# Accepted Manuscript
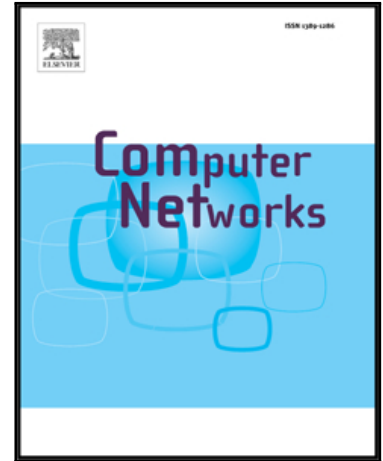
Dynamic Single Node Failure Recovery in Distributed Storage Systems

M. Itani, S. Sharafeddine, I. Elkabani

Please cite this article as: M. Itani, S. Sharafeddine, I. Elkabani, Dynamic Single Node Failure Recovery in Distributed Storage Systems, *Computer Networks* (2016), doi: 10.1016/j.comnet.2016.12.005

# Dynamic Single Node Failure Recovery in Distributed Storage Systems

M. Itani[a], S. Sharafeddine[b], I. Elkabani[a,c]

[a]*Beirut Arab University, Beirut, Lebanon*
[b]*Lebanese American University, Beirut, Lebanon*
[c]*Alexandria University, Alexandria, Egypt*

## Abstract

With the emergence of many erasure coding techniques that help provide reliability in practical distributed storage systems, we use fractional repetition coding on the given data and optimize the allocation of data blocks on system nodes in a way that minimizes the system repair cost. We selected fractional repetition coding due to its simple repair mechanism that minimizes the repair and disk access bandwidths together with the property of un-coded repair process. To minimize the system repair cost we formulate our problem using incidence matrices and solve it heuristically using genetic algorithms for all possible cases of single node failures. We then address three practical extensions that respectively account for newly arriving blocks, newly arriving nodes and variable priority files. A re-optimization mechanism for the storage allocation matrix is proposed for the first two extensions that can be easily implemented in real time without the need to redistribute original on-node blocks. The third extension is addressed by implementing variable fractional repetition codes which is shown to achieve significant cost reduction. The contributions of the paper are four fold: i. generating an optimized block distribution scheme among the nodes of a given data center for fixed and variable size blocks; ii. optimization of storage allocation under dynamic environments with data block arrivals; iii. optimization of storage allocation with newly added storage nodes; and iv. generating an effective block distribution scheme among the nodes by accounting for varying priorities among data blocks. We present a wide range of results for the various proposed algorithms and considered scenarios to quantify the achievable performance gains.

*Keywords:* Distributed storage systems, fractional repetition codes, failure recovery, genetic algorithms, variable fractional repetition codes

## 1. Introduction

Data centers are becoming a top priority for businesses and have become critical for the very functioning of big enterprises. Any interruptions in the data center operations might cause huge losses if corrective measures for interruptions or failures were not considered[1, 2]. Data centers use inexpensive hardware components that are prone to failure. In a study that examines a 3000-node production cluster of Facebook, node failures spike to more than 100 in a single day with an average failure rate of 22 nodes a day [3].

Thus there is a great need for data protection from device failures, and for mechanisms to quickly recover or at least mask the effects of node failures from users and connected devices with minimum performance cost. The easiest way for a storage system to tolerate failures and prevent data loss is to store replicas. This, however, results in decreased storage efficiency. Another alternative is to store encoded data using erasure coding [4] . Classical erasure codes transform a given message into a longer one (codeword) such that the original message can be recovered from the codeword. Although traditional erasure codes can reduce the storage overhead as compared to replication, extensive network resources are needed to repair a lost node. This is due to the fact that a surviving node should read and process all its data, then send a linear combination of them to the replacement node causing huge bandwidth consumption. To minimize the bandwidth consumed during the repair process, regenerating codes were introduced in the literature where a failed node can be recovered by connecting to a subset of surviving nodes and downloading one block of data from each [5].

In this work, we consider a family of regenerating erasure codes that provide exact and uncoded repair where a surviving node reads the exact amount of data it needs to send to a replacement node without any processing. This allows for a low complexity repair process which is achieved by fractional repetition (FR) codes that were first intro-