

# MapReduce和YARN 技术原理

[www.huawei.com](http://www.huawei.com)





# 目标

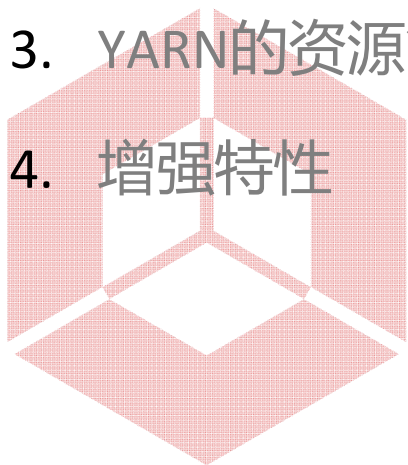
- 学完本课程后，您将能够：
  - 熟悉MapReduce和YARN是什么
  - 掌握MapReduce使用的场景及其原理
  - 掌握MapReduce和YARN功能与架构
  - 熟悉YARN的新特性

泰克教育  
TECH EDUCATION



# 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



泰克教育  
TECH EDUCATION

# MapReduce概述

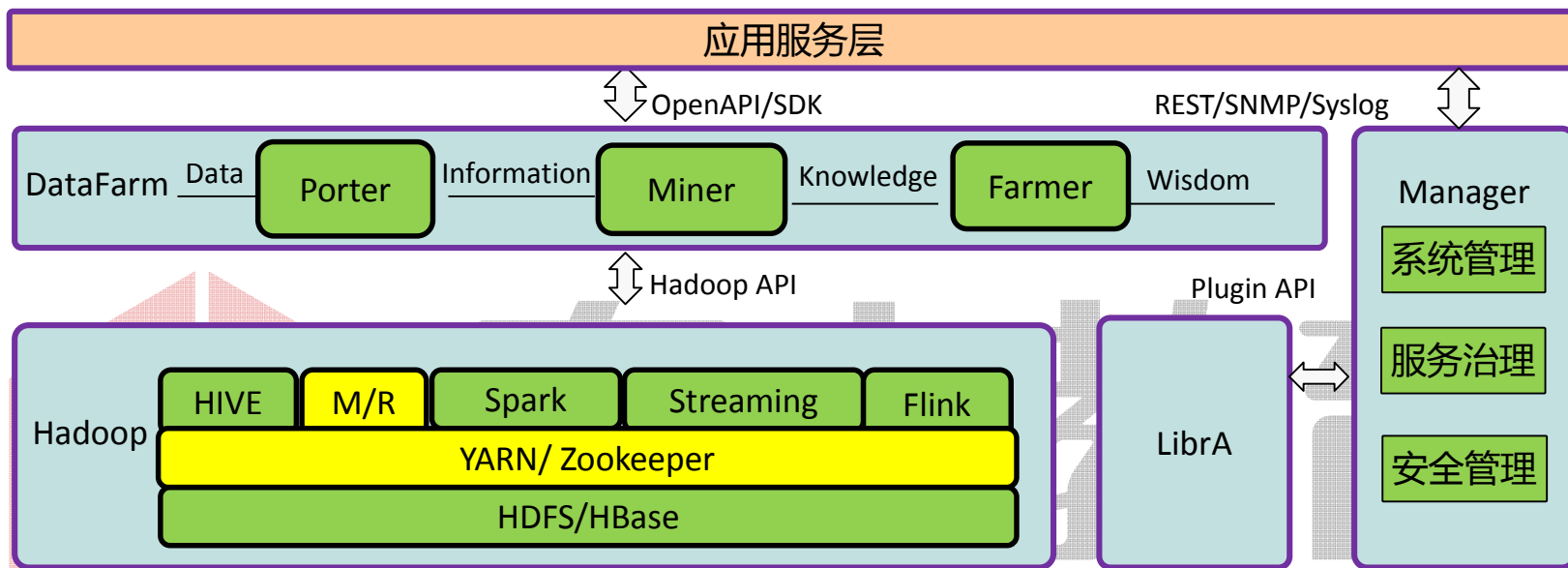
- MapReduce基于Google发布的MapReduce论文设计开发，用于大规模数据集（大于1TB）的并行计算，具有如下特点：
  - 易于编程：程序员仅需描述做什么，具体怎么做交由系统的执行框架处理。
  - 良好的扩展性：可通过添加节点以扩展集群能力。
  - 高容错性：通过计算迁移或数据迁移等策略提高集群的可用性与容错性。

# YARN概述

- 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。

泰克教育  
TECH EDUCATION

# YARN在FusionInsight产品的位置

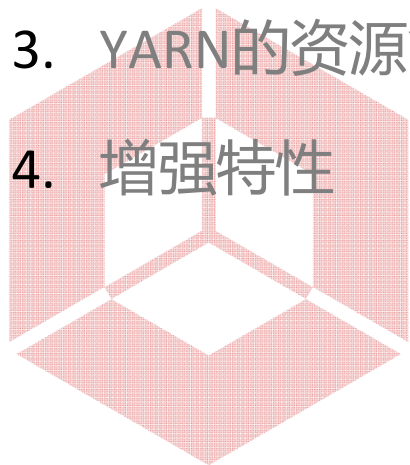


YARN是Hadoop2.0中的资源管理系统，它是一个通用的资源管理模块，可为各类应用程序提供资源管理和调度功能。



# 目录

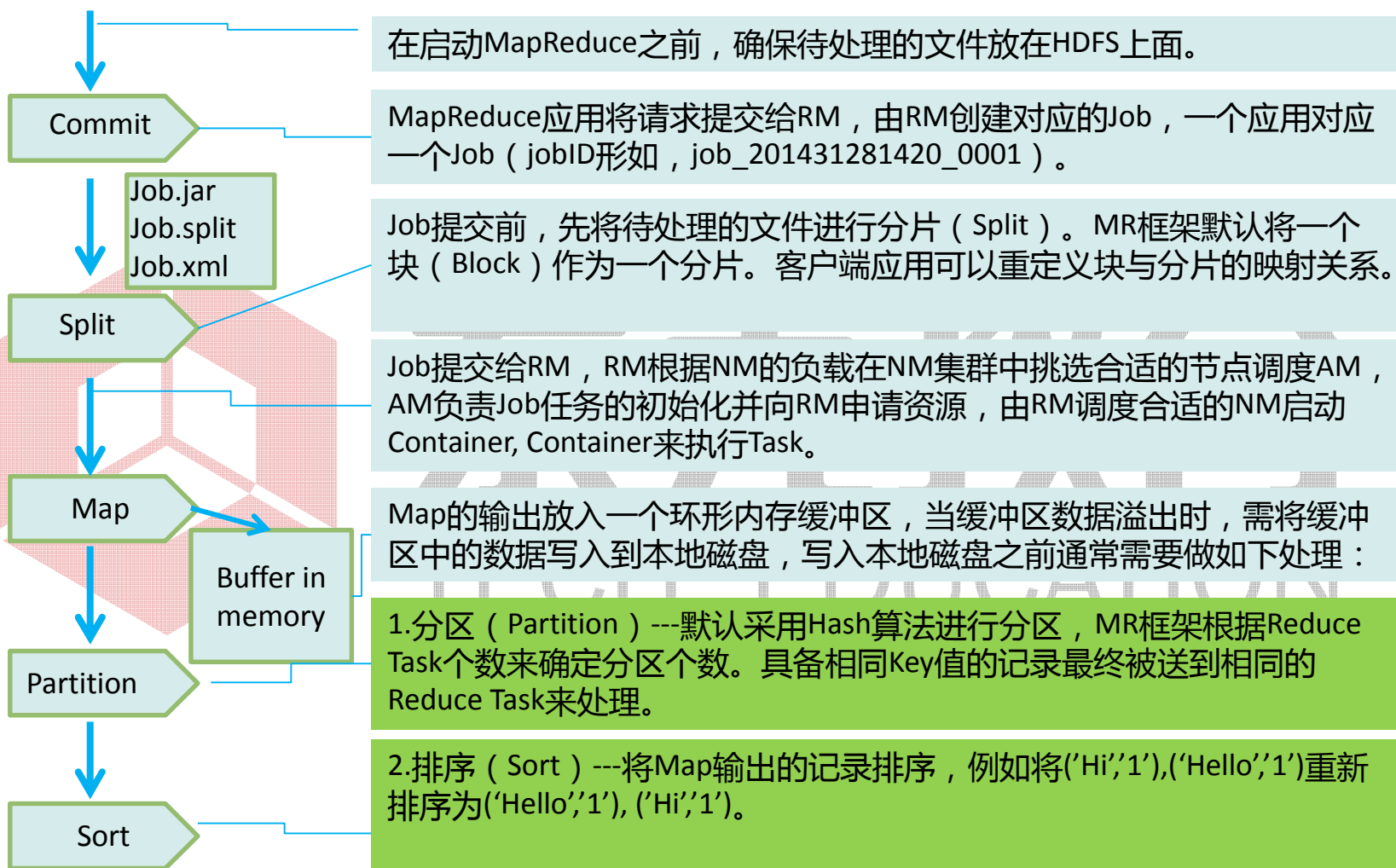
1. MapReduce和YARN基本介绍
- 2. MapReduce和YARN功能与架构**
3. YARN的资源管理和任务调度
4. 增强特性



泰克教育  
TECH EDUCATION

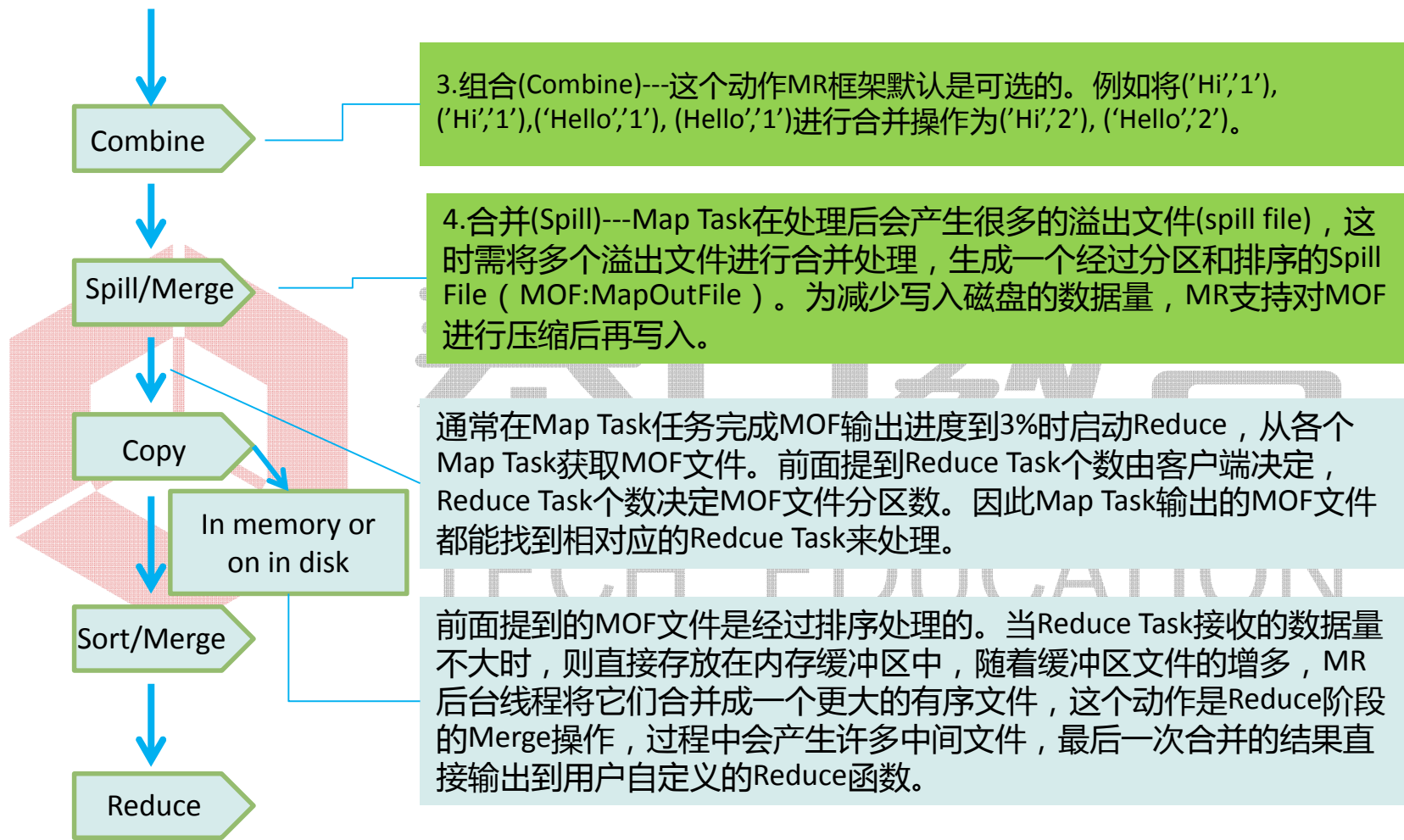


# MapReduce过程详解(1)

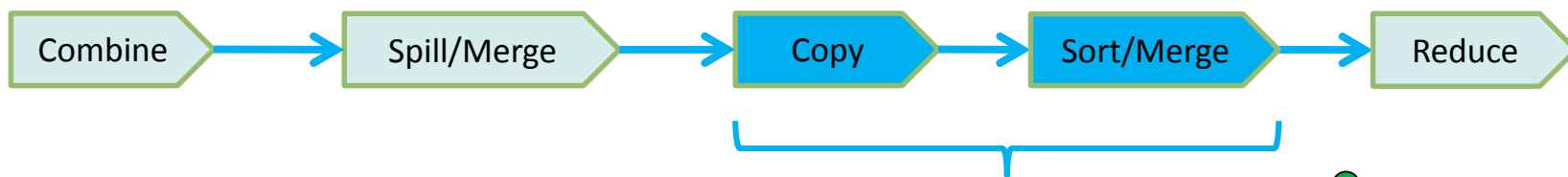




# MapReduce过程详解 (2)

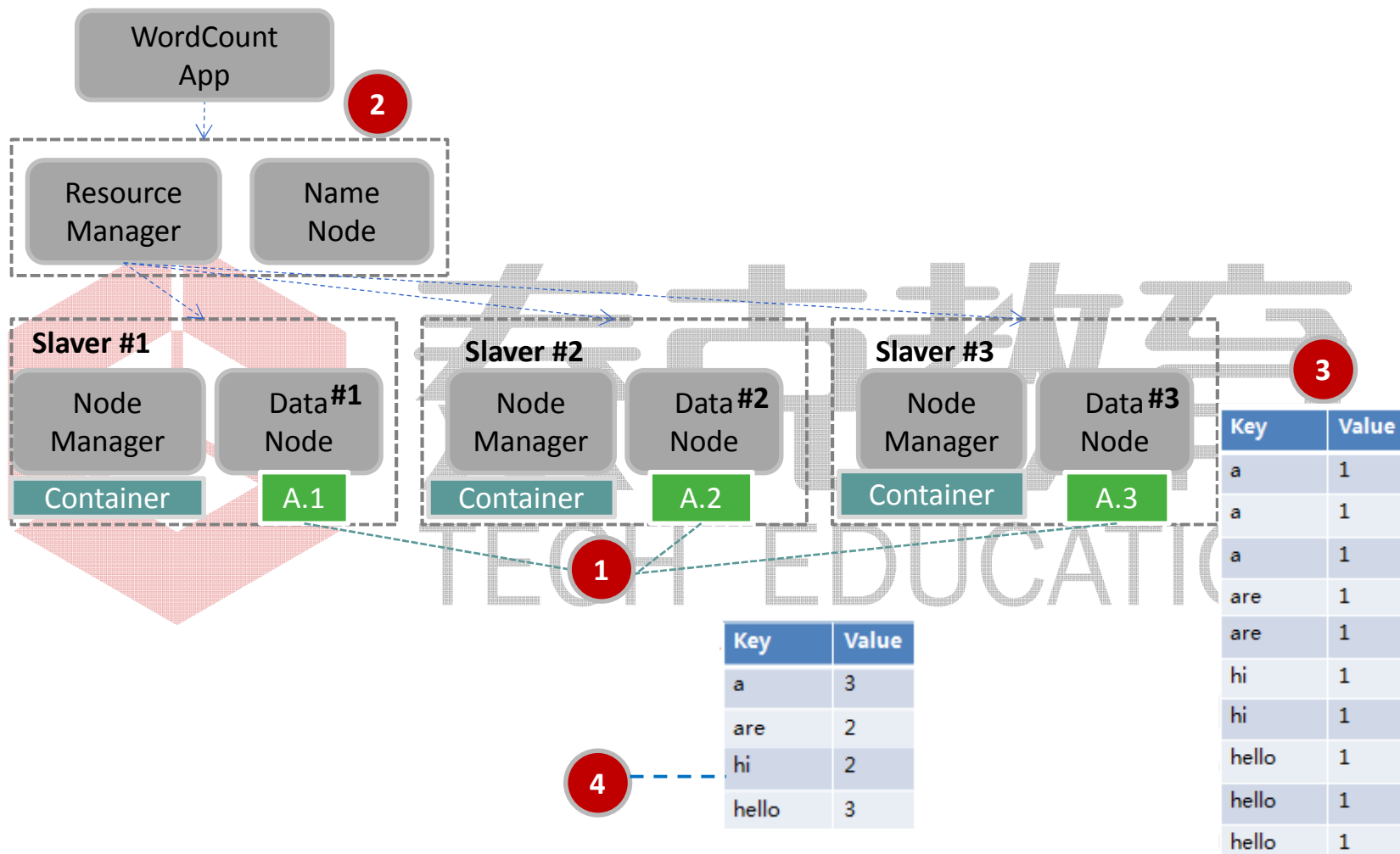


# Shuffle机制

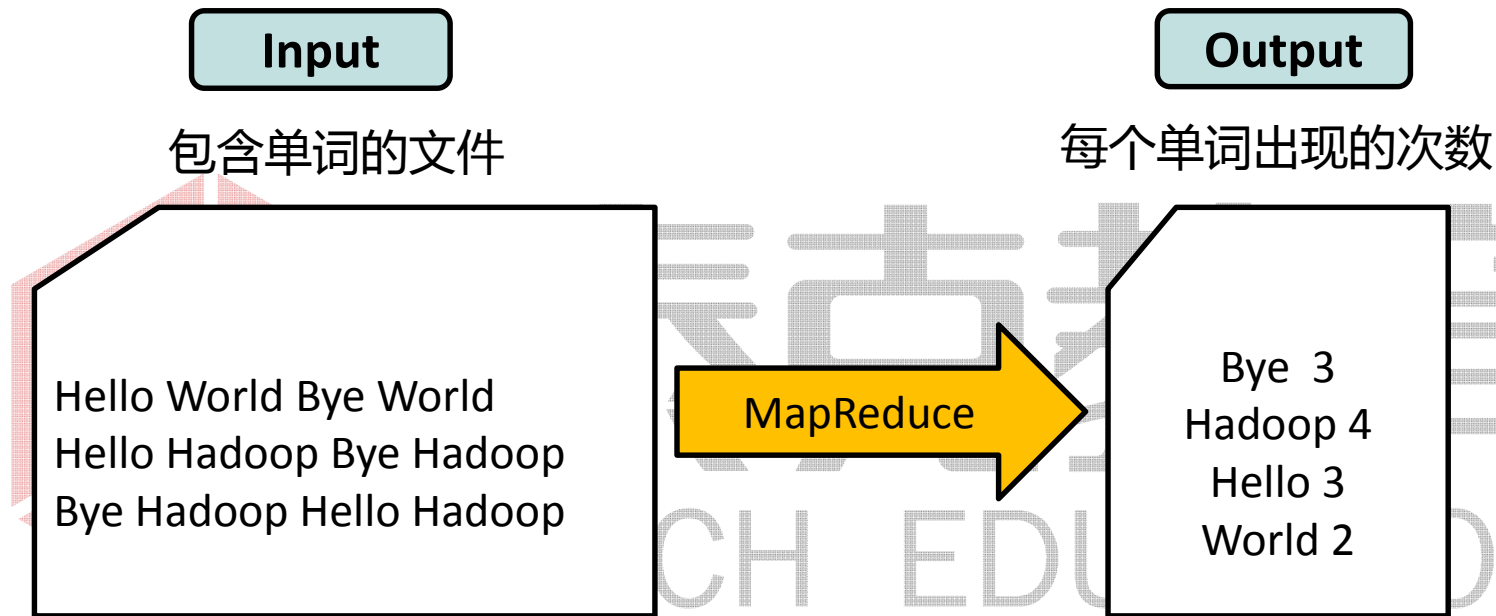


Shuffle的定义：Map阶段和Reduce阶段之间传递中间数据的过程，包括Reduce Task从各个Map Task获取MOF文件的过程，以及对MOF的排序与合并处理。

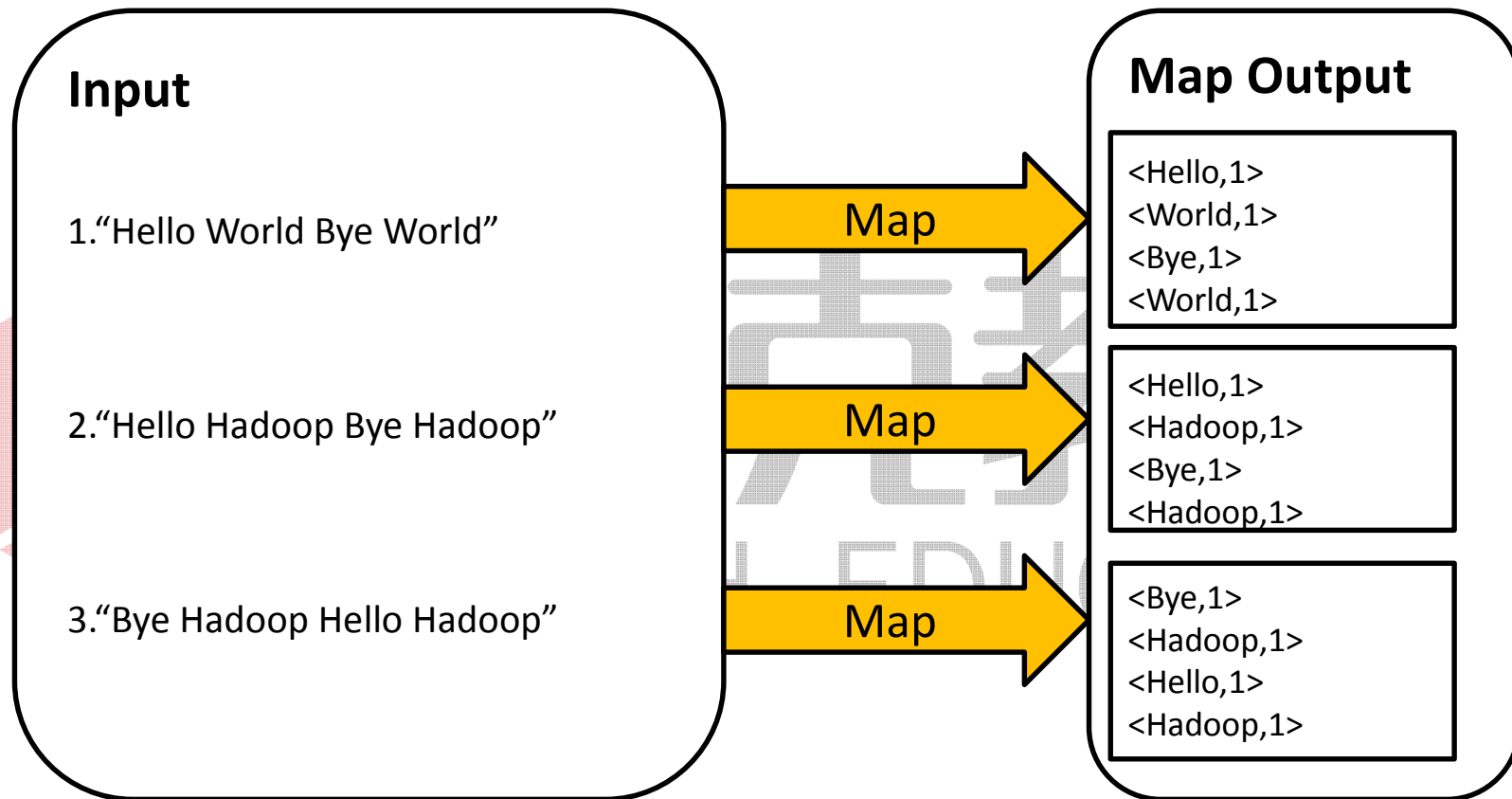
# 典型程序WordCount举例



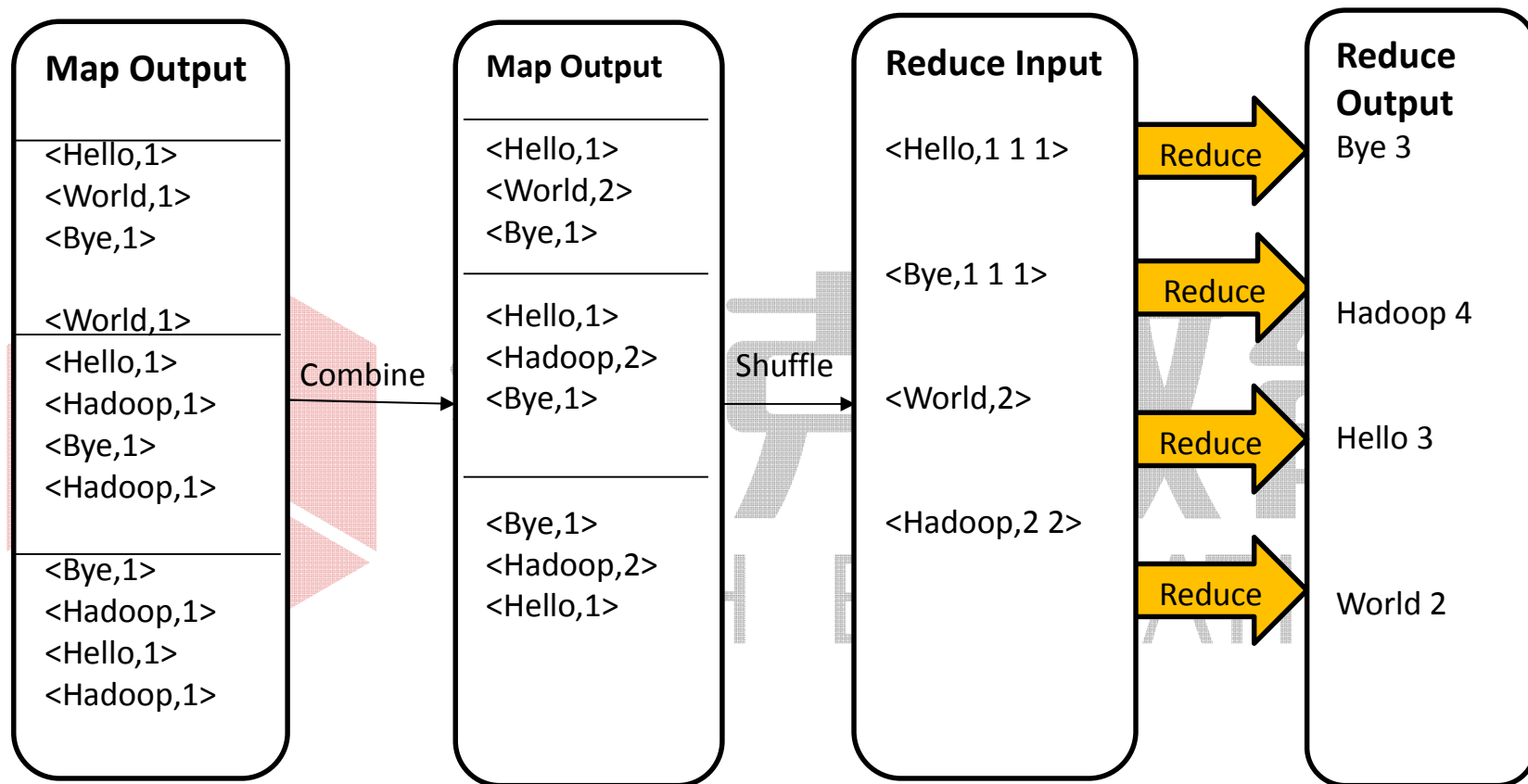
# WordCount程序功能



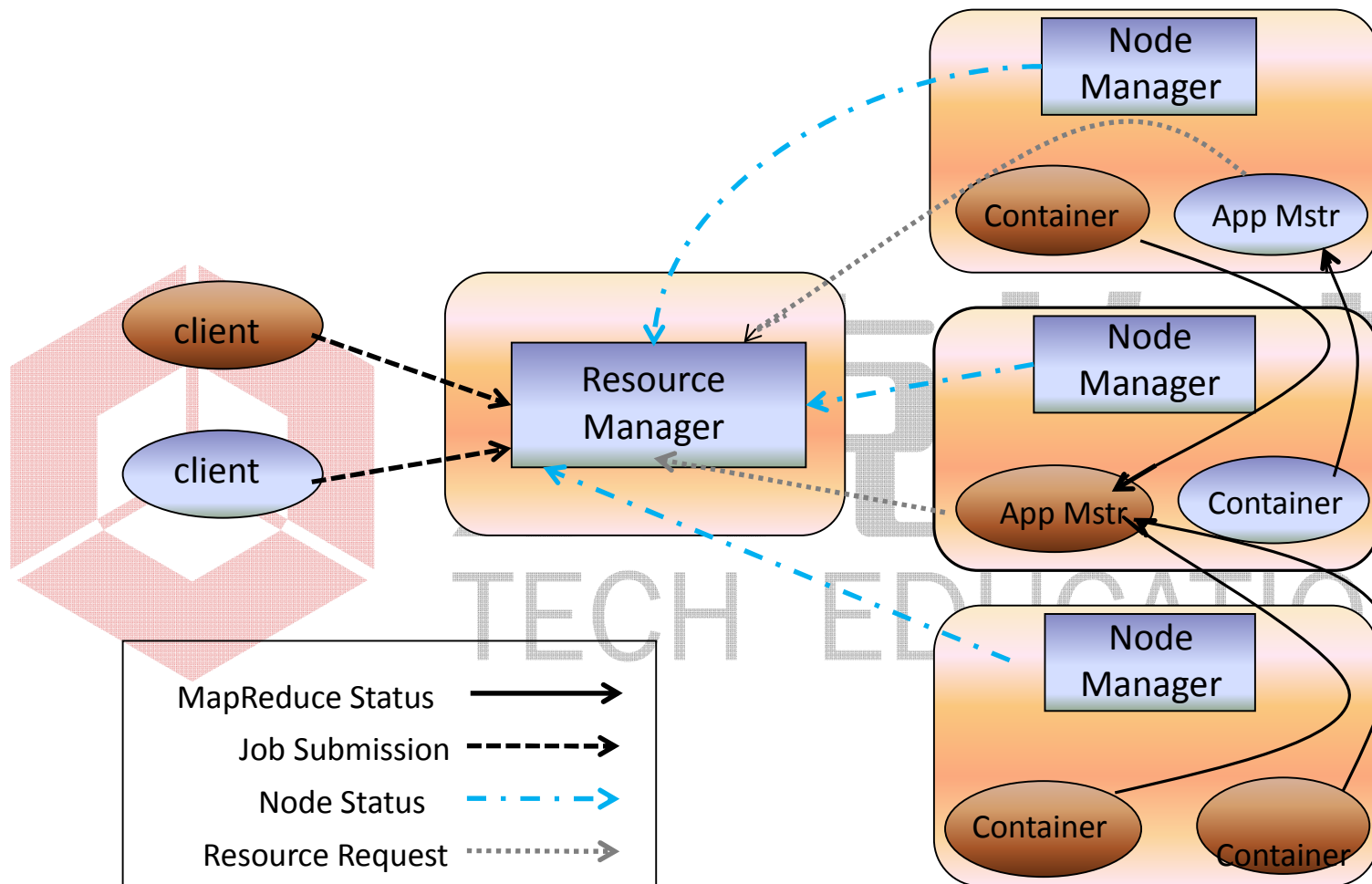
# WordCount的Map过程



# WordCount的Reduce过程

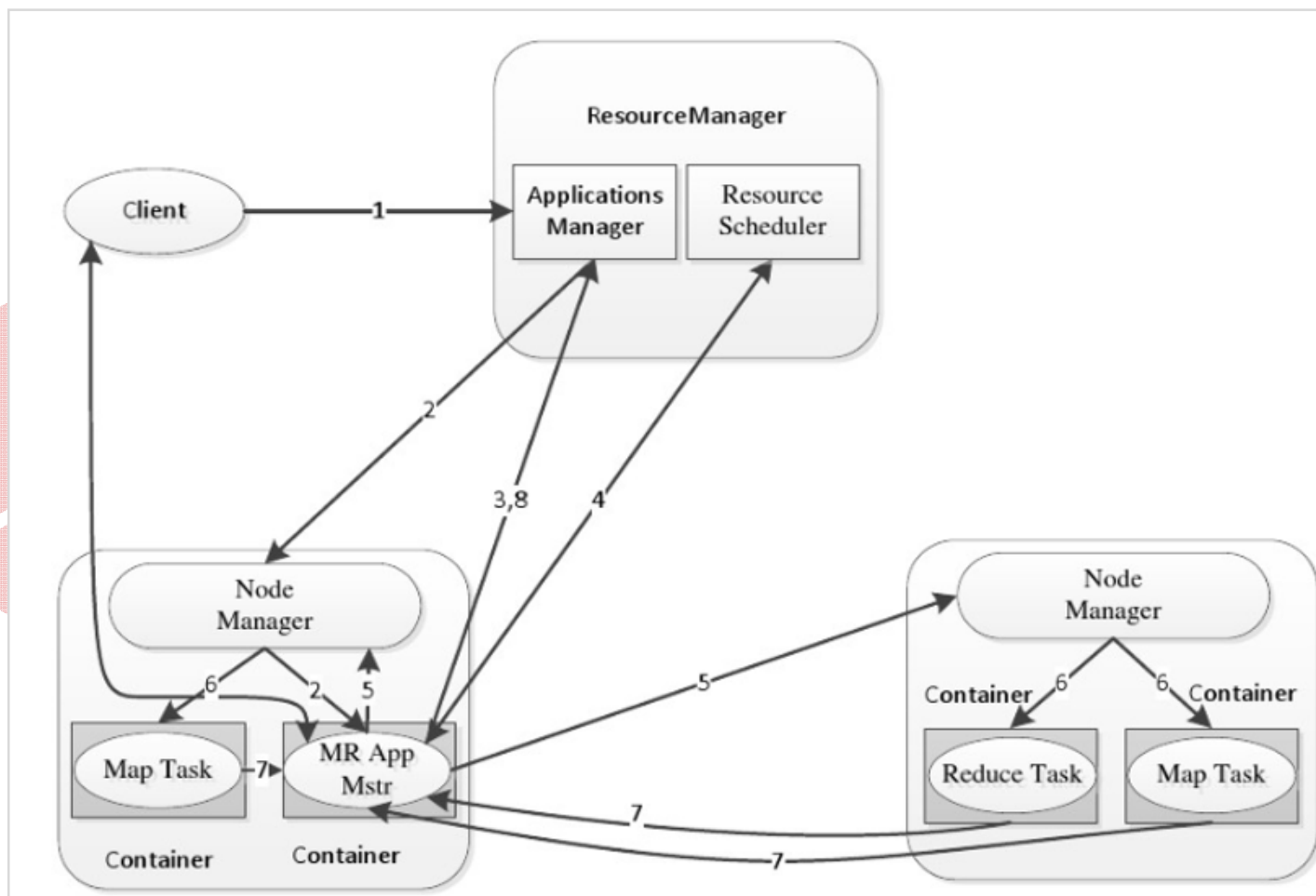


# YARN的组件架构



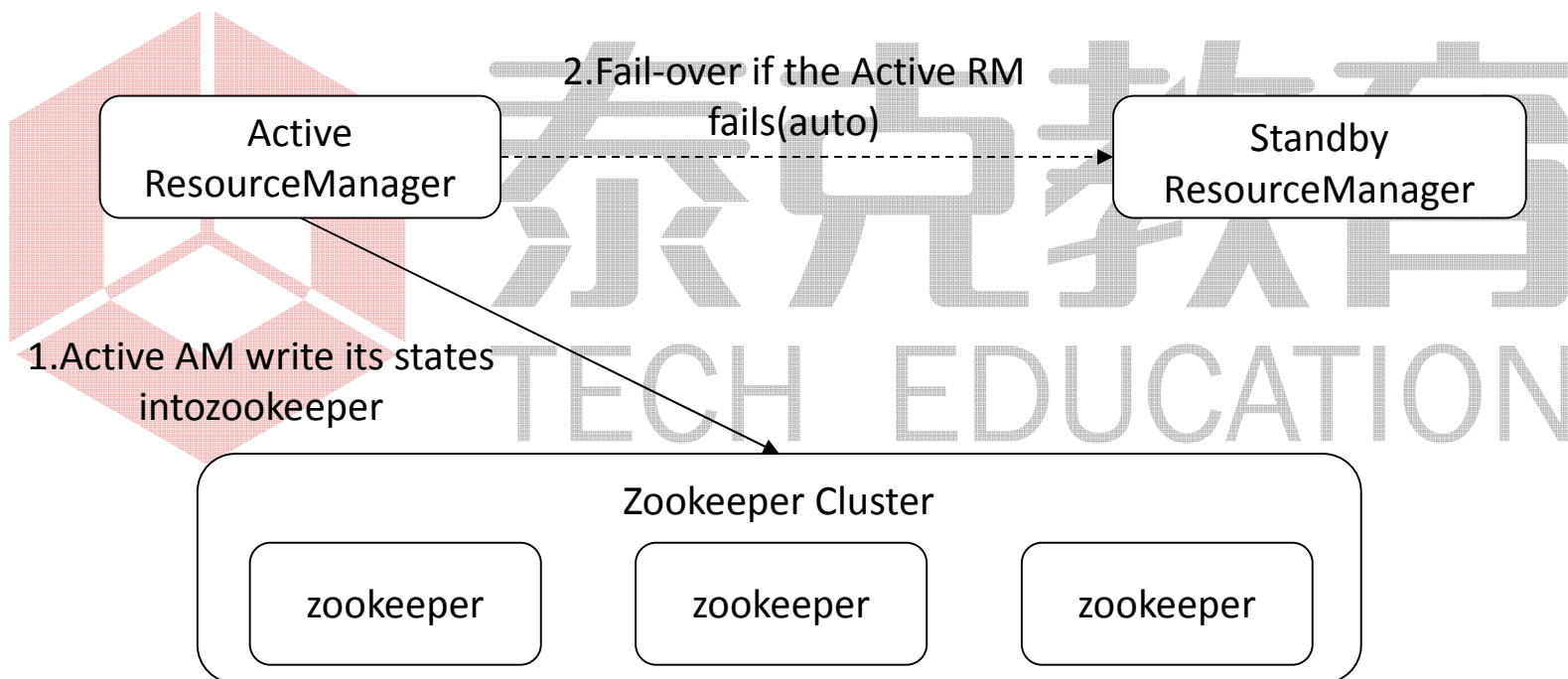


# MapReduce On YARN任务调度流程

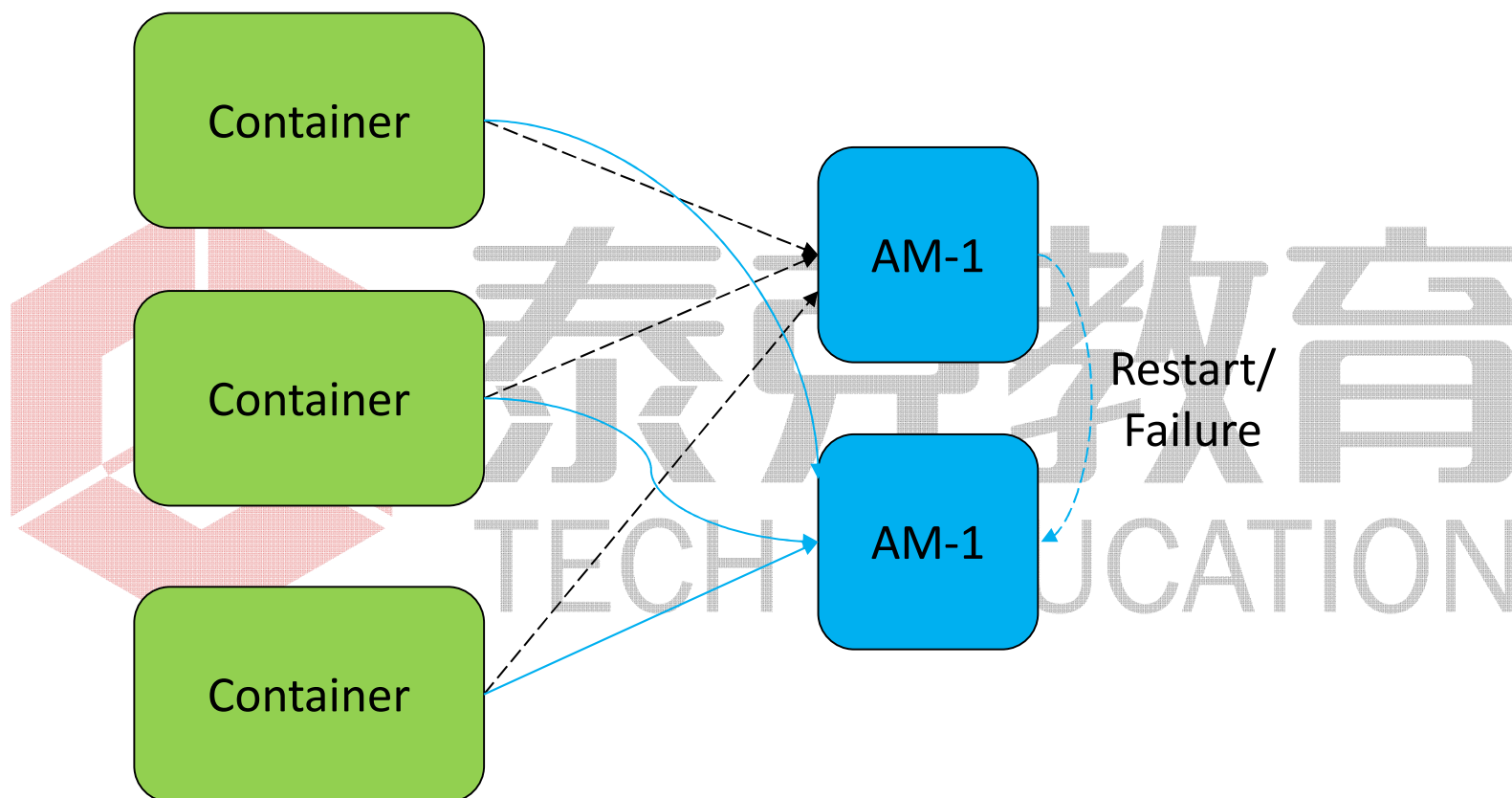


# YARN HA方案

- YARN中的ResourceManager负责整个集群的资源管理和任务调度，YARN高可用性方案通过引入冗余的ResourceManager节点的方式，解决了ResourceManager单点故障问题。



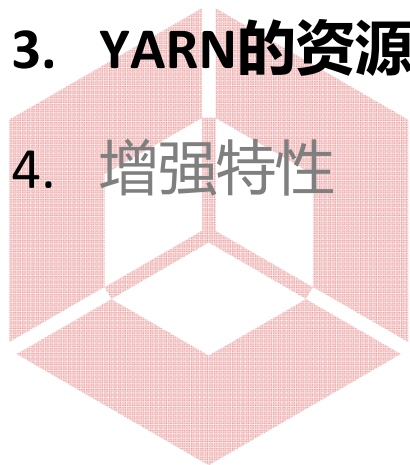
# YARN APPMaster容错机制





# 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

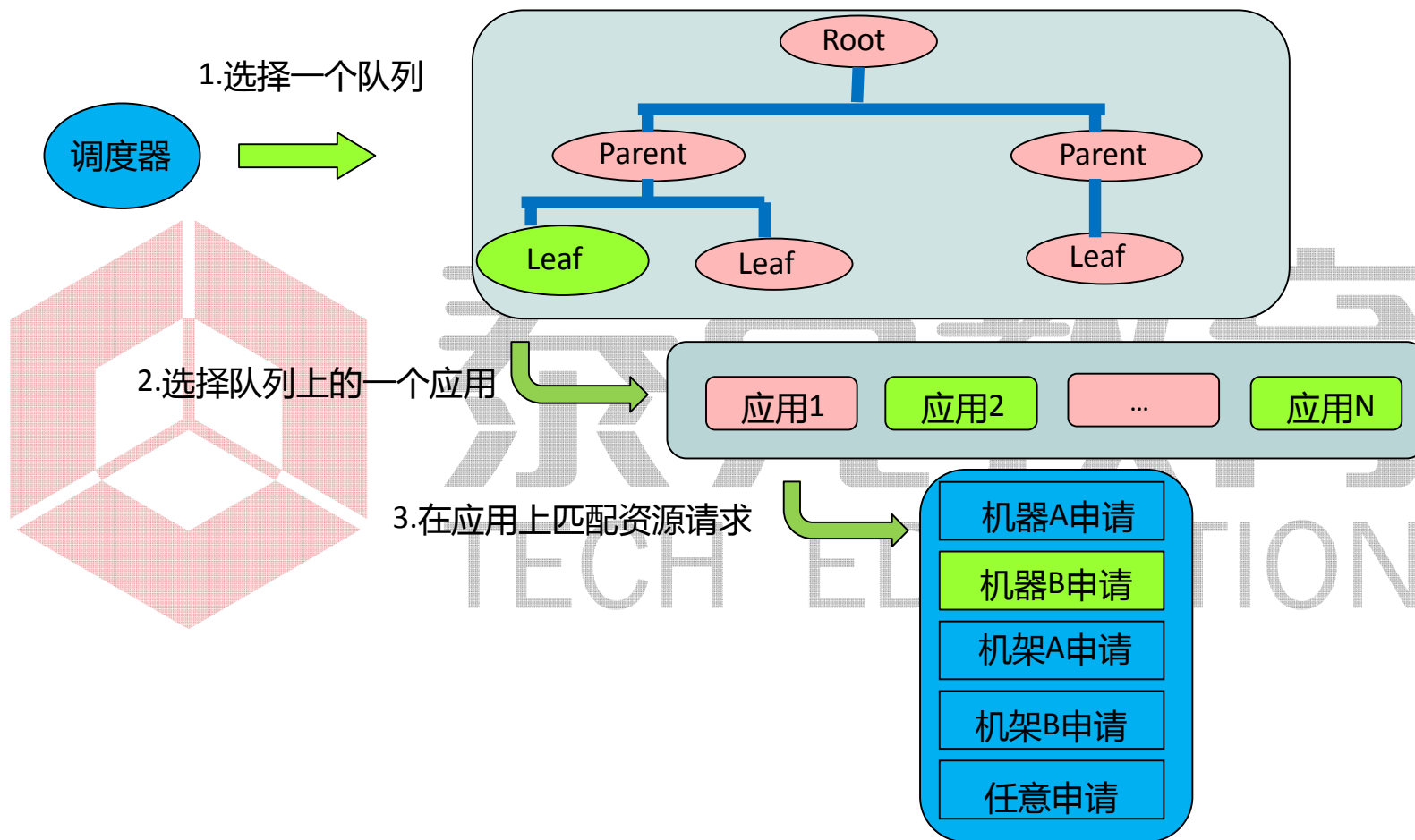


泰克教育  
TECH EDUCATION

# 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





# 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 (1)

- 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 (2)

- 队列的资源容量 (百分比) , 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%。

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%

# 队列资源限制 (3)

- **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），而最大资源使用量可通过参数限制。
- 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	
队列配置 <span>●</span>			
租户名 (队列)	最大应用数	AM最大资源百分比	
QueueA(root.QueueA)	1000	0.1	
QueueB(root.QueueB)	1000	0.1	
<input type="checkbox"/> TestParent(root.TestParent)	1000	0.1	
<input type="checkbox"/> testchild(root.TestParent.t...	1000	0.1	
default(root.default)	1000	0.1	



# 用户限制(1)

- 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`，那么随着提交任务的增加，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(2)

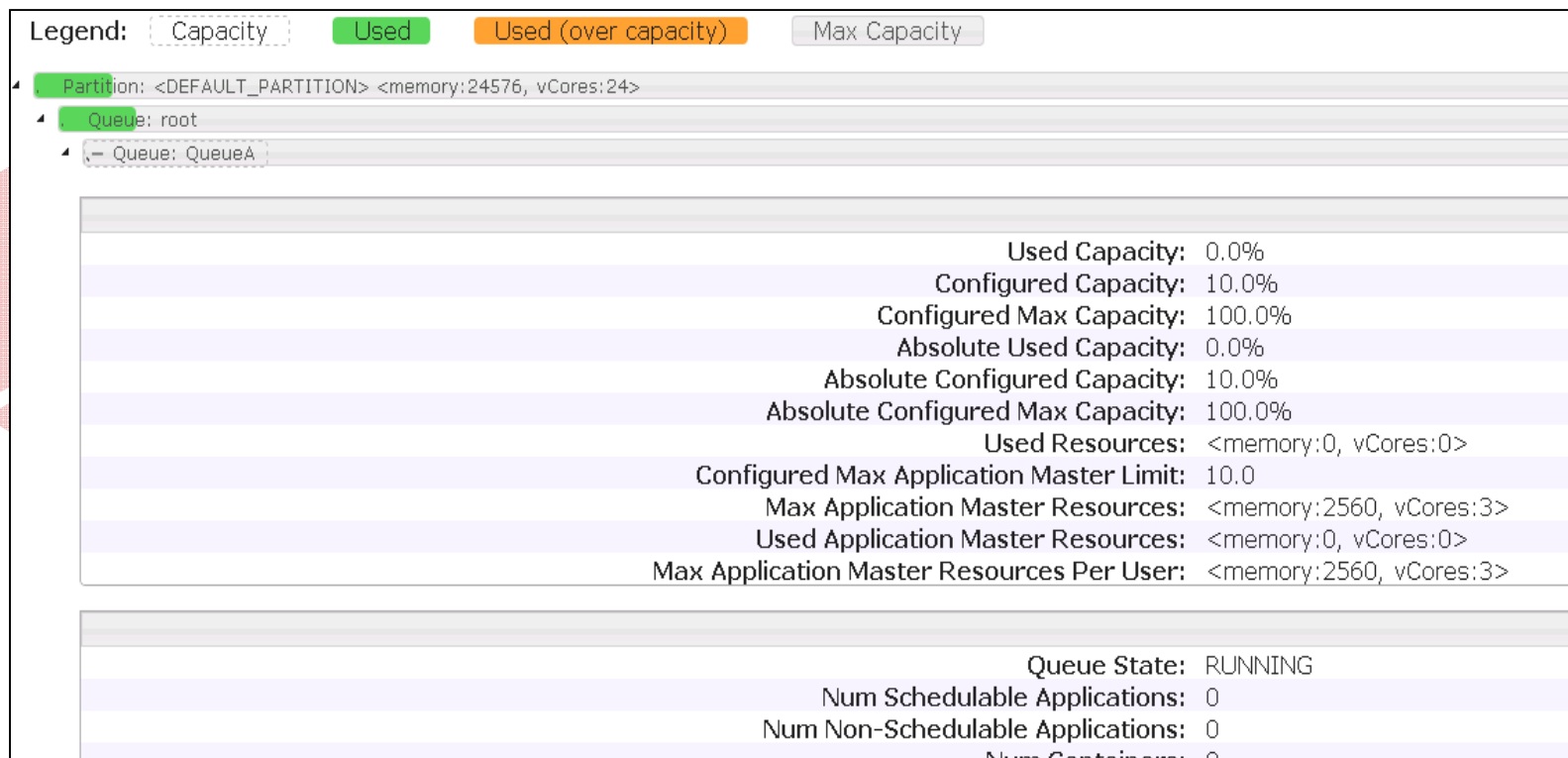
- 每个用户最多可使用的资源量（所在队列容量的倍数）：
  - queue容量的倍数，用来设置一个user可以获取更多的资源。  
Yarn.scheduler.capacity.root.QueueD.user-limit-factor=1。默认值为1，表示一个user获取的资源容量不能超过queue配置的capacity，无论集群有多少空闲资源，最多不超过maximum-capacity。

# 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。默认值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。根据上面的数据，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 $1000 * \text{用户最低资源保障率} ( \text{假设} 25\% ) * \text{用户可使用队列资源的倍数} ( \text{假设} 1 )$ 。

# 查看队列信息

- 队列的信息可以通过YARN webUI进行查看，进入方法是“服务管理” > “YARN” > “ResourceManager (主)” > “Scheduler”。





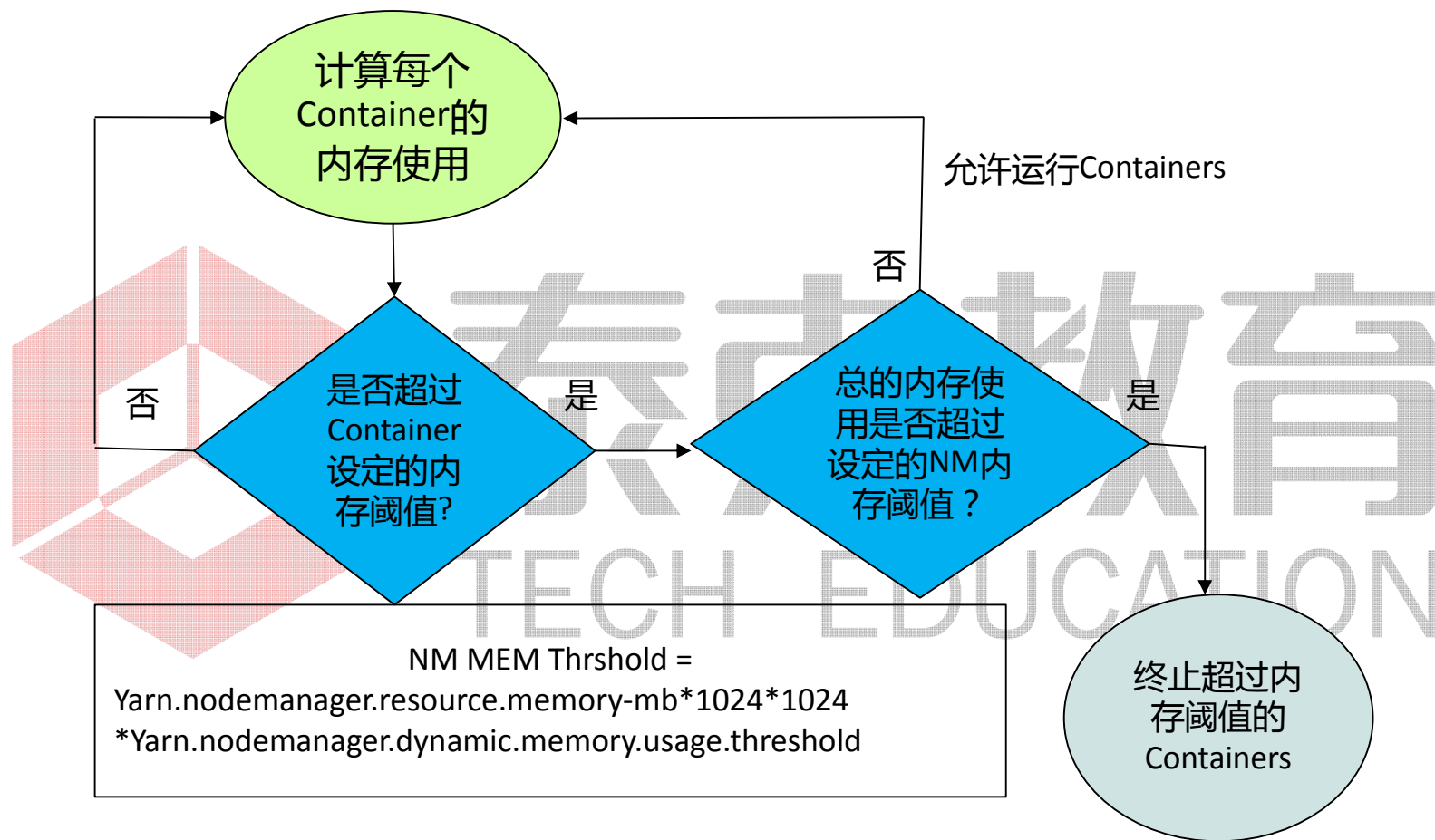
# 目录

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

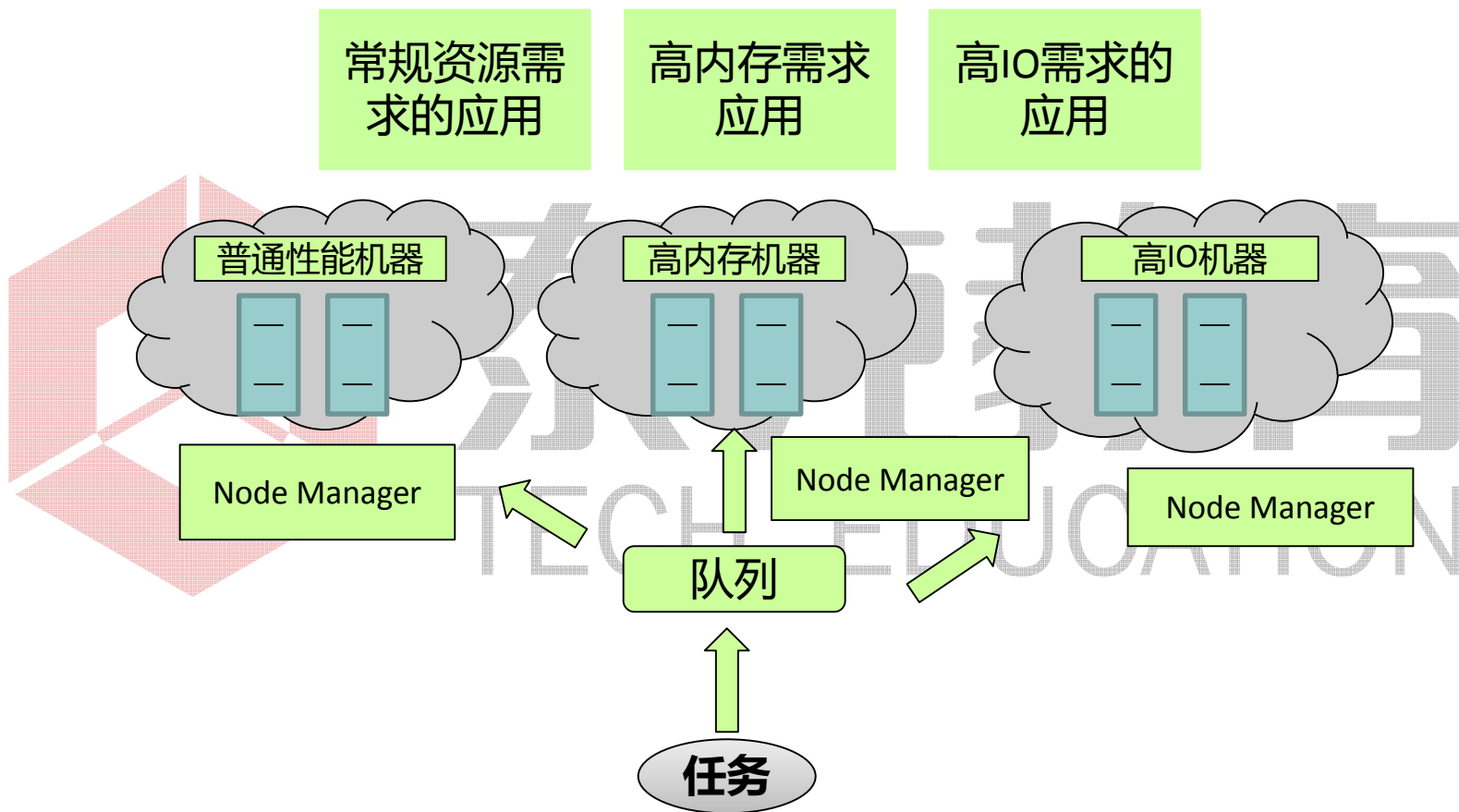
## 4. 增强特性

泰克教育  
TECH EDUCATION

# 增强特性 - YARN动态内存管理



# 增强特性 - YARN基于标签调度

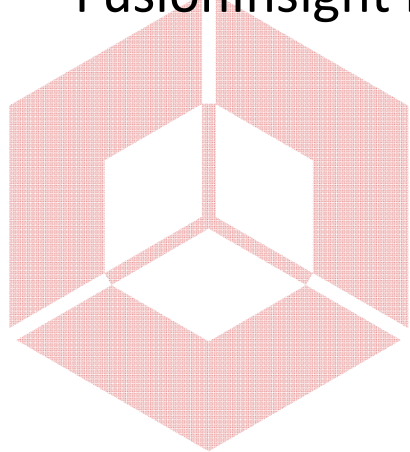






## 本章总结

- 本章首先讲述了MR和YARN的应用场景和基本架构，然后讲解了YARN资源管理与任务调度的原理，最后介绍了华为FusionInsight HD中YARN的增强特性。

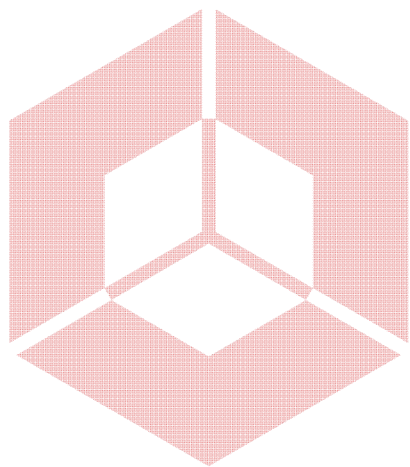


泰克教育  
TECH EDUCATION



## 思考题

1. 请简述MapReduce的工作原理。
2. 请简述YARN的工作原理。



泰克教育  
TECH EDUCATION

## 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

泰克教育  
TECH EDUCATION

## 思考题

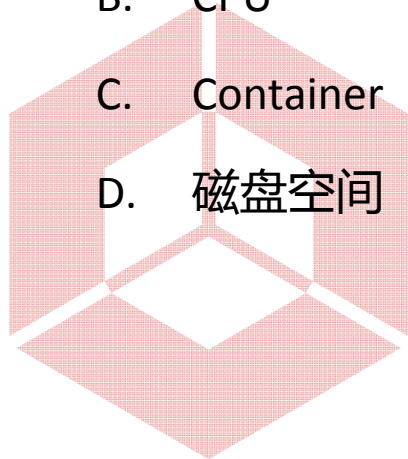
2. YARN中资源抽象用什么表示？（ ）

A. 内存

B. CPU

C. Container

D. 磁盘空间



泰克教育  
TECH EDUCATION

## 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

泰克教育  
TECH EDUCATION

## 思考题

4. 容量调度器有哪些特点？（ ）
- A. 容量保证
  - B. 灵活性
  - C. 多重租赁
  - D. 动态更新配置文件

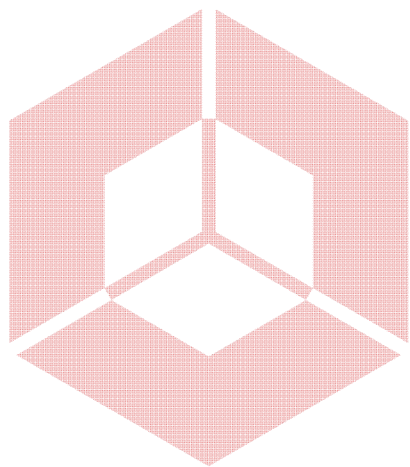
泰克教育  
TECH EDUCATION



## 更多信息

- 下载培训资料：
  - <http://support.huawei.com/learning/trainFaceDetailAction?lang=zh&pbiPath=term1000025185&courseId=Node1000009072>
- eLearning课程：
  - <http://support.huawei.com/learning/nodeQueryAction!loadTrainProjectInfo?lang=zh&pbiPath=term1000025185&courseId=Node1000009421&navId=MW000001>
- 考试大纲：
  - <http://support.huawei.com/learning/Certificate!toExamOutlineDetail?lang=zh&nodeId=Node1000003516>
- 模拟考试：
  - <http://support.huawei.com/learning/Certificate!toSimExamDetail?lang=zh&nodeId=Node1000004285>
- 认证流程：
  - [http://support.huawei.com/learning/NavigationAction!createNavi#navi\[id\]=\\_40](http://support.huawei.com/learning/NavigationAction!createNavi#navi[id]=_40)





谢谢

[www.huawei.com](http://www.huawei.com)

泰克教育  
TECH EDUCATION