

Распределенные ФС

Подкопаев Антон, podkoav239@gmail.com

СП, СПбГУ

14 октября 2013

- Модель данных, программные компоненты, персистентные структуры данных и API
- Абстракция для доступа к данным, находящимся на физическом носителе
- Традиционная модель
 - Файл — объект с именем и некоторым содержанием
 - Каталог — список файлов и подкаталогов
 - Каталоги и файлы — пространство имен
 - Файл уникально идентифицируется путем

Для ФС важна метайнформация — название файла, список блоков, время модификации, права доступа

- Некоторая интеграция с ядром ОС
- Данные на локальном HDD
- Блоки размером в несколько килобайт
- Кеширование страниц

- Компоненты распределены по разным машинам
- Распределенность существенно влияет на принимаемые решения

Ваш К.О.

- **Клиент**
 - API прикладных приложений и код для коммуникации с сервером
- **Сервер данных**
 - Содержимое файлов
- **Сервер метаданных**
 - Информация о местоположении файла и еще кое-что

- Прозрачность размещения файлов
- Совместный доступ
- Кеширование
- Репликация
- Масштабируемость

- Прикладному ПО известен только путь
- Чем меньше информации о физическом местоположении, тем лучше

- Централизованная ФС
 - Атомарные чтения и запись, блокировки, журналирование
- Распределенная ФС
 - Сетевые задержки, репликация усложняет жизнь

- Синхронные чтение/запись
- Write-through cache
- Файлы неизменяемые после создания
- Append-only
- Уведомление о изменениях для клиентов, открывших файл
- Полноценные транзакции

- Синхронная и асинхронная
- Политика согласованности реплик
- Запись в реплики

- Стремимся к линейной
 - было N дисков и K машин
 - стало на $2N$ данных — добавили N дисков, сохранили пропускную способность
 - добавили K машин — получили в два раза большую пропускную способность
- На практике есть препятствия
 - пропускная способность сети, сетевых интерфейсов серверов, производительность сервера метаданных, блокировки

Подкопаев
Антон

Основные
концепции

Особенности
РФС

Обзор

NFS

AFS

CIFS

GFS

GlusterFS

NFS у всех на слуху...

Подкопаев
Антон

Основные
концепции

Особенности
РФС

Обзор

NFS

AFS

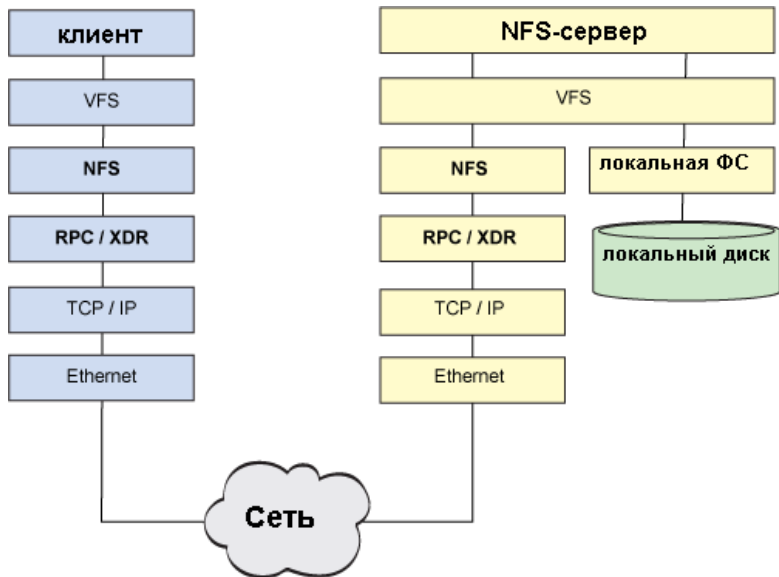
CIFS

GFS

GlusterFS

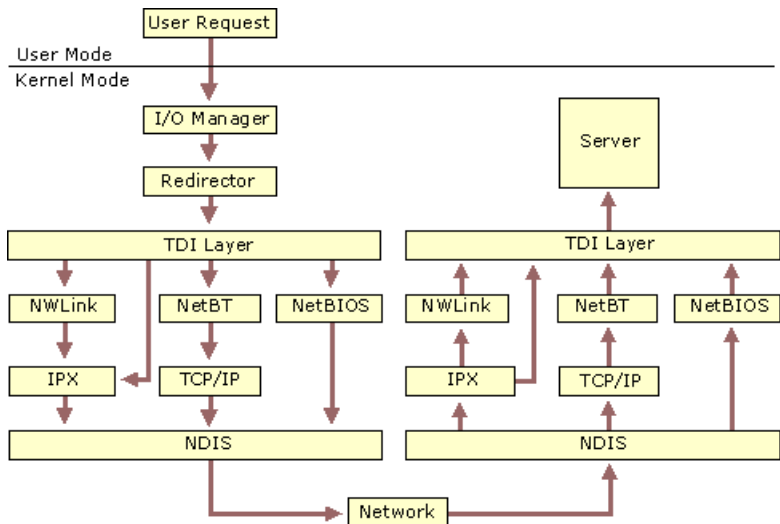


- Network File System
- Sun, 1984
- POSIX API
- На сервере NFS работает с интерфейсом файловой системы
- Поддержка блокировок и сессионного кеширования



- Andrew File System
- Carnegie Mellon, 1980-е
- Сессионное кеширование, нотификации об изменениях, блокировки файлов
- Моментальные read-only снимки томов
- Не POSIX API

- Common Internet File System
- ...aka SMB
 - Server Message Block
 - Samba
- ...aka Windows Shared Folders
- IBM, 1983
- Microsoft, 1996
- Уступающие блокировки, блокировки файлов
- Аутентификация коллективного доступа
- Уведомление об изменении каталога



- Google File System
- Начало 2000-ых

Предпосылки

- Большие файлы (N Gb) записываются и читаются пакетными процессами (creawler, indexer)
- Пропускная способность важнее случайного доступа
- Ширпотребные компьютеры

- Много файловых серверов, один активный сервер метаданных (мастер)
- Файлы хранятся фрагментами по 64 Mb
- Три реплики каждого фрагмента на различных файловых серверах
- Приоритетные операции с файлом
 - Большое последовательное чтение
 - Конкурентное наращивание
- Кеширование на клиенте не производится
- Не POSIX API

- Ячейка — единица развертывания
- В ячейке один мастер и много файловых серверов
- Ячейка GFS соответствует физическому датацентру

- Приложение собирается прочитать фрагмент
- GFS библиотека звонит мастеру, тот возвращает адреса **реплик** — файловых серверов, хранящих фрагмент
- GFS библиотека напрямую звонит одному из файловых серверов с просьбой вернуть нужный диапазон внутри данного фрагмента
- Далее прямое общение клиента и файлового сервера

- Google 2010-х годов — интерактивные приложения
- Файлы меньше в размерах и больше в количестве
- Требование во времени произвольного доступа жестче

- Google 2010-х годов — интерактивные приложения
- Файлы меньше в размерах и больше в количестве
- Требование во времени произвольного доступа жестче

GFS в Google больше не используется, на смену пришел Colossus

- Реализация Google File System закрыта
- Открытые проекты с аналогичной архитектурой
 - Apache HDFS: реализация на Java из проекта Hadoop
 - QFS: реализация на C++

- OpenSource
- Сервер данных в том числе выполняет и функции сервера метаданных

Подкопаев
Антон

Основные
концепции

Особенности
РФС

Обзор

NFS

AFS

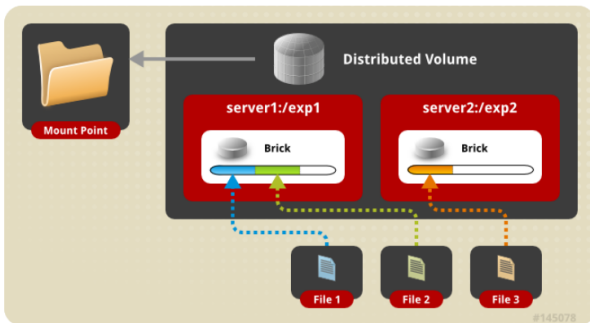
CIFS

GFS

GlusterFS

- Brick
- Логический диск
- Доверенные хранилища

- Распределенные
 - Реплицируемые
 - Разделяющие
-
- Распределенные разделяющие
 - Распределенные реплицируемые
 - Разделяющие реплицируемые

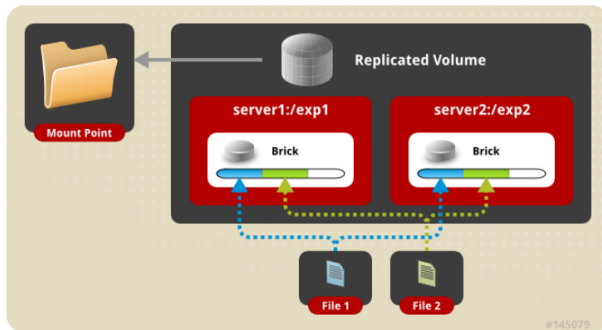


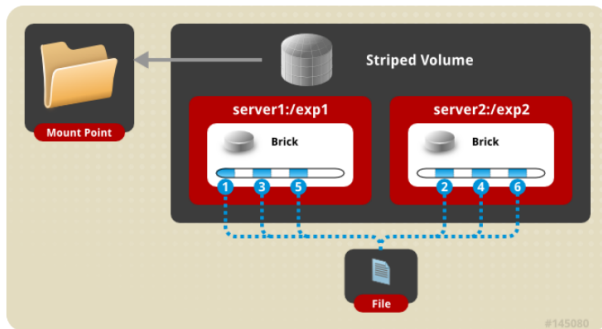
Плюсы

- Больше серверов => выше производительность при параллельном доступе
- Увеличение диска = добавление сервера (можно во время работы)

Минусы

- Потеря сервера = потеря данных на нем
- Файл не может быть больше размера узла
- Смена имени файла => дополнительное время на lookup





```
apt-get install glusterfs-server  
service glusterfs-server start  
service glusterfs-server status
```

```
apt-get install glusterfs-server  
service glusterfs-server start  
service glusterfs-server status
```

```
mkfs.xfs disk-image  
mount disk-image gluster_disk
```

```
apt-get install glusterfs-server  
service glusterfs-server start  
service glusterfs-server status
```

```
mkfs.xfs disk-image  
mount disk-image gluster_disk
```

```
gluster peer probe 192.168.1.105
```

```
gluster volume create gv0 replica 2  
192.168.1.114:/home/us1/gluster_disk  
192.168.1.105:/home/us2/gluster_disk
```

```
gluster volume start gv0  
gluster volume info  
mount -t glusterfs 192.168.1.105:/gv1 /mnt
```

Распределенные

```
gluster volume create test-volume server1:/exp1  
server2:/exp2
```

Реплицируемые

```
gluster volume create test-volume replica 2  
server1:/exp1 server2:/exp2
```

Разделяющие

```
gluster volume create test-volume stripe 2  
server1:/exp1 server2:/exp2
```

```
gluster peer probe new_server
```

```
gluster peer probe new_server  
gluster volume add-brick vol_name new_brick
```

```
gluster peer probe new_server  
gluster volume add-brick vol_name new_brick  
gluster volume info
```



```
gluster peer probe new_server  
gluster volume add-brick vol_name new_brick  
gluster volume info
```

Для расширения реплицируемых (разделяющих) дисков надо добавлять количество серверов, кратное фактору реплики (разбиения)

```
gluster volume remove-brick vol_name new_brick  
start
```

```
gluster volume remove-brick vol_name new_brick  
start  
gluster volume remove-brick vol_name new_brick  
status
```

```
gluster volume remove-brick vol_name new_brick  
start  
gluster volume remove-brick vol_name new_brick  
status  
gluster volume remove-brick vol_name new_brick  
commit
```