# Using Hadoop Distributed and Deduplicated File System (HD2FS) in Astronomy

## Paul Bartus

School of Computer Science and Mathematics
Lake Superior State University
Sault Ste. Marie, Michigan, USA, 49783
email: pbartus@lssu.edu

**Abstract.** During the last years, the amount of data has skyrocketed. As a consequence, the data has become more expensive to store than to generate. The storage needs for astronomical data are also following this trend. Storage systems in Astronomy contain redundant copies of data such as identical files or within sub-file regions. We propose the use of the Hadoop Distributed and Deduplicated File System (HD2FS) in Astronomy. HD2FS is a deduplication storage system that was created to improve data storage capacity and efficiency in distributed file systems without compromising Input/Output performance. HD2FS can be developed by modifying existing storage system environments such as the Hadoop Distributed File System. By taking advantage of deduplication technology, we can better manage the underlying redundancy of data in astronomy and reduce the space needed to store these files in the file systems, thus allowing for more capacity per volume.

**Keywords.** Astronomy, Big Data, Deduplication, Education, File System, Hadoop, Storage System

## 1. Introduction

The use of deduplication has shown potential to remove storage redundancy in similar files across file systems. The concept of a file can be adapted to refer to chunks (data blocks) and file recipes. The file recipe for a file is a synopsis that contains a list of chunk identifiers (fingerprints) that comprise the file. Each chunk identifier can be created using a collision resistant hash over the contents of the block. Once the chunk identifiers in a file recipe have been obtained, they can be combined as prescribed in the file recipe to reconstruct the file. Hadoop Distributed File System (HDFS) is used to solve the storage problem of huge data, but does not provide a handling mechanism of duplicate files. HDFS is based on Google File System (GFS) and it operates on top of the operating system. HDFS has a name node, an optional secondary name node, and several data nodes. The name node is managing access and storing all metadata, such as file names, file attributes, and block locations.

## 2. Deduplication

Deduplication systems divide files into chunks (data blocks) and identify redundant chunks by comparing their identifiers (fingerprints). The chunk index contains the chunk identifiers of the stored chunks. Every deduplication system has an additional persistent index to store the information that is necessary to rebuild file contents based on file recipes. Chunk index (fingerprint) is a unique chunk identifier for each stored chunk and can be created using a collision resistant hash function such as SHA-256 over the chunk

(a)

(b)



HD2FS cluster.

Performance Analysis.

**Figure 1.** Comparative box plots for Write and Read times for 100 MB files using HD2FS cluster for different chunk sizes.

contents. File recipe is a list of chunk indexes and it will be the new file. The deduplication ratio is defined as the ratio between the original data size and the non redundant data size. A higher *DedupRatio* value shows a high redundancy in the file content while a lower ratio shows a high number of unique chunks. The deduplication ratio is given by the following formula:

$$DedupRatio = \frac{Total\ chunks}{Distinct\ chunks}$$

Bartus & Arzuaga (2017) presented the concept of a file-aware deduplication storage system and provided a study on the relation between the percentage of duplicate content and the percentage of duplicate chunks for the most common file types.

## 3. Hadoop distributed and deduplicated file system (HD2FS)

To improve data storage capacity and efficiency in distributed file systems, the use of Hadoop Distributed and Deduplicated File System (HD2FS) is proposed. Instead of storing multiple copies of the same astronomical data, HD2FS (Figure 1a) will store only the data that is different along with a map to reconstruct all data files. Bartus (2018); Bartus & Arzuaga (2018) tested HD2FS performance. A set of user files of size 100 MB each was written and then read from HD2FS to determine how the deduplication chunk size influence the overall performance. As Figure 1b shows, for very small chunk sizes such as 512 B, or 1 KB, the write times are relatively high compared to chunk size values over 2 KB. The no-dedup experiment was done with the original unmodified HDFS cluster, therefore without deduplication. HD2FS performance is very similar to the performance of the original HDFS, but with the big advantage of improving storage space requirements.

## 4. Conclusion

HD2FS was designed and implemented using Hadoop Distributed File System (HDFS) by adding deduplication. The advantages of using HD2FS have been presented. The results show that the obtained deduplication values are superior in the case of using HD2FS when compared to HDFS. This means that there is potential for HD2FS to be effectively integrated to improve storage management of big data in Astronomy.

## References

Bartus, P. 2018, *Using Deduplication to Improve Storage Efficiency in Distributed File Systems*, PhD Dissertation, University of Puerto Rico, Mayaguez Campus

Bartus, P., & Arzuaga, E. 2018, *Gdedup: Distributed file system level deduplication for genomic big data.*, IEEE International Congress on Big Data, July 2-7, 2018, San Francisco, CA, USA.

Bartus, P., & Arzuaga, E. 2017, *Using file-aware deduplication to Improve capacity in storage systems.*, IEEE Colombian Conference on Communications and Computing (COLCOM), pages 1–6.