

# Page Rank

Uros Babic

June 14, 2024

## Abstract

This project implements and analyzes the PageRank algorithm using sequential, parallel, and distributed approaches. The goal is to compare the performance and scalability of each method in calculating the importance of pages in a randomly generated graph. The source code for this project is available on GitHub at <https://github.com/uos4658/PageRankJava>.

## 1 Introduction

PageRank is an algorithm used by search engines to rank web pages in their search results. It was developed by Larry Page and Sergey Brin at Stanford University and forms the core of Google's ranking algorithm. This project aims to implement the PageRank algorithm in three different ways: sequentially, using parallel processing, and in a distributed manner. The performance of these implementations will be compared based on computation time and scalability.

## 2 Literature Review

The original PageRank algorithm has been extensively studied and improved upon since its inception. Various optimization techniques and alternative approaches have been proposed in the literature to handle large-scale web data. Research has shown that parallel and distributed computing can significantly speed up the PageRank computation, making it feasible to process web-scale graphs efficiently.

## 3 Project Definition

### 3.1 Problem Definition

The problem addressed in this project is to calculate the PageRank of vertices in a graph representing web pages and their links. The challenge is to do this efficiently for large graphs using different computational approaches.

### 3.2 Objectives

- Implement the PageRank algorithm sequentially.

- Implement the PageRank algorithm using parallel processing.
- Implement the PageRank algorithm in a distributed computing environment using MPI.
- Compare the performance of these implementations.

### 3.3 Scope

This project focuses on the implementation and performance comparison of the PageRank algorithm. It does not cover other ranking algorithms or optimizations beyond basic parallel and distributed techniques.

## 4 System Analysis and Design

### 4.1 System Architecture

The system consists of modules for graph generation, PageRank computation (sequential, parallel, and distributed), and result evaluation. The distributed implementation uses MPI for communication between processes.

### 4.2 Use Case Diagrams

The primary use case is the calculation of PageRank values, which involves users (or automated scripts) initiating the computation and retrieving the results.

## 5 System Implementation

### 5.1 Graph Generation

A random weighted graph is generated with a specified number of vertices and edges. The weights represent the importance of the links between the vertices.

Listing 1: Graph Generation

```
public static int [][] generateRandomGraph(int numVertices, int numEdges) {
    int [][] graph = new int[numVertices][numVertices];
    Random random = new Random();

    for (int i = 0; i < numEdges; i++) {
        int source = random.nextInt(numVertices);
        int target = random.nextInt(numVertices);

        while (source == target) {
            target = random.nextInt(numVertices);
        }

        int weight = random.nextInt(10) + 1;
        graph[source][target] = weight;
    }
}
```

```

    return graph;
}

```

## 5.2 PageRank Calculation

### 5.2.1 Sequential PageRank

The sequential implementation computes the PageRank by iterating through the vertices and updating their ranks based on the incoming links.

Listing 2: Sequential PageRank Calculation

```

public class SequentialPageRank {
    public static double[] calculatePageRank(int [][] graph, double
dampingFactor, int maxIterations) {
        int numVertices = graph.length;
        double[] pageRanks = new double[numVertices];
        double[] newPageRanks = new double[numVertices];

        for (int i = 0; i < numVertices; i++) {
            pageRanks[i] = 1.0 / numVertices;
        }

        for (int iteration = 0; iteration < maxIterations; iteration++) {
            for (int i = 0; i < numVertices; i++) {
                double sum = 0.0;
                for (int j = 0; j < numVertices; j++) {
                    if (graph[j][i] != 0) {
                        sum += pageRanks[j] / countOutgoingLinks(graph, j);
                    }
                }
                newPageRanks[i] = (1 - dampingFactor) / numVertices +
dampingFactor * sum;
            }
            pageRanks = newPageRanks.clone();
        }

        return pageRanks;
    }
}

```

### 5.2.2 Parallel PageRank

The parallel implementation uses multithreading to compute the PageRank more efficiently by distributing the computation of each vertex's rank across multiple threads.

Listing 3: Parallel PageRank Calculation

```

public class ParallelPageRank {
    public static double[] calculatePageRankParallel(int [][] graph,
double dampingFactor, int maxIterations) {

```

```

    int numVertices = graph.length;
    double[] pageRanks = new double[numVertices];
    double[] newPageRanks = new double[numVertices];

    for (int i = 0; i < numVertices; i++) {
        pageRanks[i] = 1.0 / numVertices;
    }

    for (int iteration = 0; iteration < maxIterations; iteration++) {
        final double[] finalPageRanks = pageRanks.clone();
        Thread[] threads = new Thread[numVertices];
        for (int i = 0; i < numVertices; i++) {
            final int vertex = i;
            threads[i] = new Thread(() -> {
                double sum = 0.0;
                for (int j = 0; j < numVertices; j++) {
                    if (graph[j][vertex] != 0) {
                        sum += finalPageRanks[j] /
                            Main.countOutgoingLinks(graph, j);
                    }
                }
                newPageRanks[vertex] = (1 - dampingFactor) / numVertices
                    + dampingFactor * sum;
            });
            threads[i].start();
        }

        for (Thread thread : threads) {
            try {
                thread.join();
            } catch (InterruptedException e) {
                e.printStackTrace();
            }
        }

        pageRanks = newPageRanks.clone();
    }

    return pageRanks;
}
}

```

### 5.2.3 Distributed PageRank

The distributed implementation uses MPI to distribute the computation across multiple processes.

Listing 4: Distributed PageRank Calculation

```

public class DistributivePageRank {
    public static double[] calculatePageRankDistributive(int[][] graph,

```

```

    double dampingFactor,    int maxIterations, String[] args)
    throws MPIException {
        MPI.Init(args);
        int numVertices = graph.length;
        double[] pageRanks = new double[numVertices];
        double[] newPageRanks = new double[numVertices];

        for (int i = 0; i < numVertices; i++) {
            pageRanks[i] = 1.0 / numVertices;
        }

        int rank = MPI.COMM_WORLD.Rank();

        for (int iteration = 0; iteration < maxIterations; iteration++) {
            double sum = 0.0;
            for (int j = 0; j < numVertices; j++) {
                if (graph[j][rank] != 0) {
                    sum += pageRanks[j] / Main.countOutgoingLinks(graph, j);
                }
            }
            newPageRanks[rank] = (1 - dampingFactor) / numVertices +
                dampingFactor * sum;
            MPI.COMM_WORLD.Allgather(newPageRanks, rank, 1, MPI.DOUBLE,
                pageRanks, 0, numVertices, MPI.DOUBLE);

            pageRanks = newPageRanks.clone();
        }
        MPI.Finalize();
        return pageRanks;
    }
}

```

## 6 System Testing and Evaluation

### 6.1 Test Strategy

The PageRank implementations are tested by generating random graphs and verifying the correctness of the PageRank values. Each implementation is also tested for performance by measuring the computation time.

### 6.2 Test Plan

- Generate random graphs of varying sizes.
- Calculate PageRank using sequential, parallel, and distributed implementations.
- Measure and compare the computation times.
- Verify that the PageRank values are consistent across implementations.

- Compare results with known PageRank values for validation.

### 6.3 Testing Results and Error Analysis

During initial testing, an error was identified due to the inclusion of the `printGraph` function. This error significantly increased the computation time. After removing the `printGraph` function, the execution times were re-evaluated, showing significant improvements.

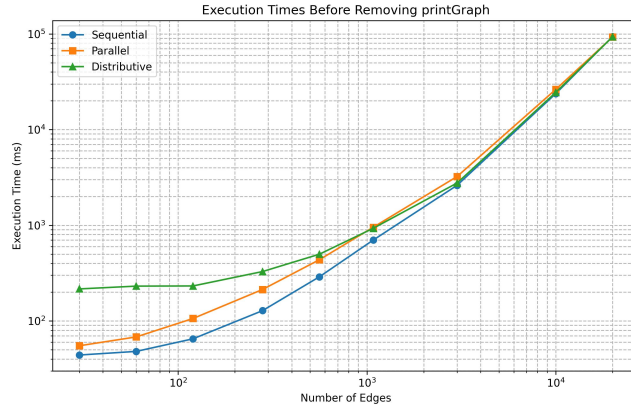


Figure 1: Execution times before removing `printGraph`

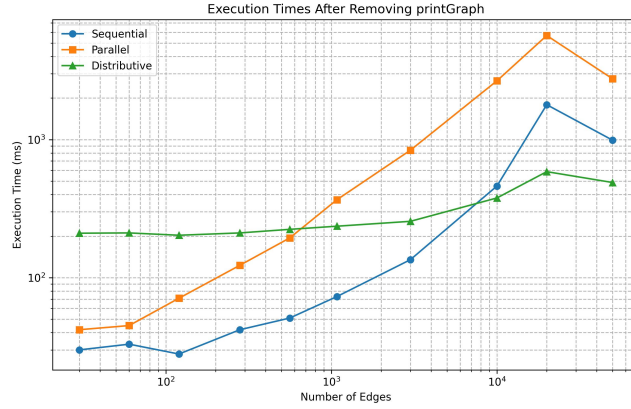


Figure 2: Execution times after removing `printGraph`

## 7 Conclusion

The PageRank algorithm was successfully implemented using sequential, parallel, and distributed approaches. Each method has its trade-offs in terms of

computation time and complexity. The choice of implementation depends on the specific requirements and resources available. The results demonstrate that parallel and distributed implementations can significantly reduce computation time, making them suitable for large-scale graph processing.

## 8 Future Work

Future improvements could include optimizing the parallel and distributed implementations for better performance and scalability. Additionally, experimenting with different graph structures and sizes could provide further insights into the algorithm's behavior. Advanced techniques such as graph partitioning and load balancing could also be explored to enhance the distributed implementation.

## 9 Literature Review

The original PageRank algorithm, developed by Larry Page and Sergey Brin, has been extensively studied and improved upon since its inception. Various optimization techniques and alternative approaches have been proposed in the literature to handle large-scale web data. Research has shown that parallel and distributed computing can significantly speed up the PageRank computation, making it feasible to process web-scale graphs efficiently.

Here is a list of references that were consulted for this report:

1. Abdul Bari. PageRank Algorithm. Available at: [https://www.youtube.com/watch?v=CsvyPNdQAHg&ab\\_channel=AbdulBari](https://www.youtube.com/watch?v=CsvyPNdQAHg&ab_channel=AbdulBari).
2. Sepandar D. Kamvar, Taher H. Haveliwala, and Gene H. Golub. "Extrapolation methods for accelerating PageRank computations." In Proceedings of the 12th international conference on World Wide Web. ACM, 2003.
3. Amy N. Langville and Carl D. Meyer. "Google's PageRank and beyond: The science of search engine rankings." Princeton University Press, 2006.
4. Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. "The PageRank citation ranking: Bringing order to the web." Stanford InfoLab, 1999.
5. Sergey Brin and Lawrence Page. "The anatomy of a large-scale hyper-textual web search engine." Computer Networks and ISDN Systems 30, no. 1-7 (1998): 107-117.
6. Taher H. Haveliwala. "Efficient computation of PageRank." Technical Report, Stanford University, 1999.

These resources provide a comprehensive overview and deep understanding of the PageRank algorithm, its computation techniques, and applications in web-scale data processing.