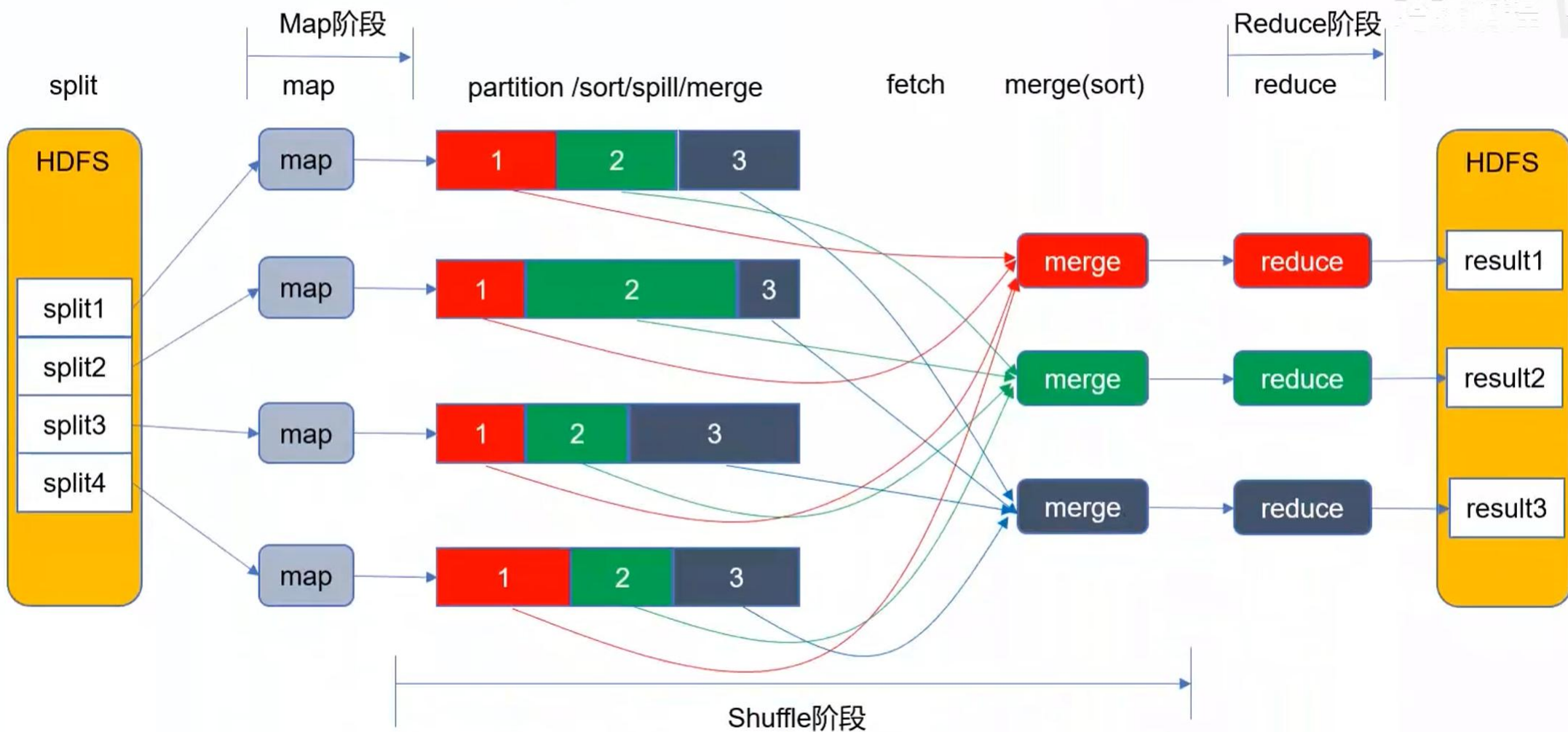# SPARK和MapReduce的性能对比实验

汇报时间：2024年12月11日

# 理论分析-MAPREDUCE

MapReduce主要适用于批量数据处理，是面向批处理的分布式计算框架。
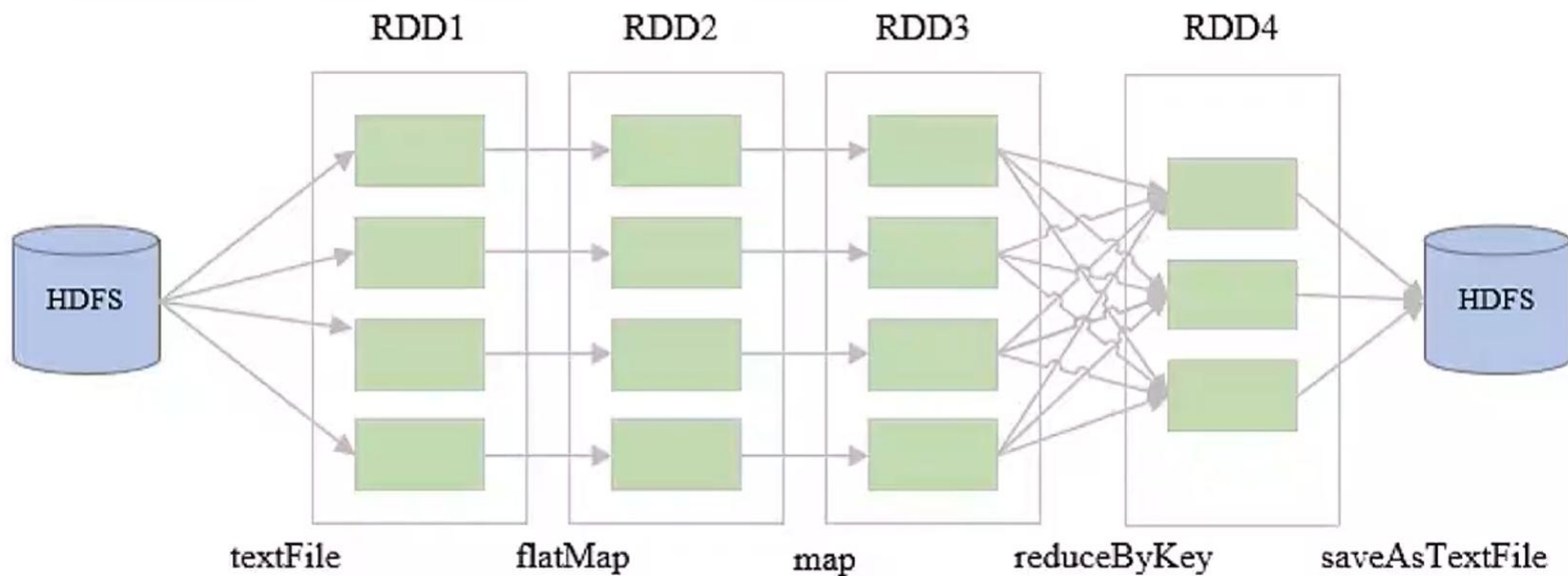MapReduce程序被分为Map（映射）阶段和Reduce（化简）阶段

# 理论分析-SPARK

弹性分布式数据集（RDD）：
- 分布在集群中的只读对象集合。
- 通过转换操作构造。
- 存储在内存或磁盘中。

- 由多个 Partition 组成。
- 失效后自动重构（弹性）。



（Spark 基于 RDD 进行计算）

# 理论分析–SPARK与MAPREDUCE对比

| MapReduce | Spark |
|---|---|
| 数据存储结构: 磁盘 HDFS 文件系统的 split | 使用内存构建弹性分布式数据集 RDD 对数据进行运算和 cache |
| 编程范式: Map + Reduce | DAG: Transformation + Action |
| 计算中间结果c到磁盘,IO及序列化、反序列化代价大 | 计算中间结果在内存中维护，存取速度比磁盘高几个数量级 |
| Task以进程的方式维护，需要数秒时间才能启动任务 | Task以线程的方式维护，对于小数据集读取能够达到亚秒级的延迟 |

# hadoop集群搭建

- 拉取Docker镜像
- 建立使用桥接模式的docker子网
- 启动master、slave1、slave2三个容器作为集群节点

```
C:\Users\ZHI>docker ps -a
CONTAINER ID    IMAGE                                                                           COMMAND     CREATED        STATUS          PORTS
                                                                          NAMES
7c241c635ddc    registry.cn-hangzhou.aliyuncs.com/hadoop_test/hadoop_base_with_spark_ui        "bash"      30 hours ago   Up 8 seconds    0.0.0.0:8080->8080/tcp, 0.0.0
.0:8088->8088/tcp, 0.0.0.0:9870->9870/tcp, 0.0.0.0:10000->10000/tcp    Master
baf8c03e2637    registry.cn-hangzhou.aliyuncs.com/hadoop_test/hadoop_base                      "bash"      2 days ago     Up 7 seconds
                                                                          Slave2
3dd2ef83397f    registry.cn-hangzhou.aliyuncs.com/hadoop_test/hadoop_base                      "bash"      2 days ago     Up 7 seconds
                                                                          Slave1
```

（查看容器状态）

ssh配置，免密登录，修改，增加如下字段/etc/ssh/sshd_config，如下图

```
# override default of no subsystems
Subsystem        sftp    /usr/lib/openssh/sftp-server
PermitRootLogin yes
PasswordAuthentication yes
PubkeyAuthentication yes
```

# hadoop集群搭建

启动Hadoop服务，如下图所示：

```
root@Master:/# start-all.sh
Starting namenodes on [Master]
Master: Warning: Permanently added 'master,172.19.0.2' (ECDSA) to the list of known hosts.
Starting datanodes
Slave2: Warning: Permanently added 'slave2,172.19.0.4' (ECDSA) to the list of known hosts.
Slave1: Warning: Permanently added 'slave1,172.19.0.3' (ECDSA) to the list of known hosts.
Slave2: WARNING: /usr/local/hadoop/logs does not exist. Creating.
Slave1: WARNING: /usr/local/hadoop/logs does not exist. Creating.
Starting secondary namenodes [Master]
Starting resourcemanager
Starting nodemanagers
```

master容器使用/usr/local/hadoop/sbin/start-all.sh命令启动hadoop集群，然后分别在三个容器中通过jps命令查看任务进程，如下三图所示：

```
root@Master:/# jps
305 NameNode
756 SecondaryNameNode
501 DataNode
1093 ResourceManager
1446 NodeManager
1817 Jps
```

```
root@Slave2:/# jps
613 Jps
152 DataNode
303 NodeManager
```

```
root@Slave1:/# jps
611 Jps
152 DataNode
301 NodeManager
```

# SPARK搭建

- 将spark压缩包拷贝到节点
- 解压并修改文件名称
- 配置环境变量
- 配置worker
- 将spark文件拷贝到从节点
- 启动Spark，如下图所示：

```
root@Master:/# /usr/local/spark/sbin/start-all.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/logs/spark--org.apache.spark.deploy.master.Master-1-Master.out
localhost: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-root-org.apache.spark.deploy.worker.Worker-1-Master.out
root@Master:/# jps
305 NameNode
2131 Worker
756 SecondaryNameNode
501 DataNode
1093 ResourceManager
1446 NodeManager
2247 Jps
1913 Master
```

# SPARK搭建

　　master容器使用/usr/local/spark/sbin/start-all.sh命令启动spark集群，slave容器使用/usr/local/spark/sbin/start-worker.sh spark://Master:7077命令将 Slave 容器作为 Worker 加入到 Spark 集群，并连接到 Master 容器上，如下图所示：

```
root@Master:/# /usr/local/spark/sbin/start-all.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/logs/spark--org.apache.spark.deploy.master.M
localhost: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-root-org.apache.spark.
root@Master:/# jps
305 NameNode
2131 Worker
756 SecondaryNameNode
501 DataNode
1093 ResourceManager
1446 NodeManager
2247 Jps
1913 Master
```

```
root@Slave1:/# /usr/local/spark/sbin/start-worker.sh spark://Master:7077
starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark--org.apache.
root@Slave1:/# jps
152 DataNode
810 Jps
698 Worker
301 NodeManager
```

```
root@Slave2:/# /usr/local/spark/sbin/start-worker.sh spark://Master:7077
starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark--org.apache.spark.deploy
root@Slave2:/# jps
152 DataNode
812 Jps
700 Worker
303 NodeManager
```

# MapReduce运算结果

所有 Map 任务的总时间为 21723ms，所有 Reduce 任务的总时间为 10339ms，总计32062ms。

```
2024-12-09 12:49:38,965 INFO mapreduce.Job:  map 0% reduce 0%
2024-12-09 12:49:52,057 INFO mapreduce.Job:  map 100% reduce 0%
2024-12-09 12:50:05,117 INFO mapreduce.Job:  map 100% reduce 100%
2024-12-09 12:50:05,127 INFO mapreduce.Job: Job job_1733744519941_0003 completed successfully
2024-12-09 12:50:05,179 INFO mapreduce.Job: Counters: 54
```

```
        File System Counters
                FILE: Number of bytes read=259739132
                FILE: Number of bytes written=390293828
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=78827139
                HDFS: Number of bytes written=15676
                HDFS: Number of read operations=11
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Launched map tasks=2
                Launched reduce tasks=1
                Data-local map tasks=2
                Total time spent by all maps in occupied slots (ms)=21723
                Total time spent by all reduces in occupied slots (ms)=10339
                Total time spent by all map tasks (ms)=21723
                Total time spent by all reduce tasks (ms)=10339
                Total vcore-milliseconds taken by all map tasks=21723
                Total vcore-milliseconds taken by all reduce tasks=10339
                Total megabyte-milliseconds taken by all map tasks=22244352
                Total megabyte-milliseconds taken by all reduce tasks=10587136
```

# SPARK展示

通过conda activate spark进入python虚拟环境，然后执行任务，spark执行时长约为9021ms。

# SPARK和MapReduce的CPU、内存占用情况对比



（MapReduce的CPU、内存占用情况）

# SPARK和MapReduce的CPU、内存占用情况对比



（SPARK的CPU、内存占用情况）

# 感 谢 指 导

BLUE THESIS PROPOSAL TEMPLATE