Kevin Keithley
CS 119, Fall 2023
5C Review, Google File System

## Category

This paper falls in between a description of a research prototype and an analysis of an existing system. It is more than a description of a research prototype because the system was already up and running when the paper was written in 2003. It also isn't quite a full analysis despite the inclusion of some measurement statistics. The paper seems like it's describing one of the first big steps forward in big data management, and the headings suggest more of a "here are some of the big things we were trying to address" more so than a systematic assessment of all the pieces of the system that you could take and try to assemble a similar system.

## Context

The paper points to three key pieces of related work. The first one is on large distributed file systems: - Scale and performance in a distributed system

This seems to have served as the background that the Google people were dealing with when creating their new system.

Two other papers seem to represent some new ways to think about implementing new systems: - Serverless network file systems - Swift: Using distributed disk striping to provide high I/O data rates

So working from a large distributed file system, they needed to incorporate elements from other systems to achieve their goals. They ended up using elements of the 2nd and 3rd papers in their design of the GFS.

## Correctness

Here is a simple list of assumptions that go into supporting the idea of a new way of thinking about implementing a large file system:

- An overview of the differences in design principles relative to earlier implementations
- Creating a new way that systems interact to better reflect their intended purposes
- Dealing with the increased number of faults because of a much larger system, and using commodity parts

I'm not an expert in any of the topics mentioned in this paper, but I believe the points made in the sections directly support the claims in the abstract, intro, and conclusion. In linear algebra I would say they seem to be both linearly independent, and they span the basis of the topic.

2003 definitely was a part of the rise in Google, and seeing what the leader in search was using behind the scenes also carries significant weight with it.

## Contributions

The main contribution of the paper appears to be a big step forward in an approach to big data management that gets so large that faults are inevitable and so large that it becomes impractically expensive to use non-commodity parts. Especially at the time of this paper in 2003, I find the methods developed to achieve probably time-honored goals of data storage and management to be quite impressive. The ideas presented in this paper appear to still be quite in line with the ideas in this course today.

## Clarity

This paper appears very well written. It is very clear about the scenario and goals that led to the development of the system. The authors explicitly state that despite having similar goals to previous file storage/management implementations, they needed to seriously rethink how to get there to create a system that will continue to work and evolve in the future. To that end, the sections and subsections seem to directly relate to the

key differentiating factors that allowed them to re-envision the future of big data storage. They reference previous works that share elements of the system that they implemented, while weaving them together into a better system.