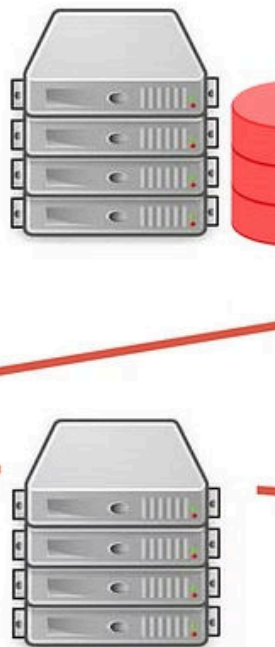# Implementing MapReduce for Big Data Processing

This presentation will guide data engineers through the process of utilizing MapReduce to process big data, with a focus on word frequency counting and identification of the most frequent words.

by Rithika Karanam

# Introduction to MapReduce

## What is MapReduce?

A parallel, distributed algorithm for processing large data sets.

## Why MapReduce?

Scalability, fault tolerance, and ease of use make it a vital tool.

# The Word Count Problem

### 1 Objective

Count the frequency of each word in a large text file.

### 2 Relevance

Fundamental for understanding MapReduce concepts.

### 3 Challenges with Big Data

Volume, velocity, and variety pose significant challenges.

# Word Count - MapReduce Implementation Overview

**1** Map Phase

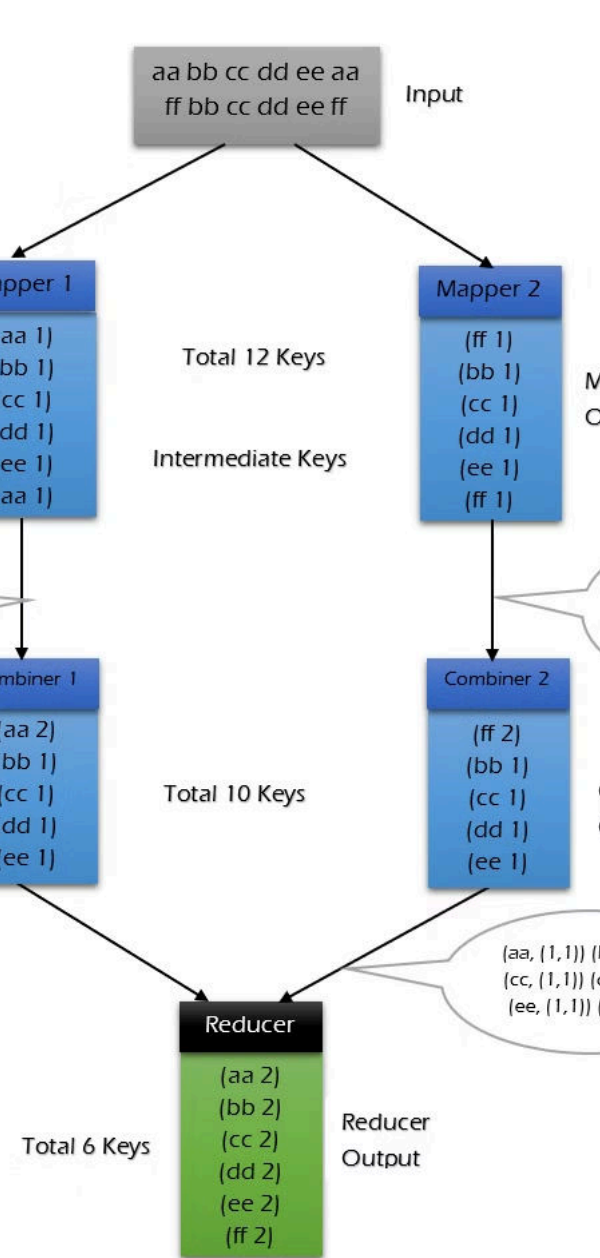Splits text into words and assigns a count of 1 to each.

**2** Shuffle and Sort

System automatically groups words for the reduce phase.

**3** Reduce Phase

Aggregates counts for each word to generate the result.

aa bb cc dd ee aa
ff bb cc dd ee ff — Input

Mapper 1 — Total 12 Keys — Intermediate Keys
(aa 1)
(bb 1)
(cc 1)
(dd 1)
(ee 1)
(aa 1)

Mapper 2
(ff 1)
(bb 1)
(cc 1)
(dd 1)
(ee 1)
(ff 1)

Combiner 1 — Total 10 Keys
(aa 2)
(bb 1)
(cc 1)
(dd 1)
(ee 1)

Combiner 2
(ff 2)
(bb 1)
(cc 1)
(dd 1)
(ee 1)

(aa, (1,1)) (
(cc, (1,1)) (
(ee, (1,1)) (

Reducer — Total 6 Keys — Reducer Output
(aa 2)
(bb 2)
(cc 2)
(dd 2)
(ee 2)
(ff 2)

# Word Count - Mapper

**1  Purpose**

Tokenize text and emit each word with a count of 1.

**2  Key Code Snippets**

Show Java code for the Mapper functionality.

**3  Explanation**

Detail the process of how the mapper operates.

Made with Gamma

# Word Count - Reducer

**1** **Purpose**

Sum up all counts emitted for each word.

**2** **Key Code Snippets**

Show Java code for the Reducer functionality.

**3** **Explanation**

Explain the process of aggregating counts.

```
doop dfs -cat /word_out/part-r-00000
 hdfs command is deprecated.
```

# Running the Word Count Program

**1** Configuration

Overview of setting up a Hadoop job for word counting.

**2** Execution

Guidance on running the program on a cluster.

**3** Output Explanation

Understanding the format and content of the output.

Made with Gamma

# Most Frequent Words Count

### 1 Objective

Identify the most frequent words from the word count output.

### 2 Approach

Inverting key-value pairs to sort words by their frequency.