

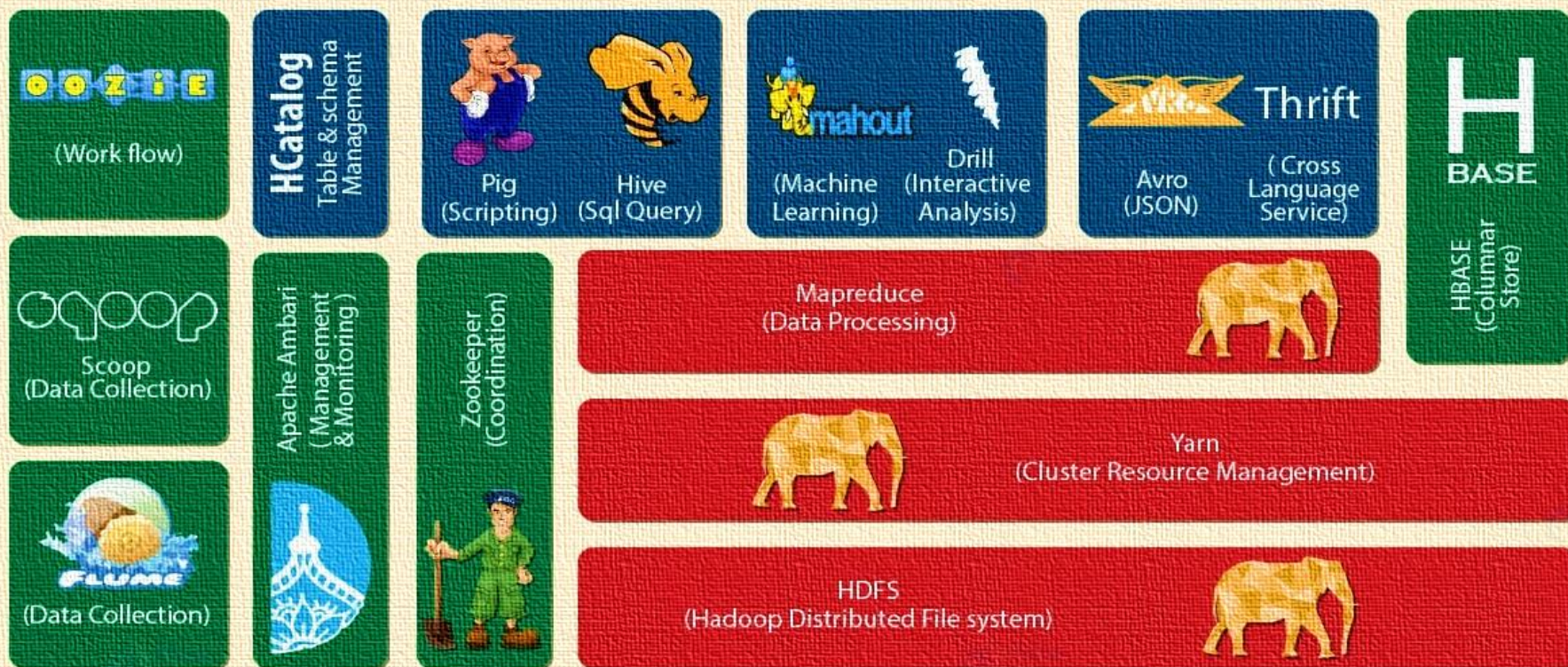


# 大数据分析技术

## Chap. 2 大数据处理架构Hadoop

王怡洋 副教授

大连海事大学 信息科学技术学院







# 内容提纲

Chap. 2.1 Hadoop 概述

Chap. 2.2 生态系统

Chap. 2.3 的安装与使用

本PPT是基于如下教材的配套讲义：

《大数据技术原理与应用——概念、存储、处理、分析与应用》

（2017年2月第2版）林子雨 编著，人民邮电出版社





## 2.1 概述

- 2.1.1 Hadoop简介
- 2.1.2 Hadoop发展简史
- 2.1.3 Hadoop的特性
- 2.1.4 Hadoop的应用现状

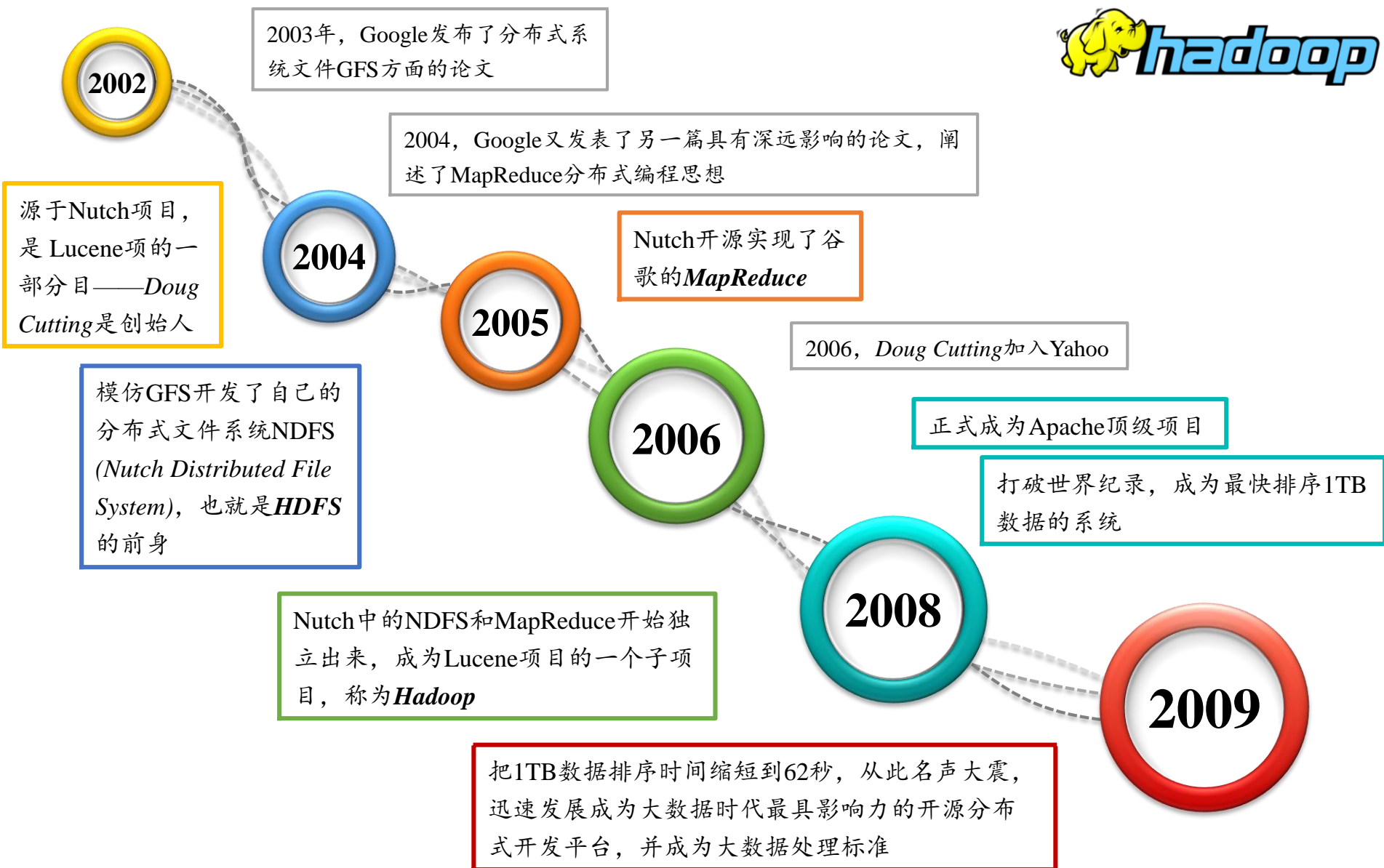


## 2.1.1 Hadoop简介

- 是*Apache*软件基金会旗下的一个开源分布式计算平台，为用户提供了系统底层细节透明的分布式基础架构
- 是基于*Java*语言开发的，具有很好的跨平台特性，并且可以部署在廉价的计算机集群中
- 其核心是分布式文件系统*HDFS* (Hadoop Distributed File System) 和分布式并行编程框架*MapReduce*
- *HDFS*是针对谷歌文件系统 (Google File System, *GFS*)的开源实现；*MapReduce*是针对谷歌*MapReduce*的开源实现
- 被公认为行业大数据标准开源软件，在分布式环境下提供了海量数据的处理能力



## 2.1.2 Hadoop发展简史







## 2.1.3 Hadoop的特性

Hadoop是一个能够对大量数据进行分布式处理的软件框架，并且是以一种可靠、高效、可伸缩的方式进行处理，它具有以下几个方面的特性：

- 高可靠性

采用冗余数据存储方式，既是一个副本发生故障，其他副本也可以保证正常对外提供服务

- 高效性

采用分布式存储和分布式处理两大核心技术，能够高效地处理PB级数据

- 高可扩展性

其设计目标是可以高效稳定地运行在廉价的计算机集群上，可以扩展到数以千计的计算机节点上

- 高容错性

采用冗余数据存储方式，自动保存数据的多个副本，并且能够自动将失败的任务进行重新分配

- 成本低

采用廉价的计算机集群，成本比较低，普通用户也很容易用自己的PC搭建Hadoop运行环境

- 运行在Linux平台上

基于Java语言开发的，可以较好地运行在Linux平台上

- 支持多种编程语言

其上的应用程序也可以使用其他语言编写，如C++



## 2.1.3 Hadoop的应用现状



- Hadoop凭借其突出的优势，已经在各个领域得到了广泛的应用，而互联网领域是其应用的主阵地



•2007年，雅虎在Sunnyvale总部建立了M45——一个包含了4000个处理器和1.5PB容量的Hadoop集群系统

•Facebook作为全球知名的社交网站，Hadoop是非常理想的选择，Facebook主要将Hadoop平台用于日志处理、推荐系统和数据仓库等方面

•国内采用Hadoop的公司主要有百度、淘宝、网易、华为、中国移动等，其中，淘宝的Hadoop集群比较大



## 2.1.3 Hadoop的应用现状

Hadoop在企业中的应用架构





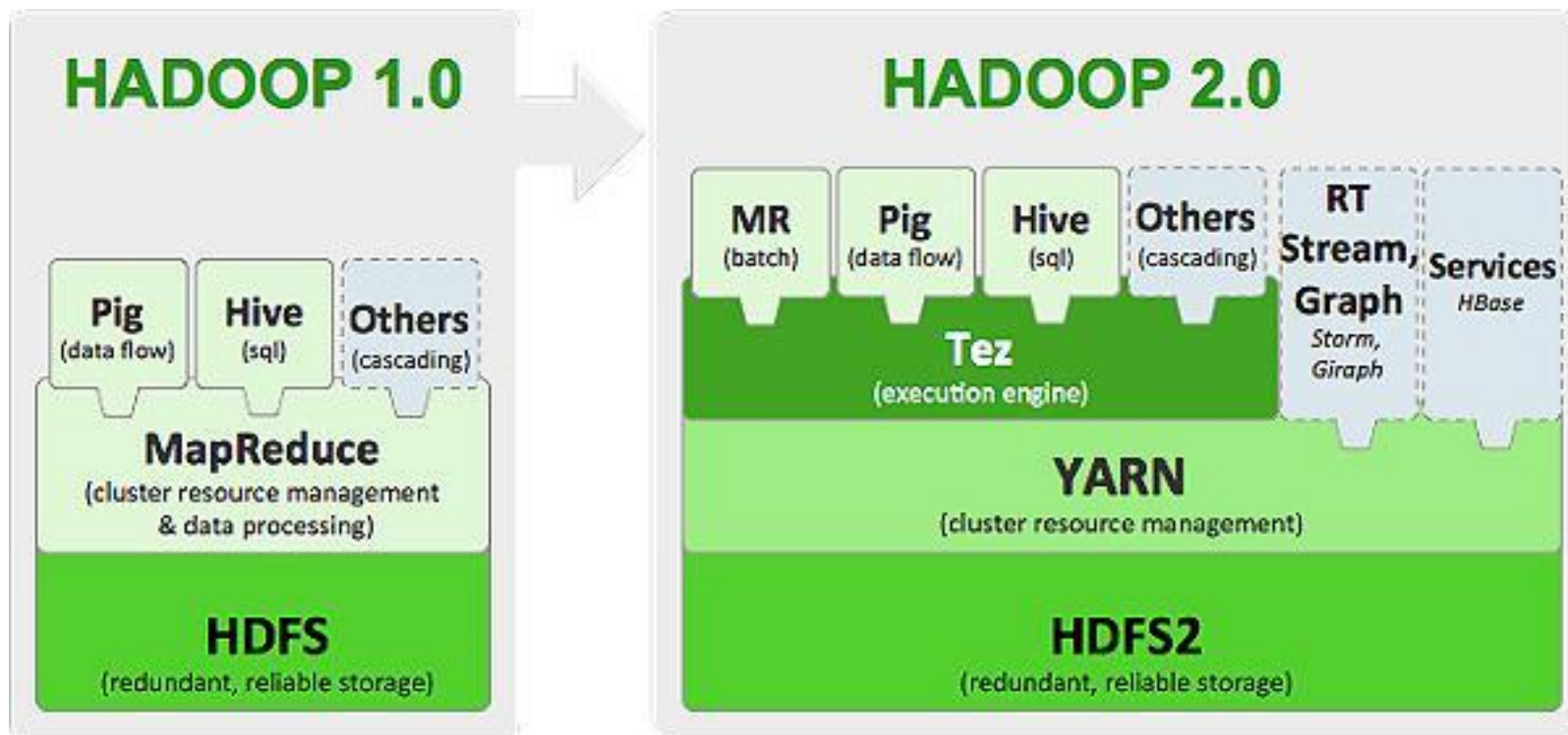


## 2.1.4 Apache Hadoop版本演变

- Apache Hadoop版本分为两代，我们将第一代Hadoop称为Hadoop 1.0，第二代Hadoop称为Hadoop 2.0
- 第一代Hadoop包含三个大版本，分别是0.20.x，0.21.x和0.22.x，其中，0.20.x最后演化成1.0.x，变成了稳定版，而0.21.x和0.22.x则增加了NameNode HA等新的重大特性
- 第二代Hadoop包含两个版本，分别是0.23.x和2.x，它们完全不同于Hadoop 1.0，是一套全新的架构，均包含HDFS Federation和YARN两个系统，相比于0.23.x，2.x增加了NameNode HA和Wire-compatibility两个重大特性



## 2.1.4 Apache Hadoop版本演变



- 拆分出YARN专门做资源调度管理
- MapReduce只做数据处理，提高了整体的效率



## 2.1.5 Hadoop各种版本

- Apache Hadoop
- Hortonworks
- Cloudera (CDH: Cloudera Distribution Hadoop)
- MapR
- .....

选择 Hadoop版本的考虑因素:

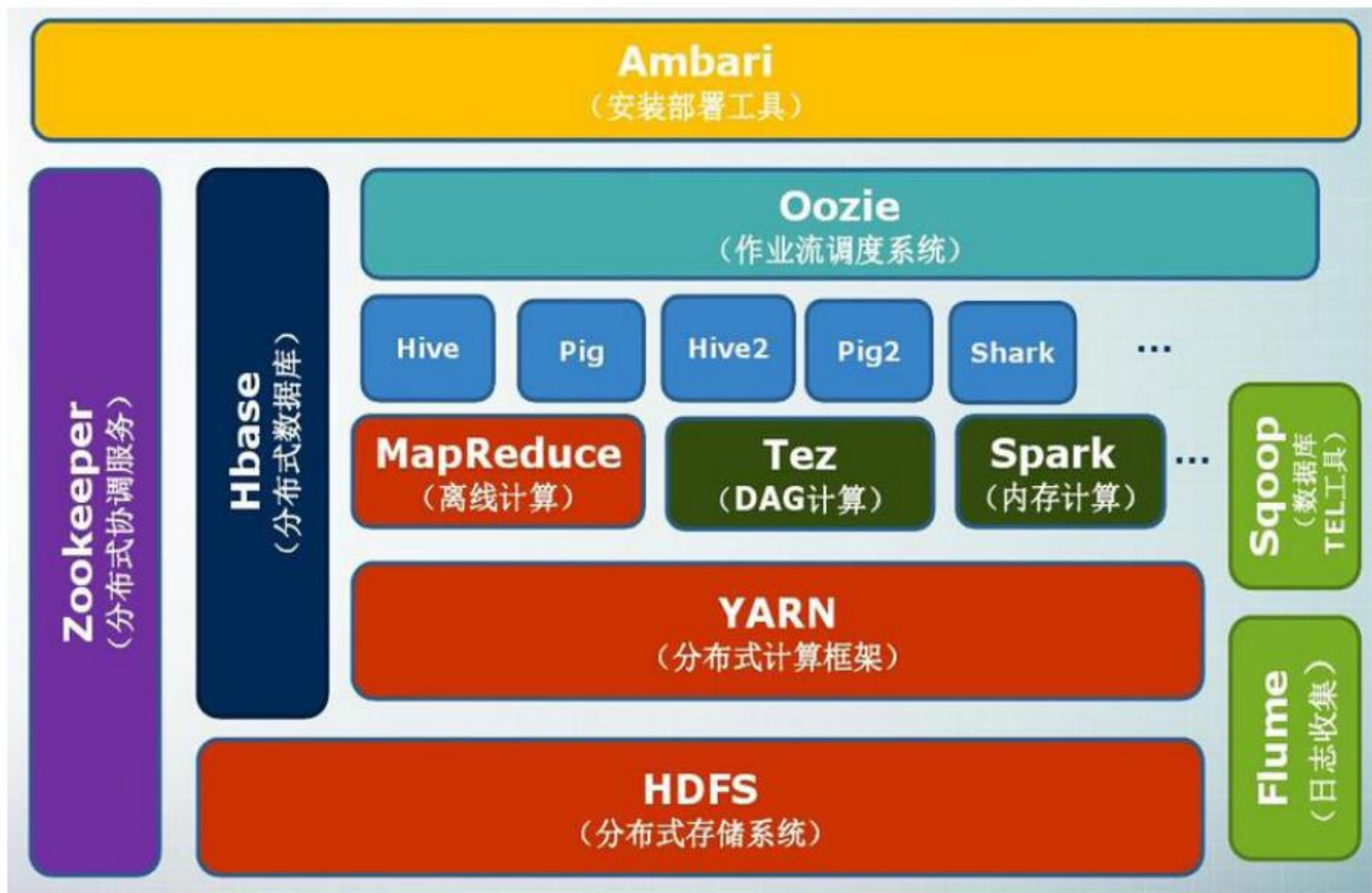
- 是否开源 (即是否免费)
- 是否有稳定版
- 是否经实践检验
- 是否有强大的社区支持





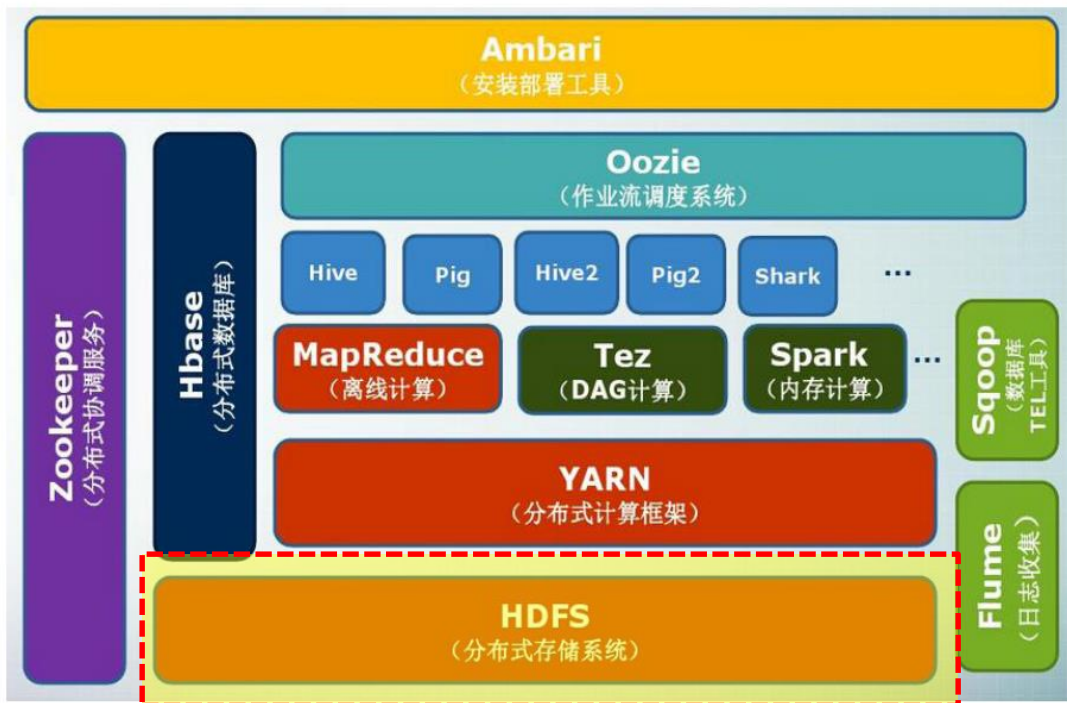
## 2.2 Hadoop项目结构

Hadoop的项目结构不断丰富发展，已经形成一个丰富的Hadoop生态系统





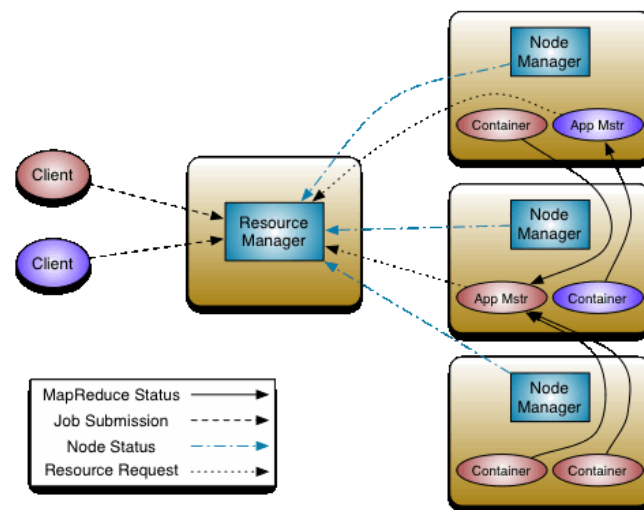
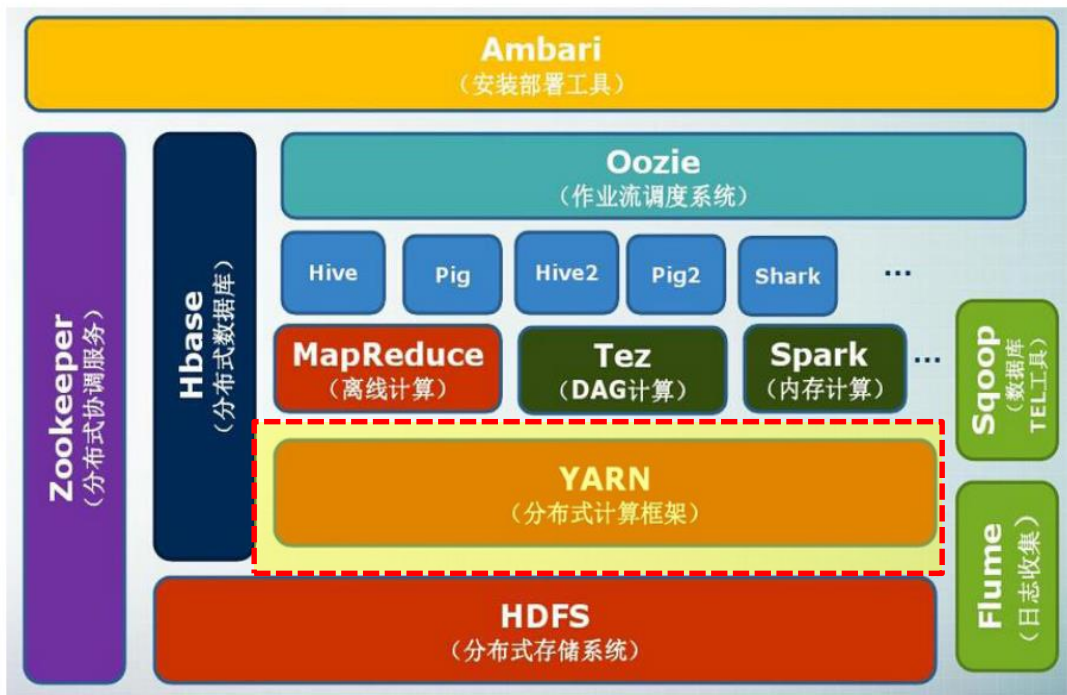
## 2.2 Hadoop项目结构



- **HDFS**: 分布式文件系统 (Hadoop Distributed File System), 是针对谷歌文件系统 (GFS) 的开源实现。



## 2.2 Hadoop项目结构

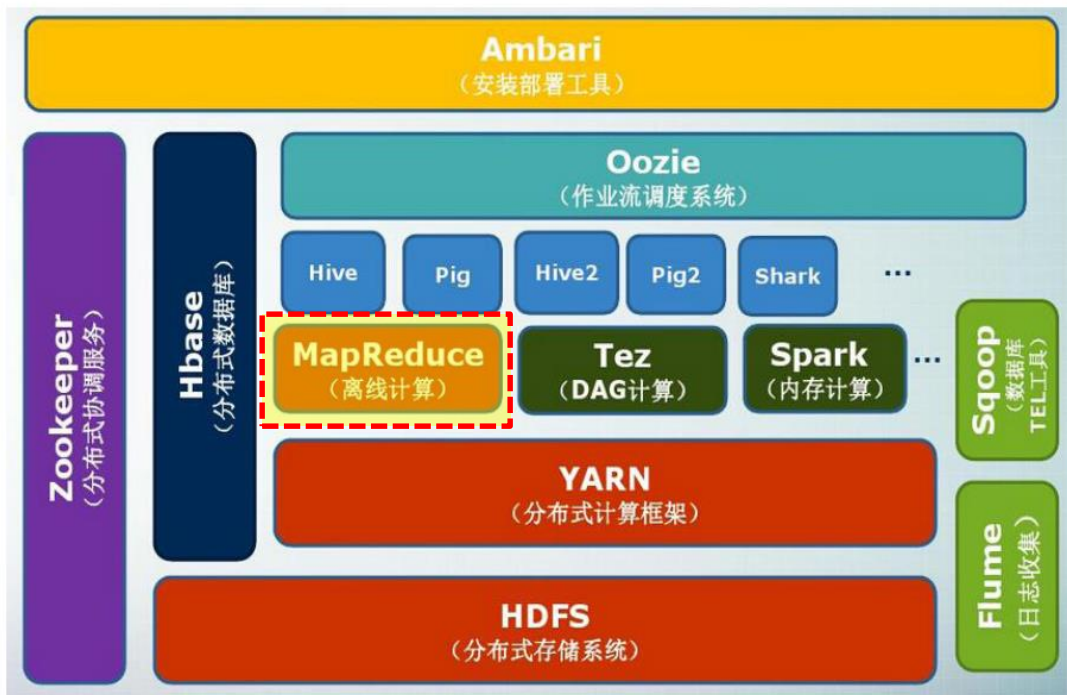


- YARN: 分布式计算框架 (Yet Another Resource Negotiator), 负责资源调度





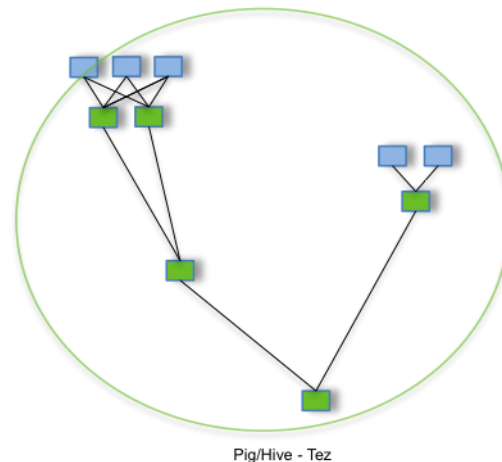
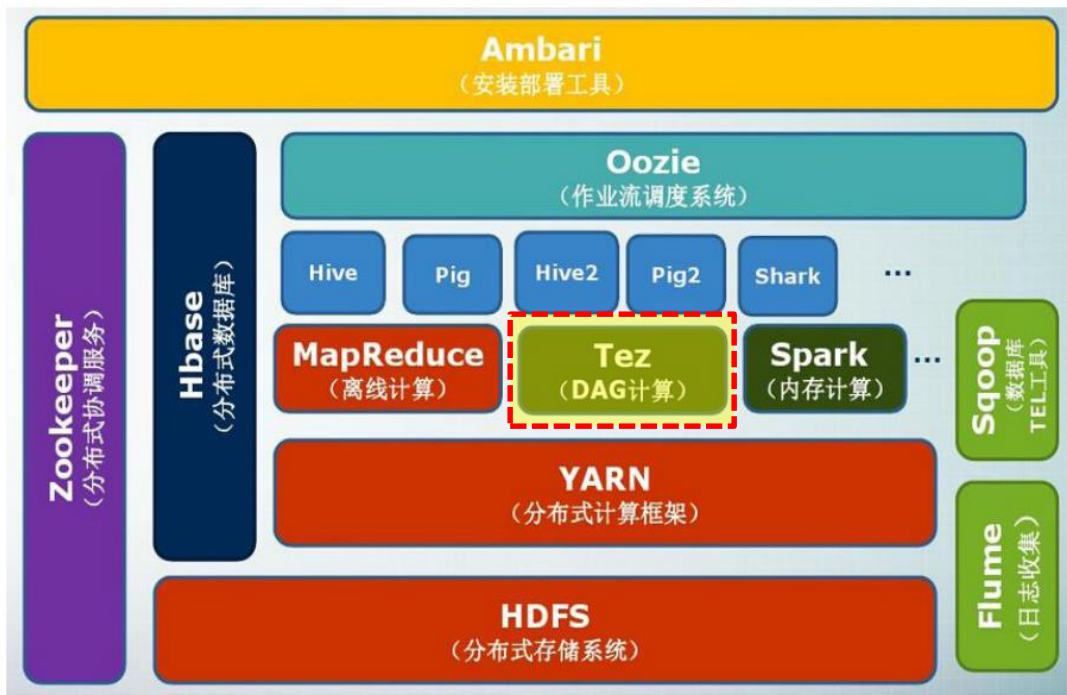
## 2.2 Hadoop项目结构



- MapReduce: 分布式编程模型，批量处理数据、采用离线的方式



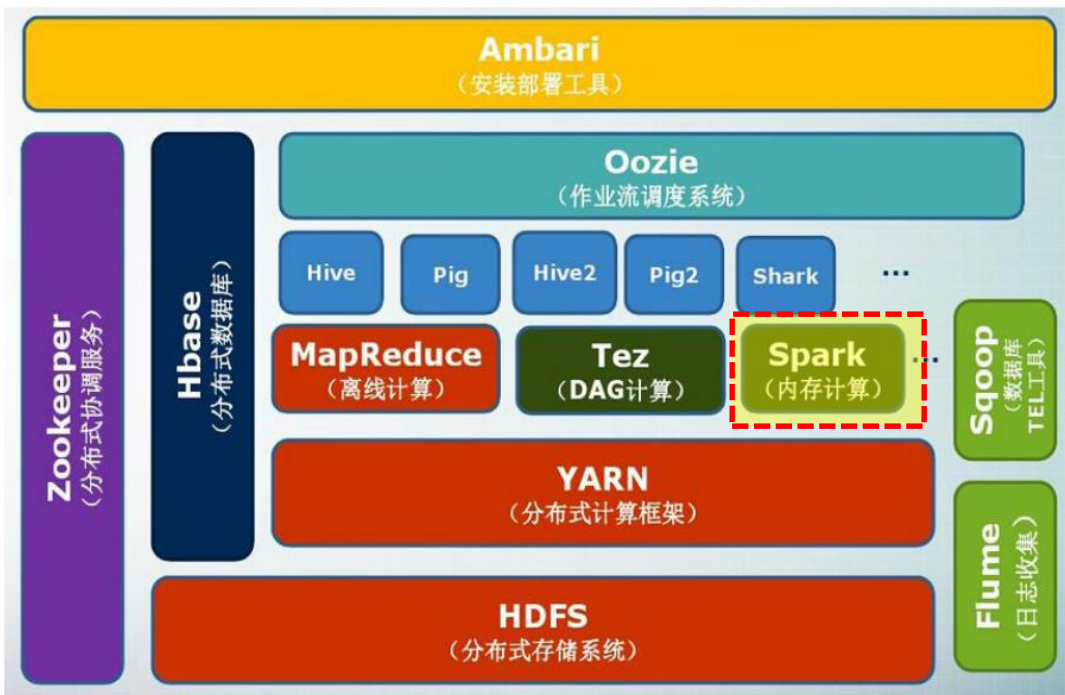
## 2.2 Hadoop项目结构



- **Tez**: 运行在YARN上的下一代Hadoop查询处理框架。它会把你很多的MapReduce作业进行分析优化之后，构成一个有向无环图，可以保证你获得最好的处理效率。



## 2.2 Hadoop项目结构

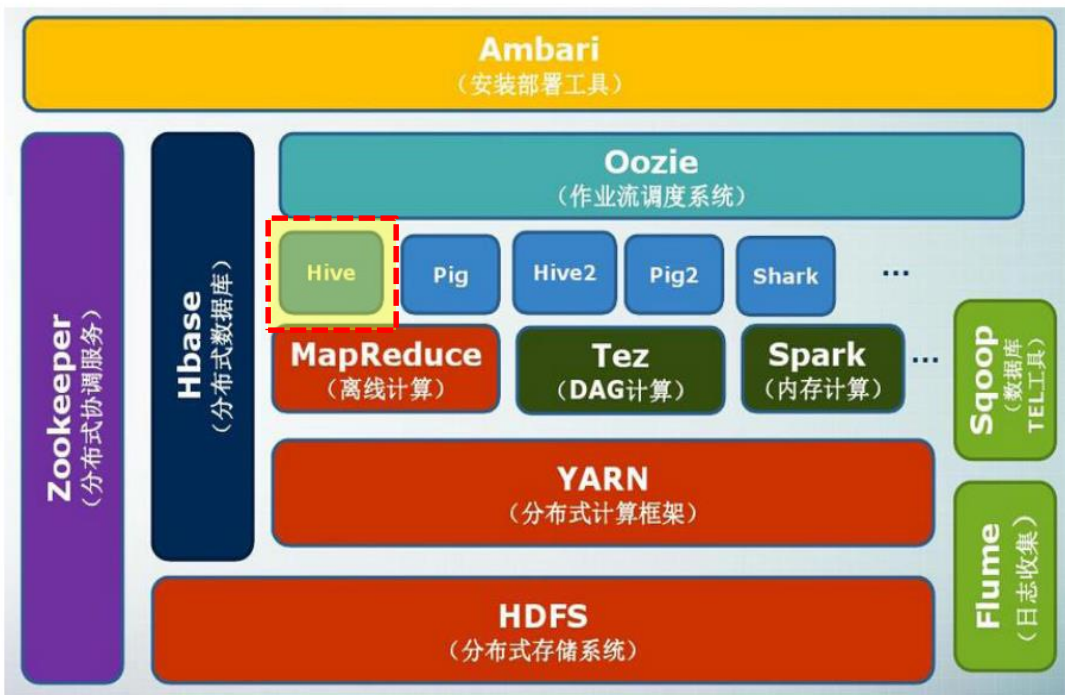


- **Spark**: 它的逻辑和MapReduce是一样的，但是它是基于内存计算，而MapReduce是基于磁盘的。Spark是直接内存中完成整个数据的处理，所以性能更高。





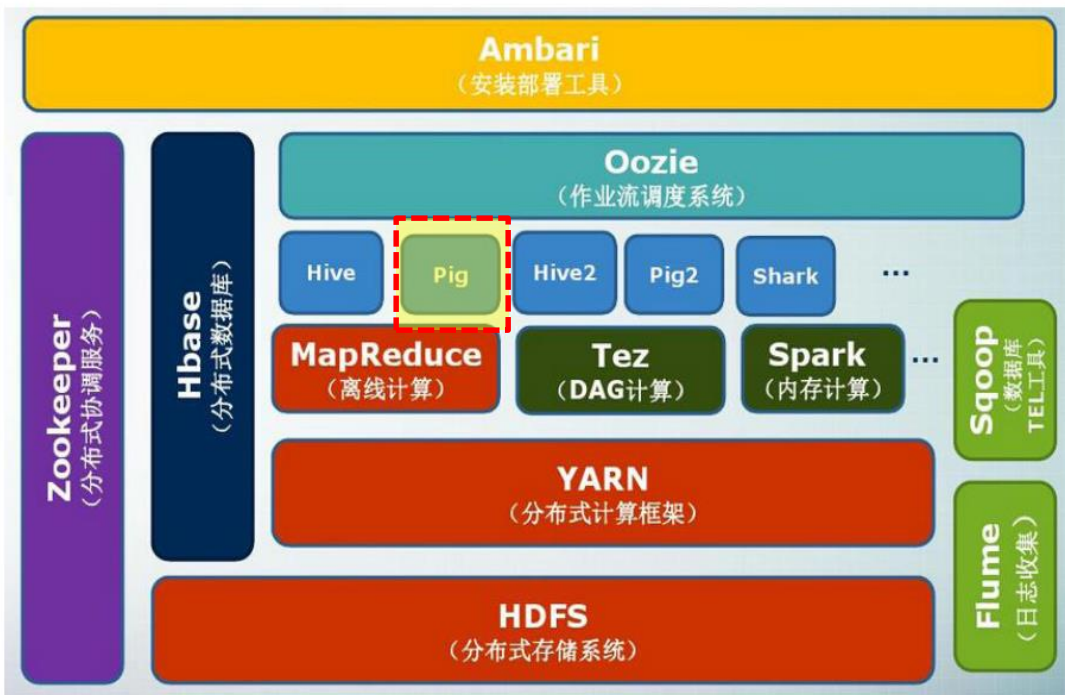
## 2.2 Hadoop项目结构



- Hive: 基于Hadoop的数据仓库，可用于对Hadoop文件中的数据集进行数据整理、特殊查询和分析存储。它提供了类似于关系数据库SQL语言的查询语言，Hive QL，可以快速实现简单的MapReduce统计，也可转换为MapReduce任务进行运行。



## 2.2 Hadoop项目结构

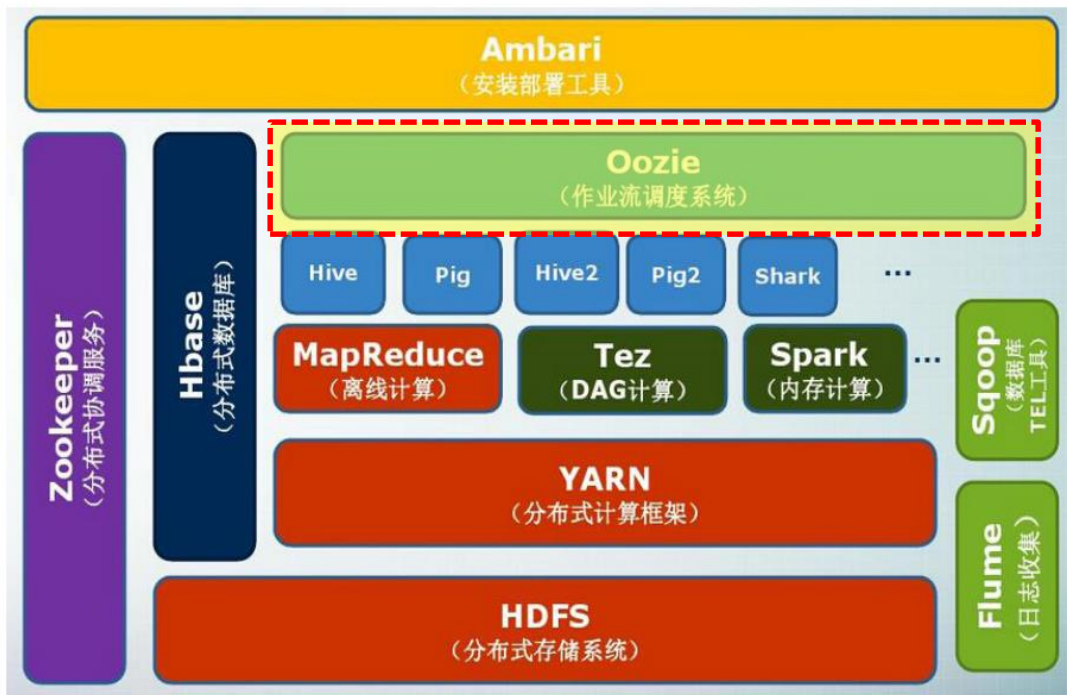


# Apache Pig

- Pig: 一种数据流语言和运行环境，适合于使用Hadoop和MapReduce平台来查询大型半结构化数据集。提供类似SQL的查询语言Pig Latin。



## 2.2 Hadoop项目结构

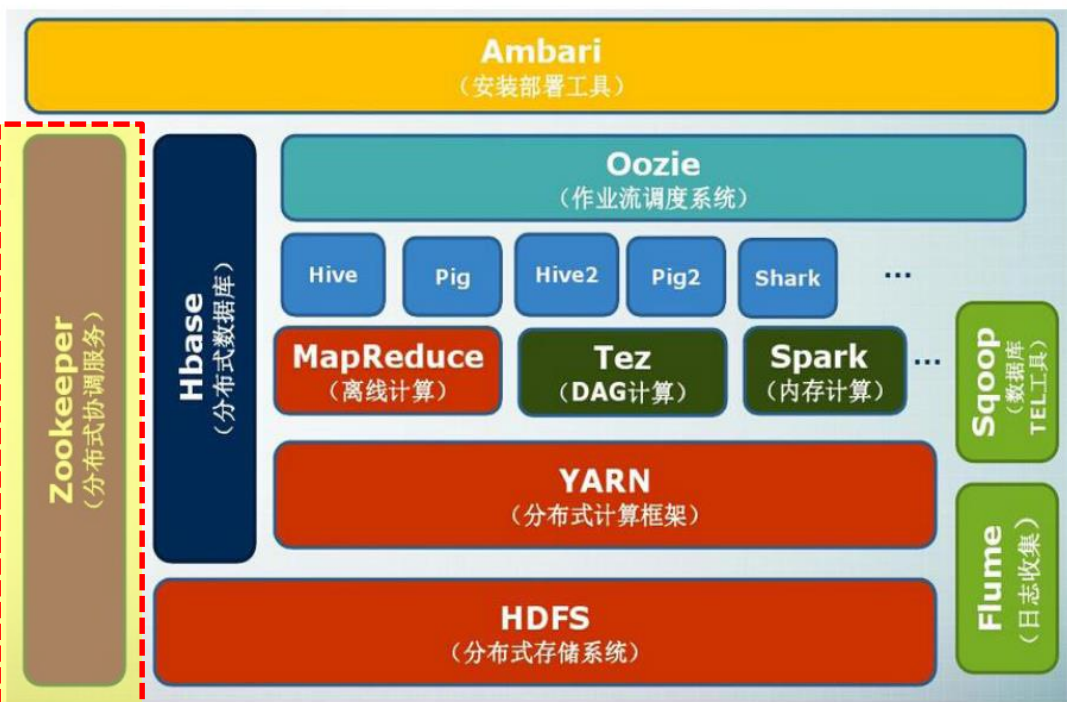


- Oozie: Hadoop上的工作流调度、管理系统。





## 2.2 Hadoop项目结构

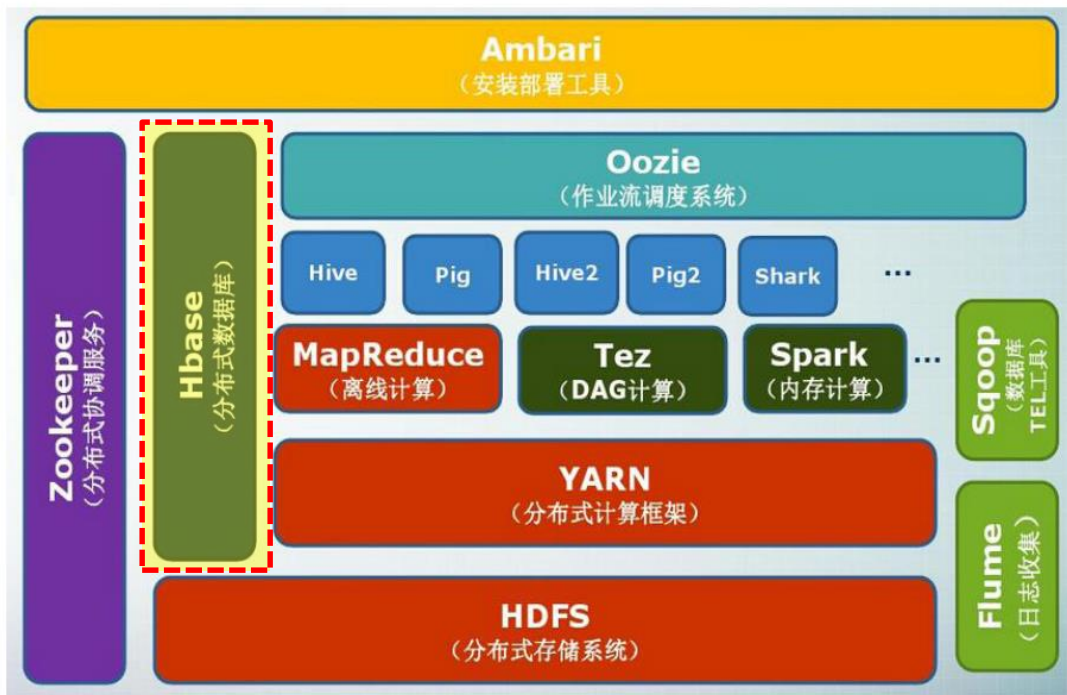


Apache ZooKeeper™

- Zookeeper: 是针对Google Chubby的一个开源实现, 是高效可靠的协同工作系统。



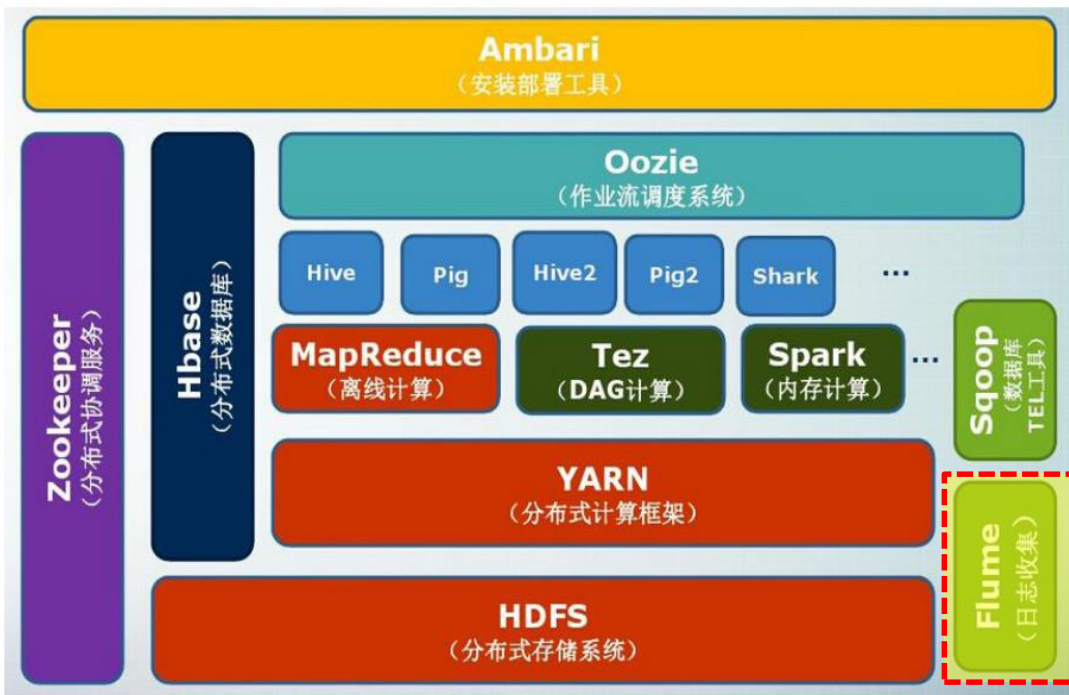
## 2.2 Hadoop项目结构



- **Hbase**: 提供可靠性、高性能、可伸缩、实时读写、分布式的列式数据库，一般采用HDFS作为其底层数据存储。它是针对Google BigTable的开源实现，具有强大的非结构化数据存储能力。



## 2.2 Hadoop项目结构

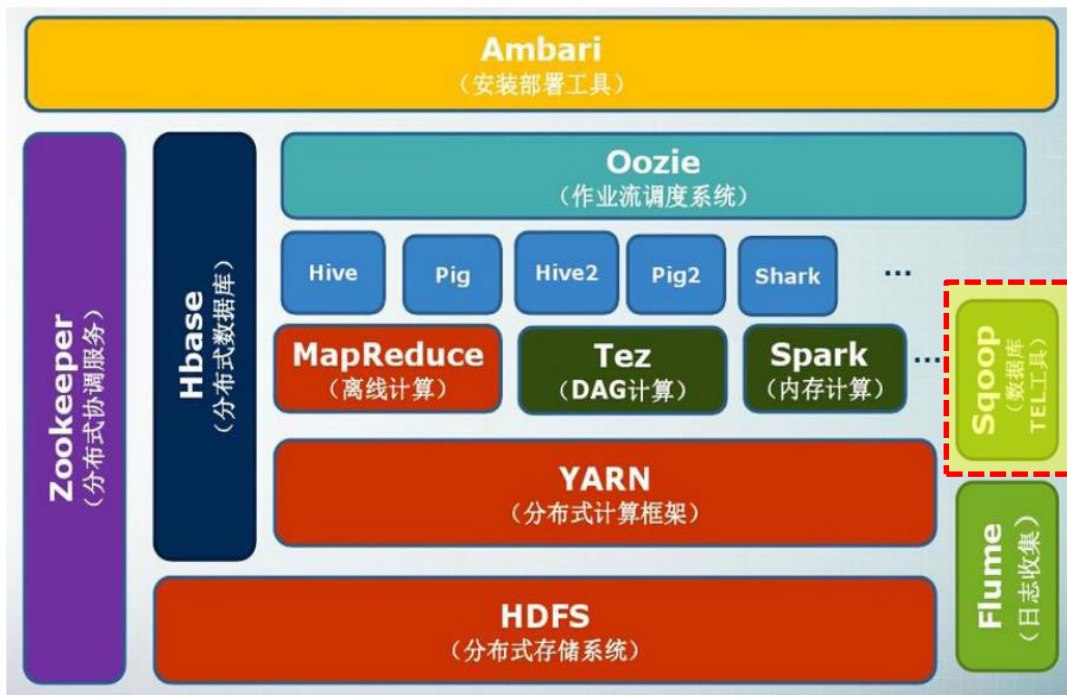


- Flume: 海量日志采集、聚合和传输的系统。





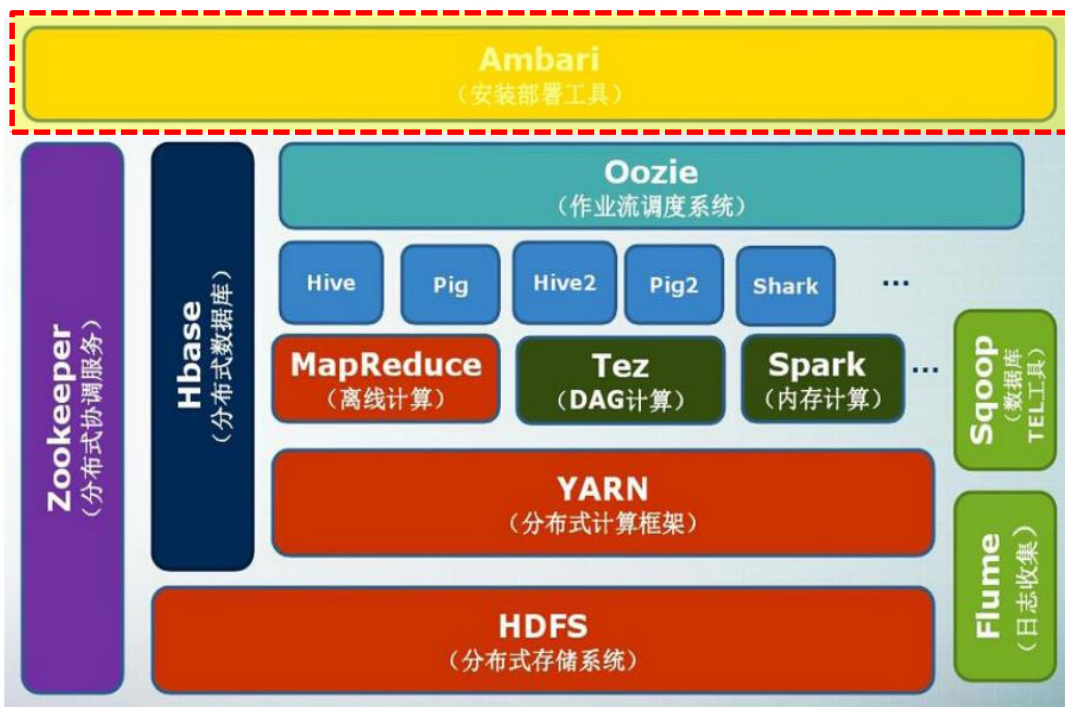
## 2.2 Hadoop项目结构



- **Sqoop**: SQL-to-Hadoop的缩写，主要用来在Hadoop和关系数据库之间交换数据，。通过Sqoop可以方便地将数据从MySQL、Oracle等关系数据库中的数据中导入Hadoop（HDFS、Hbase或Hive）；或者将数据从Hadoop导出到关系数据库，使得数据迁移变得非常方便。



## 2.2 Hadoop项目结构



- Ambari: 基于Web的工具，支持Hadoop集群的安装、部署、配置和管理。



## 2.3 Hadoop的安装与使用

- 2.3.1 Hadoop安装之前的预备知识
- 2.3.2 安装Linux虚拟机
- 2.3.3 安装双操作系统
- 2.3.4 详解Hadoop的安装与使用

➤ 更多细节内容，可以参考本书配套学习指南





## 2.3.1 Hadoop安装之前的预备知识

### (一) Linux的选择

#### (1) 选择哪个Linux发行版？

- 在Linux系统各个发行版中，CentOS系统和Ubuntu系统在服务端和桌面端使用占比最高，网络上资料最是齐全，所以建议使用CentOS 或Ubuntu
- 在学习Hadoop方面，虽然两个系统没有多大区别，但是**推荐使用Ubuntu操作系统**

#### (2) 选择32位还是64位？

- 如果电脑比较老或者内存小于2G，那么建议选择32位系统版本的Linux
- 如果内存大于4G，那么建议选择64位系统版本的Linux



## 2.3.1 Hadoop安装之前的预备知识

### (二) 系统安装方式：选择虚拟机安装还是双系统安装

- 建议电脑比较新或者配置内存4G以上的电脑可以选择虚拟机安装
- 电脑较旧或配置内存小于等于4G的电脑强烈建议选择双系统安装，否则，在配置较低的计算机上运行Linux虚拟机，系统运行速度会非常慢
- 鉴于目前教师和学生的计算机硬件配置一般不高，建议在实践教学中采用双系统安装，确保系统运行速度



## 2.3.1 Hadoop安装之前的预备知识

### (三) 关于Linux的一些基础知识

- Shell
  - 是指“提供使用者使用界面”的软件（命令解析器），类似于DOS下的command和后来的cmd.exe。它接收用户命令，然后调用相应的应用程序
- sudo命令
  - sudo是ubuntu中一种权限管理机制，管理员可以授权给一些普通用户去执行一些需要root权限执行的操作。当使用sudo命令时，就需要输入您当前用户的密码
- 输入密码
  - 在Linux的终端中输入密码，终端是不会显示任何你当前输入的密码，也不会提示你已经输入了多少字符密码，读者不要误以为键盘没有响应
- 输入法中英文切换
  - linux中英文的切换方式是使用键盘“shift”键来切换，也可以点击顶部菜单的输入法按钮进行切换。Ubuntu自带的Sunpinyin中文输入法已经足够读者使用
- Ubuntu终端复制粘贴快捷键
  - 在Ubuntu终端窗口中，复制粘贴的快捷键需要加上shift，即粘贴是ctrl+shift+v





## 2.3.1 Hadoop安装之前的预备知识

### (四) Hadoop安装方式

- 单机模式：Hadoop 默认模式为非分布式模式（本地模式），无需进行其他配置即可运行。非分布式即单 Java 进程，方便进行调试
- 伪分布式模式：Hadoop 可以在单节点上以伪分布式的方式运行，Hadoop 进程以分离的 Java 进程来运行，节点既作为 NameNode 也作为 DataNode，同时，读取的是 HDFS 中的文件
- 分布式模式：使用多个节点构成集群环境来运行Hadoop



## 2.3.2 安装Linux虚拟机

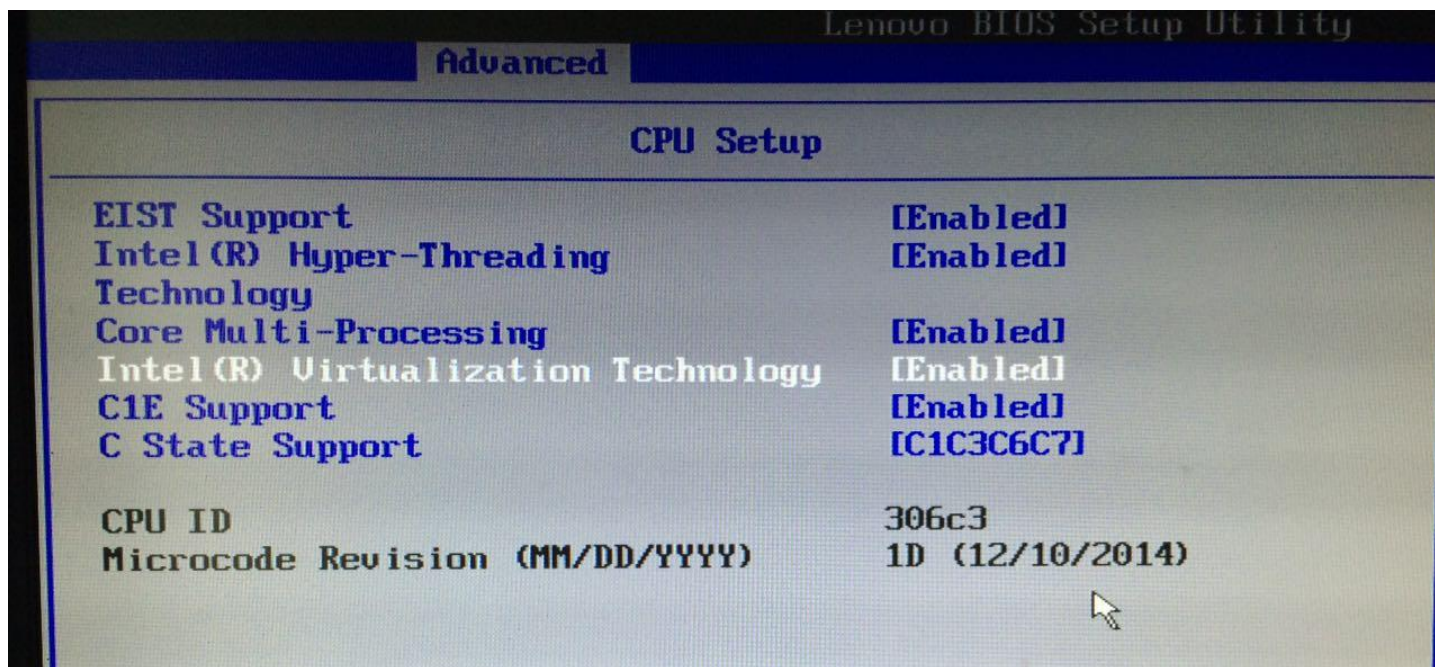
### 一、材料和工具

- 1、下载VirtualBox虚拟机软件
- 2、下载Ubuntu LTS 14.04 ISO映像文件

### 二、步骤

#### (一) 确认系统版本

如果选择的系统是64位Ubuntu系统，那么在安装虚拟机前，我们还要进入BIOS开启CPU的虚拟化

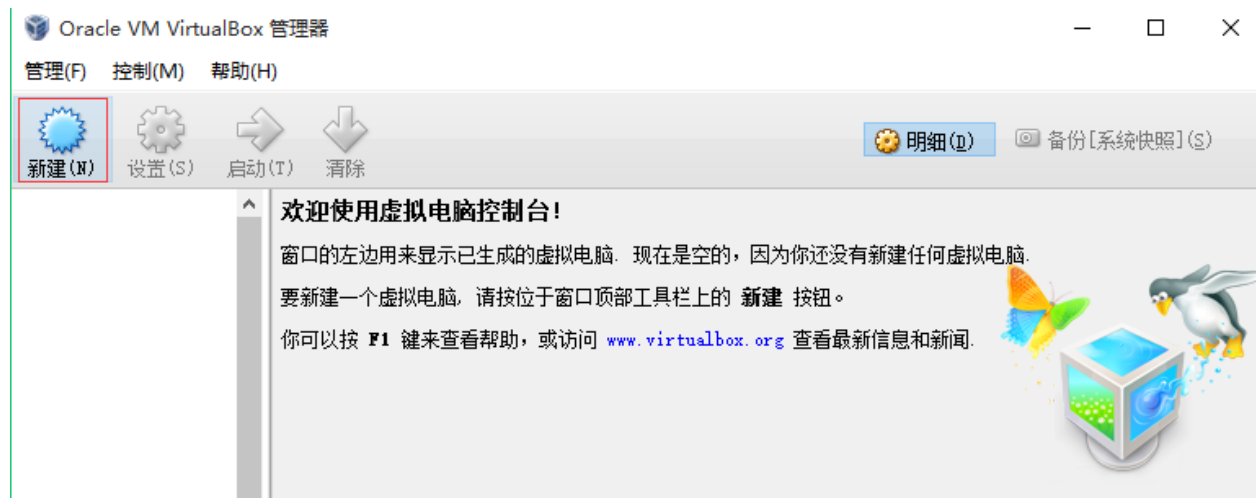




## 2.3.2 安装Linux虚拟机

### (二) 安装前的准备

1. 打开VirtualBox，点击“创建”按钮，创建一个虚拟机
2. 给虚拟机命名，选择操作系统，版本
3. 选择内存大小，这里设置的1024M
4. 创建虚拟硬盘
5. 选择虚拟硬盘文件类型VDI
6. 虚拟硬盘选择动态分配
7. 选择文件存储的位置和容量大小
8. 点击创建

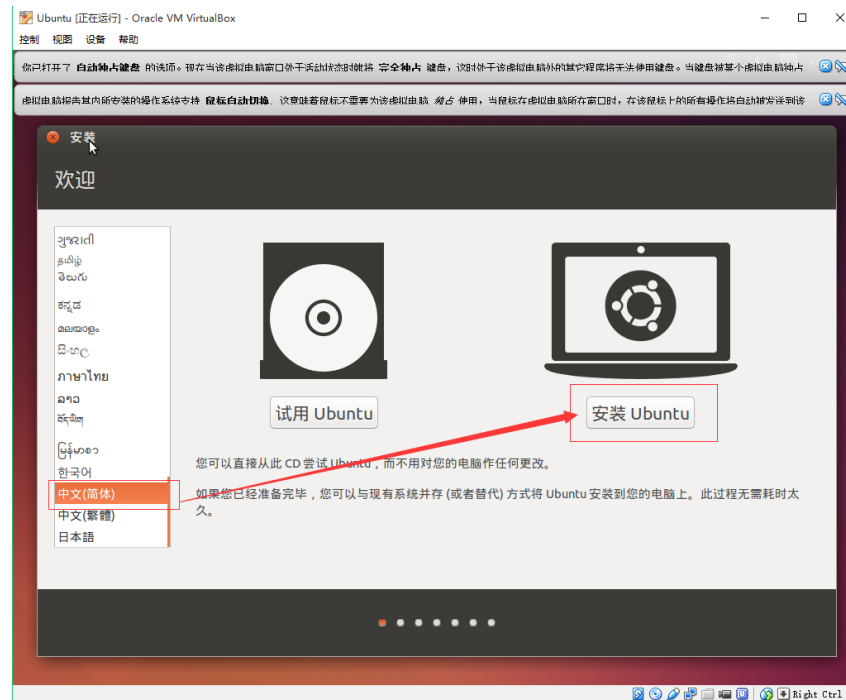
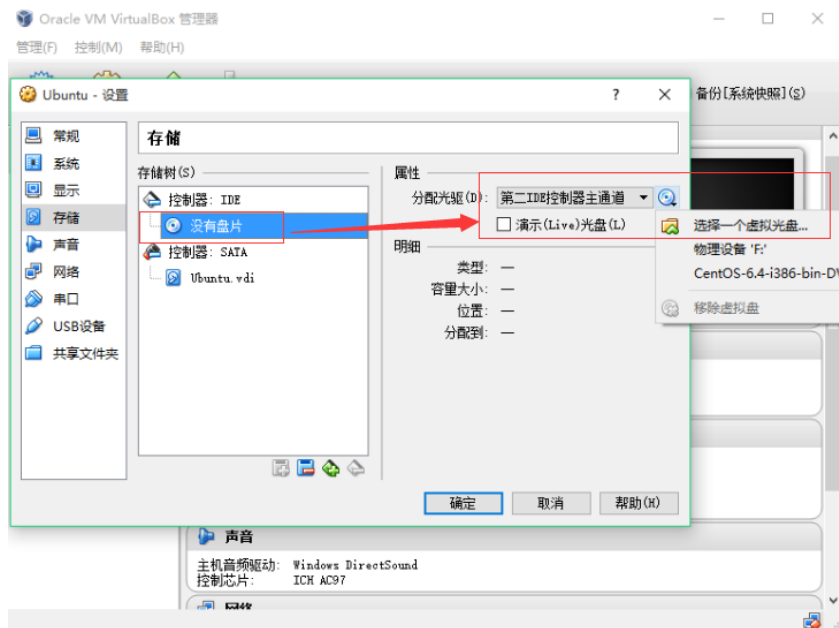






## 2.3.2 安装Linux虚拟机

### (三)安装Ubuntu





## 2.3.3 安装双操作系统

- 第一步：制作安装U盘
- 具体可参考百度经验文章
- <http://jingyan.baidu.com/article/59703552e0a6e18fc007409f.html>
- 第二步：双系统安装
- 具体可参考百度经验文章
- <http://jingyan.baidu.com/article/dca1fa6fa3b905f1a44052bd.html>

安装后Window和Ubuntu 14.04都可以用，默认windows优先启动  
可以在电脑启动时，选择进入Ubuntu系统而不是 Windows系统



## 2.3.4 Hadoop的安装与使用（单机/伪分布式）

Hadoop基本安装配置主要包括以下几个步骤：

- 创建Hadoop用户
- SSH登录权限设置
- 安装Java环境
- 单机安装配置
- 伪分布式安装配置

详细安装配置过程请参考本书配套指南教程

《**Hadoop安装教程\_单机/伪分布式配置\_Hadoop2.6.0/Ubuntu14.04**》

<http://dblab.xmu.edu.cn/blog/install-hadoop/>





# 创建Hadoop用户

如果安装 Ubuntu 的时候不是用的 “hadoop” 用户，那么需要增加一个名为 hadoop 的用户：

首先按 **ctrl+alt+t** 打开终端窗口，输入如下命令创建新用户：

```
$ sudo useradd -m hadoop -s /bin/bash
```

上面这条命令创建了可以登陆的 hadoop 用户，并使用 /bin/bash 作为 shell

接着使用如下命令设置密码，可简单设置为 hadoop，按提示输入两次密码：

```
$ sudo passwd hadoop
```

可为 hadoop 用户增加管理员权限，方便部署，避免一些对新手来说比较棘手的权限问题：

```
$ sudo adduser hadoop sudo
```



# SSH登录权限设置

## SSH是什么？

SSH 为 Secure Shell 的缩写，是建立在应用层和传输层基础上的安全协议。SSH 是目前较可靠、专为远程登录会话和其他网络服务提供安全性的协议。利用 SSH 协议可以有效防止远程管理过程中的信息泄露问题。SSH最初是UNIX系统上的一个程序，后来又迅速扩展到其他操作平台。SSH是由[客户端](#)和[服务端](#)的软件组成，服务端是一个守护进程(daemon)，它在后台运行并响应来自客户端的连接请求，客户端包含ssh程序以及像scp（远程拷贝）、slogin（远程登陆）、sftp（安全文件传输）等其他的应用程序

## 配置SSH的原因：

Hadoop名称节点（NameNode）需要启动集群中所有机器的Hadoop守护进程，这个过程需要通过SSH登录来实现。Hadoop并没有提供SSH输入密码登录的形式，因此，为了能够顺利登录每台机器，需要将所有机器配置为名称节点可以无密码登录它们



# 安装Java环境

- Java环境可选择 Oracle 的 JDK，或是 OpenJDK
- 可以在Ubuntu中直接通过命令安装 OpenJDK 7

```
$ sudo apt-get install openjdk-7-jre openjdk-7-jdk
```

- 还需要配置一下 JAVA\_HOME 环境变量
- 具体请参考网络教程：<http://dblab.xmu.edu.cn/blog/install-hadoop/>





# 单机安装配置

## Hadoop 2 安装文件的下载

Hadoop 2 可以到官网下载，需要下载 **hadoop-2.x.y.tar.gz** 这个格式的文件，这是编译好的，另一个包含 src 的则是 Hadoop 源代码，需要进行编译才可使用

- 如果读者是使用虚拟机方式安装Ubuntu系统的用户，请用虚拟机中的Ubuntu自带firefox浏览器访问本指南，再点击下载地址，才能把hadoop文件下载到虚拟机ubuntu中。请不要使用Windows系统下的浏览器下载，文件会被下载到Windows系统中，虚拟机中的Ubuntu无法访问外部Windows系统的文件，造成不必要的麻烦。
- 如果读者是使用双系统方式安装Ubuntu系统的用户，请进去Ubuntu系统，在Ubuntu系统打开firefox浏览器，再点击下载



# 单机安装配置

选择将 Hadoop 安装至 /usr/local/ 中

```
$ sudo tar -zxf ~/下载/hadoop-2.6.0.tar.gz -C /usr/local # 解压到/usr/local中  
$ cd /usr/local/  
$ sudo mv ./hadoop-2.6.0/ ./hadoop # 将文件夹名改为hadoop  
$ sudo chown -R hadoop:hadoop ./hadoop # 修改文件权限
```

Hadoop 解压后即可使用。输入如下命令来检查 Hadoop 是否可用，成功则会显示 Hadoop 版本信息：

```
$ cd /usr/local/hadoop  
$ ./bin/hadoop version
```

Hadoop 默认模式为非分布式模式（本地模式），无需进行其他配置即可运行。



# 伪分布式安装配置

- Hadoop 可以在单节点上以伪分布式的方式运行，Hadoop 进程以分离的 Java 进程来运行，节点既作为 NameNode 也作为 DataNode，同时，读取的是 HDFS 中的文件
- Hadoop 的配置文件位于 `/usr/local/hadoop/etc/hadoop/` 中，伪分布式需要修改2个配置文件 **core-site.xml** 和 **hdfs-site.xml**
- Hadoop的配置文件是 xml 格式，每个配置以声明 property 的 name 和 value 的方式来实现





# 伪分布式安装配置

## 实验步骤:

- 修改配置文件: core-site.xml, hdfs-site.xml, mapred-site.xml
- 初始化文件系统 `hadoop namenode -format`
- 启动所有进程 `start-all.sh`
- 访问web界面, 查看Hadoop信息
- 运行实例



# 伪分布式安装配置

## 修改配置文件 **core-site.xml**

```
<configuration>
  <property>
    <name>hadoop.tmp.dir</name>
    <value>file:/usr/local/hadoop/tmp</value>
    <description>Abase for other temporary directories.</description>
  </property>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://localhost:9000</value>
  </property>
</configuration>
```

- hadoop.tmp.dir表示存放临时数据的目录，即包括NameNode的数据，也包括DataNode的数据。该路径任意指定，只要实际存在该文件夹即可
- name为fs.defaultFS的值，表示hdfs路径的逻辑名称



# 伪分布式安装配置

## 修改配置文件 **hdfs-site.xml**

```
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:/usr/local/hadoop/tmp/dfs/name</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>file:/usr/local/hadoop/tmp/dfs/data</value>
  </property></configuration>
```

- dfs.replication表示副本的数量，伪分布式要设置为1
- dfs.namenode.name.dir表示本地磁盘目录，是存储fsimage文件的地方
- dfs.datanode.data.dir表示本地磁盘目录，HDFS数据存放block的地方



# 伪分布式安装配置

关于三种Shell命令方式的区别：

1. `hadoop fs`
  2. `hadoop dfs`
  3. `hdfs dfs`
- `hadoop fs`适用于任何不同的文件系统，比如本地文件系统和HDFS文件系统
  - `hadoop dfs`只能适用于HDFS文件系统
  - `hdfs dfs`跟`hadoop dfs`的命令作用一样，也只能适用于HDFS文件系统