

# CloSSer

Cloud Storage Services

Team members:

Icaro Alzuru

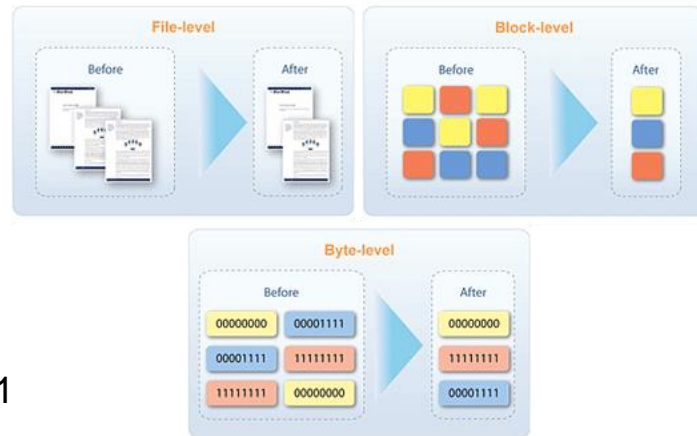
Qi Cai

Huixiang Chen

# Big Data - Storage Services

What are the storage services required by Big Data companies?

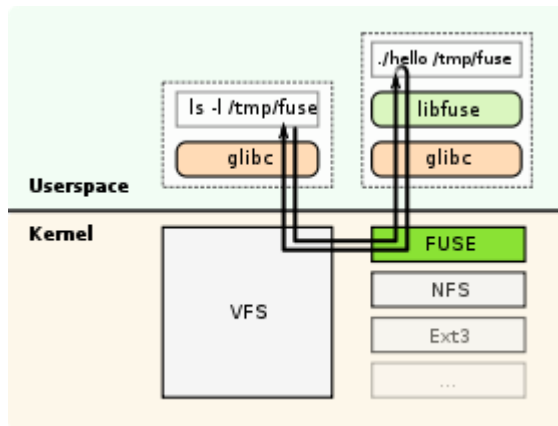
- Backup
- Replication / Synchronization. i.e. Netflix
- Deduplication: Save storage. i.e. Hotmail
  - compression save data-->:1, deduplication-->up to 25:1
- Fault tolerance
  - eg. single point failure



The solution: **CloSSer**

# Related Work

- FUSE File System



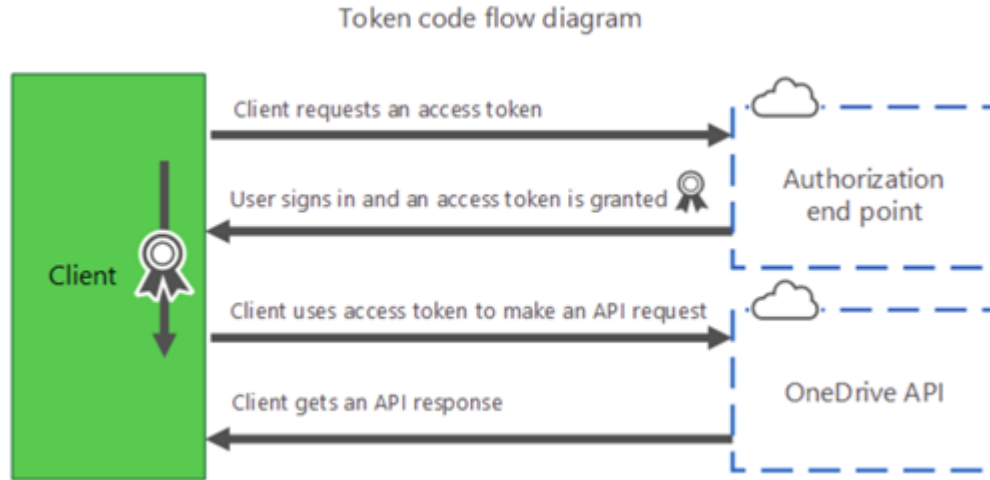
- Cloud Storages APIs

- OneDrive
- Dropbox
- Google Drive
- Box

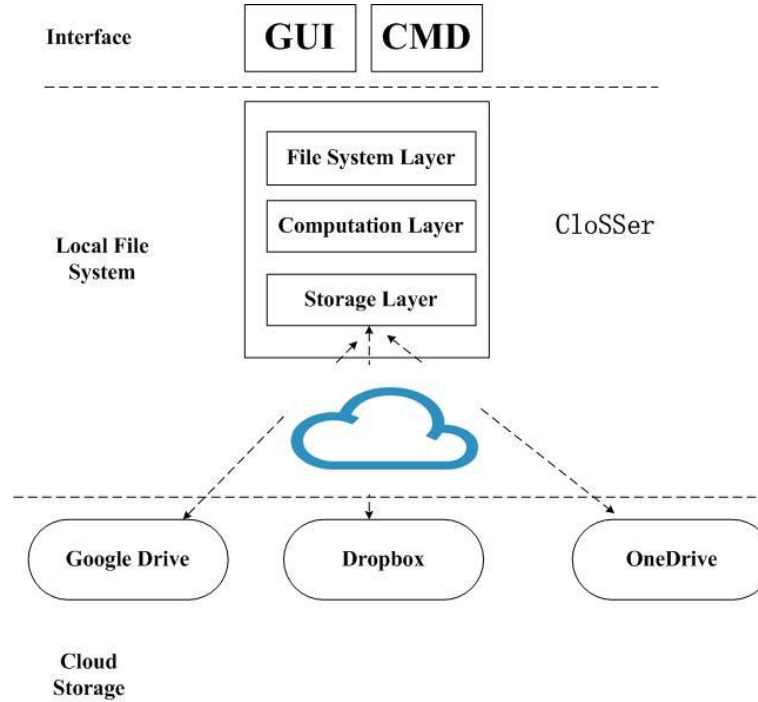


# Python API authentication

- OAuth
- OAuth2



# CloSSer Architecture



# Use Cases

- Creating several copies of some dataset for backup purposes.
  - Fault tolerance
- Saving space in environments where the same file is copied many times: Mail servers.
  - Deduplication
  - Decreases the cost of storage and communication bandwidth
- High availability environments: Netflix (many copies of the same movies)
  - Synchronized removals and copies in many servers.

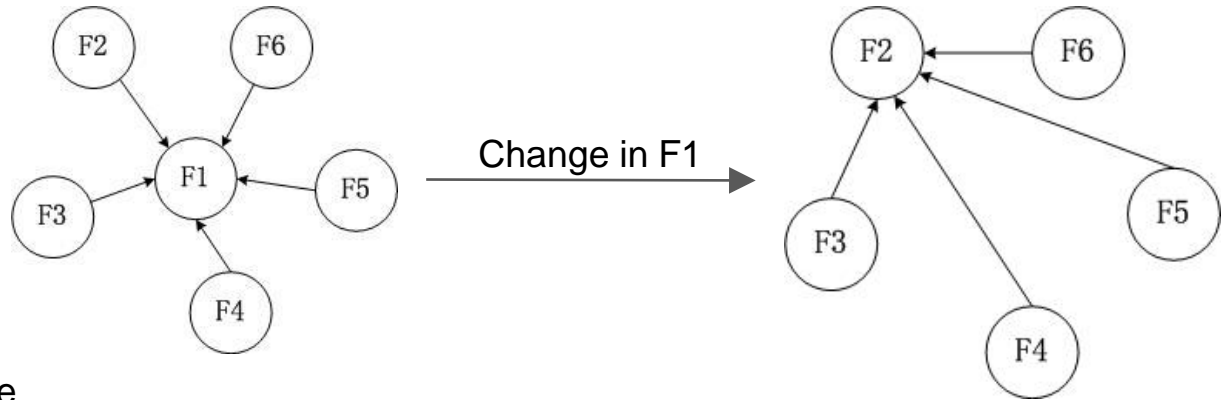
# Design and Implementation

- Storage Lager API
  - syncMirror
  - checkHealth
  - uploadFile
  - uploadMetadata
  - uploadFileAndMetadata
  - downloadFile
  - downloadMetadata
  - existsFile
  - deleteFile
  - downloadPointers
  - detectFile
  - deletePointer
  - downloadPointer

# Design and Implementation

- Computational Layer - Algorithms

- Deduplication



- Uploading a file
  - Downloading a file
  - Deleting a file
  - Rebuilding a file



# Dropbox module (Design and Implementation)

- Challenge:
  - Upload a file is a two process step: Create file/get handle, use handle to insert content
- Powerful and complete API

# Google Drive module (Design and Implementation)

- PyDrive was preferred instead of the Rest API
- Challenges:
  - PyDrive does not provide delete function
  - Folders do not exist
- Buckets are folders in the root directory

# OneDrive module (Design and Implementation)

- Main Challenge:
  - Supports only one connection to the cloud: Error "address in use"
- Buckets are folders in the root directory

# Administration Interface

(Design and Implementation)

## CloSer Administration

User

Password

Login



## CloSer: Cloud Storage Services

CloSer service state ☐ On/Off

Mounting directory

There are 3 Cloud nodes configured:

Id	Type	Bucket Name	Configure
0	mydropbox	/bucket1	<input type="button" value="cfg"/>
1	myonedrive	/bucket	<input type="button" value="cfg"/>
2	googledrive	bucket2	<input type="button" value="cfg"/>

Deduplication ☒ On/Off

Pointer directory



## Cloud Modification

Nodetype:

Nodeloc:

Accesskey:

Secretkey:

Nodeid:

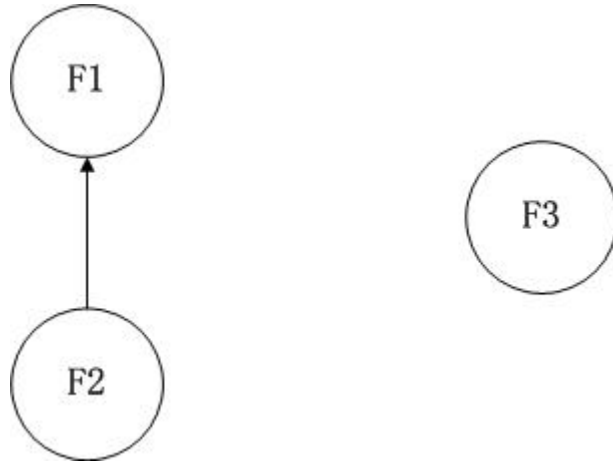
Nodekey:

Bucketname:

Modify

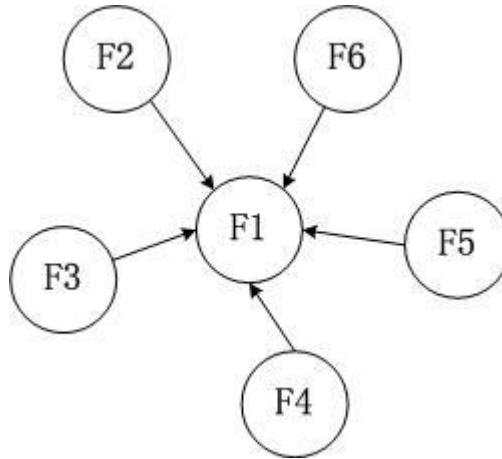
# Deduplication Algorithm

- Categorizations of file
  - non duplicated file
  - original file
  - pointer file



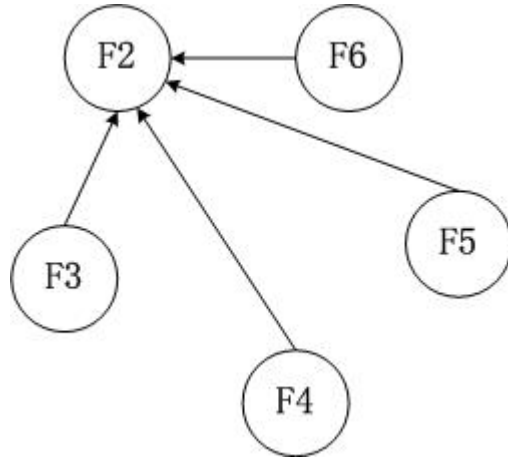
# Deduplication Algorithm

- A typical case of file modification
  - change the content of an original file



# Deduplication Algorithm

- A typical case of file modification
  - change the content of an original file



# Algorithm of downloading a file

- The workflow of downloading a file contains the following two steps:
  - Testing if the file is a pointer/leaf file on the cloud.
  - If the file is a pointer/leaf file, download the content and metadata of the original file it points to. Otherwise, the file is an original or non duplicated file. In these situations, we just download its content and metadata



# Algorithm of deleting a file

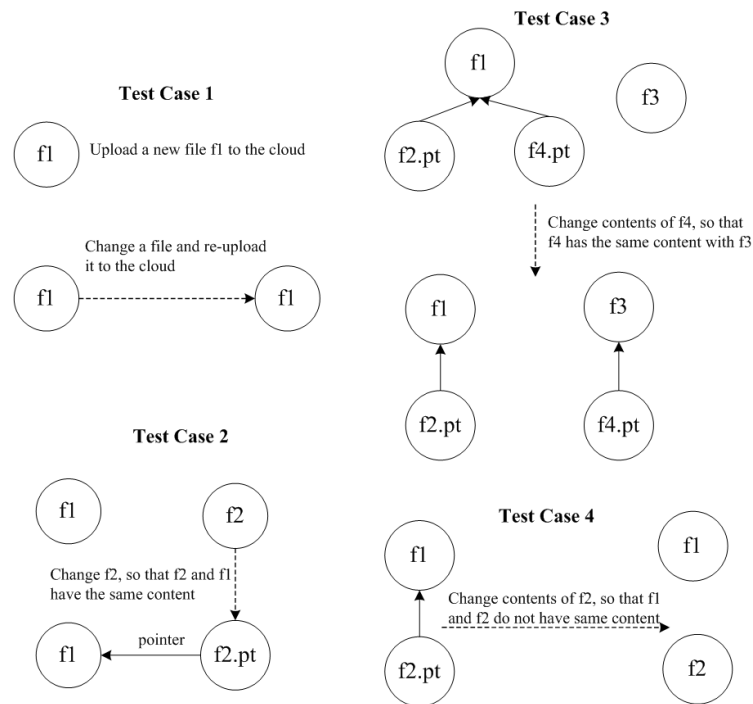
- The workflow of deleting a file can be categorized into the following three types:
  - The file is a pointer/leaf file. In this case, we just delete the corresponding pointer file on the cloud.
  - The file is a non duplicated file. In this case, we delete the content and metadata related to this file.
  - The file is an original file. In this case, we first delete the content and metadata related to this file. Then we pick up one of the pointer files as the new original file and update the contents of the rest pointer files so that they point to the new original file.

# Algorithm of rebuilding a node

- The workflow of rebuilding a failure node consists of the following steps:
  - Find an available node in the list of spare nodes to replace the failure node.
  - Find an available node in the list of survival nodes and download all the files into local directory.
  - Upload the files into the newly added node to recover the lost data

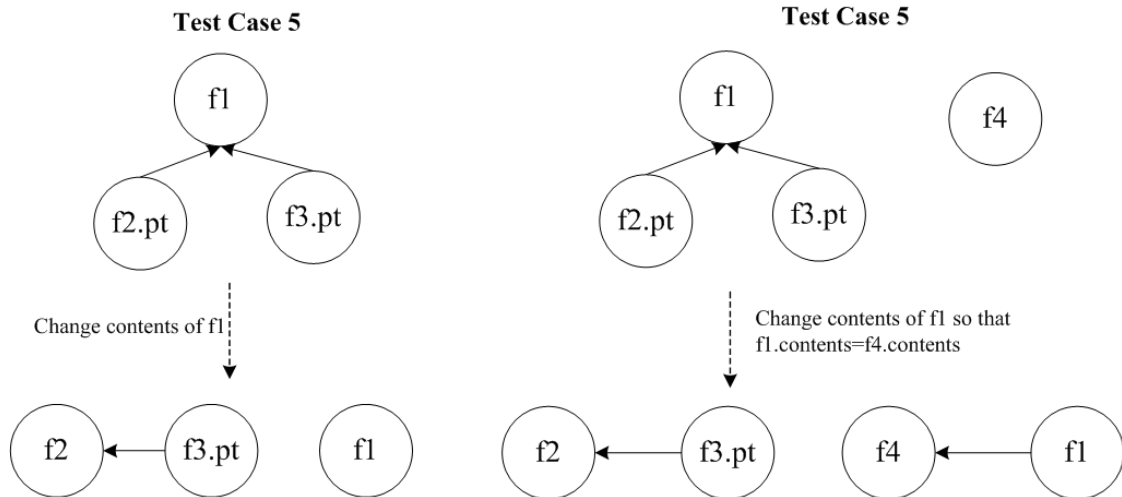
# Product Evaluation

- Standard Black Box tests for methods functionalities
- Test cases for deduplication
  - The solid arrow means two files are identical, one file points to another
  - The dashed arrow means transform from one state to another



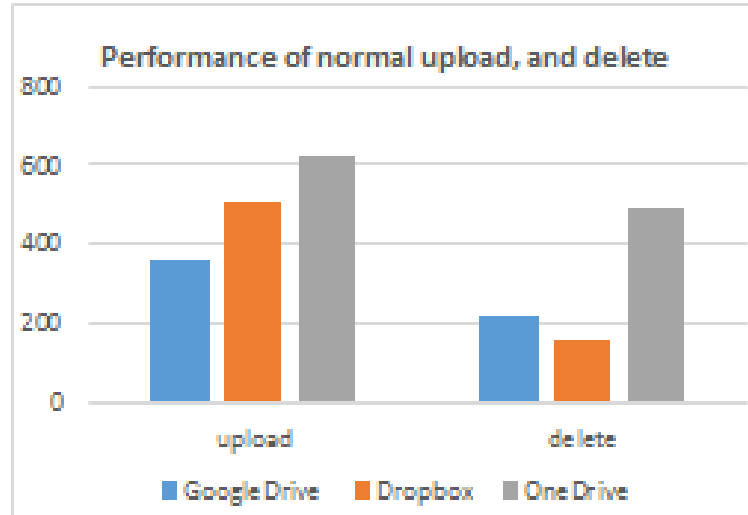
# Product Evaluation

## Test cases for deduplication (cont.)



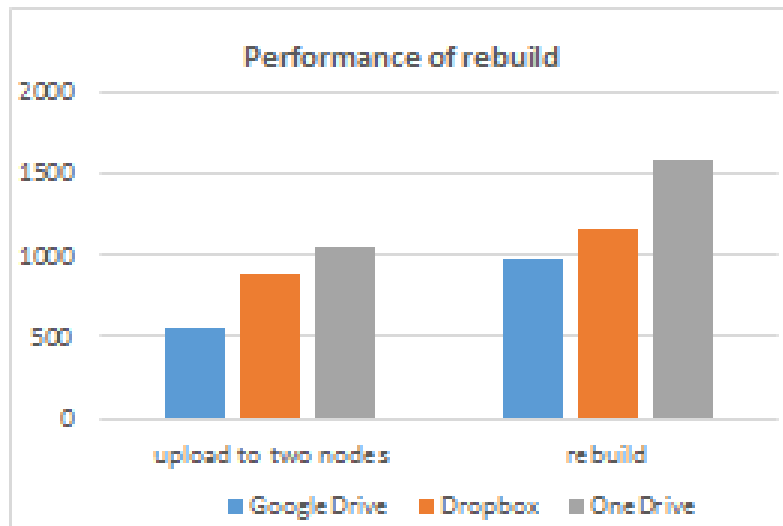
# Product Evaluation

- Performance Comparison
  - Upload 500MB of data consisted of 73 mp3 songs
  - Tested the response of uploading files and deleting files



# Performance evaluation-cont

- Run CloSSer using two nodes, and one spare node
- Fail one node on purpose by renaming the folder
- Run the reboot process and test the response time



# Future Results

- There is a lot of space for performance improvement
  - Deduplication
  - Rebuild
- Better administration interface
  - Drag and Drop
- New functionalities
  - Compression
  - Chunks distribution of big files

# Conclusions

- CloSSer can be used as a cloud backup tool
- It saves space, through deduplication, when multiple copies of the same file are present: decreasing the cost of storage and communication bandwidth
- It can be used in high availability environments to synchronize multiple servers.
- The graphical interface increases the productivity and makes easier the life of the CloSSer administrator
- Modules for OneDrive, Dropbox, and Google Drive were implemented
- Google Drive is faster for uploading files while Dropbox is faster for deleting files.
- CloSSer could be exploited commercially and for research



Thank you!

Questions?