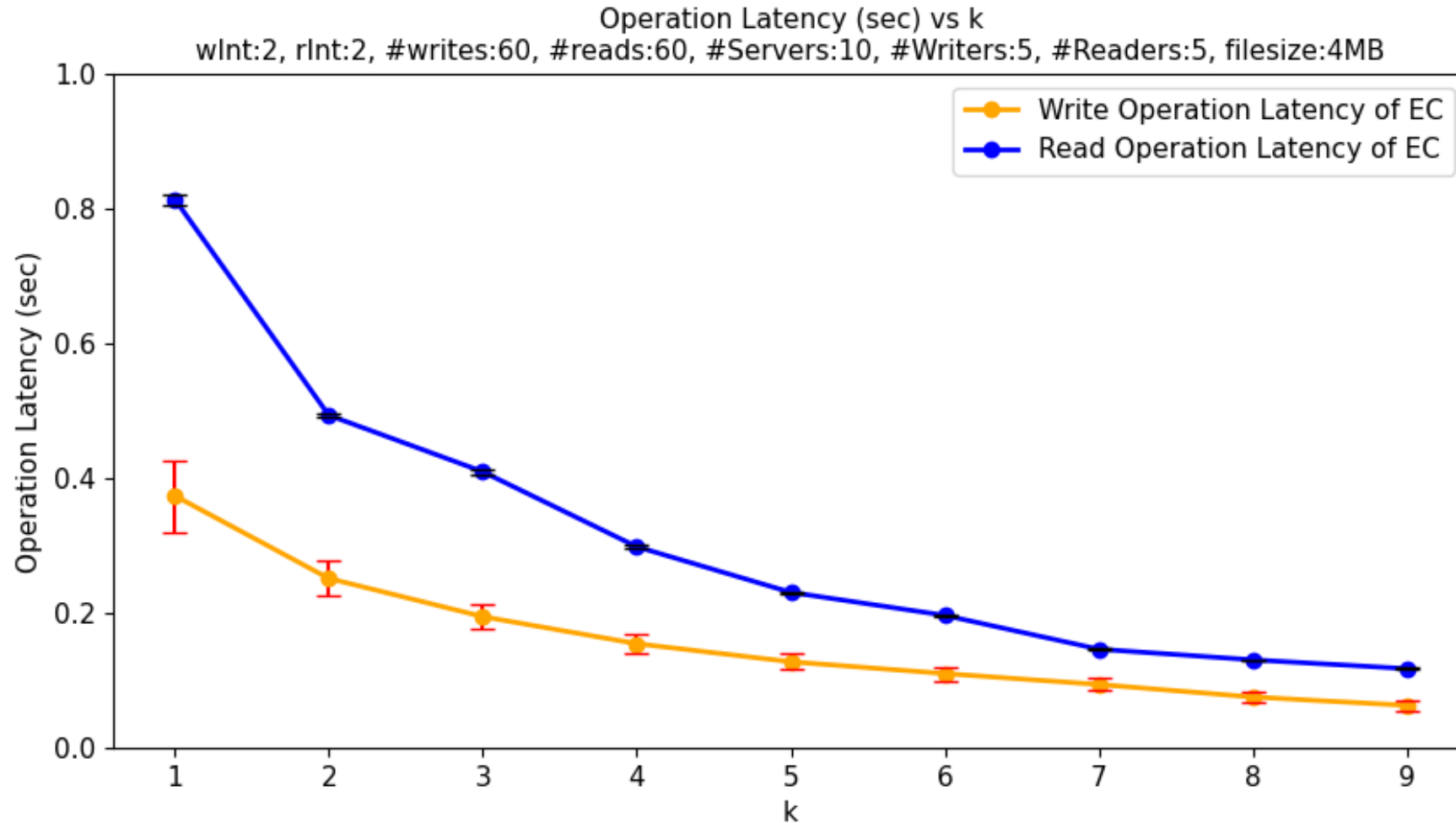
- Traverse Config Sequence `cseq`
- Find $\mu = \max(\langle c, F \rangle)$ in `cseq`
- Set $v = \text{last}(\langle c, * \rangle)$ in `cseq`
- Discover for $\mu \leq i \leq v$
`tmax=max(cseq[i].get-tag())`
- `(t,v)= (<tmax+1,wi>, val)`
- Do
 - `cseq[v].put-data(t,v)`
 - Traverse Sequence `cseq`
- `while(|cseq| > v)`

File Size (Emulab)



- the read and write latencies of both storage algorithms remain in low levels until 16 MB
- the write operation of EC algorithm is the faster
- the larger messages sent by ABD result in slower read operations

K Scalability (Emulab)



- *small k (=smaller number of data fragments) \rightarrow bigger sizes of the fragments and higher redundancy.*
- *The write latency seems to be less affected by the number of k since the first phase of write asks only for the tag*

Overview of Results

- The **write** operation of ARES_EC algorithm is the **faster** (especially in larger sizes)
- The **larger messages** sent by ARES_ABD result in **slower read** operations
- Reconfiguration (server failures & changing DAPs).
- The **reconfiguration** operation is the slower
- Trade-off between **operation latency** and **parity of EC**
- Trade-off between **operation latency** and the **number of writers**

Invited Paper: Towards Practical Atomic Distributed Shared Memory: An Experimental Evaluation

Authors: Andria Trigeorgi¹, Nicolas Nicolaou², Chryssis Georgiou¹, Theophanis Hadjistasi²,
Efsthios Stavarakis², Viveck Cadambe³, and Bhuvan Ugaonkar³

¹University of Cyprus, Nicosia, Cyprus

²Algolysis, Limassol, Cyprus

³Pennsylvania State University, PA, USA

SSS 2022

Funded by: *EU's NGIAtlantic.eu cascading grant agreement no. OC4-347*



Main Objective

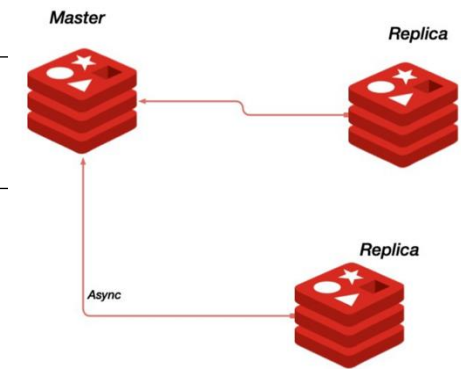
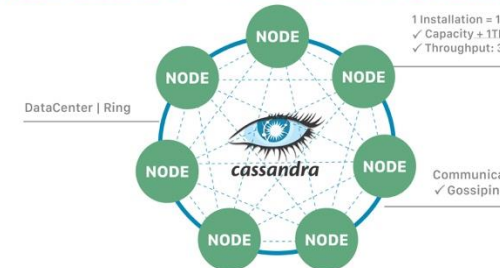
- **Primary Limitation:** Expensive DSMs that are difficult to implement in an asynchronous, fail prone, message passing environment.
- **Recent works:** Reduce high communication, storage, and computation overheads (CoBFS, ARES)

How may such algorithms compare to commercially used solutions?

Comparative Table

| Algorithm | Consistency guarantees | Data Striping | Non-Blocking Reconfiguration |
|-----------------|---------------------------|---------------|-----------------------------------|
| <i>MWMR</i> ABD | Atomic | NO | NO |
| ARES-ABD | Atomic | NO | YES |
| ARES-EC | Atomic | YES | YES |
| CASSANDRA | Tunable (eventual/atomic) | NO | NO (A single server at a time) |
| REDIS/REDIS_W | Eventual | NO | NO |

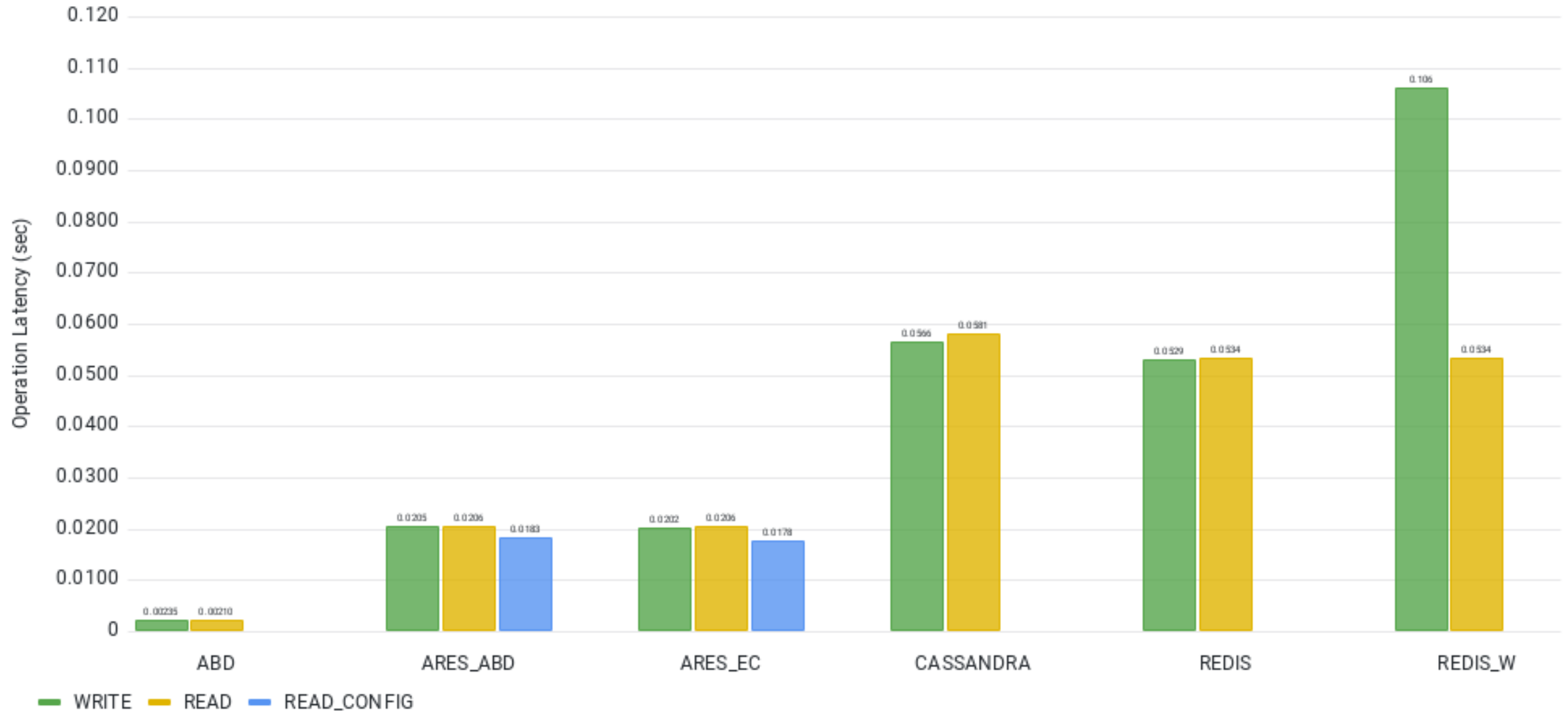
ApacheCassandra™ = NoSQL Distributed Database



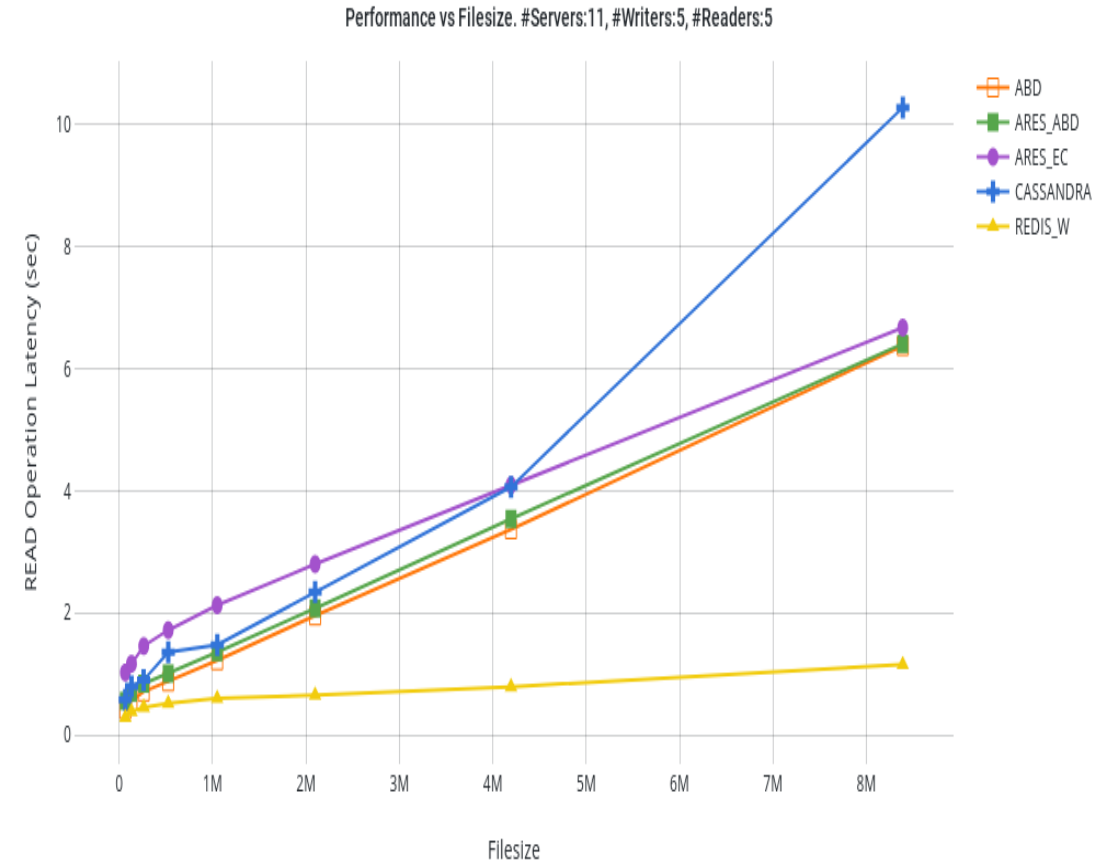
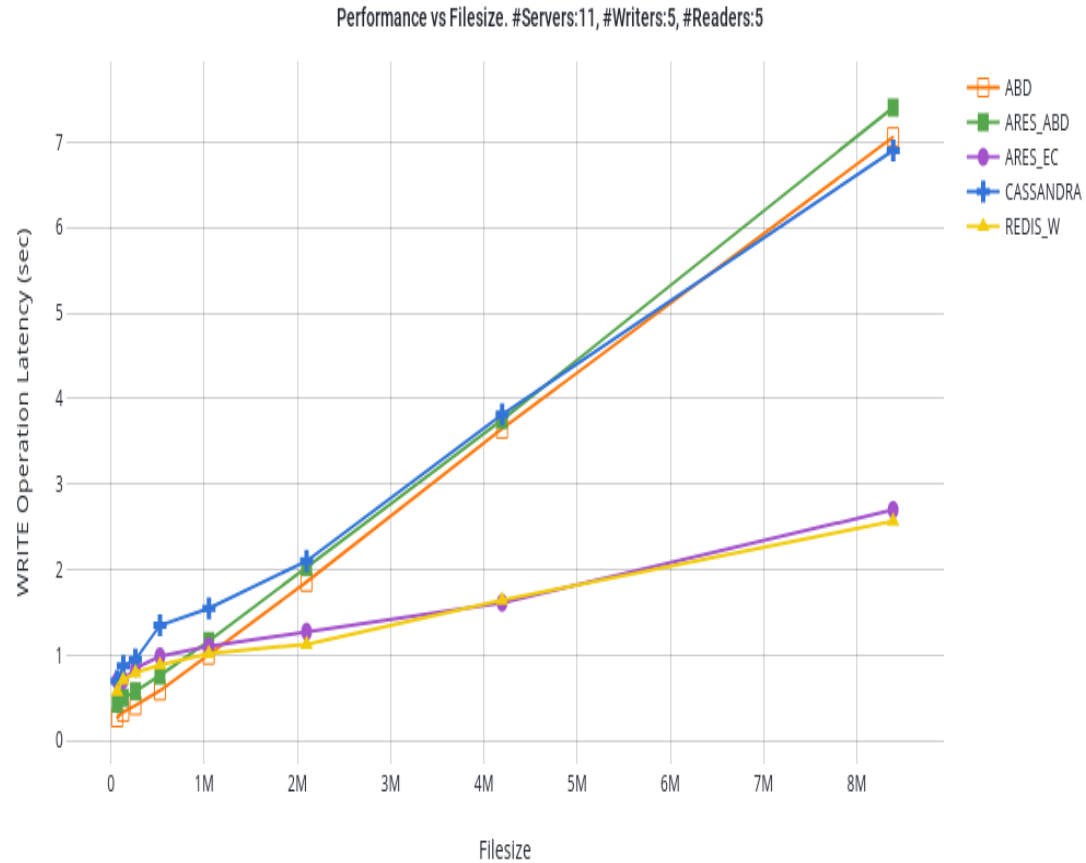
Stress Test – Topology (Fed4FIRE+)

i

Performance vs Algorithm. topology:0E+3U, clients' continent:US, S:3, W:1, R:1, fsize:32



Stress Test – Object Size (Fed4FIRE+)



Overview of Results

- ABD & ARES scale better than Cassandra.
- ARES_EC outperforms Cassandra in **larger** objects.
- ARES has non-blocking reconfiguration mechanism (server failures & changing DAPs).
- The **topology** played a major role on the performance (throughput), of all the algorithms we studied.
- Verify the **fault-tolerance** and the **responsiveness** of ARES.
- ARES trades **performance** over **consistency** with respect to Redis.

Recap and Next Goal

- ARES
 - provable guarantees
 - Operations compete closely/outperform existing DSS solutions (even when offering weaker consistency guarantees).
- Results
 - Data repository: <https://github.com/nicolaoun/ngiatlantic-public-data>
 - Graphs: <https://projects.algolysis.com/ares-ngi/results/>
- Next Goal
 - Integrate the dynamic (reconfigurable) DSM algorithm **ARES** with the DSM module in **CoBFS** → **CoARESF**

PART **III** - Fragmented ARES: Dynamic Storage for Large Objects

covering Stages 5, 7 in Methodology

Part III - Approach and Contribution

- Integrate the dynamic (reconfigurable) DSM algorithm [ARES](#) with the DSM module in [CoBFS](#) \Rightarrow [CoARES](#)
- In order to integrate ARES in CoBFS, we first needed to obtain a coverable version of ARES \Rightarrow [CoARES](#)
- Performed an in-depth experimental evaluation over [Emulab](#) and [AWS](#) comparing the different versions (with/out fragmentation, with/out EC, with/out reconfiguration).

Fragmented ARES: Dynamic Storage for Large Objects

Authors: Chryssis Georgiou¹, Nicolas Nicolaou², Andria Trigeorgi¹

¹University of Cyprus, Nicosia, Cyprus

²Algolysis, Limassol, Cyprus

DISC 2022

Augusta, GA, USA



CoARES

version: the tag of the coverable object

flag: chg when the write is successful

Query phase: find the max tag-value pair

Main difference: the condition "the writer has the latest version?"

It updates the state of the object!

Algorithm 1 Write and Read protocols for CoARES.

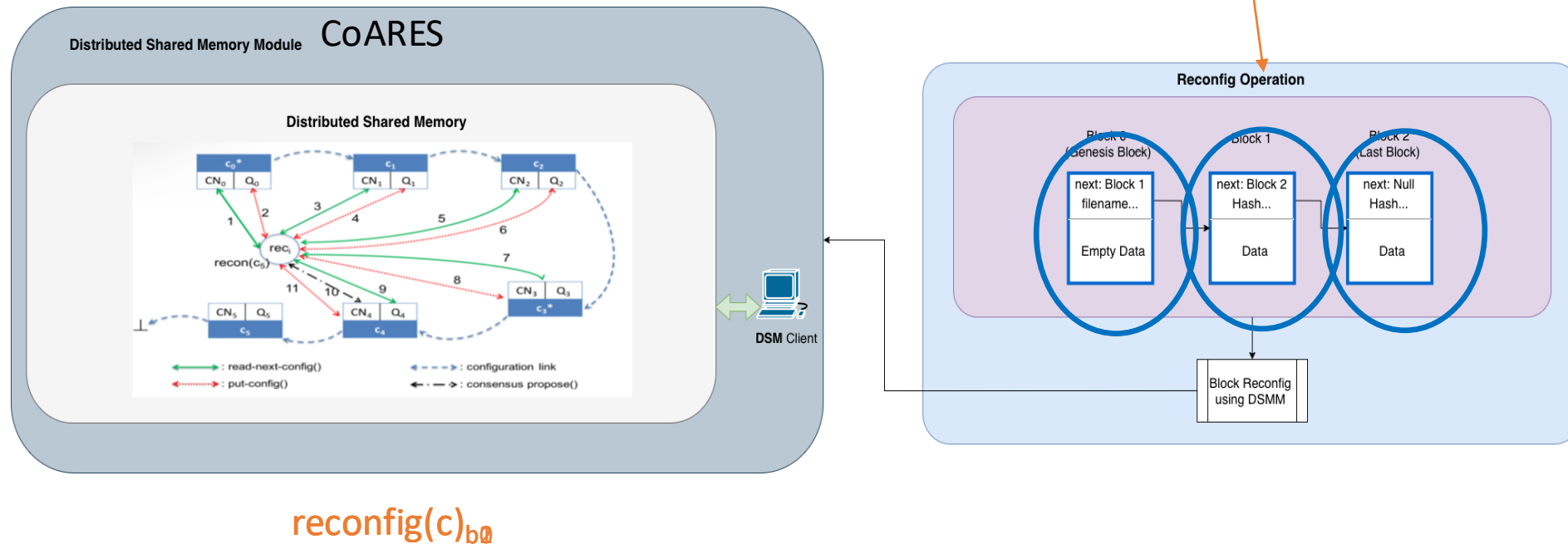
| | |
|--|--|
| <p>CVR-Write Operation:</p> <p>2: at each writer w_i</p> <p>State Variables:</p> <p>4: $cseq[i].s.t.cseq[j] \in \mathcal{C} \times \{F, P\}$</p> <p>6: $version \in \mathbb{N}^+ \times \mathcal{W} \cup \{\perp\}$ initially $\langle 0, \perp \rangle$</p> <p>Local Variables:</p> <p>8: $\mu \in \mathbb{N}^+$ initially 0, $\nu \in \mathbb{N}^+$ initially 0</p> <p>10: $\tau \in \mathbb{N}^+ \times \mathcal{W}$ initially $\langle 0, w_i \rangle$</p> <p>12: $v \in V$ initially \perp</p> <p>Initialization:</p> <p>12: $cseq[0] = \langle c_0, F \rangle$</p> <p>operation cvr-write(val), $val \in V$</p> <p>14: $cseq \leftarrow \text{read-config}(cseq)$</p> <p>16: $\mu \leftarrow \max(\{i : cseq[i].status = F\})$</p> <p>18: $\nu \leftarrow cseq$</p> <p>18: for $i = \mu : \nu$ do</p> <p>18: $\langle \tau, v \rangle \leftarrow \max(cseq[i].cfg.get-data(), \langle \tau, v \rangle)$</p> <p>20: if $version = \tau$ then</p> <p>20: $flag \leftarrow chg$</p> <p>20: $\langle \tau, v \rangle \leftarrow \langle \tau.ts + 1, \omega_i \rangle, val$</p> <p>22: else</p> <p>22: $flag \leftarrow unchg$</p> <p>24: $version \leftarrow \tau$</p> <p>24: $done \leftarrow false$</p> <p>26: while not done do</p> <p>26: $cseq[\nu].cfg.put-data(\langle \tau, v \rangle)$</p> <p>28: $cseq \leftarrow \text{read-config}(cseq)$</p> <p>28: if $cseq = \nu$ then</p> | <p>30: $done \leftarrow true$</p> <p>else</p> <p>32: $\nu \leftarrow cseq$</p> <p>end while</p> <p>34: return $\langle \tau, v \rangle, flag$</p> <p>end operation</p> <p>CVR-Read Operation:</p> <p>at each reader r_i</p> <p>State Variables:</p> <p>38: $cseq[i].s.t.cseq[j] \in \mathcal{C} \times \{F, P\}$</p> <p>Initialization:</p> <p>40: $cseq[0] = \langle c_0, F \rangle$</p> <p>operation cvr-read()</p> <p>42: $cseq \leftarrow \text{read-config}(cseq)$</p> <p>44: $\mu \leftarrow \max(\{j : cseq[j].status = F\})$</p> <p>46: $\nu \leftarrow cseq$</p> <p>46: for $i = \mu : \nu$ do</p> <p>46: $\langle \tau, v \rangle \leftarrow \max(cseq[i].cfg.get-data(), \langle \tau, v \rangle)$</p> <p>48: $done \leftarrow false$</p> <p>while not done do</p> <p>50: $cseq[\nu].cfg.put-data(\langle \tau, v \rangle)$</p> <p>52: $cseq \leftarrow \text{read-config}(cseq)$</p> <p>52: if $cseq = \nu$ then</p> <p>52: $done \leftarrow true$</p> <p>54: else</p> <p>54: $\nu \leftarrow cseq$</p> <p>end while</p> <p>56: return $\langle \tau, v \rangle$</p> <p>end operation</p> |
|--|--|

It keeps the max tag!

Difference: it returns both the value and the version

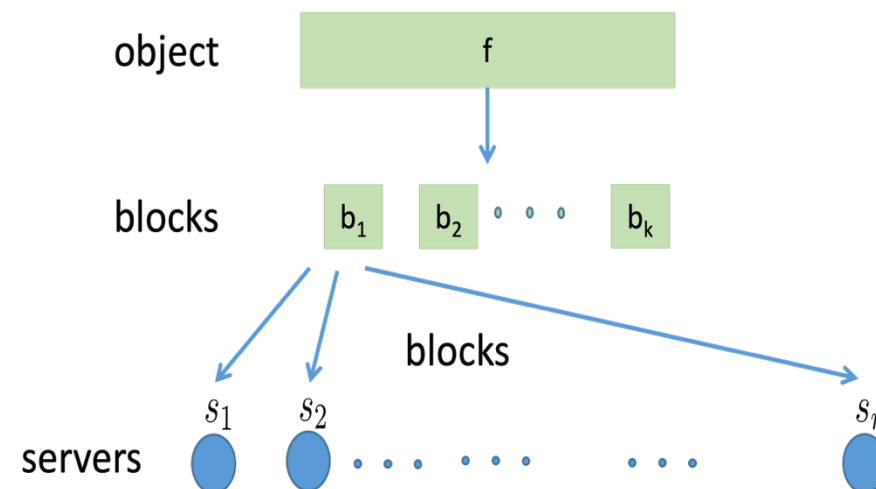
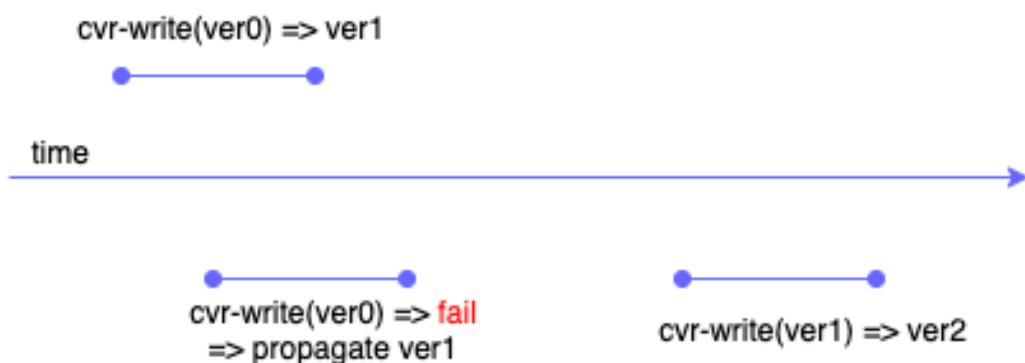
CoARES

- Integrate CoARES with CoBFS
- **Main challenge:** Enable the fragmentation approach to invoke reconfiguration operations



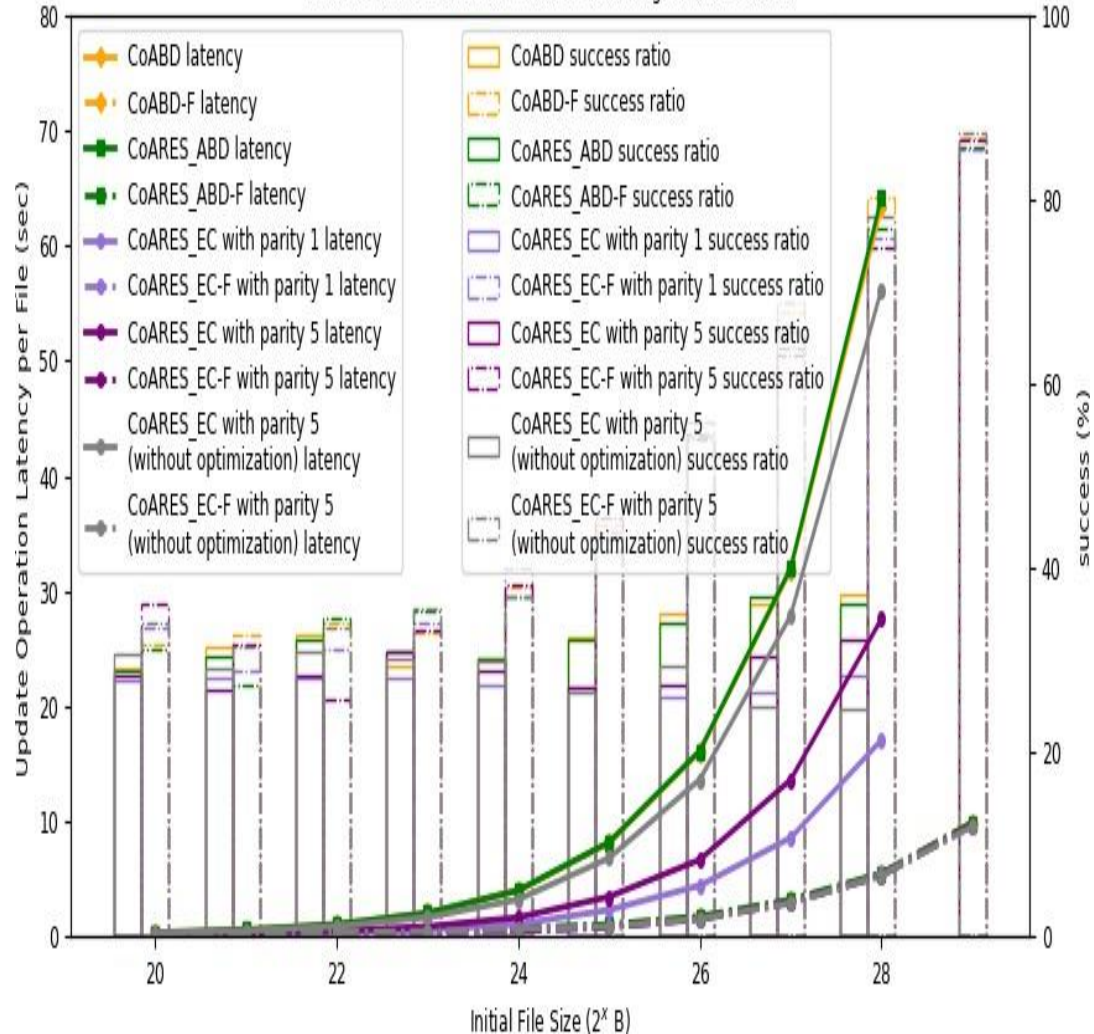
ARES VERSIONS

| | |
|-------------------|--|
| ARESABD | This is Ares that uses the ABD-DAP implementation. |
| CoARESABD | The coverable version of <i>ARESABD</i> . |
| CoARESABDF | The fragmented version of <i>CoARESABD</i> . |
| ARESEC | This is <i>ARES</i> that uses the EC-DAP implementation. |
| CoARESEC | The coverable version of <i>ARESEC</i> . |
| CoARESECF | This is the two-level data striping algorithm obtained when <i>CoARESF</i> is used with the EC-DAP implementation; i.e., it is the fragmented version of <i>CoARESEC</i> . |

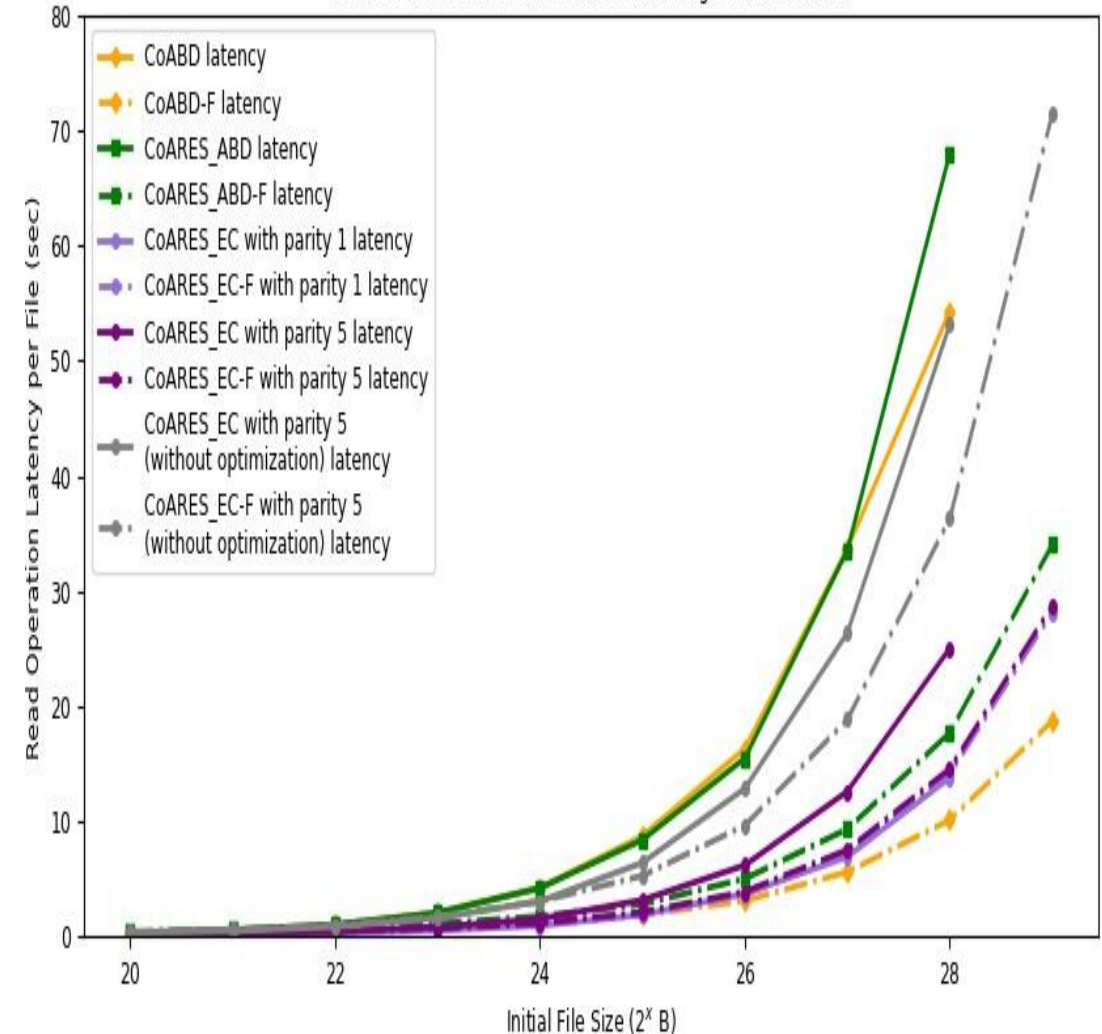


File Size (Emulab)

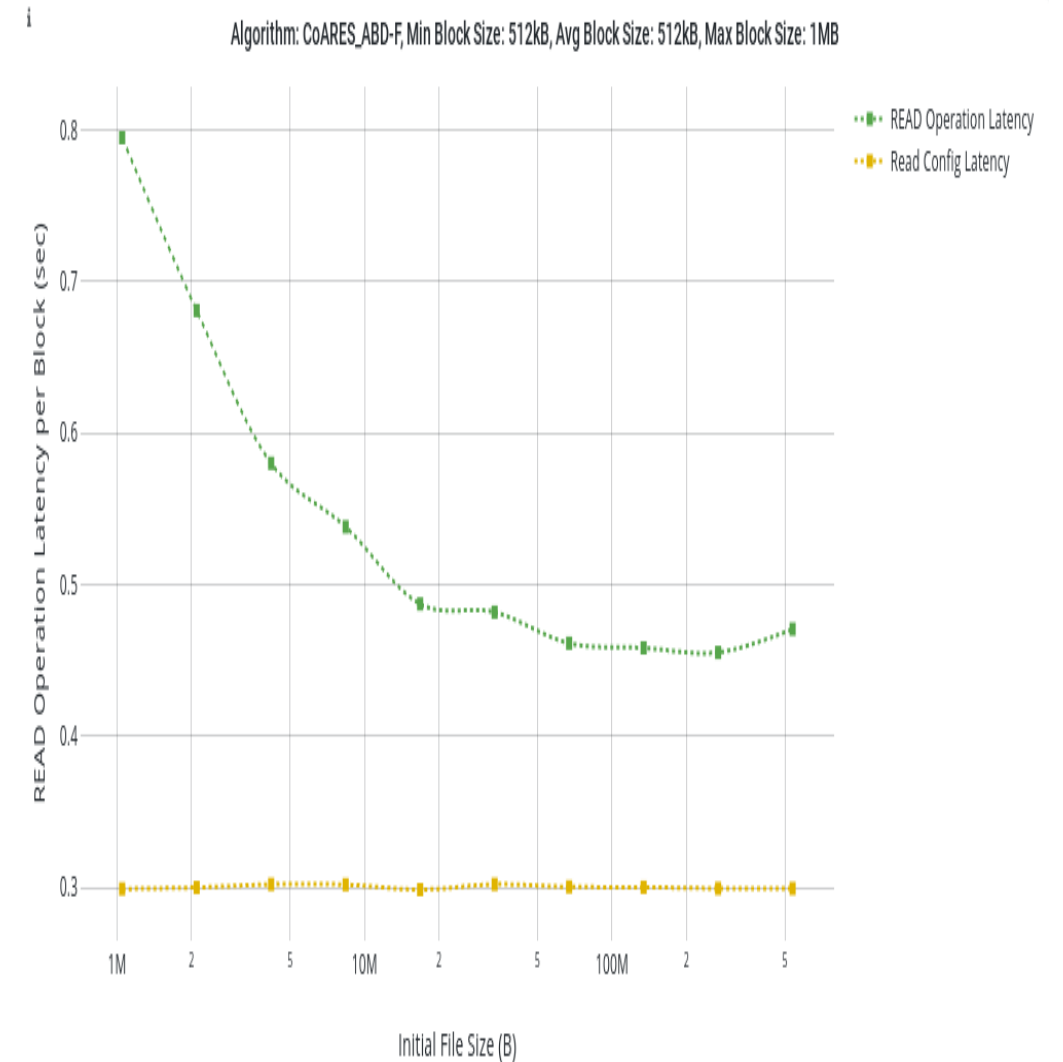
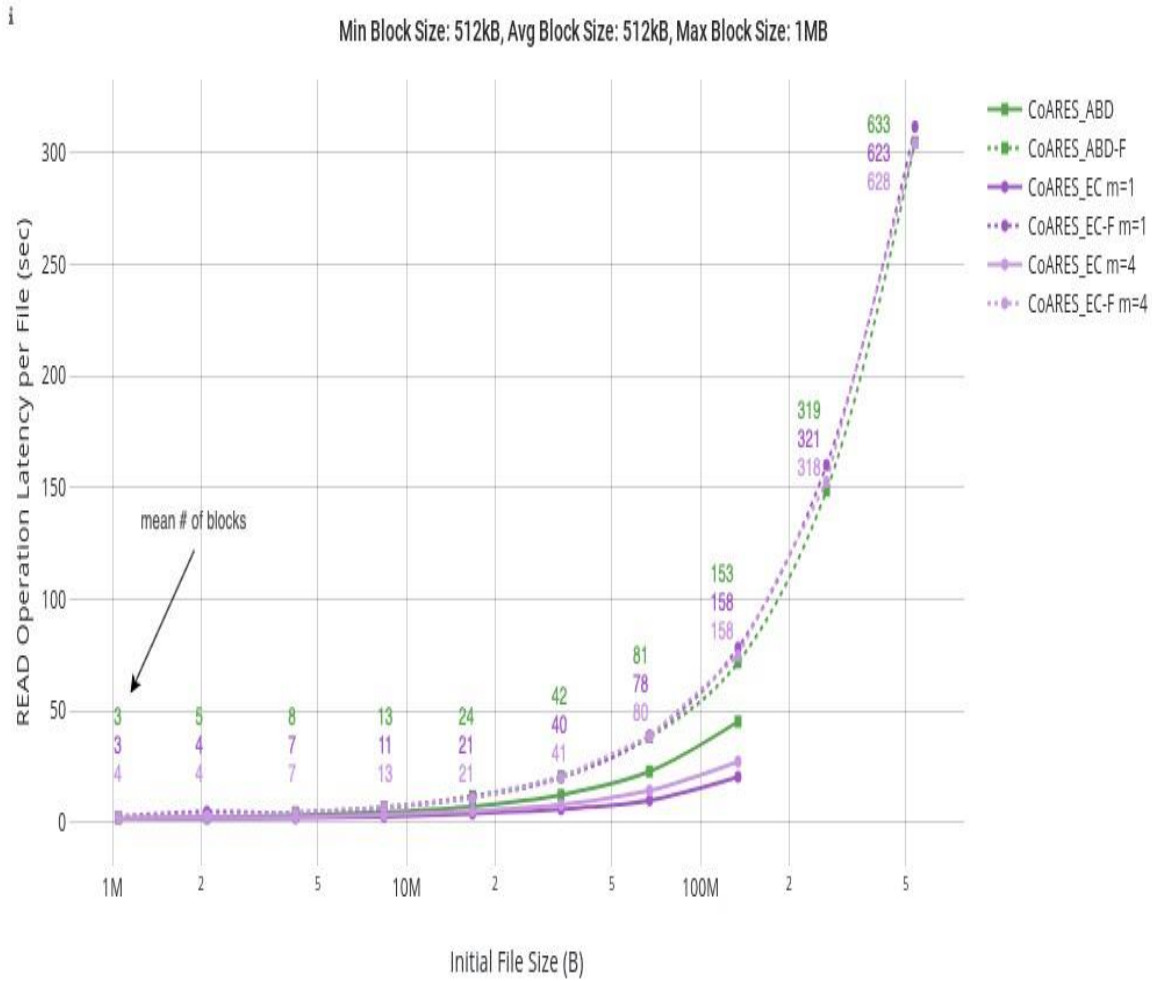
Update Operation Latency per File (sec) vs Initial File Size (2^x B)
wInt:3, rInt:3, #writes:20, #reads:20, #Servers:11, #Writers:5, #Readers:5,
maxBlockSize:1MB, minBlockSize:512KB, avgBlockSize:512KB



Read Operation Latency per File (sec) vs Initial File Size (2^x B)
wInt:3, rInt:3, #writes:20, #reads:20, #Servers:11, #Writers:5, #Readers:5,
maxBlockSize:1MB, minBlockSize:512KB, avgBlockSize:512KB



File Size (AWS)



Overview of Results

- Fragmented algorithms have **significantly lower write** operation latency both in Emulab and AWS
- For the **read** operation latency, AWS results of CoARESF suggest that there is **room for improvement**
- EC-based algorithms are the most scalable as the servers increase
- Trade-off between **block size, operation latency** and **write success rate**
- Trade-off between **operation latency** and the **parity of EC**
- Trade-off between **operation latency** and the **number of writers**

Recap

- Presented and proved correct CoARESF, the first **Dynamic, Robust, Strongly-consistent** DSM that supports **Versioning, (2-level) Data Striping, and High Access Concurrency** for Large Objects
- Data available at <https://github.com/atrigeorgi/fragmentedARES-data.git>

PART **IV** – Enhance the Performance of ARES

covering Stages 7-8 in Methodology

Part IV - Approach and Contribution

- **Identify flaws in DSMs using Distributed Tracing.**
We demonstrate this through the ***ARES*** DSM.
- Develop optimizations based on the found bottlenecks.
- The correctness of optimized ARES is rigorously proven.

Performance Analysis Challenges in DSMs

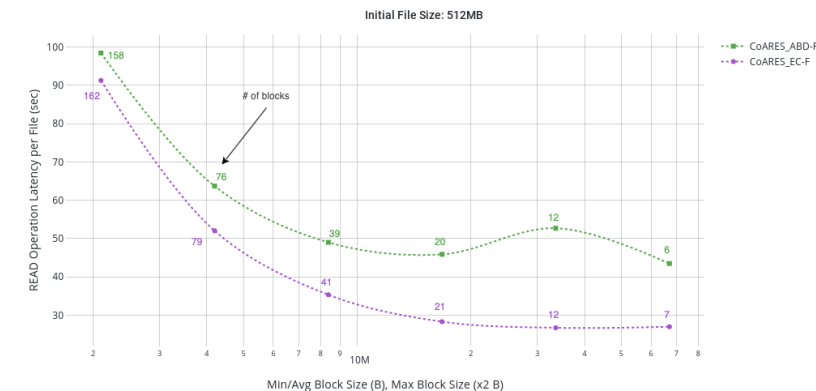
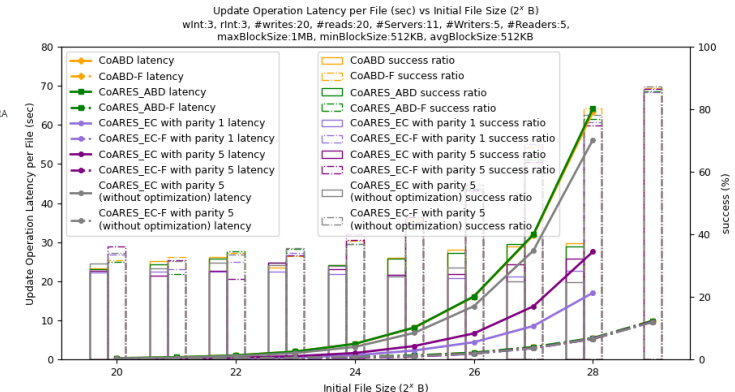
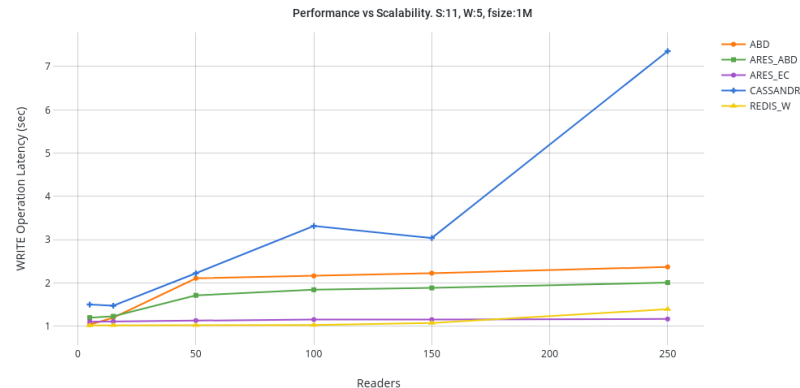
- Identifying performance bottlenecks in complex DSMs can be challenging
- Traditional logging techniques may not provide sufficient insight

```
import os
import logging
from sys import platform
from pythonjsonlogger import jsonlogger

@atrigeorgi
class SetupLogger:

    @atrigeorgi
    def setup_logger(self, logfile, level=logging.DEBUG):
        # Due to race conditions, we sometimes get error.
        # OSError: [Errno 17] File exists: 'log'
```

```
self.logger.debug('READ-COMPLETE-DSMM',
                  extra={"clientID": self.uid, "objectID": file_id, "tag": maxTag, "value": value})
```



“Distributing Tracing is a monitoring technique used to track individual requests as they move across multiple components within a distributed system. It helps to pinpoint where failures occur and what causes poor performance.”

Distributed Tracing – Terminology

- A **trace** represents the entire journey of a request.
- A **span** represents a unit of work within a trace (e.g., procedures, sections of code).
- Tracings tools: Opentemetry, Zipkin, Jaeger.

```
with self.tracer.start_as_current_span("StartWriteRequest-MEMORY"+self.phase) as parent_span:
    with self.tracer.start_as_current_span("Phase1"):
        with self.tracer.start_as_current_span("GetTag"):

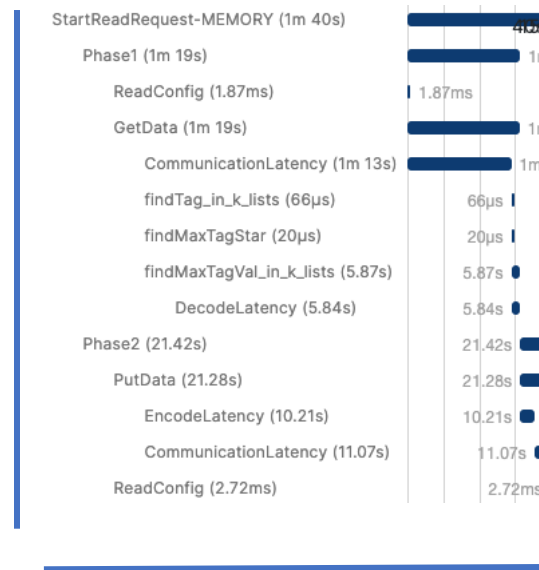
            # increase msg counter for next op
            self.msg_id += 1

            # Creating 'READ' message with the following fields
            # [type, msgID, clientID, objectID]
            message = {'type': 'READ-TAG', 'msgID': self.msg_id, 'clientID': self.uid, 'objectID': file_id}

            with self.tracer.start_as_current_span("CommunicationLatency"):
                # Broadcast it
                self.broadcastMessage(message)

                # Wait for READ-TAG-ACK messages to come from a majority
                msgs_list = self.waitReply(self.majority, "READ-TAG-ACK")
```

Trace



Spans

Tracing the Latencies of Ares: A DSM Case Study

Authors: Chryssis Georgiou¹, Nicolas Nicolaou², Andria Trigeorgi^{1,2}

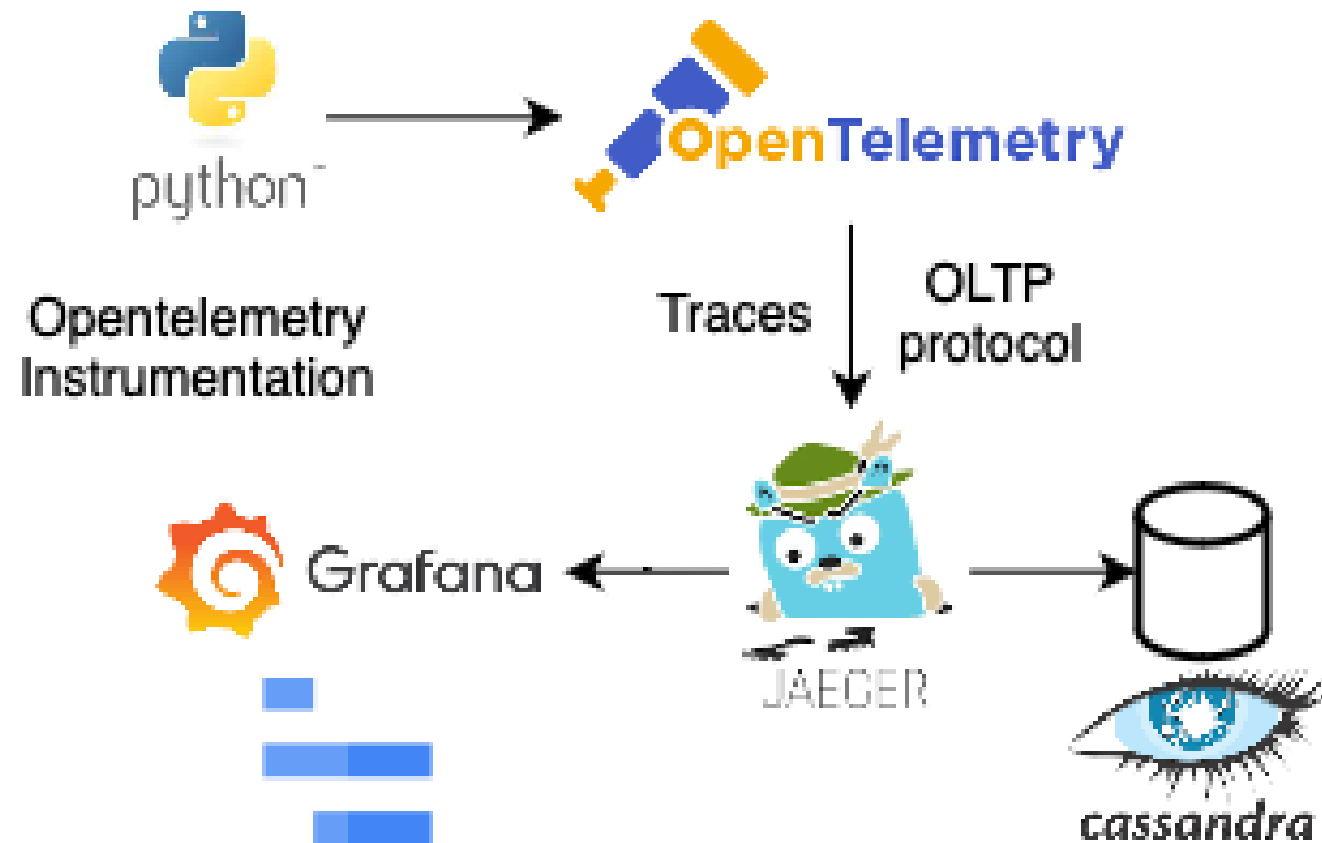
¹University of Cyprus, Nicosia, Cyprus

²Algolysis, Limassol, Cyprus

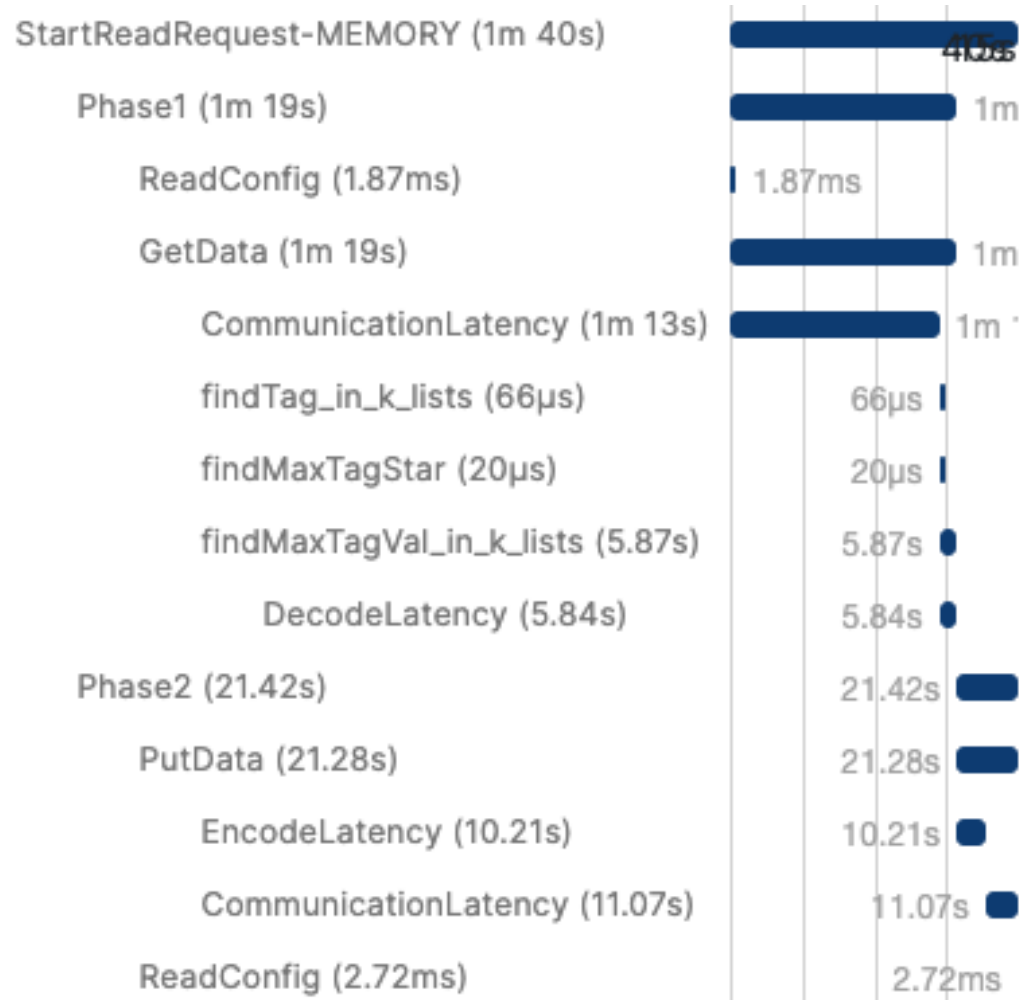
ApPLIED 2024, Nantes, France

Funded by: PHD IN INDUSTRY/1222/0121 and DUAL USE/0922/0048

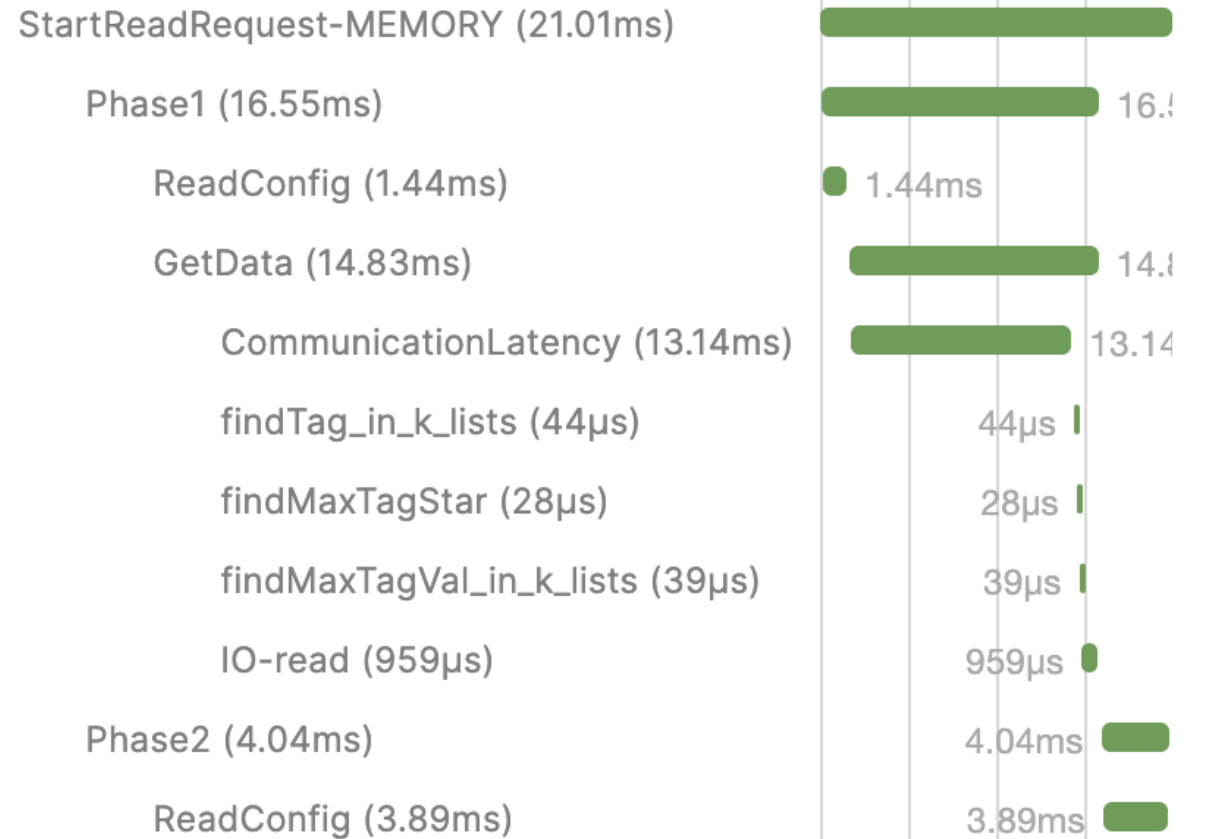
Methodology: ARES Distributed Tracing



File Size

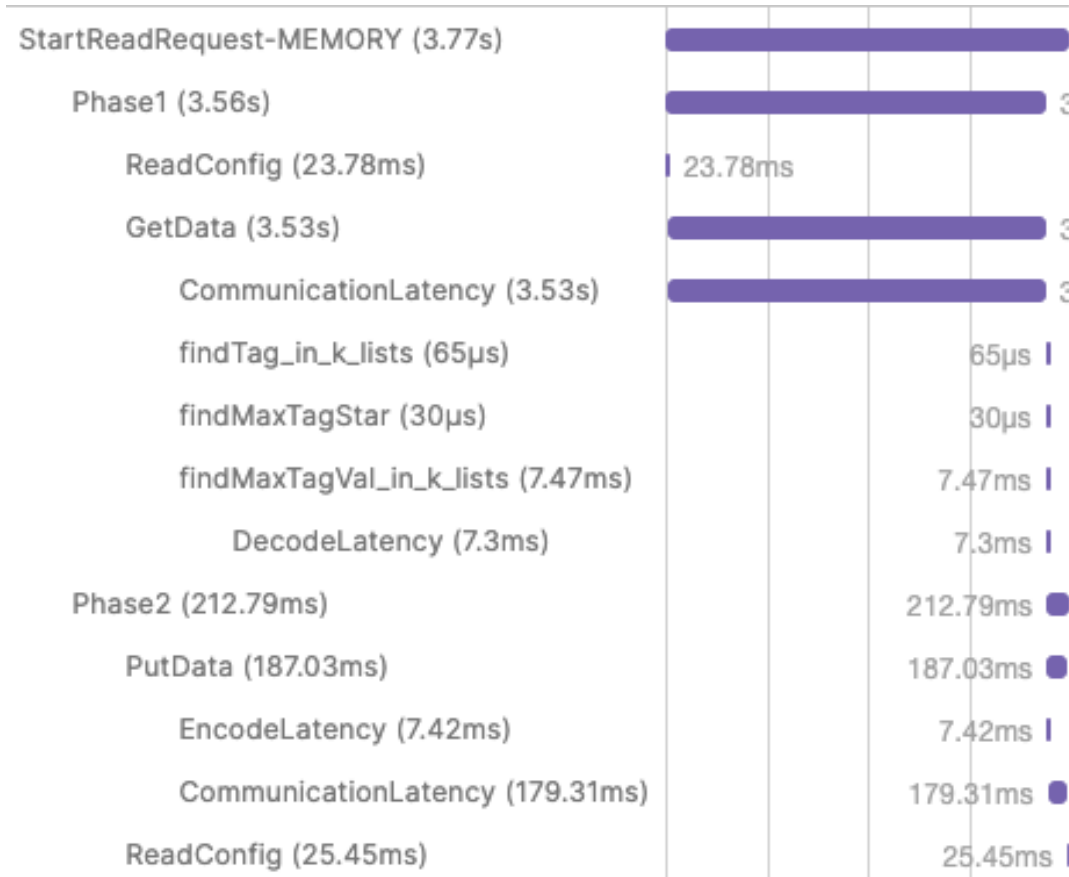


ARESEC, S:11, W:5, R:5, fsize:512MB, Debug Level:DSMM

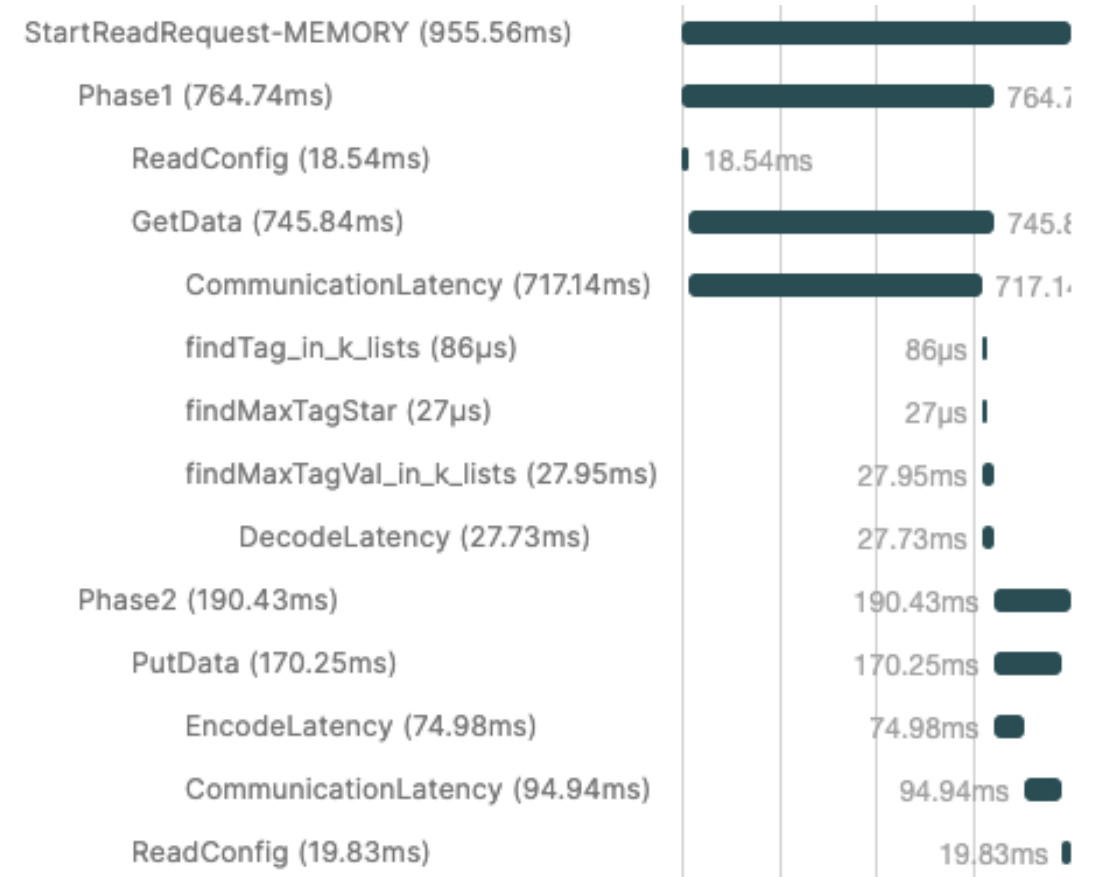


CoARESECF, S:11, W:5, R:5, init fsize:512MB, Debug Level:DSMM

Participation Scalability

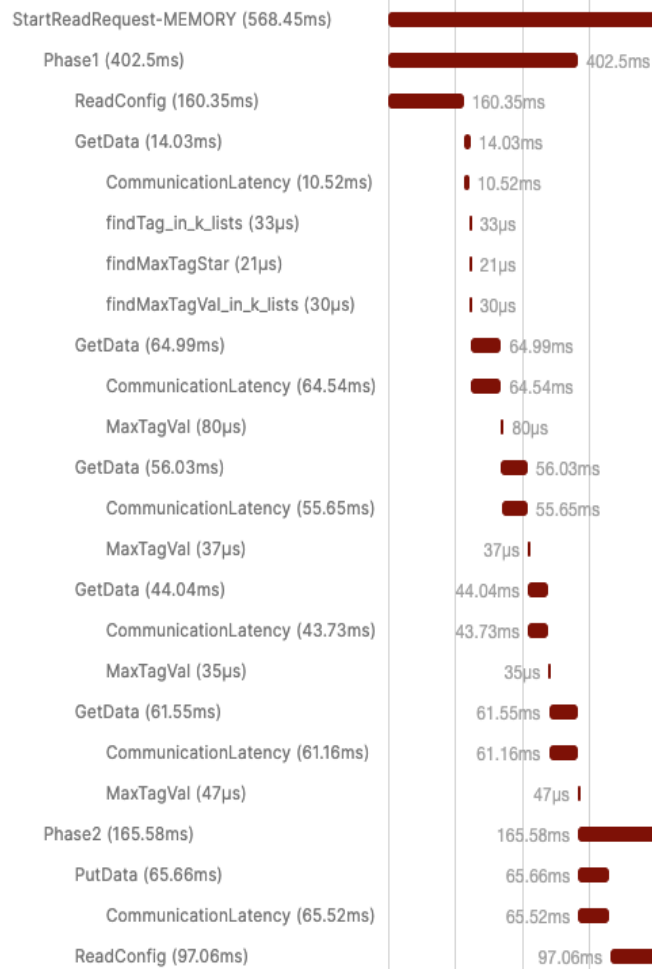


ARESEC, S:3, W:5, R:50, fsize:4MB, Debug Level:DSMM

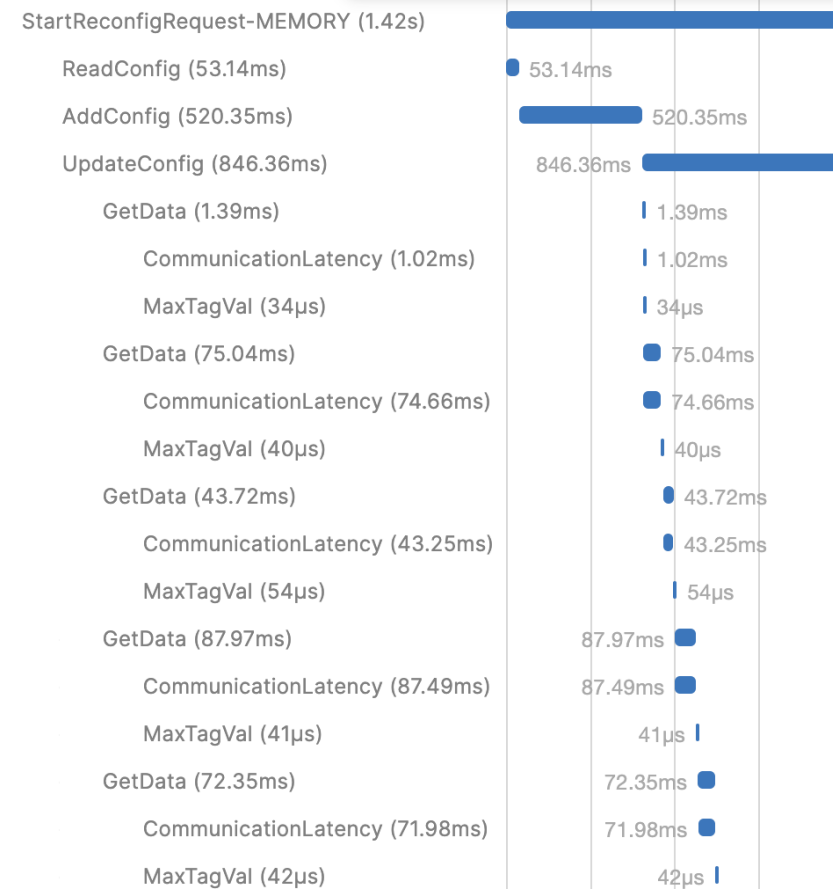


ARESEC, S:11, W:5, R:50, fsize:4MB, Debug Level:DSMM

Longevity



CoAresF, S:11, W:5, R:15, G=5, fsize:4MB,
Debug Level:DSMM



CoAresF, S:11, W:5, R:15, G=5, fsize:4MB,
Debug Level:DSMM

ARES II: Tracing the Flaws of a (Storage) God



Authors: Chryssis Georgiou¹, Nicolas Nicolaou², Andria Trigeorgi^{1,2}

¹University of Cyprus, Nicosia, Cyprus

²Algolysis, Limassol, Cyprus

SRDS 2024, Charlotte, USA

Funded by: PHD IN INDUSTRY/1222/0121 and DUAL USE/0922/0048

Optimisation 1: Piggybacking

Main difference:
skip read-config
for latest config

```
operation write(val), val ∈ V
8:  cseq ← read-config(cseq)
   μ ← max({i : cseq[i].status = F})
10:  cs ← cseq[μ]
   while cs ≠ ⊥ do
12:    τc, Cs ← cs.cfg.get-tag()
    τmax ← max(τc, τmax)
14:    cs, cseq ← find-next-config(cseq, Cs)
   end while
16:  ⟨τ, v⟩ ← ⟨⟨τmax.ts + 1, ωi⟩, val⟩

   λ ← max({i : cseq[i] ≠ ⊥})
18:  cs ← cseq[λ]
   while cs ≠ ⊥ do
20:    Cs ← cs.cfg.put-data(⟨τ, v⟩)
22:    cseq ← read-config(cseq)
    cs, cseq ← find-next-config(cseq, Cs)
   end while
24: end operation
```

Query phase:
embed latest
configs with data

Extra function:
discover the
next config

Optimisation 1: Piggybacking – EC-DAP

Algorithm 3 EC-DAP II implementation

at each process $p_i \in \mathcal{I}$

2: **procedure** c.get-tag()
 send (QUERY-TAG) to each $s \in c.Servers$
 4: **until** p_i receives $\langle t_s, nextC_s \rangle$ from $\lceil \frac{n+k}{2} \rceil$ servers in $c.Servers$
 $C_s \leftarrow \{nextC_s : \text{received } nextC_s \text{ from } s\}$
 6: $t_{max} \leftarrow \max(\{t_s : \text{received } t_s \text{ from } s\})$
 return t_{max}, C_s
 8: **end procedure**

procedure c.get-data()
 10: **send** (QUERY-LIST) to each $s \in c.Servers$
 until p_i receives $List_s, nextC_s$ from $\lceil \frac{n+k}{2} \rceil$ servers in $c.Servers$
 12: $C_s \leftarrow \{nextC_s : \text{received } nextC_s \text{ from } s\}$
 $Tags_{dec}^{\geq k} = \text{set of tags that appears in } k \text{ Lists}$
 14: $t_{max}^{dec} \leftarrow \max(Tags_{dec}^{\geq k})$
 if $Tags_{dec}^{\geq k} \neq \emptyset$ **then**

16: $fragments \leftarrow \{e : \langle \tau, e \rangle \in Lists \ \& \ \tau = t_{max}^{dec}\}$
 if $\nexists \perp \in fragments$ **then**
 18: $v \leftarrow \text{decode value for } t_{max}^{dec}$
 else
 20: $v \leftarrow \perp$
 return $\langle t_{max}^{dec}, v \rangle, C_s$
 22: **end procedure**

procedure c.put-data($\langle \tau, v \rangle$)
 24: $code_elems = [(\tau, e_1), \dots, (\tau, e_n)], e_i = \Phi_i(v)$
 send (PUT-DATA, $\langle \tau, e_i \rangle$) to each $s_i \in c.Servers$
 26: **until** p_i receives $nextC_s$ from each server $s \in S_g$ s.t. $|S_g| = \lceil \frac{n+k}{2} \rceil$
 and $S_g \subset c.Servers$
 28: $C_s \leftarrow \{nextC_s : \text{received } nextC_s \text{ from each } s \in S_g\}$
 return C_s
 30: **end procedure**

requests max τ
and nextC in one go

returns max τ
and servers' nextC

requests data list
and nextC

If nextC is finalized, servers sends
only the tag and their nextC

requests data update
and nextC

returns all
servers' nextC

Optimisation 2: Garbage Collection

```

procedure gc-config(cseq)
38:   $\mu \leftarrow \max(\{i : cseq[i].status = F\})$ 
    CID  $\leftarrow \{i : cseq[i] \neq \perp \wedge i < \mu\}$ 
40:  for id in CID do
    send (GC-CONFIG, next) to each s  $\in$  cseq[id].Servers
42:  until  $\exists Q, Q \in cseq[id].Quorums$  s.t. reci receives ACK from  $\forall s \in Q$ 
    /* remove the {id, cseq[id]} */
44:  cseq  $\leftarrow cseq \setminus \{id, cseq[id]\}$ 
    return cseq
46: end procedure
    
```

send GC request to servers

Last finalized config

older configs from last finalized config

Remove the GC configs

server updates nextC to point to the last finalized config

the server sets the data of config to \perp

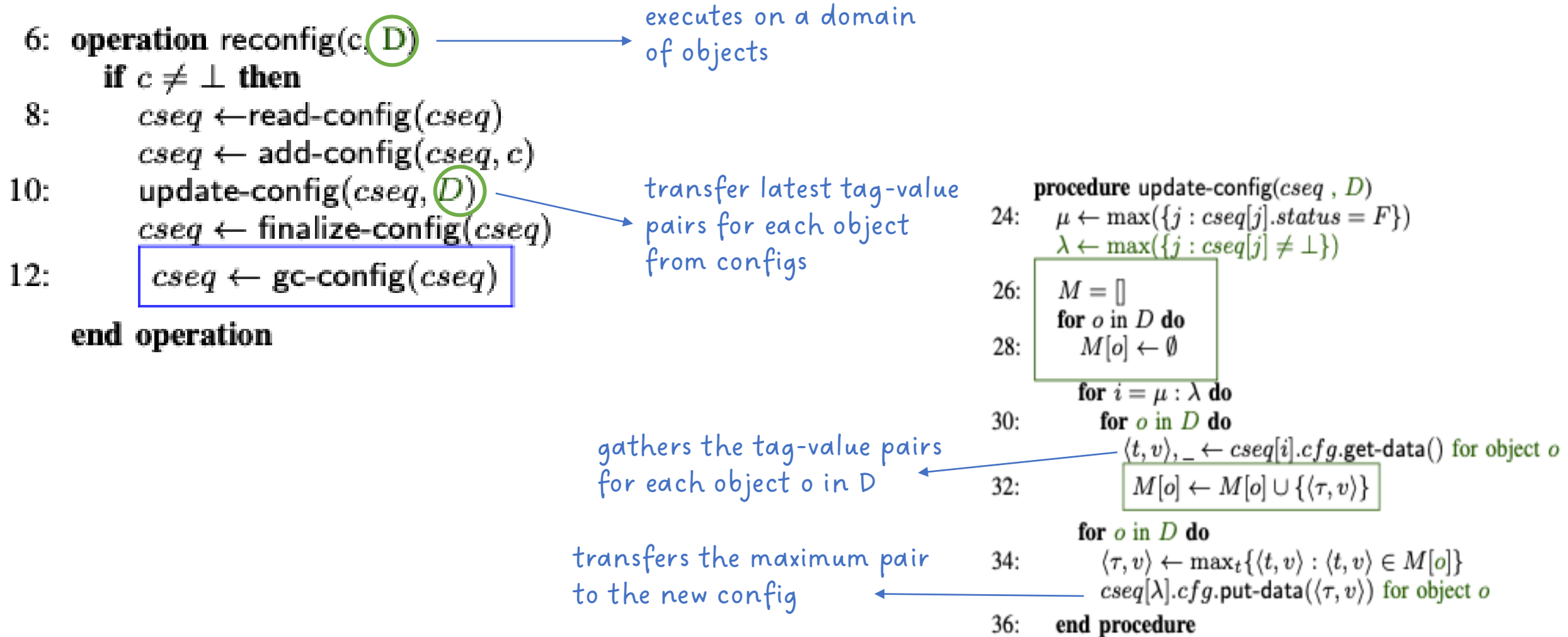
EC-DAP

```

end receive
Upon receive (GC-CONFIG, cfgTin)si, ck from q
  if cfgTin.cfg.ID > nextC.cfg.ID then
    nextC  $\leftarrow$  cfgTin
    for  $\tau, e$  in List do
      List  $\leftarrow$  List  $\setminus$  { $\langle \tau, e \rangle$ }
      List  $\leftarrow$  List  $\cup$  { $\langle \tau, \perp \rangle$ }
    send ACK to q
  end receive
    
```

Server' Response

Optimisation 3: Batching



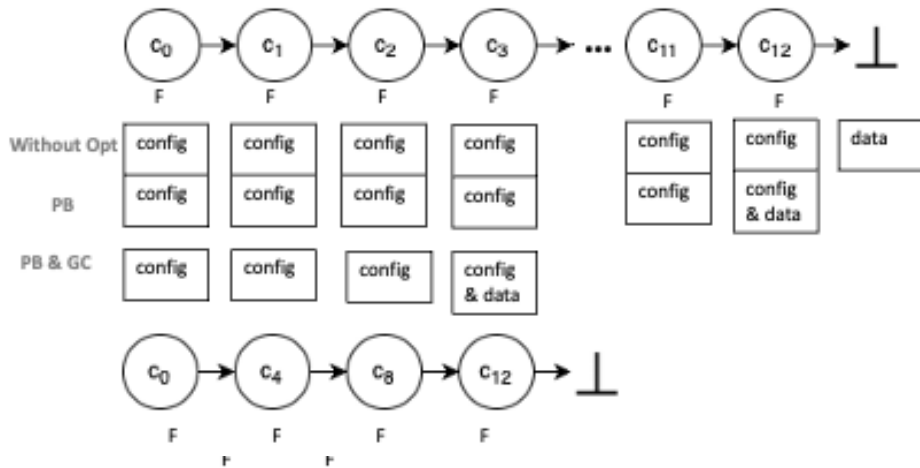
Optimization Results – Piggyback

| alg./ f_{size} | CoAREsabdf | CoAREsabdf PB | CoAREsecf | CoAREsecf PB |
|------------------|------------|-----------------|-----------|----------------|
| 1MB | 284ms | 278ms | 149ms | 142ms |
| 256MB | 9s | 5s (44%) | 9.65s | 3.82s (60%) |
| 512MB | 21.8s | 15.2s (30%) | 23.2s | 10.9s (53%) |

TABLE I: READ Operation - File Size - S:11, W:5, R:5

- **CoARESF (256 MB & 512 MB):** Significant performance drops without optimization.
 - Non-optimized: **4 rounds per block** with double read-config.
 - With PB Optimization: Reduced to **2 rounds**, lower read latency.

Optimization Results – Garbage Collection



| alg./fsize | ARES | ARES PB | ARES PB&GC | CoARES | CoARES PB | CoARES PB&GC | CoARESf | CoARESf PB | CoARESf PB&GC |
|---|-------|------------|---------------|--------|--------------|-----------------|---------|---------------|------------------|
| 1.1 Pending Reconfiguration & 1 Finalized | | | | | | | | | |
| 1MB | 159ms | 494ms | 107ms | 162ms | 506ms | 110ms | 181ms | 191ms | 127ms |
| 64MB | 5.57s | 27.4s | 5.58s | 5.81s | 26.8s | 5.73s | 6.78s | 6.62s | 6.61s |
| 1.2 Finalized Reconfiguration | | | | | | | | | |
| 1MB | 159ms | 166ms | 119ms | 163ms | 167ms | 122ms | 186ms | 193ms | 135ms |
| 64MB | 5.80s | 5.76s | 5.71s | 5.88s | 5.98s | 5.82s | 6.92s | 6.73s | 6.74s |

TABLE II: READ Operation - Reconfigurations - S:11, W:1, R:10:, G:4

Scenario 1:

- PB version has the worst latency, since transfers data and config in 11 round trips.
- PB with GC is fastest, since it updates pointers reducing actions.
- CoARESf & Larger Objects (64MB): No differences between versions since the first block finds the latest config and the next block starts from that config.

Scenario 2:

- Original vs. PB has similar performance with one extra round trip.
- PB with GC is faster, skipping every 4 configurations, fewer rounds needed.

Recap

- Used tracing to pinpoint inefficiencies by monitoring individual procedures.
- Develop optimizations, leading to ARES II.
- Show the **correctness** of ARES II and conduct performance **evaluations** to showcase its improvements over ARES.
- Distributed tracing is crucial for diagnosing and resolving performance issues in DSM algorithms.

Optimization Strategies

- **Piggy-backing**: Integrating configurations with read/write messages to expedite configuration discovery.
- **Garbage Collection**: Eliminating obsolete configurations for quicker access to the latest data.
- **Data Batching**: A single reconfiguration across multiple objects to enhance efficiency.

Conclusion

- CoBFS has the following advantages:
 - High Concurrent accesses
 - Strong Consistency
 - Large file size (tested up to 1GB file)
- CoARESF has the following advantages:
 - the first dynamic DSM with coverable fragmented objects & 2-Level of Striping
 - Optimized for High Concurrent accesses
- Theoretical principles illustrated by extensive experimental evaluation
 - atomic consistency, data striping, erasure coding, access to the same files under heavy concurrency, fault-tolerance, reconfiguration

Ongoing and Future Work

- Design Reconfiguration Orchestration Strategies (ROS) for dynamic DSM **Ongoing**
 - **When** to invoke reconfigurations
 - **How** to reconfigure
 - Ensure that the system remains operational despite server failures.
 - Improve performance by replacing older servers with more powerful ones.
- Server Failure Prediction **Ongoing**
 - Monitor environmental parameters of servers (performance, capacity, availability, and health).
 - Threshold-based approaches for determine when to reconfigure.
 - Predict failures using Machine Learning.
- Develop, Deploy, and Evaluate Web Platform **Ongoing**
 - Deploy and get access to configurations of DSM with ROS support by specifying servers.
 - Manage existing configurations.
 - Get access to existing configuration for reading and writing data objects either through the platform directly or through a third-party application using appropriate security tokens.
- Fully-functional Distributed Storage System with Security Guarantees
- Extensive Experimental Evaluation
 - Compare CoARESF against commercial solutions that employ a striping method.
 - Integrate any commercial solution into CoBFS.

The screenshot shows a web interface for configuring DSM instances. It features two main sections, one for 'DSMid 12' and one for 'DSMid 29'. Each section includes a 'Select Prometheus Instance' dropdown, a 'Select RAFT Nodes' input, and a 'Select ADSM Nodes' input. Below these, there are four columns of configuration parameters: RAFT TCP Port, RAFT HTTP Port, ADSM TCP Port, and ADSM HTTP Port. Each column has a text input for the port number, a dropdown for the protocol, and a text input for the parity blocks. The 'DSMid 12' section has a 'Save' button, while the 'DSMid 29' section has 'Add More', 'Save', and 'Save and Create' buttons. The interface is clean and modern, with a light blue background and white text.

Publications

Journal Articles

1. Ares: Adaptive, reconfigurable, erasure coded, atomic storage. N. Nicolaou, V. Cadambe, N. Prakash, **A. Trigeorgi**, K. M. Konwar, M. Medard, and N. Lynch. *Journal of ACM Trans. Storage* 2022.
2. Boosting Concurrency and Fault-Tolerance for Reconfigurable Shared Large Object. A. F. Anta, C. Georgiou, T. Hadjistasi, N. Nicolaou, E. Stavarakis, and **A. Trigeorgi**. *Under Submission* 2024.

Articles in Conferences

3. Robust and strongly consistent distributed storage systems. **A. Trigeorgi**. In *Proc. of RCIS* 2021.
4. Fragmented Object : Boosting Concurrency of Shared Large Objects. A. Anta, C. Georgiou, N. Nicolaou, T. Hadjistasi, E. Stavarakis, and **A. Trigeorgi**. In *Proc. of SIROCCO* 2021.
5. Fragmented ARES: Dynamic Storage for Large Objects. C. Georgiou, N. Nicolaou, **A. Trigeorgi**. In *Proc. of DISC* 2022. (arXiv:2201.13292)
6. Invited paper: Towards Practical Atomic Distributed Shared Memory: An Experimental Evaluation. **A. Trigeorgi**, N. Nicolaou, C. Georgiou, T. Hadjistasi, E. Stavarakis, V. Cadambe and B. Urgaonkar. In *Proc. of SSS* 2022.
7. Robust and Consistent Distributed Storage as a Service. **A. Trigeorgi**. In *Proc. of GEC* 2024.
8. Tracing the Latencies of Ares: A DSM Case Study. C. Georgiou, N. Nicolaou, and **A. Trigeorgi**. In *Proc. Of ApPLIED* 2024.
9. Ares II: Tracing the Flaws of a (Storage) God. C. Georgiou, N. Nicolaou, and **A. Trigeorgi**. In *Proc. Of SRDS* 2024.

Acknowledgments

Collaborators:

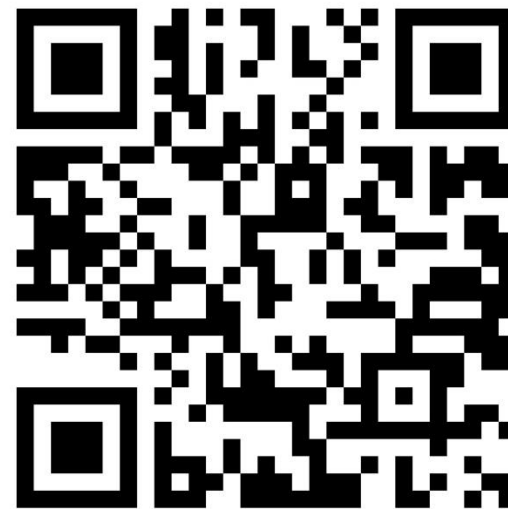
- Prof. Chryssis Georgiou, University of Cyprus (Supervisor)
- Dr. Nicolas Nicolaou, Algolysis Ltd
- Dr. Efstathios Stavrakis, Algolysis Ltd
- Prof. Antonio Fern´andez Anta, IMDEA Networks
- Dr. Theophanis Hadjistasis, Algolysis Ltd

Examination Committee:

- Prof. Anna Philippou, University of Cyprus
- Assoc. Prof. George Pallis, University of Cyprus
- Prof. Gregory Chockler, University of Surrey, UK
- Assoc. Prof. Herodotos Herodotou, Cyprus University of Technology

Thank you!

For more information you can see the websites of our related projects:



Related Projects

- **PHD IN INDUSTRY/1222/0121**
- **DUAL USE/0922/0048**
- **POST-DOC/0916/0090**
- **EU's NGIAtlantic.eu** cascading grant agreement no. OC4-347

CORRECTNESS

Data Access Primitives (DAPs)

- Operation Ordering: logical tags $\tau = \langle ts, wid \rangle$
- For a configuration c , any client p may invoke any of the following DAPs:
 - D1. $c.get - tag()$: returns a tag $\tau \in T$
 - D2. $c.get - data()$: returns a tag – value pair $(\tau, v) \in T \times V$
 - D3. $c.put - data(\langle \tau, v \rangle)$: the tag – value pair $(\tau, v) \in T \times V$ as argument
- **Property 1:** DAP consistency conditions
 - C1: if a put-data($\langle \tau, v \rangle$) precedes a get-data/get-tag operation that returns τ' , then $\tau' \geq \tau$
 - C2: if a get-data returns $\langle \tau', v' \rangle$, then there exists put-data($\langle \tau', v' \rangle$) that precedes or is concurrent to the get-data operation

DAP-ABD and DAP-EC satisfy Property 1

Correctness of CoARES

Theorem

CoARES implements a linearizable coverable object, given that the DAPs implemented in any configuration c satisfy DAP Consistency Conditions.

Proof challenges: CoARES satisfies coverability despite any reconfiguration in the system.

- New values are not overwritten (by writes associated with older versions)
- Versions are unique
- Eventually a single version path prevails

Correctness of CoARESF

Theorem

CoARESF implements a linearizable coverable fragmented object.

Proof challenges

- f remains connected and composed of the most recent blocks, despite concurrent read/write and reconfig operations
- Each block may exist in different configurations and be accessed by different DAPs
- Show that fragmented linearizable coverability cannot be violated

From ARES to ARES II

- **Piggy-backing:** Integrating configurations with read/write messages to speed up configuration discovery.
- **Garbage Collection:** Eliminating obsolete configurations for quicker access to the latest data.
- **Data Batching:** A single reconfiguration across multiple objects to enhance efficiency.

Correctness of ARES II

- The correctness of ARES II depends on the correct implementation of both **EC-DAP II** and the **Reconfiguration protocol**.

Correctness of EC-DAP II

Theorem

Revised Property 1 to accommodate the fact that get-data can return a tag associated with either a value from V or \perp .

Proof Challenges

- **Tag Guarantee:** get-data() always returns a tag \geq any previous put-data().
- **Value Origin:** Values from **get-data()** come from **put-data()**, the initial value, or \perp (due to garbage collection).

Reconfiguration Protocol Properties

Theorem

For completed actions π_1 and π_2 , where $\pi_1 \rightarrow \pi_2$:

- a. **Configuration Uniqueness**
- b. **Subsequence**
- c. **Sequence Progress**

Proof Challenges

- All processes must have the same config at index k .
- Next configs must point to a higher index than the current one.
- Future actions must start from a config at or above the last finalized config.

Atomicity

Theorem

ARES II guarantees atomicity, given that the DAPs implemented in any configuration c satisfy Property 1 and satisfy the reconfiguration properties.

Proof Challenges

- Ensure get-data returning \perp does not prevent the read operation from retrieving a valid **non- \perp** value.
- The read operation must continue until it finds a finalized (**not GC**) config with the highest tag with a non- \perp value.
- **Batching** must behave like multiple reconfigurations, maintaining the correctness of read/write operations.

Correctness of EC-DAP II

Theorem

Revised Property 1: DAP consistency conditions

C1: if a put-data($\langle \tau, v \rangle$) precedes a get-tag (or get-data) operation that returns τ' (or $\langle t', v' \rangle$ / $\langle t', \perp \rangle$), then $\tau' \geq \tau$

C2: if a get-data returns $\langle \tau', v' \rangle$ or $\langle t', \perp \rangle$, then there exists put-data($\langle \tau', v' \rangle$) that precedes or is concurrent to the get-data operation

Proof Challenges

- The tag returned by get-data() is guaranteed to be \geq the tag from any prior put-data().
- Values returned by get-data() are either from a put-data() operation, the initial object value, or \perp (due to garbage collection).

Reconfiguration Protocol Properties

Theorem

For completed actions π_1 and π_2 , where $\pi_1 \rightarrow \pi_2$:

- a. **Configuration Uniqueness:** Configuration sequences in any two processes are identical at any common index i .
- b. **Subsequence:** The configuration sequence x observed by π_1 is a subsequence of the sequence y observed by π_2 , i.e. $\lambda(x) \leq \lambda(y)$.
- c. **Sequence Progress:** Finalized configurations progress in order.

Proof Challenges

- Must have the same configuration at index k , regardless of how they received it (propose operation or DAP/read-config)
- If the last configuration has index k , subsequent configurations with index j must have a nextC pointing to a configuration with index $i \geq j+1$.
- If the last finalized configuration has index k , future actions must start from a configuration with index j such that $j \geq k$.

Atomicity

Theorem

ARES II guarantees atomicity, given that the DAPs implemented in any configuration c satisfy Property 1 and satisfy the reconfiguration properties.

Proof Challenges

- Ensure that get-data returning \perp does not prevent the read operation from retrieving a valid non- \perp value.
- The read operation must continue until it finds a finalized configuration that is not garbage collected, yielding the highest tag with a non- \perp value.
- Batching is indistinguishable from multiple reconfigurations, ensuring the correctness of read/write operations.

EXTRA SLIDES

| Algorithm/ System | Data scalability | Data access Concurrency | Consistency guarantees | Versioning | Data Striping | Non-blocking Reconfig |
|---------------------------------------|---------------------|--|-----------------------------------|------------|-------------------------|----------------------------|
| GFS | YES | concurrent appends | relaxed | YES | YES | YES (short downtime) |
| HDFS | YES | files restrict one writer at a time | atomic (centralized) | NO | YES | YES |
| Cassandra | YES | YES | tunable (default= eventual) | YES | NO | NO |
| Dropbox | YES | creates conflicting copies | eventual | YES | YES | N/A |
| Colossus | YES | concurrent appends | relaxed | YES | YES | YES |
| Blobseer | YES | YES | atomic (centralized) | YES | YES | YES |
| Tectonic | YES | files restrict one writer at a time | read-after-write | YES | YES | YES |
| CoABD | NO | NO | atomic | YES | NO | NO |
| CoBFS (using ABD) | YES | YES | atomic | YES | YES | NO |
| LDR | YES | NO | atomic | NO | NO | NO |
| RAMBO/ DYNASTORE SMSTORE/SPSNSTORE | NO | NO | atomic | NO | NO | YES |
| ARESABD | NO | NO | atomic | NO | NO | YES |
| ARESEC | NO | NO | atomic | NO | YES | YES |
| CoARESABD | NO | NO | atomic | YES | NO | YES |
| CoARESEC | NO | NO | atomic | YES | YES | YES |
| CoARESABDF | YES | YES | atomic | YES | YES | YES |
| CoARESECF | YES | YES | atomic | YES | YES (2 lvl striping) | YES |

Distributed Storage Systems

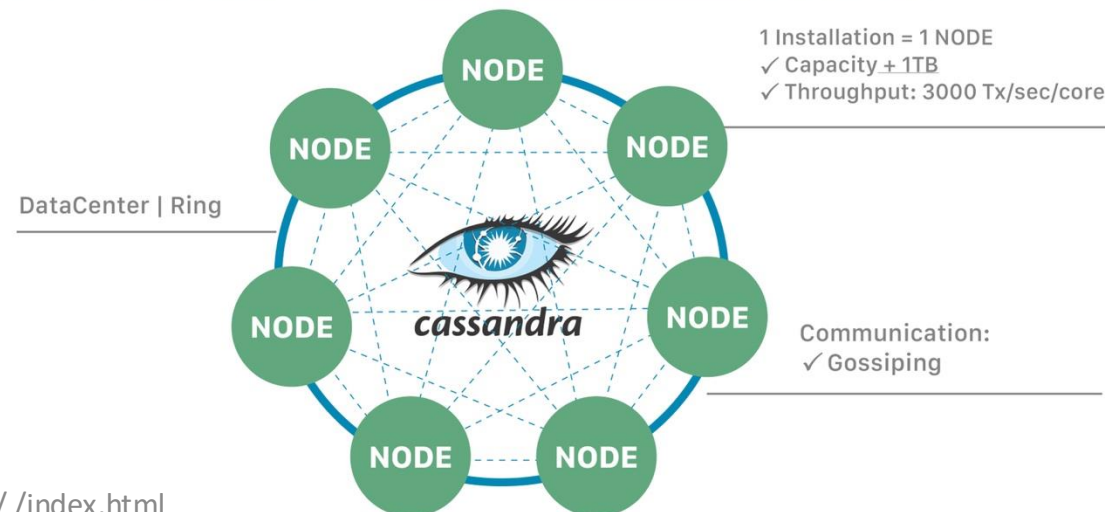
| System | Data scalability | Data access Concurr. | Consist. guarantees | Versioning | Data Striping | Non-blocking Reconfig. |
|-----------|------------------|-------------------------------------|-----------------------------|------------|---------------|------------------------|
| GFS | YES | concurrent appends | relaxed | YES | YES | YES (short downtime) |
| Colossus | YES | concurrent appends | relaxed | YES | YES | YES |
| HDFS | YES | files restrict one writer at a time | atomic centralized | NO | YES | YES |
| CASSANDRA | YES | YES | tunable (default= eventual) | YES | NO | NO |
| DROPBOX | YES | creates conflicting copies | eventual | YES | YES | N/A |
| REDIS | YES | YES | eventual | YES | NO | NO |
| BLOBSEER | YES | YES | atomic centralized | YES | YES | YES |
| TECTONIC | YES | files restrict one writer at a time | Read-your- writes | YES | YES | YES |

CASSANDRA

- **What is Cassandra?** a key-value Distributed Database.
- **Why Cassandra?** availability, high performance, horizontal scalability.
- **How it works?** Gossip protocol, Peer-to-peer communication in a Ring topology.
- Tunable Consistency = number of nodes to acknowledge an operation (default=Eventually). Can support strong consistency.
- Tunable Replication Factor (# of copies).

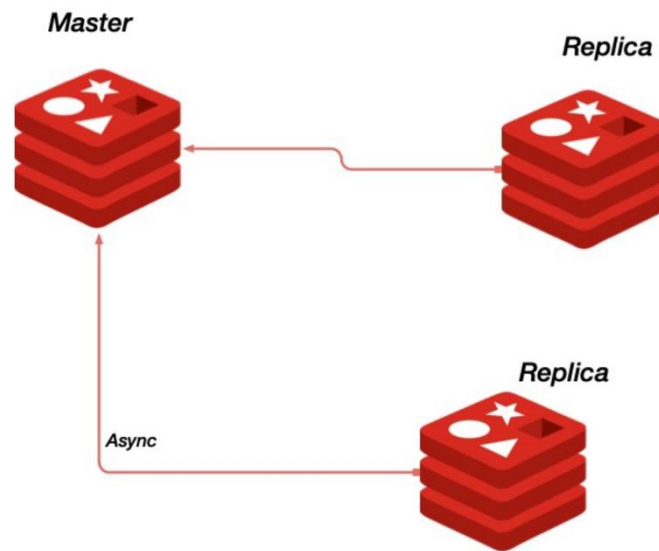
We set RF=n, CL=majority ($n/2+1$) - atomic

ApacheCassandra™ = NoSQL Distributed Database



REDIS

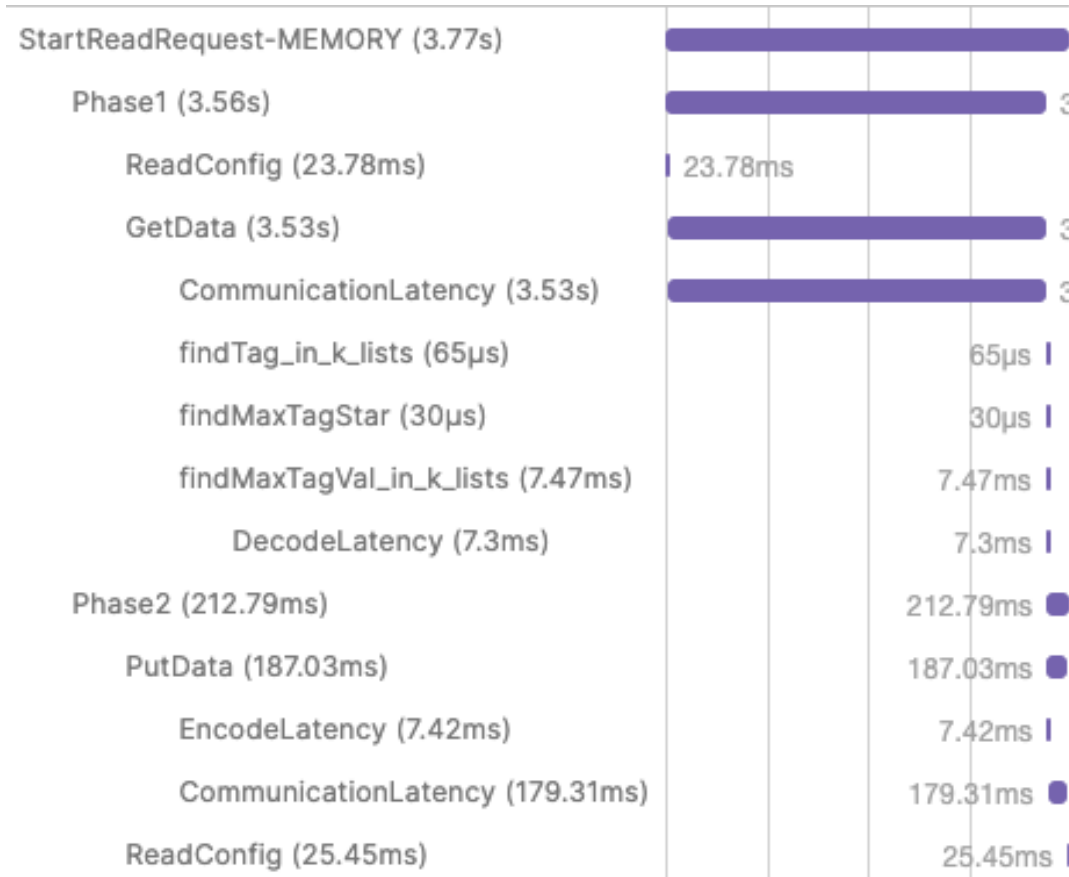
- **What is Redis?** Ultra-fast in-memory key-value store. Used as database, cache, and for simple apps.
- It has memory limitations.
- “WAIT” command for synchronous replication
(We set a majority ($n/2 + 1$) of waiting write acks).
- It provides Eventual consistency.
- **Replication:** Master ensures that one or more slaves becomes exact copies of it. Clients can connect to the master or to the slaves. Slaves are read only by default.



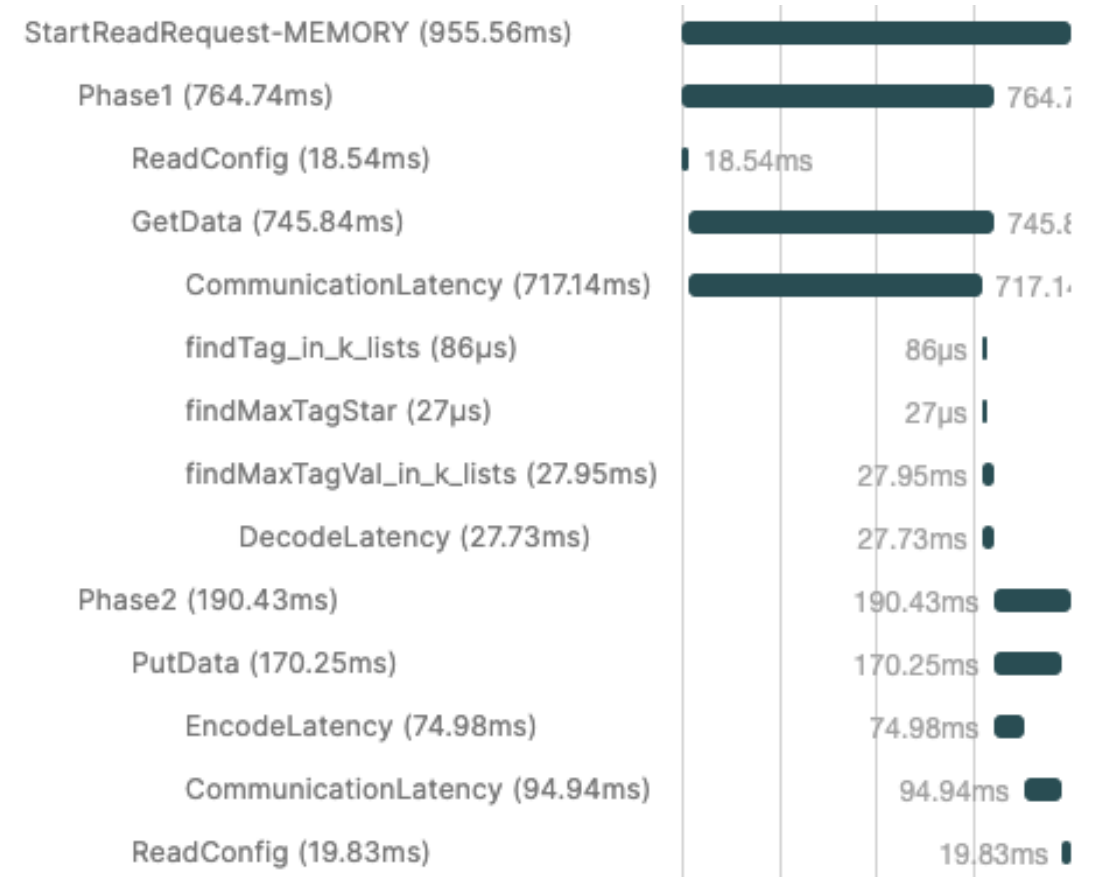
Trade-offs

- **Block size of FM.** data striping performance highly depends on the block size.
- **Parity of EC.** the further increase of the parity (and thus higher fault-tolerance) the larger the latency.
- **Parameter δ of EC.** $\delta = \text{\#writers}$ & each server sends $\delta + 1$ concurrent values in the first phase \rightarrow as the \#writers increases, the latency also increases

Participation Scalability

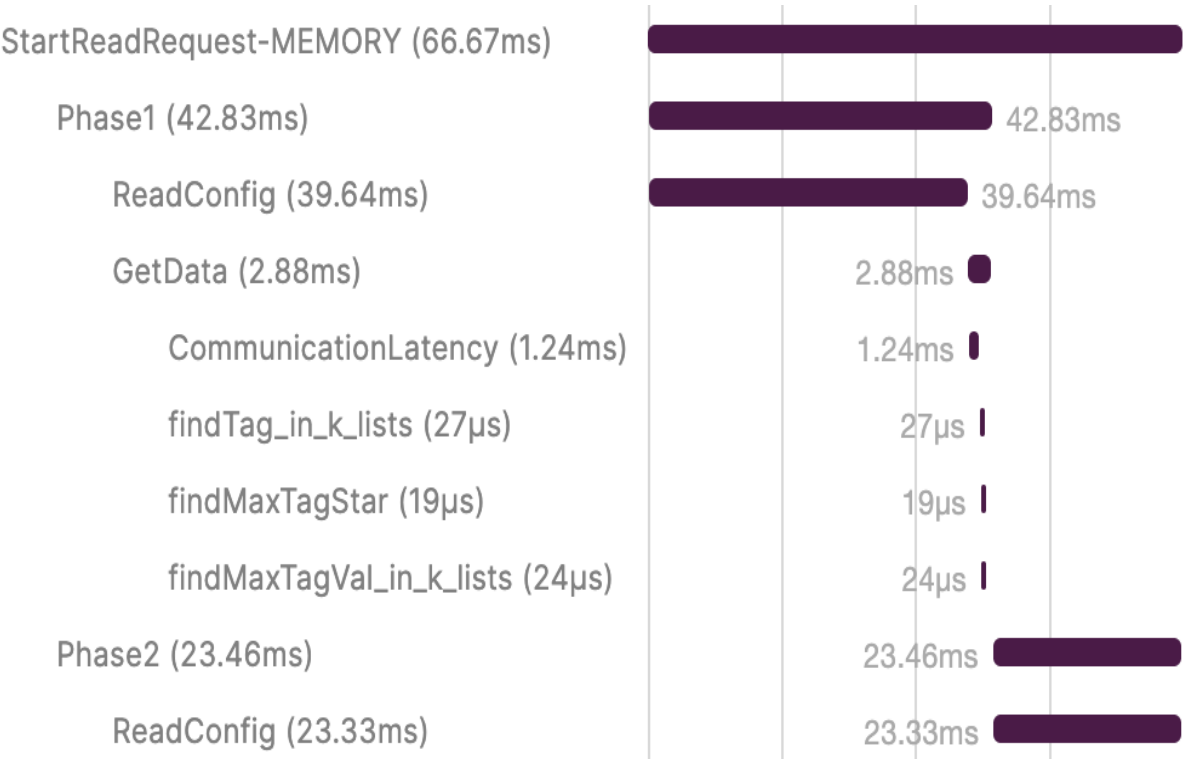


ARESEC, S:3, W:5, R:50, fsize:4MB, Debug Level:DSMM

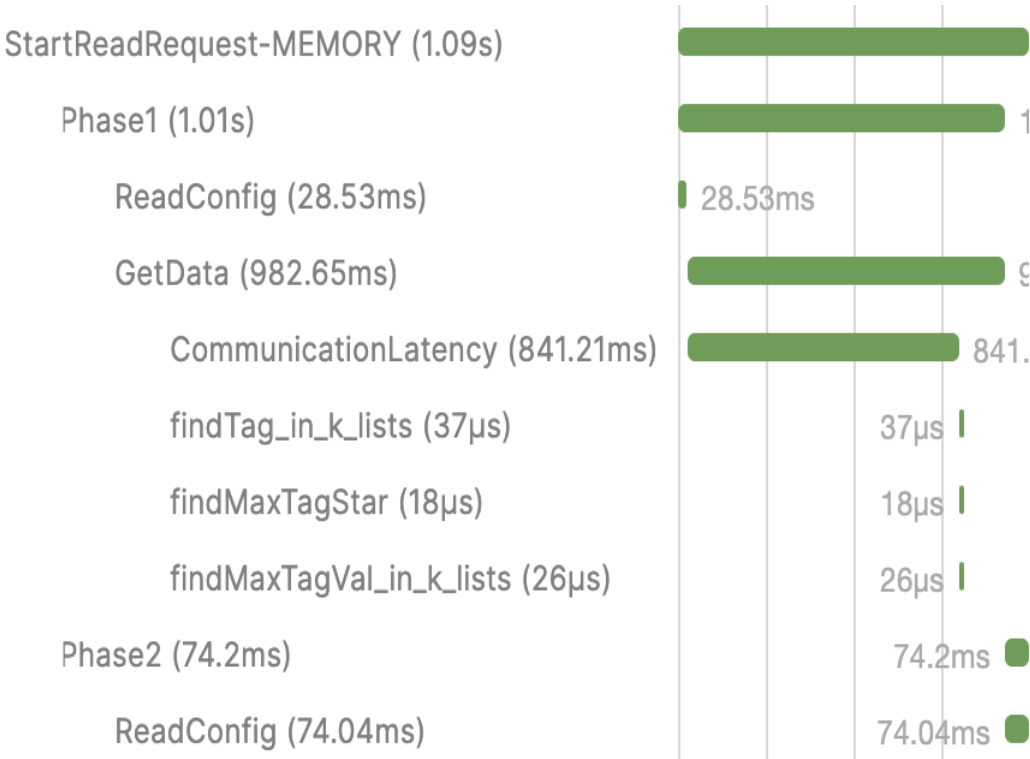


ARESEC, S:11, W:5, R:50, fsize:4MB, Debug Level:DSMM

Block Sizes

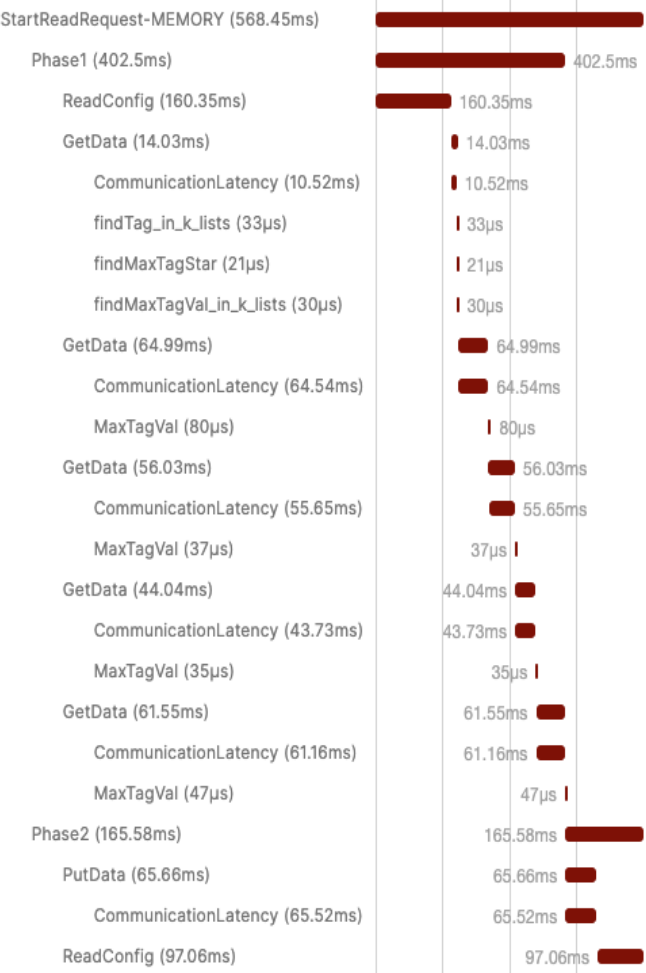


CoAresECF, S:11, W:5, R:5, fsize:512MB, Min/Avg Block Size:2MB, max Block Size:4MB, Debug Level:DSMM

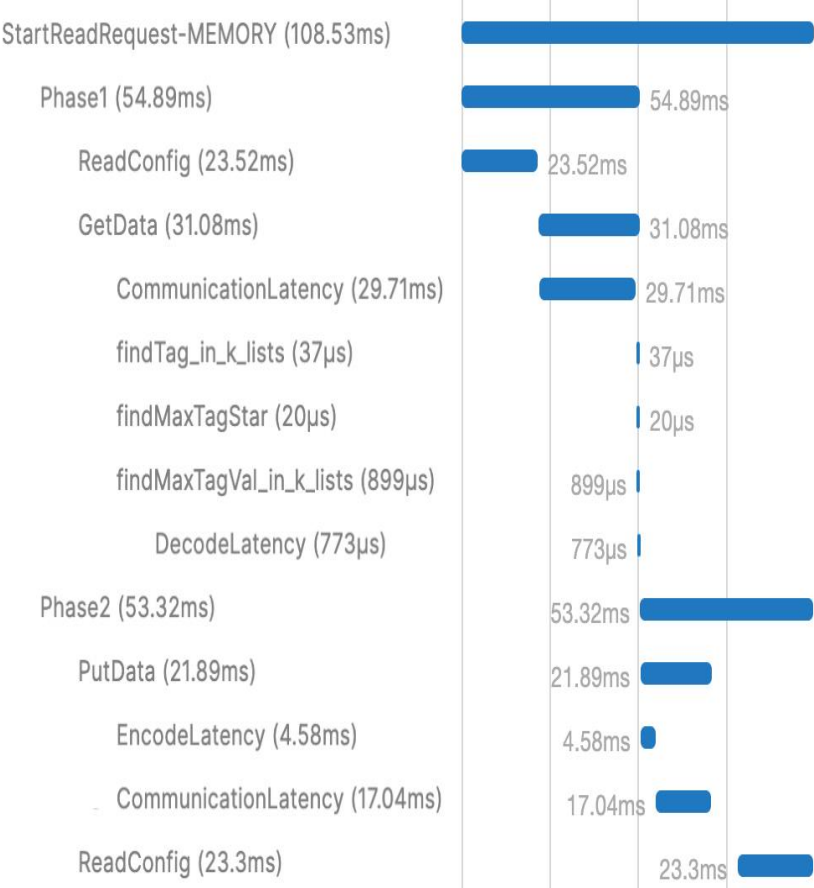


CoAresECF, S:11, W:5, R:5, fsize:512MB, Min/Avg Block Size:64MB, max Block Size:128MB, Debug Level:DSMM

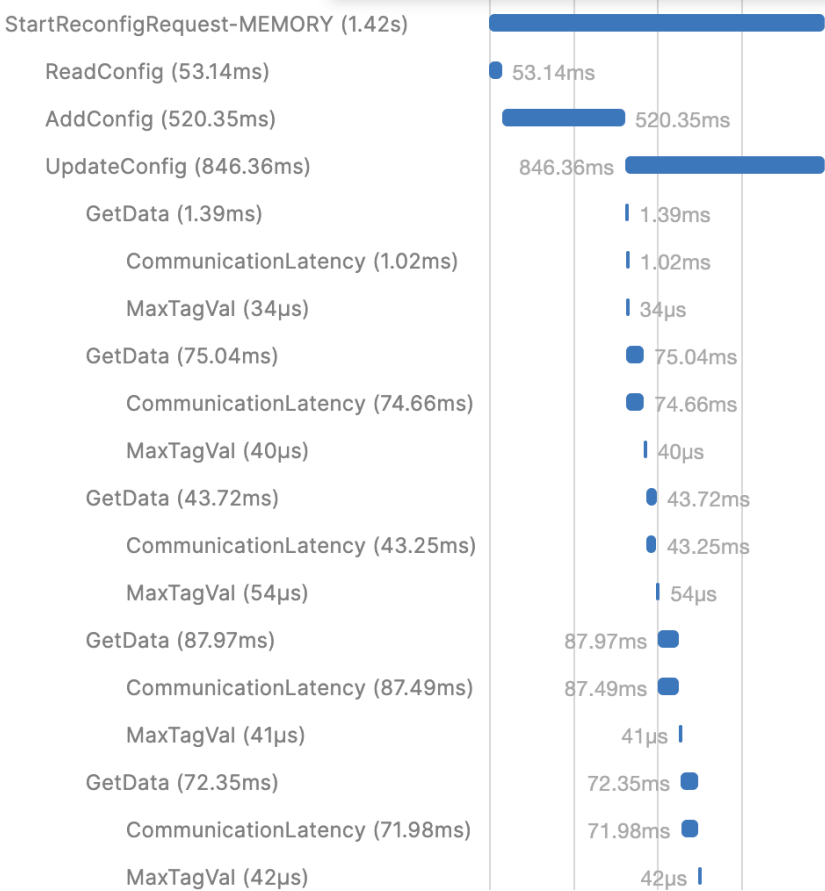
Longevity



CoAresF, S:11, W:5, R:15, G=5, fsize:4MB,
Debug Level:DSMM

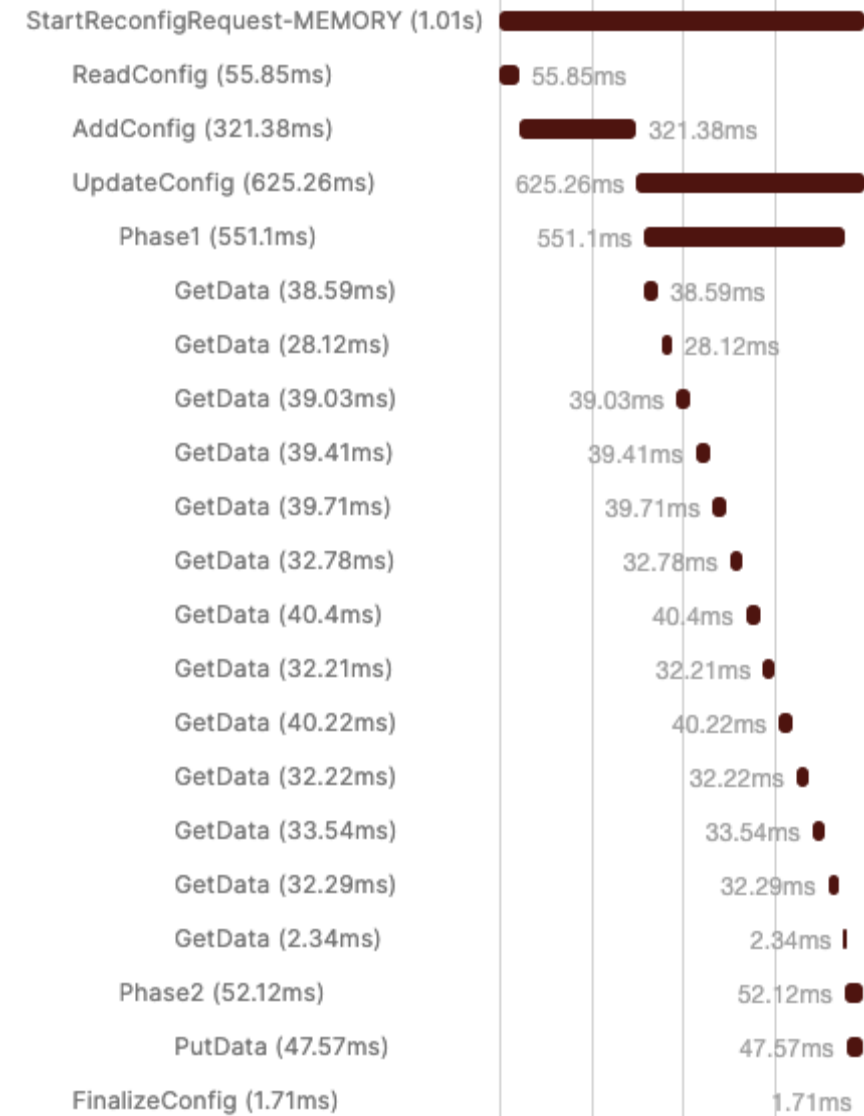
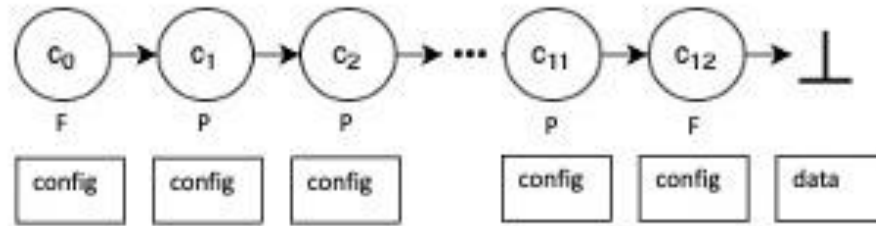


CoAresF, S:11, W:5, R:15, G=1,
fsize:4MB, Debug Level:DSMM



CoAresF, S:11, W:5, R:15, G=5, fsize:4MB,
Debug Level:DSMM

The Latencies of read-config and get-data.



The Latencies of read-config and get-data.

