

1. Naming

- (a) Consider the following URLs:

```
http://www.google.ch/maps?q=zurich&hl=en
http://map.search.ch/Zurich
http://195.141.85.93/Zurich
```

What is the naming model of these URLs? How is the name mapping performed?

Solution:

The name-mapping algorithm for a URL works in several steps, as follows:

1. The browser extracts the part before the colon (here, `http`), considers it to be the name of a network protocol to use and resolves that name to a protocol handler using a table-lookup. The name of that context is built in to the browser. The interpretation of the rest of the URL depends on the protocol handler. The remaining steps describe the interpretation for the case of the hypertext transfer protocol (`http`) handler.
2. The browser takes the part between the `//` and the following `/` (i.e. in the first case `www.google.com`) and asks the internet Domain Name System (DNS) to resolve it. The value that DNS returns is an Internet Address. In the 3rd URL, there is no need to contact the DNS server, as the URL contains the Internet Address.
3. The browser opens a connection to the server at that Internet Address, using the protocol found in step 1 and as one of the first steps of that protocol it sends the remaining part of the URL to the server. For the 1st URL for example, it will send `maps?q=zurich&hl=en`
4. The server takes the part of the given string that is before the first question mark and looks for a file in its file system that has that path name. In most cases, this file is a script. The server uses an table-lookup to identify how to execute this script and executes it. We expect that a file would not exist in the 2nd and 3rd case, since these are actually the parameters of our search.
5. If there is a part after the question mark, as in the 1st URL, then the server treats this part as parameters for the executed script. The parameter substring (in our case `q=zurich&hl=en`) is split on the `&` symbol and every token is a different parameter.
6. Every parameter is split into an attribute and a value part. The separator character is the `=` character. Every attribute – value part is passed on the script and the script is executed. The script again uses internal context to treat every parameter in a different way. In our case, the script will map `q` to query (we are querying for Zürich) and `hl` to language (we ask for English). The mapping of given values to actual values (zurich instead of Zürich and en instead of English) is done by the context that is maintained by the script.

- (b) Now consider the following URL:

```
http://spcl.inf.ethz.ch/Research/./Research/Performance/
```

How does this URL resolve?
Is it different than the following URL?

```
http://spcl.inf.ethz.ch/Research/Performance/
```

Solution:

Using the steps described before, the browser will attempt to connect to the server that resolves to `spcl.inf.ethz.ch` and sends the two substrings to the server: `/Research/..Research/Performance/` in the first case and `/Research/Performance/` in the second.

The server is responsible for mapping these paths to a content. The two given URLs refer to the same content, but this does not always hold. Cases where these two URLs refer to different content are:

- The filesystem of the server actually supports the `..` directory name.
- The server uses a lookup table to interpret this path where i.e. `../Research/Performance/` maps to a different file.
- The server uses a filesystem where `/Research` is a soft link to another directory, i.e. to `/var/tmp`. In this case, the `/Research/..` directory does not refer to `/`, but to the `/var` directory.

2. File Systems

- (a) Is the “open” system call in UNIX absolutely essential? What would be the consequences of not having it?

Solution:

If there were no “open” system call, it would be necessary to specify the name of the file to be accessed for every read operation. The system would then have to fetch the i-node for it, although it could be cached. One issue that quickly arises is when to flush the inode back to disk. It could be based on a timeout, however it would be clumsy. Overall, it may work, but with much more overhead involved.

- (b) It has been suggested that the first part of each file be kept in the same disk block as its inode. What good would this do?

Solution: Often, files are short. If the entire file fit in the same block as the inode, only one disk access would be needed to read the file, instead of two, as is presently the case. Even for longer files there would be a gain, since one fewer disk accesses would be needed.

- (c) An Operating System only supports a single directory, but allows that directory to have arbitrarily many files with arbitrarily long file names. Can something approximating a hierarchical file system be simulated? How?

Solution: One way to simulate that is to prepend each file name with the name of directory that contains it and use a distinct character to separate different directory names. For example: `usrXstudentsXtimosXSomeFile`

- (d) Systems that support sequential files always have an operation to rewind files. Do systems that support random access files need this too?

Solution: No, random access of files does not need the “rewind” operation since if you want to read the file again, you can just access byte 0.

- (e) Contiguous allocation of files leads to disk fragmentation if files are deleted. Is this internal fragmentation or external fragmentation? What if the disk is accessed in blocks and we demand that each block contains at most data of one file?

Solution: Contiguous allocation leads to external fragmentation (due to holes between files where a file was deleted, but newly created files are bigger than that hole). If additionally disk are divided into blocks there will also be internal fragmentation (space wasted due to partially empty blocks).

- (f) One way to use contiguous allocation of disk space and not suffer from holes is to compact the disk every time a file is removed. Since all files are contiguous, copying a file requires a seek and rotational delay to read the file, followed by the transfer at full speed. Writing the file back requires the same work. Assuming a seek time of 5 msec, a rotational delay of 4 msec, a transfer rate of 8MB/sec and an average file size of 8KB, how long does it take to read a file into main memory then write it back to the disk at a new location? Using these numbers, how long would it take to compact half of a 16GB disk?

Solution: It takes 9msec to start the transfer (due to 5msec seek and 4msec rotation delay). To read 2^{13} bytes (8KB) at the transfer rate of 2^{23} bytes/sec (8MB/sec) requires 2^{-10} sec (0.977msec). Hence the total time to seek, rotate and transfer is 9.977msec. Writing back takes another 9.977msec. Thus copying an average file takes 19.954msec.

To compact half of a 16GB disk would involve copying 8GB of storage, which is 2^{20} files. At 19.954 msec per file, this takes 20,923 seconds, which is 5.8 hours. Clearly, compacting the disk after every file removal is not a great idea.

- (g) Consider an inode structure with 12 direct addresses, one indirect address, one double indirect address and one triple indirect address. Each block is 4KB in size. Assuming block addresses are 32 bit values, what is the maximum file size?

Solution: Let b be the block size (4KB), then the maximum file size is: $((b/4)^3 + (b/4)^2 + b/4 + 12) \cdot b$. For the chosen block size this gives us a maximum file size of 4 TB.

- (h) Free disk space can be kept track of using free list or a bit map. Disk addresses require D bits. For a disk with B blocks, F of which are free, state the condition under which the free list uses less space than the bit map. For D having the value 16 bits, express your answer as a percentage of the disk space that must be free.

Solution: The bit map requires B bits. The free list requires DF bits. The free list requires fewer bits if $DF < B$. Alternatively, the free list is shorter if $\frac{F}{B} < \frac{1}{D}$, where $\frac{F}{B}$ is the fraction of blocks free. For 16-bit disk addresses the free list is shorter if 6 percent or less of the disk is free.

- (i) The `open()` syscall returns a filehandle, which allows us to `read()` and `write()` to that file. In order to delete a file we use the `unlink()` syscall, which takes the pathname of the file to delete as its parameter. What happens if we create/open a file, and delete it right after, can we still use the file descriptors returned from `open()`? Write a short program to try it out. How can a user access the contents of the “deleted” file without modifying your program?

Solution: The filehandles returned by `open()` remain valid until they are `close()`'d, so we can open a file, use `unlink()` to delete it and still use `read()` and `write()` on the filehandle. This technique is used by some online video streaming players. That way they can store the video file temporarily, but a user without a solid understanding of operating/file systems can not easily download the movie by saving that file. However, we can access all open filehandles of a process through the `/proc` filesystem: in `/proc/i/fd/` we will find symlinks for each open

file handle, and we can use those symlinks to e.g., copy the “hidden” file. Also see the provided code.

- (j) Implement your own version of the `ls` utility. Of course you do not need to implement all the options `ls` provides, emulating the behaviour of `ls -l --color=never` is sufficient. Hint: start by reading the man pages for `opendir()`, `readdir()` and `fstat()`.

Solution: See the provided code.

3. I/O Systems

- (a) Why might a system use interrupt-driven I/O to manage a single serial port, but polling I/O to manage a front-end processor, such as a terminal concentrator?

Solution: Polling can be more efficient than interrupt-driven I/O. This is the case when the I/O is frequent and of short duration. Even though a single serial port will perform I/O relatively infrequently and should thus use interrupts, a collection of serial ports such as those in a terminal concentrator can produce a lot of short I/O operations, and interrupting for each one could create a heavy load on the system. A well-timed polling loop could alleviate that load without wasting many resources through looping with no I/O needed.

- (b) Polling for an I/O completion can waste a large number of CPU cycles if the processor iterates a busy-waiting loop many times before the I/O completes. But if the I/O device is ready for service, polling can be much more efficient than is catching and dispatching an interrupt. Describe a hybrid strategy that combines polling, sleeping and interrupts for I/O device service. For each of these three strategies (pure polling, pure interrupts, hybrid), describe a situation in which that strategy is more efficient than is either of the others.

Solution: A hybrid approach could switch between polling and interrupts depending on the length of the I/O operation wait. For example, we could poll and loop N times and if the device is still busy at $N+1$, we could set an interrupt and sleep. This approach would avoid long busy-waiting cycles. This method would be best for very long or very short busy times. It would be inefficient if the I/O completes at $N+T$ (where T is a small number of cycles) due to the overhead of polling plus setting up and catching interrupts. Pure polling is best with very short wait times. Pure interrupts are best with known long wait times.

- (c) How does DMA increase system concurrency? How does it complicate hardware design?

Solution: DMA increases system concurrency by allowing the CPU to perform tasks while the DMA system transfers data via the system and memory buses. Hardware design is complicated because the DMA controller must be integrated into the system and the system must allow the DMA controller to be a bus master. Cycle stealing may also be necessary to allow the CPU and DMA controller to share use of the memory bus.