



Sistemas de Operação / Fundamentos de Sistemas Operativos

(Ano letivo de 2020/2021)

Guiões das aulas práticas

Quiz #IPC/02

Processes, shared memory, and semaphores

Summary

Understanding and dealing with concurrency using shared memory.

Using semaphores to control access to a shared data structure, by different processes.

Previous note

In the code provided, system calls are not used directly. Instead, equivalent functions provided by the `process.{h,cpp}` library are used. The functions in this library deal internally with error situations, either aborting execution or throwing exceptions, thus releasing the programmer of doing so. This library will be available during the practical exams.

Question 1 *Understanding race conditions in the access to a shared data structure.*

Directory *incrementer* provides an example of a simple concurrent program used to illustrate race conditions in the access to shared data by several concurrent processes. The data shared is a single pair of integer variables, which are incremented by the different processes. Three different operations are possible on the variables: set, get and increment their values. The increment is done in such a way to promote the occurrence of race conditions in one of the variables.

(a) *Generate the unsafe version (`make incrementer_unsafe`), execute it and analyse the results.*

- *If N processes increment both variables M times each, why can the final values be different from $N \times M$?*
- *Why can the two variables have different values?*
- *Why are the final value different between executions?*
- *Macros `INC_TIME` and `OTHER_TIME` represent the times taken by the increment operation and by other work. Change their values and understand what happens. Why?*

(b) *Look at the code of the unsafe version, `inc_mod_unsafe`, and analyse it.*

- *Try to understand how shared memory is used.*
- *Try to understand why race conditions can appear.*
- *What should be done to avoid race conditions?*

(c) *Generate the safe version (`make incrementer_safe`), execute it and analyse the results.*

(d) Look at the code of the safe version, `inc_mod_safe`, and analyse it.

- Try to understand how semaphores are used to avoid race conditions, thus implementing mutual exclusion.

(e) The current implementation used the System V IPC system calls (see `man ipc` and `man svipc`). Reimplement the safe version of the given code using the POSIX versions of shared memory and semaphores resources (see `man shm_overview` and `man sem_overview`).

Question 2 Implementing producer-consumer application, using a shared FIFO and semaphores.

Directory `bounded_buffer` provides an example of a simple producer-consumer application, where interaction is accomplished through a buffer with limited capacity. The application relies on a FIFO to store the items of information generated by the producers, that can be afterwards retrieved by the consumers. Each item of information is composed of a pair of integer values, one representing the id of the producer and the other a value produced. For the purpose of easily identify race conditions, the two least significant decimal digits of every value is the id of its producer. Thus the number of producers are limited to 100.

There are 2 different implementations for the fifo: `fifo_unsafe`, and `fifo_safe`.

(a) Generate the unsafe version (`make bounded_buffer_unsafe`), execute it and analyse the results. The race conditions appear in red color.

If execution does not terminate normally and you have to force it (for instance, pressing `CRTL+c`), you may need to remove the shared memory and/or semaphore resource afterwards. To do that, use command `ipcrm -M 0x1111 -S 0x1111`, because the key used to create the resources was `0x1111`. We can use command `ipcs` to see IPC system V resources in use.

(b) Look at the code of the unsafe version, `fifo_unsafe`, and analyse it.

- Try to understand how shared memory is used.
- Try to understand why race conditions can appear.
- What should be done to solve the problem?

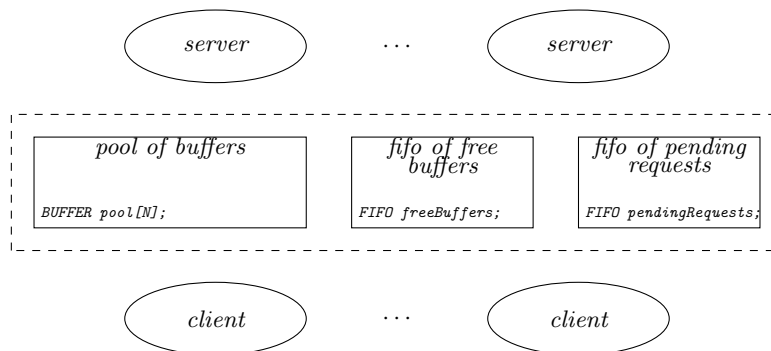
(c) Generate the safe version (`make bounded_buffer_safe`), execute it and analyse the results.

(d) Look at the code of the safe version, `fifo_safe`, and analyse it.

- Try to understand how semaphores are used to avoid both race conditions and busy waiting.
-

Question 3 *Designing and implementing a simple client-server application*

The figure below represents a simplified representation of a client-server concurrent system based on shared memory. The supporting (shared) data structure consists of a pool of N buffers of communication, individually identified by a number (between 0 and $N-1$), and two fifos, one of *ids* of buffers available and one of *ids* of buffers with pending orders. The same buffer is used for a client to place a request and the server to place the response to that request.



On the client side, interaction with the server takes place according to the following pseudo-code:

```
id = getFreeBuffer();           /* take a buffer out of fifo of free buffers */
putRequestData(data, id);       /* put request data on buffer */
addNewPendingRequest(id);       /* add buffer to fifo of pending requests */
waitForResponse(id);            /* wait (blocked) until a response is available */
resp = getResponseData(id);     /* take response out of buffer */
releaseBuffer(id);              /* buffer is free, so add it to fifo of free buffers */
```

On the server side, the interaction is described by the pseudo-code:

```
id = getPendingRequest();       /* take a buffer out of fifo of pending requests */
req = getRequestData(id);       /* take the request */
resp = produceResponse(req);     /* produce a response */
putResponseData(resp, id);      /* put response data on buffer */
signalResponseIsAvailable(id);  /* so client is waked up */
```

This is a double producer-consumer system, requiring three types of synchronization points:

- the server must block while the fifo of pending requests is empty;
- a client must block while the fifo of free buffers is empty;
- a client must block while the response to its request is not available in the buffer.

Note that in the last case there is a synchronization point per buffer. Note also that, as long as the fifos' capacities are at least the pool capacity, there is no need for a fifo full synchronization point.

Finally, consider that the purpose of the server is to process a sentence (string) to compute some statistics, specifically the number of characters, the number of digits and the number of letters.

- (a) Using the safe implementation of the fifo, used in the previous exercise, as a guideline, design and implement a solution to the data structure and its manipulation functions. Consider, for example, the following two main functions:

```
void callService(ServiceRequest & req, ServiceResponse & res);  
void processService();
```

The former is called by a client when it wants to be served; the latter is called by the server, in a cyclic way.

- (b) Implement the server process.*
 - (c) Implement the client process.*
 - (d) Does your solution work if there are more than one server?*
-