



# 大数据分析实验

## MapReduce

## 目录

一 任务背景

二 任务描述

1. 任务一

2. 任务二

三 算法流程

四 验收流程

# 一 任务背景

# 一 任务背景

- 1、理解map-reduce算法思想与流程；
- 2、应用map-reduce思想解决问题；
- 3、掌握并应用combine与shuffle过程。



## 二 任务描述



## 二 任务描述



### ◆ 任务一

- 实验数据：提供9个预处理过的源文件（source01-09）模拟9个分布式节点，每个源文件中包含大量的由英文、数字和字符（不包括逗号）构成的单词，单词由逗号与换行符分割。
- 应用map-reduce思想，模拟9个map节点与3个reduce节点实现wordCount功能。
- 输出：对应的map文件和最终的reduce结果文件。
- 要求：由于源文件较大，使用多线程来模拟分布式节点。

## 二 任务描述



### ◆ 任务二

- 掌握并应用combine与shuffle。
- Shuffle过程：map节点通过shuffle过程将任务大致均分给reduce节点。
- Combine过程：map节点通过combine过程压缩输出内容，减少map节点与reduce节点通信。

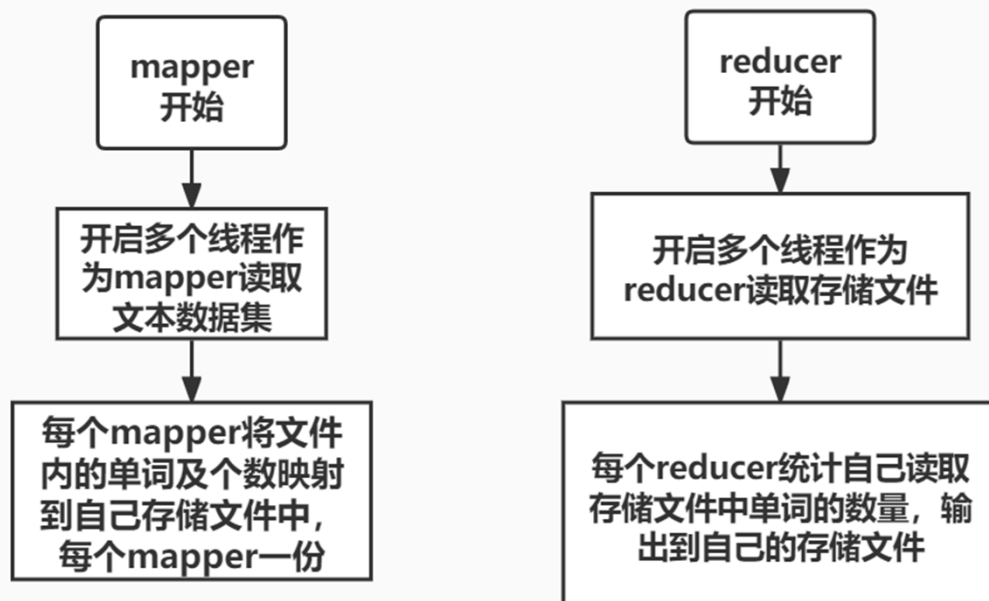
## 三 算法流程



## 三 算法流程

### ◆ MapReduce参考算法流程

MapReduce 参考算法流程



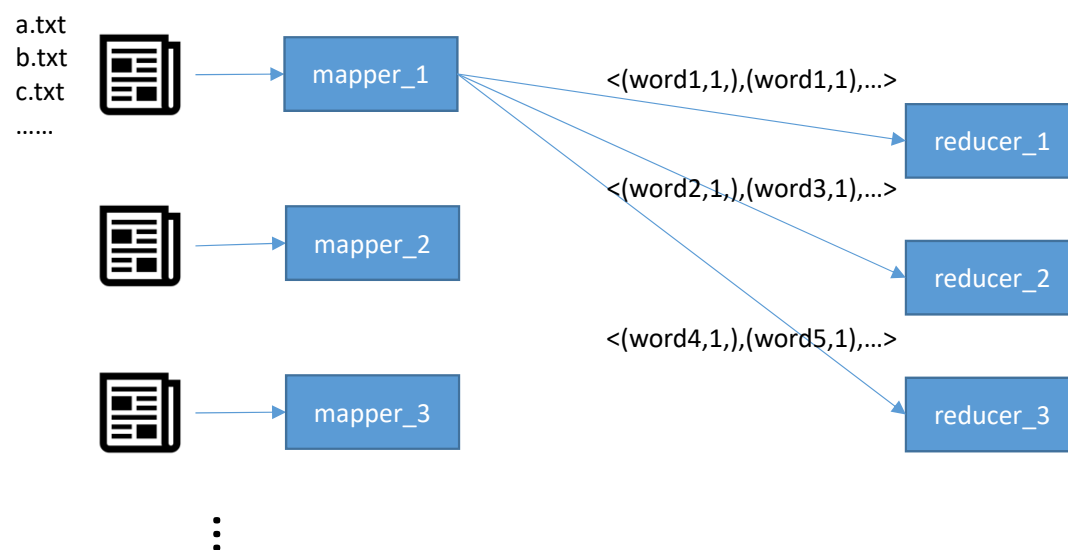
关于mapper结果的分配:  
最简单的版本可以是每3个  
mapper输出的文件作为一个  
reducer的输入文件。

但是, 当每个mapper输入  
的文件数量差距很大时,  
不同reducer的工作量差异  
可能会很大。

## 三 算法流程



### ◆ 进阶：使用shuffle过程。



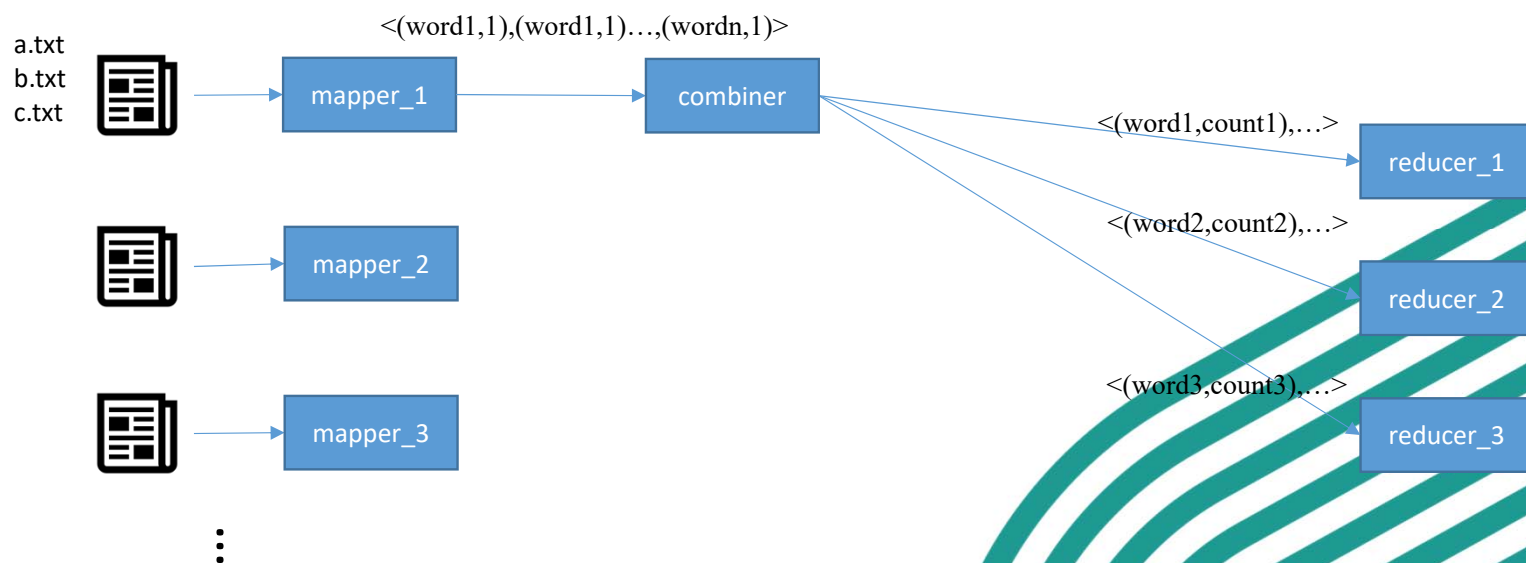
使用shuffle，一个mapper将输出平分为多份，分给多个reducer，这样每个reducer的工作量大致相同。

可以通过hash来将keyword分配到reducer上。

## 三 算法流程

### ◆ 进阶：使用combine过程。

- 注意到在上述方法中，mapper到reducer的传输开销较大，key\_list中可能包含很多重复的关键字，每个mapper可以通过combiner来压缩传输开销：



## 四 验收流程

## 四 验收流程

- 统计结果是否正确；
- 验收时对代码的大致解释；
- 验收时的提问与回答。

