

# 深入浅出Lustre Distributed Namespace特性

作者	时间	QQ技术交流群
<a href="mailto:perrynzhou@gmail.com">perrynzhou@gmail.com</a>	2022/09/17	672152841



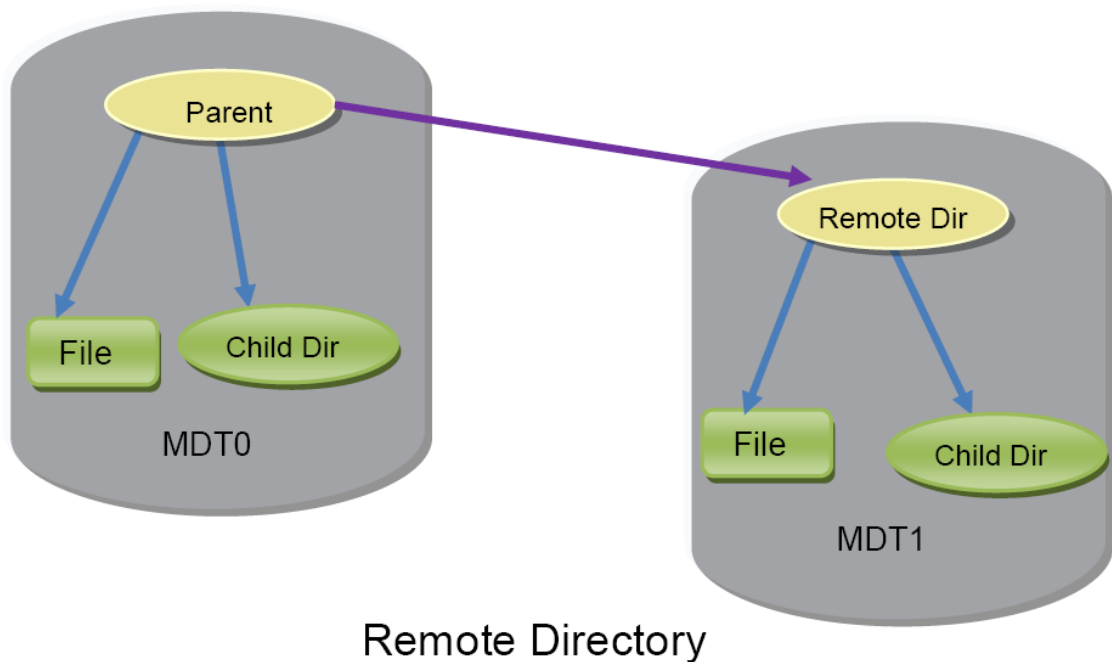
存储内核技术交流

微信扫描二维码，关注我的公众号

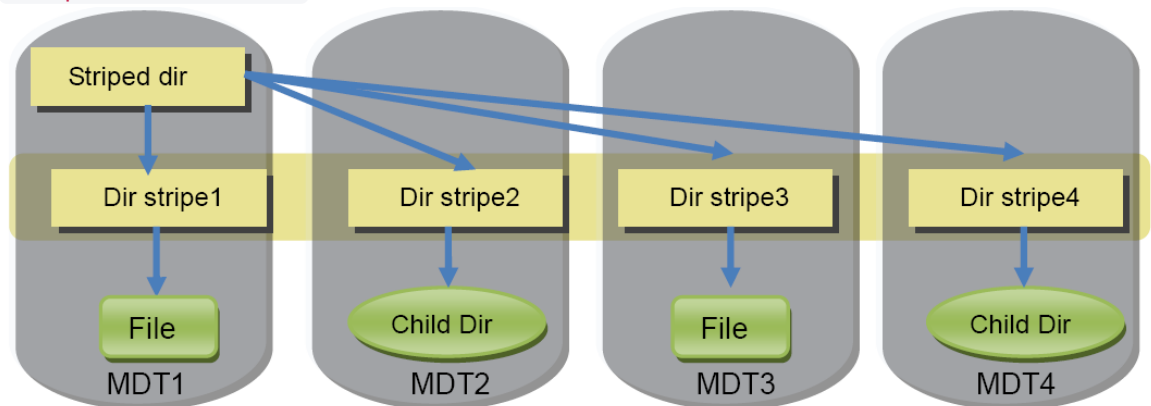


## Distributed Namespace(DNE)设计意图

- 在没有DNE功能之前Lustre文件系统元数据查询和操作仅支持在一个MDT上，如果应用的Workload对于元数据操作要求高，单个MDT不管从容量或者性能上考虑，这些都将会是一个瓶颈点。基于这些因素考虑，Lustre设计DNE来解决MDT的扩展性和性能问题。DNE核心解决了文件或者目录的元数据可以分散存储在多个MDT上，应用请求到不同的文件或者目录元数据请求会分摊到不同的MDT上，MDT0存储了Lustre文件系统的根。
- DNE设计分为2个阶段.阶段一解决根目录下的子目录可以分布到多个MDT，这个就涉及到Remote Directory.阶段二解决了给定的目录，把目录的实际存储数据分布到多个MDT上，这个就涉及到了Striped Directories。
  - Remote Directory



- **Striped Directories**



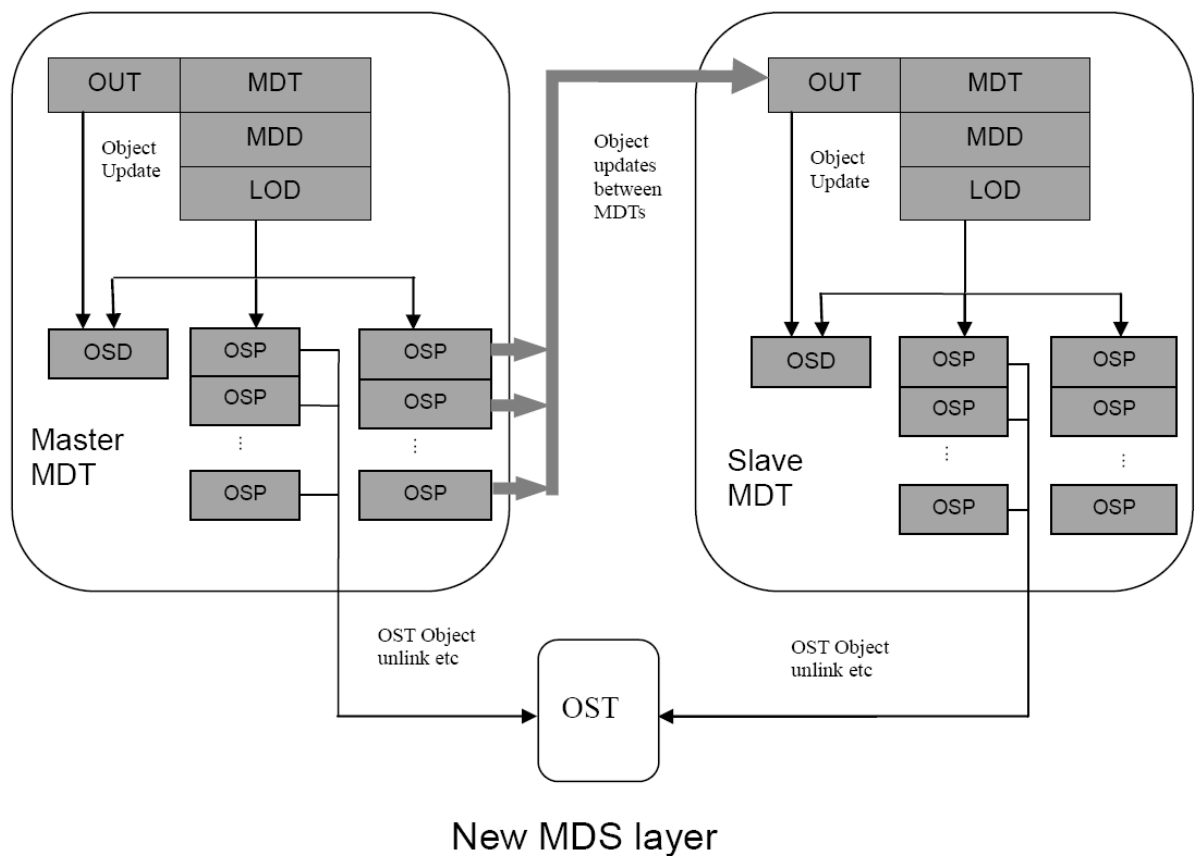
## DNE使用命令

- 管理员可以通过 `lfs mkdir -i {mdt_index} {remote_dir}` 命令在第 `{mdt_index}` mdt 上创建 Remote Directory; 通过 `rmdir {remote_dir}` 删除远程目录
- 设置普通用户在其他MDT(非MDT0)上创建远程目录需要调整如下参数

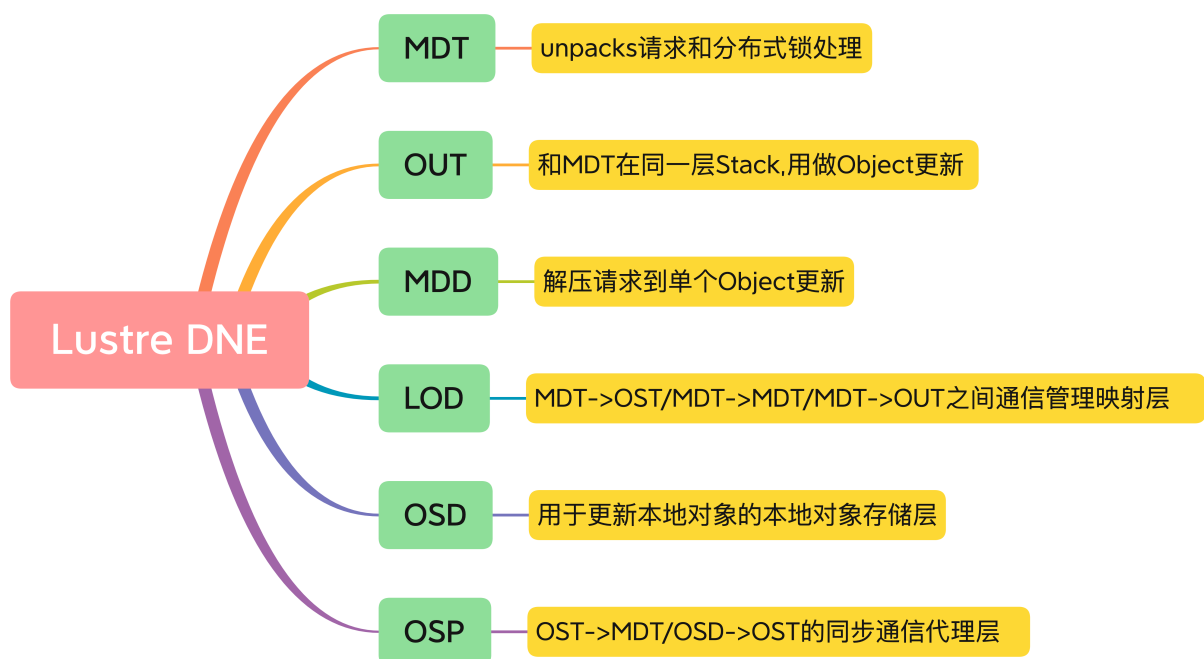
```
$ lctl set_param mdt.fsname-MDT0000.enable_remote_dir=1
$ lctl set_param.mdt.fsname-MDT0000.enable_remote_dir_gid=allowed_gid
```

## DNE架构

- DNE架构是基于原来MDS Stack基础设计，其架构中OSD是本地MDS的后端存储，一部分的OSP直接连接到OST；一部分OSP直接连接到远端的MDT上。



## • DNE架构Stack介绍



## DNE处理流程

- MDT **unpacks** 请求获取分布式锁
- MDD分割操作为多个对象的更新.MDD和LOD并不知道这些对象是本地还是远程更新,这些更新操作都会在一个事务中记性, 要么全部成功, 要么全部失败
- LOD检查此次更新时本地还是远程的更新操作, 如果查询FLD发现是远程的更新就发送给本地的OSD进行更新; 如果发现是远程更新, 则通过OSP发送给其他的远程MDT进行

更新。

- 这里检查更新的对象是本地FID还是远程FID,还是需要查询FLDB。FID一般会在Directory Entry和Inode Extended Attribute中存储。FID在FLDB中查询到后能够识别到这个FID是在哪个MDT上(MDT INDEX).下面是MDT0上的一个Remote Name Entry和远程MDT1之间的映射关系。

