# 实验材料

# MapReduce / Hadoop



大数据系统基础（B）(200)(2015-2016秋季学期)

新增文件

| 电子教案 | 实验材料 | 分类管理 |

点击带下划线的列标题，可进行排序

| 序号 | 标题 | 简要说明 | 文件大小 | 下载次数 | 上载时间 | 文件管理 |
|---|---|---|---|---|---|---|
| 1 | 实验指导书一 _mapreduce _并行编程 | 实验指导书一_mapreduce_并行编程 | 1.45M | 0 | 2015-09-26 | 修改 删除 |
| 2 | 实验一的实验数据 | | 2.91M | 0 | 2015-09-26 | 修改 删除 |
| 3 | Scala编程 | Scala是基于Java的语言，可以调用Java的库，但代码异常简洁。后面的Spark平台可以使用Scala这门函数式的编程语言，有兴趣的童鞋可以下下来看看。 | 43.64M | 0 | 2015-09-26 | 修改 删除 |

第一次实验课安排：第五周（B205）

Zhi Wang
wangzhi@sz.tsinghua.edu.cn

# Outline

- MapReduce Basis

- MapReduce Programming Model

- Algorithms in MapReduce

# Motivation: Large-scale Data Processing

- Many tasks: Process lots of data to produce other data

- Want to use hundreds or thousands of CPUs

  - but this needs to be easy

# Typical Large-Data Problems

- Iterate over a large number of records

- Extract something of interest from each

- Shuffle and sort intermediate results

- Aggregate intermediate results

- Generate final output

Key idea: provide a functional abstraction for these two operations

# What is MapReduce?

- MapReduce is a programming model Google has used successfully is processing its "big-data" sets (~ 20000 peta bytes per day)

- Users specify the computation in terms of a map and a reduce function,

- Underlying runtime system automatically parallelizes the computation across large-scale clusters of machines, and

- Underlying system also handles machine failures, efficient communications, and performance issues.

# How MapReduce is Structured

- Functional programming meets distributed computing

- A batch data processing system

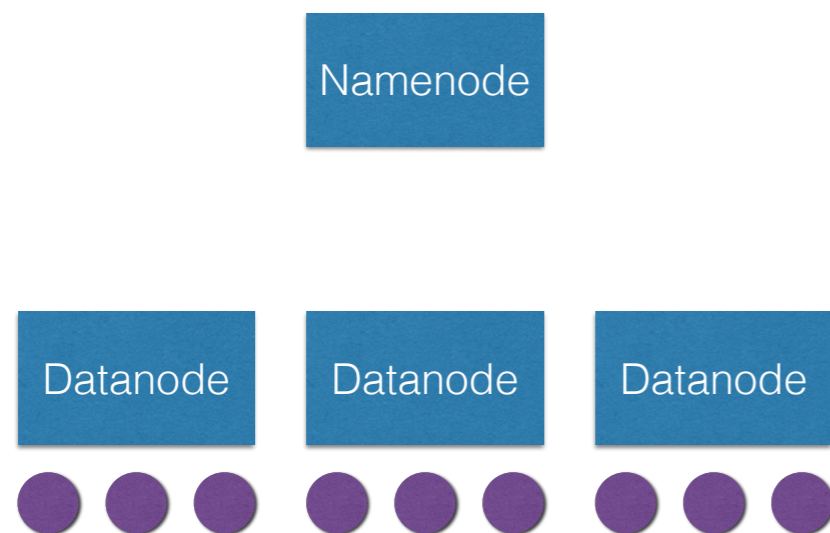- Factors out many reliability concerns from application logic

# MapReduce Provides

- Automatic parallelization & distribution

- Fault-tolerance

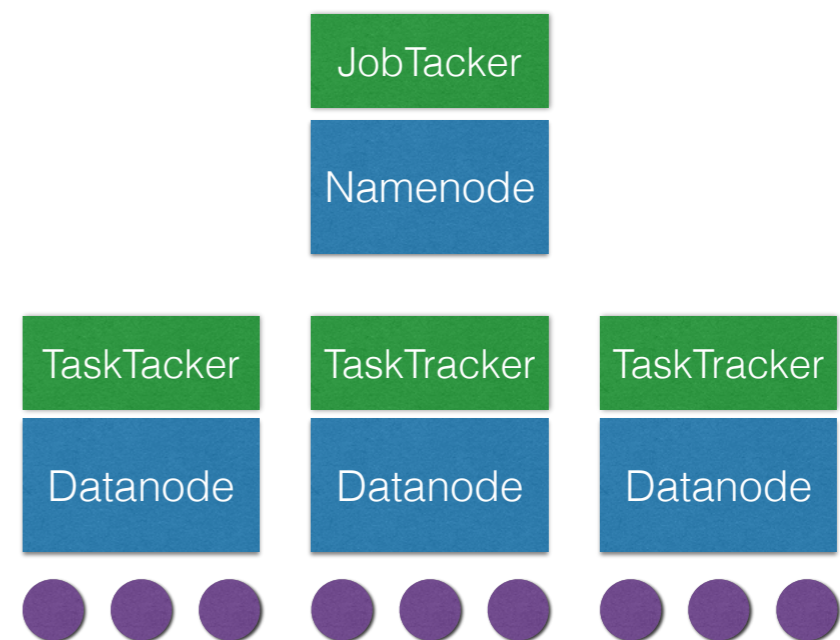- Status and monitoring tools

- A clean abstraction for programmers

# Architecture

# Architecture

# MapReduce: the complete picture



- client — forks → master
- master: *Assign tasks*
- Shard 0, Shard 1, Shard 2, Shard 3, …, Shard M-1
- Map worker → IF (×3)
- *M work items*
- Reduce worker → Output file 1
- Reduce worker → Output file 2
- *R work items*

# Step 1: Split input files into chunks (shards)

- Break up the input data into $M$ pieces (typically 64 MB)



| Shard 0 | Shard 1 | Shard 2 | Shard 3 | … | Shard M-1 |
|---------|---------|---------|---------|---|-----------|

Input files

Divided into $M$ shards

# Step 2: Fork processes

- Start up many copies of the program on a cluster of machines
  - 1 master: scheduler & coordinator
  - Lots of workers
- Idle workers are assigned either:
  - map tasks (each works on a shard) – there are $M$ map tasks
  - reduce tasks (each works on intermediate files) – there are $R$
    - $R$ = # partitions, defined by the user



- User program → master, worker, worker, worker … (*Remote fork*)

# Step 3: Map Task

- Reads contents of the input shard assigned to it
- Parses key/value pairs out of the input data
- Passes each pair to a user-defined *map* function
  - Produces intermediate key/value pairs
  - These are buffered in memory



- Shard 2 → *read* → Map worker

# Step 4: Create intermediate files

- Intermediate key/value pairs produced by the user's *map* function buffered in memory and are periodically written to the local disk
  - Partitioned into *R* regions by a partitioning function



Shard 2 → *read* → Map worker → *local write* → Intermediate file
- Partition 1
- Partition 1
- Partition *R-1*

# Step 4a. Partitioning

- Map data will be processed by Reduce workers
  - The user's *Reduce* function will be called once per unique key generated by *Map*.

- This means we will need to sort all the (key, value) data by keys and decide which Reduce worker processes which keys – the Reduce worker will do this

- Partition function: decides which of *R* reduce workers will work on which key
  - Default function: *hash(key) mod R*
  - Map worker partitions the data by keys

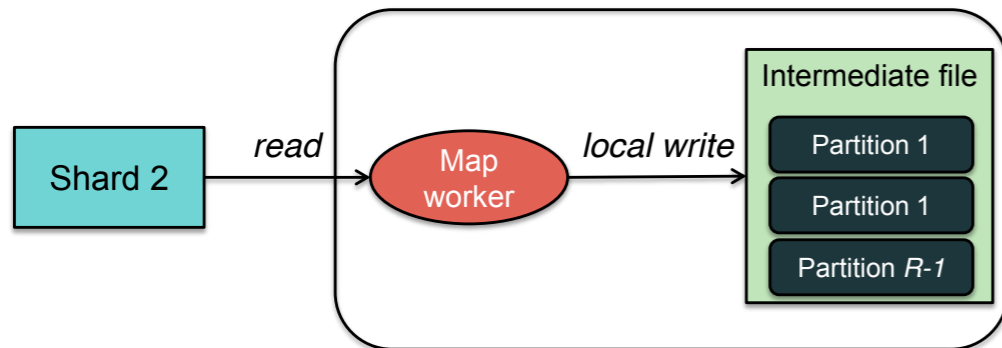- Each Reduce worker will read their partition from every Map worker

# Step 5: Reduce Task: sorting

- Reduce worker gets notified by the master about the location of intermediate files for its partition
- Uses RPCs to read the data from the local disks of the map workers
- When the *reduce* worker reads intermediate data for its partition
  - It sorts the data by the intermediate keys
  - All occurrences of the same key are grouped together



Map worker → *local write* → Intermediate file → *remote read* → Reduce worker

Map worker → *local write* → Intermediate file → *remote read* → Reduce worker

# Step 6: Reduce Task: *Reduce*

- The sort phase grouped data with a unique intermediate key

- User's *Reduce* function is given the key and the set of intermediate values for that key
  - < key, (value1, value2, value3, value4, …) >

- The output of the *Reduce* function is appended to an output file



Intermediate file / Intermediate file / Intermediate file → *remote read* → Reduce worker → *write* → Output file

## Step 7: Return to user

- When all *map* and *reduce* tasks have completed, the master wakes up the user program

- The *MapReduce* call in the user program returns and the program can resume execution.
  - Output of *MapReduce* is available in *R* output files

## MapReduce: the complete picture



# Paradigm



| | Input | Output |
|---|---|---|
| **Map** | <k1, v1> | list (<k2, v2>) |
| **Reduce** | <k2, list(v2)> | list (<k3, v3>) |

## Example



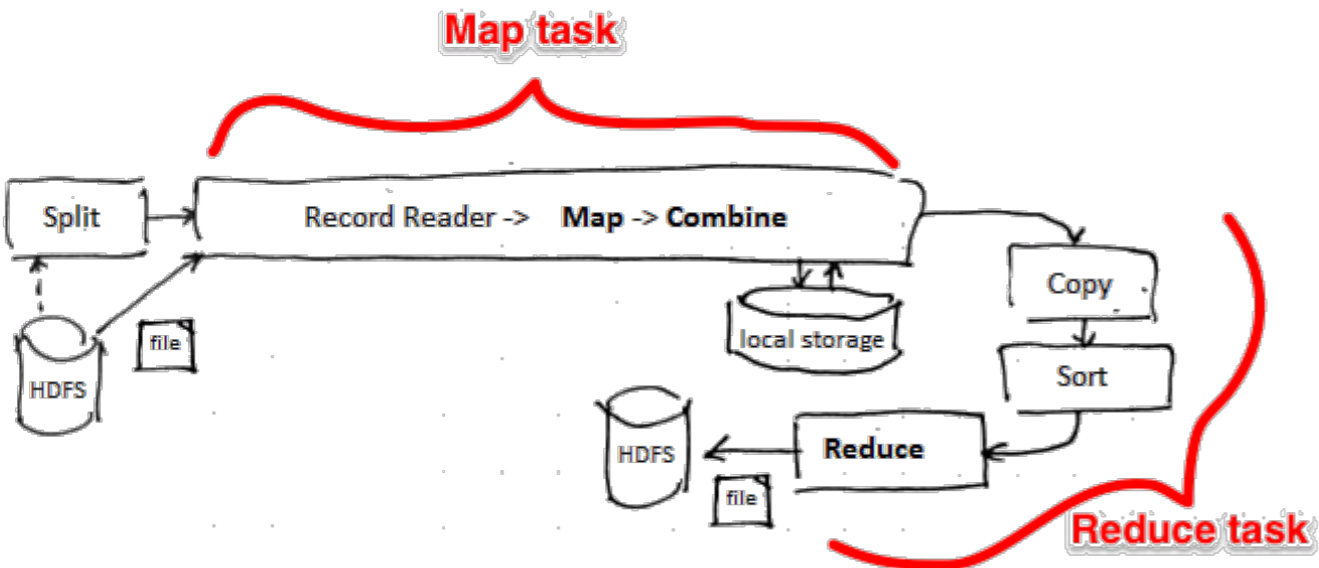| It will be seen that this mere painstaking burrower and grub-worm of a poor devil of a Sub-Sub appears to have gone through the long Vaticans and street-stalls of the earth, picking up whatever random allusions to whales he could anyways find in any book whatsoever, sacred or profane. Therefore you must not, in every case at least, take the higgledy-piggledy whale statements, however authentic, in these extracts, for veritable gospel cetology. Far from it. As touching the ancient authors generally, as well as the poets here appearing, these extracts are solely valuable or entertaining, as affording a glancing bird's eye view of what has been promiscuously said, thought, fancied, and sung of Leviathan, by many nations and generations, including our own. | it 1<br>will 1<br>be 1<br>seen 1<br>that 1<br>this 1<br>mere 1<br>painstaking 1<br>burrower 1<br>and 1<br>grub-worm 1<br>of 1<br>a 1<br>poor 1<br>devil 1<br>of 1<br>a 1<br>sub-sub 1<br>appears 1<br>to 1<br>have 1<br>gone 1 | … <br>a 1<br>a 1<br>aback 1<br>aback 1<br>abaft 1<br>abaft 1<br>abandon 1<br>abandon 1<br>abandon 1<br>abandoned 1<br>abandoned 1<br>abandoned 1<br>abandoned 1<br>abandoned 1<br>abandoned 1<br>abandoned 1<br>abandonedly 1<br>abandonment 1<br>abandonment 1<br>abased 1<br>abased 1 | a 4736<br>aback 2<br>abaft 2<br>abandon 3<br>abandoned 7<br>abandonedly 1<br>abandonment 2<br>abased 2<br>abasement 1<br>abashed 2<br>abate 1<br>abated 3<br>abatement 1<br>abating 2<br>abbreviate 1<br>abbreviation 1<br>abeam 1<br>abed 2<br>abednego 1<br>abel 1<br>abhorred 3<br>abhorrence 1 |

*Z. Wang, Foundations for Big Data Systems 2015*

# MapReduce: whole picture



*Z. Wang, Foundations for Big Data Systems 2015*

# Movie of this week

**The Trueman Show (1998)**



*Z. Wang, Foundations for Big Data Systems 2015*