

# Matrix Multiplication MapReduce

Jesus Arencibia Falcon

July 2024

## Abstract

In this study, we tackle the problem of matrix multiplication using the MapReduce programming paradigm. Matrix multiplication is a fundamental operation in various data processing and scientific computing applications. We implemented a sparse matrix multiplication algorithm on Hadoop, utilizing the Python MRJob library to distribute the workload. We evaluated the efficiency of the approach in terms of execution time and scalability.

## 1 Introduction

Matrix multiplication is a complex mathematical operation with wide-ranging applications in engineering, computational sciences, and data analysis. However, when matrices are large and sparse, the operation can become time and resource-intensive. The use of MapReduce allows us to break down the problem into manageable sub-problems that can be executed in parallel, optimizing computational efficiency.

This work focuses on the implementation of a matrix multiplication algorithm using MapReduce. The methodology involves the preparation and processing of matrices, as well as the distributed execution of the algorithm in a Hadoop environment.

## 2 Methodology

To prepare the matrices for multiplication, we first convert the matrices into a format suitable for MapReduce processing. This involves reading the matrices in .mtx file format and converting them into a text format that can be processed by our MapReduce job. The conversion process ensures that matrices are stored in a compressed format to optimize

We use the Python MRJob library to define the mapping and reducing steps. In the mapping step, elements of the matrices are distributed into key-value pairs representing the positions and values of the matrices. These pairs are sent to different processing nodes. In the reducing step, the values from the key-value pairs are combined to perform the matrix multiplication. For each entry, the corresponding matrix values are gathered, and the dot product is calculated to obtain the final element of the resultant matrix.

The MapReduce job is structured in several stages:

Mapper Initialization: Matrix dimensions are read from external files.

Mapper: Input lines are processed to generate key-value pairs representing contributions from each element of matrices A and B.

Reducer: Values from matrices A and B are combined to calculate the necessary dot products. At the end of the process, the resultant matrix from the multiplication is obtained.

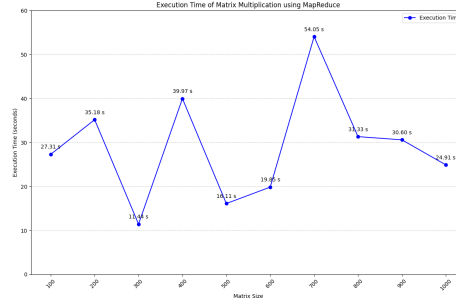
### 3 Experiments

Multiple experiments were conducted to evaluate the efficiency of the implemented matrix multiplication algorithm. We used matrices of different sizes and densities to measure execution time and scalability.

The experiments were carried out in a distributed Hadoop environment using MRJob to manage the MapReduce jobs. Some of the results obtained include:

Execution Time: The total execution time of the algorithm was measured for different matrix sizes. Results showed that the MapReduce approach efficiently handles large and sparse matrices.

Scalability: The implementation demonstrated scalability, as execution time increased linearly with matrix size.



### 4 Conclusion

Matrix multiplication using the MapReduce paradigm is shown to be an efficient solution for processing large volumes of data. The implementation of MRJob in Python allows for workload distribution, leveraging the parallel processing capabilities of Hadoop. The conducted experiments demonstrate that the approach is feasible and efficient, particularly for large and sparse matrices. The use of compressed formats and parallelization techniques significantly reduces execution time and resource requirements.

## 5 Future Work

The study on optimizing matrix multiplication techniques, particularly focusing on reduce matrices, opens avenues for further exploration and advancement. And there are several possible directions for improving and expanding this work:

**Algorithm Optimization:** Research and develop new strategies to optimize the matrix multiplication algorithm, including more advanced parallelization techniques and the use of more efficient data structures.

**Production Environment Evaluation:** Perform tests in production environments with real-world data to validate the effectiveness and efficiency of the approach in practical applications.

**Extension to Other Problems:** Explore the application of MapReduce to other computational problems related to matrices and linear algebra, thereby expanding the scope and utility of the approach.