

## 第三讲 分布式文件系统HDFS



徐辰  
cxu@dase.ecnu.edu.cn

华东师范大学



## 大纲

- 2 □ 文件系统(FS)
  - ✦ 文件系统概述
  - ✦ 文件与目录
  - ✦ 文件的物理结构
- 分布式文件系统(DFS)
- Hadoop分布式文件系统(HDFS)

## 文件系统概述

- 3 □ 文件系统出现的原因
  - ✦ 用户直接操作和管理辅助存储器上信息 (01二进制序列), 繁琐复杂、易于出错、可靠性差
- 文件系统是操作系统中负责**管理**和**存取**信息的模块
  - ✦ 统一管理用户和系统信息的存储、检索、更新、共享和保护
  - ✦ 为用户提供一整套方便有效的文件使用和操作方法

## 文件系统概述

- 4 □ 文件系统的功能:
  - ✦ 文件的按名存取 (基本功能)
  - ✦ 文件目录的建立和维护 (用于实现上述基本功能)
  - ✦ 实现逻辑文件到物理文件的转换 (核心内容)
  - ✦ 文件存储空间的分配和管理
  - ✦ 数据保密、保护和共享
  - ✦ 提供一组用户使用的操作

## 大纲

- 5 □ 文件系统(FS)
  - ✦ 文件系统概述
  - ✦ 文件与目录
  - ✦ 文件的物理结构
- 分布式文件系统(DFS)
- Hadoop分布式文件系统(HDFS)

## 文件

- 6 □ 文件
  - ✦ 文件是由文件名字标识的一组信息的集合
  - ✦ 各操作系统的文件命名规则略有不同
- 实现**按名存取**的文件系统的优点
  - ✦ 将用户从复杂的物理存储地址管理中解放出来
  - ✦ 可方便地对文件提供各种安全、保密和保护措施
  - ✦ 实现文件的共享 (同名共享、异名共享)

## 文件目录

7

### 如何实现“按名存取”？

- 当用户要求存取某个文件时，系统查找**目录文件**，获得对应的**文件目录**
- 在文件目录中，根据用户给定的文件名寻找到对应该文件的**文件控制块**（文件目录项）
- 通过文件控制块所记录的该文件的相关信息（如文件信息存放的相对位置或文件信息首块的物理位置）依次存取该文件的内容。

## 文件目录

8

### 概念：

- 文件目录：建立和维护的关于系统的所有文件的清单
- 文件控制块：每个目录项对应一个文件的信息描述
- 目录文件：目录信息也以文件的形式存放

### 文件控制块的基本内容

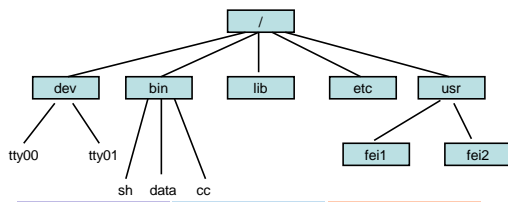
- 文件存取控制信息：如文件名、用户名、文件主存取权限等
- 文件结构信息：文件逻辑结构、文件的物理结构等
- 文件使用信息：已打开该文件的进程数、文件的修改情况等
- 文件管理信息：文件建立日期、文件访问日期等

## 文件目录

9

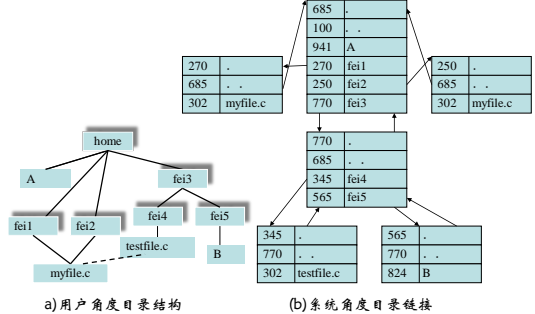
### 树形目录结构

- 文件的全名，应该从根目录开始，到该文件名为止，目录路径+文件名
- 例如，/user/include/testfile.c



## 不同角度的目录结构

10



## 大纲

11

### 文件系统(FS)

- 文件系统概述
- 文件与目录
- 文件的物理结构

### 分布式文件系统(DFS)

### Hadoop分布式文件系统(HDFS)

## 文件的物理结构

12

### 文件在物理存储中的存放方法和组织关系

- 块（物理记录）的划分
- 记录的排列
- 索引的组织
- 信息的搜索

### 常见的文件物理结构：

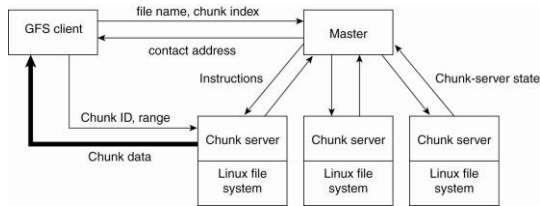
- 顺序文件，连续存储
  - 链接文件
  - 索引文件
- 非连续存储



## Cluster-Based Distributed File Systems (1)

19

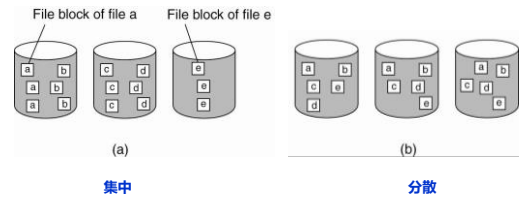
- 主节点进行管理，从节点存储数据



## Cluster-Based Distributed File Systems (2)

20

- 文件切分成块，分散存储在从节点上



## 大纲

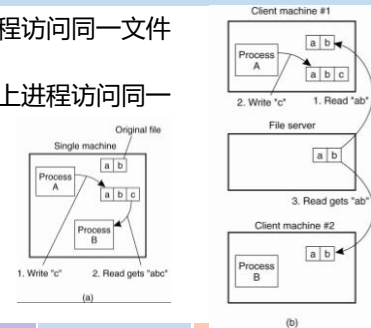
21

- 文件系统(FS)
- 分布式文件系统(DFS)
  - 体系架构
  - 文件访问
  - 备份与一致性
  - 容错管理
- Hadoop分布式文件系统(HDFS)

## 文件访问

22

- 单机多进程访问同一文件
  - 读写锁
- 不同机器上进程访问同一文件



## 大纲

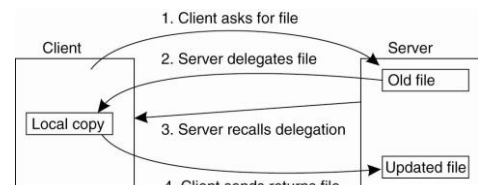
23

- 文件系统(FS)
- 分布式文件系统(DFS)
  - 体系架构
  - 文件访问
  - 备份与一致性
  - 容错管理
- Hadoop分布式文件系统(HDFS)

## 备份与一致性

24

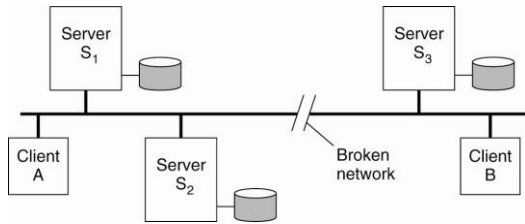
- 客户端备份: Client-Server DFS



## 备份与一致性

25

### 服务器端备份: Cluster-Based DFS



## 大纲

26

- 文件系统(FS)
- 分布式文件系统(DFS)
  - ✦ 体系架构
  - ✦ 文件访问
  - ✦ 备份与一致性
  - ✦ 容错管理
- Hadoop分布式文件系统(HDFS)

## 故障类型

27

- 发生故障即停机: Fail stop
- 拜占庭将军问题
  - ✦ 拜占庭失效指一方向另一方发送消息, 另一方没有收到, 或者收到了错误的信息的情形

## 大纲

28

- 文件系统(FS)
- 分布式文件系统(DFS)
- Hadoop分布式文件系统(HDFS)
  - ✦ 设计思想
  - ✦ 体系架构
  - ✦ 工作原理
  - ✦ 容错机制

## Hadoop发展简史

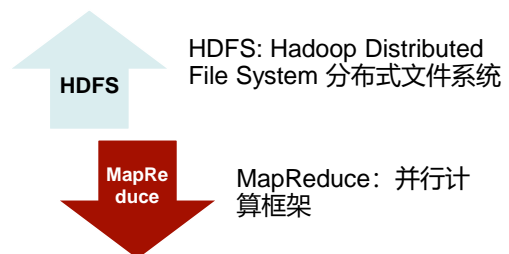


29

- Hadoop最初是由Apache Lucene项目的创始人Doug Cutting开发的文本搜索库。Hadoop源自始于2002年的Apache Nutch项目——一个开源的网络搜索引擎并且也是Lucene项目的一部分
- 2004年, Nutch项目也模仿GFS开发了自己的分布式文件系统NDFS (Nutch Distributed File System), 也就是HDFS的前身
- 2004年, 谷歌公司又发表了另一篇具有深远影响的论文, 阐述了MapReduce分布式编程思想
- 2005年, Nutch开源实现了谷歌的MapReduce
- 2006年2月, Nutch中的NDFS和MapReduce开始独立出来, 成为Lucene项目的一个子项目, 称为Hadoop, 同时, Doug Cutting加盟雅虎
- 2008年1月, Hadoop正式成为Apache顶级项目, Hadoop也逐渐开始被雅虎之外的其他公司使用
- 2008年4月, Hadoop打破世界纪录, 成为最快排序1TB数据的系统, 它采用一个由910个节点构成的集群进行运算, 排序时间只用了209秒
- 在2009年5月, Hadoop更是把1TB数据排序时间缩短到62秒。Hadoop从此名声大震, 迅速发展成为大数据时代最具影响力的开源分布式开发平台, 并成为事实上的大数据处理标准

## Hadoop核心项目

30





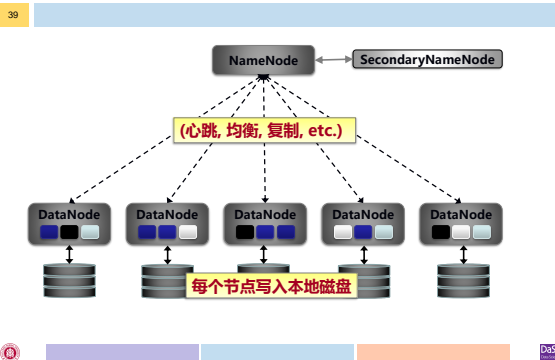
## HDFS设计假设、目标

硬件失效	流式数据访问	存储数据较大	简化的数据一致性模型	多硬件平台支持	移动计算能力比移动数据更划算
<ul style="list-style-type: none"> <li>硬件的异常比软件的异常更加常见。</li> <li>对于有上百台服务器的数据中心来说，认为总有服务器异常，<b>硬件异常是常态</b>。</li> <li>HDFS需要监测这些异常，并自动恢复数据。</li> </ul>	<ul style="list-style-type: none"> <li>基于HDFS的应用仅采用流式方式读取数据。</li> <li>运行在HDFS上的应用并非以通用业务为目的的应用程序</li> <li>应用程序关注的是<b>吞吐量</b>，而非响应时间。</li> <li>非POSIX标准接口的数据访问。</li> </ul>	<ul style="list-style-type: none"> <li>运行在HDFS的应用程序有较大的数据需要处理。</li> <li>典型的文件大小<b>为GB到TB级别</b>。</li> <li>文件仅支持追加，而不允许修改。</li> </ul>	<ul style="list-style-type: none"> <li>应用程序采用<b>WORM (Write Once Read Many)</b>的数据读写模型。</li> <li>文件仅支持追加，而不允许修改。</li> </ul>	<ul style="list-style-type: none"> <li>HDFS易于<b>运行不同的平台</b>上。</li> </ul>	<ul style="list-style-type: none"> <li>计算和存储采用就近原则，计算离数据最近。</li> <li>就近原则将有效减少网络的负载，降低网络阻塞。</li> </ul>

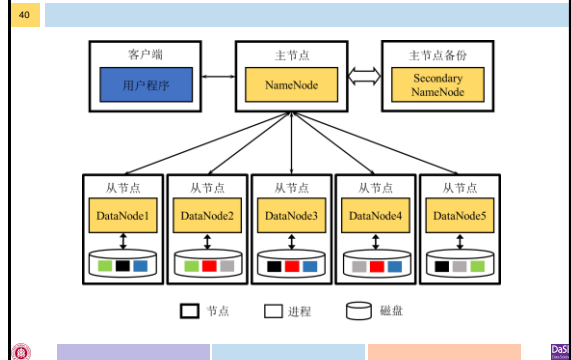
## 大纲

- 文件系统(FS)
- 分布式文件系统(DFS)
- **Hadoop分布式文件系统(HDFS)**
  - ✚ 设计思想
  - ✚ 体系架构
  - ✚ 工作原理
  - ✚ 容错机制
  - ✚ 编程使用

## HDFS架构图



## HDFS架构图



## HDFS角色

41

- **NameNode** – 每个集群一个名字节点
  - ✚ 负责文件系统元数据操作、数据块的复制和定位
- **SecondaryNameNode** – NameNode的备份节点
  - CheckpointNode or BackupNode (backups)
- **DataNodes** – 集群中每个节点一个数据节点
  - ✚ 负责数据块的存储
  - ✚ 为客户端提供实际文件数据

Diagram labels: **NameNode** (Master), **SecondaryNameNode** (backups), **DataNode** (Slaves).

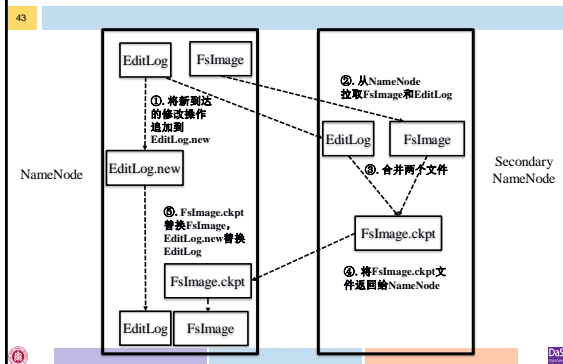
## NameNode

42

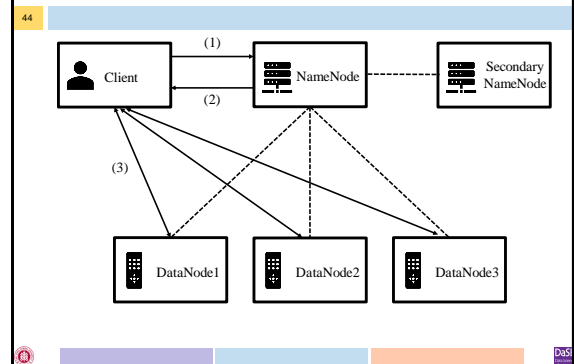
- **FsImage**: 内存在文件目录结构及其元信息在磁盘上的快照
- **EditLog**: 两次快照之间，针对目录及文件修改的操作

Diagram illustrating the NameNode internal structure. It shows a **磁盘** (Disk) containing **FsImage** and **EditLog**. The **内存** (Memory) contains a **Root** directory tree structure with **Dir1**, **Dir2**, and **Dir3** as subdirectories. **File1**, **File2**, and **File3** are listed under the directories. Below the files are **Block1**, **Block2**, **Block3**, **Block4**, **Block5**, and **Block6**.

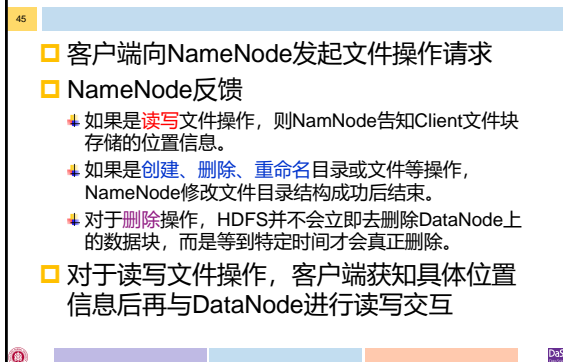
## NameNode与SecondaryNameNode



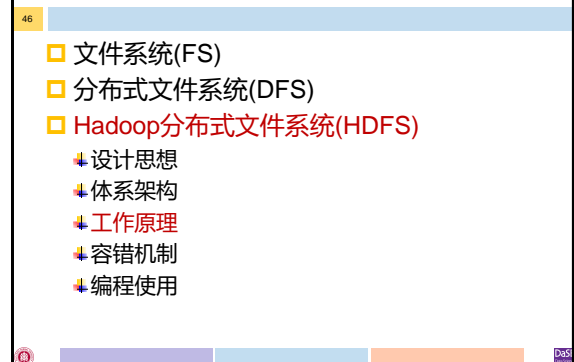
## 应用程序执行流程



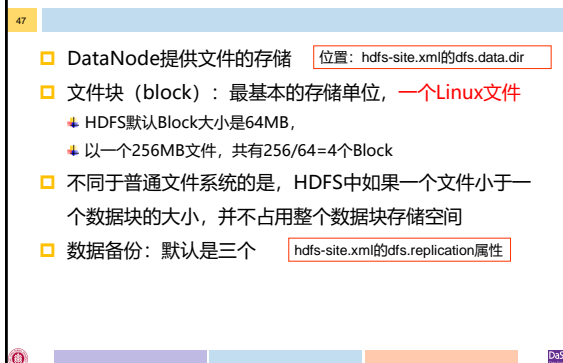
## 应用程序执行流程



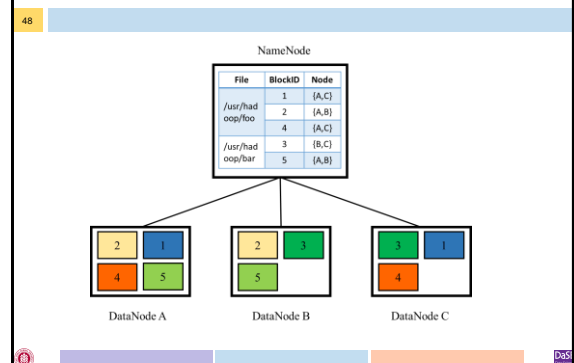
## 大纲



## 文件分块与备份



## 文件分块与备份

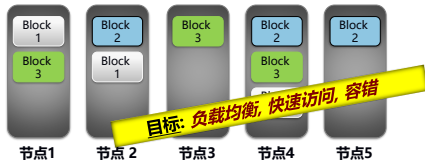




## 文件块存放策略（启发式）

49

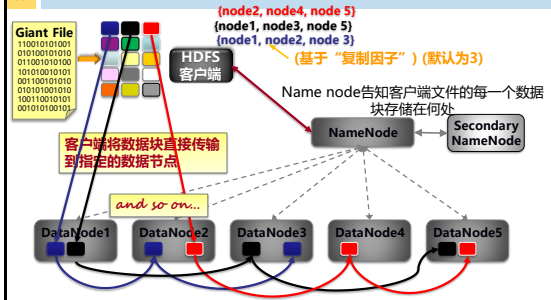
- 第一个副本：放置在上传文件的数据节点；如果是集群外提交，则随机挑选一台磁盘不太满、CPU不太忙的节点（快速写入）
- 第二个副本：放置在与第一个副本不同的机架rack的节点上（减少跨rack的网络流量）
- 第三个副本：与第一个副本相同机架的其他节点上（应对交换机故障）
- 更多副本：随机节点



e.g., 复制因子 = 3

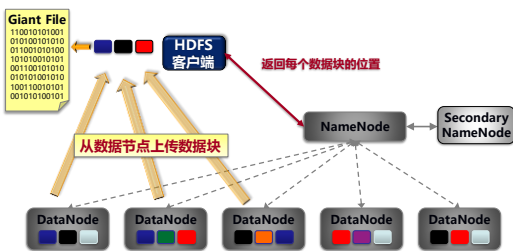
## 文件写入HDFS

50



## 从HDFS读取文件

51



## 数据读取策略

52

- 当客户端读取数据时，从NameNode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
- 最近者优先原则**：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

## 文件读写与一致性

53

- “一次写入多次读取”
  - 一个文件经过创建、写入和关闭后就不得改变文件中的内容
  - 已经写入到HDFS文件，仅容许在文件末尾追加数据，即append操作
  - 当对一个文件进行写入操作时，包括文件的追加操作，NameNode将拒绝其它针对该文件的读、写请求
  - 当对一个文件进行读取操作时，NameNode容许其它针对该文件的读请求。

## 简化的一致性模型

54

- 简化的好处
  - 避免读写冲突、用户编程无需考虑文件锁
- 问题
  - 假如用户的确需要修改已有文件中的内容，怎么办？
  - 如果HDFS容许修改文件中的已有内容，会带来哪些问题？

## 大纲

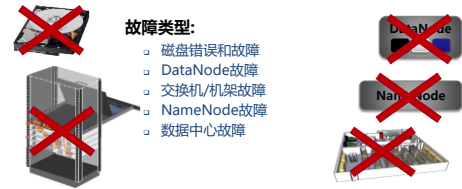
55

- 文件系统(FS)
- 分布式文件系统(DFS)
- Hadoop分布式文件系统(HDFS)
  - ✦ 设计思想
  - ✦ 体系架构
  - ✦ 工作原理
  - ✦ 容错机制
  - ✦ 编程使用

## 容错机制

56

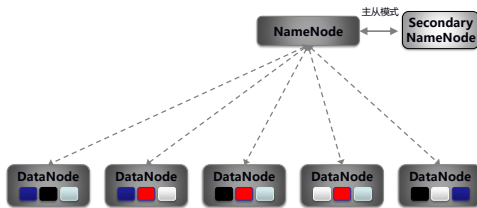
- HDFS在设计时就考虑到故障(硬件和软件)会经常出现



## NameNode故障

57

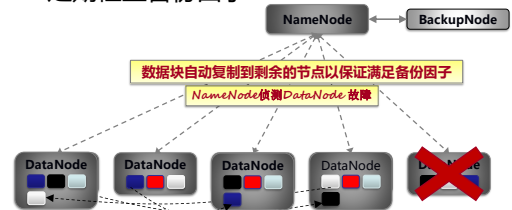
- 根据SecondaryNameNode中的FsImage和Editlog数据进行恢复



## DataNode故障

58

- “宕机”，节点上面的所有数据都会被标记为“不可读”
- 定期检查备份因子



## 其它故障

59

- 磁盘错误或故障：数据校验
- 交换机/机架故障
- 数据中心故障

## 大纲

60

- 文件系统(FS)
- 分布式文件系统(DFS)
- Hadoop分布式文件系统(HDFS)
  - ✦ 设计思想
  - ✦ 体系架构
  - ✦ 工作原理
  - ✦ 容错机制
  - ✦ 编程使用

## HDFS Shell

61

### 新建目录

✚ `./bin/hdfs dfs -mkdir input`

### 上传文件

✚ `./bin/hdfs dfs -put ./README.txt ./input`

### 查看文件

✚ `./bin/hdfs dfs -cat ./input/README.txt`

✚ `./bin/hdfs dfs -ls /`

### 拷贝

✚ `./bin/hdfs dfs -cp ./input/README.txt /`

## Java程序

62

### 设置环境

```
✚ Configuration conf = new Configuration();
✚ conf.set("fs.defaultFS", "hdfs://localhost:9000");
✚ conf.set("fs.hdfs.impl",
  "org.apache.hadoop.hdfs.DistributedFileSystem");
```

### 写文件

✚ `FSDDataOutputStream`

### 读文件

✚ `FSDDataInputStream`

## HDFS功能

63

### HDFS适合做什么

- ✚ 大文件存储
- ✚ 流式数据访问



### HDFS不适合做什么

- ✚ 大量小文件
- ✚ 随机读取
- ✚ 低延迟读取

### 谁在使用HDFS? 地球人都在用



## 课后阅读

64

### 分布式系统概念与设计, George Coulouris等著, 金蓓弘等译

✚ 第12章 12.1、12.2

✚ 第21章 21.5.1

### 论文

✚ Ghemawat, S., Gobioff, H., & Leung, S.-T. (2003). The Google File System. In SOSP (pp. 29–43).

## 本讲小结

65

### FS

### DFS

### HDFS

- ✚ 设计思想
- ✚ 体系架构
- ✚ 工作原理
- ✚ 容错机制
- ✚ 编程使用

谢谢! Q&A

