

Hadoop MapReduce "from scratch"

2019 Nov.



Presentation of the Project

Yixiao FEI

2019 Nov.

目录 Contents

- 1 Hadoop
- 2 MapReduce
- 3 Project

目录 Contents

- 1 Hadoop
- 2 MapReduce
- 3 Project



Big Data







Hadoop



• Apache Hadoop is a collection of open-source software utilities that facilitate using a network of many computers to solve problems involving massive amounts of data and computation. It provides a software framework for distributed storage and processing of big data using the MapReduce programming model.



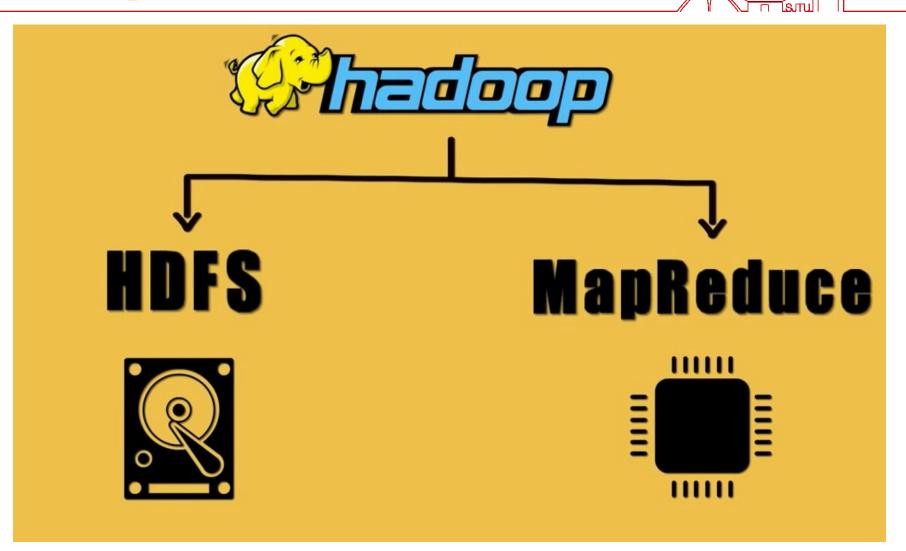


Hadoop





Hadoop

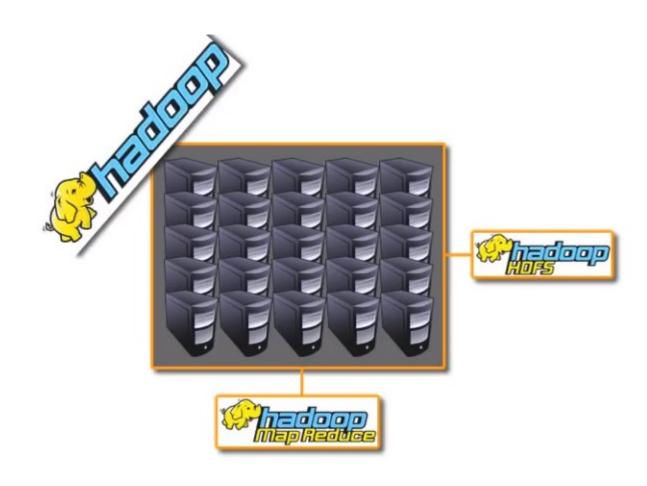


- 1 Hadoop
- 2 MapReduce
- 3 Project



Hadoop Cluster



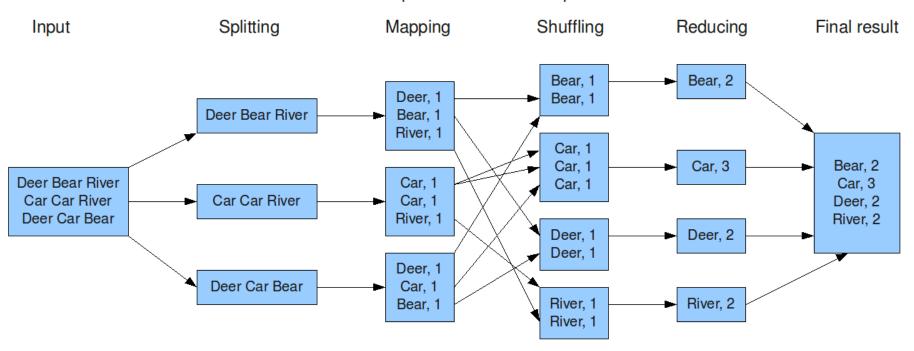




MapReduce



The overall MapReduce word count process



目录 Contents

- 1 Hadoop
- 2 MapReduce
- 3 Project



Word Count

- Trivial Method
- MapReduce
- Experiments



Important Notes

- Work: 3 people/group
- Deadline: 12:00, 24 Dec. 2019
- Filename: MP-GroupNumber.zip (eg. IM-Group1.zip) including:
 - 1) source code
 - 2) compiled runnable Jar and a Readme file
 - 3) report (your names and IDs)
- TA's contact: destiny_fyx@outlook.com



Evaluation



- Functionality (50%), Report (20%), Presentation(30%)
- Delay of submission will cost 30% of points per day
- The report includes at least:
 - ✓ System diagram
 - ✓ Function description
 - ✓ How to build/run (some demo)
 - ✓ Experiment (comparison and reflection)
 - ✓ Good written ability

No Tolerance for Plagiarism



Things that might Help

- Eclipse
- Maven
- Docker
- Linux Commands
- Git

Thanks!

