



COLLABORATE:

Building a Robust, Consistent, and Efficient Distributed Shared Storage System for Large Data Objects that Promotes Collaboration

Presented by: **Nicolas Nicolaou, Project Coordinator**

This work is supported by the Cyprus Research and Innovation Foundation under the grant agreement POSTDOC/0916/0090.



RESEARCH
& INNOVATION
FOUNDATION



European Union
European Regional
Development Fund



algolysis
algorithmic solutions

 www.algolysis.com

 research@algolysis.com

9/6/2021

In a nutshell!

- Distributed Storage – Problem Statement
- Fragmentation: How to handle Large Objects
 - Blocks
 - Erasure Coding
 - Hybrid Solutions
- Reconfiguration: dealing with failures
- DriveNest: a service to predict imminent node failures



What is a Distributed Storage System?

- How to share data **robustly** in a **message-passing** system?



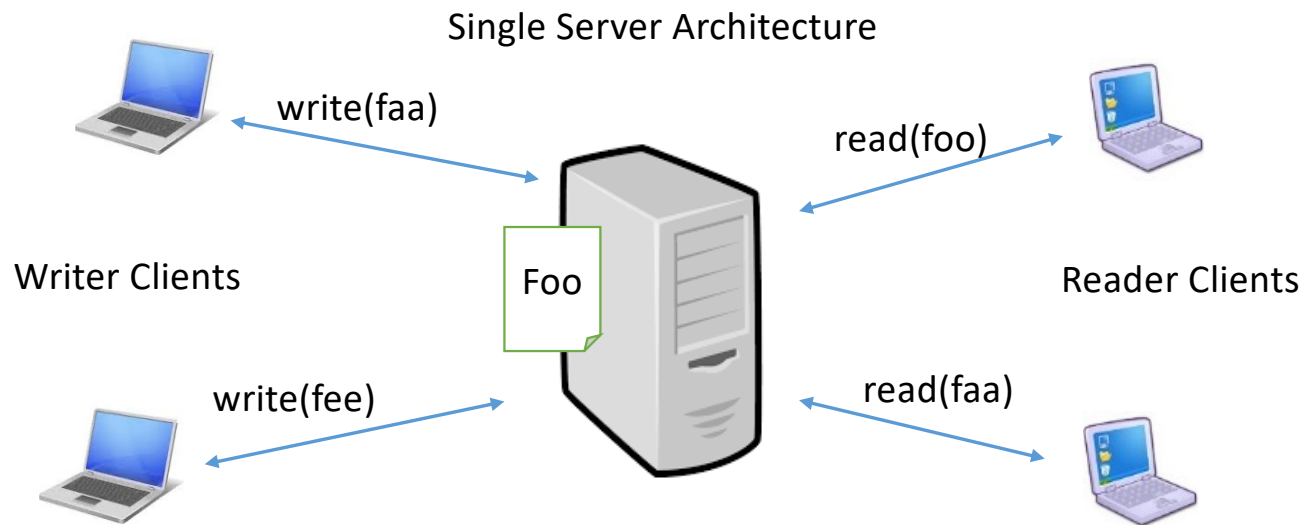
Commercial Solutions:

- Cannot provide “sufficient guarantees” when shared objects are **accesses concurrently**
- Often rely on **centralized solutions** to enable collaboration
- Not offering “suitable guarantees” for application design



What is a Distributed Storage System

- How to share data **robustly** in a **message-passing** system?



Pros:

- Easy to implement
- Strong consistency as the single server imposes the order of operations

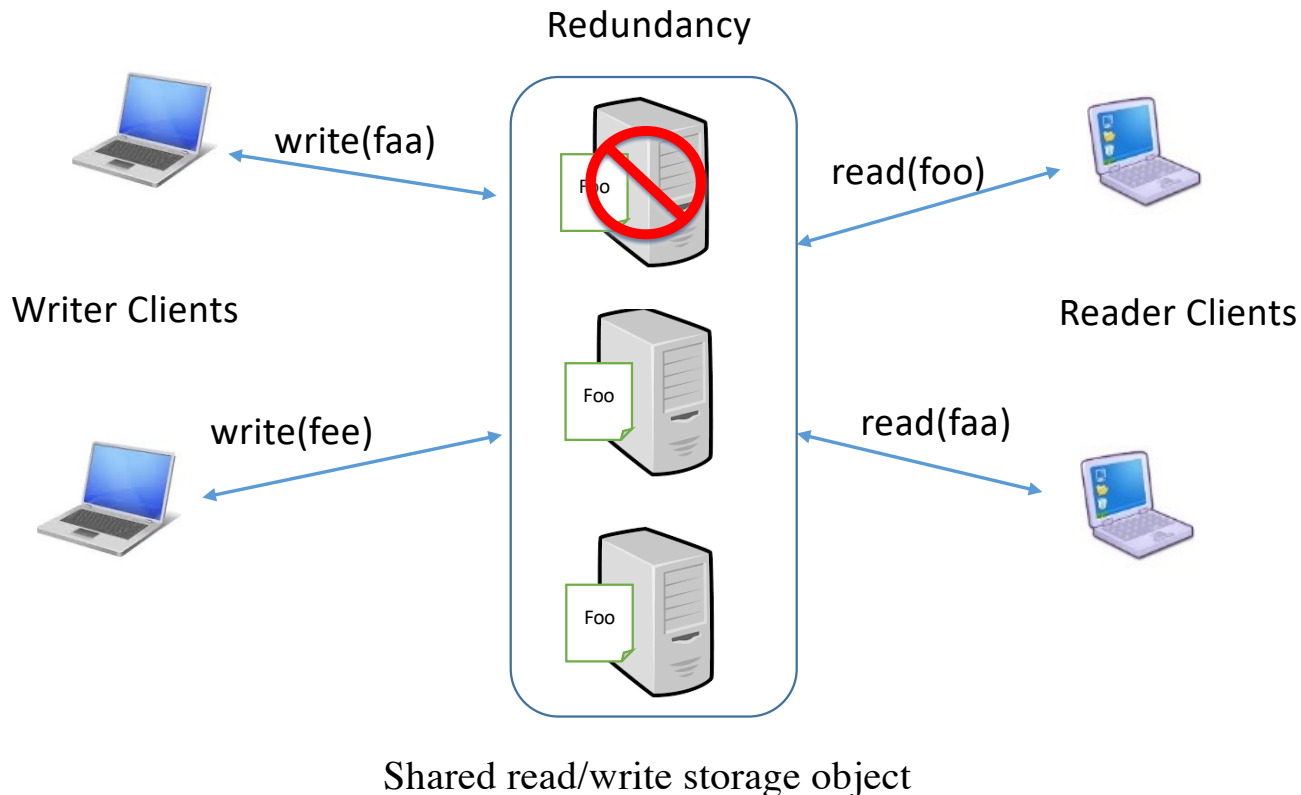
Cons:

- Performance Bottleneck
- Single Point of Failure



What is a Distributed Storage System

- How to share data **robustly** in a **message-passing** system?

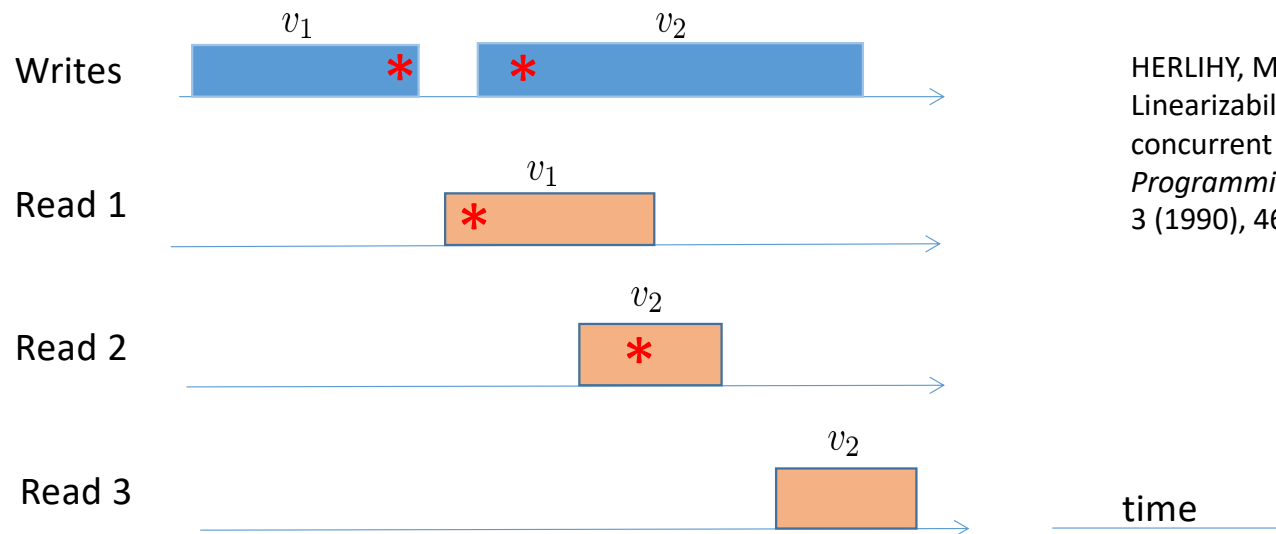


Implementing a **fault-tolerant** shared storage object in an **asynchronous, message-passing** environment:

- Availability + Survivability
=> **use redundancy**
- Asynchrony + Redundancy
=> **concurrent operations**
- Behavior of concurrent operations
=> **consistency semantics**
 - Safety, Regularity, Atomicity [Lamport86]

Atomicity/Linearizability

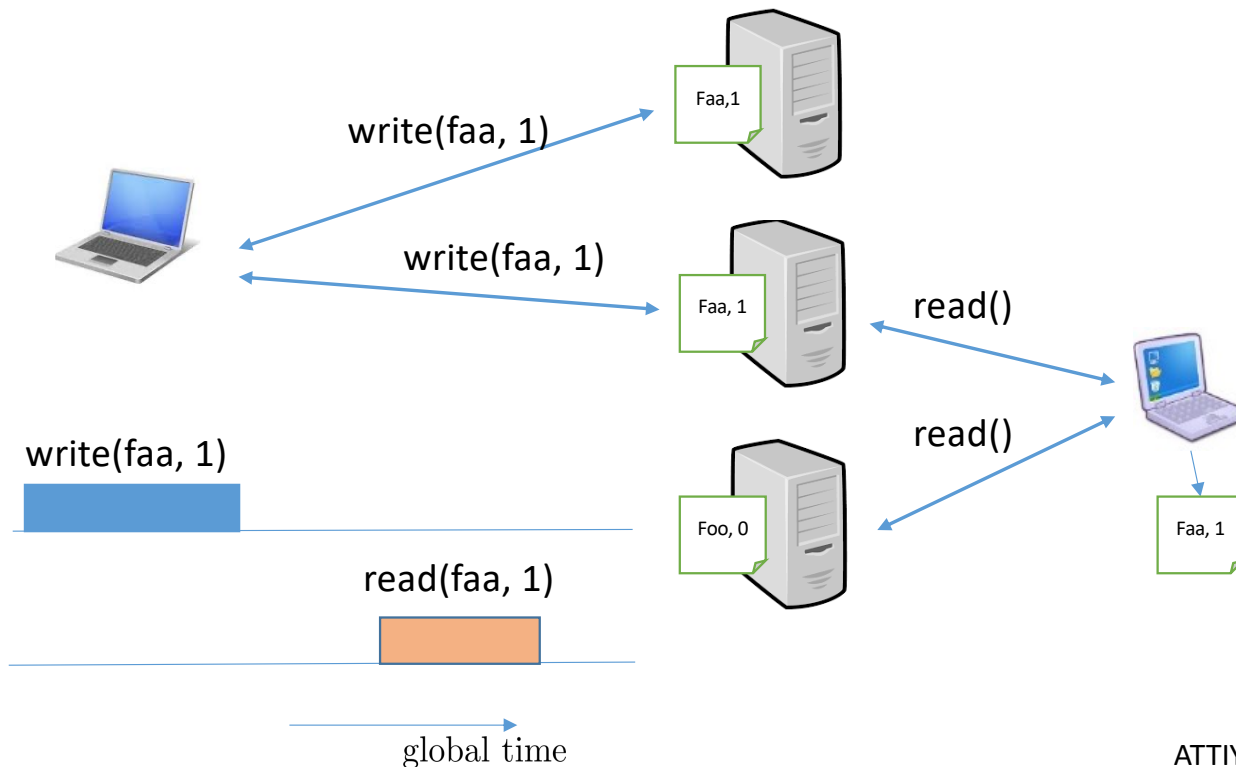
- Provides the **illusion** that operations happen in a **sequential order**
 - a read returns the **value of the preceding write**
 - a read returns a **value at least as recent** as that returned by **any preceding read**



HERLIHY, M. P., AND WING, J. M.
Linearizability: a correctness condition for
concurrent objects. *ACM Transactions on
Programming Languages and Systems* 12,
3 (1990), 463–492.



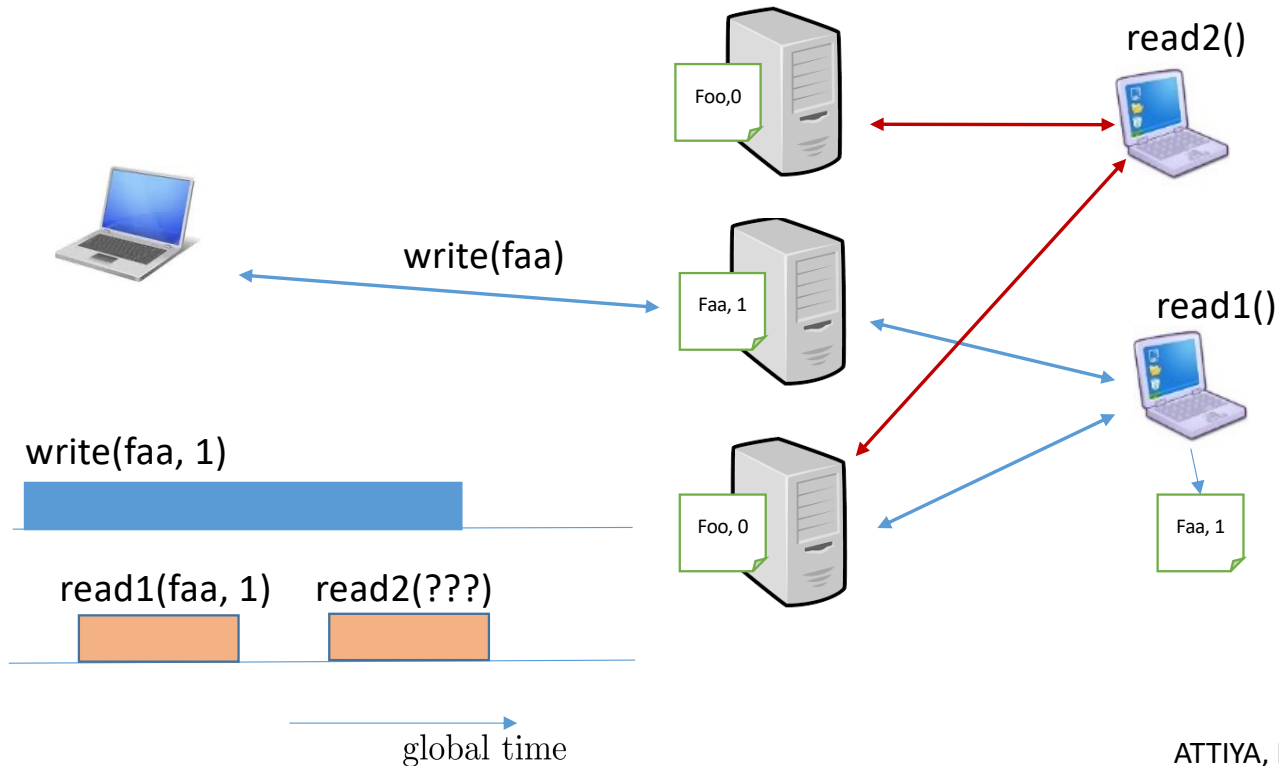
A Simple Solution - ABD



- Attiya, Bar-Noy, Dolev: an elegant, intuitive solution
 - Use the power of the majority
 - Assign logical timestamps to written values
 - **Wait-free** solution

ATTIYA, H., BAR-NOY, A., AND DOLEV, D. Sharing memory robustly in message passing systems. *Journal of the ACM* 42(1) (1996), 124–142.

A Simple Solution - ABD

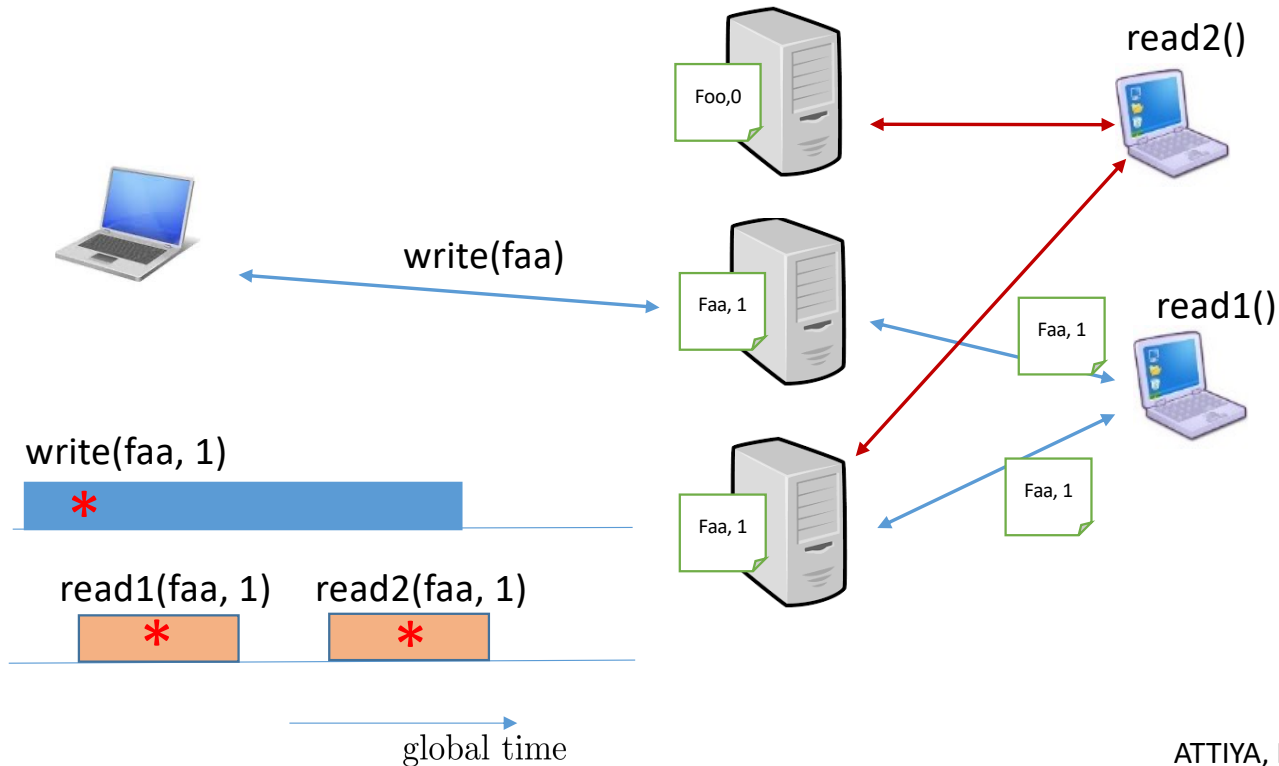


- A “read” needs to “write”
 - Phase 1: query
 - Phase 2: propagate

ATTIYA, H., BAR-NOY, A., AND DOLEV, D. Sharing memory robustly in message passing systems. *Journal of the ACM* 42(1) (1996), 124–142.



A Simple Solution - ABD

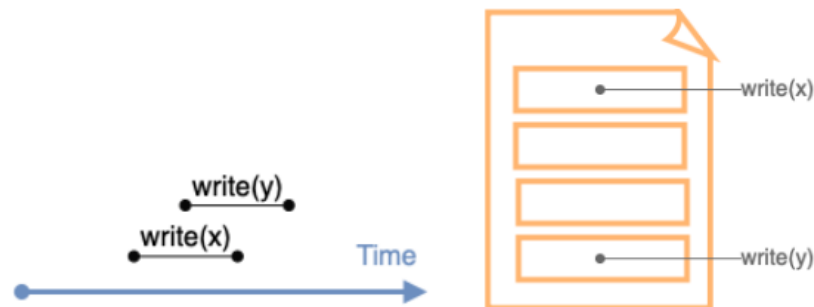


- A “read” needs to “write”
 - Phase 1: query
 - Phase 2: propagate

ATTIYA, H., BAR-NOY, A., AND DOLEV, D. Sharing memory robustly in message passing systems. *Journal of the ACM* 42(1) (1996), 124–142.

Solutions for Large Objects

- ABD efficient for small objects
 - Each write operation sends $(S/2)+1$ copies of the object
 - Each read operation sends S copies of the object
- Moreover **concurrent write operations** may **overwrite** one another
 - Unable to handle write operations working on different parts of a large object



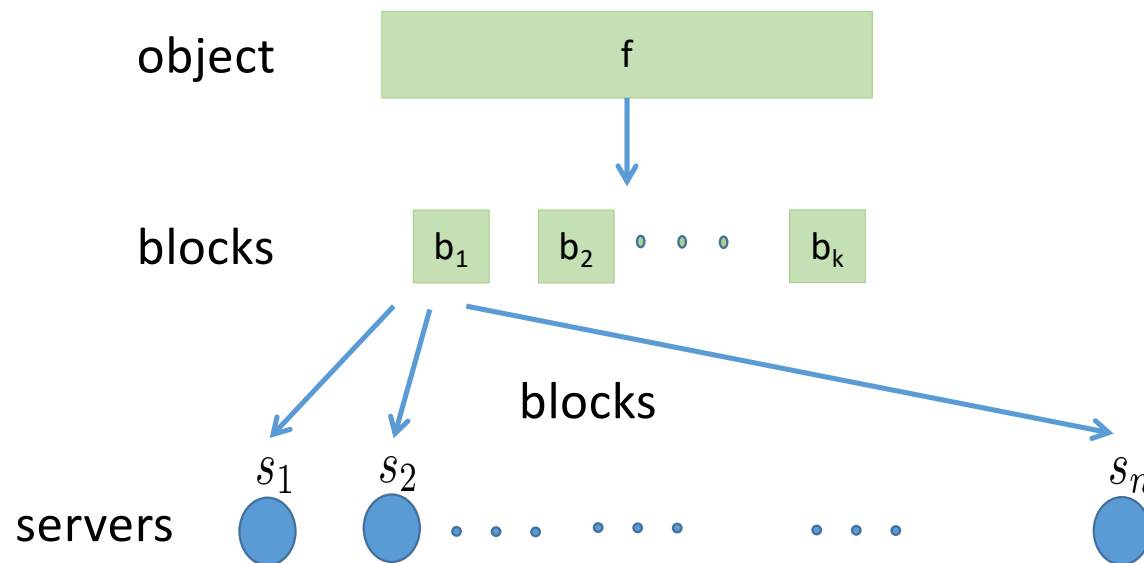
Main Project Goal

Develop **practical** and **robust** DSS in the **message-passing, asynchronous**, environment while allowing **high concurrency** and preserving **strong consistency**.



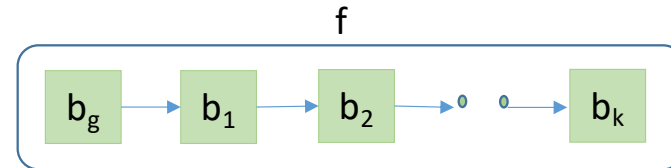
Solution 1: Fragmentation

- Most intuitive solution
 - Split large objects into smaller fragments
 - Treat each individual block as an atomic object



Solution 1: Algorithm CoBFS

- Fragmented Objects:
 - Connected list of blocks
 - Each block points to the next block
- Write Operation $\text{write}(f)$
 - Propagate only modified and new blocks
- Read Operation $\text{read}(f)$
 - Start from genesis block and read all the blocks
 - Optimization: Only blocks that have changed are send to the read

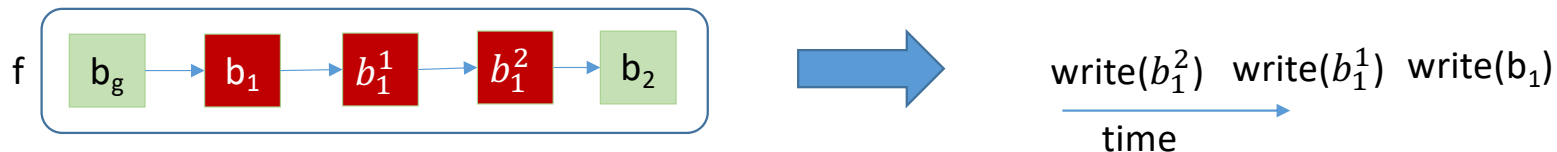


Solution 1: Algorithm CoBFS

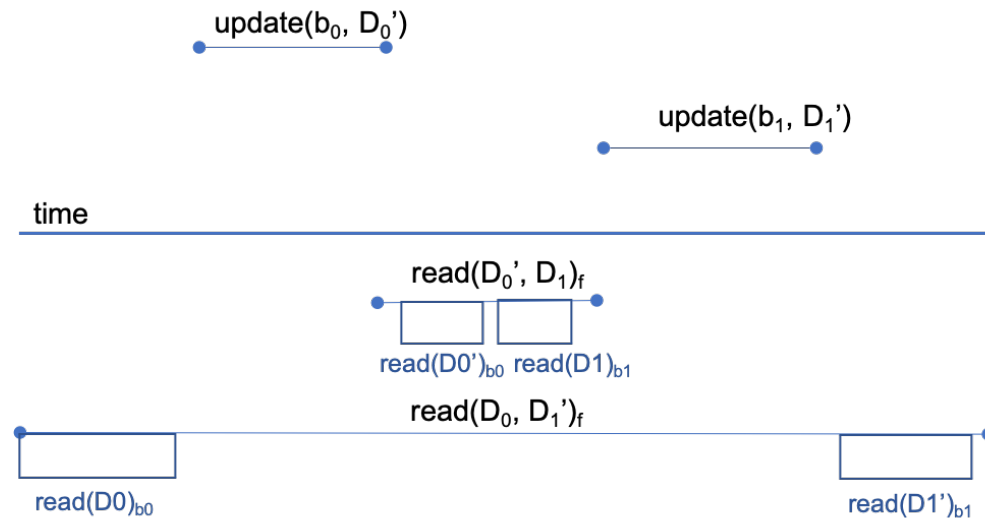
- Write/Update Operation:
 - Run fragmentation and block matching algorithms to determine
 - Modified blocks
 - New blocks
 - Case 1: Only a single block has changed



- Case 2: Changed block overflowed and new blocks introduced

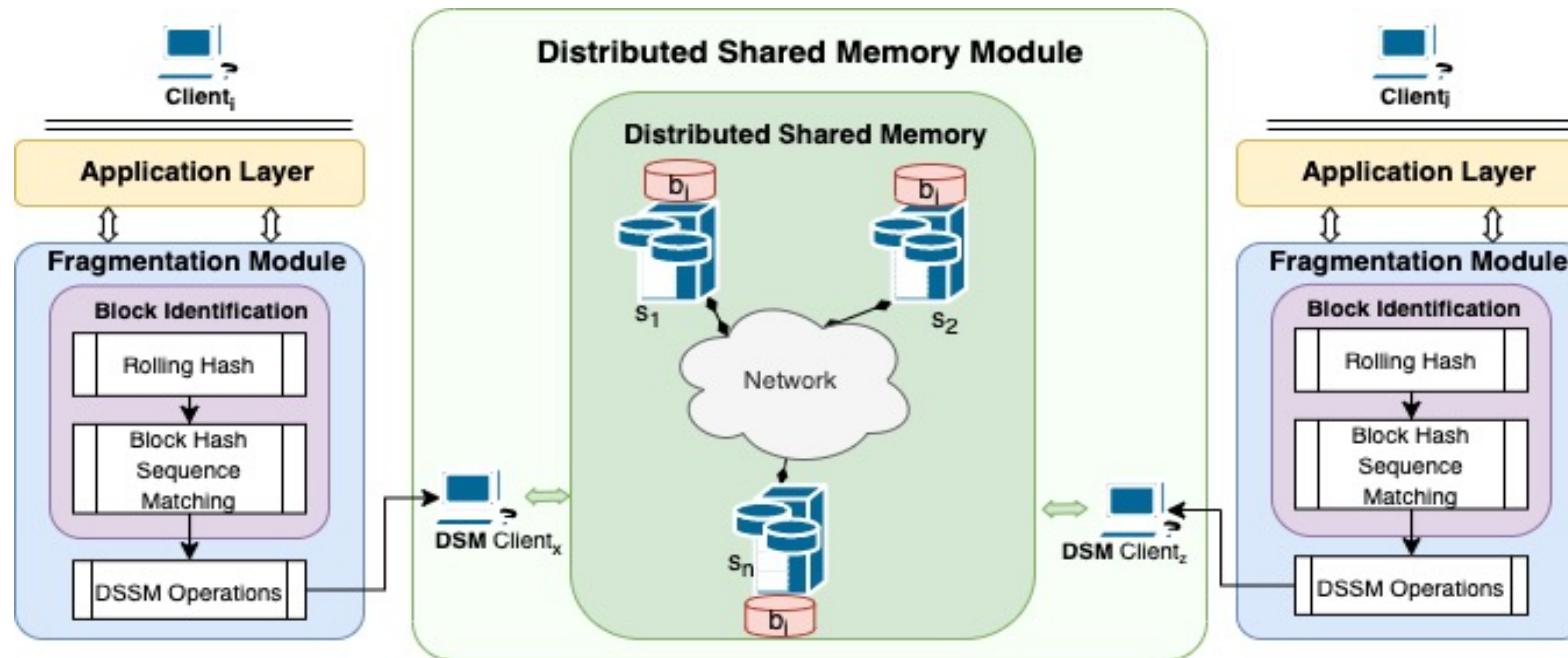


Solution 1: Fragmented Linearizability



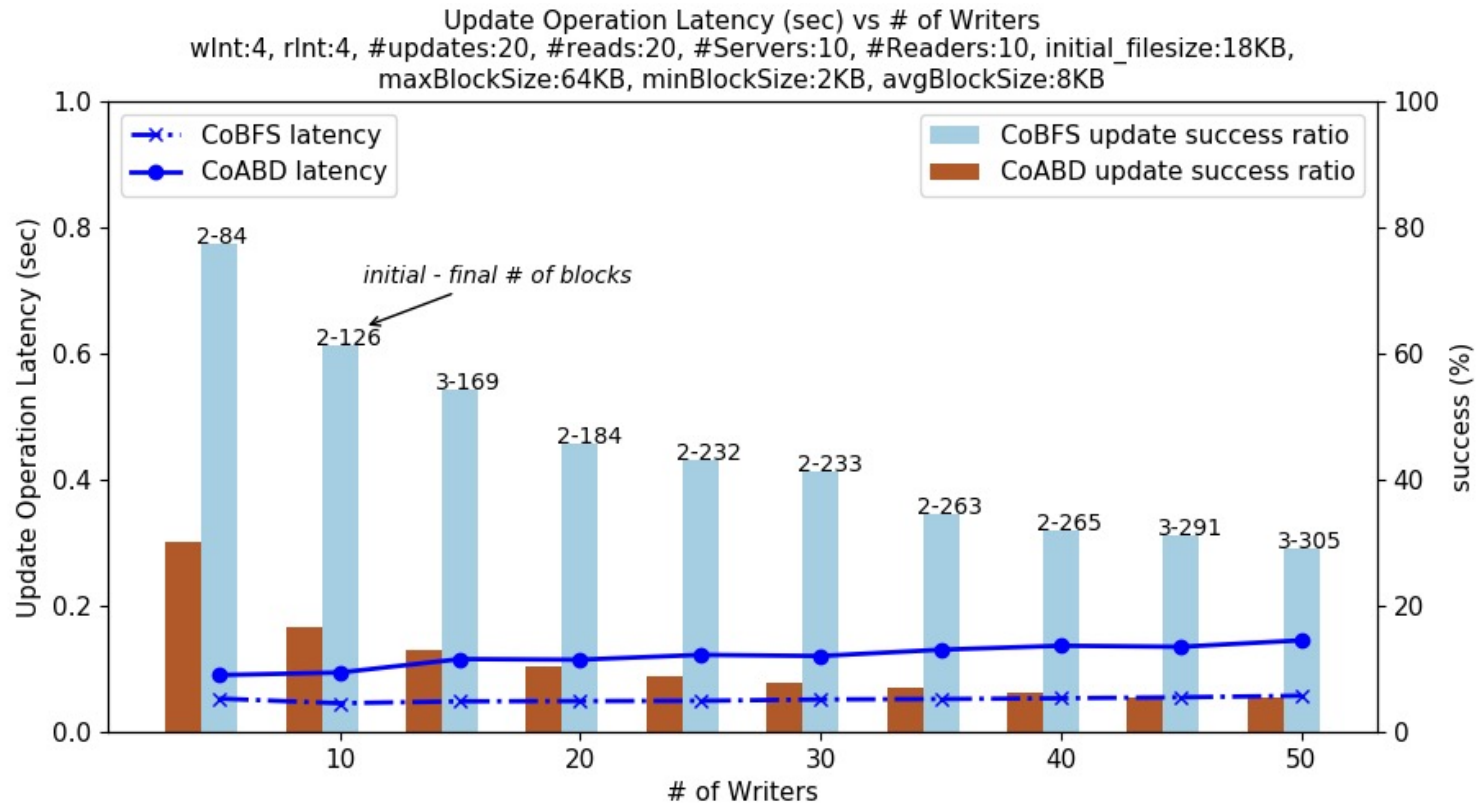
Fragmented Linearizability: all concurrent operations on different blocks prevail, and only concurrent operations on the same blocks are conflicting.

Solution 1: Basic Architecture

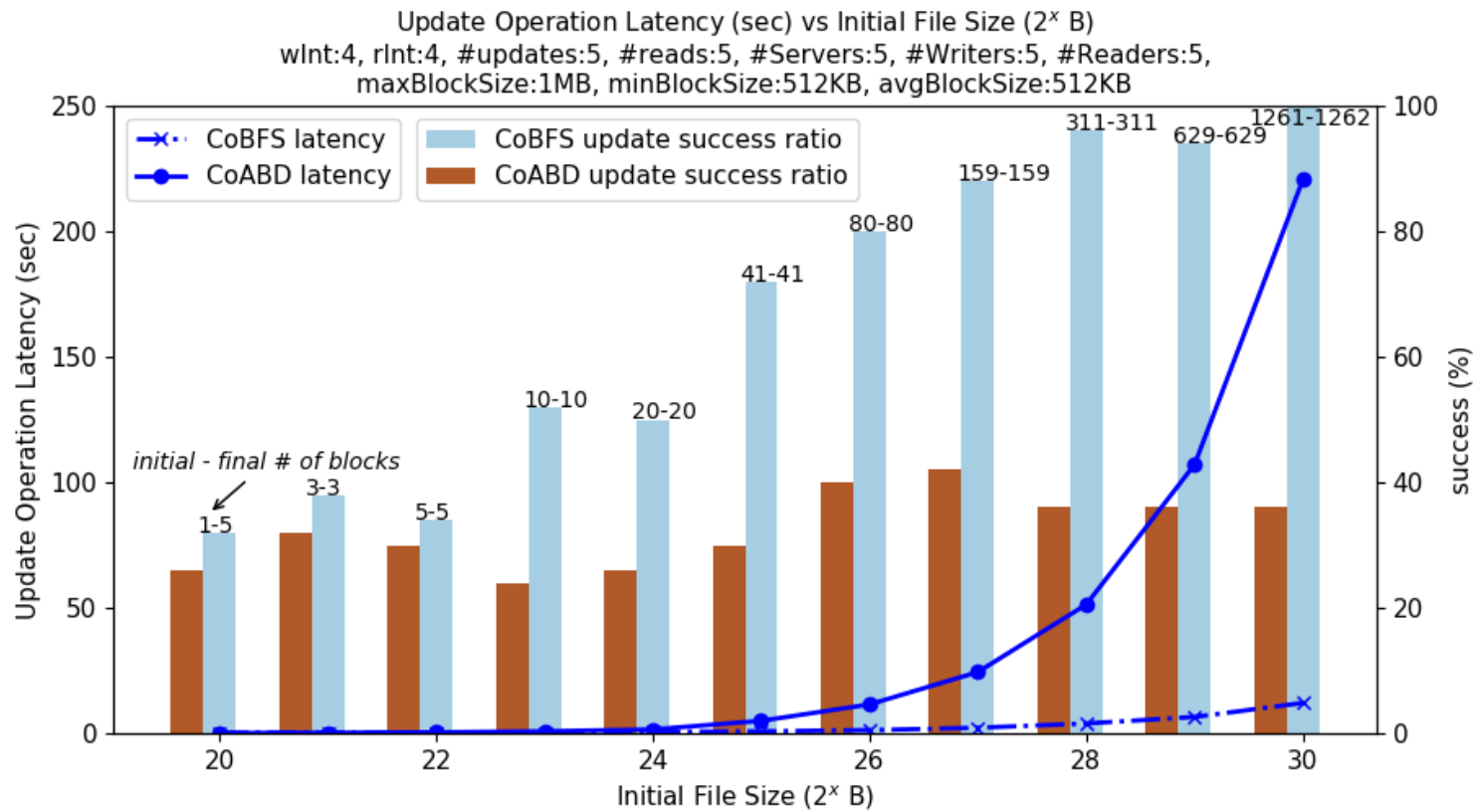


Rabin, M O. Fingerprinting by random polynomials, *Center for Research in Computing Techn., Aiken Computation Laboratory, Univ. , pp 15–18, 1981*

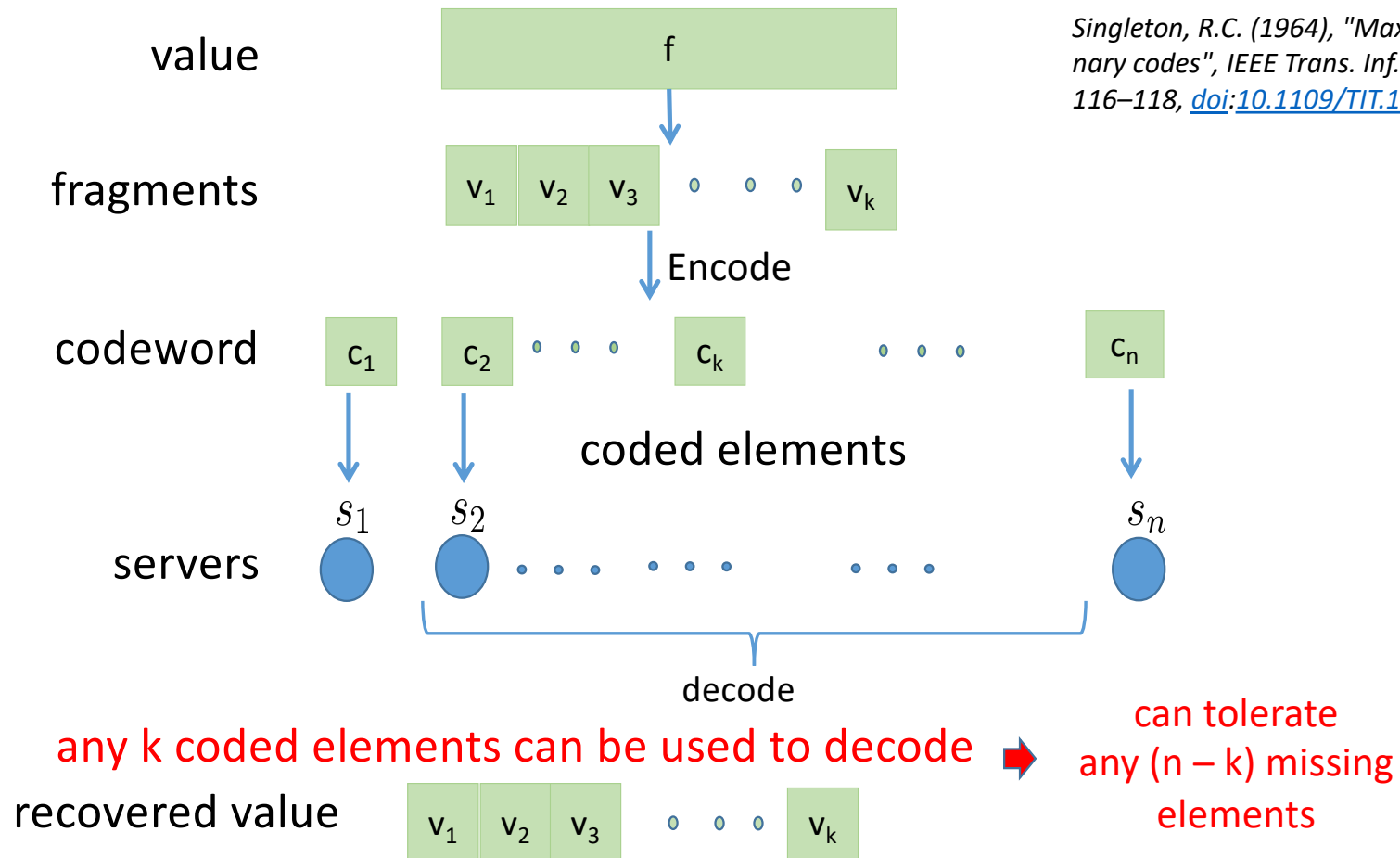
Solution 1: Experimental Results - Scalability



Solution 1: Experimental Results - Filesize

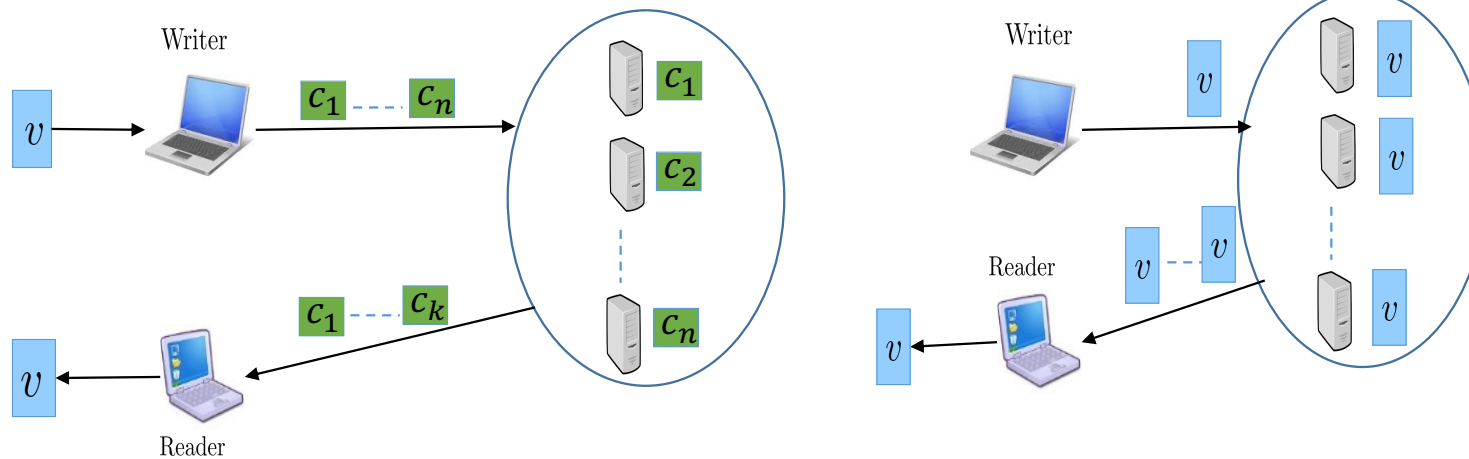


Solution 2: Erasure Coding ($[n, k]$ MDS Codes)



Singleton, R.C. (1964), "Maximum distance q -ary codes", *IEEE Trans. Inf. Theory*, **10** (2): 116–118, [doi:10.1109/TIT.1964.1053661](https://doi.org/10.1109/TIT.1964.1053661)

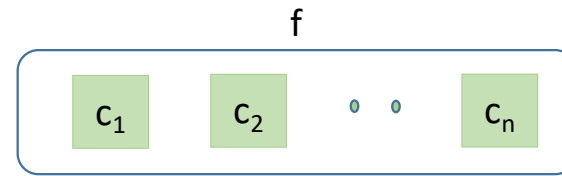
Solution 2: Erasure Coding vs Replication



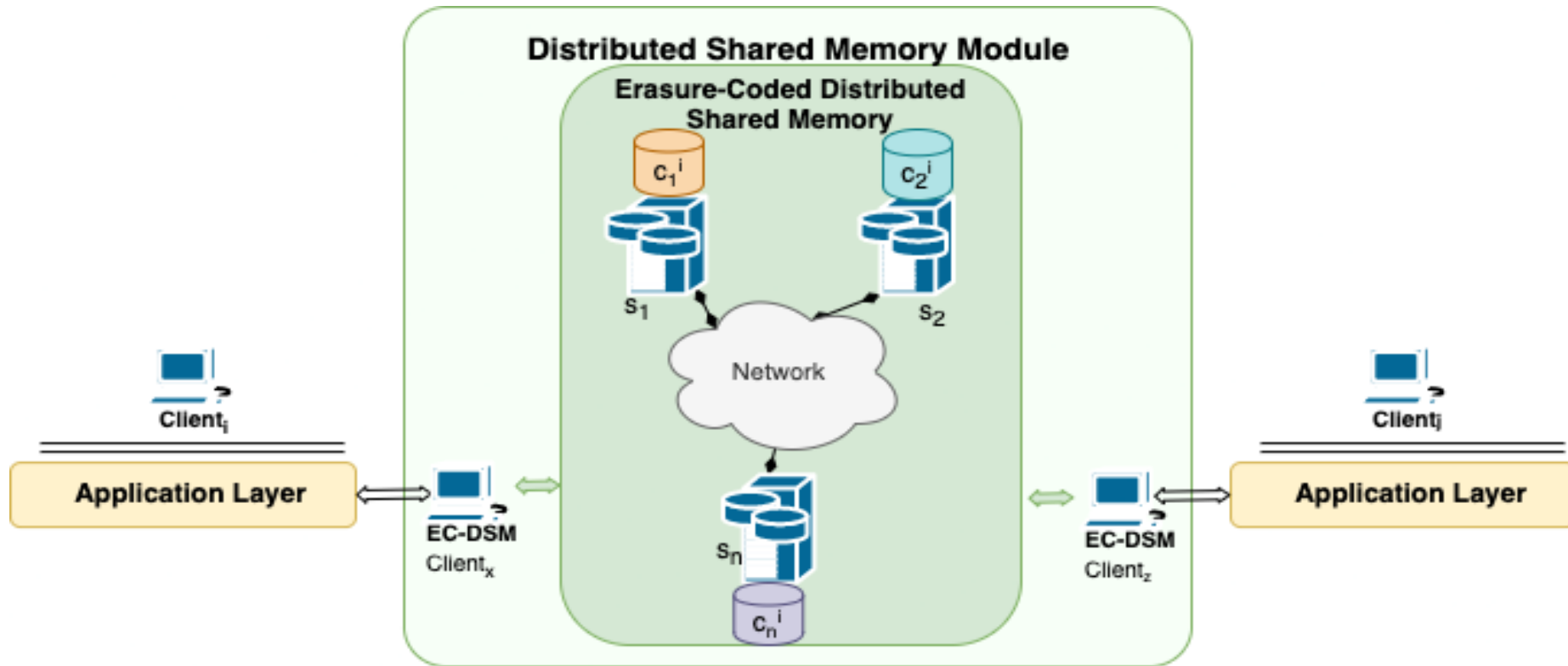
A well-designed algorithm has great potential to reduce storage and communication costs while using erasure codes

Solution 2: Algorithm CoEC

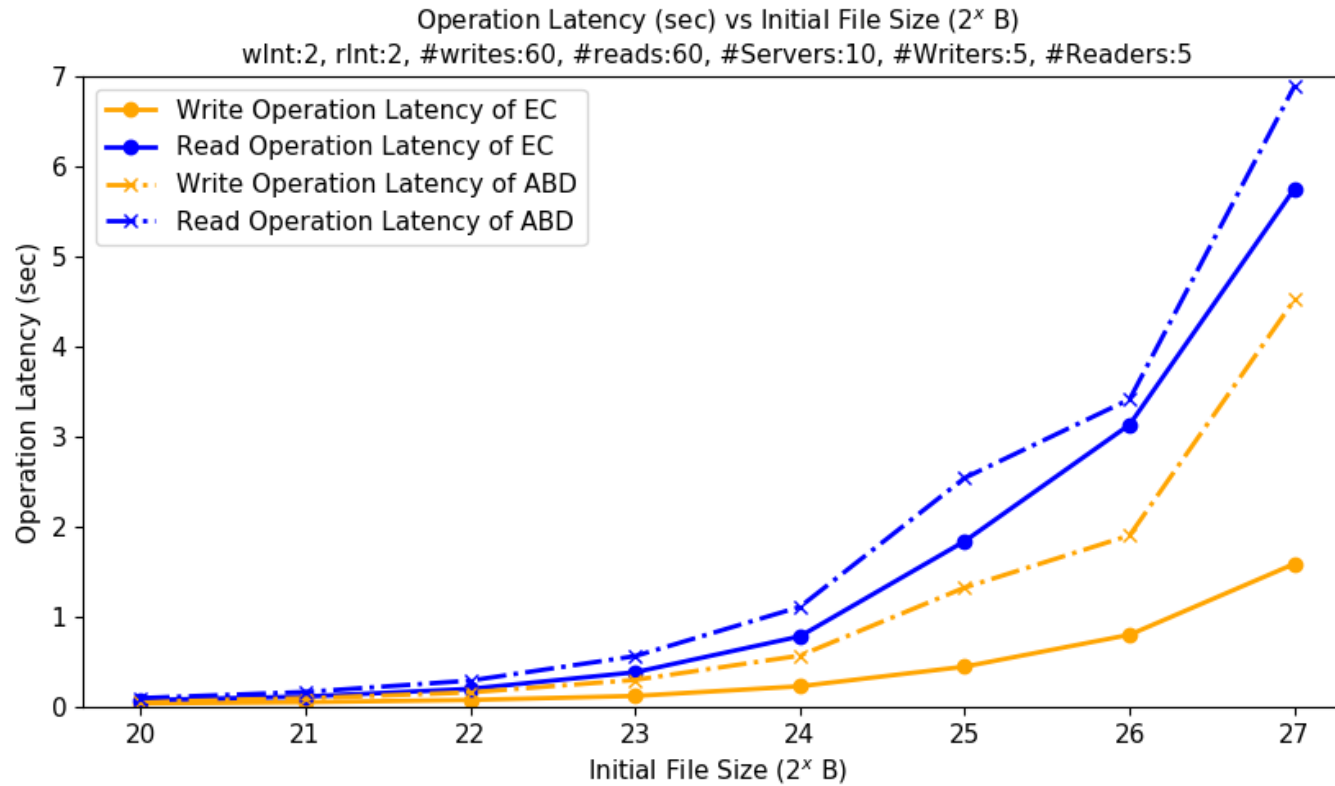
- Write Operation $\text{write}(f)$
 - Apply erasure coding on f
 - Send code c_i to server s_i
- Read Operation $\text{read}(f)$
 - Collect k codes from the servers
 - Decode and return the value of f
- Ensures Strong Consistency
- Does not prevent overwriting



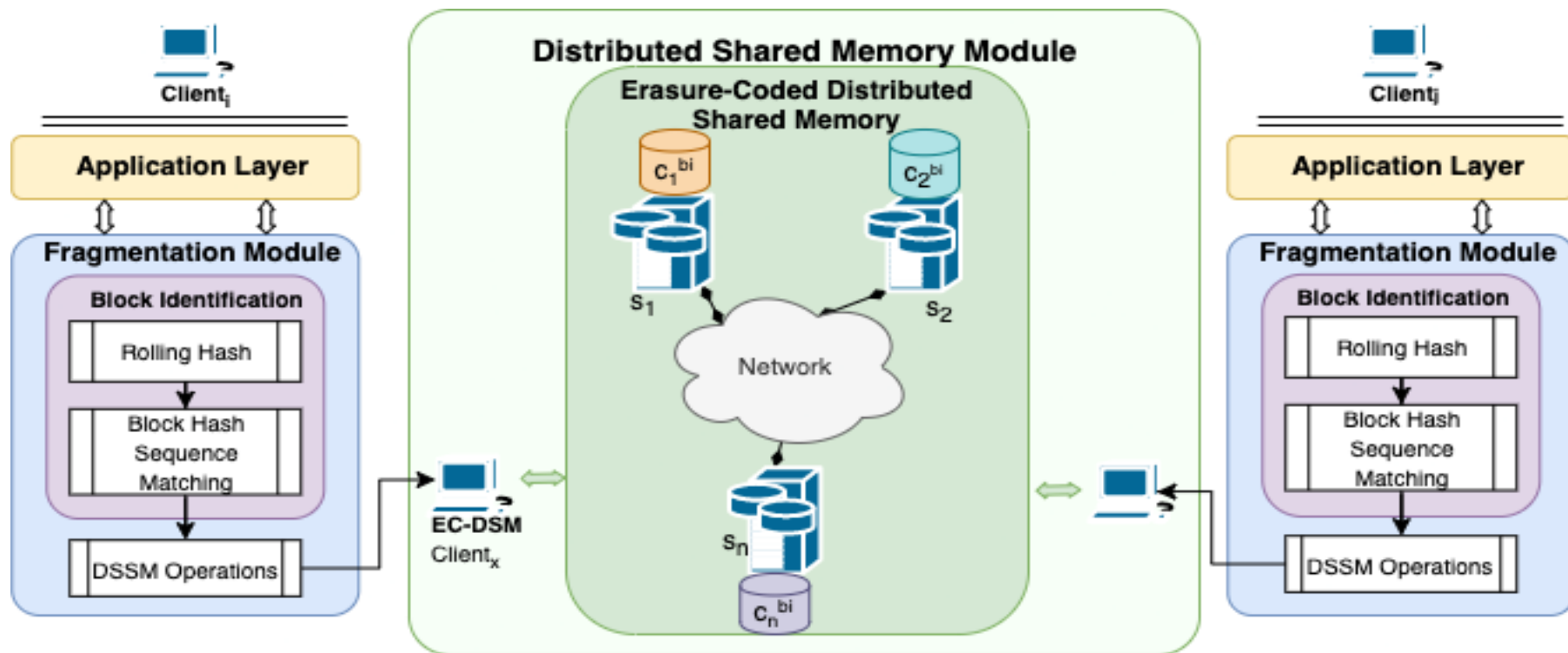
Solution 2: Erasure Coding Architecture



Solution 2: Experimental Results



Solution 3: Hybrid



What happens when things go wrong?

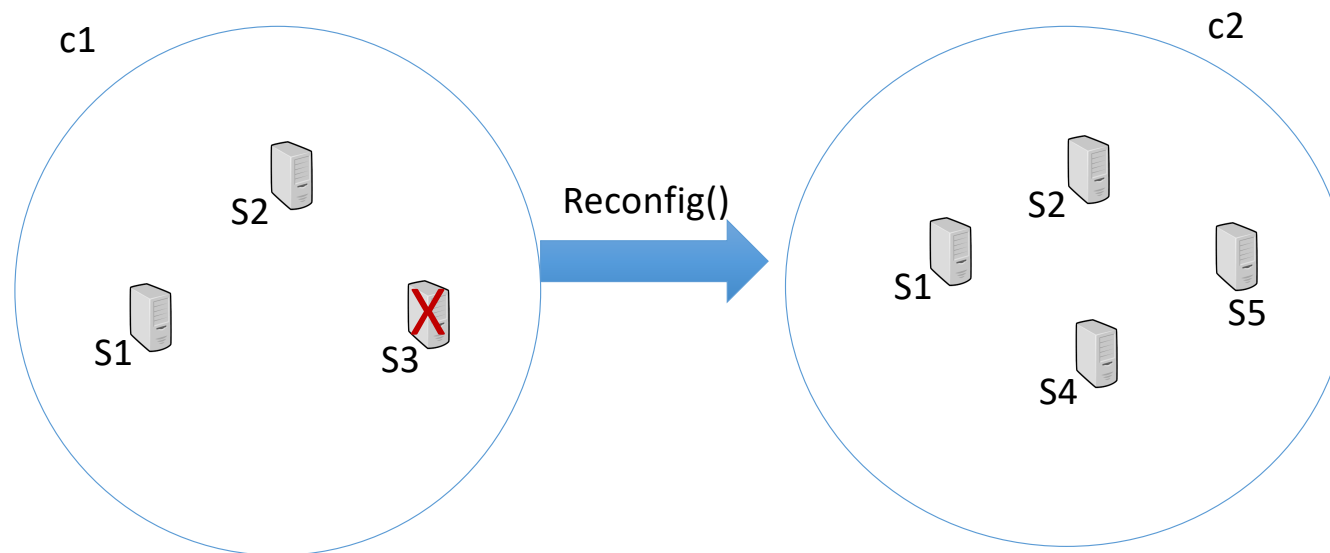
- Tolerate minority of failures
- What if more than minority fail?
- Replace failed with healthy servers => Reconfiguration

Challenge: Can we install a new configuration without stopping the service and without violating linearizability?



Re-Configuration Operation

- Change the configuration parameters (add/replace servers)
 - Due to failures
 - Due to admin maintenance



LYNCH, N., AND SHVARTSMAN, A. RAMBO: A reconfigurable atomic memory service for dynamic networks. *In Proceedings of 16th International Symposium on Distributed Computing (DISC) (2002)*, pp. 173–190.

ARES: A modular and adaptive reconfiguration protocol

Modular

- Read/Write operations are not aware of the underlying shared memory implementation
- They are using the same access primitives

Adaptive

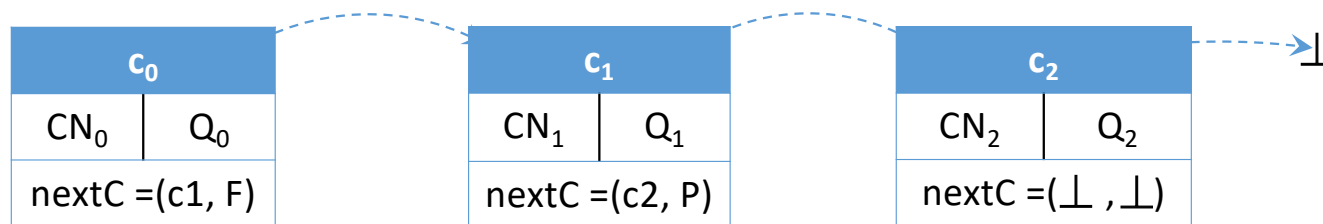
- Different shared memory algorithm may be used in every configuration
- Satisfying application demands

NICOLAOU, N., CADAMBE, V., KONWAR, K., PRAKASH, N., LYNCH, N., AND MEDARD, M. ARES: Adaptive, Reconfigurable, Erasure Coded, Atomic Storage. *In Proc. of ICDCS*, pp. 2195–2205 (2018)



Configuration Sequence

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



Reconfiguration Service

- A recon 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:    $cseq \leftarrow \text{read-config}(cseq)$ 
    $cseq \leftarrow \text{add-config}(cseq, c)$ 
10:   $\text{update-config}(cseq)$ 
    $cseq \leftarrow \text{finalize-config}(cseq)$ 
12: end operation
```

attempt get to the latest configuration } (1)
introduce the new configuration }
migrate the data to the new config } (2)
let servers know it is good to be finalized }



Reconfiguration Service Guarantees

For any two reconfig ops π_1, π_2 s.t. π_1 **before** π_2

- **Configuration Consistency**

π_1	<table border="1"><tr><td>c_0</td><td>c_1</td><td>c_2</td><td>...</td></tr></table>	c_0	c_1	c_2	...
c_0	c_1	c_2	...		
π_2	<table border="1"><tr><td>c_0</td><td>c_1</td><td>c_2</td><td>...</td></tr></table>	c_0	c_1	c_2	...
c_0	c_1	c_2	...		

- **Sequence Prefix**

π_1	<table border="1"><tr><td>c_0</td><td>c_1</td></tr></table>	c_0	c_1	
c_0	c_1			
π_2	<table border="1"><tr><td>c_0</td><td>c_1</td><td>c_2</td></tr></table>	c_0	c_1	c_2
c_0	c_1	c_2		

- **Sequence Progress**

π_1	<table border="1"><tr><td>$\langle c_0, F \rangle$</td><td>$\langle c_1, P \rangle$</td><td>$\langle c_2, P \rangle$</td></tr></table>	$\langle c_0, F \rangle$	$\langle c_1, P \rangle$	$\langle c_2, P \rangle$	
$\langle c_0, F \rangle$	$\langle c_1, P \rangle$	$\langle c_2, P \rangle$			
π_2	<table border="1"><tr><td>$\langle c_0, F \rangle$</td><td>$\langle c_1, P \rangle$</td><td>$\langle c_2, F \rangle$</td><td>...</td></tr></table>	$\langle c_0, F \rangle$	$\langle c_1, P \rangle$	$\langle c_2, F \rangle$...
$\langle c_0, F \rangle$	$\langle c_1, P \rangle$	$\langle c_2, F \rangle$...		

ARES: Experimental Results



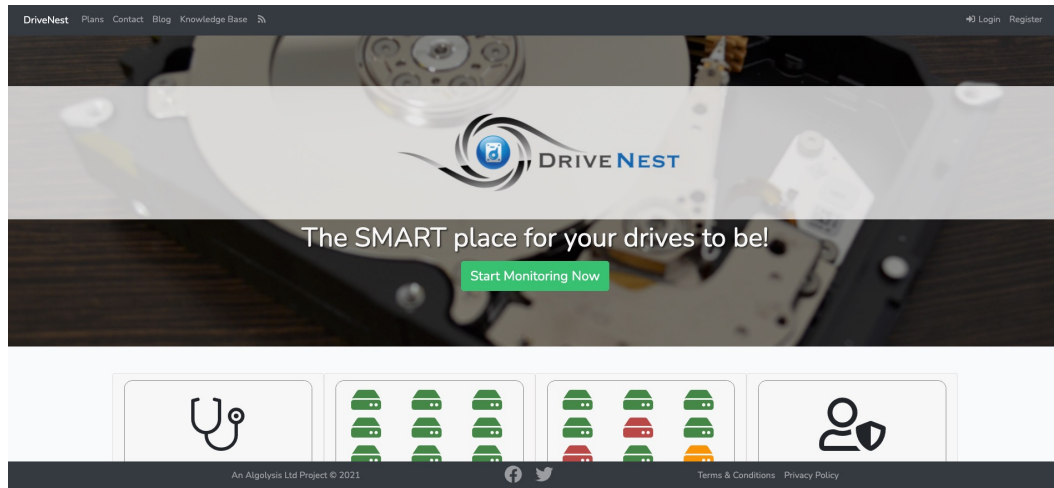
When shall we reconfigure?

- Frequent reconfigurations => Slow Down the service
- Infrequent reconfigurations => May make the service unavailable

Challenge: How can we determine when is the best time to reconfigure?



DriveNest: Monitoring Node Health

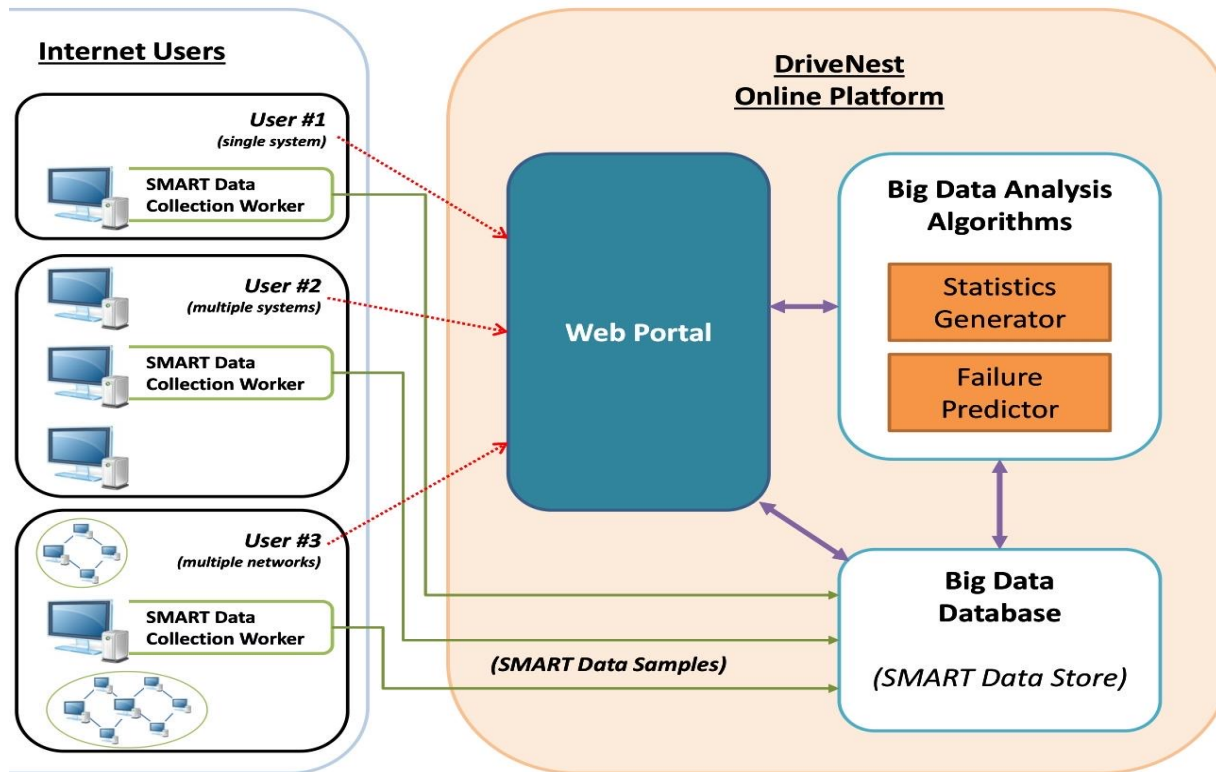


www.drivenest.com

- Crowdsourcing Platform
 - Collects data from diverse setups, locations, conditions.
- Monitor storage device health by collecting S.M.A.R.T data
- Predict soon-to-fail drives
 - Prediction performance relies on the report of failed drives
- Integration with ARES
 - Initiate recon operation to remove drives that are predicted to fail
 - Replace them with healthy nodes and migrate data



Drivenest: Architecture



- **DriveBird:** Data Collection Clients
- **Web platform:** View your drives
- **Prediction Engine:** Applies a number of machine learning/ deep learning algorithms to predict soon-to-fail drives

DriveNest: DriveBird Client

Production Submission Clients

Host a client on each machine you want to monitor for disk failures. These clients are agents that collect your drive's hardware (S.M.A.R.T) info and send it to DriveNest for analysis.
Pick the client suitable for your system and you will be up and running in seconds!

Platform	Download	checksum	Instructions
Python (cross-platform)	pyDriveBird-v1.0.7.tar.gz	pyDriveBird-v1.0.7.md5	readme

[Get Development Clients](#)

Development Submission Clients

Host a client on each machine you want to monitor for disk failures. These clients are agents that collect your drive's hardware (S.M.A.R.T) info and send it to DriveNest for analysis.
Pick the client suitable for your system and you will be up and running in seconds!

Platform	Download	checksum	Instructions
MS Windows	drivebird-1.1.7.20200603-windows.zip	drivebird-1.1.7.20200603-windows.md5	readme
Linux x64	drivebird-1.1.7.20200603-linux.tar.gz	drivebird-1.1.7.20200603-linux.md5	readme
Mac OS X 64-bit	drivebird-1.1.7.20200611-macosx.dmg	drivebird-1.1.7.20200611-macosx.md5	readme

[Get Stable Production Clients](#)

- Production Submission Clients
 - Python (cross platform)
- Development Submission Clients
 - GUI interface for all platforms
 - Tested up to Win 10 and MacOS Sierra

DriveNest: Status

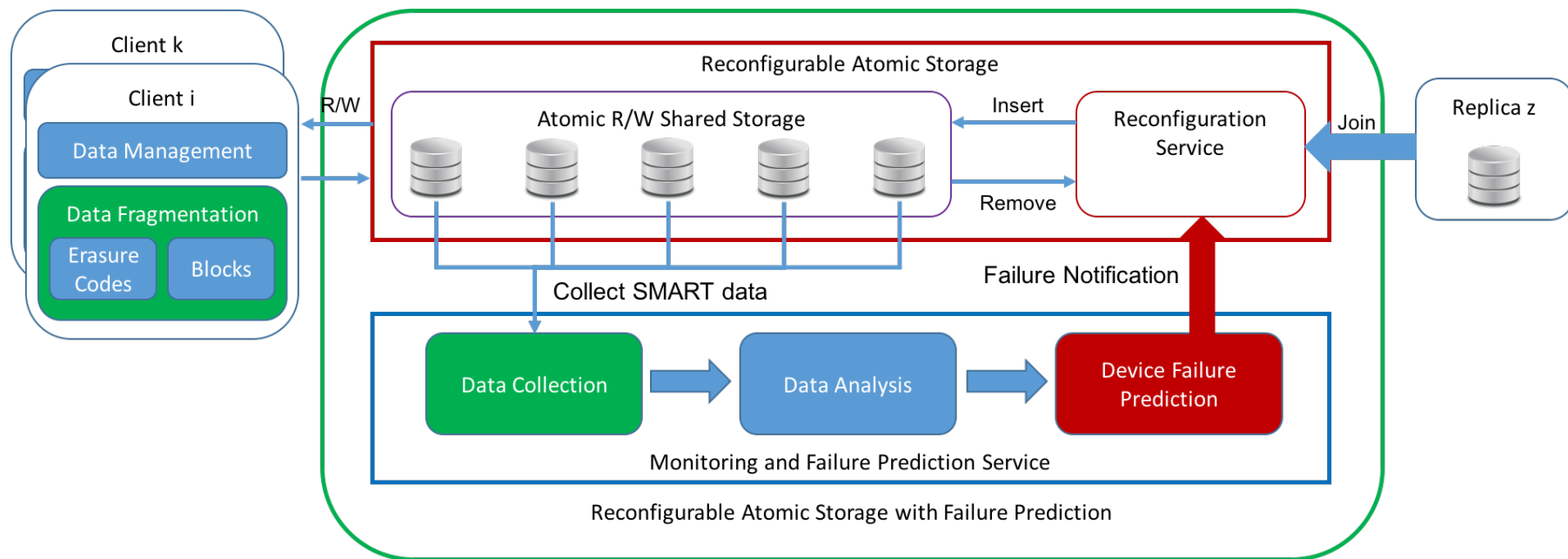
- Drivenest Collection Clients: **Alpha Testing**
- Drivenest Web Platform: **Alpha Testing**
- Drivenest Predictions: **Development Stage**

Feel free to register and give the service a test “drive”. 😊

We would be glad to hear your feedback



COLLABORATE: Overall Architecture



COLLABORATE: What's ahead

- Improve the Machine Learning outcomes of Drivenest for better failure prediction
- Embed the developed algorithms in a production level Distributed Storage Service!



Thank you!



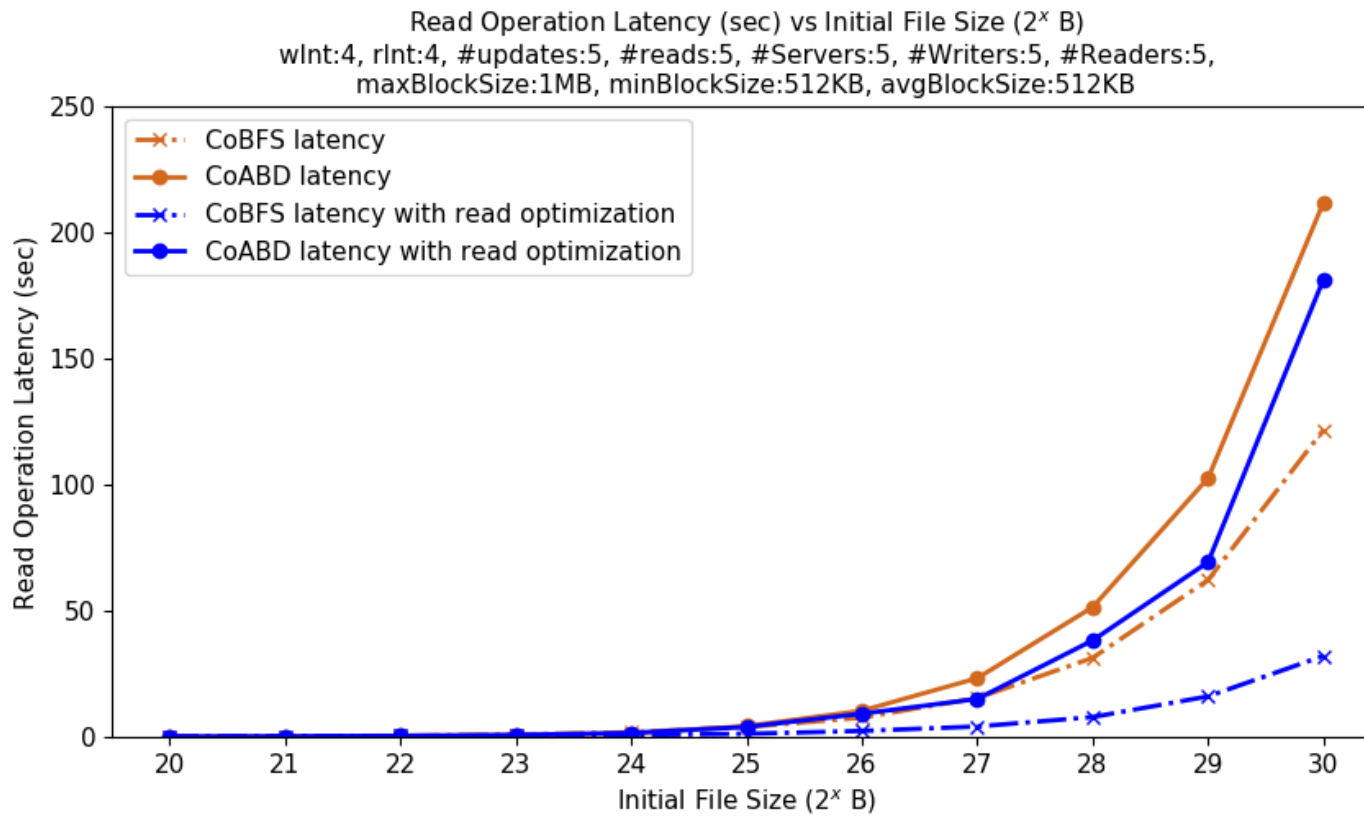
algolysis

algorithmic solutions

www.algolysis.com 

 research@algolysis.com

Solution 1: Experimental Results – Read Optimization



Solution 1: Experimental Results – Block Size

