

MapReduce和Yarn 技术原理

www.huawei.com





目标

- 学完本课程后，您将能够：
 - 熟悉**MapReduce**和**Yarn**是什么；
 - 掌握**MapReduce**使用的场景及其原理；
 - 掌握**MapReduce**和**Yarn**功能与架构；
 - 熟悉**Yarn**的新特性；



目录

1. MapReduce和Yarn基本介绍
2. MapReduce和Yarn功能与架构
3. Yarn的资源管理和任务调度
4. 增强特性

MapReduce 基本定义

MapReduce是面向大数据并行处理的计算模型、框架和平台。

它包含以下三层含义：

- 1) **MapReduce**是一个基于集群的高性能并行计算平台
(**Cluster Infrastructure**) 。
- 2) **MapReduce**是一个并行计算与运行软件框架
(**Software Framework**) 。
- 3) **MapReduce**是一个并行程序设计模型与方法
(**Programming Model & Methodology**) 。

MapReduce 应用场景

- **MapReduce**基于**Google**发布的分布式计算框架**MapReduce**论文设计开发，用于大规模数据集（大于**1TB**）的并行运算，特点如下：
 - 易于编程：程序员仅需描述做什么，具体怎么做交由系统的执行框架处理。
 - 良好的扩展性：可通过添加节点以扩展集群能力。
 - 高容错性：通过计算迁移或数据迁移等策略提高集群的可用性与容错性。

Yarn 产生背景

MRv1几个方面的缺陷:

- 扩展性受限
- 单点故障
- 不支持**MR**之外的计算

多计算框架之间无法数据共享，资源利用率低

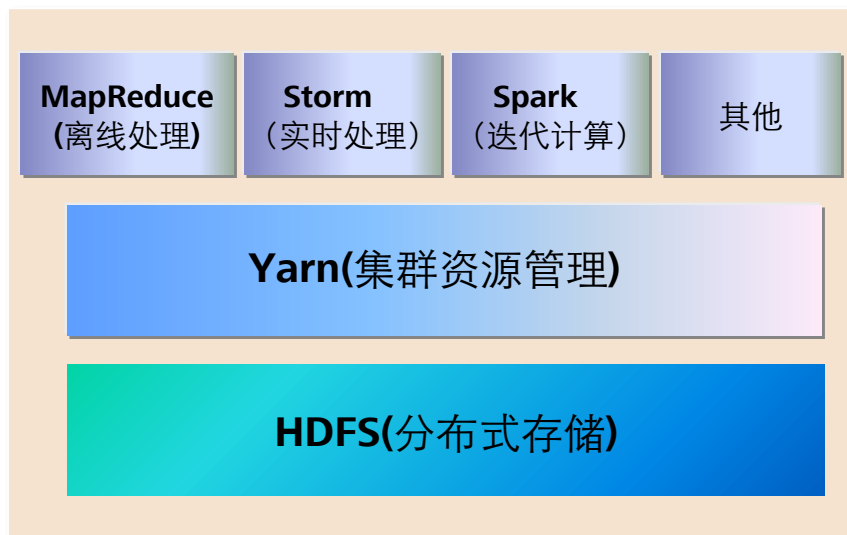
- **MR**: 离线计算框架
- **Storm**: 实时计算框架
- **Spark**: 内存计算框架

Yarn 基本定义

Apache Hadoop YARN (**Yet Another Resource Negotiator**, 另一种资源协调者) 是一种新的 **Hadoop** 资源管理器, 它是一个通用资源管理系统, 可为上层应用提供统一的资源管理和调度, 它的引入为集群在利用率、资源统一管理和数据共享等方面带来了巨大好处。

在产品中定位

- **Yarn**是**Hadoop2.0**中的资源管理系统，它是一个通用的资源管理模块，可为各类应用程序进行资源管理和调度
- **Yarn**是轻量级弹性计算平台，除了**MapReduce**框架，还可以支持其他框架，比如**Spark**、**Storm**等
- 多种框架统一管理，共享集群资源：
 - 资源利用率高
 - 运维成本低
 - 数据共享方便

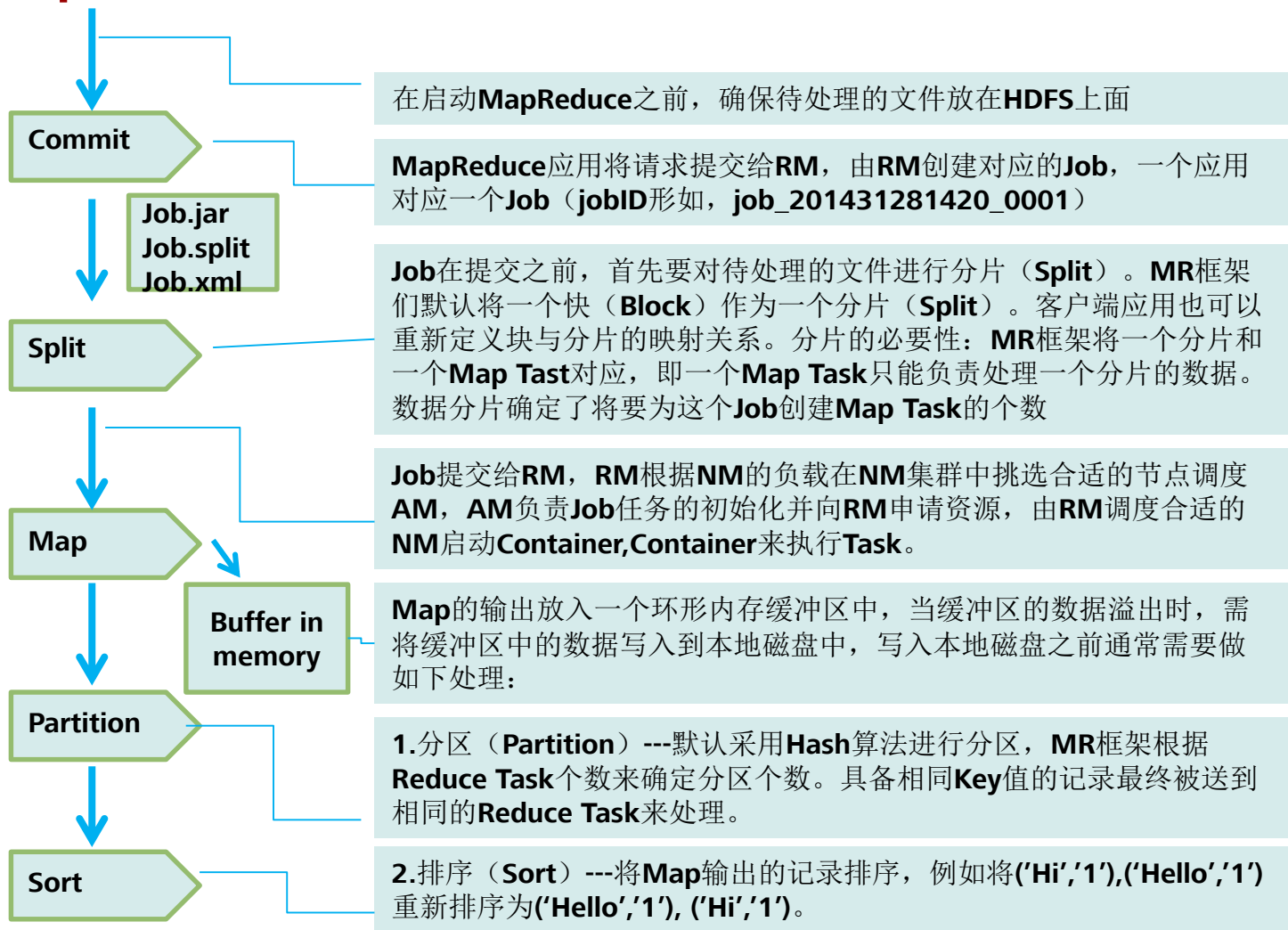




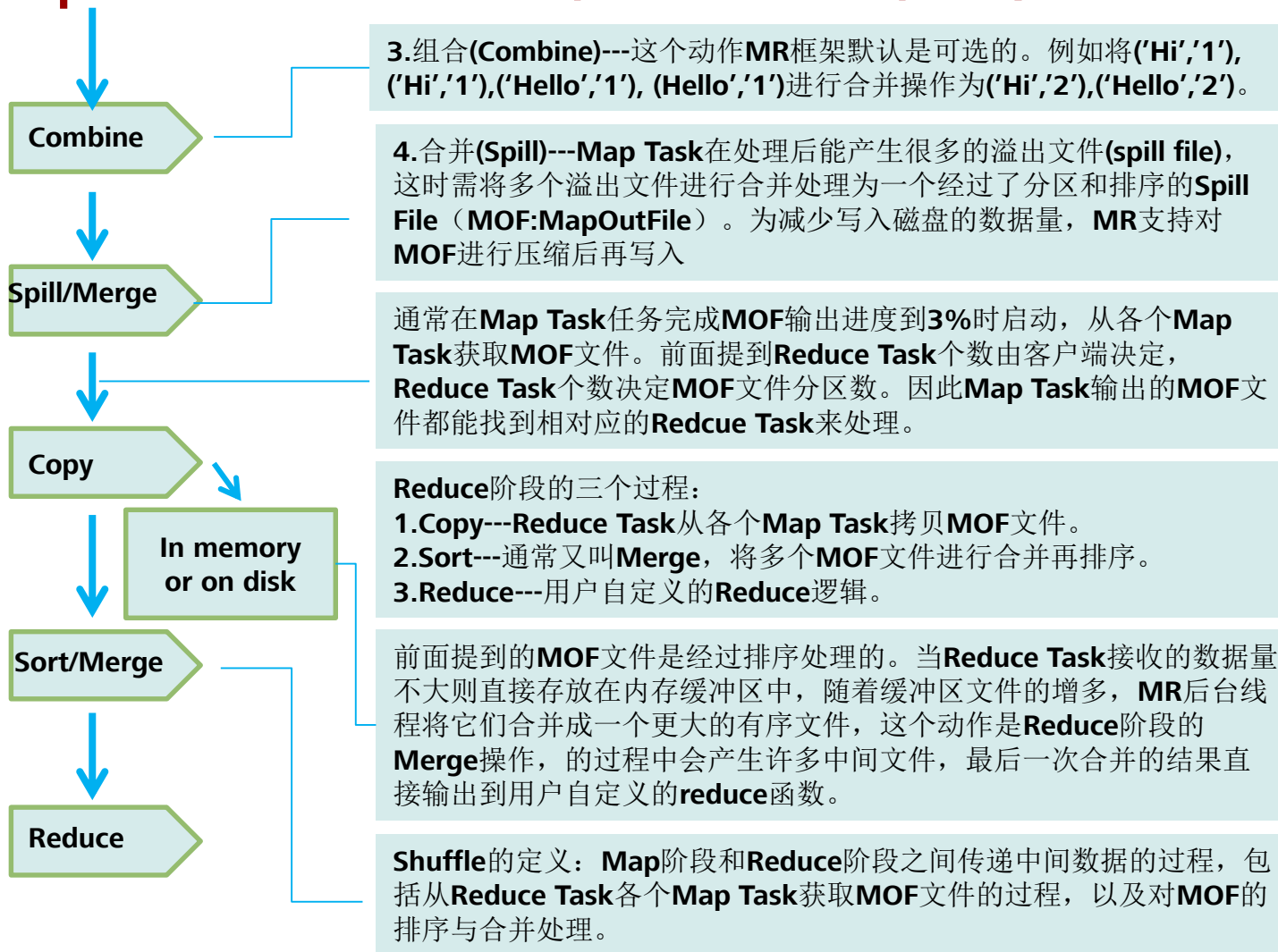
目录

1. MapReduce和Yarn基本介绍
2. MapReduce和Yarn功能与架构
3. Yarn的资源管理和任务调度
4. 增强特性

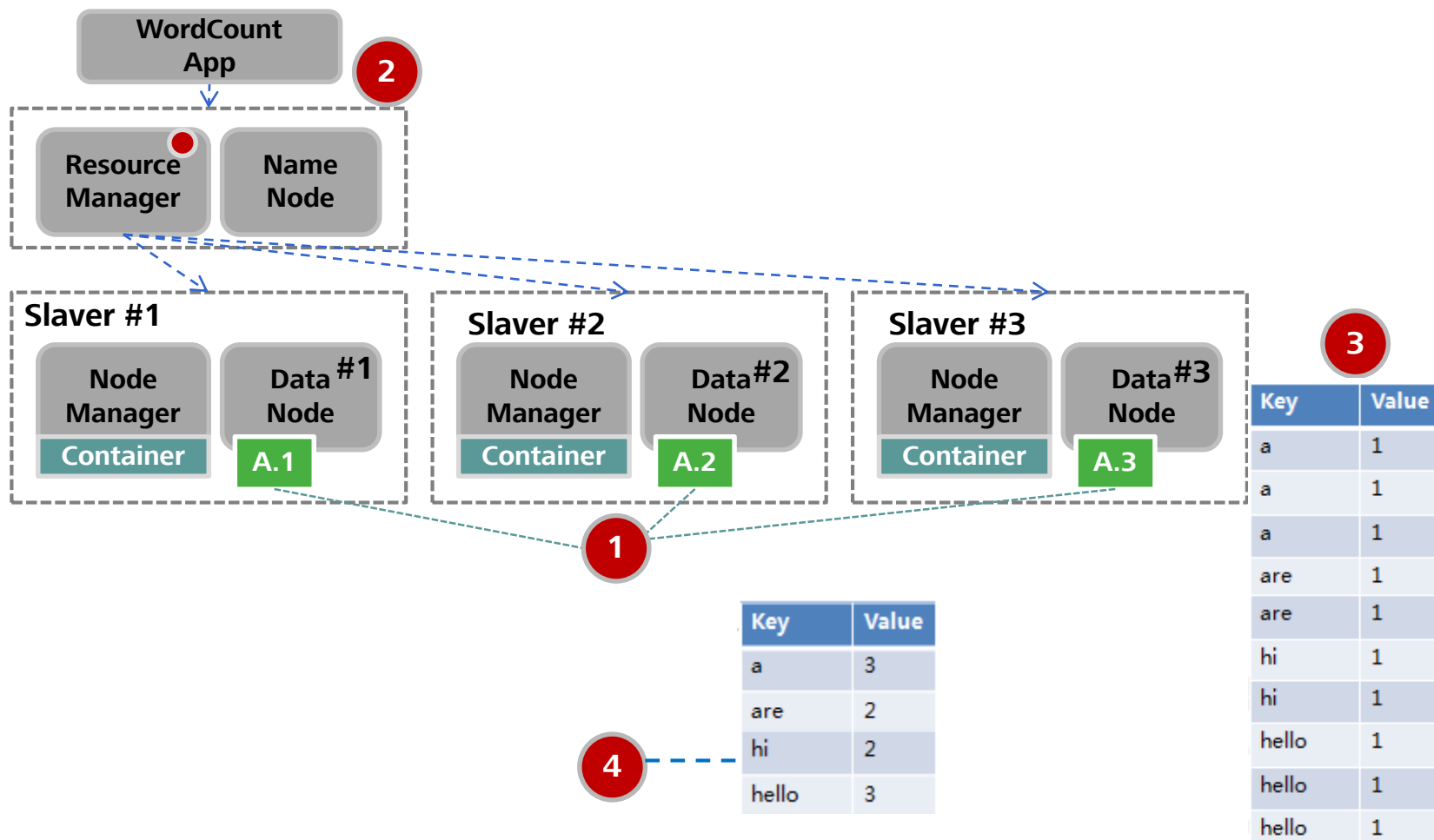
MapReduce的过程-MR过程详解



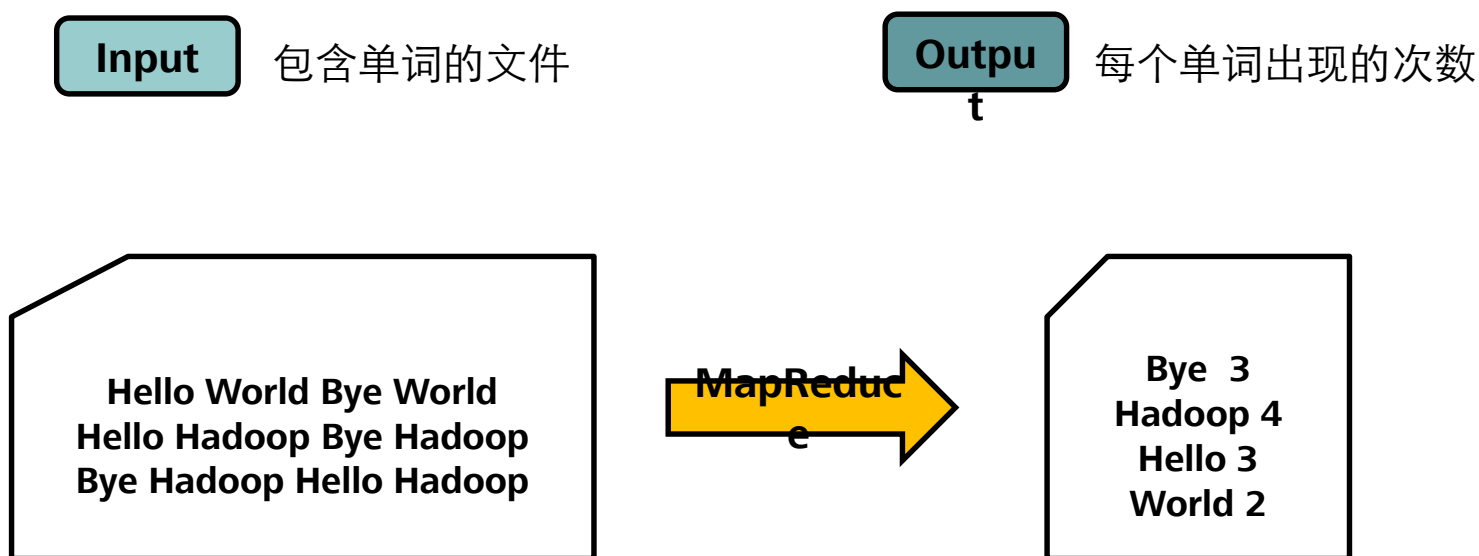
MapReduce的过程-MR过程详解



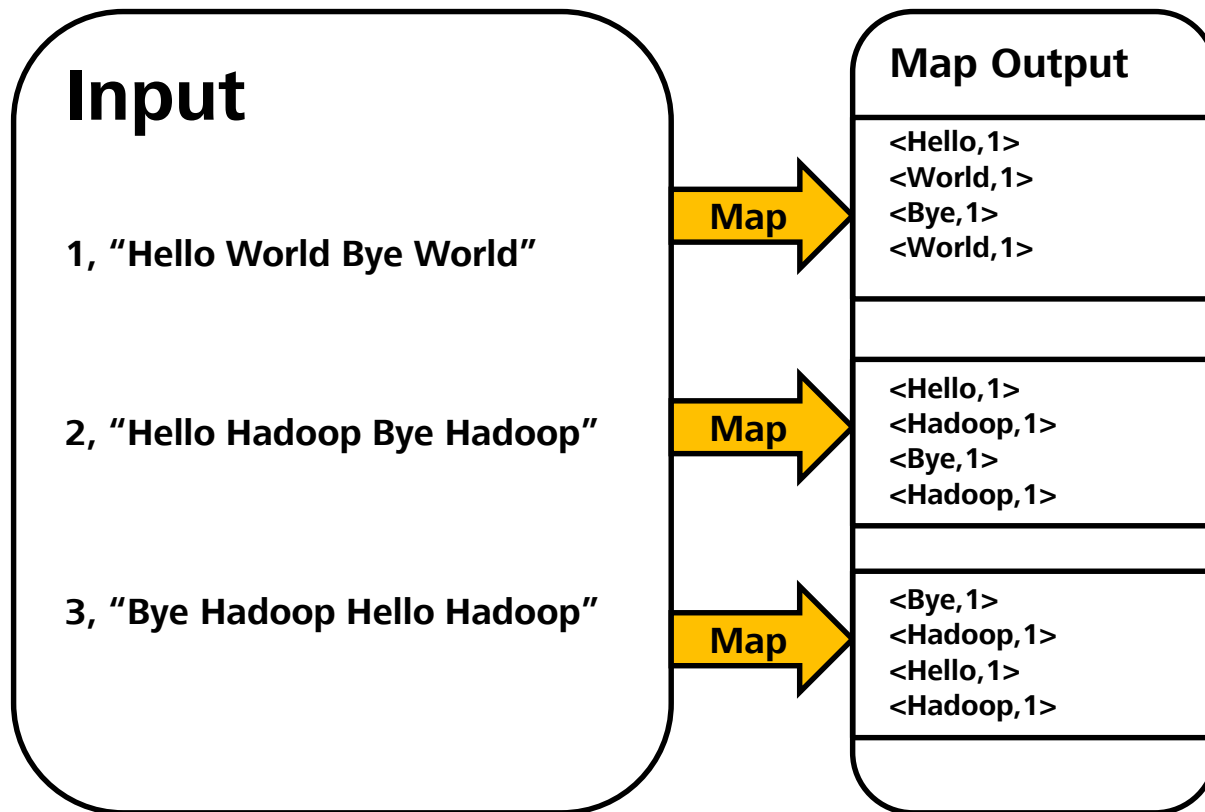
典型程序WordCount举例



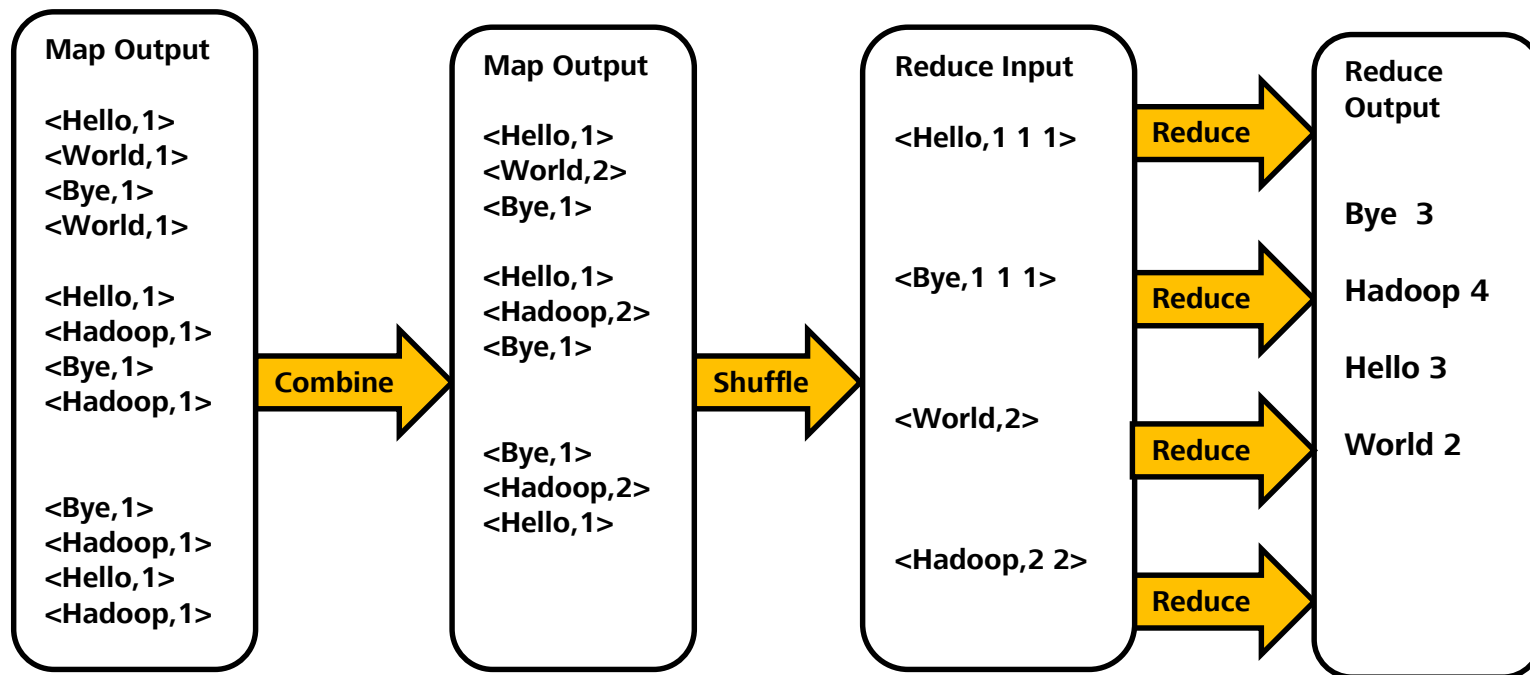
WordCount程序功能



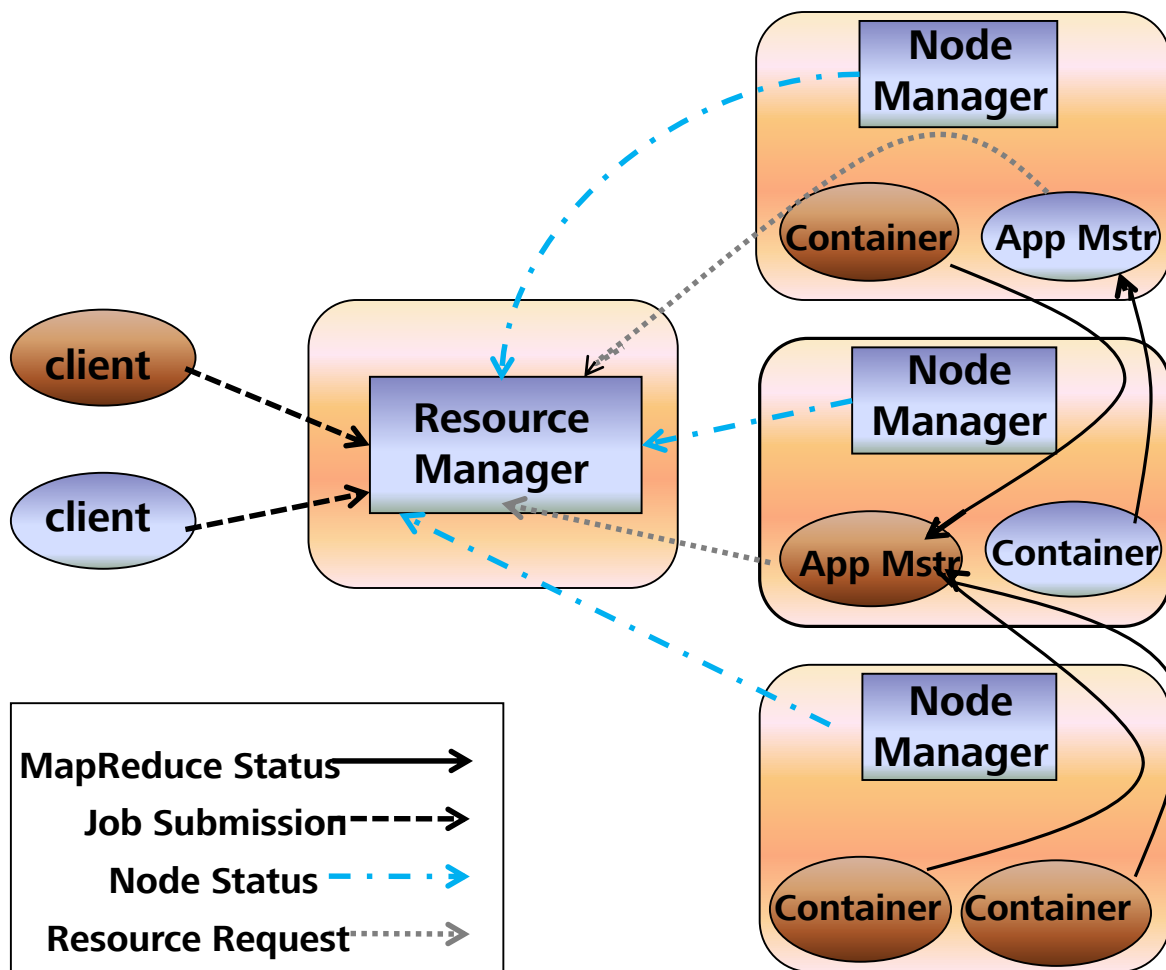
WordCount 的Map过程



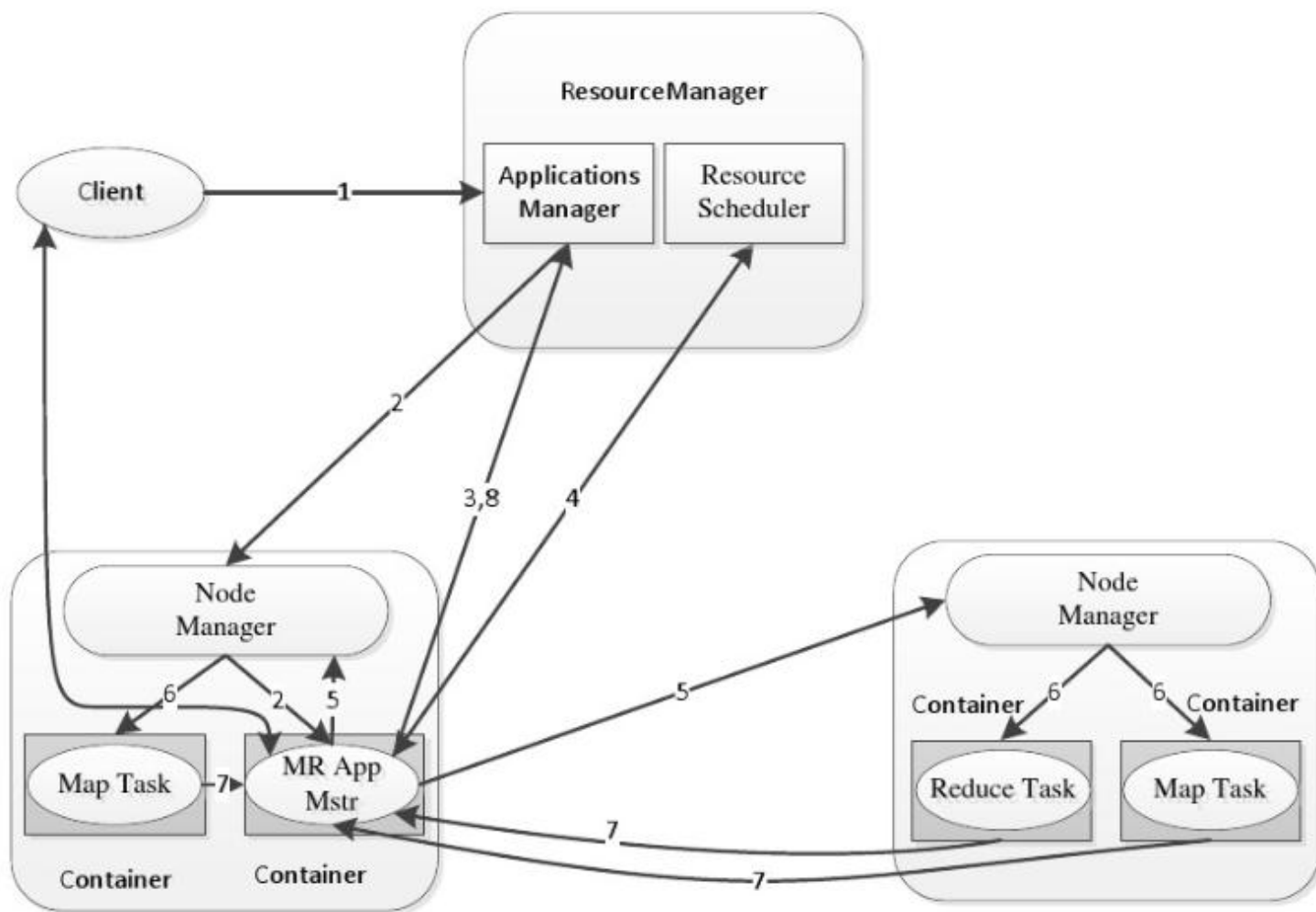
WordCount的Reduce过程



Yarn的组件架构



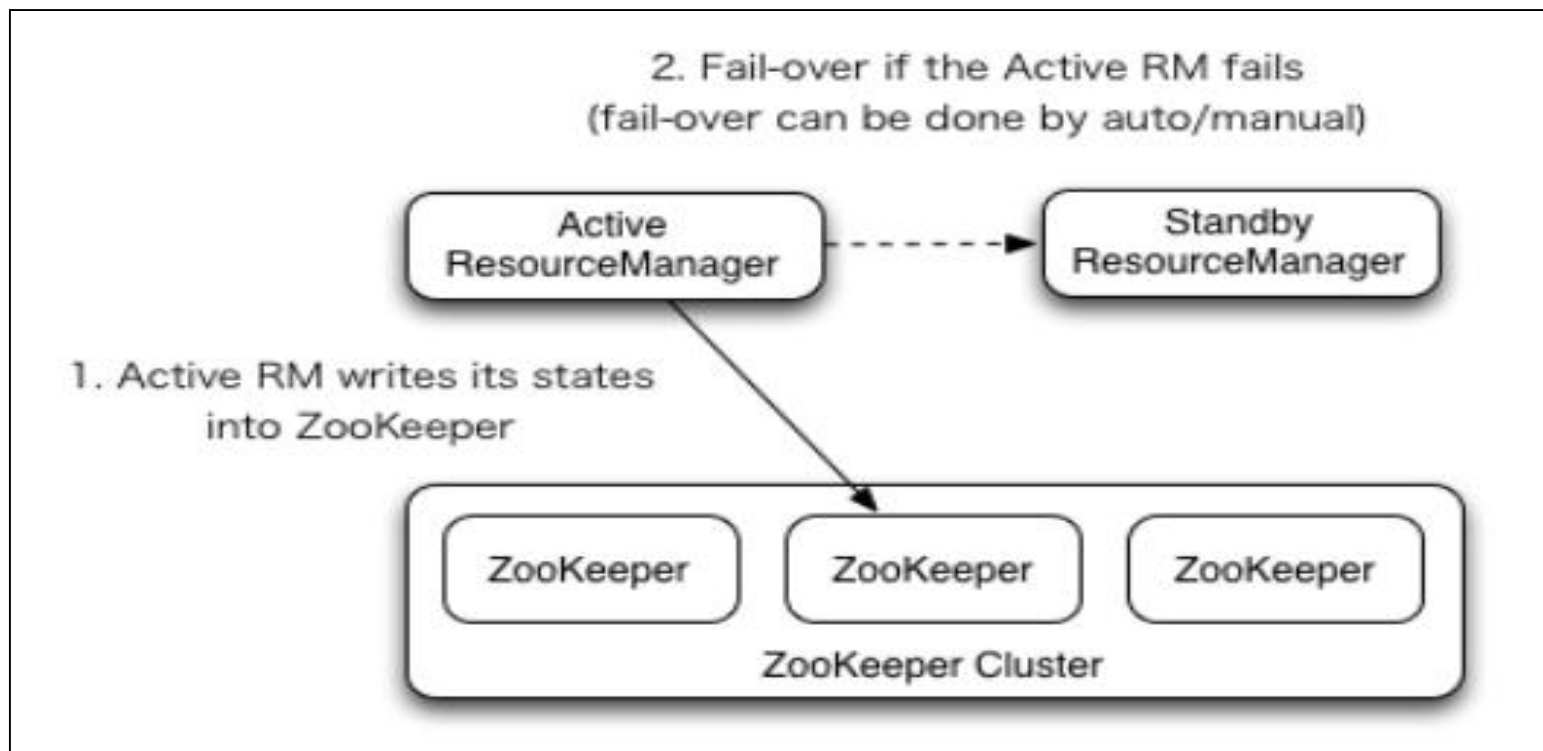
MapReduce On Yarn任务调度



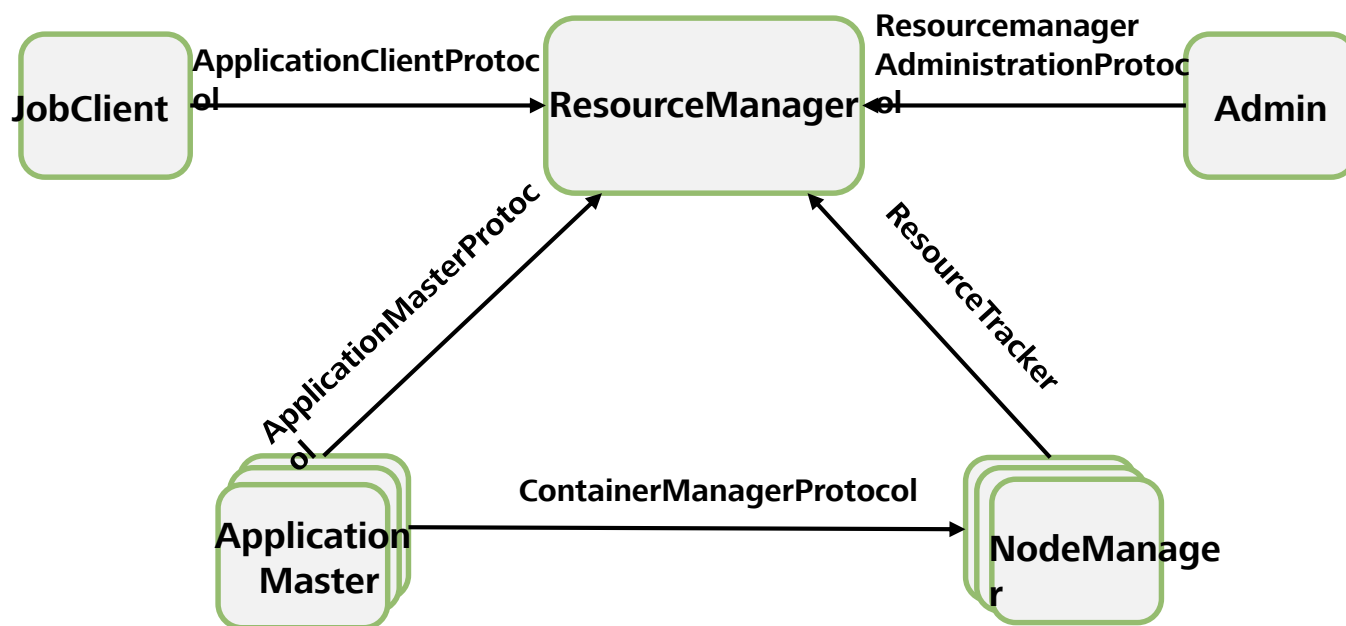
Yarn HA方案

YARN中的**ResourceManager**负责整个集群的资源管理和任务调度，在以前的版本，**ResourceManager**在**YARN**集群中存在单点故障的问题。**YARN**高可用性方案通过引入冗余的**ResourceManager**节点的方式，解决了这个基础服务的可靠性和容错性问题。

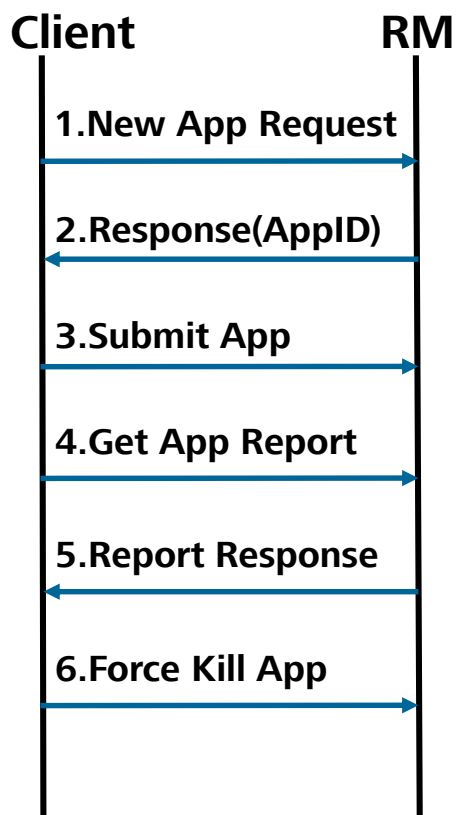
。



Yarn通信协议

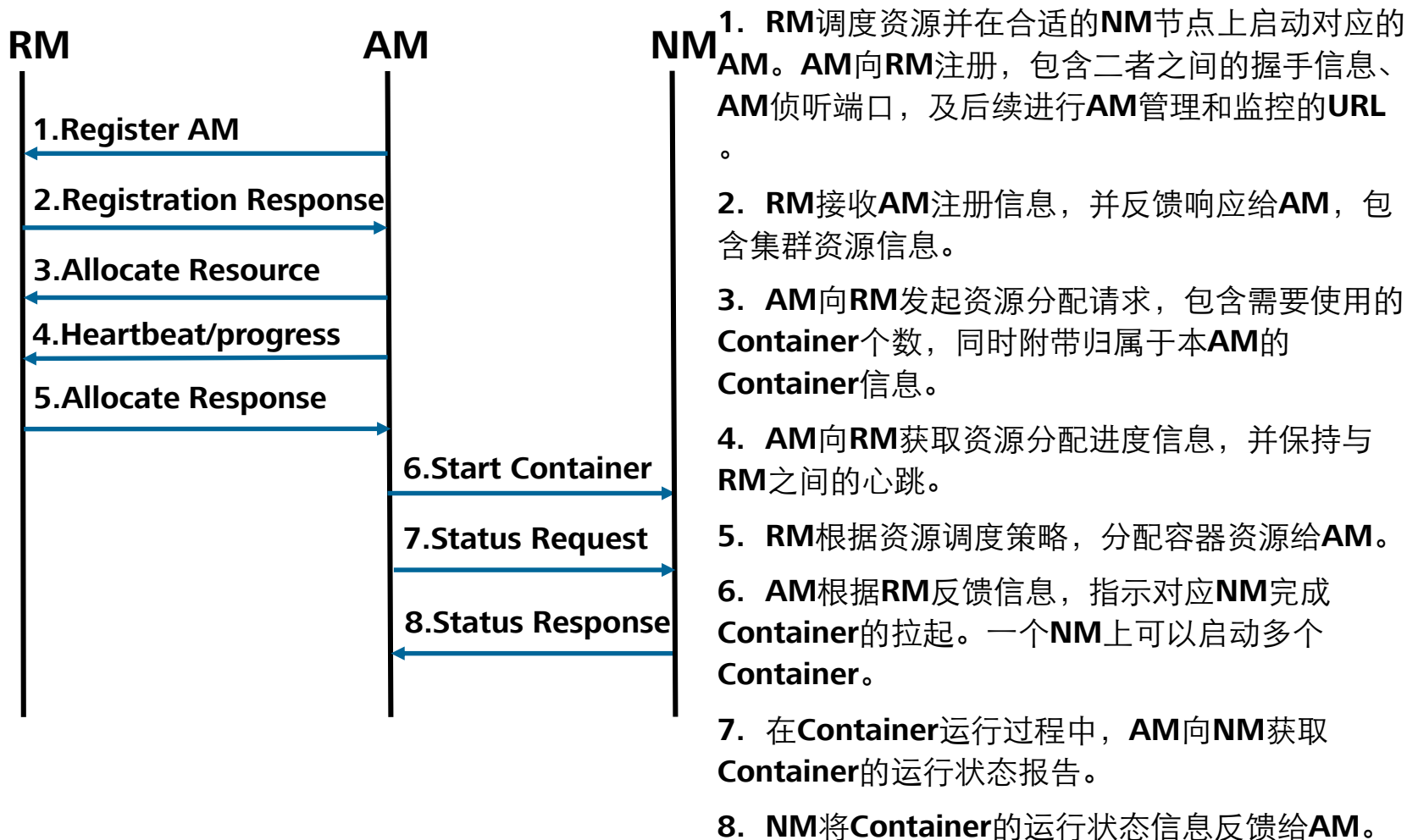


Yarn 客户端与RM的内部交互原理



1. 客户端通知**RM**提交一个应用。
2. **RM**生成一个唯一标识的应用**ID**，又叫**Job ID**，同时将当前**NM**集群的资源描述信息反馈给客户端。
3. 客户端根据**RM**的反馈信息，开始**Job**提交之前的初始化过程，包括调度队列、用户及优先级信息，和**RM**创建、启动**AM**所需的信息（例如应用**Jar**文件、**Job**资源信息、安全**Token**或其他资源描述）。
4. 客户端向**RM**查询、获取应用的执行进展报告。
5. **RM**将应用执行进展报告发送给**MR Client**。
6. 如有必要，客户端可直接通知**RM**终止**Application**的运行。

Yarn 关键组件内部交互原理





目录

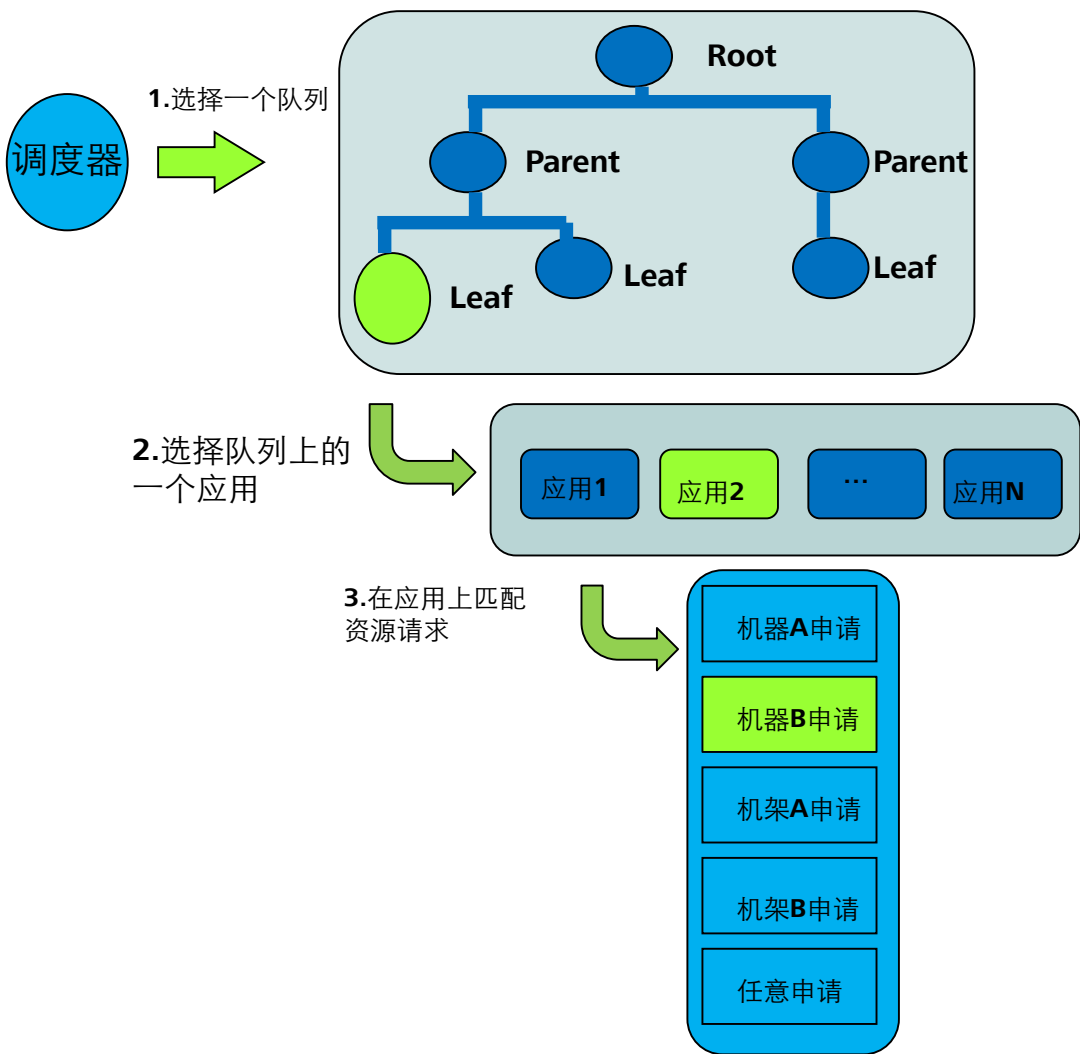
1. MapReduce和Yarn基本介绍
2. MapReduce和Yarn功能与架构
3. Yarn的资源管理和任务调度
4. 增强特性

资源管理

- 当前**YARN**支持内存和**CPU**两种资源类型的管理和分配。
- 每个**NodeManager**可分配的内存和**CPU**的数量可以通过配置选项设置（可在**yarn**服务配置页面配置）
 - **yarn.nodemanager.resource.memory-mb**
 - **yarn.nodemanager.vmem-pmem-ratio**
 - **yarn.nodemanager.resource.cpu-vcore**

资源分配模型

- 调度器维护一群队列的信息。用户可以向一个或者多个队列提交应用。
- 每次**NM**心跳的时候，调度器根据一定的规则选择一个队列，再在队列上选择一个应用，尝试在这个应用上分配资源。
- 调度器会优先匹配本地资源的申请请求，其次是同机架的，最后是任意机器的。



容量调度器的介绍

- 容量调度器使得**Hadoop**应用能够共享的、多用户的、操作简便的运行在集群上，同时最大化集群的吞吐量和利用率。
- 容量调度器以队列为单位划分资源，每个队列都有资源使用的下限和上限。每个用户可以设定资源使用上限。管理员可以约束单个队列、用户或作业的资源使用。支持作业优先级，但不支持资源抢占。

容量调度器的特点

- **容量保证**：管理员可为每个队列设置资源最低保证和资源使用上限，所有提交到该队列的应用程序共享这些资源。
- **灵活性**：如果一个队列中的资源有剩余，可以暂时共享给那些需要资源的队列，当该队列有新的应用程序提交，则其他队列释放的资源会归还给该队列。
- **支持优先级**：队列支持任务优先级调度（默认是**FIFO**）。
- **多重租赁**：支持多用户共享集群和多应用程序同时运行。为防止单个应用程序、用户或者队列独占集群资源，管理员可为之增加多重约束。
- **动态更新配置文件**：管理员可根据需要动态修改配置参数，以实现在线集群管理。

容量调度器的任务选择

- 调度时，首先按以下策略选择一个合适队列：
 - 资源利用量最低的队列优先，比如同级的两个队列**Q1**和**Q2**，它们的容量均为**30**，而**Q1**已使用**10**，**Q2**已使用**12**，则会优先将资源分配给**Q1**；
 - 最小队列层级优先，例如：**QueueA**与**QueueB.childQueueB**，则**QueueA**优先；
 - 资源回收请求队列优先。
- 然后按以下策略选择该队列中一个任务：
 - 按照任务优先级和提交时间顺序选择，同时考虑用户资源量限制和内存限制。

队列资源限制

队列的创建是在多租户页面，当创建一个租户关联Yarn服务时，会创建同名的队列。比如先创建QueueA,QueueB两个租户，即对应Yarn两个队列。



队列资源限制

队列的资源容量（百分比）：

例如，有**default**、**QueueA**、**QueueB**三个队列，每个队列都有一个[队列名].**capacity**配置；

Default队列容量为整个集群资源的**20%**，

QueueA队列容量为整个集群资源的**10%**，

QueueB队列容量为整个集群资源的**10%**，后台有一个影子队列**root-default**使队列之和达到**100%**

。

在集群的**Manager**页面点击“租户管理”»“动态资源计划”»“资源分布策略”可以看到下图页

资源分配 ●

租户名（队列）	资源容量
QueueA(root.QueueA)	10%
QueueB(root.QueueB)	10%
TestParent(root.TestParent)	10%
testchild(root.TestParent.testchild)	10%
default(root.default)	20%
testyarn(root.testyarn)	5%

队列资源限制

共享空闲资源

- 由于存在资源共享，因此一个队列使用的资源可能超过其容量（例如**QueueA.capacity**），而最多使用资源量可通过该参数限制。
 - **yarn.scheduler.capacity.root.QueueA.maximum-capacity**（此参数也是在上页胶片展示的**FusionInsight**页面配置）
- 如果某个队列任务较少，可将剩余资源共享给其他队列，例如**QueueA**的**maximum-capacity**配置为**100**，假设当前只有**QueueA**在运行任务，理论上**QueueA**可以占用整个集群**100%**的资源。

用户限制和任务限制

用户限制和任务限制的参数可通过“租户管理” > “动态资源计划” > “队列配置” 进行配置

资源分布策略				
队列配置				
租户名 (队列)	最大应用数	AM最大资源百分比	用户资源最小上限百分比	用户资源上限因子
QueueA(root.QueueA)	1000	0.1	100%	10
QueueB(root.QueueB)	1000	0.1	100%	10
TestParent(root.TestParent)	1000	0.1	100%	10
testchild(root.TestParent.t...	1000	0.1	100%	10
default(root.default)	1000	0.1	100%	10
testyam(root.testyam)	1000	0.1	100%	10

用户限制

每个用户最低资源保障（百分比）

任何时刻，一个队列中每个用户可使用的资源量均有一定的限制，当一个队列中同时运行多个用户的任务时，每个用户的可使用资源量在一个最小值与最大值之间浮动，其中，最大值取决于正在运行的任务数目，而最小值则由**minimum-user-limit-percent**决定。

例如，设置队列A的这个值为**25**，即**yarn.scheduler.capacity.root.QueueA.minimum-user-limit-percent=25**，那么随着提任务的用户增加，队列资源的调整如下：

第1个用户提交任务到QueueA	会获得QueueA的100%资源
第2个用户提交任务到QueueA	每个用户会最多获得50%的资源
第3个用户提交任务到QueueA	每个用户会最多获得33.33%的资源
第4个用户提交任务到QueueA	每个用户会最多获得25%的资源
第5个用户提交任务到QueueA	为了保障每个用户最低能获得25%的资源，第5个用户将无法再获取到QueueA的资源，必须等待资源的释放。

用户限制

- 每个用户最多可使用的资源量（所在队列容量的倍数）
- **queue**容量的倍数，用来设置一个**user**可以获取更多的资源。
yarn.scheduler.capacity.root.QueueD.user-limit-factor = 1
- 默认值为**1**，表示一个**user**获取的资源容量不能超过**queue**配置的**capacity**，无论集群有多少空闲资源，最多不超过**maximum-capacity**。

任务限制

- 最大活跃任务数

- 整个集群中允许的最大活跃任务数，包括运行或挂起状态的所有任务，当提交的任务申请数据达到限制以后，新提交的任务将会被拒绝。默认**10000**
- **yarn.scheduler.capacity.maximum-applications=10000**

- 每个队列最大任务数

- 对于每个队列，可以提交的最大任务数，以**QueueA**为例，可以在队列配置页面配置，默认是**1000**，即此队列允许最多**1000**个活跃任务。

- 每个用户可以提交的最大任务数

- 这个数值依赖每个队列最大任务数，假设根据上面的结果，**QueueA**最多可以提交**1000**个任务，那么对于每个用户而言，可以向**QueueA**提交的最大任务数为
- **1000* yarn.scheduler.capacity.root.QueueA.minimum-user-limit-percent***
yarn.scheduler.capacity.root.QueueA.user-limit-factor

队列信息

队列的信息还可以通过Yarn webUI进行查看，进入方法是“服务管理” > “Yarn” > “ResourceManager（主）” > “Scheduler”

Legend: Capacity Used Used (over capacity) Max Capacity

Partition: <DEFAULT_PARTITION> <memory:24576, vCores:24>
Queue: root
Queue: QueueA

Used Capacity:	0.0%
Configured Capacity:	10.0%
Configured Max Capacity:	100.0%
Absolute Used Capacity:	0.0%
Absolute Configured Capacity:	10.0%
Absolute Configured Max Capacity:	100.0%
Used Resources:	<memory:0, vCores:0>
Configured Max Application Master Limit:	10.0
Max Application Master Resources:	<memory:2560, vCores:3>
Used Application Master Resources:	<memory:0, vCores:0>
Max Application Master Resources Per User:	<memory:2560, vCores:3>

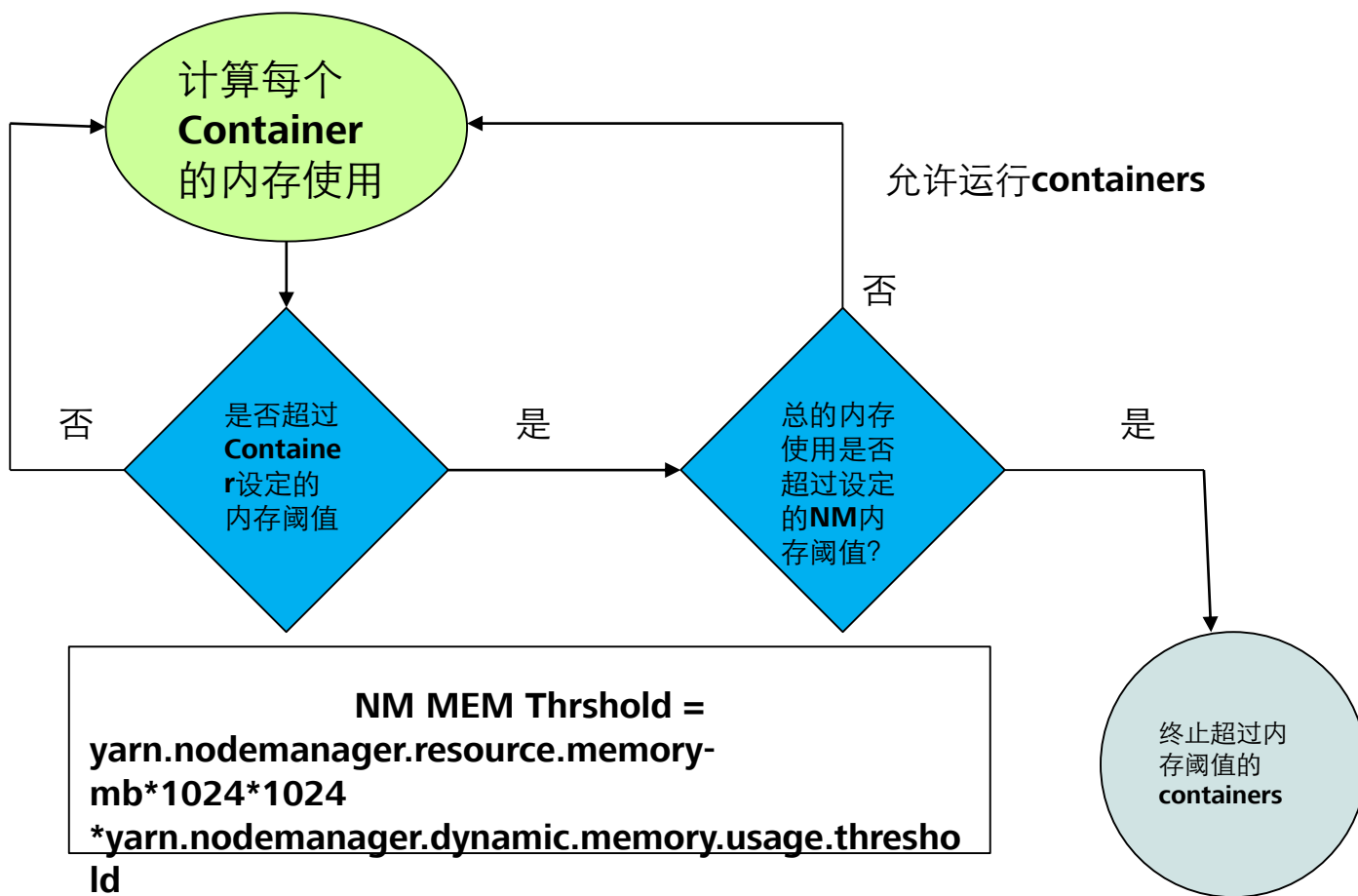
Queue State:	RUNNING
Num Schedulable Applications:	0
Num Non-Schedulable Applications:	0
Num Containers:	0



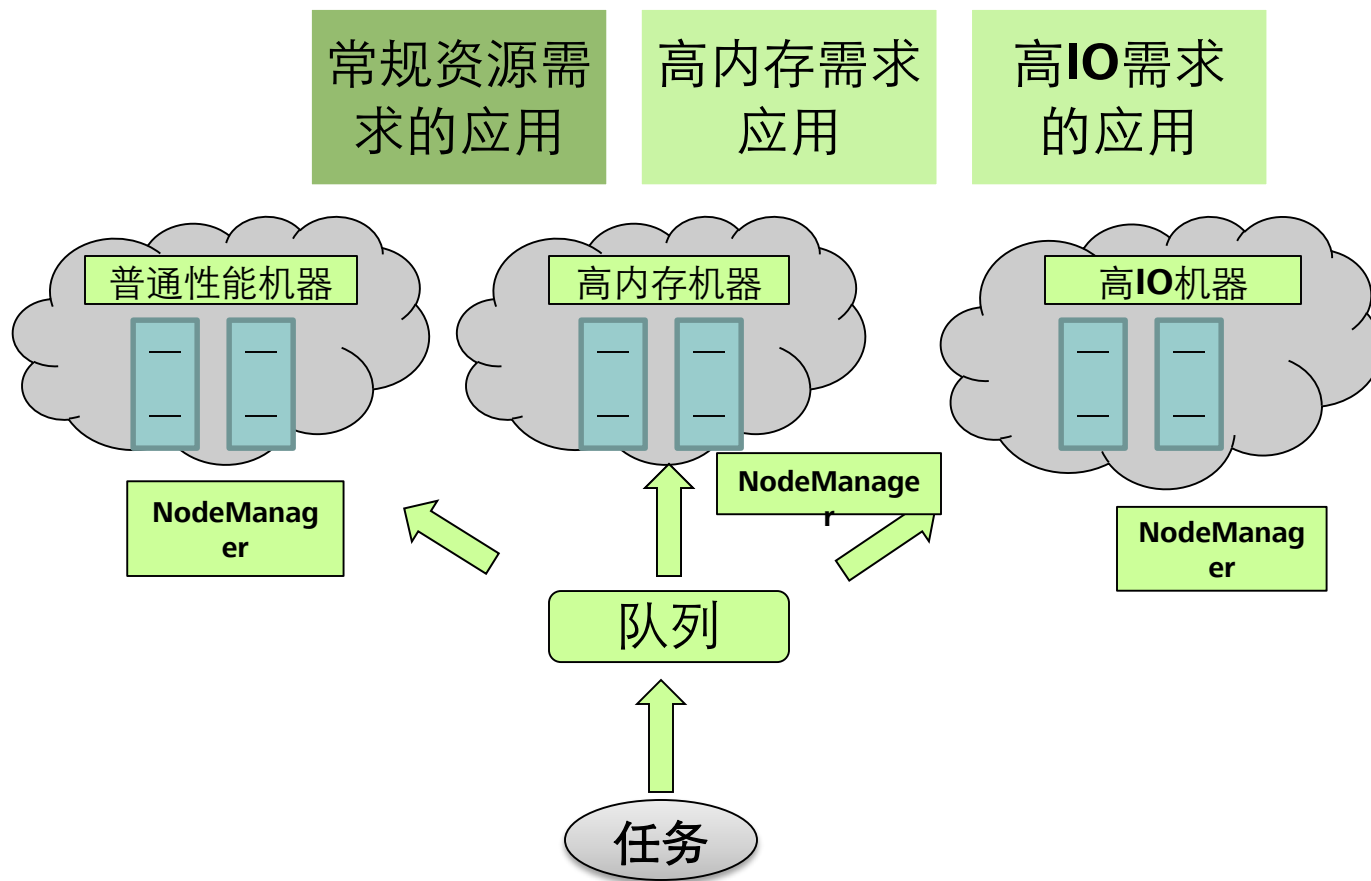
目录

1. MapReduce和Yarn基本介绍
2. MapReduce和Yarn功能与架构
3. Yarn的资源管理和任务调度
4. 增强特性

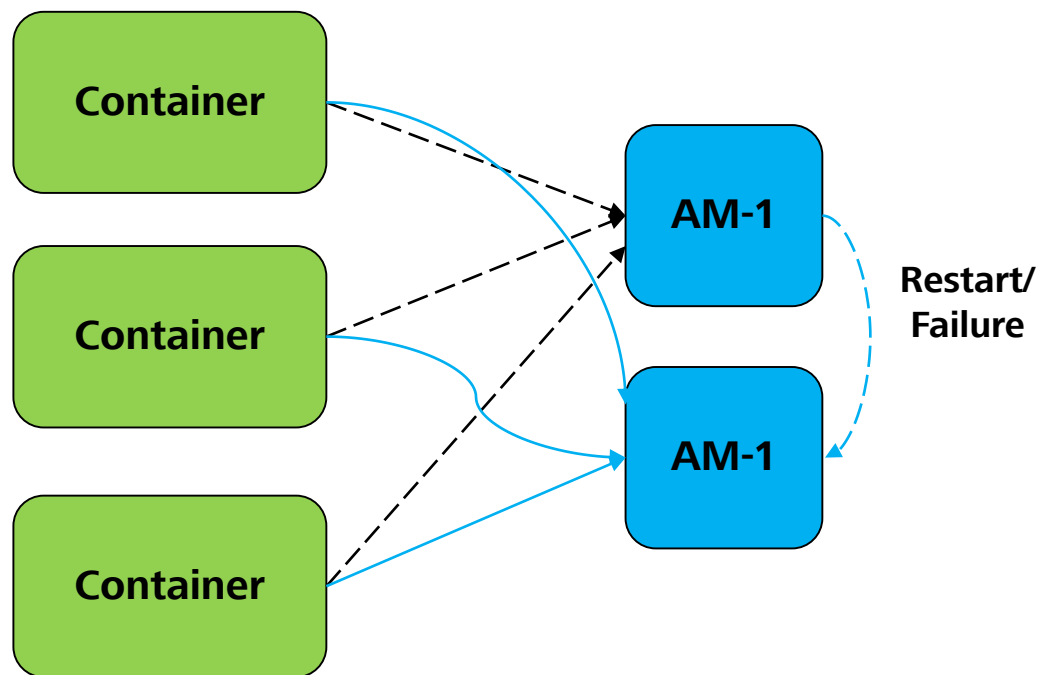
增强特性—Yarn动态内存管理



增强特性—Yarn基于标签调度



增强特性—Yarn AM作业保留





思考题

- 1.请简述**MapReduce**的工作原理。
- 2.相比于**MapReduce**，**Yarn**有哪些优点？
- 3.请简述**Yarn**的工作原理。



思考题

1. 下面哪些是**MapReduce**的特点？（ ）
 - A. 易于编程
 - B. 良好的扩展性
 - C. 实时计算
 - D. 高容错性
2. **Yarn**中资源抽象用什么表示？（ ）
 - A. 内存
 - B. CPU
 - C. Container
 - D. 磁盘空间



思考题

3. 下面哪个是**MapReduce**适合做的？（ ）

- A. 迭代计算
- B. 离线计算
- C. 实时交互计算
- D. 流式计算

4. 容量调度器的特点？（ ）

- A. 容量保证
- B. 灵活性
- C. 多重租赁
- D. 动态更新配置文件



本章总结

- 讲述了**MR**和**YARN**的应用场景
- 讲述了**MR**和**YARN**的基本架构
- 讲述了**YARN**资源管理与任务调度
- 讲述了**YARN**的增强特性



学习推荐

- 华为**Learning**网站
 - <http://support.huawei.com/learning/Index!toTrainIndex>
- 华为**Support**案例库
 - <http://support.huawei.com/enterprise/servicecenter?lang=zh>

Thank you

www.huawei.com