

# Popularity-aware Balancer in HDFS based Cloud Storage

Thanda Shwe<sup>1</sup>, and Masayoshi Aritsugi<sup>2</sup>

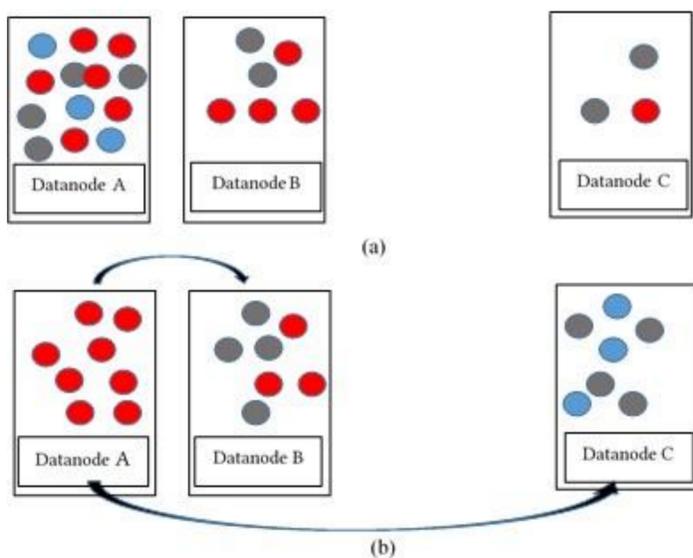
<sup>1</sup> Department of Computer Engineering and Information Technology, Mandalay Technological University, Mandalay, Myanmar

<sup>2</sup> Big Data Science and Technology, Faculty of Advanced Science and Technology, Kumamoto University, Japan, Myanmar

## Introduction

With wide adoption of distributed file systems for data-intensive applications, data storage in distributed storage systems can become storage space skew over time because of several reasons, namely, new storage nodes addition, random replica allocation policy, data block re-allocation in case of data node failures, massive file deletion load imbalance in the system because there is a guarantee that high utilized data nodes will be accessed more than least utilized data nodes. To cope with data storage skew, HDFS provides storage space balancer which is a utility for balancing data blocks across the storage devices of HDFS cluster.

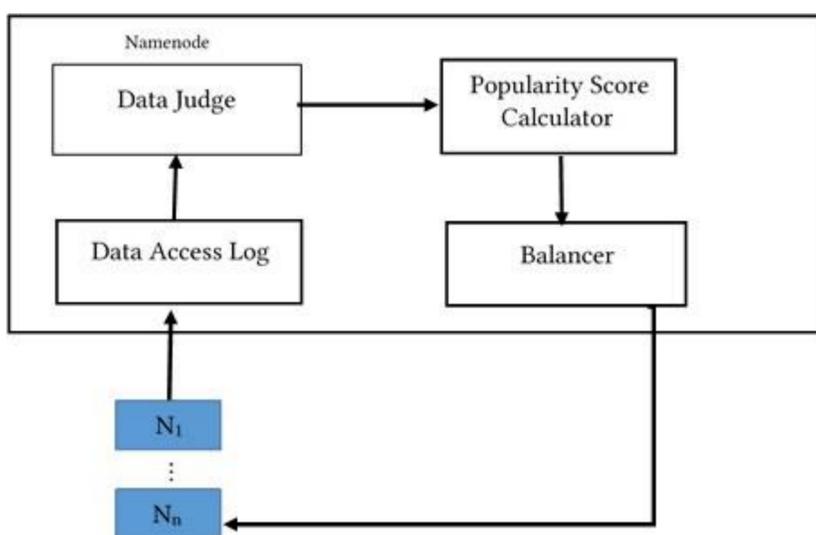
## Problem Statement



In above Figure (a), shows the data storage skew state where node A has been allocated a larger number of replicas than the others. Under this condition, intuitively, a lot of data access requests will be served by node A. To rebalance the data storage skew, HDFS's balancer migrates the data blocks from highly utilized nodes to low utilized nodes.

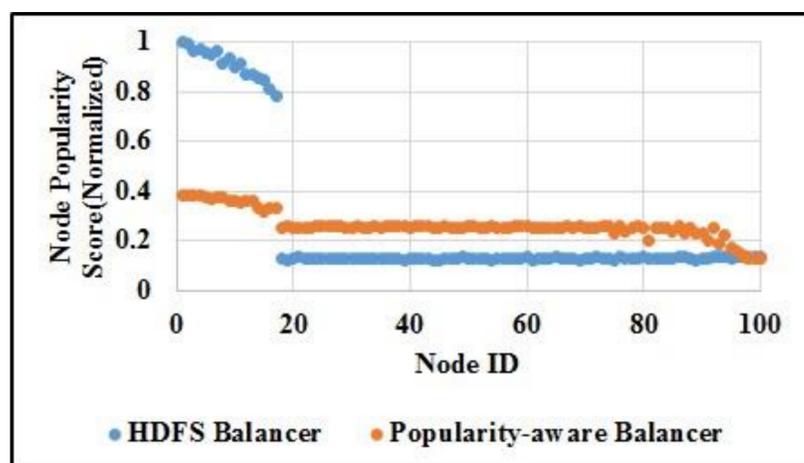
However, as shown in Figure(b), only the number of data blocks can be balanced by HDFS balancer. In datanode A, hot data are piled, resulting in creating hot spots in datanode A. Hotspots in the cluster can be reduced by appropriate placement of files. To prevent the concentration of popular data blocks, it is necessary for a balancer tool to consider the popularity of the data files and rebalance the cluster.

## Methods

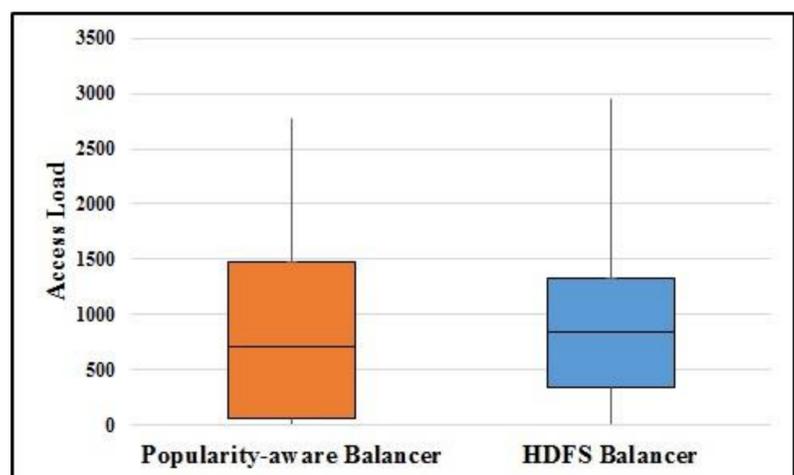


- **Data Access Log:** Record the data access for every file
- **Data Judge:** determines whether a particular data block is popular or unpopular.
- **Popularity Score Calculator:** Calculate the popularity Score based on block size and popularity level of blocks in each data node
- **Balancer:** performs the balancing activity based on node popularity score

## Evaluation Results



As expected, balancing with popularity-aware balancer shows only a few difference in node popularity score among the datanodes. This is because popularity-aware balancer balances the data blocks based not on the disk space but on the popularity level of the data and the size of each data block in the cluster.



The proposed popularity-aware balancer can distribute access load to the datanodes better than HDFS. This is because popularity-aware balancer allocated and balanced the data blocks based on node popularity score, leads to reduction in concentration of popular data in each datanode.

## Conclusion and Future Work

- The proposed system is validated through a set of several experiments and demonstrated its effectiveness using synthetic data.
- Although proposed scheme balance and allocate the data blocks based on the popularity of the data, data blocks popularity will be changing with time, and popular data at the moment may not be popular in the near future. Thus, we can predict the popularity of the data blocks in advance and allocate to balance among the cluster.