

Hello Dolly

Johan Montelius

HT2016

1 Introduction

This is an experiment where we will explore how processes are created and what is shared between them. You should have a basic understanding of what a process is but we will not assume you're an expert C programmer.

We will use the library procedure `fork()` so first take a look at the manual pages. You will probably not understand everything they talk about but we get the important information that we need to start experimenting.

```
$ man fork
:
```

We see that `fork` requires the library `unistd.h` so we need to include this in our program. We also read that `fork` will return a value of type `pid_t`. This type is defined in a header file included by `unistd.h` and is a way of making the code architecture independent. We will ignore this and assume that `pid_t` is a `int`. Further down the man pages we read that `fork` returns both the process identifier of the child process and zero. This is strange, how can a procedure return two different values? Let's give it a try, create a file called `dolly.c` and write the following:

```
#include <stdio.h>
#include <unistd.h>

int main() {
    int pid = fork();

    printf("pid = %d\n", pid);

    return 0;
}
```

The above program will call `fork` and then print the returned value. Compile and run this program, what is happening?

2 The mother and the child

So the call to `fork()` somehow creates a duplicate of the executing process and the execution then continues in both copies. By looking at the returned value we can determine if we're executing in the *mother process* or if we are in the *child process*. Try this extension to the program.

```
#include <stdio.h>
#include <unistd.h>

int main() {

    int pid = fork();

    if(pid == 0) {
        printf("I'm the child %d\n", getpid());
    } else {
        printf("My child is called %d\n", pid);
    }
    printf("That's it %d\n", getpid());
    return 0;
}
```

This is not the only way this could have been implemented. One could for example have chosen to have a construction where we would provide a function that the child process would call. Different operating systems have chosen different strategies and Windows for example have chosen to provide a procedure that creates a new process that is independent of the mother process.

2.1 wait a minute

To terminate the program in a more controlled way we can have the mother `wait()` for the child process to terminate. Try the following:

```
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/wait.h>

int main() {

    int pid = fork();

    if(pid == 0) {
        printf("I'm the child %d\n", getpid());
    }
}
```

```

    sleep(1);
} else {
    printf("My child is called %d\n", pid);
    wait(NULL);
    printf("My child has terminated \n");
}
printf("This is the end (%d)\n", getpid());

return 0;
}

```

The mother waits for its child process to terminate (actually it waits for any child it has spawned). Only then will it proceed, print out the last row and terminate.

2.2 returning a value

A process can produce a value (an integer) when it terminates and this value can be picked up by the mother process. If we change the program so that the child process returns 42 as it exists, the value can be picked up using the `wait()` procedure.

```

if(pid == 0) {
    return 42;
} else {
    int res;
    wait(&res);
    printf("the result was %d\n", WEXITSTATUS(res));
}
return 0;

```

2.3 a zombie

A *zombie* is a process that has terminated but whose parent process has not yet been informed. As long as the parent has not issued a call to `wait()` we need to keep part of the child process. When calling `wait`, the parent process should be able to pick up the exit status of the child process and possibly a return value. If the child process is completely removed from the system this information is lost.

We can see this in action if we terminate the child process but wait for a while before calling `wait()`. Do the following changes to the program, call it `zombie.c`, compile and run it in the background.

```

if(pid == 0) {
    printf("check the status\n");
    sleep(10);
}

```

```

    printf("and again\n");
    return 42;
} else {
    sleep(20);
    int res;
    wait(&res);
    printf("the result was %d\n", WEXITSTATUS(res));
    printf("and again\n");
    sleep(10);
}
return 0;
}

```

Check the status of the processes using the `ps` command. Notice how the two processes are created, how the child becomes a zombie and is then removed from the system once we have received the return value.

```

$ gcc -o zombie zombie.c
$ ./zombie&
:
$ ps -ao pid,stat,command
:

```

2.4 a clone of the process

So we have created a child process that is a *clone* of the mother process. The child is a copy of the mother with an identical memory. We can exemplify this by showing that the child has access to the same data structures but that the structures are obviously just copies of the original data structures. Extend `dolly.c` and try the following:

```

int main() {

    int pid;
    int x = 123;
    pid = fork();

    if(pid == 0) {
        printf(" child:  x is %d\n", x);
        x = 42;
        sleep(1);
        printf(" child:  x is %d\n", x);
    } else {
        printf(" mother: x is %d\n", x);
        x = 13;
    }
}

```

```

    sleep(1);
    printf("  mother: x is %d\n", x);
    wait(NULL);
}
return 0;
}

```

As you see both the mother and the child sees ha variable `x` as 123 but the changes made are only visible by themselves. If you want to see something very strange you can change the printout to also print the memory address of the variable `x`. Do this for both the mother and the child and you will see that they are actually referring to the same memory locations.

```

printf("  child:  x is %d and the address is 0x%p\n", x, &x);

```

The explanation is that processes use *virtual addresses* and they are identical, they are however mapped to different *real memory addresses*. How this is achieved is nothing that we should explore now but its fun to see that it is working.

2.5 what we do share

Since the child process is a clone of the mother process we do actually share some parts. On thing that we do share are references to open files. When a process opens a file a *file table entry* is created by the operating system. The process is given a reference to this entry and this is reference is stored in a *file descriptor table* that is owned by the process. Now when the process is cloned, this table is copied and all the references are of course pointing to the same entries in the *file table*.

The standard output is of course nothing more than a entry in the file descriptor table so this is why both processes can write to the standard output. We also read from the same standard input so we have a race condition also there.

If you look at man pages you will see a whole range of structures that the processes share or not share but most of those are not very interesting to us in this set of experiments.

3 Groups, orphans, sessions and daemons

The mother, or should we call it parent process to be gender neutral, has a special relationship to the child process. The parent process has to keep track of its child's and a child always knows the process identifier of its parent.

```

int main () {
    int pid = fork ();

```

```

if(pid == 0) {
    printf("I'm the child %d with parent %d\n", getpid(), getppid());
} else {
    printf("I'm the parent %d with parent %d\n", getpid(), getppid());
    wait(NULL);
}
return 0;
}

```

Compile and run this in a terminal, who is the parent of the parent process? The following commands might give you a clue.

```

$ ps a
:
:
$ echo $$
:

```

We could find more information about the processes using some flags to `ps`. Try `ps -fp $$` to see more information about the shell your using (`$$` will expand to the process identifier of the shell). The `PPID` field is the parent process identifier. Who is the parent of the shell? Where does it all stop?

3.1 the group

The fate of a parent and its child are not directly linked to each other but they belong to the same *process group*. Each process group has a process leader and in our simple examples the parent process has been the leader of the group. The group identifier/leader is retrieved using the system call `getpgid()`.

```

int main() {

    int pid = fork();

    if(pid == 0) {
        int child = getpid();
        printf("I'm the child %d in group %d\n", child, getpgid(child));
    } else {
        int parent = getpid();
        printf("I'm the parent %d in group %d\n", parent, getpgid(parent));
        wait(NULL);
    }
    return 0;
}

```

A group is treated as a unit by the shell, it can set a whole group to suspend, resume or run in the background (allowing the shell to use the standard input for interaction). We will however not go into how the shell is working so let's just accept that processes belong to a process group.

3.2 orphans

As a change we can try to crash the parent process and see what happens with the child process.

```
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/wait.h>

int main() {

    int pid = fork();

    if(pid == 0) {
        int child = getpid();
        printf("child: parent %d, group %d\n", getppid(), getpgid(child));
        sleep(4);
        printf("child: parent %d, group %d\n", getppid(), getpgid(child));
        sleep(4);
        printf("child: parent %d, group %d\n", getppid(), getpgid(child));
    } else {
        int parent = getpid();
        printf("parent: parent %d, group %d\n", getppid(), getpgid(parent));
        sleep(2);
        int zero = 0;
        int i = 3 / zero;
    }
    return 0;
}
```

Save the program in a file called `orphan.c`. Compile and execute the program, notice how the parent identifier of the child process changes. The process has turned into an *orphan* and adopted by the `upstart` process (or `init` or `systemd` depending on which system you using). Note the new process identifier and then check its state using the `ps` command:

```
$ps <whatever the process id was>
:
:
```

To see something fun you can take a look at the **process tree** of the process:

```
$pstree <whatever the process id was>
:
```

3.3 sessions and daemons

The origin of the notion of a session is a user attaching and logging in to the system. A session consists of a set of groups and a *session leader*. As with groups, the sessions have identifiers that are equal to the leaders process identifier. Compile and run the program below, which process is the session leader of our processes?

```
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/wait.h>

int main() {

    int pid = fork();

    if(pid == 0) {
        int child = getpid();
        printf("child: session %d\n", getsid(child));
    } else {
        int parent = getpid();
        printf("parent: session %d\n", getsid(parent));
    }
    return 0;
}
```

When you start a new terminal, a new session is created. The operating system keeps track of sessions and will terminate all groups in a session if the session leader terminates. This means that if you log in to a system and start to run processes in the background they still belong to the same session as your login shell and will be terminated if the session terminates.

If one wants to create a process that should survive the session it must form its own session. It becomes a *daemon*, a process that is running in the background detached from any *controlling terminal*.

Many of the tasks performed by the operating system are performed by daemons. They keep track of network interfaces, USB devices or schedules tasks that should run periodically etc. Your system will probably have fifty daemons running in the background but they consume very little resources.

4 Starting a program

So far we have seen how a process can be created and how the child process is related to its parent process. To understand how an operating system works there is one more very important functionality that we will take a look at - how we create a process that will execute another program.

When you use the command shell this happens (almost) every time you enter a command. Some commands are interpreted by the shell and the shell will do something for us but most “commands” are actually programs that the shell will start for us. How is a program actually started?

4.1 transforming a process

In Unix systems the execution of a program is done by transforming an existing process to run the code of the given program. As you will see, starting a program is done in two steps - creating the new process and then transforming the process into executing the program.

The mechanism that makes this possible is the family of `exec()` system calls. Look-up the man pages of `exec`, we will use the one called `execlp()`.

```
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/wait.h>

int main() {

    int pid = fork();

    if(pid == 0) {
        execlp("ls", "ls", NULL);
        printf("this will only happen if exec fails\n");
    } else {
        wait(NULL);
        printf("we're done\n");
    }
    return 0;
}
```

The call to `execlp()` will find the program `ls` and then replace the code and data segments of the process with the code and data found in the executable binary. The stack and heap areas are reset so the program starts the execution from scratch.

4.2 redirection

Even if the memory segments of the process is cleared, the process keeps the file descriptor table. By changing the table entries we can make the program read from a standard input of our choice and we can of course redirect the standard output. This allows us to control the I/O operations of the program without changing the program in any way.

To see how this works we can create a small program that does nothing but writes to standard output. Let's call this program `boba.c`.

```
#include <stdio.h>

int main() {
    printf("Don't get in my way.\n");
    return 0;
}
```

Now if we compile and run this program we will of course see the quote printed in the terminal.

```
$gcc -o boba boba.c
:
$./boba
:
```

Note that we have to write `./boba` and not simply `boba` if you have not set up your `PATH` variable to also include the current directory; more on this later.

Now if we want to redirect the output to a file called `quotes.txt` we could of course do this from the shell directly.

```
$/boba > quotes.txt
:
```

To understand how the shell achieves this we could try to write a program `jango.c`, that clones itself, redirects the standard output and then transforms the clone into `boba`. Let's go:

```
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <sys/wait.h>

int main() {
```

```

int pid = fork ();

if(pid == 0) {
    int fd = open("quotes.txt", O_RDWR | O_CREAT, S_IRUSR | S_IWUSR);
    dup2(fd, 1);
    close(fd);
    execl("boba", "boba", NULL);
    printf("this only happens on failure\n");
} else {
    wait(NULL);
}
return 0;
}

```

In the `jango.c` program we open a file `quotes.txt` (providing flags to open it in read-write mode and create it if does not exist). The operating system will create a new file table entry and add a reference to it in our file descriptor table. The table entry will be the first free entry in the table (3). We then use the system call `dup2()` to copy the entry to position 1 (the location of `stdout`). We then close the `fd` entry since we will not use this entry any more.

When we now call `execl()`, the process will turn into `boba`. The `boba` program knows nothing about what has happened but will of course direct all output to file descriptor 1 as usual. Try it and you will see that the file is created and that we will receive the output as expected.

4.3 pipes

The full beauty of how standard input and output can be redirected is shown when we introduce the concept of *pipes*. A pipe is a FIFO buffer of characters and when created we will allocate two file descriptor entries. One entry is for reading and the other for writing.

Since we are in full control over the descriptor table before we start executing a program, we can make one program send all the output to another program's input. From the shell this is a very powerful tool to combine sequences of commands.

Assume we have the commands (or rather programs) `ps axo sid` that will print the session identifier of every process in the system, `sort -u` that will sort lines and output only unique and `wc -l` that will count the number of lines. How do we combine these to find the number of sessions in the system. Using the shell this is done in one line:

```

$ ps xao sid | sort -u | wc -l
:

```

This is achieved using pipes and we can set it up ourselves in a program. There are however a lot of details to get it right and we will explore this later in the course. For now you should explore using the pipes from the command line.

5 Summary

Processes are created by cloning an existing process, the execution continues in the two duplicates and the only way of telling in which copy we are executing in, is to look at the value returned from `fork()`.

A parent and child process are in the same *process group*. If the group leader terminates all processes in the group will be sent a signal that will likely cause them to terminate.

Several groups belong to a *session* with a *controlling terminal*. If the session leader terminates or the controlling terminal closes, the whole session will be terminated.

A session that has been detached from any controlling terminal is called a *daemon*. Daemons handle many of the tasks that constitute an operating system.

Two copies have identical copies of *file descriptor tables* referring to the same *file table entries*. By changing the descriptor tables, the input and output of a process can be redirected. Two processes can use this to set up a *pipe* between them that acts as a buffered FIFO channel.

A process can be transformed to run another program using the system call `exec()`. This will reset all memory segments but the transformed process keeps the file descriptor table.