

## Metadata Load Balancing Policies and Key-Value Stores

MICHAEL SEVILLA\*

### ABSTRACT

Enter the text of your abstract here.

---

\* *Corresponding author address:* Los Alamos National Laboratory  
E-mail: msevilla@ucsc.edu

## 1. Introduction

Migrating resources is a useful tool for optimizing performance for systems that service highly accessed data (e.g. key-value stores, file system metadata servers, etc.); data<sup>1</sup> can be distributed to alleviate overloaded servers or it can be concentrated, which gives the system the ability to leverage techniques that exploit locality, like bulk operations, multiple partition strategies, secondary indexes, and caching. Unfortunately, deciding which optimization to use is difficult to reason about, especially with the scale and complexity of today's HPC architectures.

Deciding between distribution and concentration is difficult because the techniques partition data differently; distribution spreads it across a larger set of nodes while concentration keeps it on a smaller set of nodes. As a result, selecting the wrong technique can have catastrophic consequences by compromising the performance optimizations designed for a certain scenario, e.g., migrating data to an already overloaded server or increasing the network hops by spreading data across an underutilized cluster.

This paper takes an API designed to migrate file system metadata and applies it to an HPC key-value store. The API helps control distribution and concentration by letting the administrator define how to migrate load, where to migrate load, and how much load to migrate. While designed for a different domain, this API encompasses many of the same properties we need for an HPC key-value store, namely:

- services small/frequent requests
- popularity drives distribution
- locality drives concentration

The HPC key-value store has functionality for exploiting locality but could further reduce The motivating example for this integration is finding the maximum  $x$  of the neighbors in a mesh of another maximum  $y$ . For example, finding the highest temperatures of the neighbors of mesh cells with the highest pressure. Given that the hash is defined by mesh location, finding the highest pressure is one RPC to each server. Unfortunately, even with an index based on the maximum pressures, finding the highest temperatures for *neighboring* cells with the highest pressures requires an additional RPC per server. This query has a both a high computational footprint, which suggests distribution to avoid hot spots, and data locality, which alternatively encourages concentration so the system can use its functionality for secondary indices, bulk operations, and key redistribution.

Specifically, the API helps us explore the space of solutions for this problem that has locality. For example, the

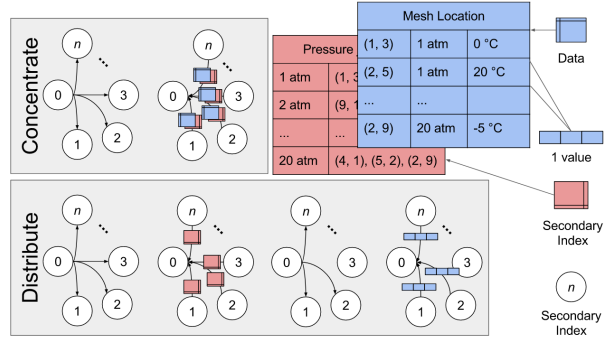


FIG. 1. Balancing for distribution takes more RPCs while migrating for concentration risk of overloading the client.

problem above has two solutions: (1) pull the index and re-query every server, or (2) pull the index and partial set of results that can be satisfied locally. Option 1 emphasizes distribution and incurs extra RPCs while option 2 opts to concentrate data at the expense of data transfer, consistency issues, and increased memory usage. Both options have advantages but inserting an API to control the mechanisms helps future programmers quickly evaluate both options and design policies for their workload.

In this paper, we make the following contributions:

1. prototype that controls concentration and distribution using the key redistribution and cursor type mechanisms from [1].
2. quantifies benefits of server/client-side caching, many small messages, and bulk operations.

<sup>1</sup>In this paper, we use the term “data” to refer to the partitioned key-value pairs AND file system metadata.

## 2. Background

Mantle speeds up file systems by making metadata access faster, leveraging the fact that file system metadata IO imposes small and frequent accesses on the underlying storage system. It uses resource migration to load balance file system metadata across a dedicated metadata cluster and

Mantle [3] is a programmable file system metadata balancer implemented on the file system above Ceph (CephFS).

system metadata IO does not scale like metadata IO [2]

HXHIM is a key-value store designed for HPC architectures and multi-dimensional data.

TABLE 1. This is a sample table caption and table layout. Enter as many tables as necessary at the end of your manuscript. Table from Lorenz (1963).

| <i>N</i> | <i>X</i> | <i>Y</i> | <i>Z</i> |
|----------|----------|----------|----------|
| 0000     | 0000     | 0010     | 0000     |
| 0005     | 0004     | 0012     | 0000     |
| 0010     | 0009     | 0020     | 0000     |
| 0015     | 0016     | 0036     | 0002     |
| 0020     | 0030     | 0066     | 0007     |
| 0025     | 0054     | 0115     | 0024     |

*Acknowledgments.* Start acknowledgments here.

References

[1] H. Greenberg, J. Bent, and G. Grider. MDHIM: A Parallel Key/Value Framework for HPC. In *7th USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage 15)*, 2015.

[2] D. Roselli, J. R. Lorch, and T. E. Anderson. A Comparison of File System Workloads. In *Proceedings of the Annual Conference on USENIX Annual Technical Conference, ATEC '00*, pages 4–4, 2000.

[3] M. A. Sevilla, N. Watkins, C. Maltzahn, I. Nassi, S. A. Brandt, S. A. Weil, G. Farnum, and S. Fineberg. Mantle: A Programmable Metadata Load Balancer for the Ceph File System. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC '15*, 2015.

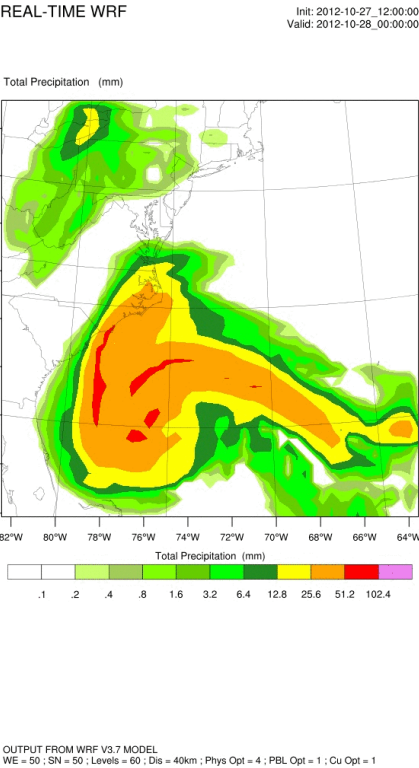


FIG. 2. Enter the caption for your figure here. Repeat as necessary for each of your figures.