

# 一、问题引导

---

- 互联网环境中的文件如何存储？(待完善)  
不能存储本地、NFS ( mount挂载 )、HDFS、FastDFS、云存储 ( 图床等 )
- 互联网环境中的文件如何进行HTTP访问？(待完善)  
web服务器：Nginx、Apache等

## 二、FastDFS介绍

---

### 1.FastDFS是什么？

---

- FastDFS 是一个使用C编写的开源的高性能**分布式文件系统 ( Distributed File System,简称DFS )**。
- 它由**淘宝**开发平台部资深架构师**余庆**开发。FastDFS孵化平台(ChinaUnix)版块<http://bbs.chinaunix.net/forum-240-1.html>。
- 它对文件进行管理，功能包括：**文件存储、文件同步、文件访问 ( 文件上传、文件下载 )**等，解决了**大容量存储和负载均衡**的问题。
- 特别适合以文件为载体的在线服务，如相册网站、视频网站、电商网站等等。特别**适合以中小文件 ( 建议范围：4KB < file\_size <500MB )**为载体的在线服务。
- FastDFS为互联网量身定制，充分考虑了**冗余备份、负载均衡、线性扩容**等机制，并注重**高可用、高性能**等指标，使用FastDFS很容易搭建一套高性能的文件服务器集群提供文件上传、下载等服务。

## 三、FastDFS架构原理分析 ( 重点 )

---

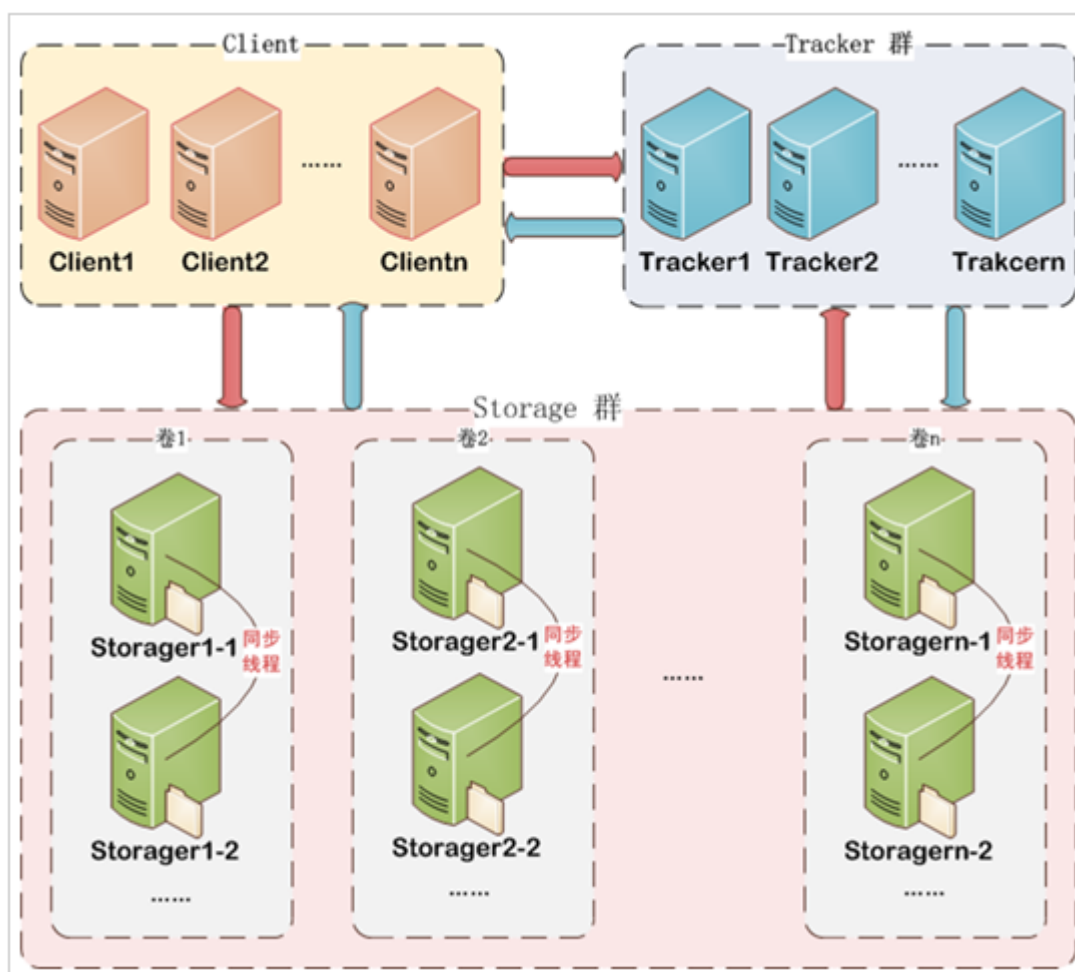
### 1.架构分析

---

FastDFS 系统有三个角色：跟踪服务器(Tracker Server)、存储服务器(Storage Server)和客户端(Client)。

- **Tracker Server**：跟踪服务器
  - 主要做调度工作，并对Storage Server起到负载均衡的作用；
  - 负责管理所有的 storage server和 group，每个 storage 在启动后会连接 Tracker，告知自己所属 group 等信息，并保持周期性心跳。
  - Tracker Server可以有多台，**Tracker Server之间是相互平等关系同时提供服务，Tracker Server不存在单点故障。客户端请求Tracker Server采用轮询方式，如果请求的Tracker无法提供服务则换另一个Tracker。**
- **Storage Server**：存储服务器
  - 主要提供容量和备份服务；
  - 以 group 为单位，不同group之间互相独立，每个 group 内可以有多台 storage server，数据互为备份。
  - **采用分组存储方式的好处是灵活、可控性较强。**比如上传文件时，可以由客户端直接指定上传到的组也可以由Tracker进行调度选择。

- 一个分组的存储服务器访问压力较大时，可以在该组增加存储服务器来扩充服务能力（纵向扩容）。当系统容量不足时，可以增加组来扩充存储容量（横向扩容）。
- **Client**：客户端
  - 上传下载数据的服务器，也就是我们自己的项目所部署在的服务器。



## 2.存储策略

为了支持大容量，存储节点（服务器）采用了**分卷（或分组）**的组织方式。存储系统由一个或多个卷组成，**卷与卷之间的文件是相互独立的**，所有卷的文件容量累加就是整个存储系统中的文件容量。一个卷可以由一台或多台存储服务器组成，**一个卷下的存储服务器中的文件都是相同的**，卷中的多台存储服务器起到了**冗余备份**和**负载均衡**的作用。

**在卷中增加服务器时，同步已有的文件由系统自动完成，同步完成后，系统自动将新增服务器切换到线上提供服务。**当存储空间不足或即将耗尽时，可以动态添加卷。只需要增加一台或多台服务器，并将它们配置为一个新的卷，这样就扩大了存储系统的容量。

## 3.Storage状态收集

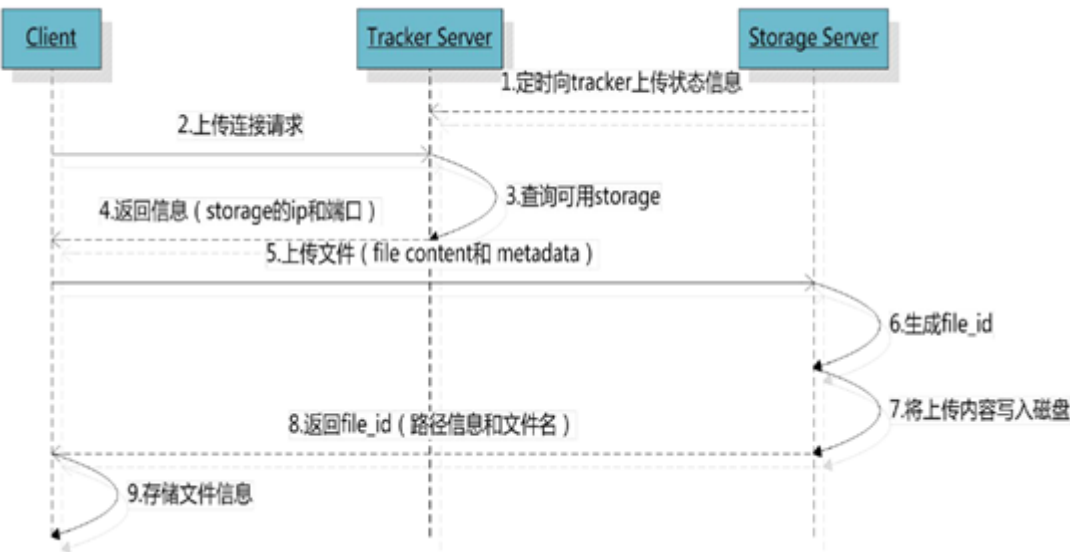
Storage Server会通过配置连接集群中所有的Tracker Server，定时向他们报告自己的状态，包括**磁盘剩余空间、文件同步状况、文件上传下载次数**等统计信息。

storage server有7个状态，如下：

|   |                               |                               |
|---|-------------------------------|-------------------------------|
| 1 | FDFS_STORAGE_STATUS_INIT      | :初始化，尚未得到同步已有数据的源服务器          |
| 2 | FDFS_STORAGE_STATUS_WAIT_SYNC | :等待同步，已得到同步已有数据的源服务器          |
| 3 | FDFS_STORAGE_STATUS_SYNCING   | :同步中                          |
| 4 | FDFS_STORAGE_STATUS_DELETED   | :已删除，该服务器从本组中摘除（注：本状态的功能尚未实现） |
| 5 | FDFS_STORAGE_STATUS_OFFLINE   | :离线                           |
| 6 | FDFS_STORAGE_STATUS_ONLINE    | :在线，尚不能提供服务                   |
| 7 | FDFS_STORAGE_STATUS_ACTIVE    | :在线，可以提供服务                    |

当storage server的状态为 FDFS\_STORAGE\_STATUS\_ONLINE 时，当该storage server向tracker server发起一次 heart beat时，tracker server将其状态更改为 FDFS\_STORAGE\_STATUS\_ACTIVE。

## 4.文件上传流程分析



流程说明：

1、tracker server收集storage server的状态信息

- 1

storage server定时向已知的tracker server（可以是多个）发送磁盘剩余空间、文件同步状况、文件上传下载次数等统计信息
- 2

storage server会连接整个集群中所有的Tracker server，向他们报告自己的状态。

2、选择tracker server

- 1

当集群中不止一个tracker server时，由于tracker之间是完全对等的关系，客户端在upload文件时可以任意选择一个trakcer。

3、选择存储的group

- 1 当tracker接收到upload file的请求时，会为该文件分配一个可以存储该文件的group，支持如下选择group的规则：
- 2 1. Round robin，所有的group间轮询
- 3 2. Specified group，指定某一个确定的group
- 4 3. Load balance，剩余存储空间多多group优先

#### 4、选择storage server

- 1 当选定group后，tracker会在group内选择一个storage server给客户端，支持如下选择storage的规则：
- 2 1. Round robin，在group内的所有storage间轮询
- 3 2. First server ordered by ip，按ip排序
- 4 3. First server ordered by priority，按优先级排序（优先级在storage上配置）

#### 5、选择storage path

- 1 当分配好storage server后，客户端将向storage发送写文件请求，storage将会为文件分配一个数据存储目录，支持如下规则（在storage配置文件可以通过配置store\_path\*参数来设置，该参数可以设置多个，通过\*来区别）：
- 2 1. Round robin，多个存储目录间轮询
- 3 2. 剩余存储空间最多的优先

#### 6、生成文件名

- 1 选定存储目录之后，storage会为文件生一个文件名称，由源storage server ip、文件创建时间、文件大小、文件crc32和一个随机数拼接而成，然后将这个二进制串进行base64编码，转换为可打印的字符串。

#### 7、选择两级目录

- 1 当选定存储目录之后，storage会为文件分配一个fileid，每个存储目录下有两级256\*256的子目录，storage会按文件名称进行两次hash（猜测），路由到其中一个子目录，然后将文件以fileid为文件名存储到该子目录下。

#### 8、生成fileid

- 1 当文件存储到某个子目录后，即认为该文件存储成功，接下来会为该文件生成一个文件名，文件名由group、存储目录、两级子目录、文件名、文件后缀名（由客户端指定，主要用于区分文件类型）拼接而成。

## 5.文件同步分析

写文件时，客户端将文件写至group内一个storage server即认为**写文件成功**，storage server写完文件后，会由后台线程将文件同步至同group内其他的storage server。

**同步规则总结如下：**

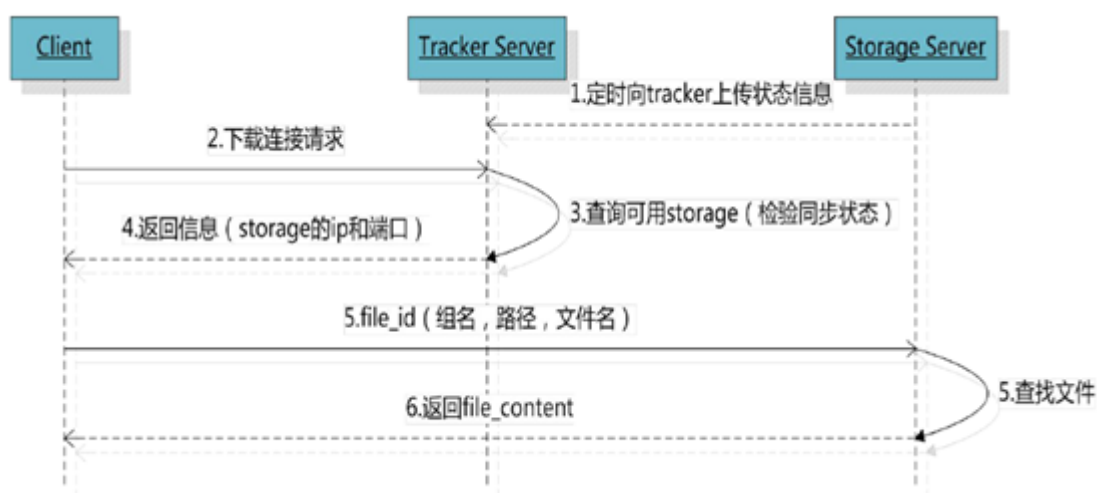
1. 只在本组内的storage server之间进行同步；
2. 源头数据才需要同步，备份数据不需要再次同步，否则就构成环路了；
3. 上述第二条规则有个例外，就是新增加一台storage server时，由已有的一台storage server将已有的所有数据（包括源头数据和备份数据）同步给该新增服务器。

每个storage写文件后，同时会写一份binlog，binlog里不包含文件数据，只包含文件名等元信息，**这份binlog用于后台同步**，storage会记录向group内其他storage同步的进度，以便重启后能接上次的进度继续同步；**进度以时间戳的方式进行记录**，所以最好能保证集群内所有server的时钟保持同步。

storage的同步进度会作为元数据的一部分汇报到tracker上，tracker在选择读storage的时候会以同步进度作为参考。

比如一个group内有A、B、C三个storage server，A向C同步到进度为T1（T1以前写的文件都已经同步到B上了），B向C同步到时间戳为T2（ $T2 > T1$ ），tracker接收到这些同步进度信息时，就会进行整理，将最小的那个做为C的同步时间戳，本例中T1即为C的同步时间戳为T1（即所有T1以前写的文件都已经同步到C上了）；同理，根据上述规则，tracker会为A、B生成一个同步时间戳。

## 6.文件下载流程分析



客户端upload file成功后，会拿到一个storage生成的文件名，接下来客户端根据这个文件名即可访问到该文件。

### 流程说明：

#### 1、tracker server收集storage server的状态信息

- 1 storage server定时向已知的tracker server（可以是多个）发送磁盘剩余空间、文件同步状况、文件上传下载次数等统计信息
- 2 storage server会连接整个集群中所有的Tracker server，向他们报告自己的状态。

#### 2、选择tracker server

1 跟upload file一样，在download file时客户端可以选择任意tracker server。

### 3、选择可用的storage server

- 1 client发送download请求给某个tracker，必须带上文件名信息，tracker从文件名中解析出文件的group、路径信息、文件大小、创建时间、源storage server ip等信息，然后为该请求选择一个storage用来服务读请求。
- 2 由于group内的文件同步是在后台异步进行的，所以有可能出现在读的时候，文件还没有同步到某些storage server上，为了尽量避免访问到这样的storage，tracker按照如下规则选择group内可读的storage：
  - 3 1. 该文件上传到的源头storage - 源头storage只要存活着，肯定包含这个文件，源头的地址被编码在文件名中。
  - 4 2. 文件创建时间戳==storage被同步到的时间戳 且(当前时间-文件创建时间戳) > 文件同步最大时间(如5分钟) - 文件创建后，认为经过最大同步时间后，肯定已经同步到其他storage了。
  - 5 3. 文件创建时间戳 < storage被同步到的时间戳。 - 同步时间戳之前的文件确定已经同步了
  - 6 4. (当前时间-文件创建时间戳) > 同步延迟阈值(如一天)。 - 经过同步延迟阈值时间，认为文件肯定已经同步了。

## 7.新增Storage server分析

组内新增加一台 storage server A 时，由系统自动完成已有数据同步，处理逻辑如下： 1. storage server A连接tracker server，tracker server将storage server A的状态设置为 FDFS\_STORAGE\_STATUS\_INIT。storage server A询问追加同步的源服务器和追加同步截至时间点，如果该组内只有storage server A或该组内已成功上传的文件数为0，则没有数据需要同步，storage server A就可以提供在线服务，此时tracker将其状态设置为 FDFS\_STORAGE\_STATUS\_ONLINE，否则tracker server将其状态设置为FDFS\_STORAGE\_STATUS\_WAIT\_SYNC，进入第二步的处理；

2. 假设tracker server分配向storage server A同步已有数据的源storage server为B。同组的storage server和tracker server通讯得知新增了storage server A，将启动同步线程，并向tracker server询问向storage server A追加同步的源服务器和截至时间点。storage server B将把截至时间点之前的所有数据同步给storage server A；而其余的storage server从截至时间点之后进行正常同步，只把源头数据同步给storage server A。到了截至时间点之后，storage server B对storage server A的同步将由追加同步切换为正常同步，只同步源头数据； 3. storage server B向storage server A同步完所有数据，暂时没有数据要同步时，storage server B请求tracker server将storage server A的状态设置为FDFS\_STORAGE\_STATUS\_ONLINE； 4 当storage server A向tracker server发起heart beat时，tracker server将其状态更改为FDFS\_STORAGE\_STATUS\_ACTIVE。

## 四、FastDFS安装

### 1.下载

- 下载libfastcommon包
  - github源码

```
1 | https://github.com/happyfish100/libfastcommon/releases
```

- linux wget下载

```
1 | wget https://github.com/happyfish100/libfastcommon/archive/v1.0.39.tar.gz
```

- 下载fastdfs源码包
  - github源码

```
1 | https://github.com/happyfish100/fastdfs/releases
```

- linux wget下载

```
1 | wget https://github.com/happyfish100/fastdfs/archive/v5.11.tar.gz
```

## 2.需求

- Tracker Server : 192.168.10.135
- Storage Server : 192.168.10.136、192.168.10.137

## 3.安装

### tracker和storage安装

tracker server和storage server都有一些相同的安装操作，如下：

- 安装gcc环境

```
1 | yum install -y gcc-c++
```

- 安装libevent，FastDFS依赖libevent库（暂不安装）

```
1 | yum install -y libevent
```

- 安装libfastcommon，libfastcommon是FastDFS官方提供的包，包含了FastDFS运行所需要的一些基础库。

```
1 | wget https://github.com/happyfish100/libfastcommon/archive/v1.0.39.tar.gz
2 | tar -zxvf v1.0.39.tar.gz
3 | cd libfastcommon-1.0.39
4 | ./make.sh && ./make.sh install
```

- 拷贝libfastcommon.so文件至/usr/lib目录（新版本不需要此步）

```
1 | cp /usr/lib64/libfastcommon.so /usr/lib/
```

**注：**

libfastcommon安装好后会自动将库文件拷贝至/usr/lib64下，由于FastDFS程序引用usr/lib目录，所以需要将/usr/lib64下的库文件拷贝至/usr/lib下。

- 下载安装FastDFS，进入FastDFS目录，编译安装

```
1 | wget https://github.com/happyfish100/fastdfs/archive/v5.11.tar.gz
2 | tar -zxvf v5.11.tar.gz
3 | cd cd fastdfs-5.11
4 | ./make.sh && ./make.sh install
```

- 拷贝/root/fastdfs-5.11/conf目录下的文件到/etc/fdfs目录下

```
1 | cp /kbb/soft/fastdfs-5.11/conf/* /etc/fdfs
```

## tracker server配置

**注意：**base\_path目录要存在

- 修改/etc/fdfs/tracker.conf

```
1 | vim /etc/fdfs/tracker.conf
```

修改内容如下：

```
1 | base_path=/kbb/server/fastdfs/tracker
```

- 创建tracker服务器上面的目录

```
1 | mkdir /kbb/server/fastdfs/tracker -p
```

## storage server配置

**注意：**base\_path和store\_path0目录要存在

- 修改/etc/fdfs/storage.conf

```
1 | vim /etc/fdfs/storage.conf
```

修改内容如下：



```
1 #指定storage的组名
2 group_name=group1
3 base_path=/kbb/server/fastdfs/storage
4 store_path0=/kbb/server/fastdfs/storage
5 #如果有多个挂载磁盘则定义多个store_path, 如下
6 #store_path1=.....
7 #store_path2=.....
8 #配置tracker服务器IP和端口
9 tracker_Server=111.231.106.221:22122
10 #如果有多个则配置多个tracker
11 #tracker_Server=192.168.101.4:22122
```

- 创建storage服务器上面的目录

```
1 | mkdir /kbb/server/fastdfs/storage -p
```

## 4.启动

### 4.1 Tracker启动命令

```
1 | /usr/bin/fdfs_trackerd /etc/fdfs/tracker.conf
```

### 4.2 Storage启动命令

```
1 | /usr/bin/fdfs_storaged /etc/fdfs/storage.conf
```

### 4.3 Tracker开机自启动

```
1 | vim /etc/rc.d/rc.local
```

将运行命令行添加进文件：

```
1 | /usr/bin/fdfs_trackerd /etc/fdfs/tracker.conf
```

### 4.4 Storage开机自启动

```
1 | vim /etc/rc.d/rc.local
```

将运行命令行添加进文件：

```
1 | /usr/bin/fdfs_storaged /etc/fdfs/storage.conf
```

## 5.上传图片测试

FastDFS安装成功后可通过【fdfs\_test】命令测试上传、下载等操作。

- 修改client.conf

```
1 | vim /etc/fdfs/client.conf
```

修改内容如下：

```
1 | base_path=/kkb/server/fastdfs/client
2 | tracker_server=111.231.106.221:22122
```

- 创建client的数据目录

```
1 | mkdir -p /kkb/server/fastdfs/client
```

- 使用fdfs\_test命令将/home下的tomcat.png上传到FastDFS中

```
1 | /usr/bin/fdfs_test /etc/fdfs/client.conf upload /etc/fdfs/anti-steal.jpg
```

说明：

[http://192.168.10.135/group1/M00/00/00/wKh1BVVY2M-AM\\_9DAAAT7-0xdqM485\\_big.png](http://192.168.10.135/group1/M00/00/00/wKh1BVVY2M-AM_9DAAAT7-0xdqM485_big.png)就是文件的访问路径。

对应storage服务器上的磁盘路径：

/home/fastdfs/fdfs\_storage/data/00/00/wkh1BVVY2M-AM\_9DAAAT7-0xdqM485\_big.png 文件。

## 6.tracker.conf

### 6.1 基本配置

```
1 | disable
2 | #func：配置是否生效
3 | #value:true、false
4 | disable=false
5 | bind_addr
6 | #func：绑定IP
7 | #value：IP地址
8 | bind_addr=192.168.6.102
9 | port
10 | #func：服务端口
11 | #value：端口整数
12 | port=22122
13 | connect_timeout
14 | #func：连接超时
15 | #value：秒单位正整数
16 | connect_timeout=30
```

```
17 network_timeout
18 #func : 网络超时
19 #valu : 秒单位正整数值
20 network_timeout=60
21 base_path
22 #func : Tracker数据/日志目录地址
23 #valu : 路径
24 base_path=/home/michael/fdfs/base4tracker
25 max_connections
26 #func : 最大连接数
27 #valu : 正整数值
28 max_connections=256
29 work_threads
30 #func : 线程数, 通常设置CPU数
31 #valu : 正整数值
32 work_threads=4
33 store_lookup
34 #func : 上传文件的选组方式。
35 #valu : 0、1或2。
36 # 0 : 表示轮询
37 # 1 : 表示指定组
38 # 2 : 表示存储负载均衡 ( 选择剩余空间最大的组 )
39 store_lookup=2
40 store_group
41 #func : 指定上传的组, 如果在应用层指定了具体的组, 那么这个参数将不会起效。另外如果store_lookup如果是0或2, 则此参数无效。
42 #valu : group1等
43 store_group=group1
44 store_server
45 #func : 上传服务器的选择方式。(一个文件被上传后, 这个storage Server就相当于这个文件的storage Server源, 会对同组的storage Server推送这个文件达到同步效果)
46 #valu : 0、1或2
47 # 0 : 轮询方式 ( 默认 )
48 # 1 : 根据ip 地址进行排序选择第一个服务器 ( IP地址最小者 )
49 # 2 : 根据优先级进行排序 ( 上传优先级由storage Server来设置, 参数名为upload_priority ), 优先级值越小优先级越高。
50 store_server=0
51 store_path
52 #func : 上传路径的选择方式。storage Server可以有多个存放文件的base path ( 可以理解为多个磁盘 )。
53 #valu :
54 # 0 : 轮流方式, 多个目录依次存放文件
55 # 2 : 存储负载均衡。选择剩余空间最大的目录存放文件 ( 注意 : 剩余磁盘空间是动态的, 因此存储到的目录或磁盘可能也是变化的 )
56 store_path=0
57 download_server
58 #func : 下载服务器的选择方式。
59 #valu :
60 # 0 : 轮询 ( 默认 )
61 # 1 : IP最小者
62 # 2 : 优先级排序 ( 值最小的, 优先级最高。 )
63 download_server=0
64 reserved_storage_space
```

```

65 #func：保留空间值。如果某个组中的某个服务器的剩余自由空间小于设定值，则文件不会被上传到这个组。
66 #valu：
67 # G or g for gigabyte
68 # M or m for megabyte
69 # K or k for kilobyte
70 reserved_storage_space=1GB
71 log_level
72 #func：日志级别
73 #valu：
74 # emerg for emergency
75 # alert
76 # crit for critical
77 # error
78 # warn for warning
79 # notice
80 # info for information
81 # debug for debugging
82 log_level=info
83 run_by_group / run_by_user
84 #func：指定运行该程序的用户组
85 #valu：用户组名或空
86 run_by_group=
87
88 #func：
89 #valu：
90 run_by_user=
91 allow_hosts
92 #func：可以连接到tracker Server的ip范围。可设定多个值。
93 #valu
94 allow_hosts=
95 check_active_interval
96 #func：检测 storage Server 存活的时间间隔，单位为秒。
97 #      storage Server定期向tracker Server 发心跳，
98 #      如果tracker Server在一个check_active_interval内还没有收到storage server的一次心跳，
99 #      那边将认为该storage Server已经下线。所以本参数值必须大于storage Server配置的心跳时间间隔。
100 #      通常配置为storage Server心跳时间间隔的2倍或3倍。
101 check_active_interval=120
102 thread_stack_size
103 #func：设定线程栈的大小。 线程栈越大，一个线程占用的系统资源就越多。
104 #      如果要启动更多的线程（V1.x对应的参数为max_connections，V2.0为work_threads），可以适当降低本参数值。
105 #valu：如64KB，默认值为64，tracker Server线程栈不应小于64KB
106 thread_stack_size=64KB
107 storage_ip_changed_auto_adjust
108 #func：这个参数控制当storage Server IP地址改变时，集群是否自动调整。注：只有在storage Server进程重启时才完成自动调整。
109 #valu：true或false
110 storage_ip_changed_auto_adjust=true

```

## 6.2 同步

```

1 storage_sync_file_max_delay
2 #func : 同组storage服务器之间同步的最大延迟时间。存储服务器之间同步文件的最大延迟时间，根据实际情况
   进行调整
3 #valu : 秒为单位，默认值为1天 ( 24*3600 )
4 #sinc : v2.0
5 storage_sync_file_max_delay=86400
6 storage_sync_file_max_time
7 #func : 存储服务器同步一个文件需要消耗的最大时间，缺省为300s，即5分钟。
8 #sinc : v2.0
9 storage_sync_file_max_time=300
10 sync_log_buff_interval
11 #func : 同步或刷新日志信息到硬盘的时间间隔。注意：tracker Server 的日志不是时时写硬盘的，而是先写
   内存。
12 #valu : 以秒为单位
13 sync_log_buff_interval=10

```

## 6.3 trunk 和 slot

```

1 #func : 是否使用trunk文件来存储几个小文件
2 #valu : true或false
3 #sinc : v3.0
4 use_trunk_file=false
5
6 #func : 最小slot大小
7 #valu : <= 4KB，默认为256字节
8 #sinc : v3.0
9 slot_min_size=256
10
11 #func : 最大slot大小
12 #valu : >= slot_min_size，当小于这个值的时候就存储到trunk file中。默认为16MB。
13 #sinc : v3.0
14 slot_max_size=16MB
15
16 #func : trunk file的size
17 #valu : >= 4MB，默认为64MB
18 #sinc : v3.0
19 trunk_file_size=64MB

```

## 6.4 HTTP相关

```

1 是否启用 HTTP
2 #func : HTTP是否生效
3 #valu : true或false
4 http.disabled=false
5 HTTP 服务器端口号
6 #func : tracker Server上的http port
7 #valu :
8 #note : 只有http.disabled=false时才生效
9 http.Server_port=7271
10 检查Storage存活状态的间隔时间 ( 心跳检测 )
11 #func : 检查storage http Server存活的间隔时间

```

```
12 #valu : 单位为秒
13 #note : 只有http.disabled=false时才生效
14 http.check_alive_interval=30
15 心跳检测使用的协议方式
16 #func : 检查storage http Server存活的方式
17 #valu :
18 # tcp : 连接到storage Server的http端口 , 不进行request和response。
19 # http : storage check alive url must return http status 200.
20 #note : 只有http.disabled=false时才生效
21 http.check_alive_type=tcp
22 检查 Storage 状态的 URI
23 #func : 检查storage http Server是否alive的uri/url
24 #note : 只有http.disabled=false时才生效
25 http.check_alive_uri=/status.html
```

## 五、FastDFS-Nginx扩展模块源码分析

### 1. 背景

在大多数业务场景中，往往需要为FastDFS存储的文件提供http下载服务，而尽管FastDFS在其storage及tracker都内置了http服务，但性能表现却不尽如人意；(余老师在 V4.05 以后的版本就把内置 HTTP服务去掉了) 作者余庆在后来的版本中增加了基于当前主流web服务器的扩展模块(包括nginx/apache)，其用意在于利用web服务器直接对本机storage数据文件提供http服务，以提高文件下载的性能。

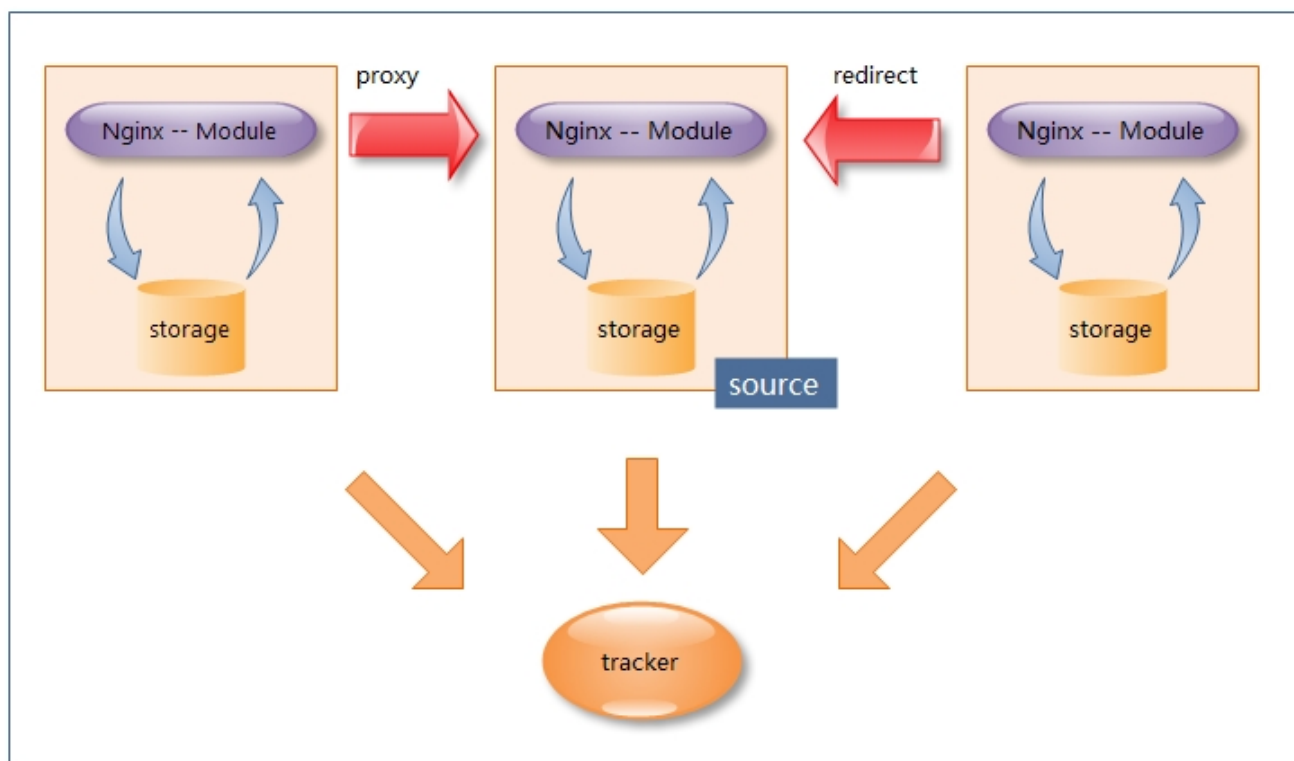
不使用Nginx的扩展模块，只安装web服务器（Nginx和Apache），也可以对存储的文件进行访问。

为什么要使用Nginx的扩展模块来访问存储的文件，原因有两个：

- 如果进行文件合并，那么不使用FastDFS的nginx扩展模块，是无法访问到具体的文件的，因为文件合并之后，多个小文件都是存储在一个trunk文件中的，在存储目录下，是看不到具体的小文件的。
- 如果文件未同步成功，那么不使用FastDFS的nginx扩展模块，是无法正常访问到指定的文件的，而使用了FastDFS的nginx扩展模块之后，如果要访问的文件未同步成功，那么会解析出来该文件的源存储服务器ip，然后将该访问请求重定向或者代理到源存储服务器中进行访问。

### 2. 概要介绍

#### 2.1 FastDFS整合Nginx的参考架构



#### 说明：

在每一台storage服务器主机上部署Nginx及FastDFS扩展模块，由Nginx模块对storage存储的文件提供http下载服务，仅当当前storage节点找不到文件时会向源storage主机发起 `redirect` 或 `proxy` 动作。

#### 注：

图中的tracker可能为多个tracker组成的集群；

且当前FastDFS的Nginx扩展模块支持单机多个group的情况

## 2.2 几个概念

- **storage\_id**：指storage server的id，从FastDFS4.x版本开始，tracker可以对storage定义一组ip到id的映射，以id的形式对storage进行管理。而文件名写入的不再是storage的ip而是id，这样的方式对于数据迁移十分有利。
- **storage\_sync\_file\_max\_delay**：指storage节点同步一个文件最大的时间延迟，是一个阈值；如果当前时间与文件创建时间的差距超过该值则认为同步已经完成。
- **anti\_steal\_token**：指文件ID防盗链的方式，FastDFS采用 `token认证` 的方式进行文件防盗链检查。

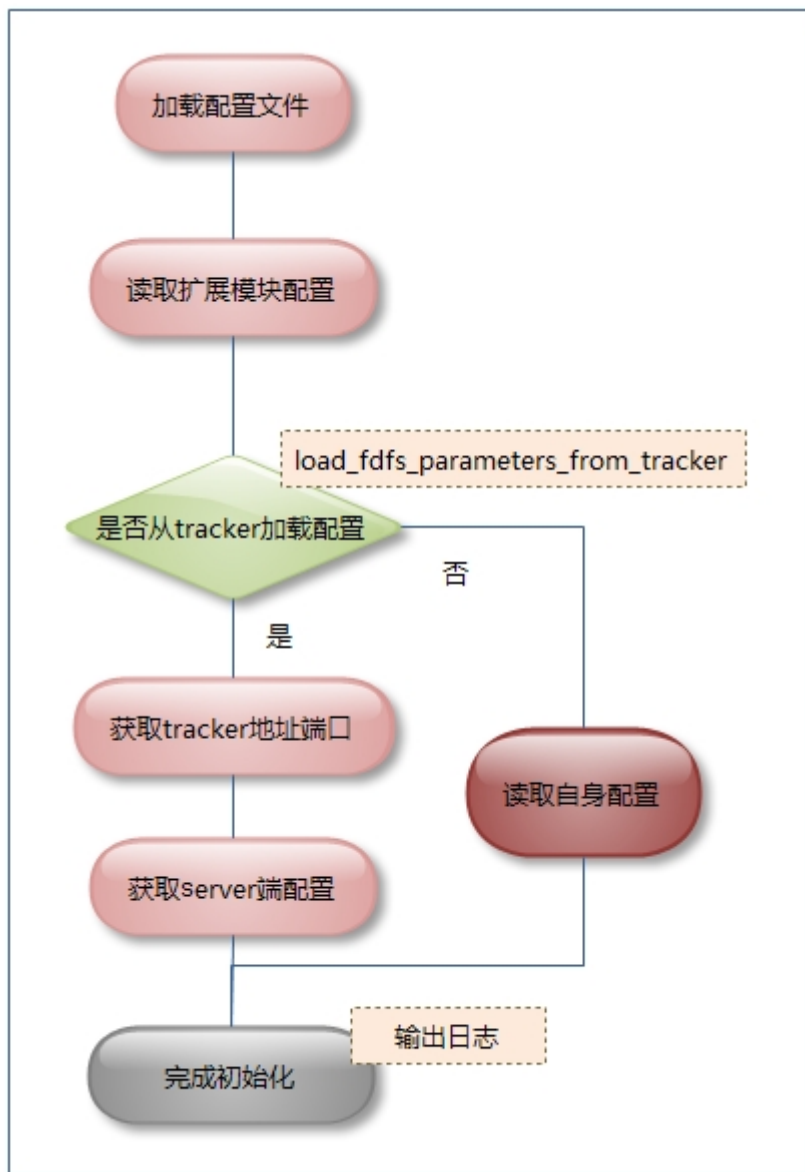
## 3. 实现原理

### 3.1 源码包说明

下载后的源码包很小，仅包括以下文件：

|   |                           |  |
|---|---------------------------|--|
| 1 | ngx_http_fastdfs_module.c | //nginx-module接口实现文件，用于接入fastdfs-module核心模块逻辑  |
| 2 | common.c                  | //fastdfs-module核心模块，实现了初始化、文件下载的主要逻辑          |
| 3 | common.h                  | //对应于common.c的头文件                              |
| 4 | config                    | //编译模块所用的配置，里面定义了一些重要的常量，如扩展配置文件路径、文件下载chunk大小 |
| 5 | mod_fastdfs.conf          | //扩展配置文件的demo                                  |

## 3.2 初始化



### 3.2.1 加载配置文件

目标文件：/etc/fdfs/mod\_fastdfs.conf

### 3.2.2 读取扩展模块配置

一些重要参数包括：



```
1 group_count           //group个数
2 url_have_group_name   //url中是否包含group
3 group.store_path      //group对应的存储路径
4 connect_timeout       //连接超时
5 network_timeout       //接收或发送超时
6 storage_server_port   //storage_server端口，用于在找不到文件情况下连接源storage下载文件(该做法已过时)
7 response_mode         //响应模式，proxy或redirect
8 load_fdfs_parameters_from_tracker //是否从tracker下载服务端配置
```

### 3.2.3 加载服务端配置

根据 `load_fdfs_parameters_from_tracker` 参数确定是否从tracker获取server端的配置信息

- **load\_fdfs\_parameters\_from\_tracker=true:**

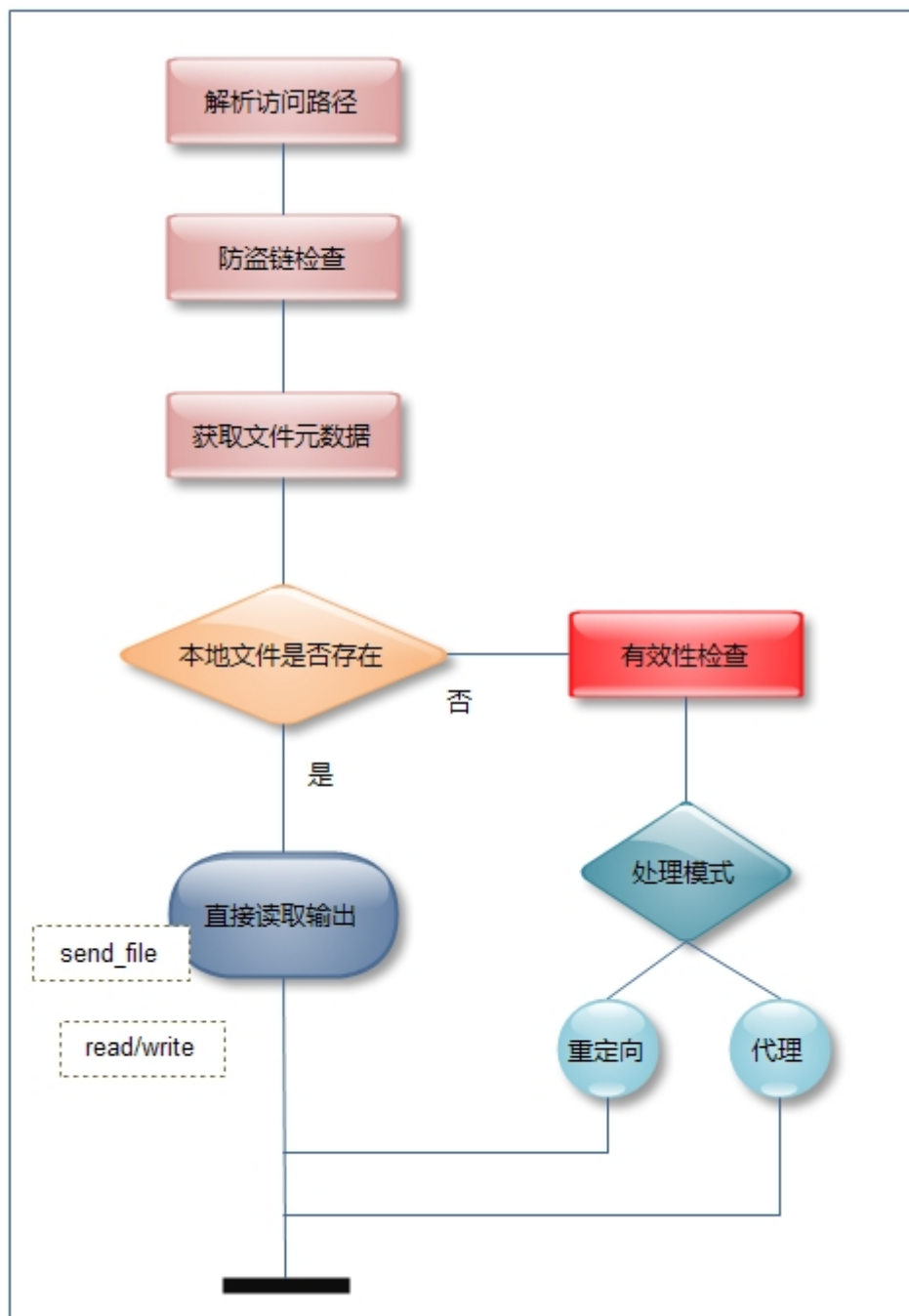
1. 调用 `fdfs_load_tracker_group_ex` 解析tracker连接配置；
2. 调用 `fdfs_get_ini_context_from_tracker` 连接tracker获取配置信息；
3. 获取 `storage_sync_file_max_delay` 阈值
4. 获取 `use_storage_id`
5. 如果 `use_storage_id` 为 true，则连接tracker获取 `storage_ids` 映射表(调用方法：  
**`fdfs_get_storage_ids_from_tracker_group`**)

- **load\_fdfs\_parameters\_from\_tracker=false:**

1. 从 `mod_fastdfs.conf` 加载所需配置：`storage_sync_file_max_delay`、`use_storage_id`;
2. 如果 `use_storage_id` 为 true，则根据 `storage_ids_filename` 获取 `storage_ids` 映射表(调用方法：  
**`fdfs_load_storage_ids_from_file`**)

---

## 3.3 下载过程(重点)



### 3.3.1 解析访问路径

得到 `group` 和 `file_id_without_group` 两个参数;

### 3.3.2 防盗链检查

- 根据 `g_http_params.anti_steal_token` 配置(见 `http.conf` 文件), 判断是否进行防盗链检查;
- 采用 token 的方式实现防盗链, 该方式要求下载地址带上 token, 且 token 具有时效性(由 `ts` 参数指明);

检查方式:

1 | `md5(fileid_without_group + privKey + ts) = token`; 同时 `ts` 没有超过 `ttd` 范围 (可参考 `JavaClient CommonProtocol`)

调用方法：**fdfs\_http\_check\_token** 关于FastDFS的防盗链可参考：<http://bbs.chinaunix.net/thread-1916999-1-1.html>

### 3.3.3 获取文件元数据

根据文件ID 获取元数据信息, 包括：**源storage ip,文件路径、名称，大小** 代码：

```
1 | if ((result=fdfs_get_file_info_ex1(file_id, false, &file_info)) != 0)...
```

在**fdfs\_get\_file\_info\_ex1**的实现中，存在一个取巧的逻辑：当获得文件的ip段之后，仍然需要确定该段落是storage的id还是ip。 代码：

```
1 | fdfs_shared.func.c
2 | -> fdfs_get_server_id_type(ip_addr.s_addr) == FDFS_ID_TYPE_SERVER_ID
3 | ...
4 |     if (id > 0 && id <= FDFS_MAX_SERVER_ID) {
5 |         return FDFS_ID_TYPE_SERVER_ID;
6 |     } else {
7 |         return FDFS_ID_TYPE_IP_ADDRESS;
8 |     }
```

判断标准为ip段的整数值是否在 0 到 -> **FDFS\_MAX\_SERVER\_ID**(见tracker\_types.h)之间； 其中  $FDFS\_MAX\_SERVER\_ID = (1 \ll 24) - 1$ ，该做法利用了ipv4地址的特点(由4\*8个二进制位组成)，即ipv4地址数值务必大于该阈值

### 3.3.4 检查本地文件是否存在

调用**trunk\_file\_stat\_ex1**获取本地文件信息，该方法将实现：

1. 辨别当前文件是trunkfile还是singlefile
2. 获得文件句柄fd
3. 如果文件是trunk形式则同时也将相关信息(偏移量/长度)一并获得

代码：

```
1 |     if (bSameGroup)
2 |     {
3 |         FDFS_TrunkHeader trunkHeader;
4 |         if ((result=trunk_file_stat_ex1(pStorePaths, store_path_index, \
5 |             true_filename, filename_len, &file_stat, \
6 |             &trunkInfo, &trunkHeader, &fd)) != 0)
7 |         {
8 |             bFileExists = false;
9 |         }
10 |        else
11 |        {
12 |            bFileExists = true;
13 |        }
14 |    }
15 |    else
```

```

16     {
17         bFileExists = false;
18         memset(&trunkInfo, 0, sizeof(trunkInfo));
19     }

```

### 3.3.5 文件不存在的处理

- 进行有效性检查

检查项有二：

**\*代码\*：**

```

1     if (is_local_host_ip(file_info.source_ip_addr) || \
2         (file_info.create_timestamp > 0 && (time(NULL) - \
3             file_info.create_timestamp > "'storage_sync_file_max_delay'")))

```

在通过有效性检查之后将进行代理或重定向处理

- 重定向模式

配置项 `response_mode = redirect`，此时服务端返回返回302响应码，url如下：

```

1 http:// {源storage地址} : {当前port} {当前url} {参数"redirect=1"}(标记已重定向过)

```

**代码：**

```

1     response.redirect_url_len = snprintf( \
2         response.redirect_url, \
3         sizeof(response.redirect_url), \
4         "http://%s%s%s%s%c%s", \
5         file_info.source_ip_addr, port_part, \
6         path_split_str, url, \
7         param_split_char, "redirect=1");

```

注：该模式下要求源storage配备公开访问的webserver、同样的端口(一般是80)、同样的path配置。

- 代理模式

配置项 `response_mode = proxy`，该模式的工作原理如同反向代理的做法，而仅仅使用源storage地址作为代理proxy的host，其余部分保持不变。 **代码：**

```

1     if (pContext->proxy_handler != NULL)
2     {
3         return pContext->proxy_handler(pContext->arg, \
4             file_info.source_ip_addr);
5     }
6     //其中proxy_handler方法来自ngx_http_fastdfs_module.c文件的
    ngx_http_fastdfs_proxy_handler方法
7     //其实现中设置了大量回调、变量，并最终调用代理请求方法，返回结果：
8     rc = ngx_http_read_client_request_body(r, ngx_http_upstream_init); //执行代理
    请求，并返回结果

```

### 3.3.6 输出本地文件

- 根据是否**trunkfile**获取文件名，文件名长度、文件offset;

代码：

```

1     bTrunkFile = IS_TRUNK_FILE_BY_ID(trunkInfo);
2     if (bTrunkFile)
3     {
4         trunk_get_full_filename_ex(pStorePaths, &trunkInfo, \
5             full_filename, sizeof(full_filename));
6         full_filename_len = strlen(full_filename);
7         file_offset = TRUNK_FILE_START_OFFSET(trunkInfo) + \
8             pContext->range.start;
9     }
10    else
11    {
12        full_filename_len = snprintf(full_filename, \
13            sizeof(full_filename), "%s/data/%s", \
14            pStorePaths->paths[store_path_index], \
15            true_filename);
16        file_offset = pContext->range.start;
17    }

```

- 若nginx开启了send\_file开关而且当前为非chunkFile的情况下尝试**使用sendfile方法以优化性能**;

代码：

```

1     if (pContext->send_file != NULL && !bTrunkFile)
2     {
3         http_status = pContext->if_range ? \
4             HTTP_PARTIAL_CONTENT : HTTP_OK;
5         OUTPUT_HEADERS(pContext, (&response), http_status)
6         .....
7         return pContext->send_file(pContext->arg, full_filename, \
8             full_filename_len, file_offset, download_bytes);
9     }

```

- 否则使用lseek 方式随机访问文件，并输出相应的段;

做法：使用chunk方式循环读，输出... 代码：

```
1  while (remain_bytes > 0)
2  {
3      read_bytes = remain_bytes <= FDFS_OUTPUT_CHUNK_SIZE ? \
4          remain_bytes : FDFS_OUTPUT_CHUNK_SIZE;
5      if (read(fd, file_trunk_buff, read_bytes) != read_bytes)
6      {
7          close(fd);
8          .....
9          return HTTP_INTERNAL_SERVER_ERROR;
10     }
11
12     remain_bytes -= read_bytes;
13     if (pContext->send_reply_chunk(pContext->arg, \
14         (remain_bytes == 0) ? 1: 0, file_trunk_buff, \
15         read_bytes) != 0)
16     {
17         close(fd);
18         return HTTP_INTERNAL_SERVER_ERROR;
19     }
20 }
```

其中chunk大小见config文件配置：`-DFDFS_OUTPUT_CHUNK_SIZE='256*1024'`

## 六、配置FastDFS的Nginx模块

FastDFS的Nginx模块需要安装到每一台storage server中

### 1. 下载文件

```
1 | wget https://github.com/happyfish100/fastdfs-nginx-module/archive/v1.20.tar.gz
```

### 2. 解压缩

```
1 | tar -zxvf v1.20.tar.gz
```

### 3. 修改config文件（特别关键的一步）

修改 `fastdfs-nginx-module/src/config` 文件

```
1 | cd fastdfs-nginx-module-1.20/src
2 | vim config
```

- 修改前的内容如下：

```

1 ngx_addon_name=ngx_http_fastdfs_module
2
3 if test -n "${ngx_module_link}"; then
4     ngx_module_type=HTTP
5     ngx_module_name=$ngx_addon_name
6     ngx_module_incs="/usr/local/include"
7     ngx_module_libs="-lfastcommon -lfdescclient"
8     ngx_module_srcs="$ngx_addon_dir/ngx_http_fastdfs_module.c"
9     ngx_module_deps=
10    CFLAGS="$CFLAGS -D_FILE_OFFSET_BITS=64 -DFDFS_OUTPUT_CHUNK_SIZE='256*1024' -
DFDFS_MOD_CONF_FILENAME='\"/etc/fdfs/mod_fastdfs.conf\"'"
11    . auto/module
12 else
13    HTTP_MODULES="$HTTP_MODULES ngx_http_fastdfs_module"
14    NGX_ADDON_SRCS="$NGX_ADDON_SRCS $ngx_addon_dir/ngx_http_fastdfs_module.c"
15    CORE_INCS="$CORE_INCS /usr/local/include"
16    CORE_LIBS="$CORE_LIBS -lfastcommon -lfdescclient"
17    CFLAGS="$CFLAGS -D_FILE_OFFSET_BITS=64 -DFDFS_OUTPUT_CHUNK_SIZE='256*1024' -
DFDFS_MOD_CONF_FILENAME='\"/etc/fdfs/mod_fastdfs.conf\"'"
18 fi

```

- 其中第6行和第15行要进行修改，修改后的内容如下：

```

1 ngx_addon_name=ngx_http_fastdfs_module
2
3 if test -n "${ngx_module_link}"; then
4     ngx_module_type=HTTP
5     ngx_module_name=$ngx_addon_name
6     ngx_module_incs="/usr/include/fastdfs /usr/include/fastcommon/"
7     ngx_module_libs="-lfastcommon -lfdescclient"
8     ngx_module_srcs="$ngx_addon_dir/ngx_http_fastdfs_module.c"
9     ngx_module_deps=
10    CFLAGS="$CFLAGS -D_FILE_OFFSET_BITS=64 -DFDFS_OUTPUT_CHUNK_SIZE='256*1024' -
DFDFS_MOD_CONF_FILENAME='\"/etc/fdfs/mod_fastdfs.conf\"'"
11    . auto/module
12 else
13    HTTP_MODULES="$HTTP_MODULES ngx_http_fastdfs_module"
14    NGX_ADDON_SRCS="$NGX_ADDON_SRCS $ngx_addon_dir/ngx_http_fastdfs_module.c"
15    CORE_INCS="$CORE_INCS /usr/include/fastdfs /usr/include/fastcommon/"
16    CORE_LIBS="$CORE_LIBS -lfastcommon -lfdescclient"
17    CFLAGS="$CFLAGS -D_FILE_OFFSET_BITS=64 -DFDFS_OUTPUT_CHUNK_SIZE='256*1024' -
DFDFS_MOD_CONF_FILENAME='\"/etc/fdfs/mod_fastdfs.conf\"'"
18 fi

```

## 4.拷贝mod\_fastdfs.conf

将 fastdfs-nginx-module-1.20/src/mod\_fastdfs.conf 拷贝至 /etc/fdfs/ 下

```
1 | cp mod_fastdfs.conf /etc/fdfs/
```

## 5.修改mod\_fastdfs.conf

```
1 | vim /etc/fdfs/mod_fastdfs.conf
```

文件内容如下：

```
1 | base_path=/kkb/server/fastdfs/storage
2 | tracker_server=192.168.10.135:22122 #url中是否包含group名称
3 | url_have_group_name=true #指定文件存储路径，访问时使用该路径
4 | store_path0=/kkb/server/fastdfs/storage
```

## 6.拷贝libfdfsclient.so(新版不需要)

将 libfdfsclient.so 拷贝至/usr/lib下

```
1 | cp /usr/lib64/libfdfsclient.so /usr/lib/
```

## 7.创建nginx/client目录

```
1 | mkdir -p /var/temp/nginx/client
```

# 七、安装Nginx ( Apache )

Nginx需要每一台storage server都需要安装。

## 1.下载文件

- 查看GitHub上面最新的nginx release 版本

<https://github.com/nginx/nginx/releases>

- 下载各个版本的nginx的地址：

<http://nginx.org/download/>

- 上传nginx压缩包

```
1 | yum install -y lrzsz
2 | rz
```

通过rz命令上传文件：nginx-1.15.6.tar.gz



## 2.安装第三方软件

### 2.1 安装PCRE

PCRE(Perl Compatible Regular Expressions)是一个Perl库，包括 perl 兼容的正则表达式库。Nginx的http模块使用pcre来解析正则表达式，所以需要在linux上安装pcre库。

```
1 | yum install -y pcre-devel
```

注：pcre-devel是使用pcre开发的一个二次开发库。Nginx也需要此库。

### 2.2 安装ZLIB

zlib库提供了很多种压缩和解压缩的方式，Nginx使用zlib对http包的内容进行gzip，所以需要在linux上安装zlib库。

```
1 | yum install -y zlib-devel
```

### 2.3 安装OPENSSL

OpenSSL 是一个强大的安全套接字层密码库，囊括主要的密码算法、常用的密钥和证书封装管理功能及SSL协议，并提供丰富的应用程序供测试或其它目的使用。

Nginx不仅支持http协议，还支持https（即在ssl协议上传输http），所以需要在linux安装openssl库。

```
1 | yum install -y openssl-devel
```

## 3.解压缩

```
1 | tar -xf nginx-1.15.6.tar.gz
```

## 4.执行configure配置

```
1 | cd nginx-1.15.6/
2
3 | ./configure \
4 | --prefix=/kbb/server/nginx \
5 | --pid-path=/var/run/nginx/nginx.pid \
6 | --lock-path=/var/lock/nginx.lock \
7 | --error-log-path=/var/log/nginx/error.log \
```

```
8 --http-log-path=/var/log/nginx/access.log \  
9 --http-client-body-temp-path=/var/temp/nginx/client \  
10 --http-proxy-temp-path=/var/temp/nginx/proxy \  
11 --http-fastcgi-temp-path=/var/temp/nginx/fastcgi \  
12 --http-uwsgi-temp-path=/var/temp/nginx/uwsgi \  
13 --http-scgi-temp-path=/var/temp/nginx/scgi \  
14 --with-http_gzip_static_module \  
15 --add-module=/kbb/soft/fastdfs-nginx-module-1.20/src
```

#### 注意：

prefix= /kbb/server/nginx 中的 /kbb/server/nginx 指的是要安装的nginx的路径

add-module= /opt/fastdfs-nginx-module/src 中的路径指的是 fastdfs-nginx-module模块 的解压缩路径

## 5.创建临时目录

上面执行的configure命令，设置了一些配置参数，其中的一些参数指定的目录一定要存在。

```
1 | mkdir /var/temp/nginx -p
```

## 6.编译安装

```
1 | make && make install
```

如果出现以下错误，请修改fastdfs-nginx-module-1.20/src/config配置文件

```
/kbb/server/fastdfs-nginx-module-1.20/src/nginx_http_fastdfs_module.c  
In file included from /kbb/server/fastdfs-nginx-module-1.20/src/common.c:26:0,  
                 from /kbb/server/fastdfs-nginx-module-1.20/src/nginx_http_fastdfs_module.c:6:  
/usr/include/fastdfs/fdfs_define.h:15:27: fatal error: common_define.h: No such file or directory  
#include "common_define.h"  
^  
compilation terminated.  
make[1]: *** [objs/addon/src/nginx_http_fastdfs_module.o] Error 1  
make[1]: Leaving directory `/kbb/soft/nginx-1.15.6'  
make: *** [build] Error 2
```

## 7.修改nginx.conf

```
1 | vim /kbb/server/nginx/conf/nginx.conf
```

修改内容如下：

```
1 server {
2     listen      80;
3     server_name localhost;
4     location /group1/M00/{
5         ngx_fastdfs_module;
6     }
7 }
```

说明：

`location /group1/M00/`：以`/group1/M00/`开头的请求，才会正常使用Nginx模块`ngx_fastdfs_module`下载访问图片。

## 8.启动Nginx

切换到nginx/bin目录

```
1 | ./nginx
```

# 八、合并存储

## 1.合并存储简介

在处理海量小文件问题上，文件系统处理性能会受到显著的影响，在 `读写次数 ( IOPS )` 与 `吞吐量 ( Throughput )` 这两个指标上会有不少的下降。主要需要面对如下几个问题：

- **元数据管理低效**，磁盘文件系统中，目录项(dentry)、索引节点(inode)和数据(data)保存在存储介质的不同位置上。因此，**访问一个文件需要经历至少3次独立的访问**。这样，并发的小文件访问就转变成了大量的随机访问，而这种访问对于广泛使用的磁盘来说是非常低效的；
- **数据布局低效**；
- **IO访问流程复杂**；因此一种解决途径就是将小文件合并存储成大文件，使用seek来定位到大文件的指定位置来访问该小文件。

海量文件描述参考自：<http://blog.csdn.net/liuaigui/article/details/9981135>

注：

FastDFS提供的合并存储功能，默认创建的大文件为64MB，然后在该大文件中存储很多小文件。大文件中容纳一个小文件的空间称为一个Slot，规定Slot最小值为256字节，最大为16MB，也就是小于256字节的文件也需要占用256字节，超过16MB的文件不会合并存储而是创建独立的文件。

## 2.合并存储配置

FastDFS提供了合并存储功能的实现，所有的配置都在 `tracker.conf` 文件之中，具体摘录如下：

注意：开启合并存储只需要设置 `use_trunk_file = true` 和 `store_server=1`

```
1 # which storage server to upload file
2 # 0: round robin (default)
3 # 1: the first server order by ip address
4 # 2: the first server order by priority (the minimal)
5 # Note: if use_trunk_file set to true, must set store_server to 1 or 2
6 store_server = 0
7 #是否启用trunk存储
8 use_trunk_file = false
9 #trunk文件最小分配单元
10 slot_min_size = 256
11 #trunk内部存储的最大文件，超过该值会被独立存储
12 slot_max_size = 16MB
13 #trunk文件大小
14 trunk_file_size = 64MB
15 #是否预先创建trunk文件
16 trunk_create_file_advance = false
17 #预先创建trunk文件的基准时间
18 trunk_create_file_time_base = 02:00
19 #预先创建trunk文件的时间间隔
20 trunk_create_file_interval = 86400
21 #trunk创建文件的最大空闲空间
22 trunk_create_file_space_threshold = 20G
23 #启动时是否检查每个空闲空间列表项已经被使用
24 trunk_init_check_occupying = false
25 #是否纯粹从trunk-binlog重建空闲空间列表
26 trunk_init_reload_from_binlog = false
27 #对trunk-binlog进行压缩的时间间隔
28 trunk_compress_binlog_min_interval = 0
```

配置说明：

### 3.合并存储文件命名与文件结构

我们知道向FastDFS上传文件成功时，服务器返回该文件的存取ID叫做fileid，当没有启动合并存储时该fileid和磁盘上实际存储的文件一一对应，当采用合并存储时就不再一一对应而是多个fileid对应的文件被存储成一个大文件。

注：下面将采用合并存储后的大文件统称为 `Trunk文件`，没有合并存储的文件统称为 `源文件`；

请大家注意区分三个概念：

- **Trunk文件**：storage服务器磁盘上存储的实际文件，默认大小为64MB

Trunk文件文件名格式：`fdfs_storage1/data/00/00/000001` 文件名从1开始递增，类型为int；

- **合并存储文件的fileid**：表示服务器启用合并存储后，每次上传返回给客户端的fileid，注意此时该fileid与磁盘上的文件没有一一对应关系；
- **没有合并存储的fileid**：表示服务器未启用合并存储时，Upload时返回的fileid

### 3.1 合并存储后fileid的变化

在启动合并存储时服务返回给客户端的fileid也会有所变化，具体如下：

- **合并存储时fileid**：

```
1 | group1/M00/00/00/CgAEbFQWwbyIPCu1AAAFr1bq36EAAAAAQAAAAAAXH82.conf
```

可以看出合并存储的fileid更长，因为其中需要加入保存的 **大文件id** 以及 **偏移量**，具体包括了如下信息：

**file\_size**：占用大文件的空间（注意按照最小slot-256字节进行对齐）

**mtime**：文件修改时间

**crc32**：文件内容的crc32码

**formatted\_ext\_name**：文件扩展名

**alloc\_size**：文件大小与size相等

**id**：大文件ID如000001

**offset**：文件内容在trunk文件中的偏移量

**size**：文件大小

- **没有合并存储时fileid**：

```
1 | group1/M00/00/00/CmQPR1P0T4-AA9_ECDsoXi21HR0.tar.gz
```

**CmQPR1P0T4-AA9\_ECDsoXi21HR0.tar.gz**，这个文件名中，除了 **.tar.gz** 为文件后缀，**CmQPR1P0T4-AA9\_ECDsoXi21HR0** 这部分是一个base64编码缓冲区，组成如下：

**storage\_id**：ip的数值型

**timestamp**：创建时间

**file\_size**：若原始值为32位则前面加入一个随机值填充，最终为64位

**crc32**：文件内容的检验码

### 3.2 Trunk文件内部结构

trunk内部是由多个小文件组成，每个小文件都会有一个trunkHeader，以及紧跟在其后的真实数据，结构如下：

```
|||----- 24bytes-----|||
|—1byte —|—4bytes —|—4bytes —|—4bytes—|—4bytes —|—7bytes —|
|—filetype—|—alloc_size—|—filesize—|—crc32 —|—mtime —|—formatted_ext_name—|
|||-----file_data filesize bytes-----|||
|-----file_data-----|
```

## 4.合并存储空间闲空间管理

### 4.1 概述

Trunk文件为64MB（默认），因此每次创建一次Trunk文件总是会产生空余空间，比如为存储一个10MB文件，创建一个Trunk文件，那么就会剩下接近54MB的空间（TrunkHeader 会24字节，后面为了方便叙述暂时忽略其所占空间），下次要想再次存储10MB文件时就不需要创建新的文件，存储在已经创建的Trunk文件中即可。另外当删除一个存储的文件时，也会产生空余空间。

在Storage内部会为每个store\_path构造一颗以空闲块大小作为关键字的空闲平衡树，相同大小的空闲块保存在链表之中。每当需要存储一个文件时会首先到空闲平衡树中查找大于并且最接近的空闲块，然后试着从该空闲块中分割出多余的部分作为一个新的空闲块，加入到空闲平衡树中。例如：

要求存储文件为300KB，通过空闲平衡树找到一个350KB的空闲块，那么就会将350KB的空闲块分裂成两块，前面300KB返回用于存储，后面50KB则继续放置到空闲平衡树之中。

假若此时找不到可满足的空闲块，那么就会创建一个新的trunk文件64MB，将其加入到空闲平衡树之中，再次执行上面的查找操作（此时总是能够满足了）。

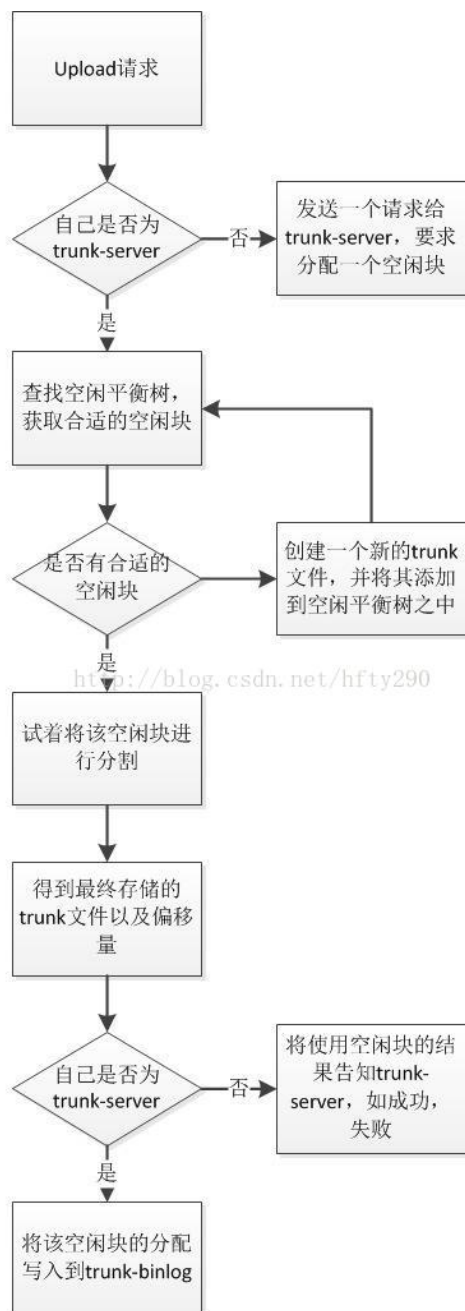
### 4.2 TrunkServer

假若所有的Storage都具有分配空闲空间的能力（upload文件时自主决定存储到哪个TrunkFile之中），那么可能会由于同步延迟导致数据冲突，例如：

Storage-A：Upload一个文件A.txt 100KB，将其保存到000001这个TrunkFile的开头，与此同时，Storage-B也接受Upload一个文件B.txt 200KB，将其保存在000001这个TrunkFile文件的开头，当Storage-B收到Storage-A的同步信息时，他无法将A.txt 保存在000001这个trunk文件的开头，因此这个位置已经被B.txt占用。

为了处理这种冲突，引入了TrunkServer概念，只有TrunkServer才有权限分配空闲空间，决定文件应该保存到哪个TrunkFile的什么位置。TrunkServer由Tracker指定，并且在心跳信息中通知所有的Storage。

引入TrunkServer之后，一次Upload请求，Storage的处理流程图如下：



## 4.3 TrunkFile同步

开启了合并存储服务后，除了原本的源文件同步之外，TrunkServer还多了TrunkBinlog的同步（非TrunkServer没有TrunkBinlog同步）。源文件的同步与没有开启合并存储时过程完全一样，都是从binlog触发同步文件。

TrunkBinlog记录了TrunkServer所有分配与回收空闲块的操作，由TrunkServer同步给同组中的其他Storage。TrunkServer为同组中的其他Storage各创建一个同步线程，每秒将TrunkBinlog的变化同步出去。同组的Storage接收到TrunkBinlog只是保存到文件中，不做其他任何操作。

TrunkBinlog文件记录如下：

```

1 | 1410750754 A 0 0 0 1 0 67108864
2 |
3 | 1410750754 D 0 0 0 1 0 67108864

```

各字段含义如下（按照顺序）：

时间戳

操作类型（A：增加，D：删除）

store\_path\_index

sub\_path\_high

sub\_path\_low

file.id（TrunkFile文件名，比如000001）

offset（在TrunkFile文件中的偏移量）

size（占用的大小，按照slot对齐）

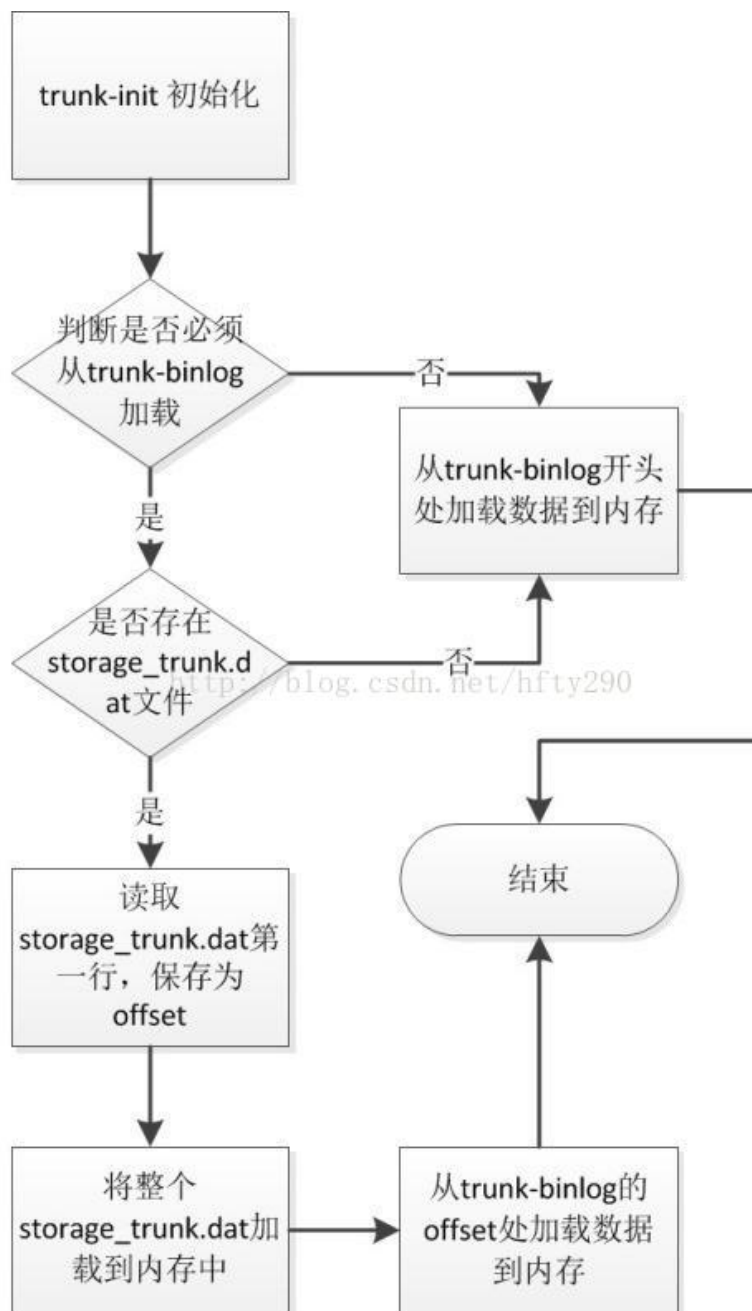
## 4.4 空闲平衡树重建

当作为TrunkServer的Storage启动时可以从TrunkBinlog文件中加载所有的空闲块分配与加入操作，这个过程就可以实现空闲平衡树的重建。

当长期运行时，随着空闲块的不断删除添加会导致TrunkBinlog文件很大，那么加载时间会很长，FastDFS引入了检查点文件storage\_trunk.dat，每次TrunkServer进程退出时会将当前内存里的空闲平衡树导出为storage\_trunk.dat文件，该文件的第一行为TrunkBinlog的offset，也就是该检查点文件负责到这个offset为止的TrunkBinlog。也就是说下次TrunkServer启动的时候，先加载storage\_trunk.dat文件，然后继续加载这个offset之后的TrunkBinlog文件内容。

下面为TrunkServer初始化的流程图：





## 4.5 TrunkBinlog压缩

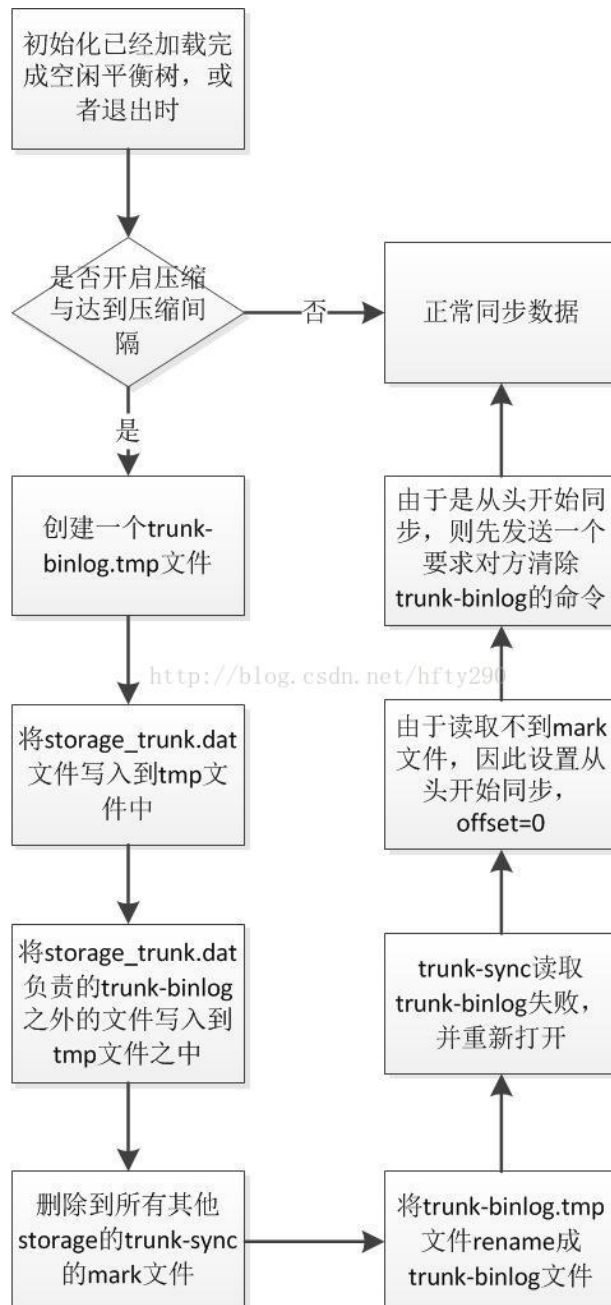
上文提到的storage\_trunk.dat既是检查点文件，其实也是一个压缩文件，因为从内存中将整个空闲平衡树直接导出，没有了中间步骤，因此文件就很小。这种方式虽然实现了TrunkServer自身重启时快速加载空闲平衡树的目的，但是并没有实际上缩小TrunkBinlog文件的大小。假如这台TrunkServer宕机后，Tracker会选择另外一台机器作为新的TrunkServer，这台新的TrunkServer就必须从很庞大的TrunkBinlog中加载空闲平衡树，由于TrunkBinlog文件很大，这将是一个很漫长的过程。

为了减少TrunkBinlog，可以选择压缩文件，在TrunkServer初始化完成后，或者退出时，可以将storage\_trunk.dat与其负责偏移量之后的TrunkBinlog进行合并，产生一个新的TrunkBinlog。由于此时的TrunkBinlog已经从头到尾整个修改了，就需要将该文件完成的同步给同组内的其他Storage，为了达到该目的，FastDFS使用了如下方法：

- 1) TrunkServer将TrunkBinlog同步给组内其他Storage时会将同步的最后状态记录到一个mark文件之中，比如同步给A，则记录到A.mark文件（其中包括最后同步成功的TrunkBinlog偏移量）。
- 2) TrunkServer在将storage\_trunk.dat与TrunkBinlog合并之后，就将本地记录TrunkBinlog最后同步状态的所有mark文件删除，如，一组有A、B、C，其中A为TrunkServer则A此时删除B.mark、C.mark。
- 3) 当下次TrunkServer要同步TrunkBinlog到B、C时，发现找不到B.mark、C.mark文件，就会自动从头转换成从头开始同步文件。
- 4) 当TrunkServer判断需要从头开始同步TrunkBinlog，由于担心B、C已经有旧的文件，因此就需要向B、C发送一个删除旧的TrunkBinlog的命令。
- 5) 发送删除命令成功之后，就可以从头开始将TrunkBinlog同步给B、C了。

大家发现了么，这里的删除TrunkBinlog文件，会有一个时间窗口，就是删除B、C的TrunkBinlog文件之后，与将TrunkBinlog同步给他们之前，假如TrunkBinlog宕机了，那么组内的B、C都会没有TrunkBinlog可使用。

流程图如下：



## 5.Tracker-Leader选择TrunkServer

### 5.1 描述

在一个FastDFS集群之中，在开启合并存储时，为了分配空间引入了一个TrunkServer角色，该TrunkServer是该Group中的一个Storage，只是该Storage要负责为该组内的所有Upload操作分配空间。为了避免不同Tracker为该Group选择了不同的TrunkServer，此时引入了Tracker-Leader角色，也就是TrunkServer最终是由Tracker-Leader来选择的，然后通知给该组内的所有Storage。

### 5.2 Tracker -Leader选择TrunkServer时机

- 1) 当组内的一个Storage状态变成Active时，并且该组还没有指定TrunkServer，在tracker\_mem\_active\_store\_server函数中触发。
- 2) 某个Tracker在经过选择，被设置成Leader时，则为当前还没有指定TrunkServer的组选择TrunkServer，在relationship\_select\_leader函数中触发。
- 3) 在定期指定的任务tracker\_mem\_check\_alive函数中，默认该函数100秒指定一次。会尝试着为每个当前还没有指定TrunkServer的组选择TrunkServer。对已经指定的组检查其TrunkServer是否还处于活动状态（根据TrunkServer与Tracker的最后心跳时间计算），若不处于活动状态，则会尝试着给该组选择一个新的TrunkServer。在tracker\_mem\_check\_alive函数中触发。

## 5.3 Tracker-Leader选择TrunkServer的过程

该过程由tracker\_mem\_find\_trunk\_server函数负责，具体操作步骤如下：

- 1) 依次向组内当前状态为ACTIVE的Storage发送TRUNK\_GET\_BINLOG\_SIZE命令的消息，来查询每个Storage当前保存的Trunk-Binlog的文件大小。来找到Trunk-Binlog文件最大的Storage。
- 2) 若该Group的最后一个TrunkServer与要设置的新的TrunkServer并非同一个，则像该新的TrunkServer发送TRUNK\_DELETE\_BINLOG\_MARKS命令，让其删除trunk-binlog同步的状态mark文件。（既然这个TrunkServer是新的，那么就要清除同步trunk-binlog的状态，使其从头同步trunk-binlog给组内的其他Storage）。
- 3) 变更该组的TrunkServer，并将修改写入到storage\_groups\_new.dat文件之中，更新该组的最后TrunkServer设置，设置TrunkServer已经变更标志，该标志使得在与Storage的心跳中通知对方TrunkServer已经变更。

# 九、存储缩略图

## 1.FastDFS主从文件

### 应用背景

使用FastDFS存储一个图片的多个分辨率的备份时，**希望只记录源图的fileid**，并能将其它分辨率的图片与源图关联。可以使用从文件方法。

### 解决办法

**名词注解**：主从文件是指文件ID有关联的文件，一个主文件可以对应多个从文件。

- 主文件ID = 主文件名 + 主文件扩展名
- 从文件ID = 主文件名 + **从文件后缀名（比如\_200\*200）** + 从文件扩展名

以本场景为例：主文件为原始图片，从文件为该图片的一张或多张缩略图。

**流程说明**：

1. 先上传主文件（即：原文件），得到主文件FID
2. 然后上传从文件（即：缩略图），指定主文件FID和从文件后缀名，上传后得到从文件FID。

**Java代码**：

```

1 public static String uploadFile(String filePath) throws Exception {
2     String fileId = "";
3     String fileExtName = "";
4     if (filePath.contains(".")) {
5         fileExtName = filePath.substring(filePath.lastIndexOf(".") + 1);
6     } else {
7         logger.warn("上传失败, 无效的扩展名");
8         return fileId;
9     }
10
11     try {
12         fileId = client.upload_file1(filePath, fileExtName, null);
13     } catch (Exception e) {
14         logger.warn("Upload file \"" + filePath + "\" fails");
15     } finally {
16         trackerServer.close();
17     }
18     return fileId;
19 }
20
21 /**
22  *
23  * @param masterFileId      FastDFS服务器返回的主文件的fileid
24  * @param prefixName        从文件后缀名 ( 比如_200_200 )
25  * @param slaveFilePath      从文件所在路径 ( 主从文件在本地都需要有对应的文件 )
26  * @return
27  * @throws Exception
28  */
29 public static String uploadSlaveFile(String masterFileId, String prefixName,
String slaveFilePath) throws Exception {
30     String slaveFileId = "";
31     String slaveFileExtName = "";
32     if (slaveFilePath.contains(".")) {
33         slaveFileExtName =
34             slaveFilePath.substring(slaveFilePath.lastIndexOf(".") + 1);
35     } else {
36         logger.warn("上传失败, 无效的扩展名");
37         return slaveFileId;
38     }
39     try {
40         slaveFileId = client.upload_file1(masterFileId, prefixName,
41                                             slaveFilePath, slaveFileExtName, null);
42     } catch (Exception e) {
43         logger.warn("Upload file \"" + slaveFilePath + "\" fails");
44     } finally {
45         trackerServer.close();
46     }
47     return slaveFileId;
48 }
49
50 public static int download(String fileId, String localFile) throws Exception {
51     int result = 0;
52     // 建立连接

```

```

53 TrackerClient tracker = new TrackerClient();
54 TrackerServer trackerServer = tracker.getConnection();
55 StorageServer storageServer = null;
56 StorageClient1 client = new StorageClient1(trackerServer, storageServer);
57 // 上传文件
58 try {
59     result = client.download_file1(fileId, localFile);
60 } catch (Exception e) {
61     logger.warn("Download file \"\" + localFile + "\"fails");
62 } finally {
63     trackerServer.close();
64 }
65 return result;
66 }
67
68 public static void main(String[] args) {
69     try {
70         String masterFileId = uploadFile("E:/1.jpg");
71         System.out.println("主文件:" + masterFileId);
72         //下载上传成功的主文件
73         download(masterFileId, "E:/master.jpg");
74         //第三个参数：待上传的从文件（由此可知，必须先本地生成从文件，再上传）
75         String slaveFileId = uploadSlaveFile(masterFileId, "_120x120",
"E:/2.jpg");
76         System.out.println("从文件:" + slaveFileId);
77         //下载上传成功的从文件
78         download(slaveFileId, "E:/slave.jpg");
79     } catch (Exception e) {
80         logger.error("upload file to FastDFS failed.", e);
81     }
82 }

```

## 2.Nginx生成缩略图

### 2.1 image\_filter模块

nginx\_http\_image\_filter\_module在nginx 0.7.54以后才出现的，用于对 JPEG、GIF和PNG 图片进行转换处理（**压缩图片、裁剪图片、旋转图片**）。这个模块默认不被编译，所以要在编译nginx源码的时候，加入相关配置信息（下面会给出相关配置）。

#### 2.1.1 检测nginx模块安装情况

```
1 [root@localhost img]# /kbb/server/nginx/sbin/nginx -v
2 nginx version: nginx/1.15.6
3 built by gcc 4.8.5 20150623 (Red Hat 4.8.5-28) (GCC)
4 built with OpenSSL 1.0.2k-fips 26 Jan 2017
5 TLS SNI support enabled
6 configure arguments: --prefix=/kbb/server/nginx --with-http_stub_status_module --
  with-http_ssl_module --with-http_realip_module
7 [root@localhost img]#
```

### 2.1.2 安装步骤

- 安装gd，HttpImageFilterModule模块需要依赖gd-devel的支持

```
1 yum -y install gd-devel
```

- 将http\_image\_filter\_module包含进来

```
1 cd /root/nginx-1.15.6
2
3 ./configure \
4 --prefix=/kbb/server/nginx \
5 --with-http_stub_status_module \
6 --with-http_ssl_module \
7 --with-http_realip_module \
8 --with-http_image_filter_module
9
10 make && make install
```

### 2.1.3 访问普通图片

需求：

假设我们图片的真实路径是在本地 `/kbb/data/img/xxxx.jpg`，下面有很多jpg格式的图片，我们希望通过访问 `/img/xxxx_100x100.jpg` 这样的请求路径可以生成宽为100，高也为100的小图，并且请求的宽和高是可变的，那么这时候需要在nginx模块中拦截该请求并返回转换后的小图，在对应的server {}段中进行配置。

nginx.conf配置如下：

```

1      location ~* /img/(.*)_(\d+)x(\d+)\.(jpg|gif|png)$ {
2          root /;
3          set $s $1;
4          set $w $2;
5          set $h $3;
6          set $t $4;
7          image_filter resize $w $h;
8          image_filter_buffer 10M;
9          rewrite ^/img/(.*)$ /kkb/data/img/$s.$t break;
10     }

```

## 重启Nginx

```
1 | /kkb/server/nginx/sbin/nginx -s reload
```

## 2.1.4 访问FastDFS图片

```

1  location ~ group1/M00/(.+)_([0-9]+)x([0-9]+)\.(jpg|gif|png) {
2      # 设置别名 (类似于root的用法)
3      alias /kkb/server/fastdfs/storage/data/;
4      # fastdfs中的ngx_fastdfs_module模块
5      ngx_fastdfs_module;
6      set $w $2;
7      set $h $3;
8
9      if ($w != "0") {
10         rewrite group1/M00(.+)_(\d+)x(\d+)\.(jpg|gif|png)$ group1/M00$1.$4
11     break;
12     }
13
14     if ($h != "0") {
15         rewrite group1/M00(.+)_(\d+)x(\d+)\.(jpg|gif|png)$ group1/M00$1.$4
16     break;
17     }
18
19     #根据给定的长宽生成缩略图
20     image_filter resize $w $h;
21
22     #原图最大2M, 要裁剪的图片超过2M返回415错误, 需要调节参数image_filter_buffer
23     image_filter_buffer 2M;
24
25     #try_files group1/M00$1.$4 $1.jpg;
26 }

```

## 重启Nginx

```
1 | /kkb/server/nginx/sbin/nginx -s reload
```



## 2.2 Nginx Image 缩略图 模块

- 本nginx模块主要功能是对请求的图片进行缩略/水印处理，支持文字水印和图片水印。
- 支持自定义字体，文字大小，水印透明度，水印位置。
- 判断原图是否是大于指定尺寸才处理。 ....等等
- 支持 jpeg / png / gif (Gif生成后变成静态图片)

### 2.2.1 安装nginx image模块

- 编译nginx前请确认您的系统已经安装了libcurl-dev libgd2-dev libpcre-dev 依赖库

```
1 | yum install -y gd-devel pcre-devel libcurl-devel
```

- 下载nginx ( 略 )
- 下载nginx image模块，并解压缩

```
1 | wget https://github.com/oupula/nginx_image_thumb/archive/master.tar.gz
2 | tar -xf master.tar.gz
```

- 添加nginx image模块

```
1 | cd /root/nginx-1.15.6
2 |
3 | ./configure \
4 | --prefix=/kbb/server/nginx \
5 | --pid-path=/var/run/nginx/nginx.pid \
6 | --lock-path=/var/lock/nginx.lock \
7 | --error-log-path=/var/log/nginx/error.log \
8 | --http-log-path=/var/log/nginx/access.log \
9 | --http-client-body-temp-path=/var/temp/nginx/client \
10 | --http-proxy-temp-path=/var/temp/nginx/proxy \
11 | --http-fastcgi-temp-path=/var/temp/nginx/fastcgi \
12 | --http-uwsgi-temp-path=/var/temp/nginx/uwsgi \
13 | --http-scgi-temp-path=/var/temp/nginx/scgi \
14 | --with-http_gzip_static_module \
15 | --with-http_stub_status_module \
16 | --with-http_ssl_module \
17 | --add-module=/root/nginx-1.15.6/nginx_image_thumb-master \
18 | --add-module=/kbb/server/fastdfs-nginx-module-1.20/src
19 |
20 | make && make install
```

### 2.2.2 访问普通图片

#### 修改nginx.conf配置文件

```
1 |         location /img/ {
```

```

2      root /kkb/data/;
3      # 开启压缩功能
4      image on;
5      # 是否不生成图片而直接处理后输出
6      image_output on;
7
8      image_water on;
9      # 水印类型：0为图片水印，1位文字水印
10     image_water_type 0;
11     # 水印出现位置
12     image_water_pos 9;
13     # 水印透明度
14     image_water_transparent 80;
15     # 水印文件
16     image_water_file "/kkb/data/logo.png";
17 }

```

### 访问方式

- nginx服务器地址：<http://192.168.10.136>
- 访问的源文件地址：<http://192.168.10.136/img/1.jpg>
- 访问的压缩文件地址：<http://192.168.10.136/img/1.jpg!c300x200.jpg>

其中 `c` 是生成图片缩略图的参数，`300` 是生成缩略图的宽度，`200` 是生成缩略图的高度

### 参数说明：

一共可以生成四种不同类型的缩略图：

- 1 `C` 参数按请求宽高比例从图片高度 10% 处开始截取图片，然后缩放/放大到指定尺寸（图片缩略图大小等于请求的宽高）
- 2
- 3 `M` 参数按请求宽高比例居中截图图片，然后缩放/放大到指定尺寸（图片缩略图大小等于请求的宽高）
- 4
- 5 `T` 参数按请求宽高比例按比例缩放/放大到指定尺寸（图片缩略图大小可能小于请求的宽高）
- 6
- 7 `W` 参数按请求宽高比例缩放/放大到指定尺寸，空白处填充白色背景颜色（图片缩略图大小等于请求的宽高）

### 2.2.3 访问FastDFS图片

```

1      location /group1/M00/{
2          alias /kkb/server/fastdfs/storage/data/;
3
4          image on;
5          image_output on;
6          image_jpeg_quality 75;
7
8          image_water on;
9          image_water_type 0;
10         image_water_pos 9;
11         image_water_transparent 80;

```

```

12         image_water_file "/kkb/data/logo.png";
13
14         # image_backend off;
15         #配置一个不存在的图片地址，防止查看缩略图时图片不存在，服务器响应慢
16         # image_backend_server http://www.baidu.com/img/baidu_jgylogo3.gif;
17     }

```

## 2.2.4 配置参数说明

```

1  image on/off 是否开启缩略图功能,默认关闭
2
3  image_backend on/off 是否开启镜像服务，当开启该功能时，请求目录不存在的图片（判断原图），将自动从
   镜像服务器地址下载原图
4
5  image_backend_server 镜像服务器地址
6
7  image_output on/off 是否不生成图片而直接处理后输出 默认off
8
9  image_jpeg_quality 75 生成JPEG图片的质量 默认值75
10
11 image_water on/off 是否开启水印功能
12
13 image_water_type 0/1 水印类型 0:图片水印 1:文字水印
14
15 image_water_min 300 300 图片宽度 300 高度 300 的情况才添加水印
16
17 image_water_pos 0-9 水印位置 默认值9 0为随机位置,1为顶端居左,2为顶端居中,3为顶端居右,4为中部居
   左,5为中部居中,6为中部居右,7为底端居左,8为底端居中,9为底端居右
18
19 image_water_file 水印文件(jpg/png/gif),绝对路径或者相对路径的水印图片
20
21 image_water_transparent 水印透明度,默认20
22
23 image_water_text 水印文字 "Power By Vampire"
24
25 image_water_font_size 水印大小 默认 5
26
27 image_water_font 文字水印字体文件路径
28
29 image_water_color 水印文字颜色,默认 #000000

```

# 十、Java客户端

## 将jar包install到maven仓库

如果通过eclipse来执行以下命令的时候，是不需要加上 `mvn` 这个命令的。

```
1 | mvn install:install-file -DgroupId=org.csource.fastdfs -DartifactId=fastdfs -
  | -Dversion=1.25 -Dpackaging=jar -Dfile=E:\01-java-soft\FastDFS\fastdfs_client-1.25.jar
```

## 配置文件fdfs\_client.conf

该配置文件要放到fdfs子目录下，名字叫fdfs\_client.conf

```
1 | tracker_server=192.168.10.135:22122
```

## 编写测试代码

```
1 | public class FastDFSClient {
2 |     private static Logger logger = Logger.getLogger(FastDFSClient.class);
3 |
4 |     private static TrackerClient trackerClient = null;
5 |     private static TrackerServer trackerServer = null;
6 |     private static StorageServer storageServer = null;
7 |     private static StorageClient1 client = null;
8 |     // fdfsclient的配置文件的路径
9 |     private static String CONF_NAME = "/fdfs/fdfs_client.conf";
10 |
11 |     static {
12 |         try {
13 |             // 配置文件必须指定全路径
14 |             String confName = FastDFSClient.class.getResource(CONF_NAME).getPath();
15 |             // 配置文件全路径中如果有中文，需要进行utf8转码
16 |             confName = URLDecoder.decode(confName, "utf8");
17 |
18 |             ClientGlobal.init(confName);
19 |             trackerClient = new TrackerClient();
20 |             trackerServer = trackerClient.getConnection();
21 |             storageServer = null;
22 |             client = new StorageClient1(trackerServer, storageServer);
23 |         } catch (Exception e) {
24 |             e.printStackTrace();
25 |         }
26 |     }
27 |     /**
28 |      * 上传文件方法
29 |      * <p>
30 |      * Title: uploadFile
31 |      * </p>
32 |      * <p>
33 |      * Description:
34 |      * </p>
35 |      *
```

```

36      * @param fileName
37      *          文件全路径
38      * @param extName
39      *          文件扩展名, 不包含 (.)
40      * @param metas
41      *          文件扩展信息
42      * @return
43      * @throws Exception
44      */
45      public static String uploadFile(String fileName, String extName, NameValuePair[]
metas) throws Exception {
46          String result = client.upload_file1(fileName, extName, metas);
47          return result;
48      }
49
50      public static void main(String[] args) {
51          try {
52              uploadFile("E:/1.jpg", "jpg", null);
53          } catch (Exception e) {
54              logger.error("upload file to FastDFS failed.", e);
55          }
56      }

```

## 十一、查看错误日志

- tracker日志

```
1 | tail -100f /kbb/server/fastdfs/tracker/logs/tracker.log
```

其中 `/kbb/server/fastdfs/tracker` 是在 `/etc/fdfs/tracker.conf` 中配置的base\_path路径

- storage日志

```
1 | tail -100f /kbb/server/fastdfs/storage/logs/storage.log
```

其中 `/kbb/server/fastdfs/storage` 是在 `/etc/fdfs/storage.conf` 中配置的base\_path路径

- nginx日志

```
1 | tail -100f /var/log/nginx/error.log
```

`/var/log/nginx/error.log` 是通过 `configure`脚本 指定的, 或者默认在 `nginx/conf` 目录下。