



WIKIPEDIA
The Free Encyclopedia

WIKIPEDIA

MapReduce

MapReduce is a programming model and an associated implementation for processing and generating big data sets with a parallel and distributed algorithm on a cluster.^{[1][2][3]}

A MapReduce program is composed of a *map* procedure, which performs filtering and sorting (such as sorting students by first name into queues, one queue for each name), and a *reduce* method, which performs a summary operation (such as counting the number of students in each queue, yielding name frequencies). The "MapReduce System" (also called "infrastructure" or "framework") orchestrates the processing by marshalling the distributed servers, running the various tasks in parallel, managing all communications and data transfers between the various parts of the system, and providing for redundancy and fault tolerance.

The model is a specialization of the *split-apply-combine* strategy for data analysis.^[4] It is inspired by the map and reduce functions commonly used in functional programming,^[5] although their purpose in the MapReduce framework is not the same as in their original forms.^[6] The key contributions of the MapReduce framework are not the actual map and reduce functions (which, for example, resemble the 1995 Message Passing Interface standard's^[7] *reduce*^[8] and *scatter*^[9] operations), but the scalability and fault-tolerance achieved for a variety of applications due to parallelization. As such, a single-threaded implementation of MapReduce is usually not faster than a traditional (non-MapReduce) implementation; any gains are usually only seen with multi-threaded implementations on multi-processor hardware.^[10] The use of this model is beneficial only when the optimized distributed shuffle operation (which reduces network communication cost) and fault tolerance features of the MapReduce framework come into play. Optimizing the communication cost is essential to a good MapReduce algorithm.^[11]

MapReduce libraries have been written in many programming languages, with different levels of optimization. A popular open-source implementation that has support for distributed shuffles is part of Apache Hadoop. The name MapReduce originally referred to the proprietary Google technology, but has since become a generic trademark. By 2014, Google was no longer using MapReduce as its primary *big data* processing model,^[12] and development on Apache Mahout had moved on to more capable and less disk-oriented mechanisms that incorporated full map and reduce capabilities.^[13]

Overview

MapReduce is a framework for processing parallelizable problems across large datasets using a large number of computers (nodes), collectively referred to as a cluster (if all nodes are on the same local network and use similar hardware) or a grid (if the nodes are shared across geographically and administratively distributed systems, and use more heterogeneous hardware). Processing can occur