

API Best Practices and MapReduce

The background of the slide features a photograph of the Alma Mater statue at the University of Illinois, which is a central figure in a long, flowing gown with her arms outstretched. The entire image is overlaid with a semi-transparent red color, creating a monochromatic effect.

CS 240 - The University of Illinois

Wade Fagen-Ulmschneider

October 26, 2021



APIS

API Notes:

API Notes:

Best Web API Design Practices:

1.

2.

3.

4.

5.

Data Stores



Data Stores:

- Optimized for data storage and lookup.
- Focus on queries for data:
 - `SELECT * FROM table WHERE date="2021-11-26";`
 - `find_one({date: {$eq: "2021-11-26"}})`
 - ...
- Developer is responsible for processing the data.

Big Data Processing

A photograph of a crowd of people gathered around a statue of Alma Mater, overlaid with a semi-transparent orange filter. The text "Big Data Processing" is centered in white. The background shows a large group of people, some looking towards the camera and others looking towards the statue. The statue is a large, classical-style figure standing on a pedestal. The overall scene is outdoors, likely at a university or public square.

Big Data Processing:

- Optimized for data processing on large data sets.
- Queries involve functions that perform logic as part of processing the data.
- Key Technology: _____.

MapReduce

A photograph of a crowd of people gathered around a statue of a woman in a long dress, with the text 'MapReduce' overlaid in white. The image is heavily filtered with an orange-red color. The statue is the central focus, and the crowd is visible in the foreground and background. The text 'MapReduce' is prominently displayed in the center of the image.

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

Google, Inc.

Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

Programs written in this functional style are automatically parallelized and executed on a large cluster of commodity machines. The run-time system takes care of the details of partitioning the input data, scheduling the program's execution across a set of machines, handling machine failures, and managing the required inter-machine communication. This allows programmers without any experience with parallel and distributed systems to easily utilize the resources of a large distributed system.

Our implementation of MapReduce runs on a large cluster of commodity machines and is highly scalable: a typical MapReduce computation processes many terabytes of data on thousands of machines. Programmers find the system easy to use: hundreds of MapReduce programs have been implemented and upwards of one thousand MapReduce jobs are executed on Google's clusters every day.

1 Introduction

Over the past five years, the authors and many others at Google have implemented hundreds of special-purpose computations that process large amounts of raw data, such as crawled documents, web request logs, etc., to compute various kinds of derived data, such as inverted indices, various representations of the graph structure of web documents, summaries of the number of pages crawled per host, the set of most frequent queries in a

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with these issues.

As a reaction to this complexity, we designed a new abstraction that allows us to express the simple computations we were trying to perform but hides the messy details of parallelization, fault-tolerance, data distribution and load balancing in a library. Our abstraction is inspired by the *map* and *reduce* primitives present in Lisp and many other functional languages. We realized that most of our computations involved applying a *map* operation to each logical "record" in our input in order to compute a set of intermediate key/value pairs, and then applying a *reduce* operation to all the values that shared the same key, in order to combine the derived data appropriately. Our use of a functional model with user-specified map and reduce operations allows us to parallelize large computations easily and to use re-execution as the primary mechanism for fault tolerance.

The major contributions of this work are a simple and powerful interface that enables automatic parallelization and distribution of large-scale computations, combined with an implementation of this interface that achieves high performance on large clusters of commodity PCs.

Section 2 describes the basic programming model and gives several examples. Section 3 describes an implementation of the MapReduce interface tailored towards our cluster-based computing environment. Section 4 describes several refinements of the programming model that we have found useful. Section 5 has performance measurements of our implementation for a variety of tasks. Section 6 explores the use of MapReduce within Google including our experiences in using it as the basis

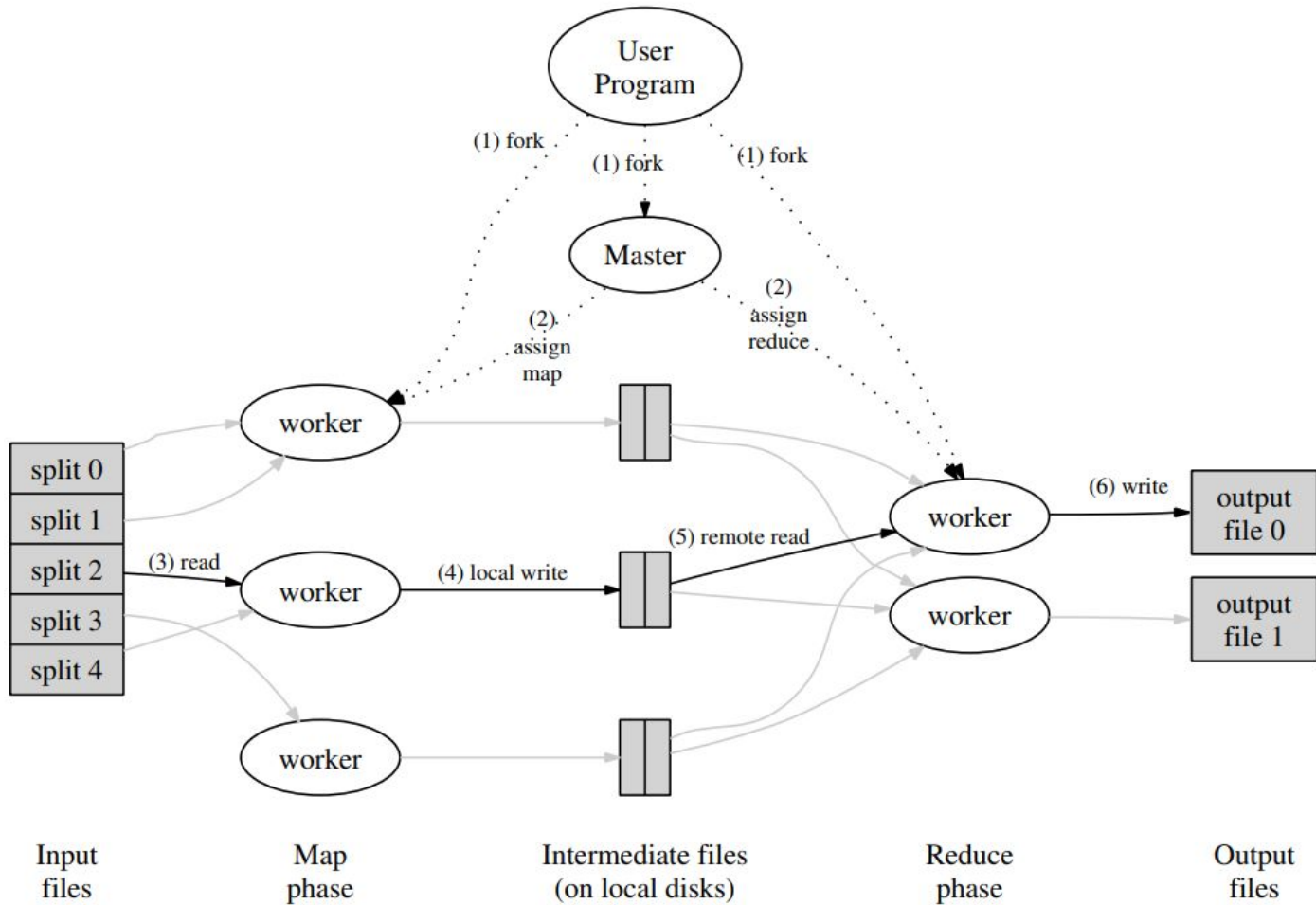


Figure 1, <https://static.googleusercontent.com/media/research.google.com/en//archive/mapreduce-osdi04.pdf>

Map Function

Reduce Function

MapReduce Examples

The image features a central photograph of a crowd of people gathered around a statue of a woman in a long, flowing dress. The statue is mounted on a pedestal with the inscription 'ALMA MATER' visible. The background consists of bare tree branches. The entire image is overlaid with a semi-transparent orange filter. The text 'MapReduce Examples' is written in a large, white, sans-serif font across the middle of the image.

<i>The</i>	<i>quick</i>	<i>brown</i>	<i>fox</i>	<i>jumps</i>	<i>over</i>	<i>the</i>	<i>lazy</i>	<i>dog</i>
[0]	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]

A → **B, C**
B → **A, C, D**
C → **A, B, D**
D → **B, C**

A → **B**

A → **C**

B → **A**

B → **C**

B → **D**

C → **A**

C → **B**

C → **D**

D → **B**

D → **C**