



大数据知识体系型公开课 《大数据概念、技术与应用》

第3讲 分布式文件系统HDFS

林子雨 博士/助理教授

厦门大学计算机科学系

厦门大学云计算与大数据研究中心

E-mail: ziyulin@xmu.edu.cn ▶▶

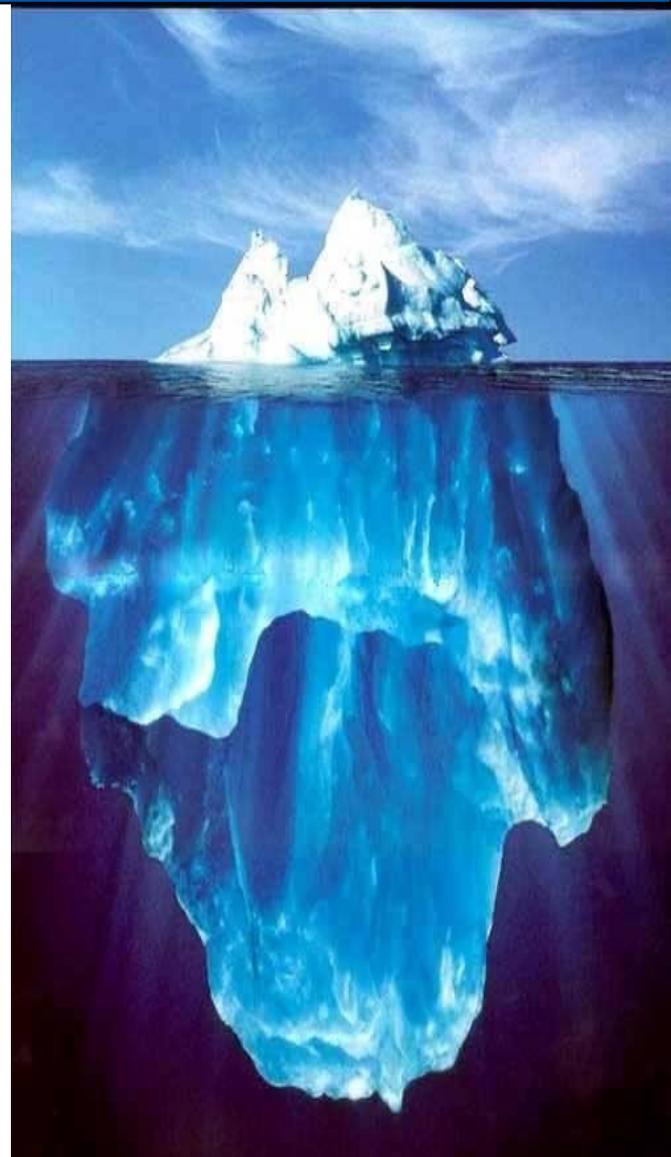
主页: <http://www.cs.xmu.edu.cn/linziyu>





提纲

- 3.1 HDFS简介
- 3.2 HDFS体系结构
- 3.3 HDFS存储原理





3.1 HDFS简介

总体而言，HDFS要实现以下目标：

- 兼容廉价的硬件设备
- 流数据读写
- 大数据集
- 简单的文件模型
- 强大的跨平台兼容性

HDFS特殊的设计，在实现上述优良特性的同时，也使得自身具有一些应用局限性，主要包括以下几个方面：

- 不适合低延迟数据访问
- 无法高效存储大量小文件
- 不支持多用户写入及任意修改文件



3.1.1块

HDFS采用抽象的块概念可以带来以下几个明显的好处：

- **支持大规模文件存储：**文件以块为单位进行存储，一个大规模文件可以被分拆成若干个文件块，不同的文件块可以被分发到不同的节点上，因此，一个文件的大小不会受到单个节点的存储容量的限制，可以远远大于网络中任意节点的存储容量
- **简化系统设计：**首先，大大简化了存储管理，因为文件块大小是固定的，这样就可以很容易计算出一个节点可以存储多少文件块；其次，方便了元数据的管理，元数据不需要和文件块一起存储，可以由其他系统负责管理元数据
- **适合数据备份：**每个文件块都可以冗余存储到多个节点上，大大提高了系统的容错性和可用性



3.1.2 名称节点和数据节点

在HDFS中，名称节点（NameNode）负责管理分布式文件系统的命名空间（Namespace），保存了两个核心的数据结构，即FsImage和EditLog，FsImage用于维护文件系统树以及文件树中所有的文件和文件夹的元数据，操作日志文件EditLog中记录了所有针对文件的创建、删除、重命名等操作。名称节点记录了每个文件中各个块所在的数据节点的位置信息。下图展示了名称节点的数据结构。

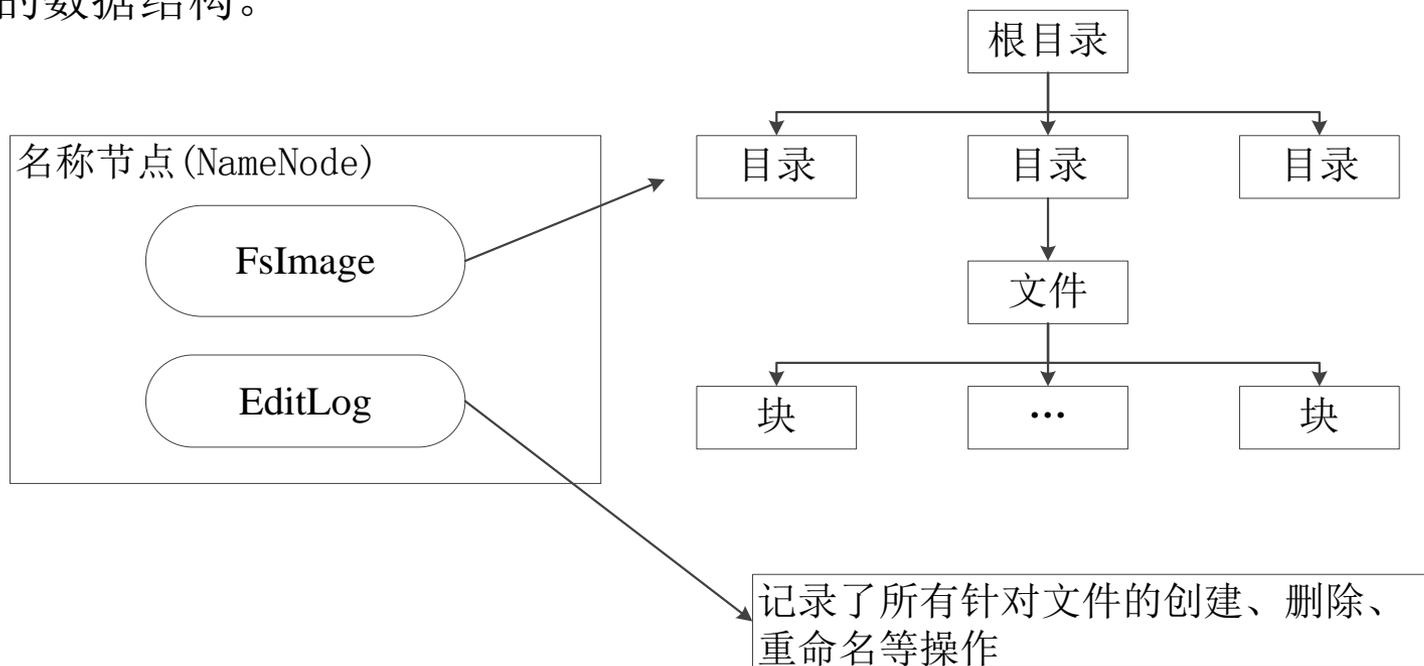


图3-3 名称节点的数据结构



3.1.2 名称节点和数据节点

数据节点（DataNode）是分布式文件系统HDFS的工作节点，负责数据的存储和读取，会根据客户端或者是名称节点的调度来进行数据的存储和检索，并且向名称节点定期发送自己所存储的块的列表。每个数据节点中的数据会被保存在各自节点的本地Linux文件系统中



3.2 HDFS体系结构

- 3.2.1 计算机集群结构
- 3.2.2 HDFS体系结构概述
- 3.2.3 HDFS体系结构的局限性



3.2.1 计算机集群结构

- 分布式文件系统把文件分布存储到多个计算机节点上，成千上万的计算机节点构成计算机集群
- 与之前使用多个处理器和专用高级硬件的并行化处理装置不同的是，目前的分布式文件系统所采用的计算机集群，都是由普通硬件构成的，这就大大降低了硬件上的开销

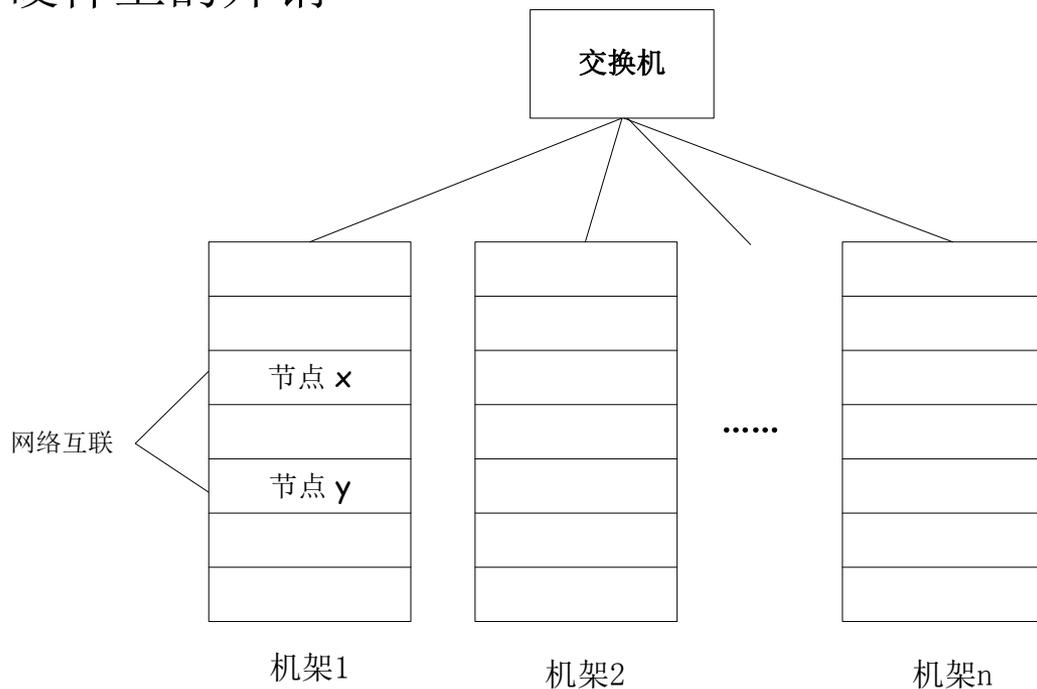


图3-1 计算机集群的基本架构



3.2.2 HDFS体系结构概述

HDFS采用了主从（Master/Slave）结构模型，一个HDFS集群包括一个名称节点（NameNode）和若干个数据节点（DataNode）（如图3-4所示）。名称节点作为中心服务器，负责管理文件系统的命名空间及客户端对文件的访问。集群中的数据节点一般是一个节点运行一个数据节点进程，负责处理文件系统客户端的读/写请求，在名称节点的统一调度下进行数据块的创建、删除和复制等操作。每个数据节点的数据实际上是保存在本地Linux文件系统中的

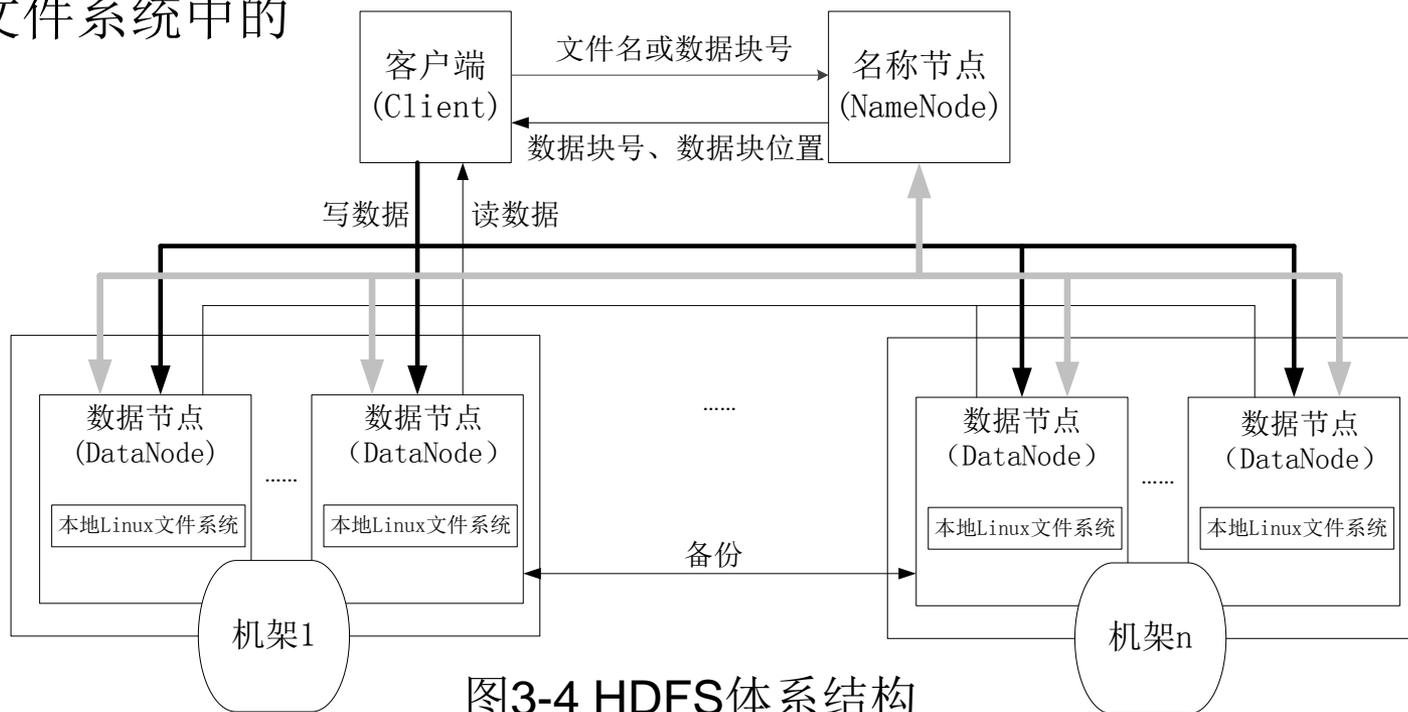


图3-4 HDFS体系结构



3.2.3 HDFS体系结构的局限性

HDFS只设置唯一一个名称节点，这样做虽然大大简化了系统设计，但也带来了一些明显的局限性，具体如下：

- (1) **命名空间的限制：**名称节点是保存在内存中的，因此，名称节点能够容纳的对象（文件、块）的个数会受到内存空间大小的限制。
- (2) **性能的瓶颈：**整个分布式文件系统的吞吐量，受限于单个名称节点的吞吐量。
- (3) **隔离问题：**由于集群中只有一个名称节点，只有一个命名空间，因此，无法对不同应用程序进行隔离。
- (4) **集群的可用性：**一旦这个唯一的名称节点发生故障，会导致整个集群变得不可用。



3.3 HDFS存储原理

- 3.3.1 冗余数据保存
- 3.3.2 数据存取策略
- 3.3.3 数据错误与恢复



3.3.1 冗余数据保存

作为一个分布式文件系统，为了保证系统的容错性和可用性，HDFS采用了多副本方式对数据进行冗余存储，通常一个数据块的多个副本会被分布到不同的数据节点上，如图3-5所示，数据块1被分别存放到数据节点A和C上，数据块2被存放在数据节点A和B上。这种多副本方式具有以下几个优点：

- (1) 加快数据传输速度
- (2) 容易检查数据错误
- (3) 保证数据可靠性

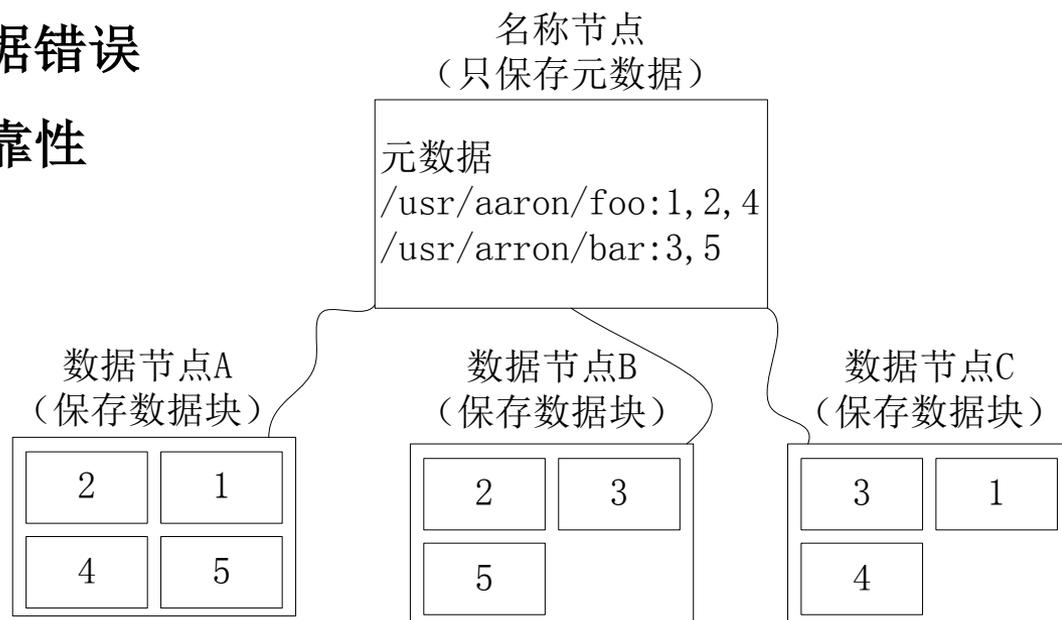


图3-5 HDFS数据块多副本存储



3.3.2 数据存取策略

数据存取策略包括数据存放、数据读取和数据复制等方面，它在很大程度上会影响到整个分布式文件系统的读写性能，是分布式文件系统的核心内容。

1. 数据存放

- 为了提高数据可靠性与系统可用性，以及充分利用网络带宽，**HDFS**采用了以机架（**Rack**）为基础的数据存放策略
- **HDFS**默认每个数据节点都是在不同的机架上，缺点是写入数据的时候不能充分利用同一机架内部机器之间的带宽。优点：首先，可以获得很高的数据可靠性，即使一个机架发生故障，位于其他机架上的数据副本仍然是可用的；其次，在读取数据的时候，可以在多个机架并行读取数据，大大提高了数据读取速度；再次，可以更容易实现系统内部负载均衡和错误处理。



3.3.2 数据存取策略

2. 数据读取

- HDFS提供了一个API可以确定一个数据节点所属的机架ID，客户端也可以调用API获取自己所属的机架ID
- 当客户端读取数据时，从名称节点获得数据块不同副本的存放位置列表，列表中包含了副本所在的数据节点，可以调用API来确定客户端和这些数据节点所属的机架ID，当发现某个数据块副本对应的机架ID和客户端对应的机架ID相同时，就优先选择该副本读取数据，如果没有发现，就随机选择一个副本读取数据



本讲小结

- 分布式文件系统是大数据时代解决大规模数据存储问题的有效解决方案，**HDFS**开源实现了**GFS**，可以利用由廉价硬件构成的计算机集群实现海量数据的分布式存储
- **HDFS**具有兼容廉价的硬件设备、流数据读写、大数据集、简单的文件模型、强大的跨平台兼容性等特点。但是，也要注意，**HDFS**也有自身的局限性，比如不适合低延迟数据访问、无法高效存储大量小文件和不支持多用户写入及任意修改文件等
- 块是**HDFS**核心的概念，一个大的文件会被拆分成很多个块。**HDFS**采用抽象的块概念，具有支持大规模文件存储、简化系统设计、适合数据备份等优点
- **HDFS**采用了主从（**Master/Slave**）结构模型，一个**HDFS**集群包括一个名称节点和若干个数据节点。名称节点负责管理分布式文件系统的命名空间；数据节点是分布式文件系统**HDFS**的工作节点，负责数据的存储和读取
- **HDFS**采用了冗余数据存储，增强了数据可靠性，加快了数据传输速度。**HDFS**还采用了相应的数据存放、数据读取策略，来提升系统整体读写响应性能。**HDFS**把硬件出错看作一种常态，设计了错误恢复机制



主讲教师



主讲教师：林子雨

单位：厦门大学计算机科学系

E-mail: ziyulin@xmu.edu.cn

个人网页: <http://www.cs.xmu.edu.cn/linziyu>

数据库实验室网站: <http://dblabb.xmu.edu.cn>



扫一扫访问个人主页

林子雨，男，1978年出生，博士（毕业于北京大学），现为厦门大学计算机科学系助理教授（讲师），曾任厦门大学信息科学与技术学院院长助理、晋江市发展和改革委员会副局长。中国高校首个“数字教师”提出者和建设者，厦门大学数据库实验室负责人，厦门大学云计算与大数据研究中心主要建设者和骨干成员，2013年度厦门大学奖教金获得者。主要研究方向为数据库、数据仓库、数据挖掘、大数据、云计算和物联网，编著出版中国高校第一本系统介绍大数据知识的专业教材《大数据技术原理与应用》并成为畅销书籍；主讲厦门大学计算机系本科生课程《数据库系统原理》和研究生课程《分布式数据库》《大数据技术基础》。具有丰富的政府和企业信息化培训经验，曾先后给中国移动通信集团公司、福州马尾区政府、福建省物联网科学研究院、石狮市物流协会、厦门市物流协会等多家单位和个人开展信息化培训，累计培训人数达2000人以上。



大数据学习教材推荐



扫一扫访问教材官网

《大数据技术原理与应用——概念、存储、处理、分析与应用》，由厦门大学计算机科学系林子雨博士编著，是中国高校第一本系统介绍大数据知识的专业教材。

全书共有13章，系统地论述了大数据的基本概念、大数据处理架构Hadoop、分布式文件系统HDFS、分布式数据库HBase、NoSQL数据库、云数据库、分布式并行编程模型MapReduce、流计算、图计算、数据可视化以及大数据在互联网、生物医学和物流等各个领域的应用。在Hadoop、HDFS、HBase和MapReduce等重要章节，安排了入门级的实践操作，让读者更好地学习和掌握大数据关键技术。

本书可以作为高等院校计算机专业、信息管理等相关专业的大数据课程教材，也可供相关技术人员参考、学习、培训之用。

欢迎访问《大数据技术原理与应用——概念、存储、处理、分析与应用》教材官方网站：
<http://dblab.xmu.edu.cn/post/bigdata>



Principles and Applications of Big Data Technology - Big Data Conception, Storage, Processing, Analysis and Application

林子雨 编著



中国工信出版集团

人民邮电出版社
POSTS & TELECOM PRESS

The background of the slide features several faint, light-blue silhouettes of people. At the top, there are two groups of people standing and holding hands. On the right side, a person is shown in profile, looking towards the center. On the left side, two people are shown in profile, facing each other. The overall scene suggests a group of people in a meeting or a presentation.

Thank You!

Department of Computer Science, Xiamen University