

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



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

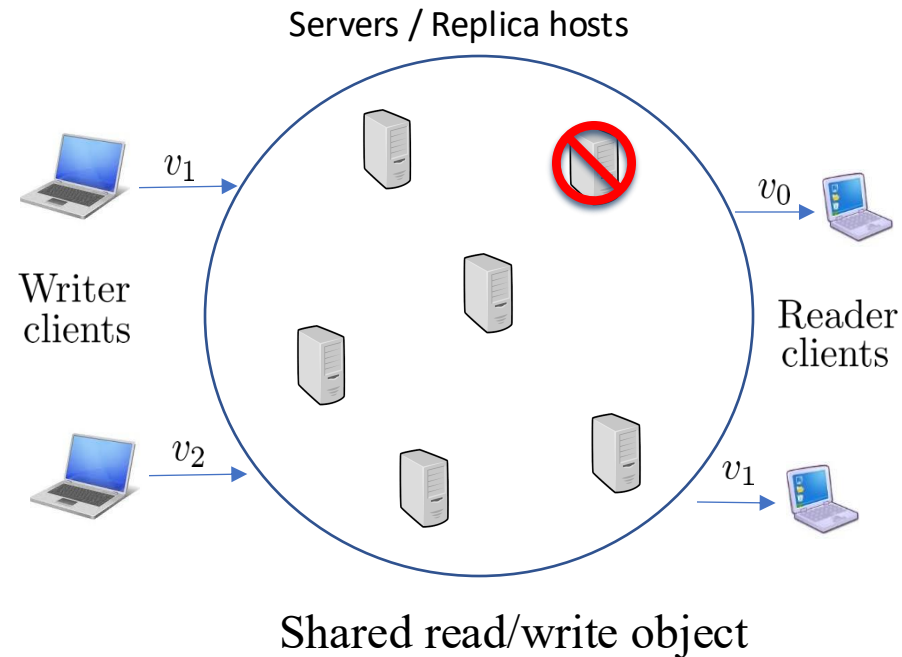
¹University of Cyprus, Nicosia, Cyprus

²Algolysis, Limassol, Cyprus

SRDS 2024, Charlotte, USA

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

Distributed Shared Memory Emulations (DSMs)



- A set of **servers (configuration)** maintain replicas of the same data object.
- Clients (**readers/writers**) access the object by sending messages to these servers.
- Read/Write operations are structured in terms of **phases**.
- Each phase consists of **two** communication exchanges (broadcast & convergecast).
- Fixed Configuration -> **Static** environment, Reconfiguration -> **Dynamic** environment
- Consistency guarantees
 - Safety, Regularity, **Atomicity** (Atomic DSMs) [Lamport 1986]

Seminal Algorithm - ABD

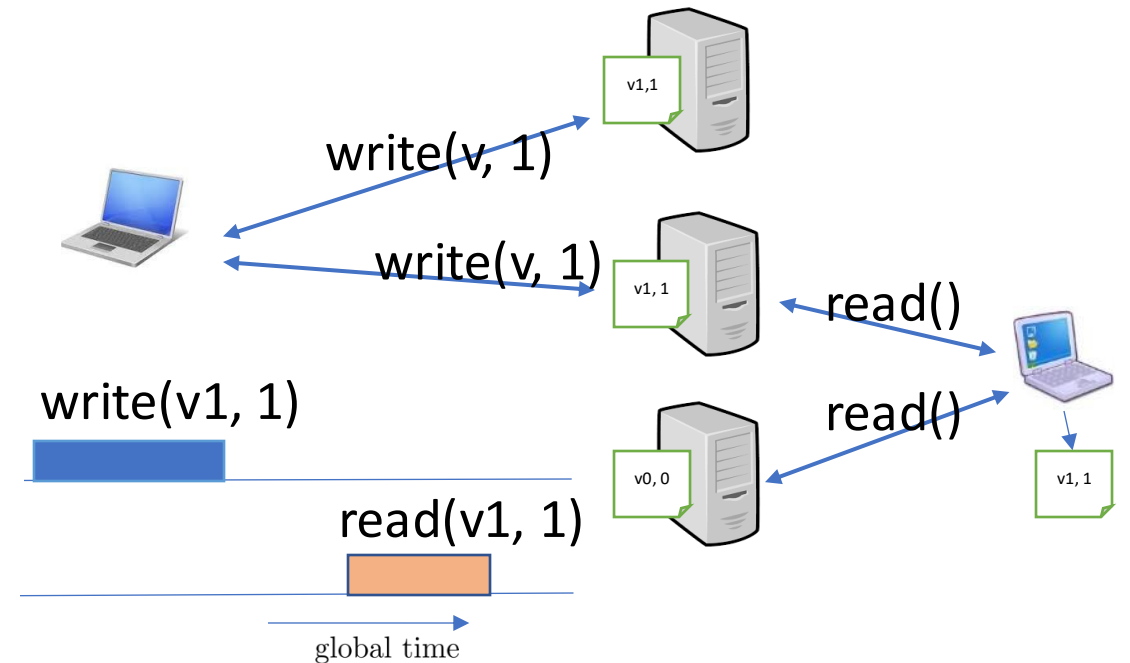
An elegant, intuitive solution that

- 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

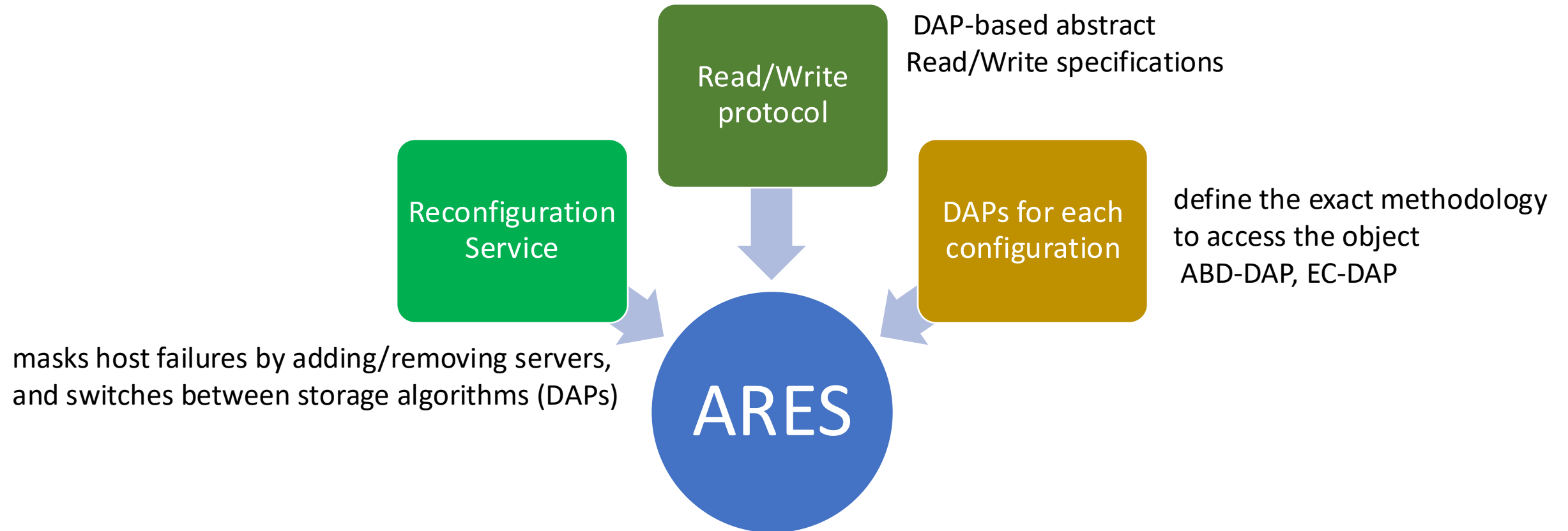
in 1997 for **MWMMR**, assigning tags $\langle ts, wid \rangle$ – **MW-ABD**

- Many more complex ABD-like protocols were developed over the years to address various challenges such as fault-tolerance, efficiency, and scalability.

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.

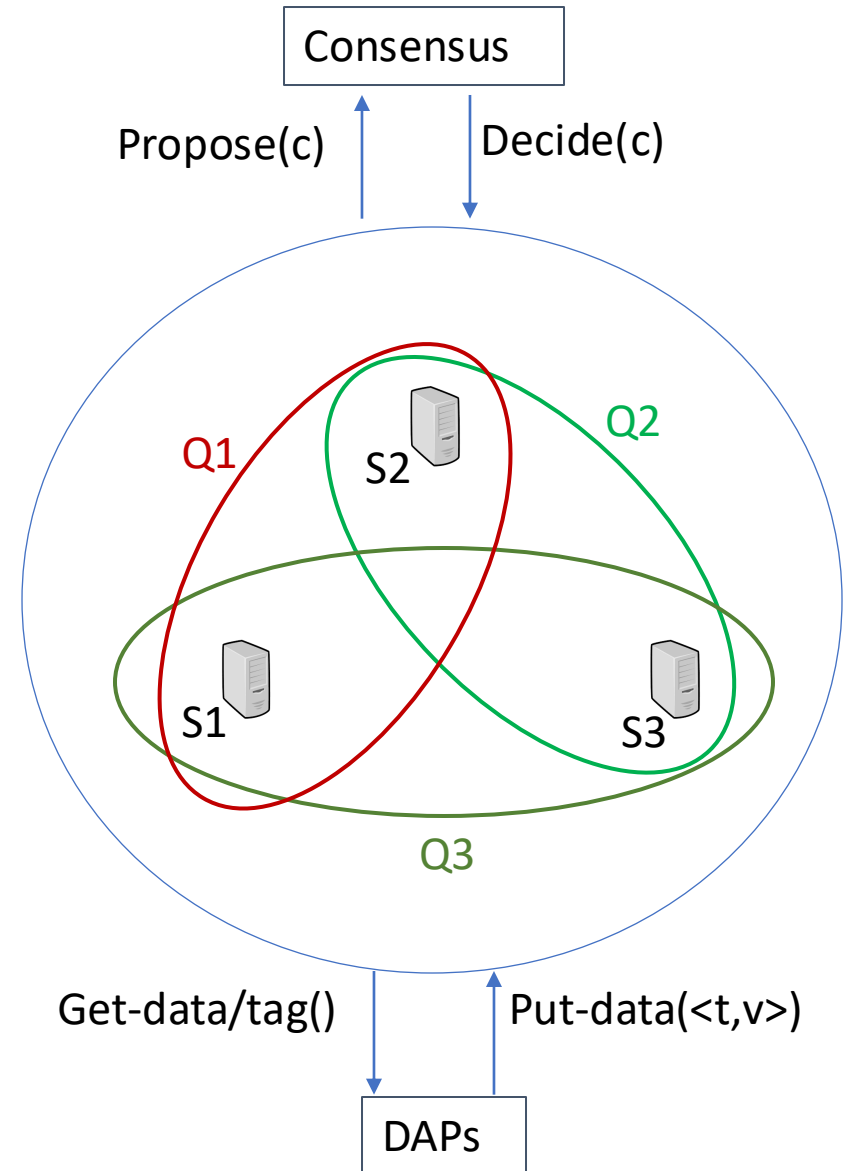
N. Lynch, A. Shvartsman.. "Robust emulation of shared memory using dynamic quorum-acknowledged broadcasts," In Proc. of FTCS pp. 272–281 (1997).

ARES - Adaptive, Reconfigurable, Erasure Code, Atomic Storage

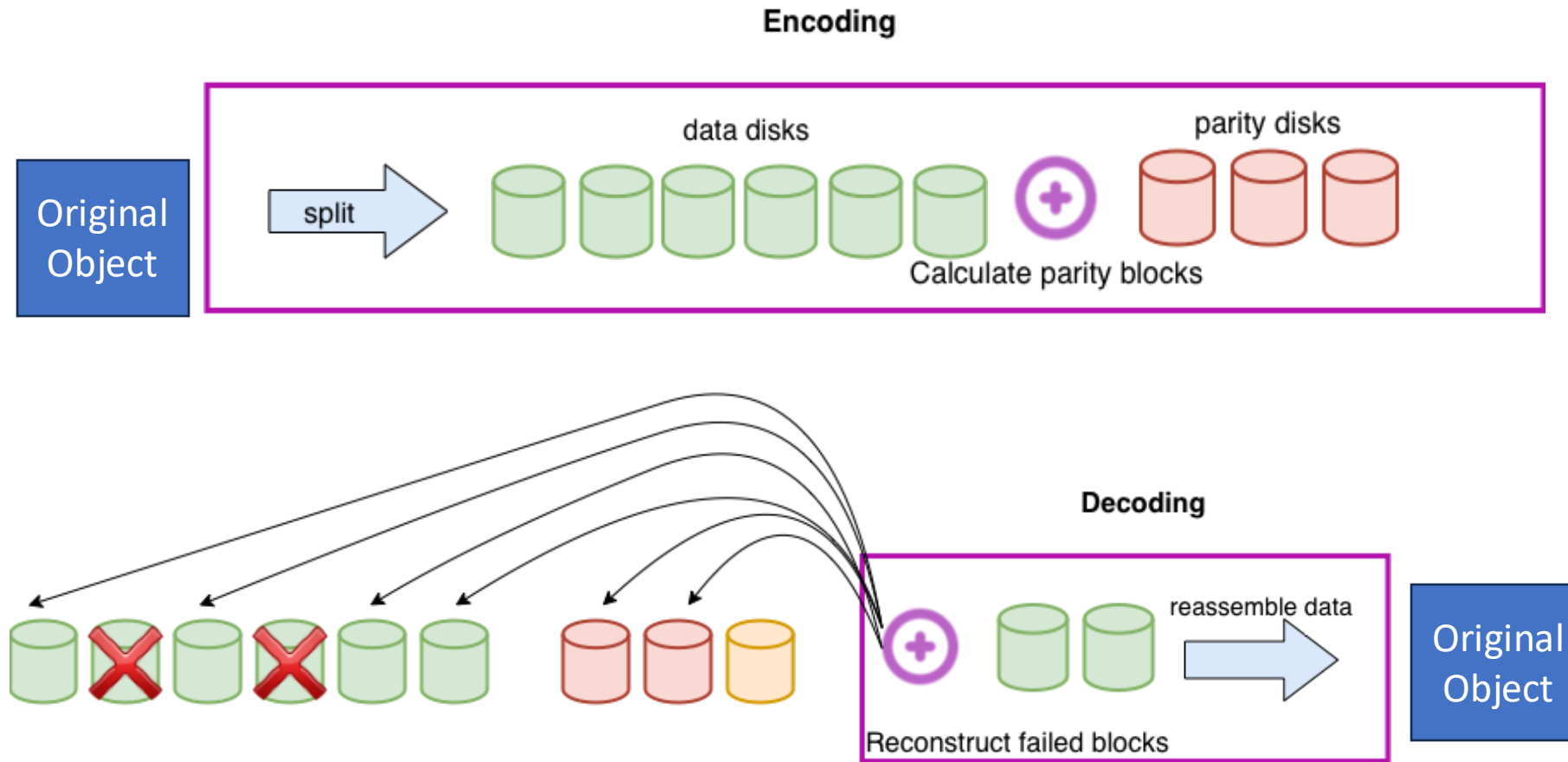


Configurations

- 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



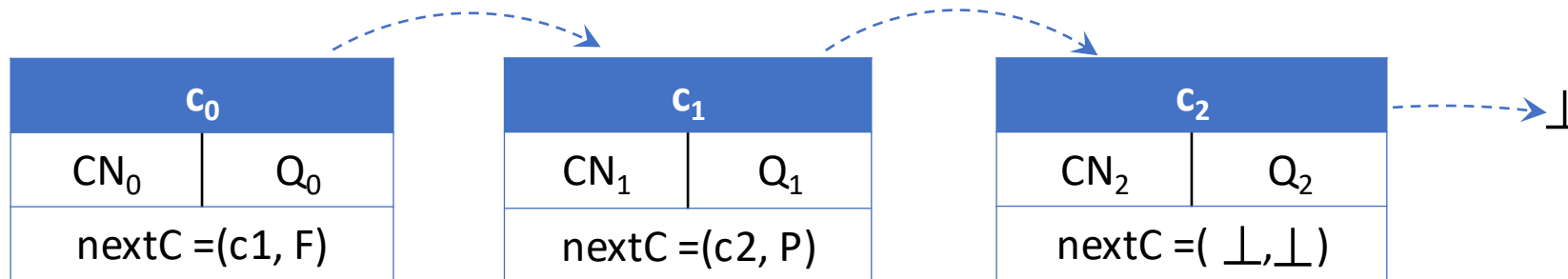
Erasure Code (EC)



(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
    $cseq \leftarrow \text{add-config}(cseq, c)$          introduce the new configuration
10:   $\text{update-config}(cseq)$                        move the latest value to the new config
    $cseq \leftarrow \text{finalize-config}(cseq)$       let servers know it is good to be finalized
12: end operation
```

(1) }
(2) }

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) = (\langle tmax+1, wi \rangle, val)`
- Do
 - `cseq[v].put-data(t,v)`
 - Traverse Sequence `cseq`
- `while(|cseq| > v)`

Main Objective

The primary goal is to **identify flaws** in **DSMs** and **guide their optimization**.
We demonstrate this through the ***ARES*** DSM.

Performance Analysis Challenges in DSMs

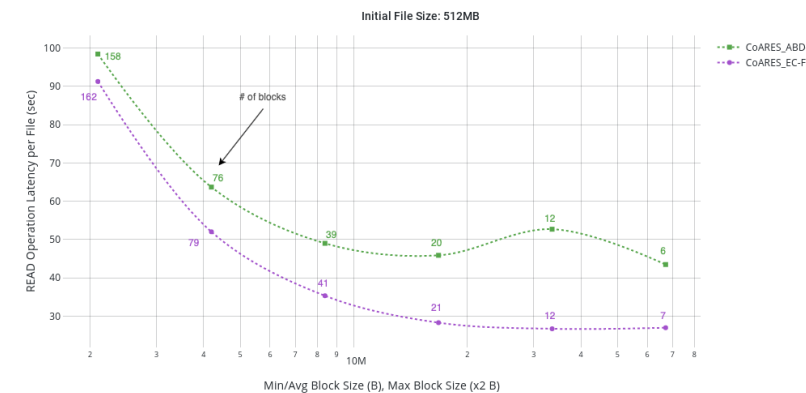
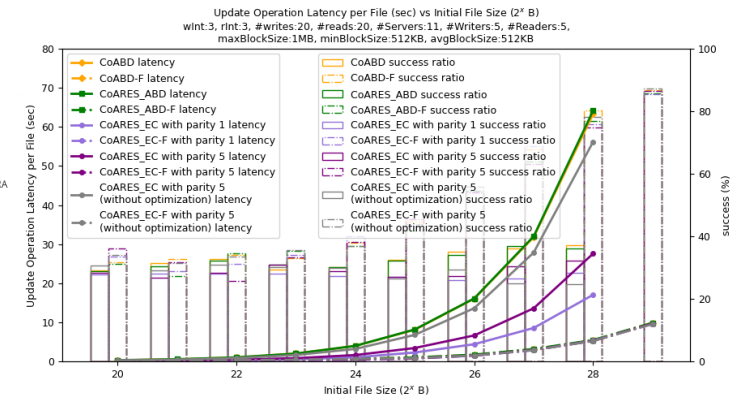
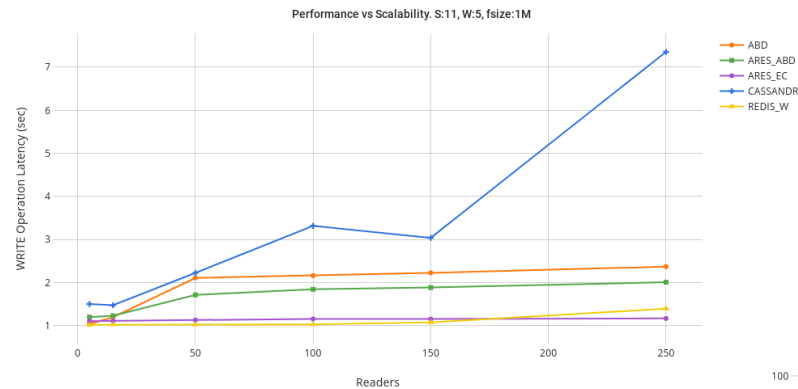
- Identifying performance bottlenecks in complex DSMs can be challenging
- Traditional logging techniques may not provide sufficient insight

```
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"):

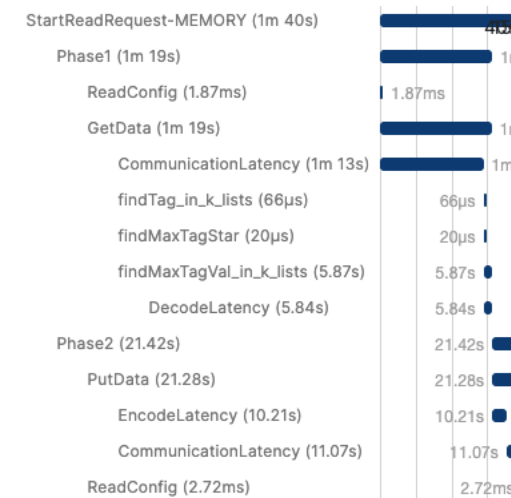
            # Increase msg counter for next op
            self.msg_id += 1

            # 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")
```

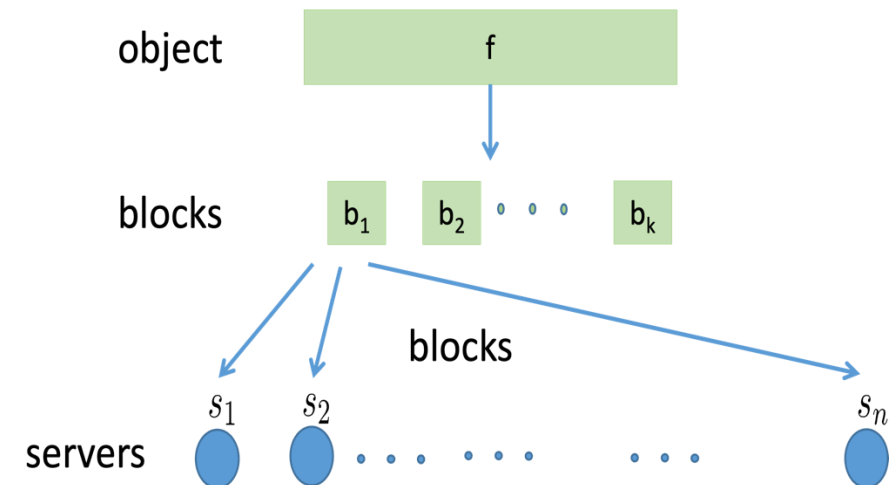
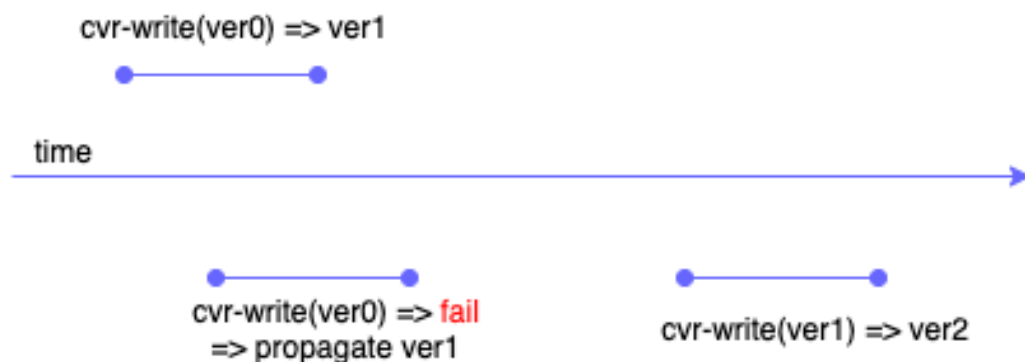
Trace



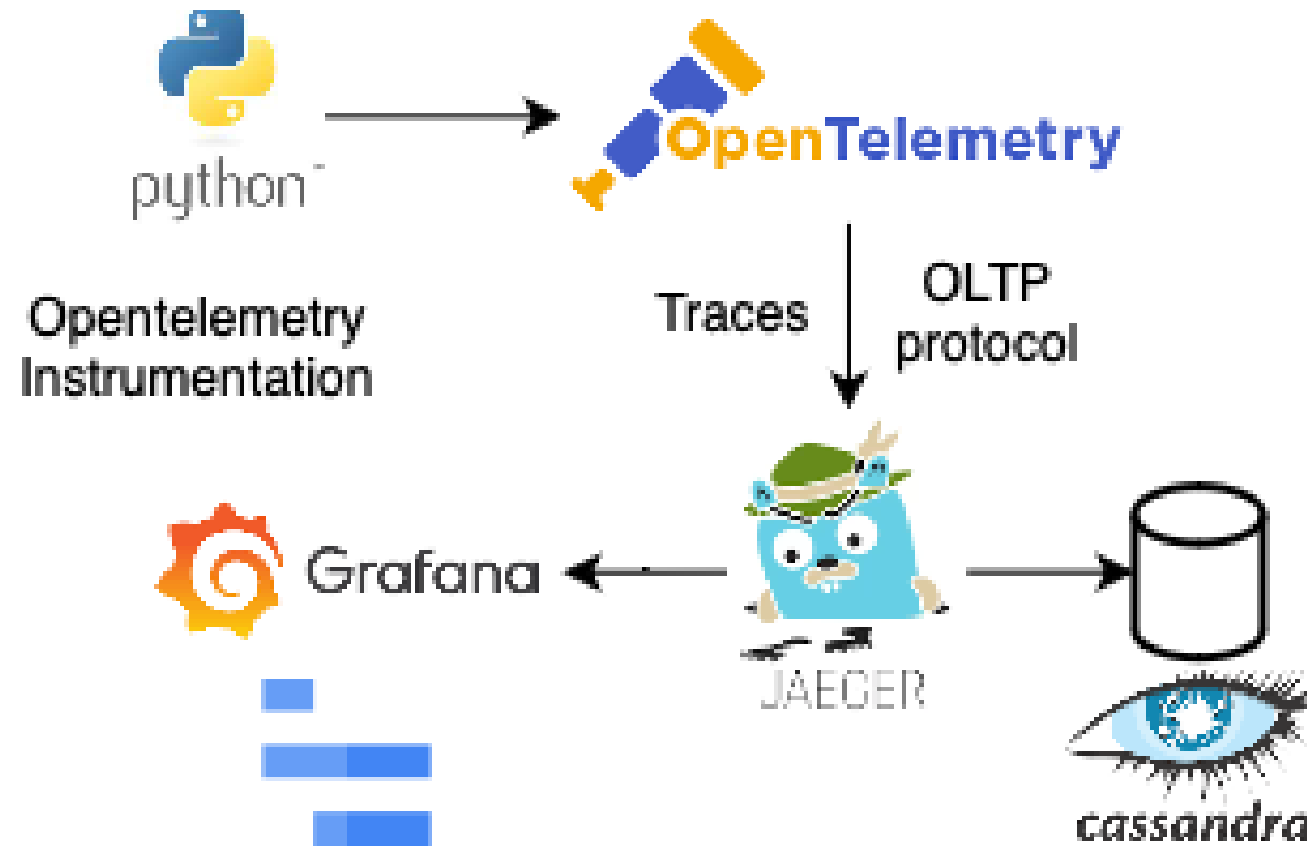
Spans

Evaluated Algorithms

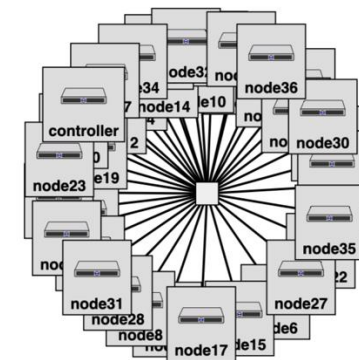
<i>ARESABD</i>	This is Ares that uses the ABD-DAP implementation.
<i>CoARESABD</i>	The coverable version of <i>ARESABD</i> .
<i>CoARESABDF</i>	The fragmented version of <i>CoARESABD</i> .
<i>ARESEC</i>	This is <i>ARES</i> that uses the EC-DAP implementation.
<i>CoARESEC</i>	The coverable version of <i>ARESEC</i> .
<i>CoARESECF</i>	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> .



Methodology: ARES Distributed Tracing



Experimental Setup

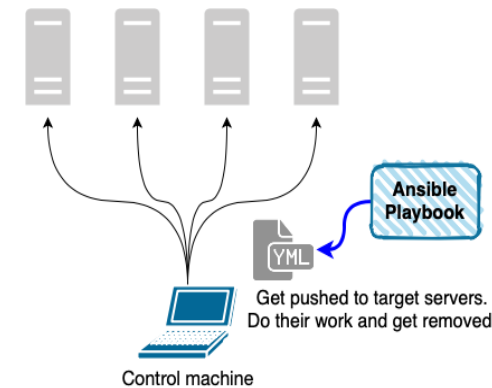


We used two main tools to run the experiments:

- **Emulab:** an emulated WAN environment testbed.

- 39 machines with 100 Mb/s bandwidth
- Each server is deployed on a different machine.
- Clients are all deployed in the remaining machines in a round robin fashion.

- **Ansible:** a tool to automate different IT tasks.

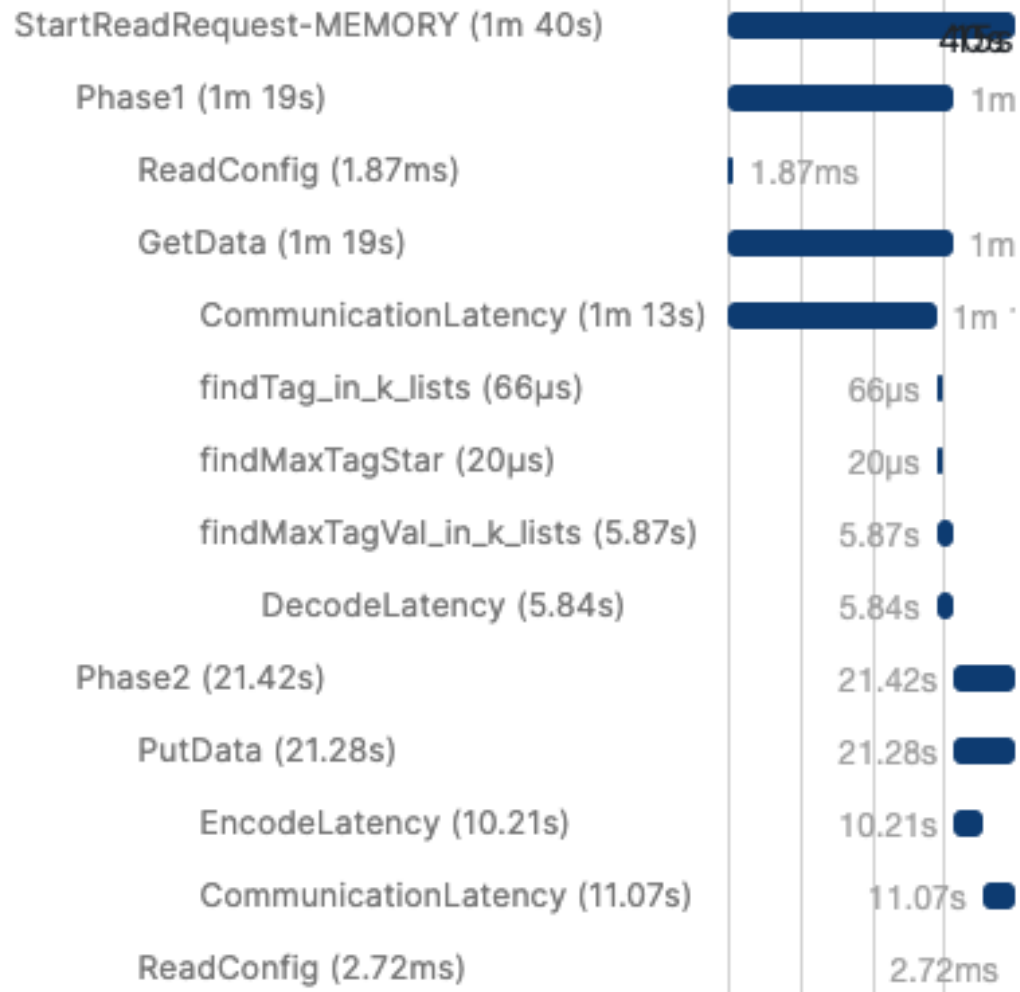


- **Performance Metric**

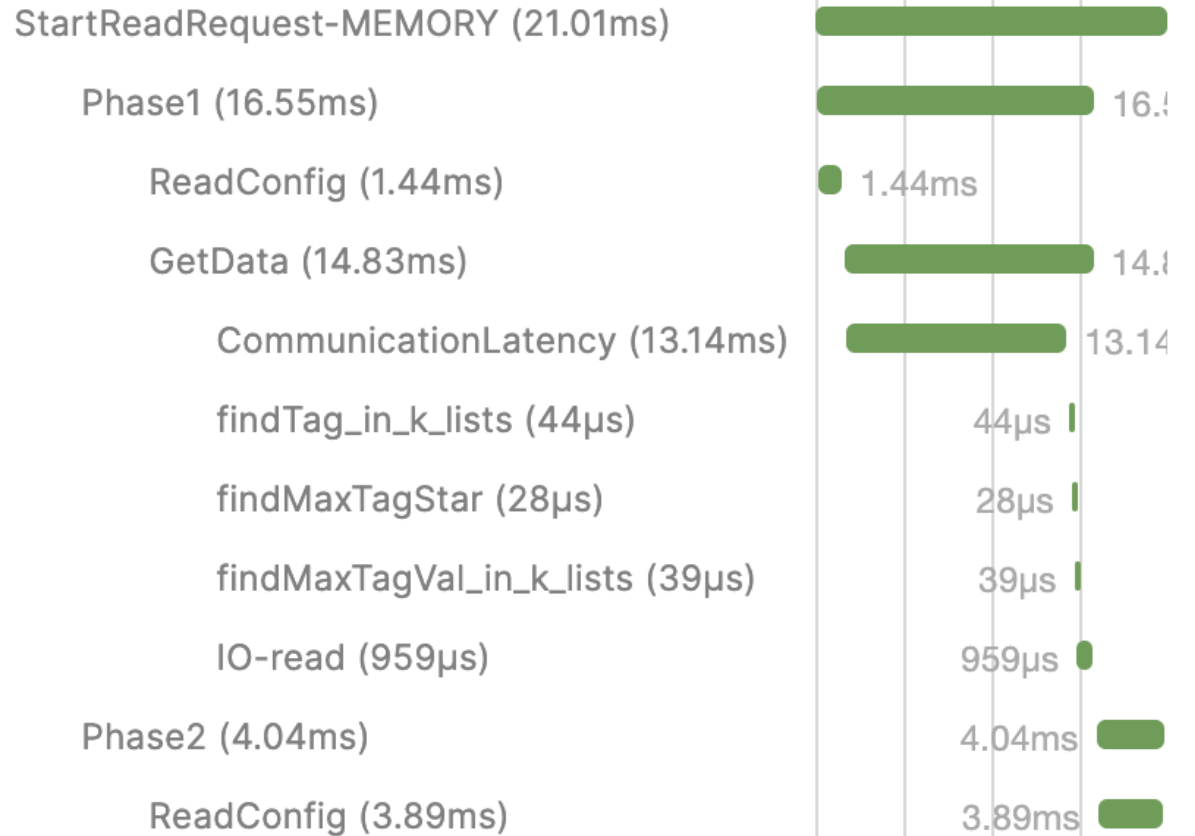
- Operation latency of clients (Communication + Computation Overhead).
- Sample traces near the average duration for each scenario.
- Three executions.



Object Size

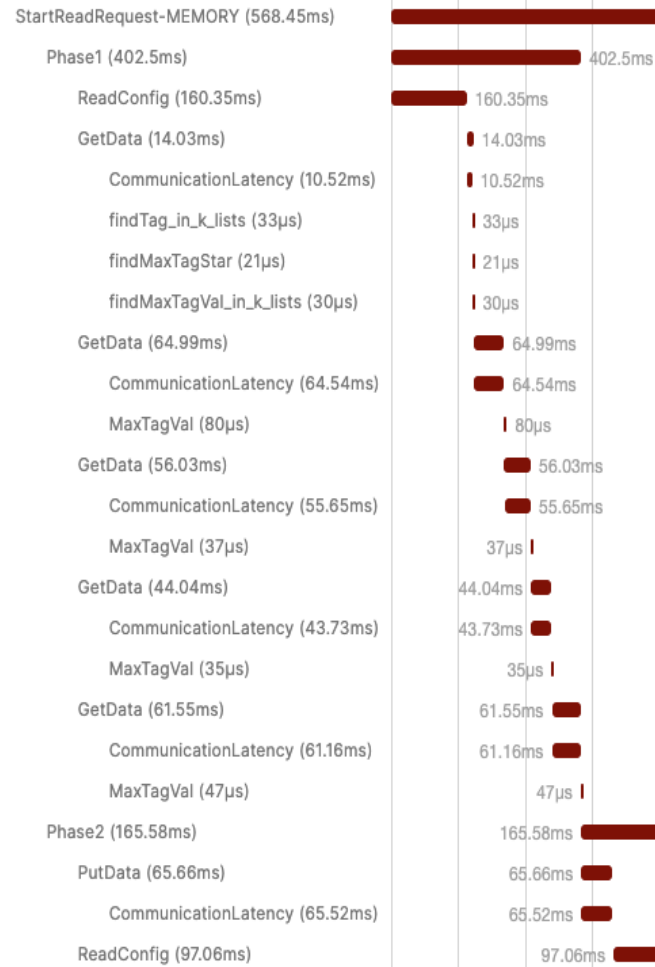


ARESEC, S:11, W:5, R:5, fsize:512MB, Debug Level:DSMM

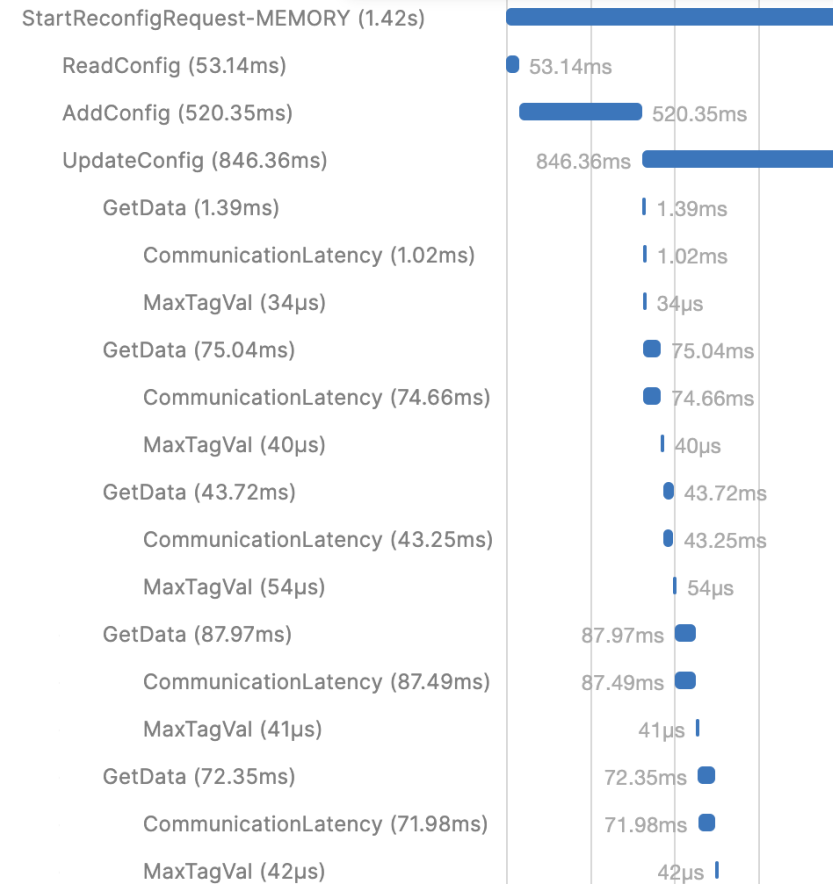


CoARESECF, S:11, W:5, R:5, init fsize:512MB, Debug Level:DSMM

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

From ARES to ARES II

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⟩)
    cseq ← read-config(cseq)
22:    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:  $C_s \leftarrow \{nextC_s : \text{received } nextC_s \text{ from } s\}$ 
    $Tags_{dec}^{\geq k}$  = set of tags that appears in  $k$  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:  $code\text{-}elems = [(\tau, e_1), \dots, (\tau, e_n)], e_i = \Phi_i(v)$ 
   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  $s \in \mathcal{S}_g$  s.t.  $|\mathcal{S}_g| = \lceil \frac{n+k}{2} \rceil$ 
   and  $\mathcal{S}_g \subset c.Servers$ 
28:  $C_s \leftarrow \{nextC_s : \text{received } nextC_s \text{ from each } s \in \mathcal{S}_g\}$ 
   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
  
```

Last finalized config

older configs from last finalized config

EC-DAP

send GC request to servers

Remove the GC configs

server updates nextC to point to the last finalized config

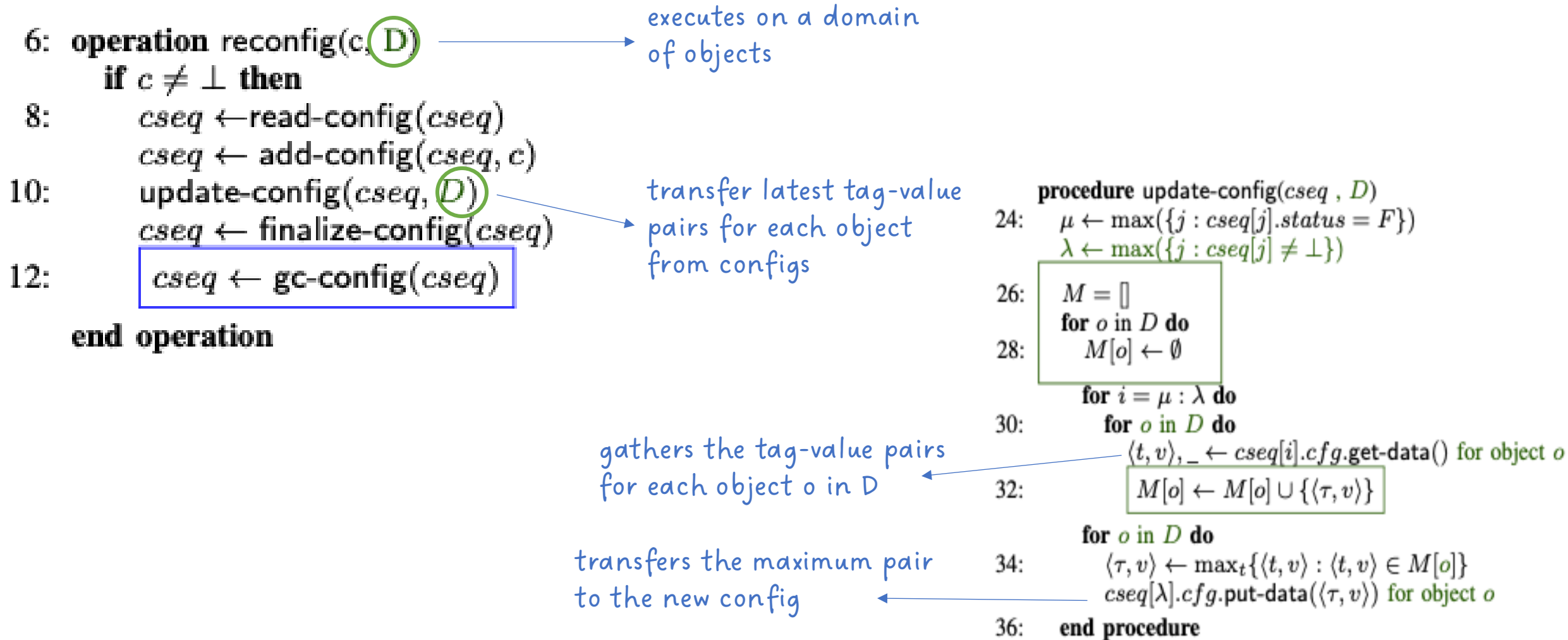
the server sets the data of config to \perp

```

end receive
Upon receive (GC-CONFIG,  $cfgT_{in}$ )si, ck from q
if  $cfgT_{in}.cfg.ID > nextC.cfg.ID$  then
   $nextC \leftarrow cfgT_{in}$ 
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



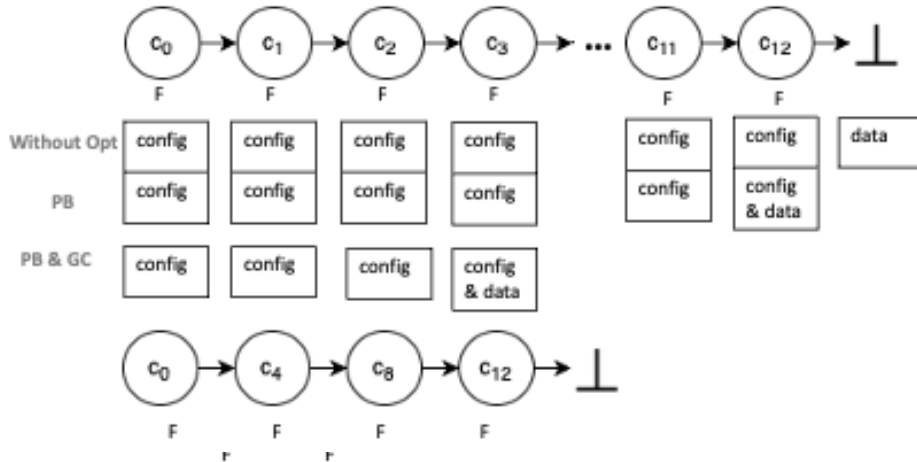
Optimization Results – Piggyback

alg./ f_{size}	CoARESABDF	CoARESABDF <i>PB</i>	CoARESECF	CoARESECF <i>PB</i>
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

- **Non-Fragmented Algorithms:** No notable improvements for medium-sized objects.
 - Removal of **read-config** occurs only twice, so impact on latency is minimal.
- **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	COARES _F	COARES _F PB	COARES _F PB&GC
11 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
12 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.
- CoARES_F & 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.

Conclusions and Future Work

- 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.

Future Work - Devise strategies on **when** and **how** to introduce new configurations.

- Ensure that the system remains operational despite server failures.
- Improve performance by replacing older servers with more powerful ones.
- * Monitoring tools to collect health metrics, threshold-based approaches for determine when to reconfigure, machine learning algorithms for anomaly detection, server rebalancing policy

Thank you!

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

