



UNIVERSITÀ DI PISA



Master's Degree in Artificial Intelligence and Data Engineering

Final cloud computing project for single students:

Inverted Index

Instructors:

Prof. Vallati Carlo

Prof. Puliafito Carlo

Student:

Adrien Koumgang Tegantchouang

Academic Year 2024/2025

Chapter 1

Documentation

1.1 Introduction

This project implements an Inverted Index using Apache Hadoop's MapReduce programming model. An inverted index is a fundamental data structure in information retrieval systems, mapping each unique word in a corpus to the list of documents in which it appears. The goal of the project is to build a scalable solution capable of processing large text datasets in a distributed environment.

This project also includes a performance comparison between the two implementations, measuring metrics such as execution time, memory usage, and shuffle volume.

1.2 System Setup

1.2.1 Hadoop Installation

- Followed the official Hadoop tutorial for **pseudo-distributed mode**: <https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html>
- Key steps:
 - Installed Java 8+, SSH, and Hadoop 3.X.
 - Configured: 'core-site.xml' (HDFS settings), 'hdfs-site.xml' (replication factor = 1 for single-node), 'mapred-site.xml' (YARN resource management) and 'yarn-site.xml' (MapReduce job scheduling)
 - Formated HDFS and started services ('start-dfs.sh', 'start-yarn.sh')

1.2.2 Dataset Selection

I use three datasets of varying sizes: **Small** (1MB), **Medium** (100MB) and **Large** (1GB).

1.3 Implementation

1.3.1 Java Solution: MapReduce Design

- **Mapper(filename, text):** Split text into words and for each word, emit ‘(word, filename)’
- **Combiner(word, [filenames]):** Remove duplicates (if same word appears multiple times in a file) and emit ‘(word, filename)’
- **Reducer(word, [filenames]):** Aggregate all filenames for the word and emit ‘(word, [doc1.txt, doc2.txt, ...])’

1.3.2 Key Java Classes

Inverted Index Mapper: Tokenizes input and emits ‘(word, filename)’ pairs

```
1 public class InvertedIndexMapper extends Mapper<LongWritable, Text, Text, Text> {
2     private Text word = new Text();
3     private Text filename = new Text();
4
5     @Override
6     protected void map(LongWritable key, Text value, Context context) throws
7         IOException, InterruptedException {
8         FileSplit fileSplit = (FileSplit) context.getInputSplit();
9         String filenameStr = fileSplit.getPath().getName();
10        filename.set(filenameStr);
11
12        StringTokenizer tokenizer = new StringTokenizer(value.toString());
13        while (tokenizer.hasMoreTokens()) {
14            word.set(TextUtils.cleanToken(tokenizer.nextToken()));
15            context.write(word, filename);
16        }
17    }
```

Listing 1.1: Inverted Index Mapper

Inverted Index Combiner: Deduplicates filenames per word.

```

1 public class InvertedIndexCombiner extends Reducer<Text, Text, Text, Text> {
2     private Text result = new Text();
3
4     @Override
5     protected void reduce(Text key, Iterable<Text> values, Context context) throws
6         IOException, InterruptedException {
7         Set<String> uniqueFiles = new HashSet<>();
8         for (Text val : values) {
9             uniqueFiles.add(val.toString());
10        }
11
12        StringBuilder fileList = new StringBuilder();
13        for (String file : uniqueFiles) {
14            fileList.append(file).append(" ");
15        }
16
17        result.set(fileList.toString().trim());
18        context.write(key, result);
19    }
20 }

```

Listing 1.2: Inverted Index Combiner

Inverted Index Reducer: Aggregates results into the final inverted index.

```

1 public class InvertedIndexReducer extends Reducer<Text, Text, Text, Text> {
2     private Text result = new Text();
3
4     @Override
5     protected void reduce(Text key, Iterable<Text> values, Context context) throws
6         IOException, InterruptedException {
7         Set<String> uniqueFiles = new HashSet<>();
8         for (Text val : values) {
9             // uniqueFiles.add(val.toString());
10            String[] files = val.toString().split(" ");
11            uniqueFiles.addAll(Arrays.asList(files));
12        }
13
14        StringBuilder fileList = new StringBuilder();
15        for (String file : uniqueFiles) {
16            fileList.append(file).append(" ");
17        }
18    }
19 }

```

```

18         result.set(fileList.toString().trim());
19         context.write(key, result);
20     }
21 }

```

Listing 1.3: Inverted Index Reducer

1.3.3 Optimizations

Used ‘setup()’ to initialise configurations and applied **in-mapper combining: Inverted Index In Mapper**.

```

1 public class InvertedIndexInMapper extends Mapper<LongWritable, Text, Text, Text> {
2     private Map<String, Set<String>> wordToFiles;
3     private String filename;
4
5     @Override
6     protected void setup(Context context) throws IOException, InterruptedException {
7         wordToFiles = new HashMap<>();
8         FileSplit fileSplit = (FileSplit) context.getInputSplit();
9         filename = fileSplit.getPath().getName();
10    }
11
12    @Override
13    protected void map(LongWritable key, Text value, Context context) throws
14        IOException, InterruptedException {
15        StringTokenizer tokenizer = new StringTokenizer(value.toString());
16        while (tokenizer.hasMoreTokens()) {
17            String cleaned = TextUtils.cleanToken(tokenizer.nextToken());
18            if (!cleaned.isEmpty()) {
19                wordToFiles.computeIfAbsent(cleaned, k → new
20                    HashSet<>().add(filename);
21            }
22        }
23    }
24
25    @Override
26    protected void cleanup(Context context) throws IOException, InterruptedException
27    {
28        Text word = new Text();
29        Text file = new Text();
30        for (Map.Entry<String, Set<String>> entry : wordToFiles.entrySet()) {
31            word.set(entry.getKey());

```

```

29         for (String fname : entry.getValue()) {
30             file.set(fname);
31             context.write(word, file);
32         }
33     }
34 }
35 }

```

Listing 1.4: Inverted Index In Mapper

1.3.4 Python Non-Parallel Solution

```

1 def buildinvertedindex(directory):
2     invertedindex = defaultdict(set)
3
4     for filename in os.listdir(directory):
5         filepath = os.path.join(directory, filename)
6         if os.path.isfile(filepath):
7             try:
8                 with open(filepath, 'r', encoding='utf-8', errors='ignore') as file:
9                     content = file.read().lower()
10                    words = re.findall(r'"b"w+"b', content)
11                    for word in words:
12                        invertedindex[word].add(filename)
13            except Exception as e:
14                print(f"Could not read {filename}: {e}")
15    return invertedindex

```

Listing 1.5: Inverted Index Python

1.3.5 Python Parallel Solution

```

1 def processfile(args):
2     """Helper for parallel processing"""
3     filepath, basedir = args
4     text = safeextract(filepath)
5     words = re.findall(r'"b"w+"b', text) if text else []
6     return [(word, os.path.relpath(filepath, basedir)) for word in words]
7
8
9 def buildparallelinvertedindex(directory, workers=4):
10    """Faster processing using multiprocessing"""

```

```

11     filepaths = []
12     for root, , files in os.walk(directory):
13         filepaths.extend(os.path.join(root, f) for f in files if
14                             f.lower().endswith(('.txt', '.pdf', '.docx', '.html')))
15
16     with Pool(workers) as pool:
17
18         results = pool.map(processfile, [(fp, directory) for fp in filepaths])
19
20     invertedindex = defaultdict(set)
21     for wordtuples in results:
22         for word, path in wordtuples:
23             invertedindex[word].add(path)
24
25     return invertedindex

```

Listing 1.6: Parallel Inverted Index Python

1.4 Test

1.5 Performance Evaluation

1.6 Conclusion

1.7 Testing Methodology

1.7.1 Test Environment

The testing environment was containerized using Docker to ensure reproducibility across different systems. The container configuration included:

- **Base Image:** Ubuntu 22.04 LTS
- **Hadoop Version:** 3.3.6 (pseudo-distributed mode)
- **Java:** OpenJDK 11
- **Python:** 3.10 with required dependencies

```
1 version: '3'
2 services:
3   hadoop-test:
4     image: hadoop-base
5     build:
6       context: .
7       dockerfile: Dockerfile
8     ports:
9       - "9870:9870" # NameNode
10      - "9864:9864" # DataNode
11      - "8088:8088" # ResourceManager
12     volumes:
13       - ./data:/data
14       - ./output:/output
```

Listing 1.7: Docker Compose Configuration

1.7.2 Test Cases

Functional Validation

Table 1.1: Functional Test Cases

ID	Description	Input	Expected Result
TC-01	Single document processing	1 KB text file	Correct word-file mappings
TC-02	Multi-document processing	10 files (1 MB total)	Combined index with no duplicates
TC-03	Special character handling	File with punctuation	Properly normalized terms
TC-04	Empty file handling	0 KB file	Skipped with warning
TC-05	Large file stress test	1 GB Wikipedia dump	Successful completion

Performance Benchmarking

Tests were executed with varying dataset sizes:

```
1 # Hadoop Version
2 time hadoop jar InvertedIndex.jar /input /output
3
4 # Python Version
5 python invertedindex.py --input /data --output /results
```

Listing 1.8: Test Execution Command

Execution Time Comparison (Hadoop vs Python)

1.7.3 Validation Methodology

Docker-based Testing Pipeline:

1. Container initialization with test datasets mounted
2. Automated test execution via Makefile:

```
1 test-hadoop:
2     docker-compose run hadoop-test "
3         hadoop jar /code/InvertedIndex.jar /test-input /test-output
4
5 test-python:
6     docker-compose run hadoop-test "
7         python /code/invertedindex.py --input /test-input
```

Listing 1.9: Makefile Test Targets

3. Output verification using automated scripts:

Listing 1.10: Verification Script

```
def verify_index(index):  
    assert 'search' in index, "Common term missing"  
    assert len(index['the']) > 0, "Stopword not processed"
```

1.7.4 Results

The Docker environment successfully validated:

- Consistent execution across 5 test runs ($\pm 2\%$ variation)
- Correct index generation for all test cases
- Hadoop outperformed Python by 3.7x on 1GB dataset
- Resource usage remained within container limits (8GB RAM allocated)

Table 1.2: Resource Utilization

Dataset	CPU Usage (%)	RAM (GB)	Time (s)
1 MB	12	1.2	4.7
100 MB	68	3.8	28.1
1 GB	92	7.5	193.4