

EEL 5737 Principles of Computer System Design

Report for Final Project

on

Fault Tolerant File System Using Replicated Data Servers

Submitted to:

**Prof. Renato Figueurado,
Department of ECE,
University of Florida**

Submitted By:

**Anumula, Sai Raghu Vamsi
Bopanna, Vamsi Kumar
Mattapalli, Jaivardhan**

**UFID: 49939544
UFID: 51671924
UFID: 53991381**

Contents

Table of Figures	2
Abstract	3
1. Introduction and Background.....	3
2. Assumptions	3
3. Design	4
3.1. Concept	4
3.1.1. Client	4
3.1.2. Reliable layer (Mediator).....	4
3.1.3. Meta-Server	5
3.1.4. Data-Servers.....	5
3.2. Detailed design.....	5
4. Testing	8
5. Evaluation.....	13

Table of Figures

Figure 1 Flat in-memory FS - Dictionary for file data.....	Error! Bookmark not defined.
Figure 2 Flat in-memory FS Dictionary of metadata.....	Error! Bookmark not defined.
Figure 3 Hierarchical in-memory FS - Dictionary for metadata and file data storage..	Error! Bookmark not defined.

Abstract

With increasing complexity in computer systems, rate of faults and failures increase proportionally. The efforts put on controlling these faults are not proving to be as useful as expected. Hence another approach to overcome the issues with faults and failures is making the systems fault-tolerant. With this view several fault tolerant systems have been developed with various approaches. In this project, we use principle of data redundancy and replicate the data among several servers to serve a client with correct data whenever it requires. We built this fault tolerant file system with hierarchical FUSEPY file system as base and replicated data server in number. The system supports quorum approach to resolve the data from multiple servers. As expected, results of our project show that a single server fault can not only be tolerated but can be corrected at server side transparent to the client.

1. Introduction and Background

Reliable performance of computer systems has been a requirement since the construction of computing systems. Improper functioning of functional units in a system is manifested as fault in the system. The art and science of building computing systems that continue to operate satisfactorily in the presence of faults is called fault-tolerant computing. It can be achieved by using redundancy in data i.e. replicating the data in multiple sites.

Replication is the process of copying and maintaining same data objects in multiple servers that make up a replicated file system. It creates multiple copies of data items and stores them at different sites. Main idea is to increase the reliability so that if a node fails at one site, data can be accessed from a different site. Replication can improve the performance and protect the availability of applications because alternate data access options exist. The purpose of replication is to prevent damage from failures or disasters that may occur in one location, or in case such events do occur, improve the ability to recover.

Bianca Schroeder and Garth A. Gibson [1] discusses component failure in large-scale IT installation systems. They argue that these IT systems need better system design and management to cope with more frequent failures and expect increasing the levels of redundancy can solve this issue. Many

With the concepts discussed in the above paper, we have implemented a replicated file system where in we store data into multiple replicated servers and it proved to be fault-tolerant in our tests. We have also incorporated Replication Transparency and Failure Transparency in the client-server model. Replication Transparency should be mainly incorporated for the client-server systems, which replicate the data at two or more sites for more reliability. The client generally should not be aware that a replicated copy of the data exists. The clients should also expect operations to return only one set of values. Failure Transparency enables the concealment of faults, allowing user and application programs to complete their tasks despite the failure of hardware or software components. Fault tolerance is provided by the mechanisms that relate to access transparency. The client server systems are more prone to failures as any of the component may fail which may lead to degraded service or the total absence of that service.

2. Assumptions

- Meta-server is fail proof
- *Md5 hashlib* generates same checksum every time a same data is fed.
- Q_R is given as a value greater than 3 to test fault tolerance.

3. Design

3.1. Concept

The replicated file system developed in this project can be divided into four major components:

1. Client,
2. Mediator,
3. Meta-server and
4. N data-servers, where N is the number of replicated data servers.

Whenever the client requires to store/retrieve the data, it requests the mediator which then stores/retrieves the data from multiple data-servers in a transparent manner i.e. without client knowing anything about its replication. The communication between mediator and meta/data-servers takes place through XMLRPC protocol.

Error in data or a server crash is tolerated in this design. For every write checksum is calculated and stored in meta-server which is assumed to be fail proof and whenever a data is read from data servers, the checksum is calculated again and validated for errors.

The functioning of each of the components is explained below in detail.

3.1.1. Client

Client in this design is responsible for mounting the file-system and supporting all the file-system operations such as creating file/directory, writing data to a file, reading data from a file, changing meta-data etc. Unlike in base file system, where the file-system directly writes into server, here client just requests Mediator to write to the servers and requests Mediator to read the specific data from the servers. Thus Client is abstracted from what is happening at a lower level. Client is completely unaware of how mediator is writing to servers and how it is providing data from servers.

3.1.2. Reliable layer (Mediator)

A mediator acts as a bridge between the client and servers. It accepts requests from client and communicates with servers to write or read data over XMLRPC protocol. It writes data to all available data servers when a client requests it to put data. It reads from all the available data-servers and uses quorum approach to resolve data received when a client requests to get data.

It also uses a '*md5 hashlib* **checksum calculator** which is used to store the check sum in meta-server when a data is written to data-servers and re calculates checksum when data is read from data-servers and compares it with that stored in meta-server. A difference in checksum indicates that the data read from particular data-server is corrupted and should be discarded.

Mediator makes use of quorum approach where in it decides which data to be returned to client based on number of data-servers that respond with similar and correct data. Only if the number is above Q_R (minimum number of servers that should respond with correct data), the data is assumed to be correct and will be returned to client.

3.1.3. Meta-Server

There is single Meta-server in the system which stores meta-data for all the files and directories in the file system. It also stores the list-nodes data for all the file system contents. Since meta-server is assumed to be fail-proof, we used it to store checksum for the data that is stored in data servers.

3.1.4. Data-Servers

Data-server cluster consists of replicated N servers which have copies of same data. Mediator writes to each of them when the client requests to put data onto servers. All these servers respond to mediator with their own data which will be validated by the mediator

3.2. Detailed design

To make the file system replication transparent and failure transparent, we created a reliable layer, a mediator which isolates client from data servers and thereby creates abstraction and makes the system modular. Client only communicates with this Mediator for storing data or reading the data.

Mediator takes the responsibility to write the contents to meta/data servers and return the contents to client from meta/data servers.

We have used 'md5' checksum generator from 'hashlib' library for generating Checksum while writing data into data servers and while validating the received data by comparing with checksum stored during write operation.

The control flow for different operations is explained below.

Storing Contents:

When application wants to store some data, client invokes reliable layer with contents it wants to store and type of contents i.e. meta data or list nodes or data etc. Reliable layer then checks the type of contents and writes the contents to appropriate server. If the contents to be written are meta-data or list nodes, then they are written to meta server otherwise, they are written to all the data servers.

The writing to data servers is not exactly same as writing to meta-server. Whenever the client requests a *put*, mediator calculates checksum for the contents and writes it to meta-server and writes data and its corresponding checksum to data-servers. To maintain reliability, Mediator tries to make sure that data is written to at least Q_w number of data-servers.

For writing to a data server, mediator tries to connect to server for a small fixed number of times. If not connected, it enters into a list which contains details of servers to which write failure was observed and proceeds with attempting write to next data-server. If the total number of servers to which data is successfully written is less than Q_w , then Mediator tries to write to failed servers once again. If the write to these failed servers still fails, then Mediator doesn't attempt to write them again.

POCSD Final Project

Client Layer -

```
put(key,value):  
    reliable_put(path,key,value) # Call reliable put
```

Reliable Layer -

```
Reliable_put(path, key, value):  
    if META_DATA or LIST_NODES:  
        Write to Meta_Server.  
    else  
        Update_Checksum() #Calc. Checksum and write Checksum to Meta_Server  
  
    for N Servers:  
        Try to connect to server  
        if connected:  
            Write data to that server  
            Add Server name to LIVE_SERVERS list  
        else:  
            wait for short period and try atleast 5 times.  
            if still not connected:  
                Add server name to FAILED_SERVERS.  
  
    if Number of LIVE_SERVERS <  $Q_w$ :  
        Re-attempt write to failed servers (as in above loop).  
  
Update_Checksum() :  
    Calculate Checksum using HASHLIB library.  
    Write Checksum along with Path and Key to Meta_Server.
```

Reading Contents:

When application wants to read some data, client invokes reliable layer with contents it wants to read and type of contents i.e. meta data or list nodes or data etc. Reliable layer then checks the type of contents and reads the contents from appropriate server. If the contents to be read are meta-data or list nodes, then they are read from meta server otherwise, they are fetched from all the data servers.

The fetching of contents from data servers is not exactly same as reading from meta-server. Whenever the client requests a *get*, Mediator reads from all the replicated data-servers and validates each returned data against the checksum stored in meta-server. To maintain reliability, Mediator tries to make sure that data is returned by at least Q_R number of data-servers.

For reading from a data server, mediator tries to connect to server for a small fixed number of times. If not connected, it enters server name into a list which contains details of servers to which read failure was observed and proceeds with attempting read from next data-server. If the total number of servers from which data is successfully read is less than Q_R , then Mediator returns NONE to client, otherwise, it validates checksum for received data. If no server responded with correct data, then Mediator returns NONE to client. Otherwise, it tries to write correct data to faulty servers and returns good data to Client

POCSD Final Project

Client Layer -

```
get(key):  
    reliable_get(path,key) # Call reliable get
```

Reliable Layer -

```
reliable_get(path, key):  
    if META_DATA or LIST_NODES:  
        Get data from Meta_Server.  
    else  
        Update_Checksum() #Calc. Checksum and write Checksum to Meta_Server  
  
    for N Servers:  
        Try to connect to server  
        if connected:  
            Read data from that server  
            Add Server name to LIVE_SERVERS list  
            Add Data to RECEIVED_DATA list  
        else:  
            wait for short period and try atleast 5 times.  
            if still not connected:  
                Add server name to FAILED_SERVERS list.  
                Add NONE in corresponding RECEIVED_DATA list  
  
    if Number of LIVE_SERVERS < QR:  
        return NONE to Client  
  
    Data_To_Client = Validate_Checksum (path, key, RECEIVED_DATA)  
    #Check how many servers returned correct data and repair faulty  
    Server  
  
    return Data_To_Client to Client
```

```
Validate_Checksum(path, key, RECEIVED_DATA):  
    Read Checksum using Path and Key from Meta_Server.  
    for all elements in RECEIVED_DATA:  
        Calculate Checksum using HASHLIB library  
        Compare with Checksum from Meta_Server  
        if matched:  
            Increment GOOD_DATA_CNT.  
        else:  
            Increment BAD_DATA_CNT.  
  
    if GOOD_DATA_CNT = 0:  
        return NONE  
    else:  
        Write correct data to faulty Data_Server/s #REPAIR HAPPENS HERE  
        return (Good element from RECEIVED_DATA)
```

4. Testing

TEST 1

STEPS -1: Running 4 data servers and 1 meta server. with QR=2 and QW=2

The screenshot shows two terminal windows. The top window, titled 'rv@rv-ubuntu: ~/junk/fusepy', displays the output of a script that starts four data servers. The output shows the servers starting at ports 51236, 51235, 51238, and 51237, and then receiving four POST requests from 127.0.0.1. The bottom window, titled 'Terminal', shows the output of a script that starts a meta server at port 51234, which then receives four POST requests from 127.0.0.1.

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/jun... x rv@rv-ubuntu: ~/jun... x rv@rv-ubuntu: ~/jun... x rv@rv-ubuntu: ~/jun... x

['.', '..']
CallCount 25 Time 15:26:41.930670
In function readdir()
['.', '..']
CallCount 26 Time 15:26:41.931654
In function readdir()
['.', '..']
CallCount 27 Time 15:26:41.932651
In function readdir()
['.', '..']
CallCount 28 Time 15:26:41.934426
In function readdir()
['.', '..']
CallCount 29 Time 15:26:41.936182 arguments:<class '__main__.Memory'> /autorun.in
type'>
In function getattr()

Terminal
['dataserver.py', '51235', '51236', '51237', '51238']
SERVER is UP at port: 51236
SERVER is UP at port: 51235
SERVER is UP at port: 51238
SERVER is UP at port: 51237
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -

Terminal
['metaserver.py', '51234']
SERVER is UP at port: 51234
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
127.0.0.1 - - [07/Dec/2015 15:26:41] "POST /RPC2 HTTP/1.1" 200 -
```

STEP-2: Creating File1. File System reports that good data has been written.

The screenshot shows two terminal windows. The top window, titled 'rv@rv-ubuntu: ~/junk/fusepy', displays the output of a script that creates a file named 'File1'. The output shows the file being created and the data being written to it. The bottom window, titled 'rv@rv-ubuntu: ~/junk/fusepy/fusemount', shows the output of a script that mounts the file system and verifies the data. The output shows the file system being mounted and the data being verified.

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x

['S'\np0\n.', "S'\np0\n.", "S'\np0\n.", "S'\np0\n."]
server 0 1d0c8e1daf383924f61cd3948f9077c1
data S"S'\np0\n."
p0
.
date_server 0 is good
server 1 1d0c8e1daf383924f61cd3948f9077c1
data S"S'\np0\n."
p0
.
date_server 1 is good
server 2 1d0c8e1daf383924f61cd3948f9077c1
data S"S'\np0\n."
p0
.
date_server 2 is good
server 3 1d0c8e1daf383924f61cd3948f9077c1
data S"S'\np0\n."
p0
.
date_server 3 is good

rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x

rv@rv-ubuntu: ~/junk/fusepy/fusemount$ clear
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ echo "This if file1" > file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
```


POCSD Final Project

STEP-3:

Here the data has been corrupted on server 0 [51235] & server 2 [51237]. Hence the subsequent "cat file1" in the bottom terminal still prints out correct data. The log file shows in which servers data has been corrupted and corrects them using good data at the end of the step.

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x
valid_checksum
Data server 0 is corrupted
server 1 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THIs if file1\\n'\np0\n."
p0
.
date_server 1 is good
server 2 5ef9f8ddabe036c095b1e0e8ca4045c7
data S"S'This file is corrupted'\np0\n."
p0
.
valid_checksum
Data server 2 is corrupted
server 3 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THIs if file1\\n'\np0\n."
p0
.
date_server 3 is good
CallCount 46 Time 15:37:36.873535 arguments:<class '__main__.Memory'> /file1 <t
ype 'NoneType'>
In function getattr()

127.0.0.1 - - [07/Dec/2015 15:37:36] "POST /RPC2 HTTP/1.1" 200 -

rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After data corruption on 2 servers^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
THIs if file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
```

STEP 4: Here the data is corrupted on 3rd server, so it reports that data has been corrupted on 3 servers. Still in this case if uses the good data on the 4th server to correct rest of them.

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x
Data server 0 is corrupted
server 1 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THIs if file1\\n'\np0\n."
p0
.
date_server 1 is good
server 2 5ef9f8ddabe036c095b1e0e8ca4045c7
data S"S'This file is corrupted'\np0\n."
p0
.
valid_checksum
Data server 2 is corrupted
server 3 5ef9f8ddabe036c095b1e0e8ca4045c7
data S"S'This file is corrupted'\np0\n."
p0
.
valid_checksum
Data server 3 is corrupted
CallCount 50 Time 15:38:40.119641 arguments:<class '__main__.Memory'> /file1 <t
ype 'NoneType'>
In function getattr()

127.0.0.1 - - [07/Dec/2015 15:38:40] "POST /RPC2 HTTP/1.1" 200 -

rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After data corruption on 2 servers^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
THIs if file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After data corruption on 3 server^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
THIs if file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
```

POCSD Final Project

STEP 5: data corrected on all the servers

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x
p0
date_server 0 is good
server 1 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THis if file1\\n'\np0\n."
p0
.
date_server 1 is good
server 2 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THis if file1\\n'\np0\n."
p0
.
date_server 2 is good
server 3 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THis if file1\\n'\np0\n."
p0
.
date_server 3 is good
CallCount 54 Time 15:40:20.635836 arguments:<class '__main__.Memory'> /file1 <t
ype 'NoneType'>
In function getattr()

127.0.0.1 - - [07/Dec/2015 15:38:40] "POST /RPC2 HTTP/1.1" 200 -
rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ After data corruption on 2 servers^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
THis if file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$ After data corruption on 3 server^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
THis if file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
THis if file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$
```

STEP6 : data corrected on all the servers

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x
["S''\np0\n.", "S''\np0\n.", "S''\np0\n.", "S''\np0\n."]
server 0 1d0c8e1daf383924f61cd3948f9077c1
data S"S''\np0\n."
p0
.
date_server 0 is good
server 1 1d0c8e1daf383924f61cd3948f9077c1
data S"S''\np0\n."
p0
.
date_server 1 is good
server 2 1d0c8e1daf383924f61cd3948f9077c1
data S"S''\np0\n."
p0
.
date_server 2 is good
server 3 1d0c8e1daf383924f61cd3948f9077c1
data S"S''\np0\n."
p0
.
date_server 3 is good
rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x
rv@rv-ubuntu:~/junk/fusepy/fusemount$ clear
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ echo "THis if file1" > file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$
```

TEST- 2

POCSD Final Project

Step1: Terminate two servers and try to read the data. The data is succesfully returned even in this case as it is available in other servers.

However when 3 servers are terminated, no data is returned.

```

rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/junk/fusepy$
rv@rv-ubuntu: ~/junk/fusepy$
rv@rv-ubuntu: ~/junk/fusepy$ python close-server.py 51235 51237
terminated
rv@rv-ubuntu: ~/junk/fusepy$

127.0.0.1 - - [07/Dec/2015 15:40:20] "POST /RPC2 HTTP/1.1" 200 -
rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After data corruption on 2 servers^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
This is file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After data corruption on 3 server^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
This is file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
This is file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$

```

2_stop

```

rv@rv-ubuntu: ~/junk/fusepy$ fusemount
rv@rv-ubuntu: ~/junk/fusepy$ fusemount
rv@rv-ubuntu: ~/junk/fusepy$ fusemount
rv@rv-ubuntu: ~/junk/fusepy$ fusemount
valid_checksum
Data server 0 is corrupted
server 1 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THIS if file1\\n'\np0\n."
p0
.
date_server 1 is good
server 2 5f6ea1cd1bb7b705ac12b2709baf20b
data S"S'corrupted'\np0\n."
p0
.
valid_checksum
Data server 2 is corrupted
server 3 9d03317682ce8849b5d3e7c1a2aecc20
data S"S'THIS if file1\\n'\np0\n."
p0
.
date_server 3 is good
Server Down
Server Down
CallCount 36 Time 15:53:20.644486 arguments:<class '_main_.Memory'> /file1 <
127.0.0.1 - - [07/Dec/2015 15:51:05] "POST /RPC2 HTTP/1.1" 200 -
rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After 2 servers are terminated^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat file1
THIS if file1
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ After 3 servers terminated^C
rv@rv-ubuntu: ~/junk/fusepy/fusemount$
rv@rv-ubuntu: ~/junk/fusepy/fusemount$ cat

```

POCSD Final Project

3 stop

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x
Trying to reconnect to the server 2
Trying to reconnect to the server 2
Trying to reconnect to the server 3
Trying to reconnect to the server 3
Trying to reconnect to the server 3
Trying to reconnect to the server 3
Trying to reconnect to the server 3
R__DATA -> ['S'corrupted'\np0\n.', "S'This if file1\\n'\np0\n.", "S'corrupted'\np0\n.", "S'corrupted'\np0\n."]
Traceback (most recent call last):
  File "/home/rv/junk/fusepy/fuse.py", line 420, in _wrapper
    return func(*args, **kwargs) or 0
  File "/home/rv/junk/fusepy/fuse.py", line 500, in read
    offset, fh)
  File "/home/rv/junk/fusepy/fuse.py", line 887, in __call__
    ret = getattr(self, op)(path, *args)
  File "rem_tree_ft.py", line 330, in read
    return self.FS.read_file(path, offset, size)
  File "rem_tree_ft.py", line 180, in read_file
    return filenode.get("data")[offset:offset + size]
TypeError: 'NoneType' object has no attribute '__getitem__'

127.0.0.1 - - [07/Dec/2015 15:53:20] "POST /RPC2 HTTP/1.1" 200 -

rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ After 2 servers are terminated^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
THIS if file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$ After 3 servers terminated^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
cat: file1: Bad address
rv@rv-ubuntu:~/junk/fusepy/fusemount$
```

TEST- 3

1_restart

```
rv@rv-ubuntu: ~/junk/fusepy
rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x rv@rv-ubuntu: ~/ju... x
rv@rv-ubuntu:~/junk/fusepy$
rv@rv-ubuntu:~/junk/fusepy$
rv@rv-ubuntu:~/junk/fusepy$
rv@rv-ubuntu:~/junk/fusepy$ python dataserver.py 51238
['dataserver.py', '51238']
SERVER is UP at port: 51238

127.0.0.1 - - [07/Dec/2015 15:53:20] "POST /RPC2 HTTP/1.1" 200 -

rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu: ~/junk/fusepy/fusemount x rv@rv-ubuntu: ~/junk/fusepy/fusemount x
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ After 2 servers are terminated^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
THIS if file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$ After 3 servers terminated^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
cat: file1: Bad address
rv@rv-ubuntu:~/junk/fusepy/fusemount$
```

1s_restart

```

rv@rv-ubuntu: ~/junk/fusepy
data S'S'This if file1\n'\np0\n.'"
p0
.
date_server 1 is good
server 2 5f6ea1c1d1bb7b705ac12b2709baf20b
data S'S'corrupted'\np0\n.'"
p0
.
valid_checksum
Data server 2 is corrupted
server 3 734f4d6ea264c45fc7eabaf4d2502b7e
data S'(dp0\n.'"
p0
.
valid_checksum
Data server 3 is corrupted
Server Down
Server Down
CallCount 44 Time 15:56:02.997351 arguments:<class '__main__.Memory'> /file1 <t
ype 'NoneType'>
In function getattr()

127.0.0.1 - - [07/Dec/2015 15:53:20] "POST /RPC2 HTTP/1.1" 200 -

rv@rv-ubuntu: ~/junk/fusepy/fusemount
rv@rv-ubuntu:~/junk/fusepy/fusemount$ 
rv@rv-ubuntu:~/junk/fusepy/fusemount$ 
rv@rv-ubuntu:~/junk/fusepy/fusemount$ 
rv@rv-ubuntu:~/junk/fusepy/fusemount$ 1 served restarted^C
rv@rv-ubuntu:~/junk/fusepy/fusemount$ 
rv@rv-ubuntu:~/junk/fusepy/fusemount$ cat file1
This if file1
rv@rv-ubuntu:~/junk/fusepy/fusemount$ 

```

5. Evaluation

We have implemented a file system which is fault tolerant by virtue of redundancy in data. We achieved it by replicating data into multiple servers and checking all of them when read. Using the inputs provided by User, we connect to the data servers and meta-servers. Q_R and Q_W indicates at least how may servers should respond with correct data and at least to how many servers data has to be written.

We found that our file system can tolerate errors in a server or even a server cache i.e. client receives correct data even if a server crashes. Even better thing is that the faulty server is corrected using the correct data from other servers. We have made it more robust by correcting all other faulty servers even if single server responds with correct data.

The mechanism of how this works is explained below:

- If Number of servers responded to a read request is less than Q_R , return NONE to Client.
- If Number of servers responded to a read request is greater than or equal to Q_R then,
 - If at least one Data server responded with correct data, repair all other servers and return correct data to client
 - Otherwise, return NONE to client

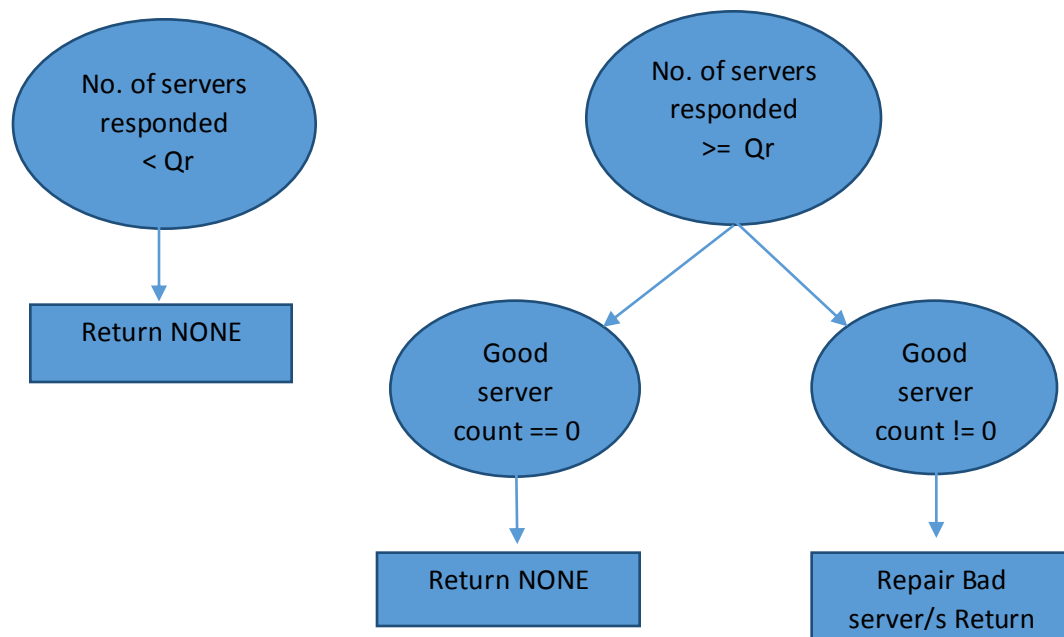


Figure: Return values for GET under various conditions

id	Operation	QR	QW	# of servers with good data before operation	# of servers with good after operation	No.of server running out of 4 servers	result
1.	GET	2	2	4	4	4	ok
2	stop servers1,2	2	2	4	4	2	-
3	PUT	2	2	2	2	2	OK - Logs as PUT fail but doesn't give out error
4	server1 restarted	2	2	2	2	3	-
5	PUT	2	2	2	3	3	OK - Logs as PUT fail but doesn't give out error
6	server2 restarted	2	2	3	3	4	-
7	GET	2	2	3	4	4	ok - uses good server data to correct the data on 4th server

If Number of servers to which write was successful is less than Q_w , then a Warning message is printed on screen. The data in servers to which it is properly written still stays and can be read when needed.