

101 – 5

GlusterFS 系统原理剖析



讲师介绍

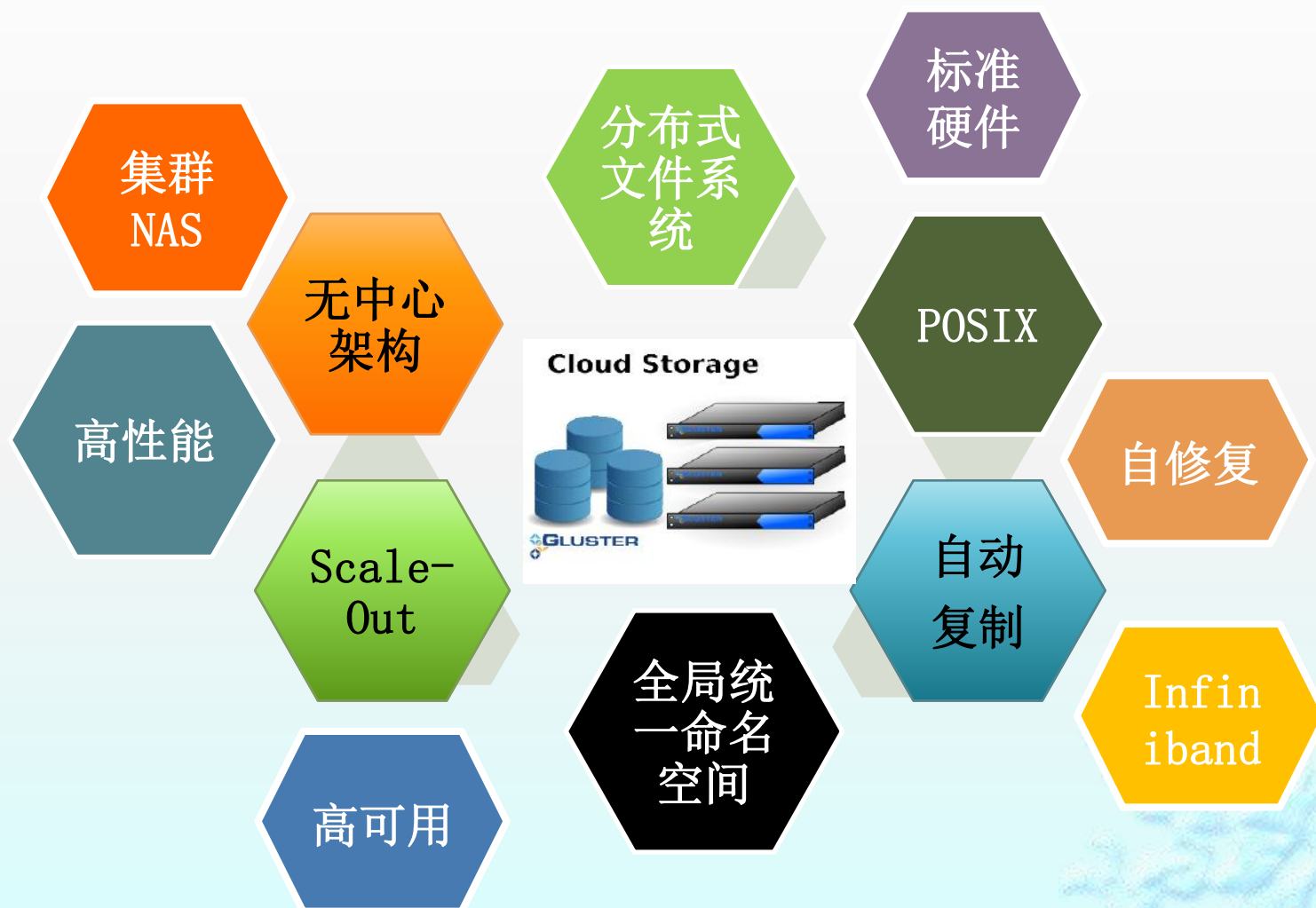
- ◆ 刘爱贵，博士
- ◆ 毕业于中科院高能物理研究所
- ◆ 主要研究方向：分布式存储、高性能计算、数据挖掘
- ◆ 专注于存储技术的研究与开发，GlusterFS等分布式文件系统资深理论研究与实践者
- ◆ 博客：<http://blog.csdn.net/liuaigui>
- ◆ 微博：<http://weibo.com/liuag>
- ◆ Email：aigui.liu@gmail.com
- ◆ QQ：9187434



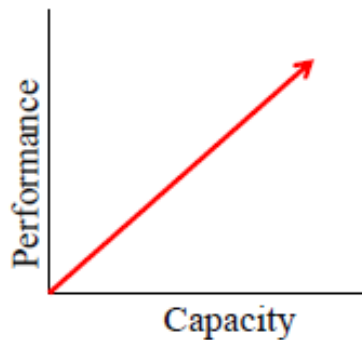
培 训 提 纲

- ◆ GlusterFS 架 构 特 点
- ◆ GlusterFS 核 心 工 作 原 理
- ◆ GlusterFS 典 型 功 能 剖 析

GlusterFS 是什么？



GlusterFS 架构设计目标



◆ Elasticity

- Flexibility adapt to growth/reduction
- Add, delete volumes & users
- Without disruption

◆ Scale linearly

- Multiple dimensions
 - Performance
 - Capacity
- Aggregated resources

◆ Eliminate metadata

- Improve file access speed

◆ Simplicity

- Ease of management
- No complex Kernel patches
- Run in user space

GlusterFS 架构特点

软件定义

无中心架构

全局命名空间

高性能

用户空间实现

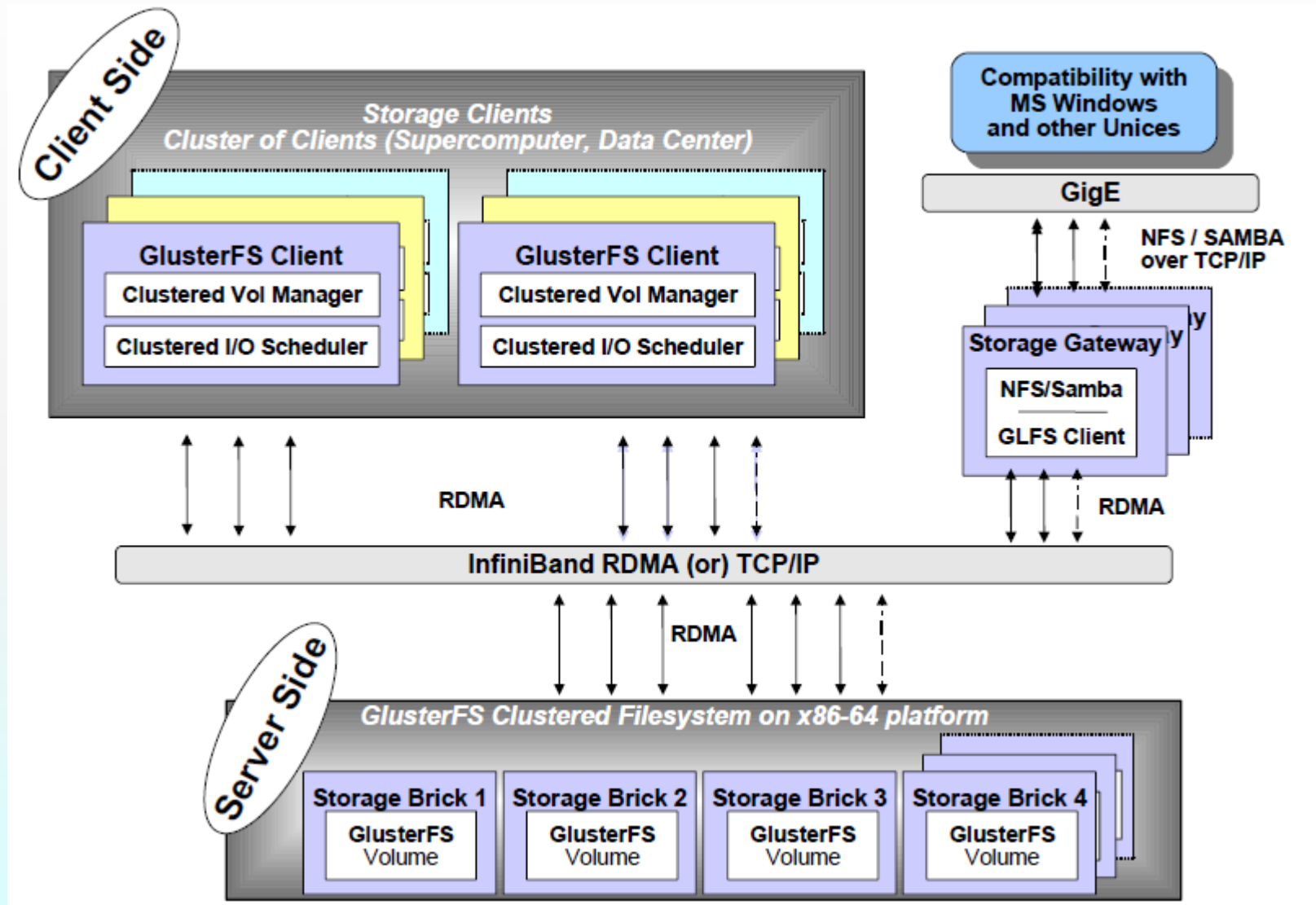
堆栈式设计

弹性横向扩展

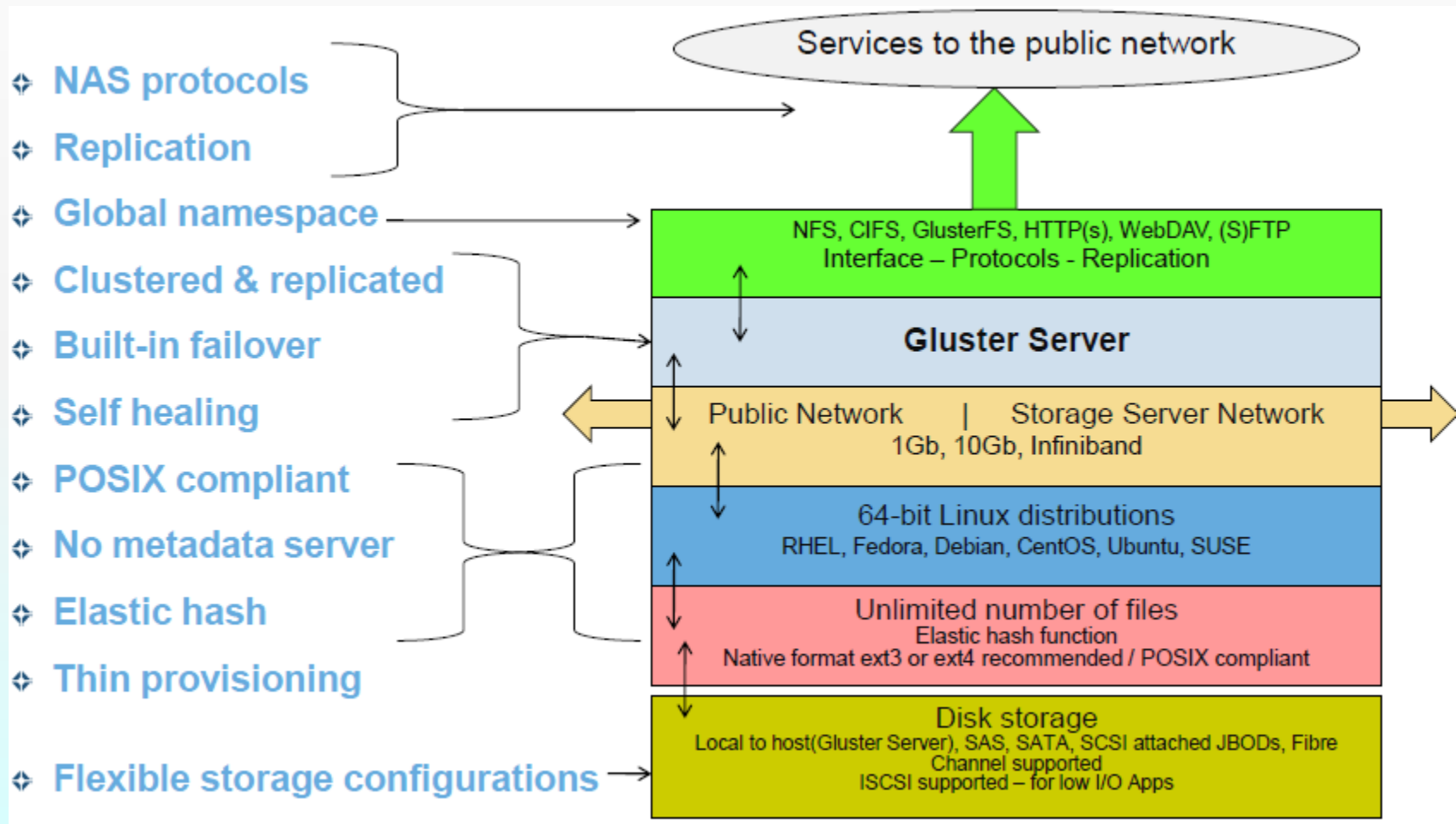
高速网络通信

数据自动修复

GlusterFS 总体架构

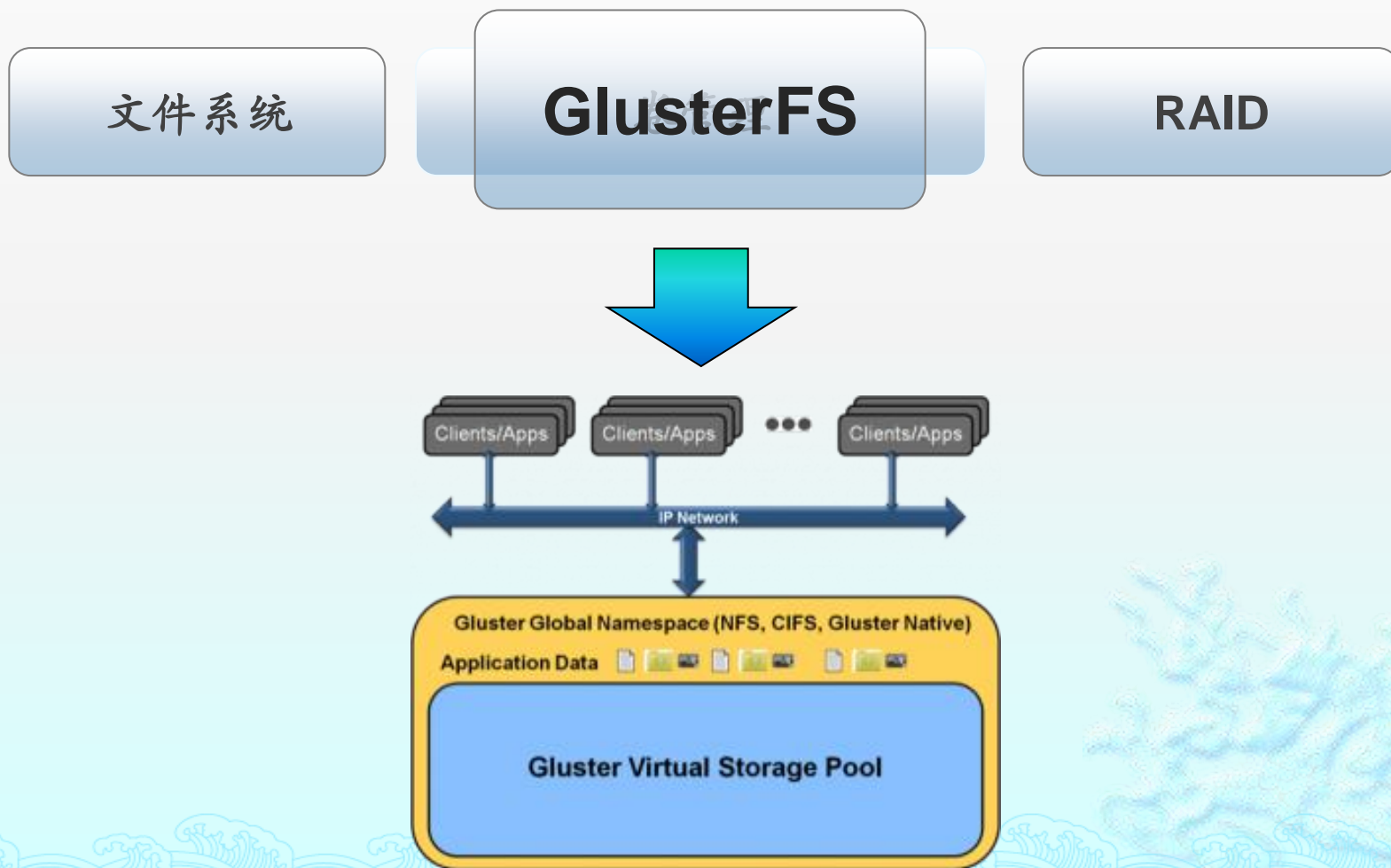


模块化/堆栈式存储OS架构

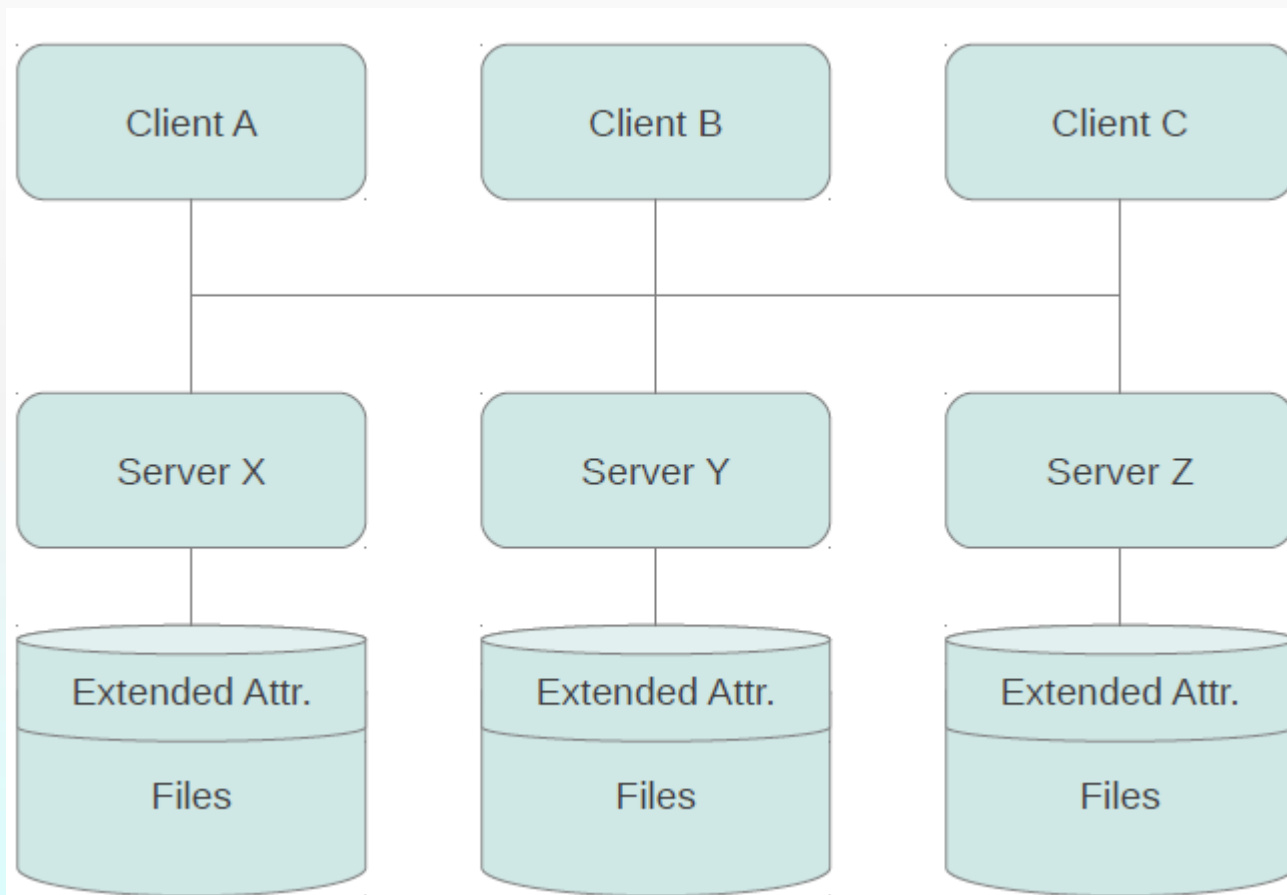


全局统一命名空间

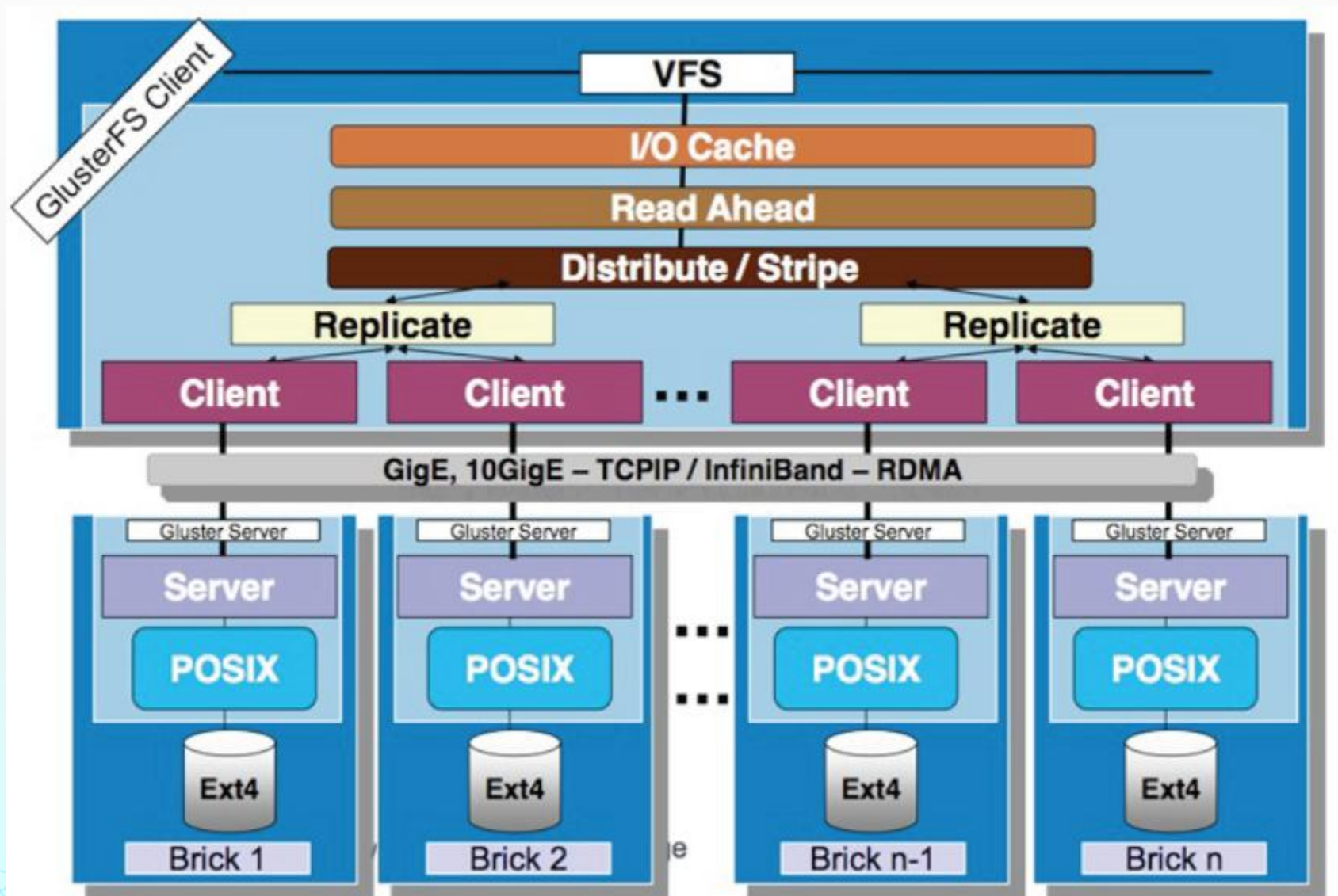
通过分布式文件系统将物理分散的存储资源虚拟化成统一的存储池



无集中元数据服务



GlusterFS堆栈式软件架构



GlusterFS 基本概念

◆ Brick

- A filesystem mountpoint
- A unit of storage used as a GlusterFS building block

◆ Translator

- Logic between the bits and the Global Namespace
- Layered to provide GlusterFS functionality

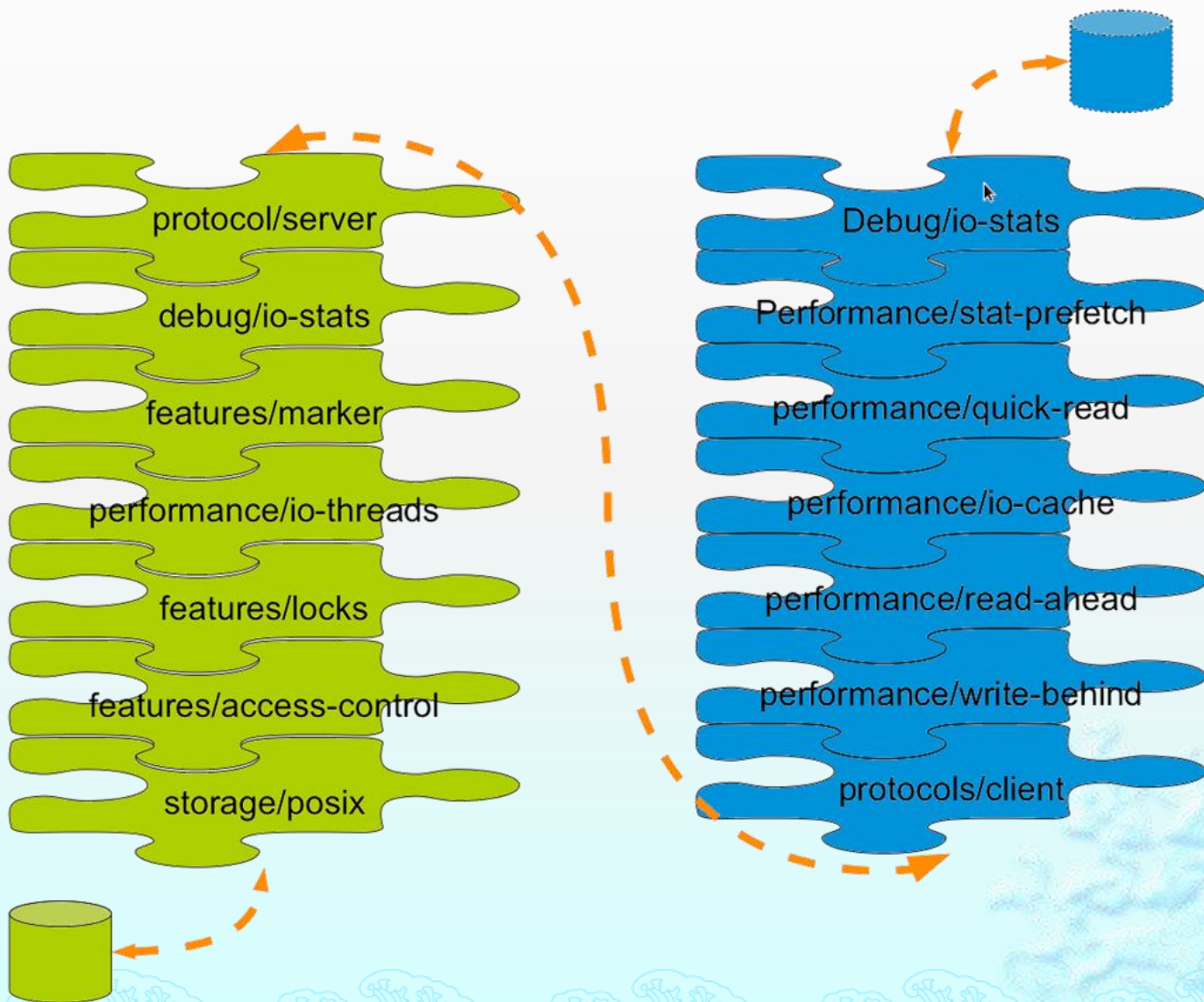
◆ Volume

- Bricks combined and passed through translators

◆ Node / Peer

- Server running the gluster daemon and sharing volumes

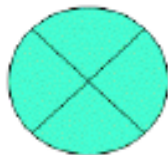
Translators



弹性hash算法

Elastic Hashing Algorithm

Goal: Systematically locate files based solely on their name



Algorithm

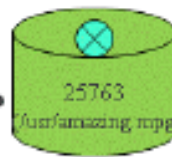
- 1) Run path/filename through hash
- 2) Assign to logical disk based on numerical result
- 3) Separate logical storage from physical storage



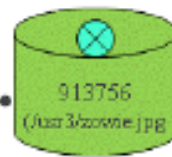
Logical Disk 1
10-11k



Logical Disk 64
19k-20k



Logical Disk 256
25k-30k



Logical Disk 30,000
91k-92k

Gluster Mgt Functions: Add, subtract, replicate, heal, recover, etc.

...



Disk 1



Disk 2

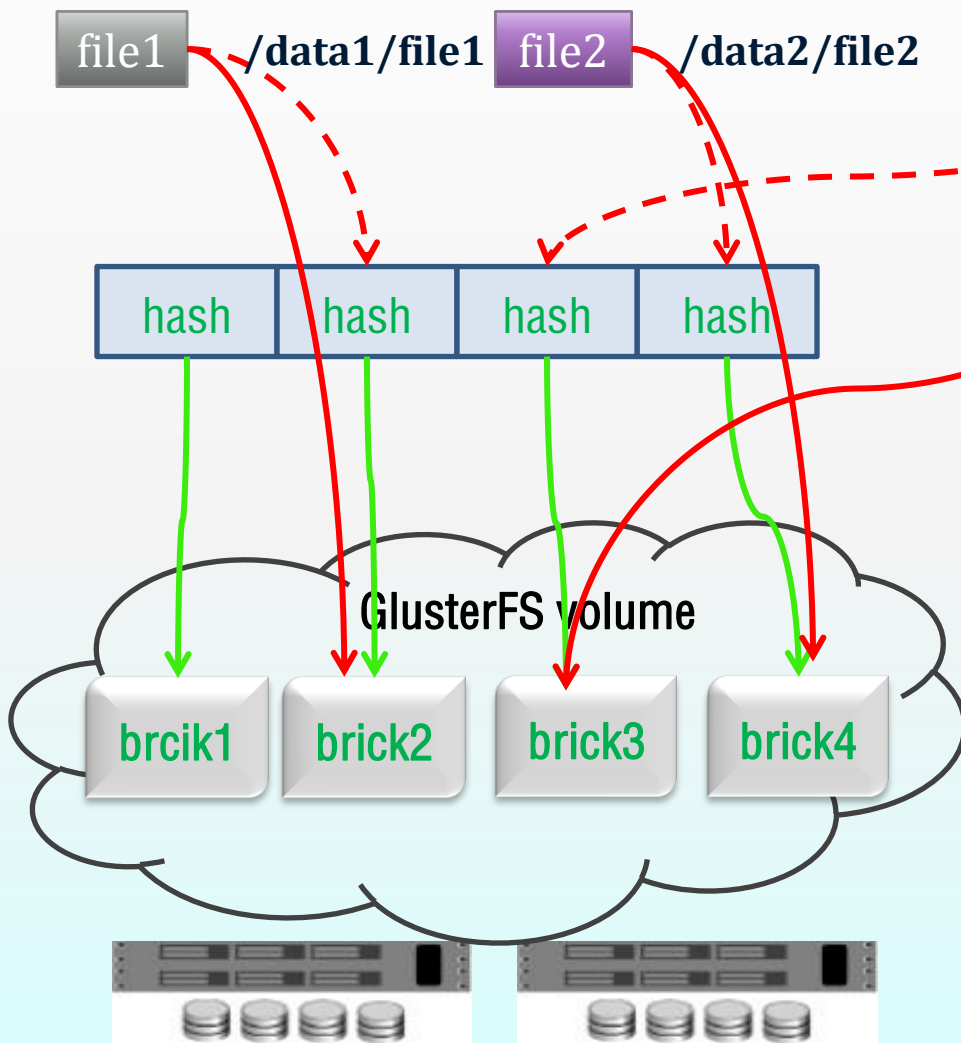
...



Disk 10

- Provision huge # of virtual disks
- Use hash algorithm to assign systematically & in distributed fashion to *virtual* storage locations
- Virtualization lets you deal flexibly with physical disks
 - ...add, subtract, deal with different disk performance or capacity parameters, etc.
- **WE HAVE MET OUR GOAL!**

弹性Hash算法流程



- 1、使用Davies-Meyer算法计算32位hash值，输入参数为文件名；
- 2、根据hash值在集群中选择子卷（存储服务器），进行文件定位；
- 3、对所选择的子卷进行数据访问。

Brick	Hash range
Brick1	00000000 ~ 3FFFFFFF
Brick2	4FFFFFFF ~ 7FFFFFFF
Brick3	8FFFFFFF ~ BFFFFFFF
Brick4	CFFFFFFF ~ FFFFFFFF

GlusterFS 卷 类 型

◆ 基本卷

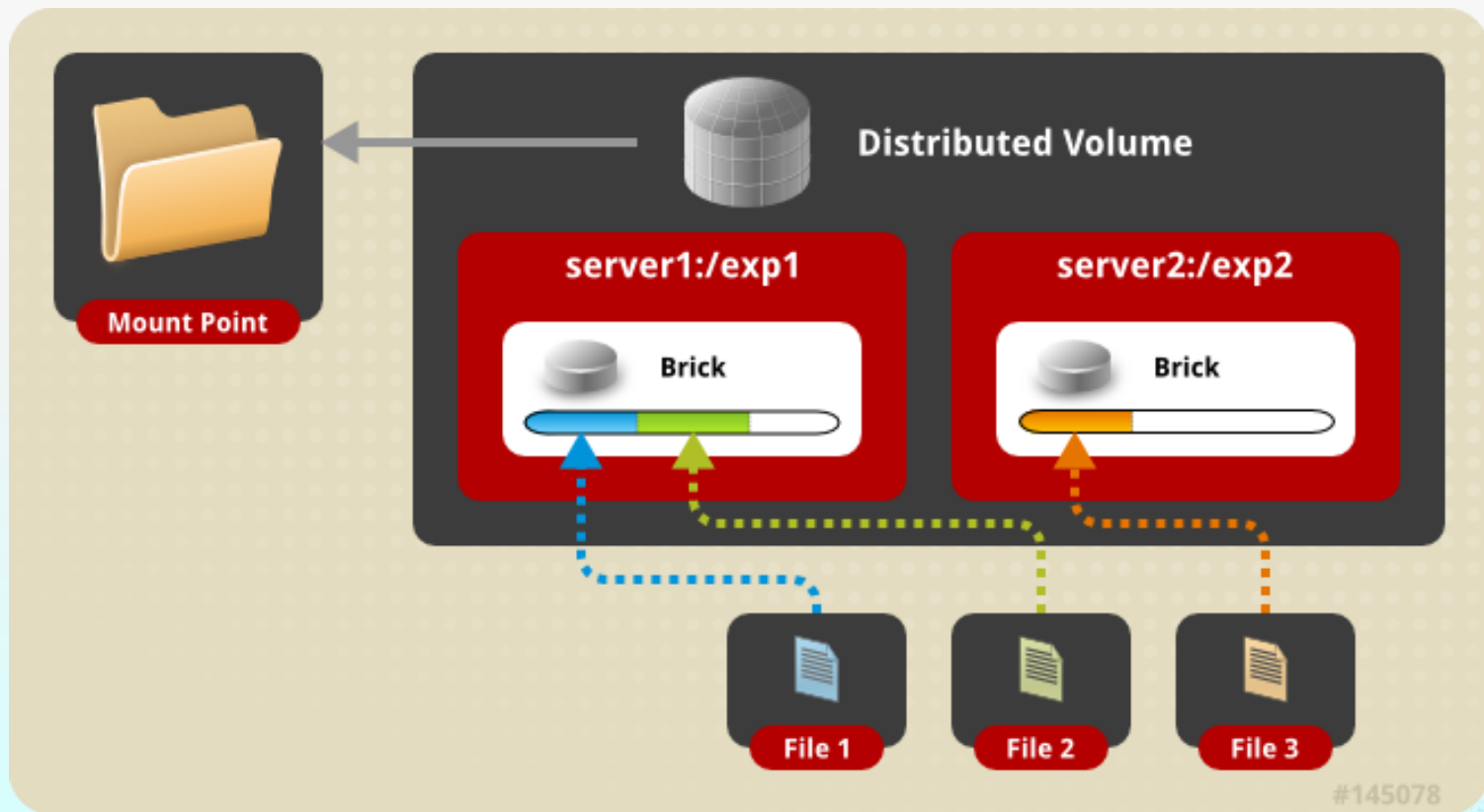
- ◆ 哈希卷 (Distributed Volume)
- ◆ 复制卷 (Replicated Volume)
- ◆ 条带卷 (Striped Volumes)

◆ 复合卷

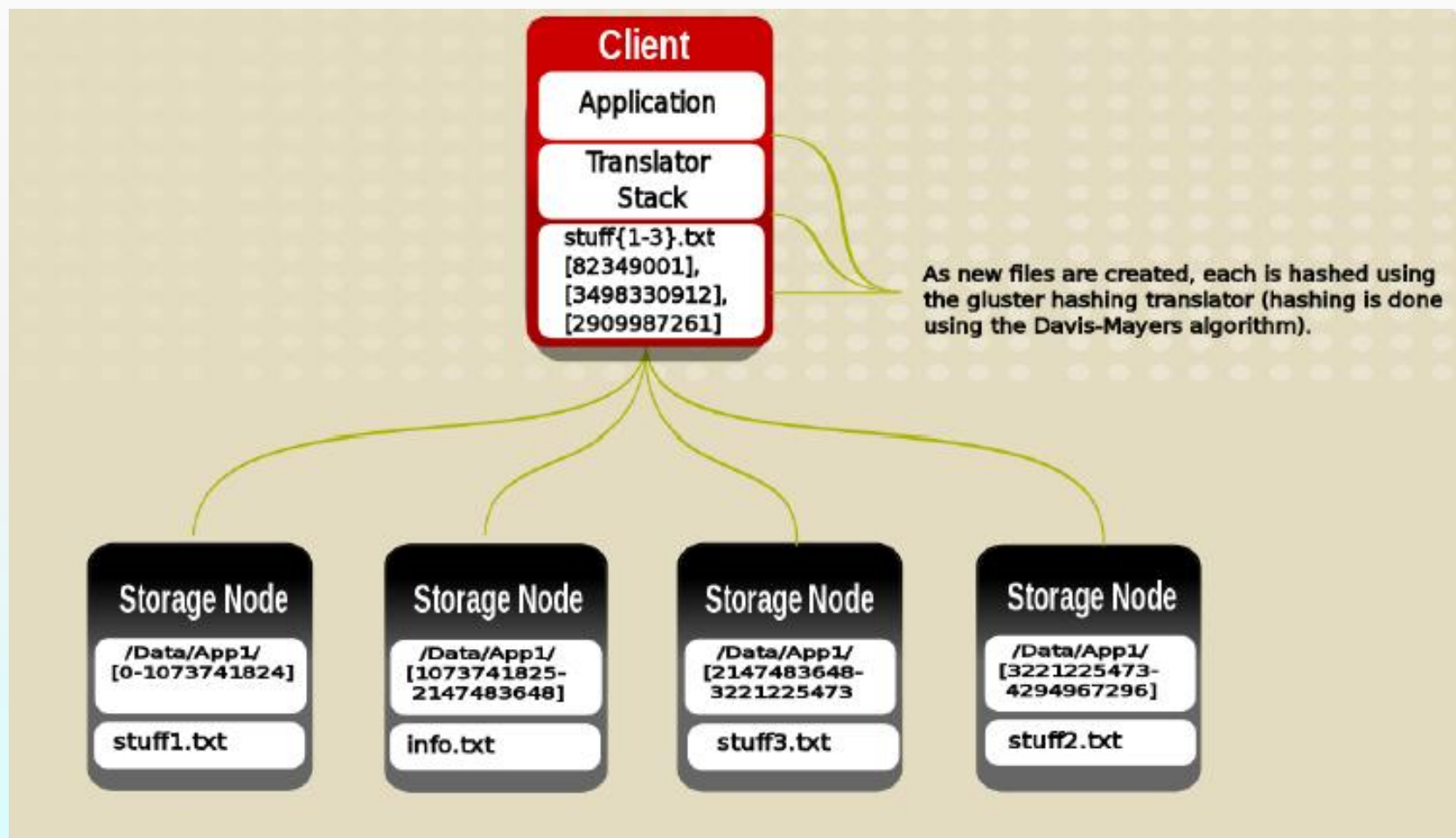
- ◆ 哈希复制卷 (Distributed Replicated Volume)
- ◆ 哈希条带卷 (Distributed Striped Volume)
- ◆ 复制条带卷 (Replicated Striped Volume)
- ◆ 哈希复制条带卷 (Distributed Replicated Striped Volume)

哈希卷 (Distributed Volume)

- 文件通过hash算法在所有brick上分布
- 文件级RAID 0，不具有容错能力

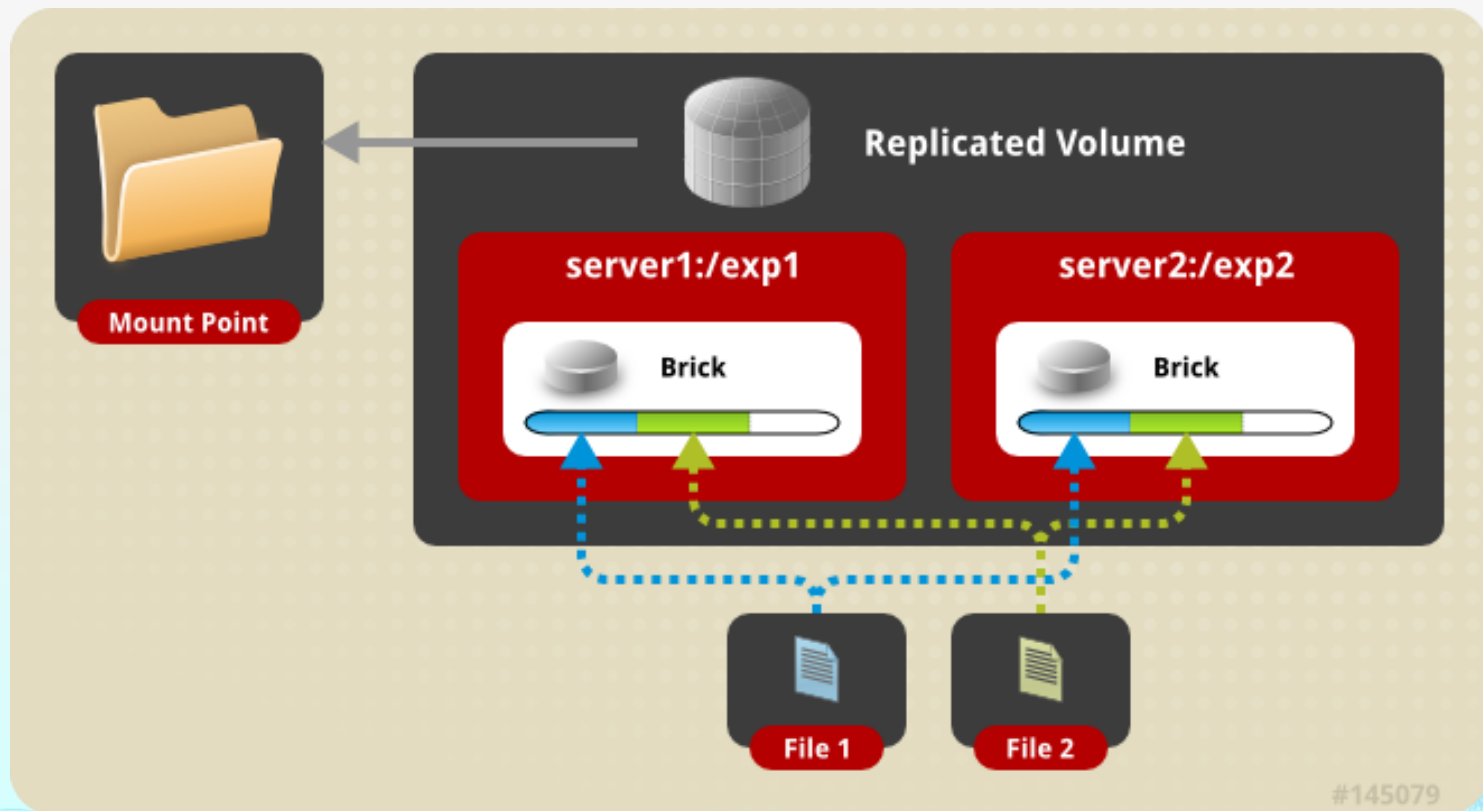


哈希卷工作原理

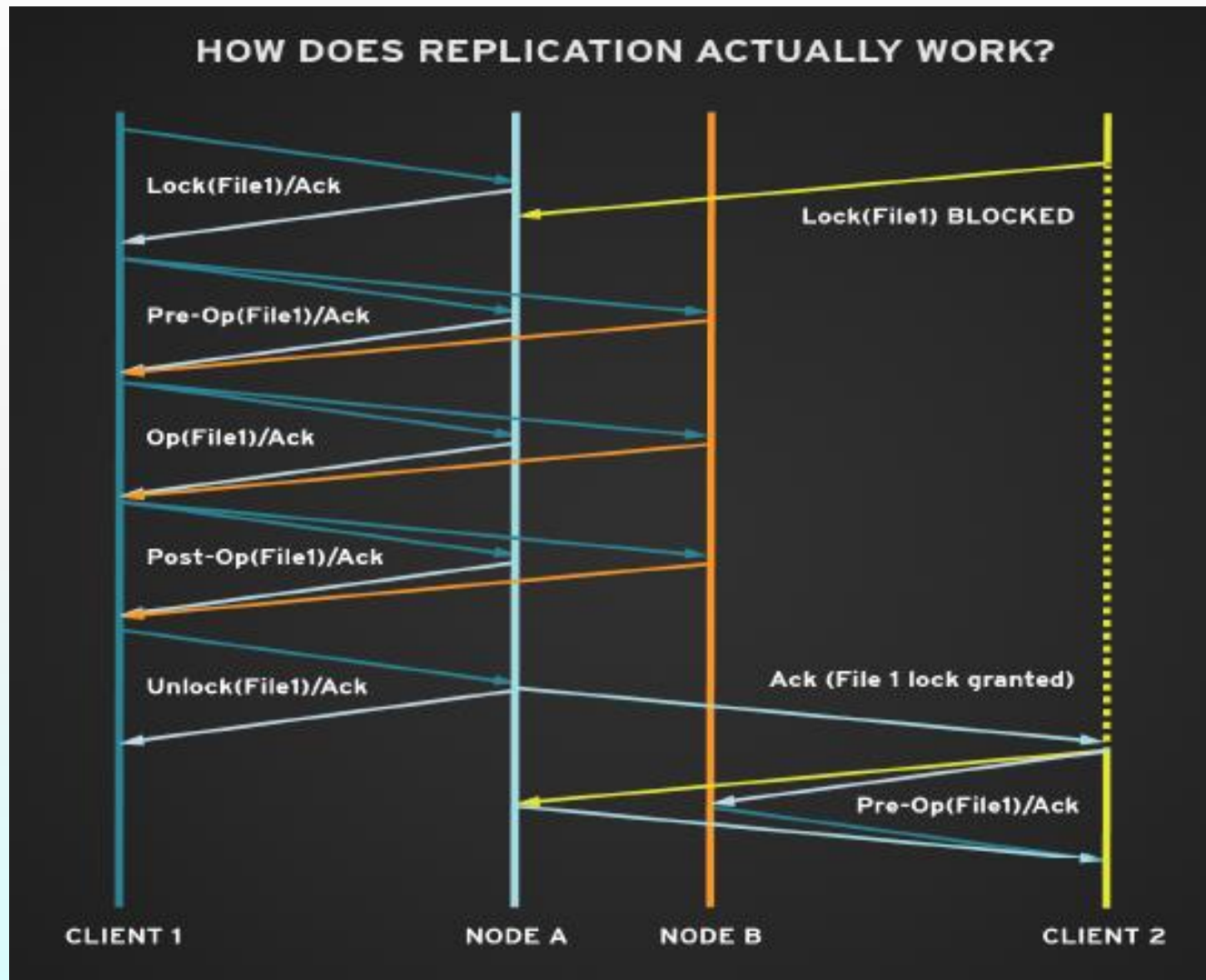


复制卷 (Replicated Volume)

- 文件同步复制到多个brick上
- 文件级RAID 1，具有容错能力
- 写性能下降，读性能提升

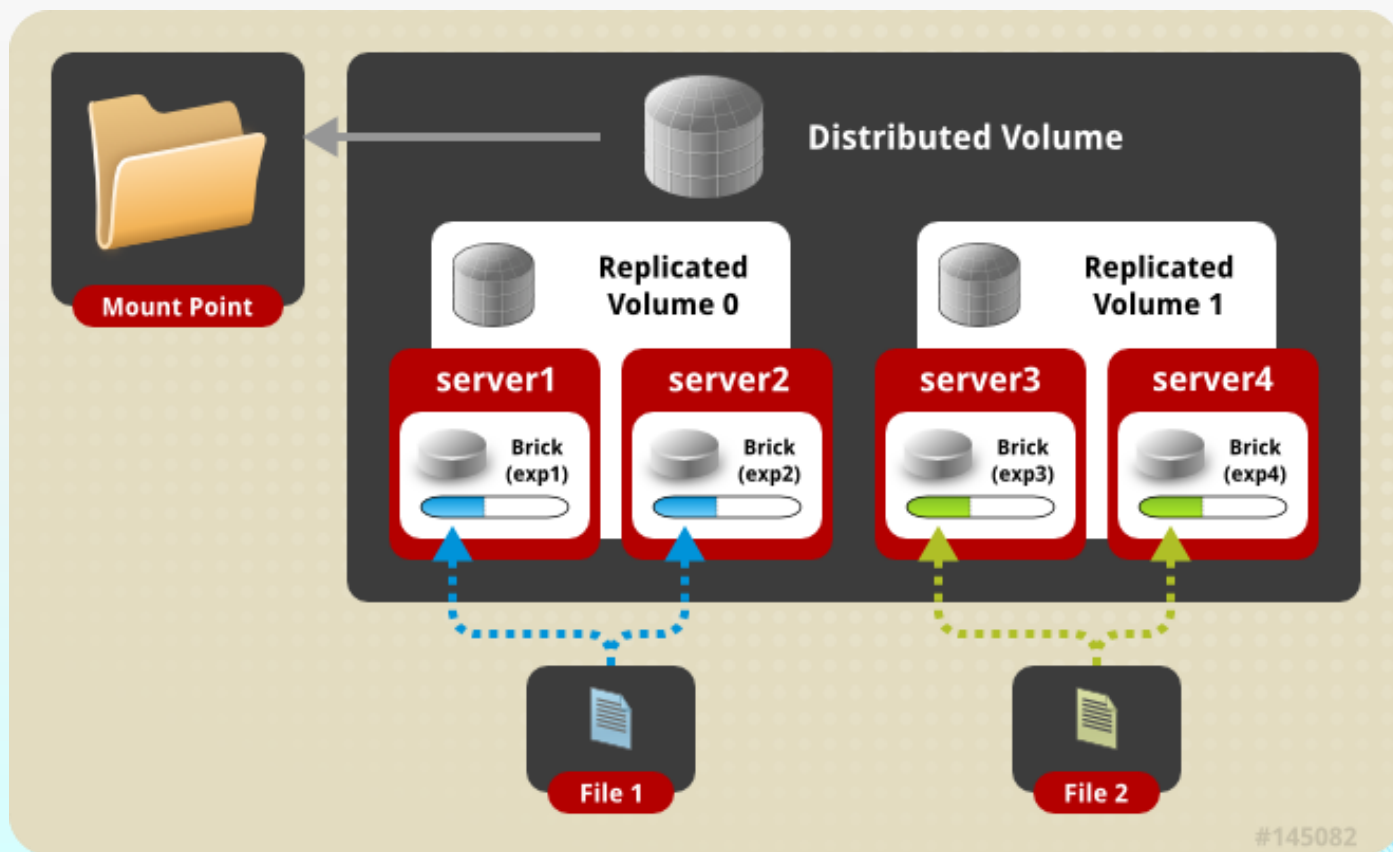


复制卷工作原理



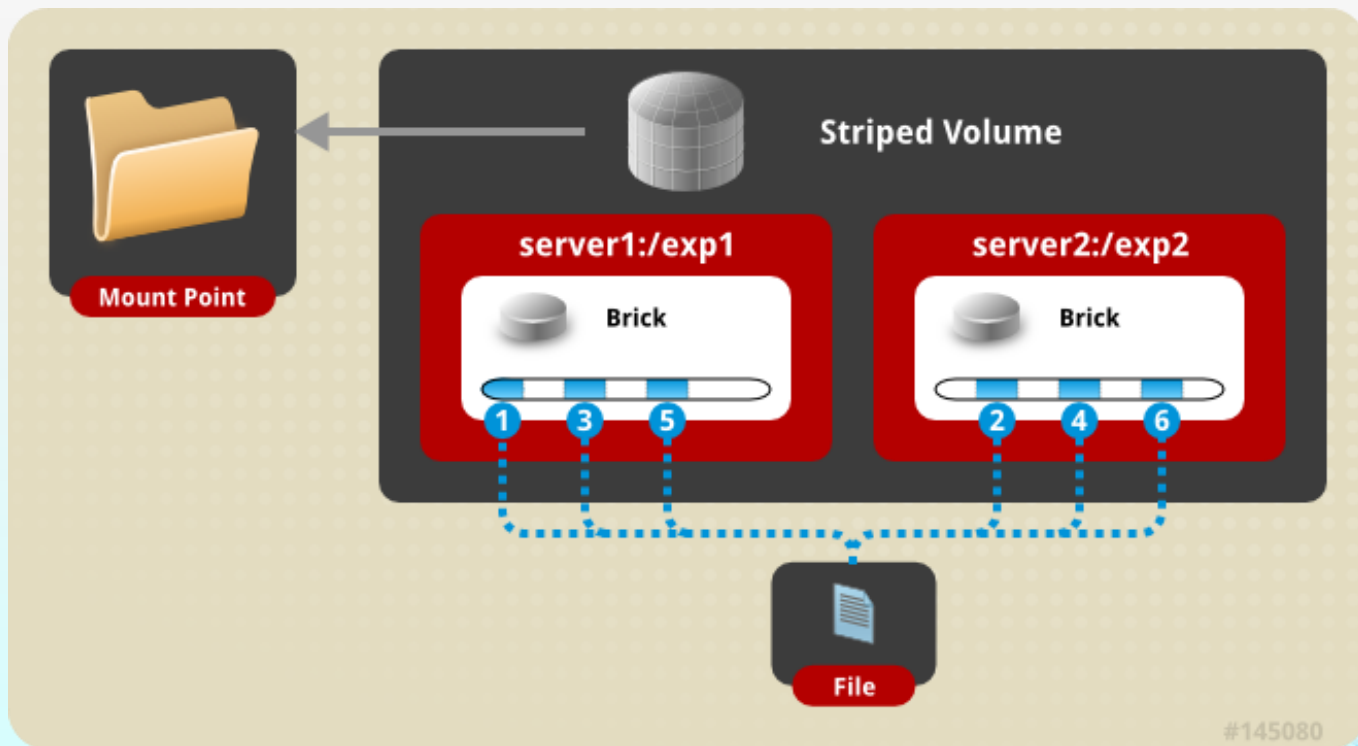
复合卷：哈希 + 复制

- 哈希卷和复制卷的复合方式
- 同时具有哈希卷和复制卷的特点



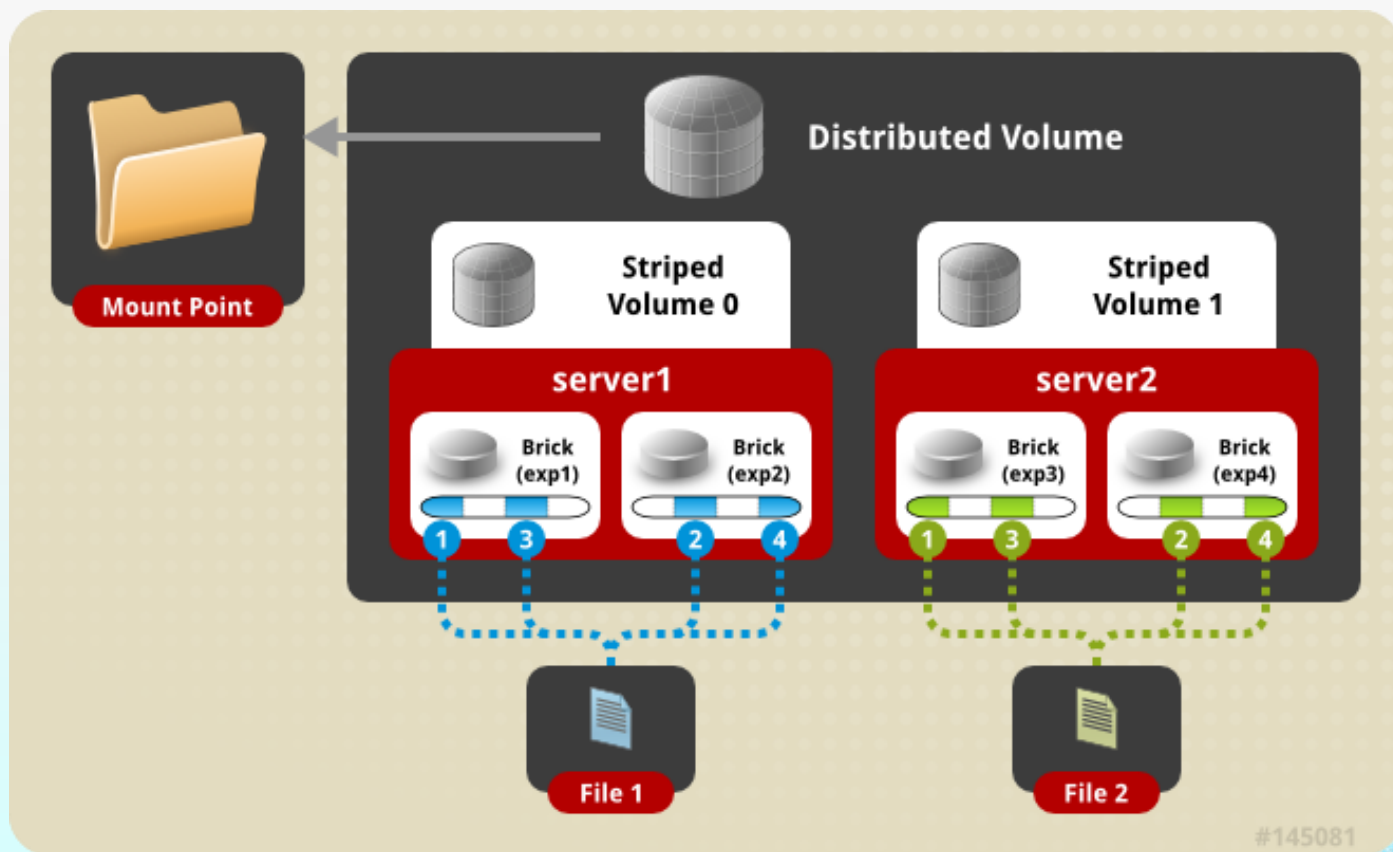
条带卷 (Striped Volumes)

- 单个文件分布到多个brick上，支持超大文件
- 类似RAID 0，以Round-Robin方式
- 通常用于HPC中的超大文件高并发访问



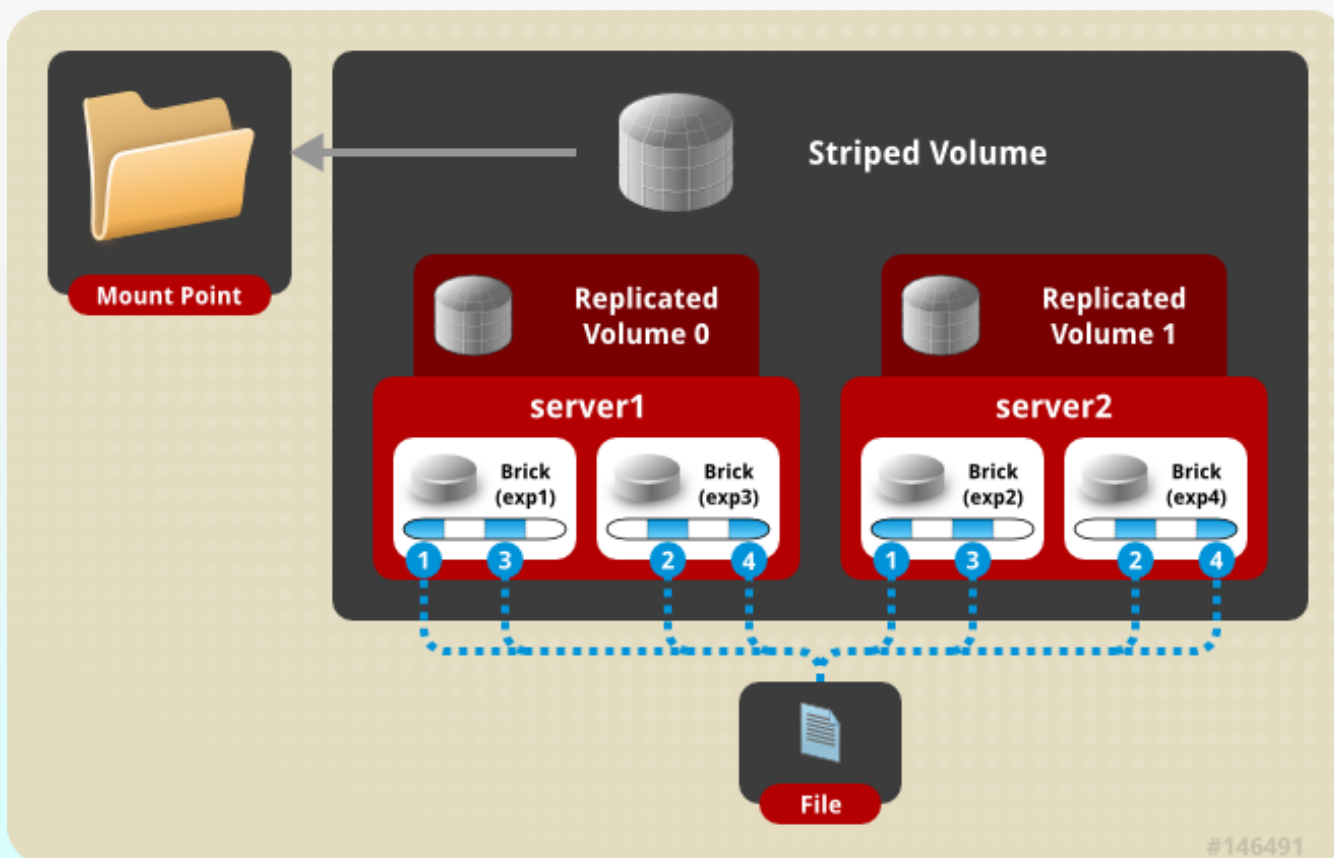
复合卷：哈希 + 条带

- 哈希卷和条带卷的复合方式
- 同时具有哈希卷和条带卷的特点



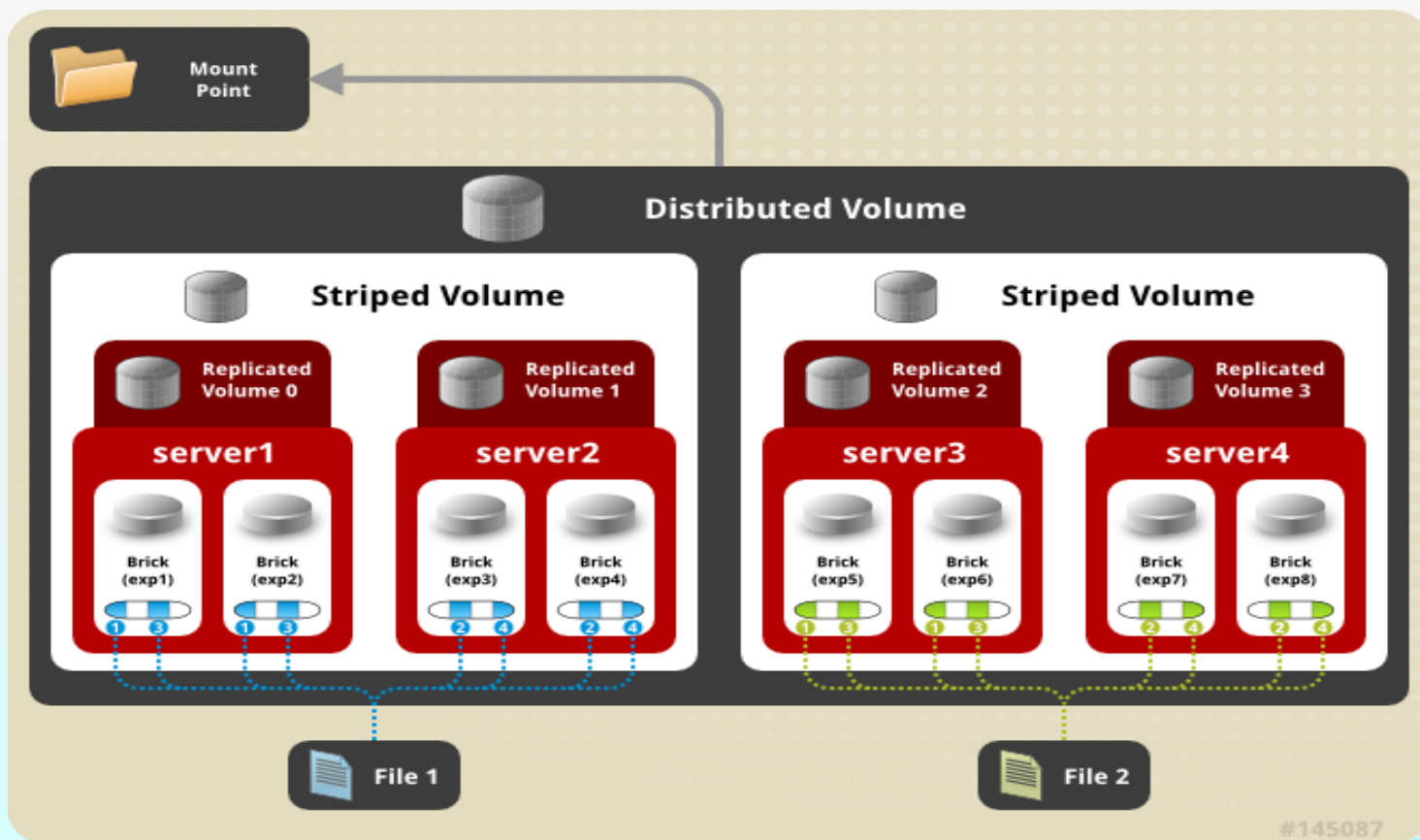
复合卷：条带 + 复制

- 类似RAID 10
- 同时具有条带卷和复制卷的特点

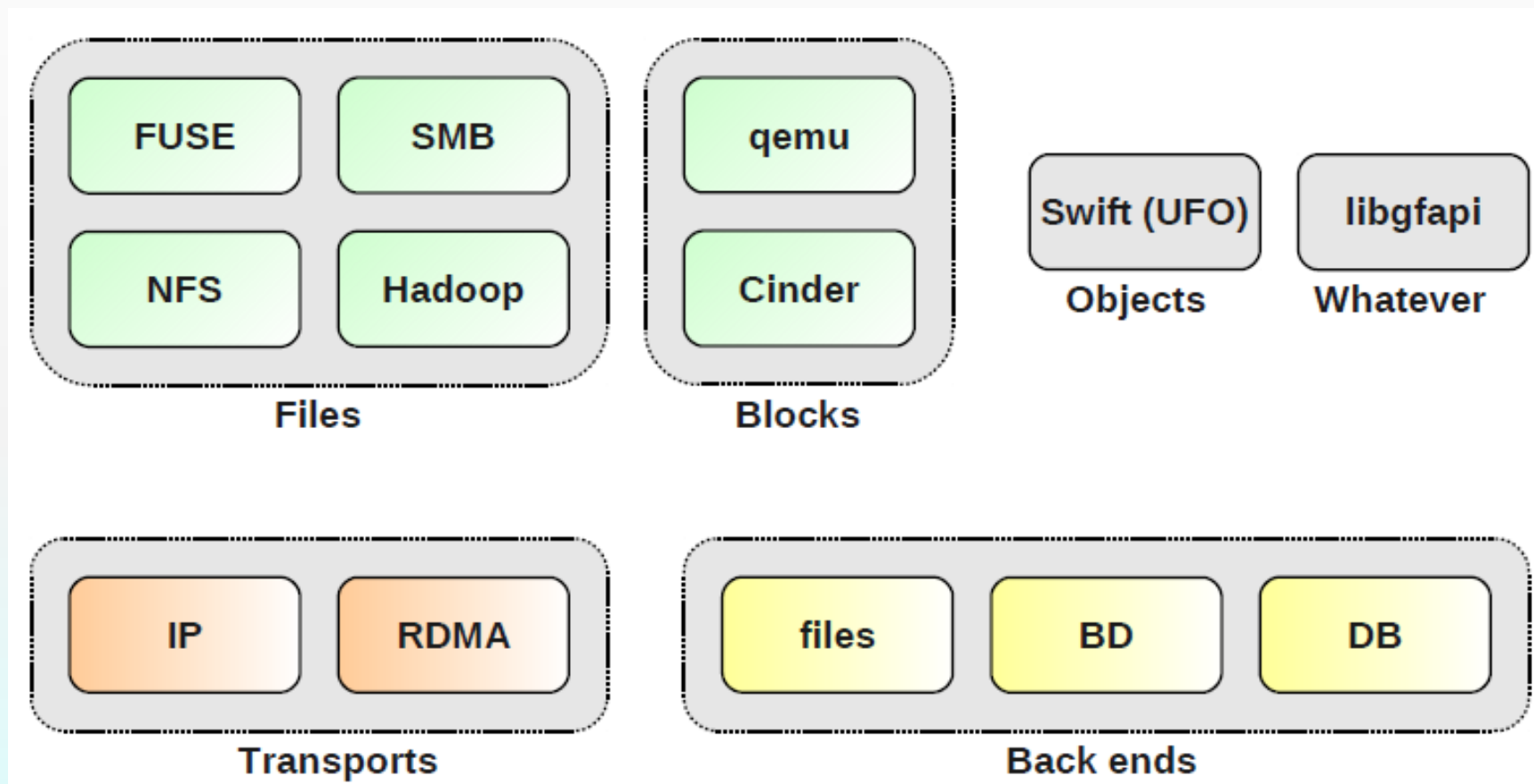


复合卷：哈希 + 条带 + 复制

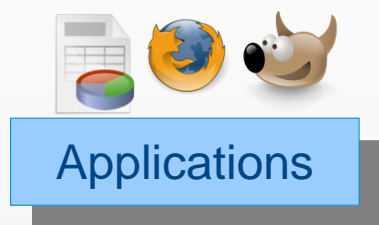
- 三种基本卷的复合卷
- 通常用于类Map Reduce应用



GlusterFS 访 问 接 口



GlusterFS – FUSE Architecture



Glibc

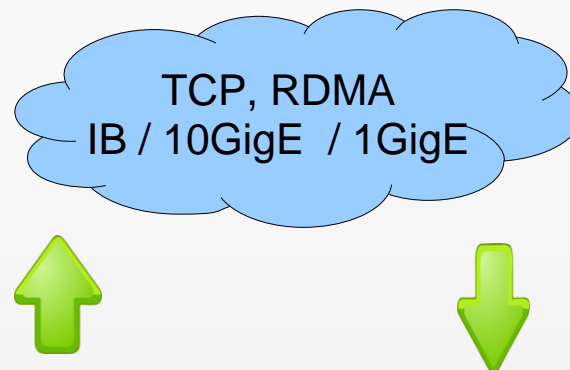
Userland

GlusterFS

Kernel

VFS

FUSE



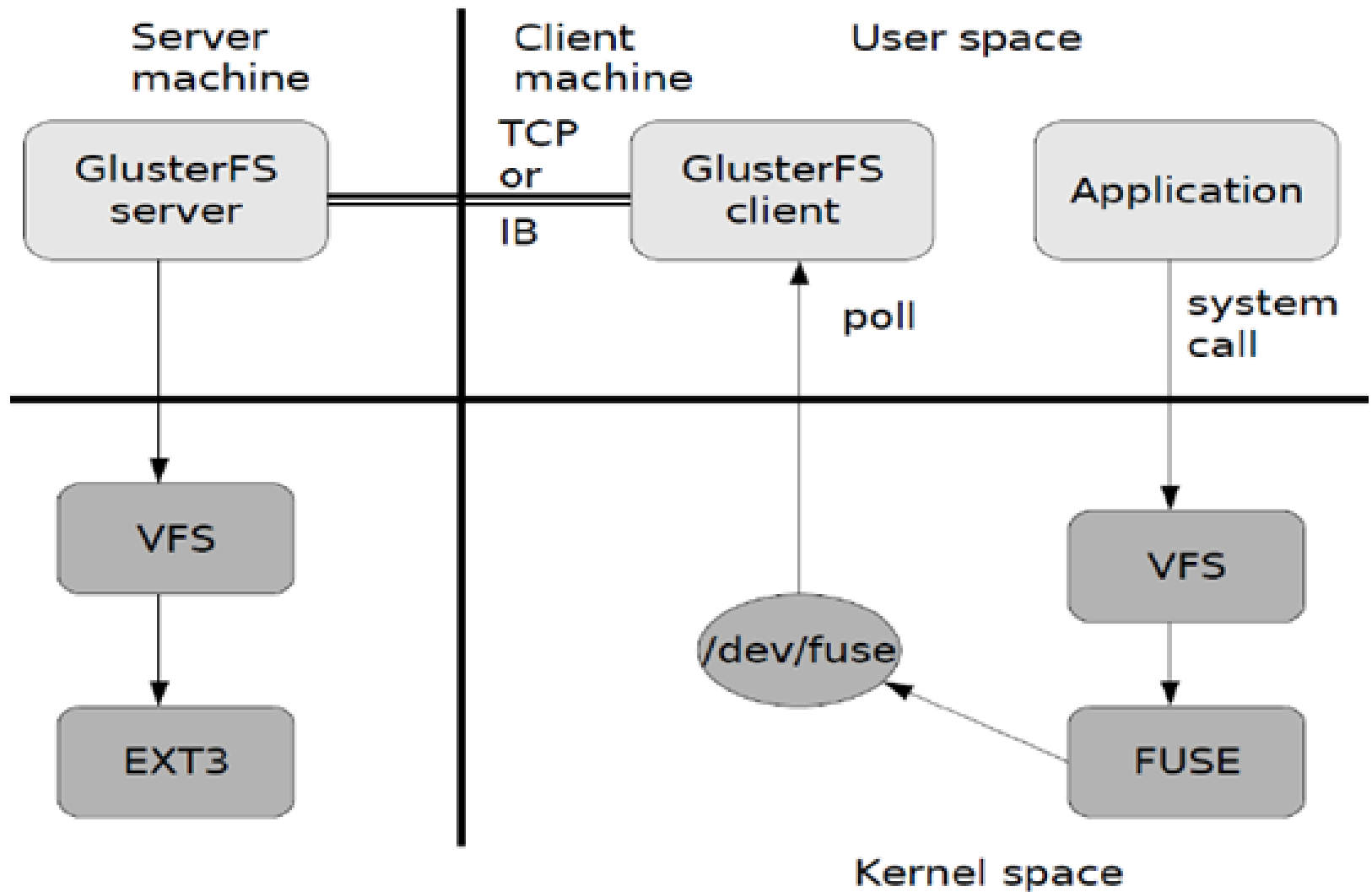
Userland

GlusterFS

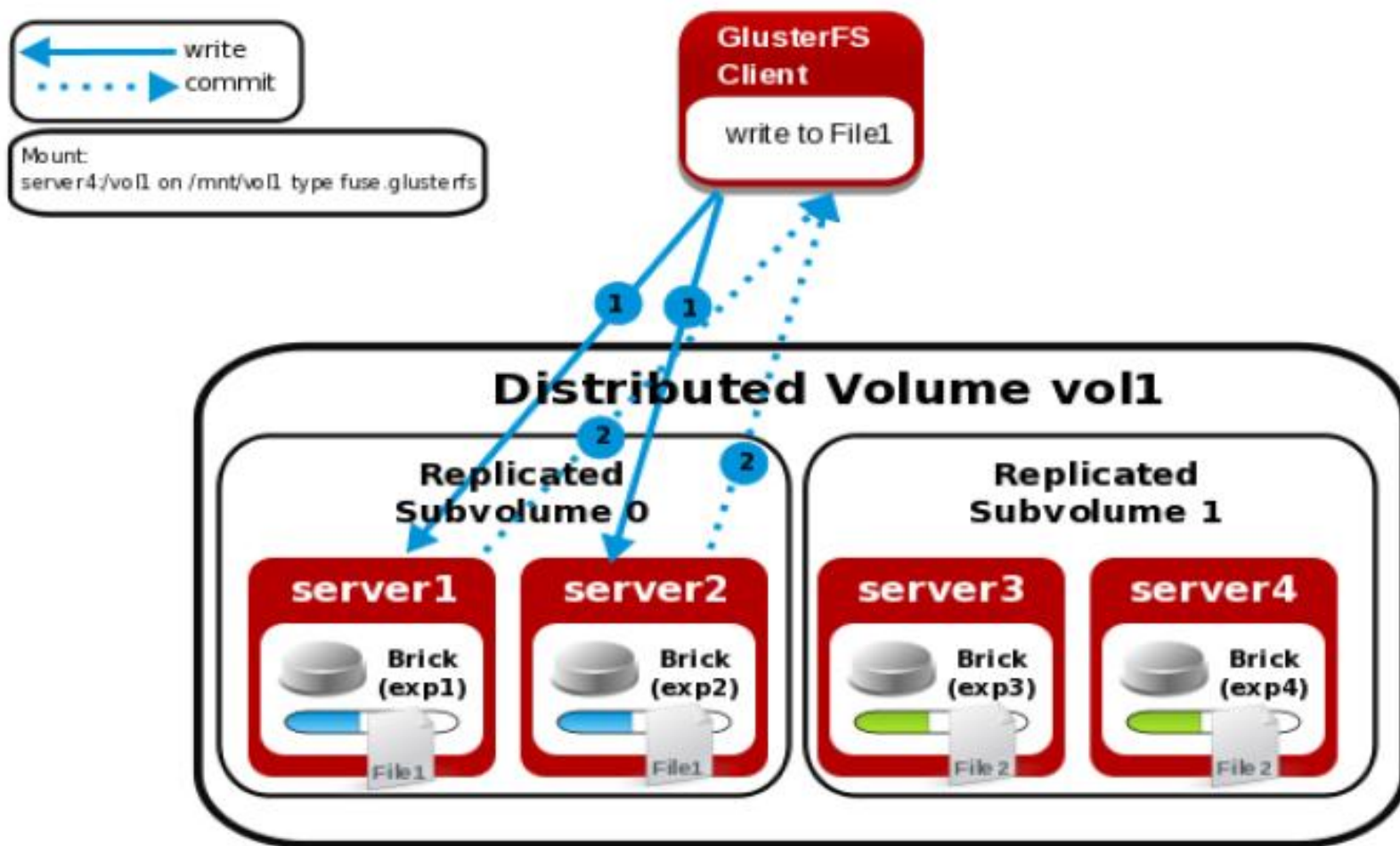
Kernel

Disk FS
(ZFS / Ext3 / XFS)

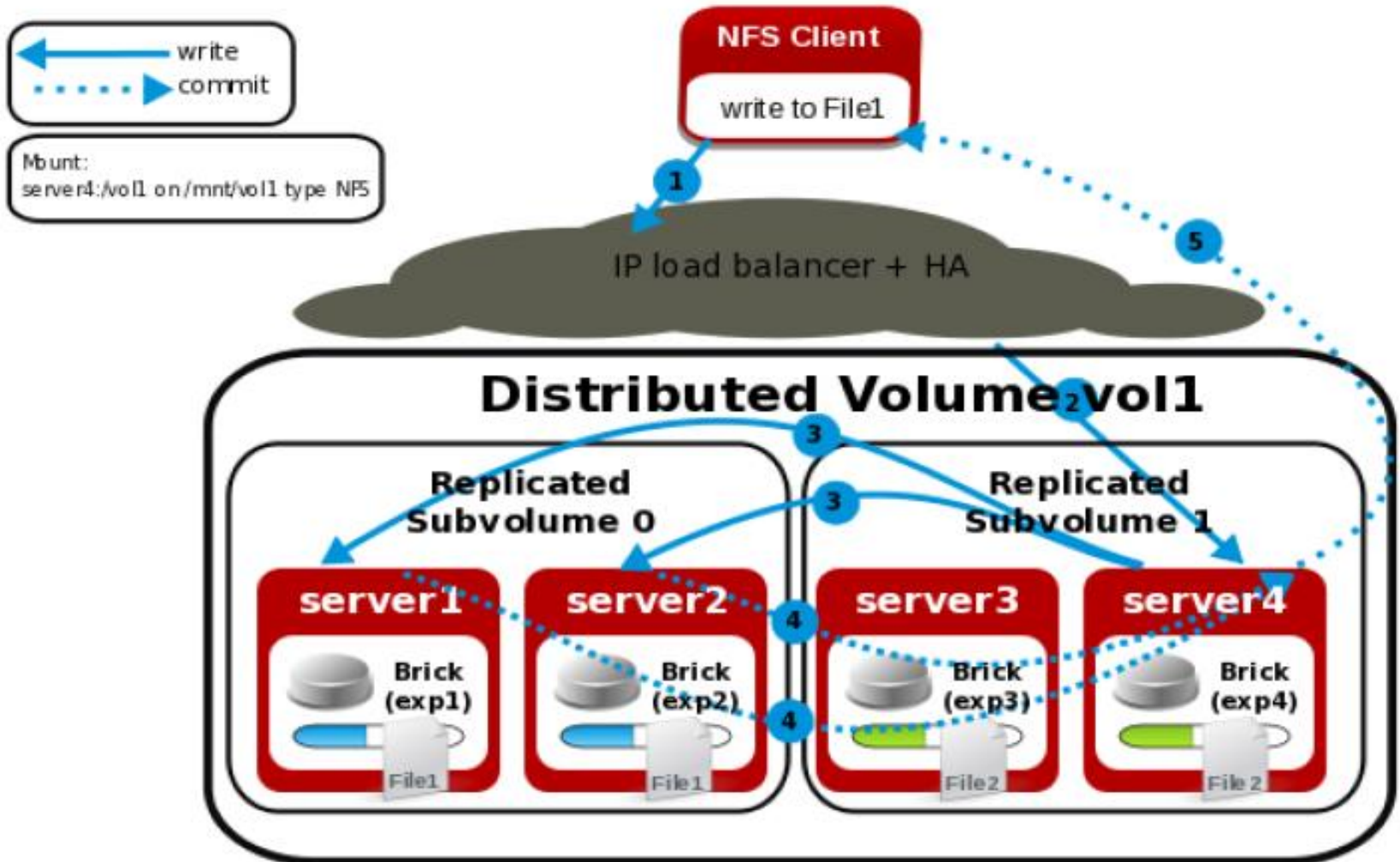
GlusterFS 数据流



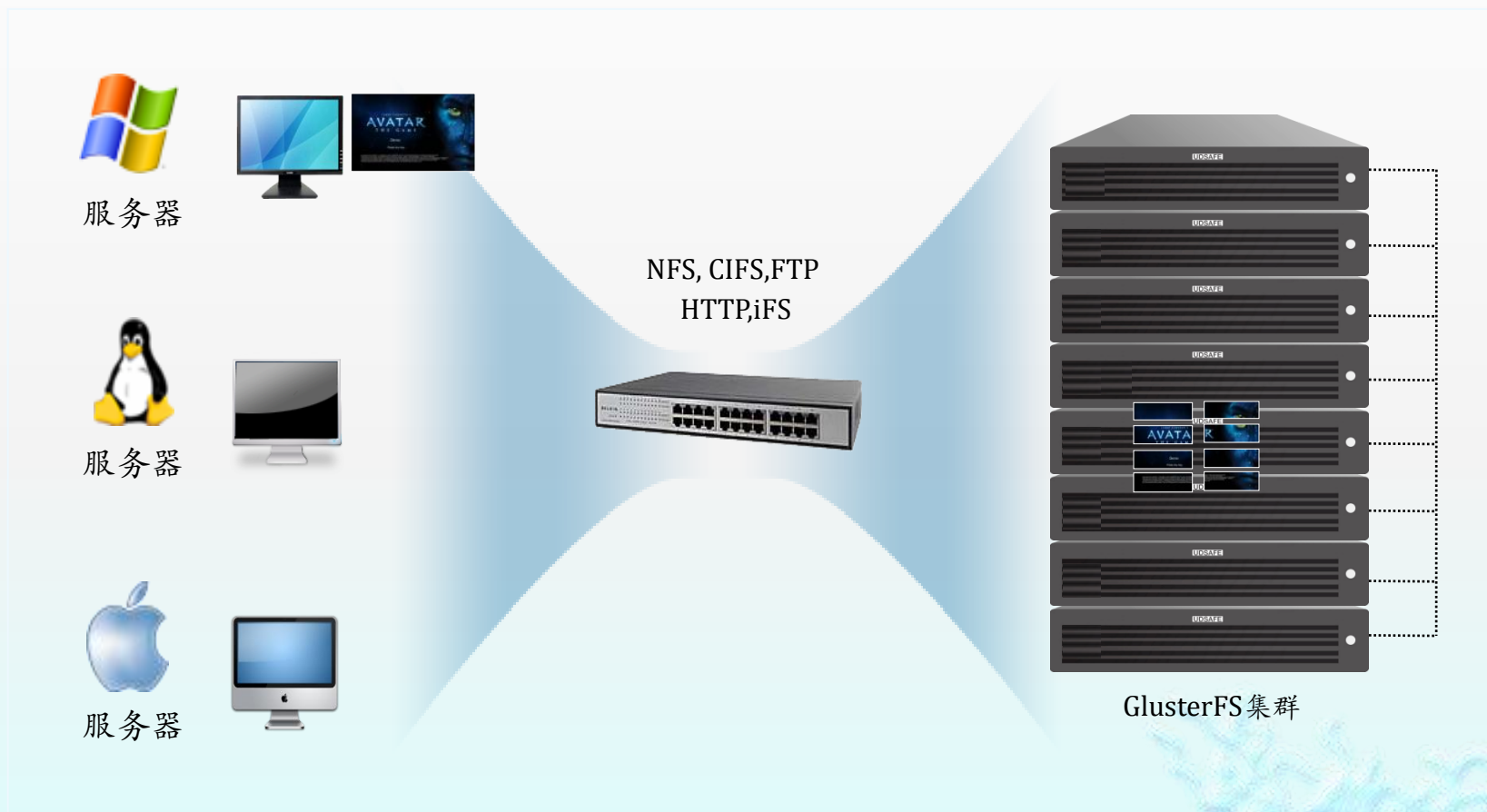
FUSE 访问



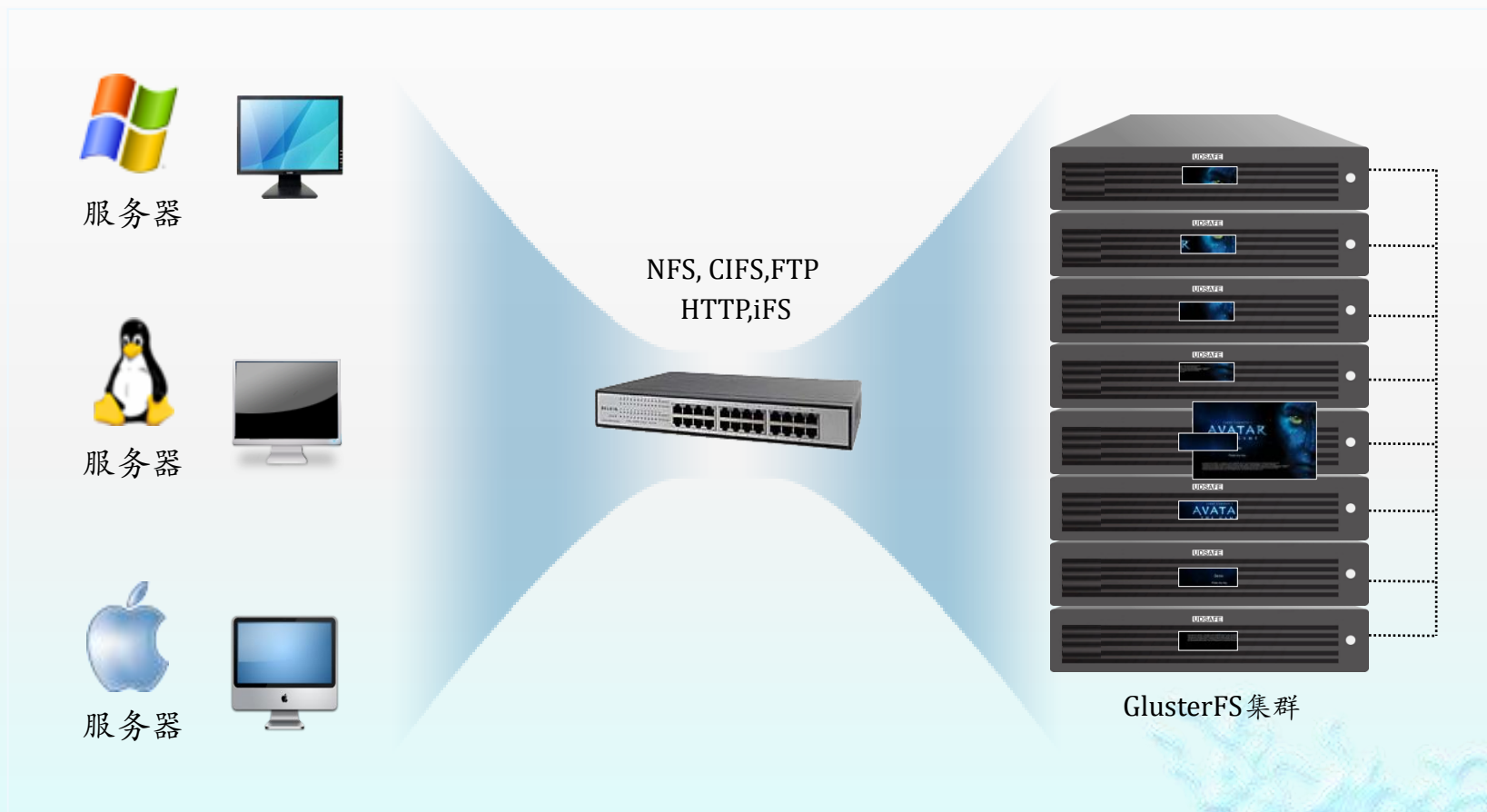
NFS/CIFS 访问



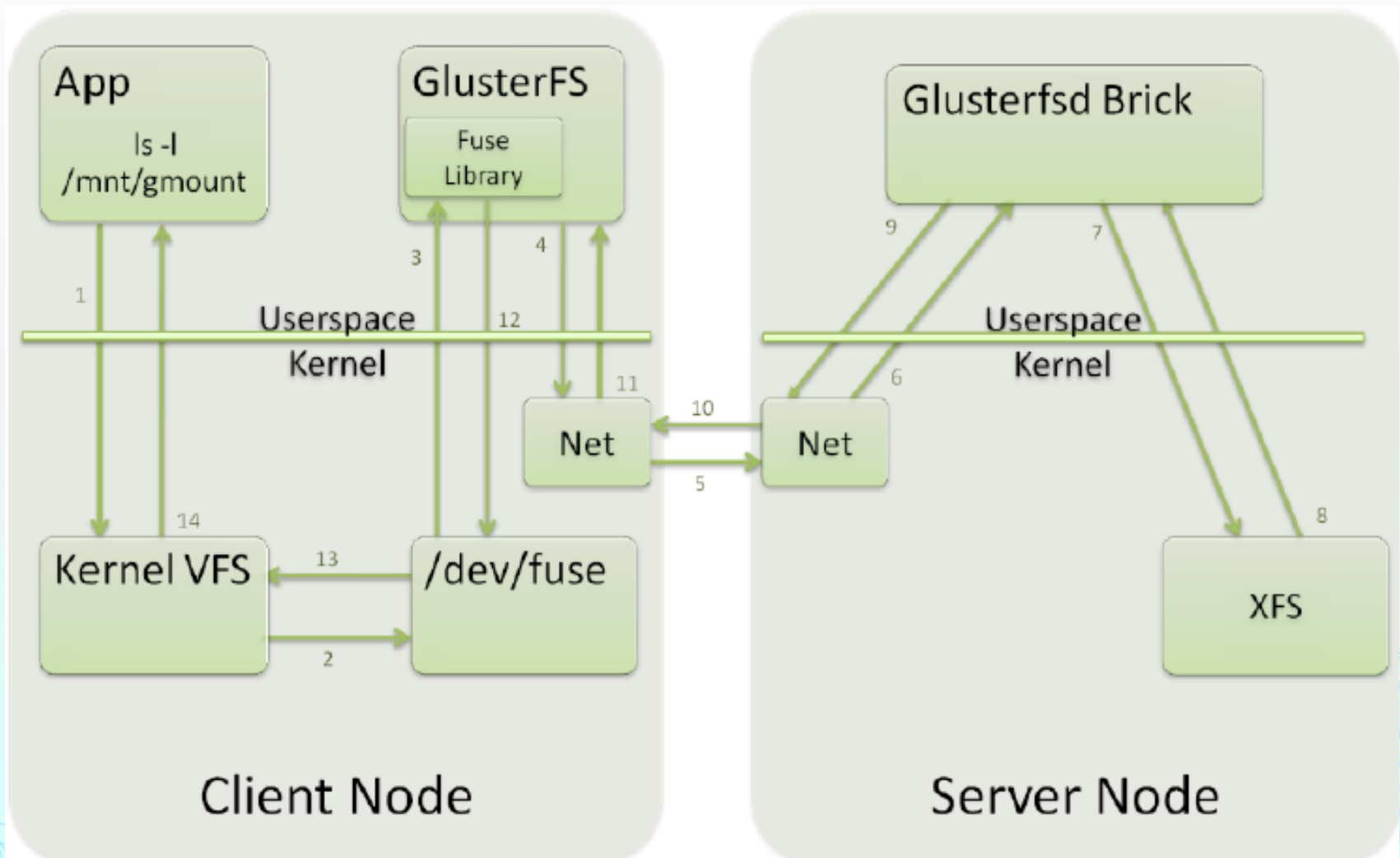
写入一个文件



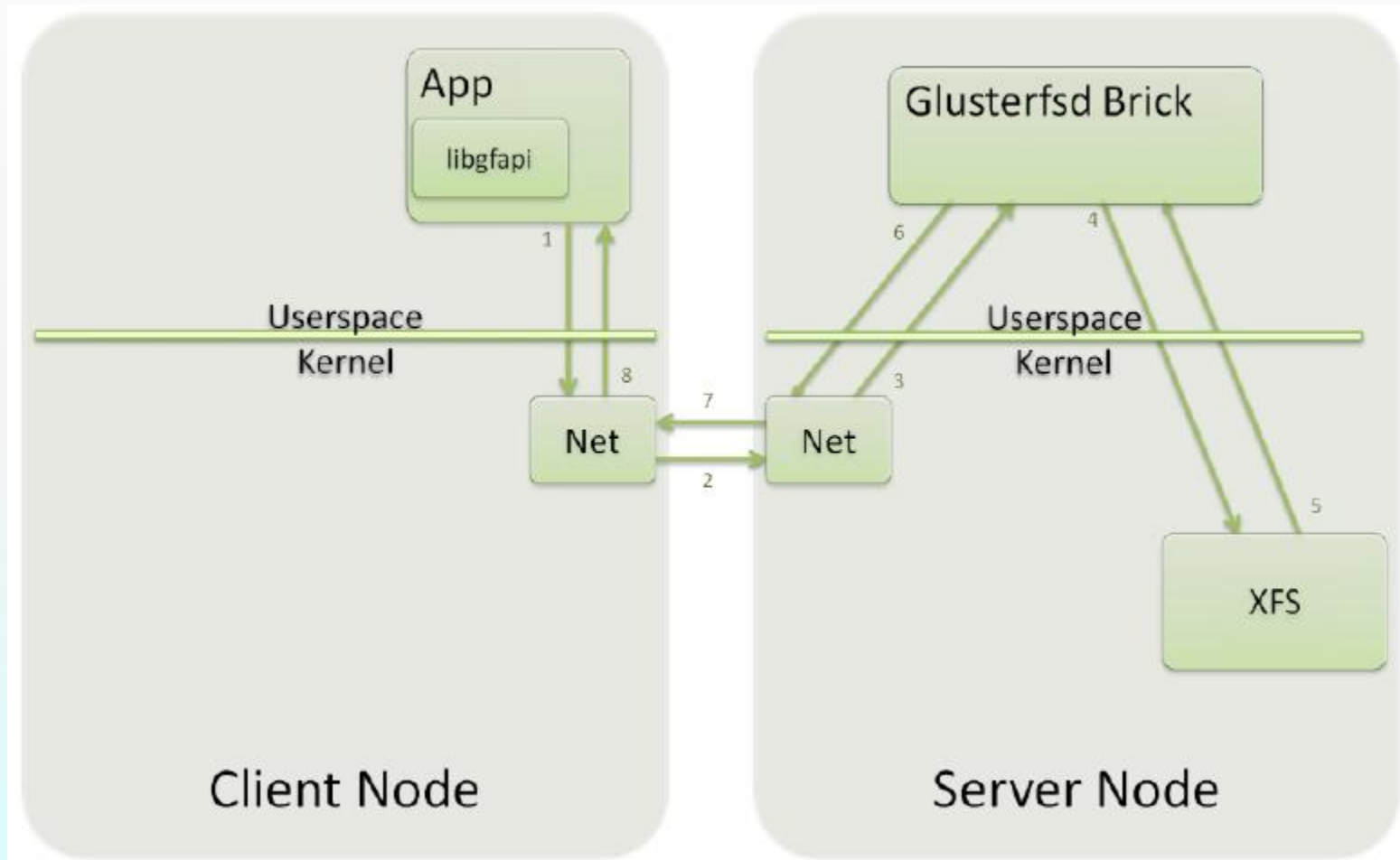
读取一个文件



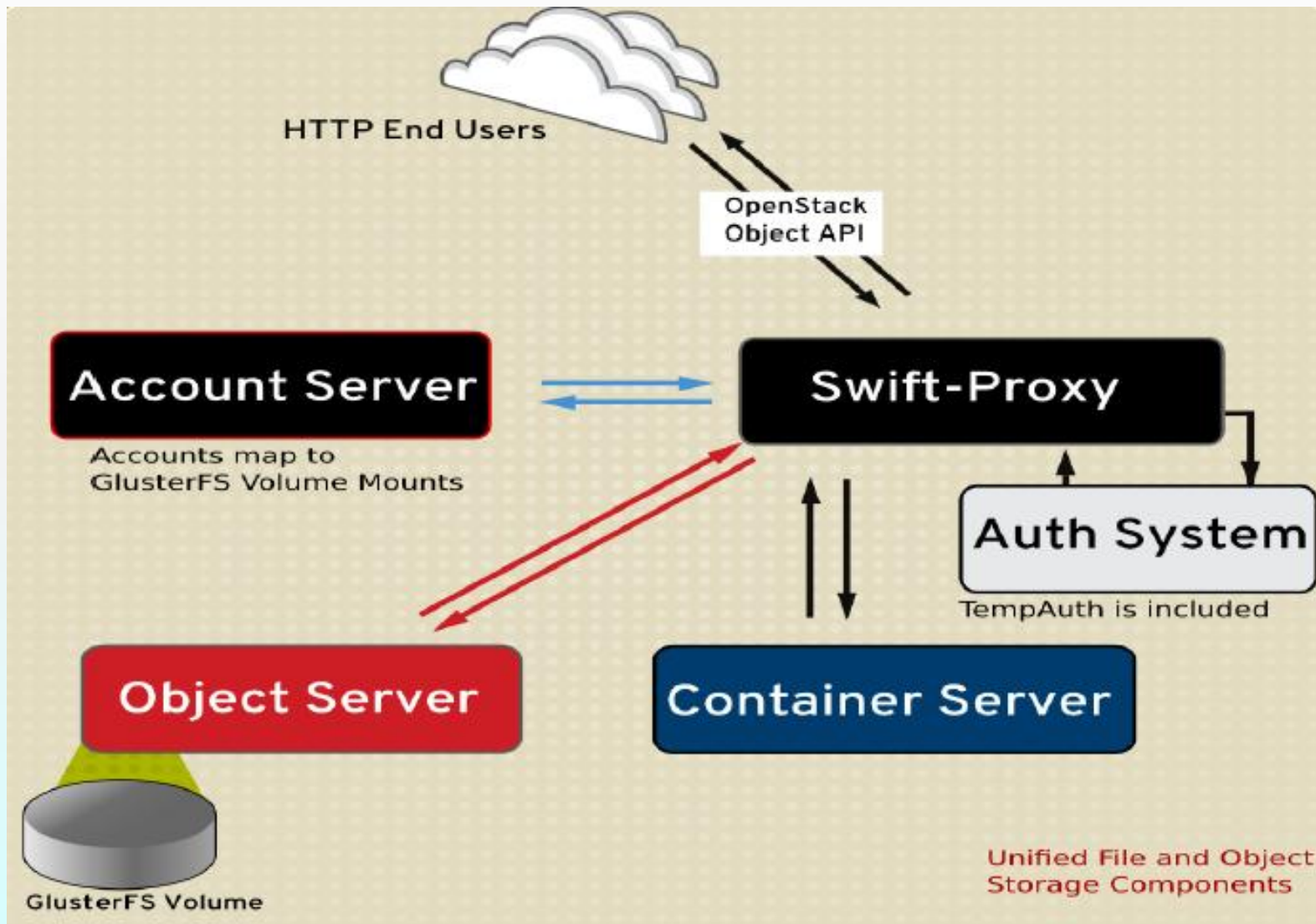
FUSE w/ Libgfapi 访问



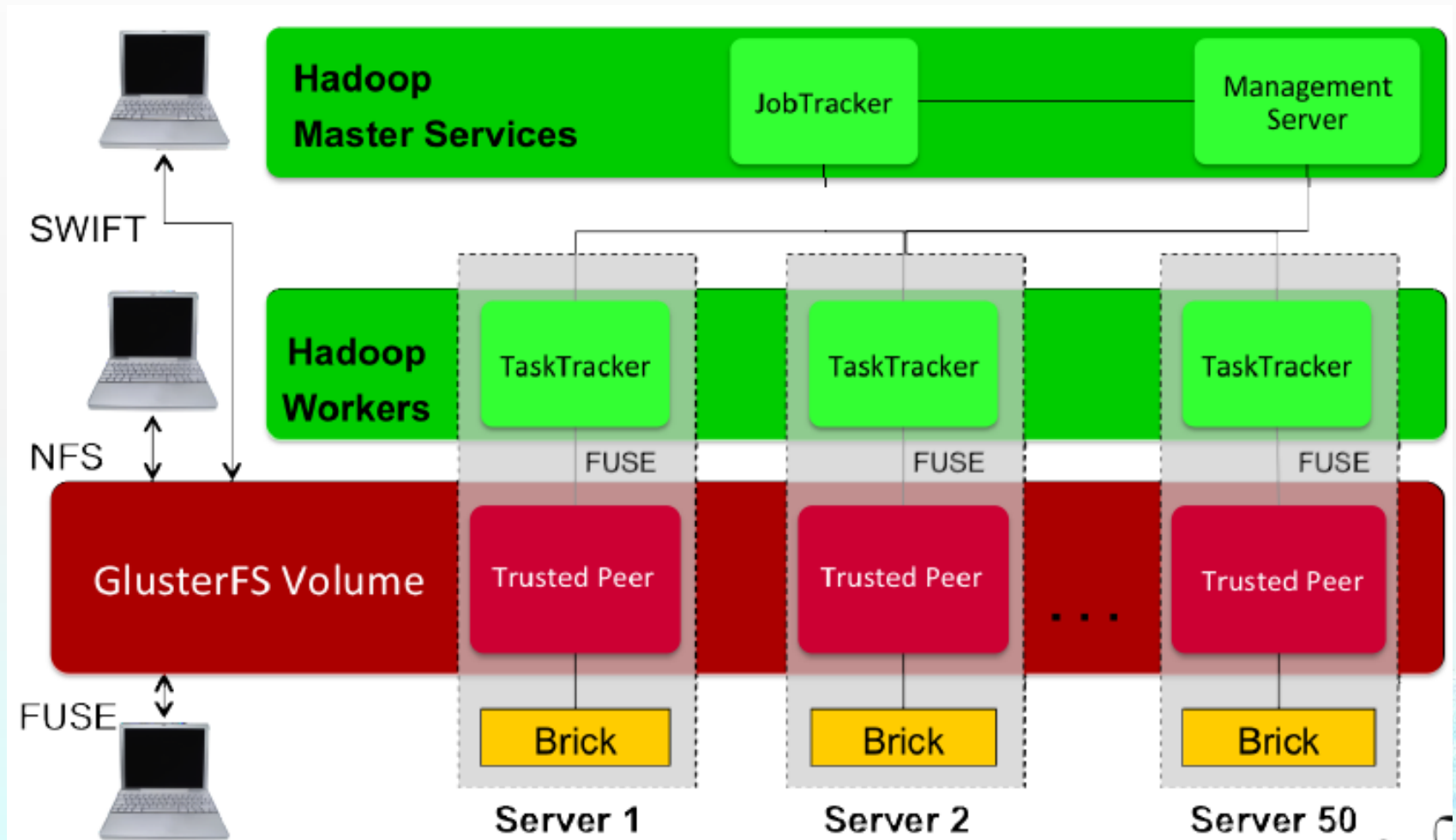
libgfapi 访问



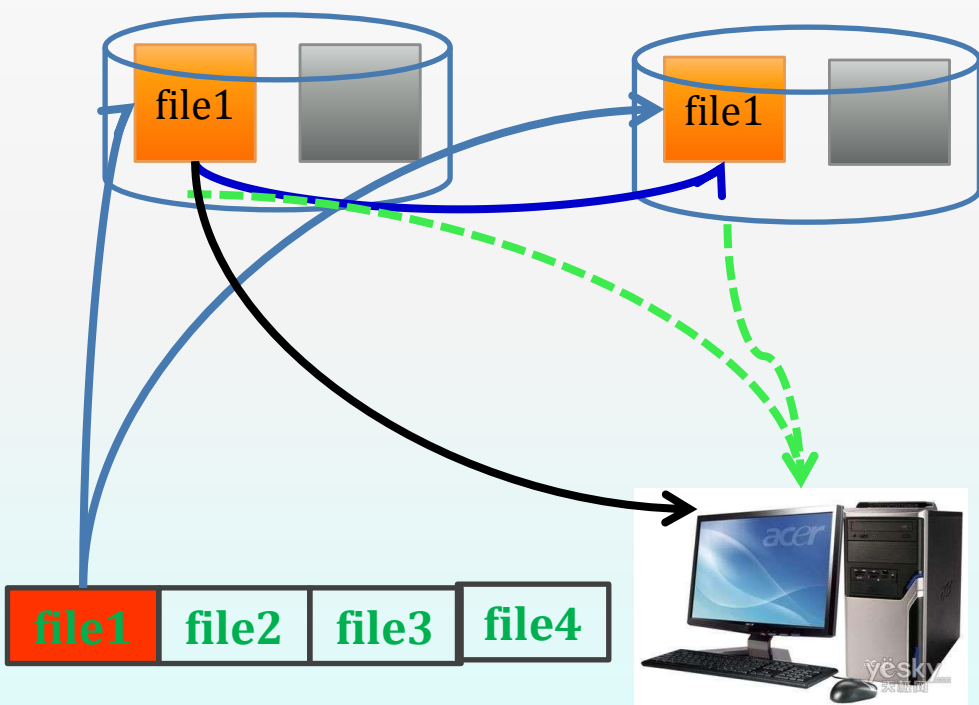
REST API 访问



Hadoop 访问



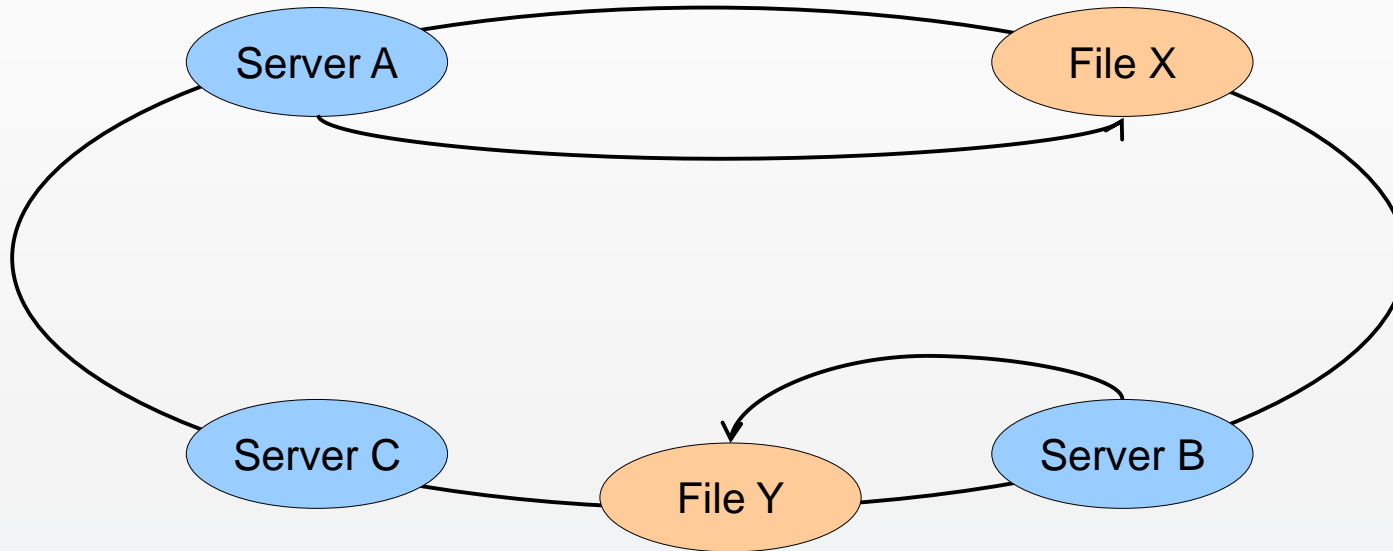
数据自修复 Self-heal



Self-heal发展

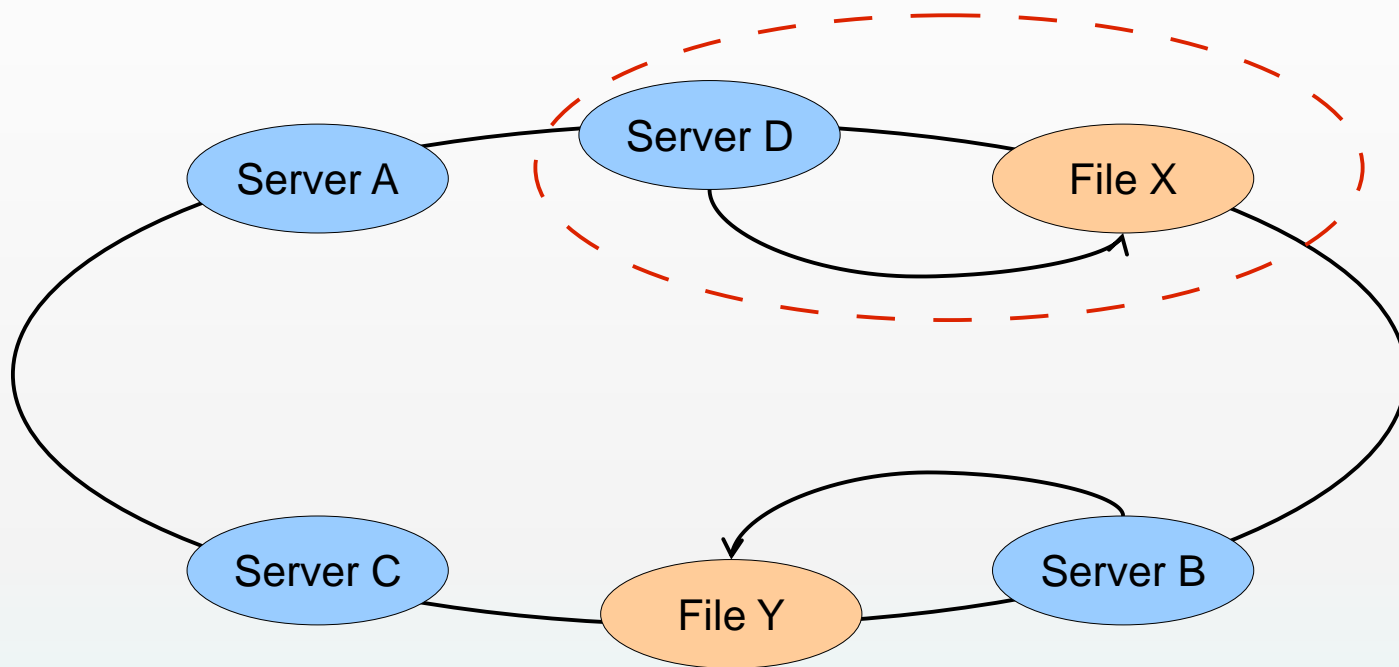
- ◆ 第一代：按需同步进行
 - ◆ 第二代：完全人工扫描
 - ◆ 第三代：并发自动修复(3.3)
 - ◆ 第四代：基于日志
-
- ◆ 镜像卷文件副本保持一致性
 - ◆ 触发时机：访问文件目录时
 - ◆ 判断依据：扩展属性
 - ◆ 脑裂问题：报错或按规则处理

Distributed Hash Table (DHT)



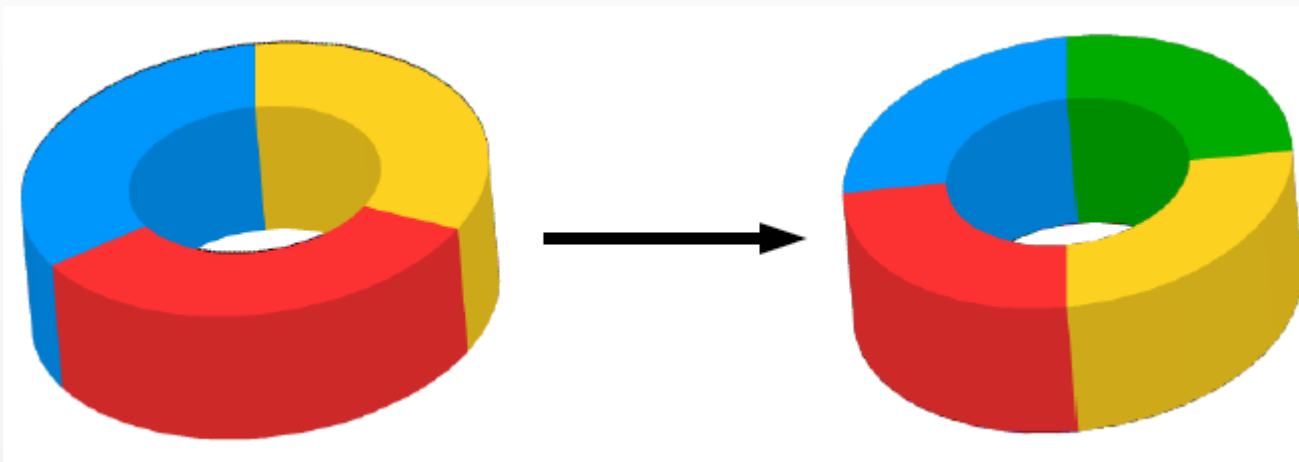
- GlusterFS弹性扩展的基础
- 确定目标hash和brick之间的映射关系

添加节点



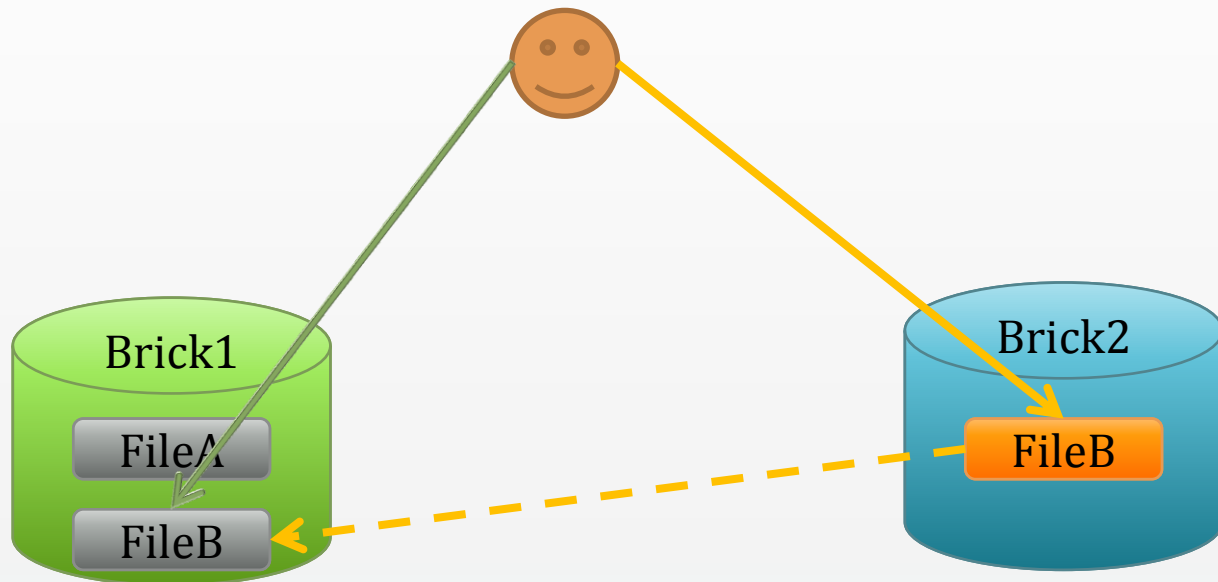
- 添加新节点，最小化数据重新分配
- 老数据分布模式不变，新数据分布到所有节点上
- 执行rebalance，数据重新分布

容量负载均衡



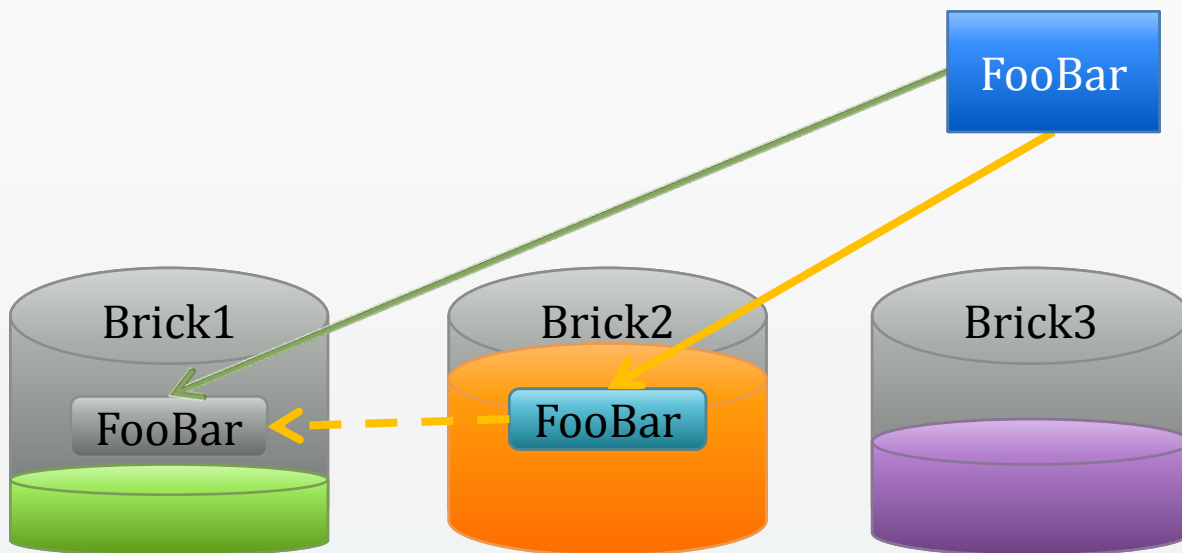
- Hash范围均衡分布，节点一变动全局
- 目标：优化数据分布，最小化数据迁移
- 数据迁移自动化、智能化、并行化

文件更名



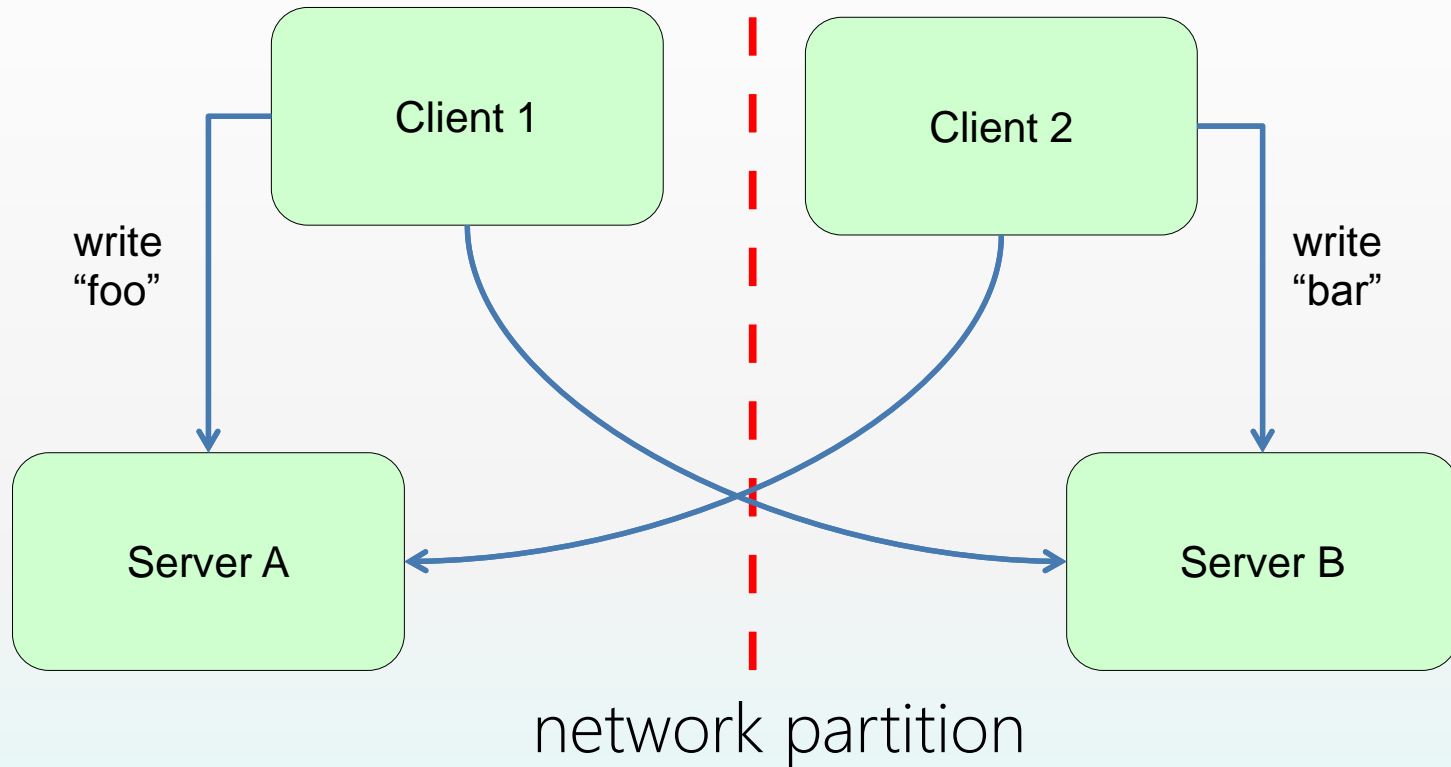
- 文件更名: FileA \rightarrow FileB
- 原先的hash映射关系失效, 大文件难以实时迁移
- 采用文件符号链接, 访问时解析重定向

容量负载优先



- 设置容量阈值，优先选择可用容量充足brick
- Hash目标brick上创建文件符号链接
- 访问时解析重定向

Split Brain



- 裂脑如何产生的？
- 解决方法：1、报错处理；2、Quorum方法 ($N=2?$)；3、仲裁机制

Q & A

