# OS 2014    Exercise 4

# A Caching File System

## Caching algorithms

Cache replacement algorithms are algorithms that a computer program or hardware component can use to manage a cache of information stored on the computer. When the cache is full and new information arrives, the algorithm must choose which items to discard to make room for the new ones.

There are several common algorithms to choose which item to discard:
- Least Recently Used (LRU) discards the least recently used item first. This is one of the most used algorithms.
- Most Recently Used (MRU) discard the most recently used item first (opposite of LRU).
- Random Replacement (RU) randomly selects a candidate item and discards it.
- **Least Frequently Used (LFU)** counts how often each item is used, and discards the least frequently used one first.

In this exercise you will implement the LFU algorithm. There are a few ways to implement it, we will mention two of them:
1. **Array based implementation** – on each access to the item, we increase its access counter. When an item should be discarded from the cache, we search for the item with the smallest counter and remove it.
2. **Sorted List based implementation** – the idea is to maintain a list sorted according to the counter's value. This means that the list should be updated on each access to an item. However, the removal is simple and done by popping the last item and pushing the new item.

## Background Reading and Resources

1. Read and understand the relevant TA class is a first step
2. FUSE API documentation
3. The Big Brother File System (later will be referred as "bbfs") is a good tutorial.

## Assignment

In this exercise you will implement a simple single threaded caching file system using FUSE, and a simple algorithm to manage the file system's buffer cache.  This file system will actually become part of the file system you see and can use on your machine, but you'll be responsible for how it works!

Caching file systems saves files or parts of them to a memory segment called the buffer cached, which is stored in the main memory. That way, if the user wishes to access information that is saved in the cache, it can be retrieved from the memory instead of the disk, thus reducing the number of disk accesses required.

In order to make this system efficient, we want to ensure that the information that will be requested the most is saved to the cache, thus ensuring the disk will be accessed as seldom as possible. Therefore, we will use the LFU algorithm to choose which item to discard from the cache when we need more space.

Managing a complete cached file system is an interesting task, but away too big for an exercise in our course. Therefore, the supported functions in our file system will be read-only (except rename), and you may assume that the files won't be changed at all after mounting your file system. Moreover, while in reality we would cache all the files' inode data (including the stat, size, etc), here you are requested to cache only the content of the file (e.g. the data received from "read" request).

To manage a caching file system, you will need two relevant parameters – the number of blocks (*numberOfBlocks*) to save in the cache, and the size of each such block (*blockSize*). This enables you to cache  *numberOfBlocks\* blockSize* bytes. You may assume that this is not too big a value (so you can allocate it on the heap without problems). The blocks in the cache must be aligned according to the *blockSize*. This means that when you retrieve file content from the disk (following a read request), you will always be asked for complete blocks (except at the end of the file) with an offset that is a multiple of *blockSize*.

Let's demonstrate the behavior of such a system, when *numberOfBlocks* is 10 and *blockSize* is 1024. We start with an empty buffer cache. Assume that we also have a file named "tmp" with 1500 bytes in it.

1. If a user tries to read the last 500 bytes, we must retrieve all the file and save both its blocks in the buffer cache (we can't read only the required 24 bytes from the first block because we can't read partial blocks).
2. If a user tries to read the last 50 bytes of the file, we will retrieve from the disk and cache the second block of the file (bytes 1024-1499). However, we will return only its last 50 bytes to the user.
3. If a user does 2 and then tries to read all the file:
   In the first read, the file system retrieves and caches the second block, returning only its last 50 bytes. In the second read, the file system retrieves and caches only the first block, because the second block is already in the buffer cache and shouldn't need to be read again. Of course, the function will return all the file.

You may see from this example that the function always returns to the user all the information he requested, but retrieves from the disk only the blocks that are not already in the buffer cache.

## Running your file system

Your file system will be implemented in a file called CachingFileSystem.cpp, and should be activated using the following command:

**$ CachingFileSystem *rootdir mountdir numberOfBlocks blockSize***

Where:
- *mountdir* is the mount point directory - which should be empty. The files in your file system will appear to be in this directory.
- *rootdir* is a directory (folder) containing files, which will be mounted by your filesystem at *mountdir*. This means that when using your file system, *mountdir* should respond as if it contains all the files in rootdir, and *rootdir* is where they really are. Thus the user should be able to open and read the files in *rootdir* (through your filesystem of course) using a path through *mountdir*. All the actions done on a file in *mountdir* (again, using your file system) should actually be done on the corresponding file in *rootdir* (renaming for example).
- *NumberOfBlocks* is the number of blocks in the buffer cache
- *BlockSize* is the number of bytes in each block.

## Supported functions

Your FUSE file system implementation must override the following functions (note that some of these functions may have an empty implementation):

- *getattr*
- *access*
- *open*
- *read*
- *fgetattr*
- *flush*
- *release*
- *opendir*
- *readdir*
- *releasedir*
- *rename*
- *init*
- *destroy*
- *ioctl*

The functions are described in the supplied CachingFileSystem.cpp file. While you need to implement several functions, in many of them the main task of your function is only to find an equivalent function in Linux, and to call to this function. For example, the *getattr* function initializes the stat structure with the file attributes. Pay attention that in your file system, the files (in *mountdir*) have exactly the same attributes as the real files in the root dir. Therefore, you only need to initialize the stat of the requested file in the root directory.

However, there are several functions that are a bit more complicated and we will describe them here.
- *open*. While open is a simple function, pay attention that you should use the given flags (in the *fuse_file_info*).
- *readdir*. This function implementation is a bit more complicated than a simple forward. Pay attention to the difference between what you are requested to do and what Linux's "*readdir*" does.
- *read*. The read function is the most important function in your file system. It is actually the only function that uses your cache. Most of the logic of the program, as described previously, should be in here.
- *rename*. The rename function renames a file. Pay attention that this file may have cached blocks in the file system's buffer cache. In this situation, the cache must be updated such that when there is an open and read of the renamed file, the information that was in the cache before the "rename" should be identified and should not be retrieved from the disk again. Of course, the rename function doesn't change the accesses counter to each block.
- *ioctl*. This function will write the current status of the cache to the log file (see details later). The function prints a line for every block that was used in the cache (meaning, with at least one access). Each line contains the following values separated by a single space.
    ◦ The name of the file - relative path from *mountdir*.
    ◦ Number of the block (not the number in the cache, but the enumeration within the file itself - starting with 1).
    ◦ The number of times the block was accessed (including the first time when it was retrieved, and therefore a zero value is invalid).

The order of the entries is not important and may depend on your implementation.

**For example**, the following lines could be appended to the log file after iocl is called:
1430334011 ioctl (see details about this line, later)
SomeFile.txt 2 1
myFolder/SomeOtherFile.txt 1 5

## Logging your files system behavior

Your program will log its activity to a hidden log file named "*.filesystem.log*". This file will be positioned under *rootdir*. You need to create the log in the main function, before invoking FUSE. If the file already exists, then append to it and don't recreate (empty) it.

The log should contain exactly the following information (without anything else). In each file system function that you implement (described above), when the function is called you need to write to the log: "UNIX_TIME FUNCTION_NAME" where FUNCTION_NAME is the name of the function and the UNIX_TIME is the unix time (use "time" function to obtain it). For example, when a user uses "*open*", our "*caching_open*" will be called. Our first step within the function is to write "1430374011 open" (and not *caching_open*!) to the log. Pay attention that you don't log your private functions, but only the file system functions. You should check before submitting your code that you have exactly 14 such print instructions in the code (one for each function). In addition to that, the log also contains the LFU information printed by the *ioctl* function as it was described above.

The *.filesystem.log* is designed to help us improve or check the system. The users that use your system (via the mounted directory) shouldn't be aware of the existence of this file. To do that, you must block them from using this file in any of your functions. For example, an attempt to open the log file will fail and return an error. Specifically, any attempt to access the log file should return -*ENOENT*, which is an error defined in *errno.h* meaning "no such file or directory". We will speak more about errors and the functions' return values later.

Pay attention that in order to cause the users not see the file via the *ls* command, it's not enough to block an access to the .filesystem.log. Think why it happens, and how you can solve this problem.

Note – while we previously said that *mountdir* should reflect the content of the *rootdir*, the log is an exception. It exists in the *rootdir* but is invisible using *mountdir*. This simple example demonstrates that the files that we see are not the real data in the disk. In this way, it is possible to create hidden files (e.g. all the files started with "." in Linux).

## Errors handling in the main function

When your program starts, you first must check the received parameters in main. If the number of parameters is wrong, or the parameters values are invalid, you must print (to *stdout*) the following message:
"usage: CachingFileSystem rootdir mountdir numberOfBlocks blockSize\n".

Invalid values are:
- The *rootdir* and the *mountdir* are invalid if the directories don't exist.
- The *numberOfBlocks* and the *blockSize* are invalid if they are not positive numbers (zero is invalid too).

**If there was another error in the main function (e.g. failure in alloc), you must write to the stderr a message starting with "System Error" and exit.**

**In any error in the main function, the program should exit, and FUSE shouldn't be invoked in the first place (this means that after running FUSE's main, you shouldn't make any system call in the main function [except free if you need], because they may fail after invoking FUSE).**

## Errors in the file system functions and their return values

Each system call and library function may fail. When they fail, they update the *errno* value, which contains the last error number. *errno.h* is a library file, which defines an int value for each possible reason for an error.

When you use these functions in your file system, you must check the return value of each such call in order to understand if an error occurred. If there was an error, also your function must fail. In this case, you must return "–*errno*" (you need the minus before the *errno*, because *errno.h* defines positive values, and a minus value indicates an error). Pay attention, if you won't preserve this behavior and return the correct error, commands like "mv" in the shell will fail (you are welcome to try it yourself!).

You are required to stick with this behavior for each failure in your file system functions. For example, a failure to open the log file, to print into the log file, *malloc* failure, etc. However, there is a single failure that is not related to external function – when a function is requested to access the log file. In this case, return a *-ENOENT* error (no such file or directory).

## Assumptions

- You can assume that in case the *mountdir* exists, it is empty.
- You can also assume we'll not be supplying a cache size (block size * number of blocks) that's too big.
- The files won't be changed after running your system (remember that the requested file-system commands doesn't change the files' content).
- You may assume that we won't use "links", so each file will be defined by a single path.

## Guidelines and Tips
**This part is extremely important. Read the tips carefully, they might help you and save your precious time.**
**If you have problems of are unsure what to do, read the tips again and make sure that you follow them.**

- *CachingFileSystem*.**cpp**. We've provided you with a basic file that you should start working on. Your goal is to implement all the functions in this file.

- *fuse_main* **arguments**. This function may get a few arguments (using its first two parameters). In the *CachingFileSystem* that we provided you with, we used the "-*s*" flag, which makes FUSE run in a single thread. It's extremely recommended that you won't change it (to avoid multi-thread problems).
  Another important flag is "-*f*". It is used when running FUSE to bring the program to the foreground - this means you'll be able to see the cout (for debugging purpose), but also means you'll need a second shell to access your folder. **Don't forget to cancel the -f flag before submitting your exercise.**

- **fuse_file_info**. The *struct fuse_file_info *fi* that you receive in many functions contains several fields. However, there are only two fields that are relevant for you. The first is the "flags" which must be used in the "*open*" command. The other field is "*fh*", which you are responsible to init, delete and use in multiple situations.

- **Global variables**. In FUSE, global variables should be handled through *fuse_get_context()->private_data*  - see the fuse.h documentation and bbfs system (in the background) for more details.
  That being said, global **static** variables are valid and work as expected. While it is considered as bad design in

applications, globals are widely used in operating systems, and for our purpose they are enough and it will help you to avoid unnecessary complications.

- *mountdir* and *rootdir* **folders**. The file system for your user name is the Network File System (NFS). Therefore, the *mountdir* and *rootdir* can't be located there. Instead, **you must use directories for the *mountdir* and *rootdir* under the */tmp* directory.**
  Moreover, **don't use the same folder for mountdir and rootdir, and make sure that mountdir is empty**.

- **Running FUSE**. FUSE changes the file system, which may lead to very strange behaviors, including crashing the operating system if you do something incorrectly. When encountering a problem on the lab computers, restart the machine.
  The systems' servers (e.g. river) don't support FUSE. This is a lesson from last year, when River collapsed multiple times. We strongly suggest you to work from the Aquarium or Rotbergs' open space. It is possible to run your code at home using a Linux virtual machine. In this way, a bug won't affect your operating system, but only the virtual machine, and a restart won't be needed. However, you may then need to make some adaptations to run your code here (we had such cases last year).

- **Code Design**. As usual, make your code readable and as simple as possible, and you shouldn't have any memory leaks. In this exercise especially I would suggest to invest in a good design. Keep the logic of the LFU (and maybe the calculations of the blocks to bring) in separate files.

- **Working step by step**. Don't write all the code and then try to run it! Write a single file-system function and check it, only when it works move on to the next one.

- **Understand each single line in your code**. We supplied you with a file containing empty functions in purpose. This is meant to encourage you to start with a clean slate. We also provided you with references to two FUSE projects: the hello world from the class's presentation and the bbfs. Use them in order to <u>understand</u> what your functions need to do, and how to work with FUSE's structs.
  We strongly advise you to minimize the copy-and-paste from these examples and use them mainly for learning. In any case, under any circumstances, don't use functions that you don't understand. This will lead you to have bugs that you don't know how to solve because you don't understand your own program. If you do such things, neither we nor your friends will be able to help you.

- In contrast to other exercises that you did (in the university and in our course), here you use an interface with a low-level, complicated library. Be careful, understand what you need to do and what your code does. You will need to google to understand the Linux functions. We probably won't be able to solve each problem you encounter, because we didn't have the same problems. Therefore, we strongly encourage you to help each other and reply to your friends' posts in the Q&A forum.

- **Unmmounting your files system**. Pay attention that when you run your program, you mount the file system. Therefore, you must unmmount it (use "*fusermount -u*") before running your program again.

- After running your file system, its working directory may be changed to "/" (it happens when you don't use the "*f*" flag). If you need the working directory for some reason, you should save when the main is called (possible as a global static char*).

- **Testing your program**. We suggest you to do the following two steps in order to test your program.

  o **Basic test using shell commands.** Pay attention that the commands that you implement aren't equivalent to the common shell commands. However, the shell commands will use your functions. You can cd into directories, ls to see the files, open the files by gedit, mv files within the mount directory (which uses your rename). However, be aware than each shell function may use multiply functions of yours. For example, if "mv" doesn't work, it may be caused by due to a bug outside from "rename" function.

  However, the "ioctl" function can't be called in this way. If you want to use this function, you may use the following line in your shell:
  ```
  python -c "import os,fcntl; fd = os.open('/tmp/fuse_mount', os.O_RDONLY);
  fcntl.ioctl(fd, 0); os.close(fd)"
  ```
  where "/tmp/fuse_mount" is the mount point of your file system.

  o **More sophisticated tests** should be done via scripts. There are things that are harder (but still possible) to do via the shell. To do them, we suggest you to create testers. These testers may check automatically that the log file can't be

accessed, that you always preserve the correct information in the log file (by reading specific blocks from files), and that the rename function works as expected.

## *Theoretical part (10 points)*

Don't write more than a few lines per question.

1. In this exercise you cached files' blocks in the heap in order to enable fast access to this data.
   Does this always provides faster response then accessing the disk?
   Hint: where is the cache saved?
2. In the beginning of the exercise we described two different ways to implement LFU. Which one is better? If there is no one answer, explain under which circumstances each one of them would be preferred.
3. In the class you saw different algorithms for choosing which memory page will be moved to the disk. The most common approach is the clock-algorithm, which is LRU-like. Also our blocks-caching algorithms tries to minimize the accesses to disks by saving data in the memory. However, when we manage the buffer cache, we may actually use more sophisticated algorithms (such as LRU), which will be much harder to manage for swapping pages. Why?
   Hint – who handles accesses to memory and who accesses to file blocks?
4. Give a simple working pattern with files when LRU is better than LFU and another working pattern when LFU is better. Finally, give a working pattern when both of them don't help at all.
5. What will probably be the ideal block-size in this exercise? What happens when we use a smaller/bigger value than it?

## *Submission*

Submit a tar file on-line containing the following:

- o A README file. The README should be structured according to the course guidelines, and contains an explanation on how and why your library functions are built the way they are, as well as answers to the theoretical questions.
- o CachingFileSystem.cpp, containing your implementation of the file system, and all other relevant files you implemented.
- o Your Makefile with all the requested command