# Sistemas Operativos 72.11

Entrando en calor

# Entorno de desarrollo

# Entorno de desarrollo

## more bash tricks

Julia Evans
@b0rk

### cd -

changes to the directory you were last in

pushd & popd let you keep a stack

### ctrl+z

suspends (SIGSTOP) the running program

### fg

brings backgrounded or suspended program to the foreground

### ♡ shellcheck ♡

an amazing shell script linter! it can help you write safe & correct scripts

### <( )

process substitution

treat process output like a file (no more temp files!)
eg:

$ diff <(ls) <(ls -a)

### fc

"fix command"

open the last command you ran in an editor

then run the edited version

### type

tells you if something is a builtin, program, or alias

try running type on

{time} [ping] (pushd)

(they're all different types)

# Entorno de desarrollo
## Docker

- ¿Por qué?
- Se solicitan dadores de conocimiento, factor "Docker en Windows" positivo
- Demo https://github.com/alejoaquili/ITBA-72.11-SO.git

# Manuales

## man

- ¿Por qué?
- Demo https://github.com/alejoaquili/ITBA-72.11-SO.git

# Manuales

man pages = awesome

JULIA EVANS
@b0rk

## man pages are split up into 8 sections

① ② ③ ④ ⑤ ⑥ ⑦ ⑧

$ man 2 read

means "get me the man page for read from section 2"

There's both
→ a program called "read"
→ and a system call called "read"

so

$ man 1 read

gives you a different man page from

$ man 2 read

If you don't specify a section, man will look through all the sections & show the first one it finds

## man page sections

① programs
$ man grep
$ man ls

② system calls
$ man sendfile
$ man ptrace

③ C functions
$ man printf
$ man fopen

④ devices
$ man null
    for /dev/null docs

⑤ file formats
$ man sudoers
    for /etc/sudoers
$ man proc
    files in /proc !

⑥ games
not super useful.
$ man sl
    is funny if you have
    sl though.

⑦ miscellaneous
explains concepts!
$ man 7 pipe
$ man 7 symlink
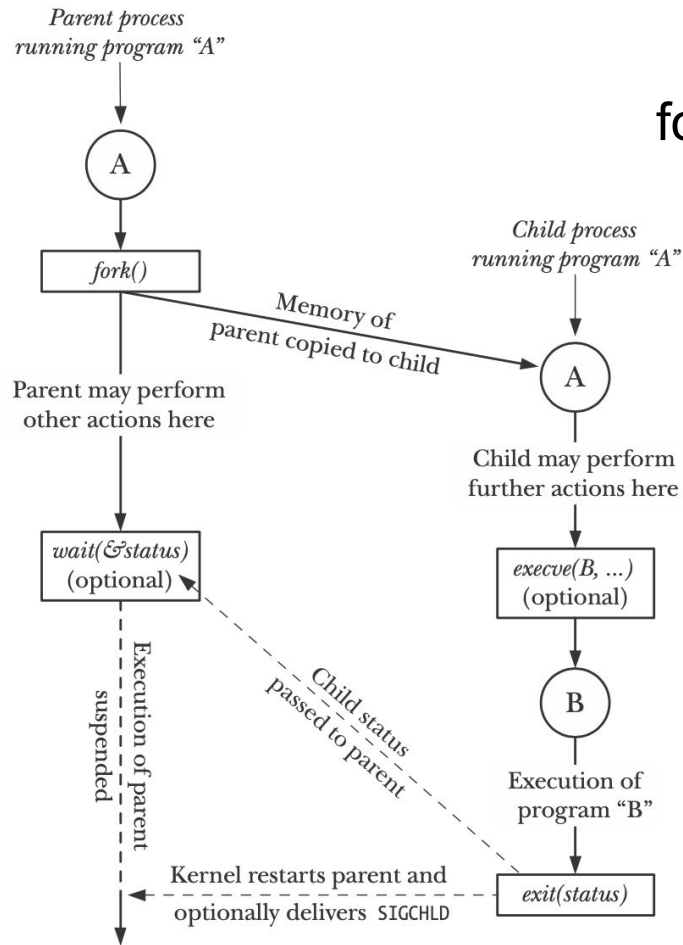
⑧ sysadmin programs
$ man apt
$ man chroot

ITBA
Instituto Tecnológico
de Buenos Aires

# POSIX
## fork execve wait exit



*Parent process running program "A"*

A

fork()

Memory of parent copied to child

*Child process running program "A"*

A

Parent may perform other actions here

Child may perform further actions here

wait(&status) (optional)

execve(B, ...) (optional)

Execution of parent suspended

Child status passed to parent

B

Execution of program "B"

Kernel restarts parent and optionally delivers SIGCHLD

exit(status)

**Figure 24-1:** Overview of the use of *fork()*, *exit()*, *wait()*, and *execve()*

The Linux programming interface - Kerrisk, Michael

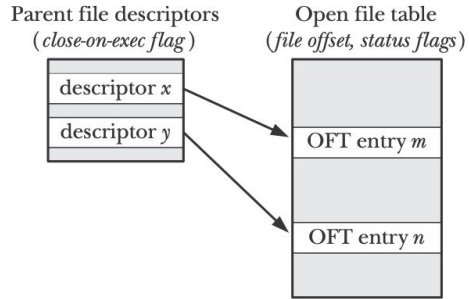### Process management

| Call | Description |
|------|-------------|
| pid = fork( ) | Create a child process identical to the parent |
| pid = waitpid(pid, &statloc, options) | Wait for a child to terminate |
| s = execve(name, argv, environp) | Replace a process' core image |
| exit(status) | Terminate process execution and return status |

# POSIX
## fork execve wait exit



a) **Descriptors and open file table entries before *fork()***

Parent file descriptors ( *close-on-exec flag* )

Open file table ( *file offset, status flags* )

descriptor x
descriptor y

OFT entry m

OFT entry n

b) **Descriptors after *fork()***

Parent file descriptors

Open file table

descriptor x
descriptor y

*Descriptors duplicated in child*

Child file descriptors

descriptor x
descriptor y

OFT entry m

OFT entry n

c) **After closing unused descriptors in parent (y) and child (x)**

Parent file descriptors

Open file table

descriptor x
descriptor y

Child file descriptors

descriptor x
descriptor y
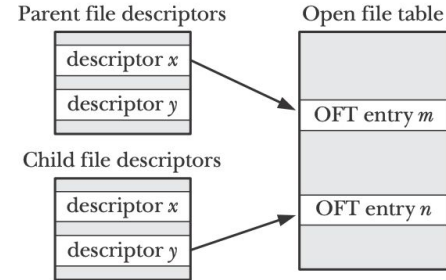
OFT entry m

OFT entry n

**Figure 24-2:** Duplication of file descriptors during *fork()*, and closing of unused descriptors

# POSIX
## Pair programming

- Programar una shell funcional a partir de este esqueleto

```
#define TRUE 1

while (TRUE) {                                  /* repeat forever */
    type_prompt( );                             /* display prompt on the screen */
    read_command(command, parameters);          /* read input from terminal */

    if (fork( ) != 0) {                         /* fork off child process */
        /* Parent code. */
        waitpid(-1, &status, 0);                /* wait for child to exit */
    } else {
        /* Child code. */
        execve(command, parameters, 0);         /* execute command */
    }
}
```

# Comandos UNIX

## ps

JULIA EVANS
@b0rk

### ps

ps shows which processes are running

I usually run ps like this:

$ ps aux

u means include username column

a+x together show all process

(ps -ef works too)

### w

is for _wide_. ps auxwww will show all the command line args for each process

### e

is for _environment_. ps auxe will show the environment vars!

### wchan

you can choose which columns to show with ps (ps -eo ...)

One cool column is 'wchan' which tells you the name of the kernel function if the process is sleeping

try it:

ps -eo user, pid, wchan, cmd

### ★ process state ★

Here's what the letters in ps's STATE column mean:

R: running
S/D: asleep
Z: zombie
l: multithreaded
+: in the foreground

### f

is for "forest" ☺. ps auxf will show you an ASCII art process tree !

pstree can display a process tree too

ps has 3 different sets of command line arguments 💔

1. UNIX ( 1 dash)
2. BSD (no dash)
3. GNU ( 2 dashes)

you can write monstrosities like:
$ ps f -f

forest (BSD)

full format (UNIX)

ITBA
Instituto Tecnológico
de Buenos Aires

# copy on write

Julia Evans
@b0rk

**On Linux, you start new processes using the fork() or clone() system call**

calling fork gives you a child process that's a copy of you

parent    child

---

the cloned process has EXACTLY the same memory

→ same heap

→ same stack

→ same memory maps

if the parent has 3GB of memory, the child will too

---

copying all that memory every time we fork would be slow and a waste of RAM

often processes call exec right after fork which means they don't use the parent process's memory basically at all!

---

so Linux lets them share physical RAM and only copies the memory when one of them tries to write.

process — I'd like to change that memory

ok I'll make you your own copy! — Linux

---

Linux does this by giving both the processes identical page tables

same RAM

but marks every page as read only

---

when a process tries to write to a shared memory address

① there's a page fault

② Linux makes a copy of the page & updates the page table

③ the process continues, blissfully ignorant

It's just like I have my own copy

# Copy on Write



**Figure 24-3:** Page tables before and after modification of a shared copy-on-write page

The Linux programming interface - Kerrisk, Michael

# POSIX
## open read write close

- Universalidad de I/O en UNIX
- Eficiencia read write
- Protocolo
  - End of file
  - Read 0
  - Read block

**File management**

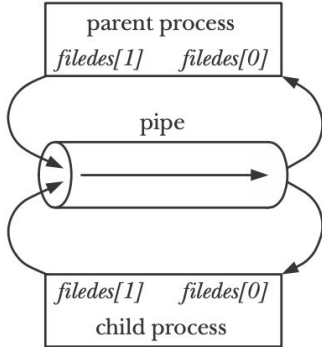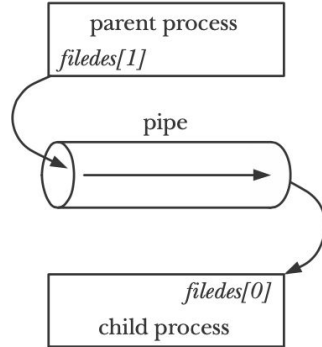| Call | Description |
|---|---|
| fd = open(file, how, ...) | Open a file for reading, writing, or both |
| s = close(fd) | Close an open file |
| n = read(fd, buffer, nbytes) | Read data from a file into a buffer |
| n = write(fd, buffer, nbytes) | Write data from a buffer into a file |
| position = lseek(fd, offset, whence) | Move the file pointer |
| s = stat(name, &buf) | Get a file's status information |

# POSIX

# POSIX
## pipe



**Figure 44-2:** Process file descriptors after creating a pipe



a) After *fork()*

b) After closing unused descriptors

**Figure 44-3:** Setting up a pipe to transfer data from a parent to a child
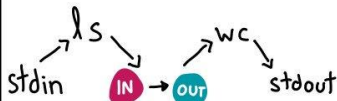
# POSIX
## pipes

# POSIX
## dup

- DEMO shell, cat y strace

| shell | |
|---|---|
| fd | device |
| 0 | /dev/tty |
| 1 | /dev/tty |
| 2 | /dev/tty |

| shell | |
|---|---|
| fd | device |
| 0 | /dev/tty |
| 1 | /dev/tty |
| 2 | /dev/tty |
| **3** | **r-end** |
| **4** | **w-end** |

| shell | | child1 | | child2 | |
|---|---|---|---|---|---|
| fd | device | **fd** | **device** | **fd** | **device** |
| 0 | /dev/tty | **0** | **/dev/tty** | **0** | **/dev/tty** |
| 1 | /dev/tty | **1** | **/dev/tty** | **1** | **/dev/tty** |
| 2 | /dev/tty | **2** | **/dev/tty** | **2** | **/dev/tty** |
| 3 | r-end | **3** | **r-end** | **3** | **r-end** |
| 4 | w-end | **4** | **w-end** | **4** | **w-end** |

shell → pipe → shell → fork x2 →

w-end  r-end

shell

w-end  r-end

child1  child2

# POSIX
## dup

Left side:

| shell | | child1 | | child2 | |
|---|---|---|---|---|---|
| **fd** | **device** | **fd** | **device** | **fd** | **device** |
| 0 | /dev/tty | 0 | /dev/tty | 0 | /dev/tty |
| 1 | /dev/tty | 1 | /dev/tty | 1 | /dev/tty |
| 2 | /dev/tty | 2 | /dev/tty | 2 | /dev/tty |
| | | **3** | **r-end** | 3 | r-end |
| | | 4 | w-end | **4** | **w-end** |

shell

w-end → r-end

child1   child2

child1

close(3)

child2

close(4)

→

Right side:

| shell | | child1 | | child2 | |
|---|---|---|---|---|---|
| **fd** | **device** | **fd** | **device** | **fd** | **device** |
| 0 | /dev/tty | 0 | /dev/tty | 0 | /dev/tty |
| 1 | /dev/tty | 1 | /dev/tty | 1 | /dev/tty |
| 2 | /dev/tty | 2 | /dev/tty | 2 | /dev/tty |
| | | 4 | w-end | 3 | r-end |

shell

w-end → r-end

child1   child2

# POSIX
## dup

| shell | |
|---|---|
| **fd** | **device** |
| 0 | /dev/tty |
| 1 | /dev/tty |
| 2 | /dev/tty |

| child1 | |
|---|---|
| **fd** | **device** |
| 0 | /dev/tty |
| **1** | **/dev/tty** |
| 2 | /dev/tty |
| 4 | w-end |

| child2 | |
|---|---|
| **fd** | **device** |
| **0** | **/dev/tty** |
| 1 | /dev/tty |
| 2 | /dev/tty |
| 3 | r-end |

| shell | |
|---|---|
| **fd** | **device** |
| 0 | /dev/tty |
| 1 | /dev/tty |
| 2 | /dev/tty |

| child1 | |
|---|---|
| **fd** | **device** |
| 0 | /dev/tty |
| 2 | /dev/tty |
| 4 | w-end |

| child2 | |
|---|---|
| **fd** | **device** |
| 1 | /dev/tty |
| 2 | /dev/tty |
| 3 | r-end |

child1

close(1)

child2

close(0)

# POSIX
## dup

| shell | | child1 | | child2 | |
|---|---|---|---|---|---|
| **fd** | **device** | **fd** | **device** | **fd** | **device** |
| 0 | /dev/tty | 0 | /dev/tty | 1 | /dev/tty |
| 1 | /dev/tty | 2 | /dev/tty | 2 | /dev/tty |
| 2 | /dev/tty | 4 | w-end | 3 | r-end |

child1

dup(4)

child2

dup(3)

| shell | | child1 | | child2 | |
|---|---|---|---|---|---|
| **fd** | **device** | **fd** | **device** | **fd** | **device** |
| 0 | /dev/tty | 0 | /dev/tty | **0** | **r-end** |
| 1 | /dev/tty | **1** | **w-end** | 1 | /dev/tty |
| 2 | /dev/tty | 2 | /dev/tty | 2 | /dev/tty |
| | | 4 | w-end | 3 | r-end |

# POSIX
## dup

| shell | |
|---|---|
| fd | device |
| 0 | /dev/tty |
| 1 | /dev/tty |
| 2 | /dev/tty |

| child1 | |
|---|---|
| fd | device |
| 0 | /dev/tty |
| 1 | w-end |
| 2 | /dev/tty |
| **4** | **w-end** |

| child2 | |
|---|---|
| fd | device |
| 0 | r-end |
| 1 | /dev/tty |
| 2 | /dev/tty |
| **3** | **r-end** |

| shell | |
|---|---|
| fd | device |
| 0 | /dev/tty |
| 1 | /dev/tty |
| 2 | /dev/tty |

| child1 | |
|---|---|
| fd | device |
| 0 | /dev/tty |
| 1 | w-end |
| 2 | /dev/tty |

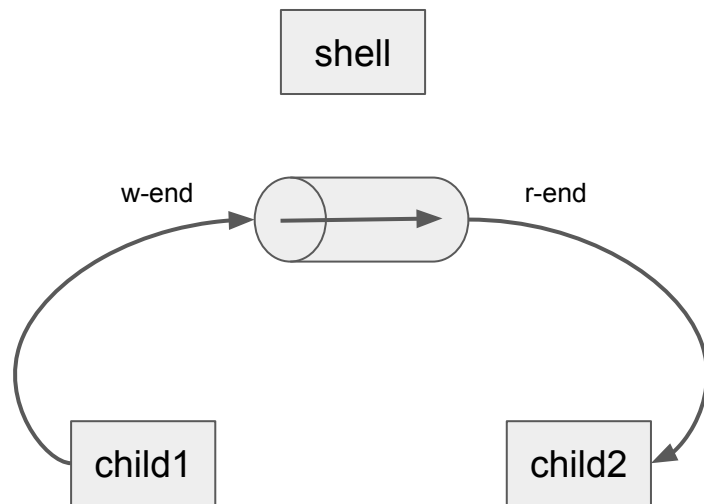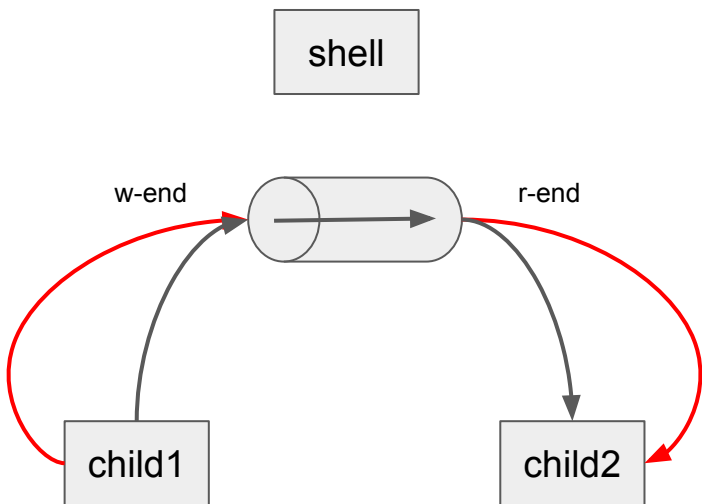| child2 | |
|---|---|
| fd | device |
| 0 | r-end |
| 1 | /dev/tty |
| 2 | /dev/tty |

child1

close(4)

child2

close(3)

# POSIX
## dup

```
//...
 case PIPE:
   pcmd = (struct pipecmd*)cmd;
   if(pipe(p) < 0)
     panic("pipe");
   if(fork1() == 0){
     close(1);
     dup(p[1]);
     close(p[0]);
     close(p[1]);
     runcmd(pcmd->left);
   }
   if(fork1() == 0){
     close(0);
     dup(p[0]);
     close(p[0]);
     close(p[1]);
     runcmd(pcmd->right);
   }
   close(p[0]);
   close(p[1]);
   wait();
   wait();
   break;
//...
```

# Ejercicio 1

Escribir un programa en C que liste recursivamente el contenido de un directorio que recibe como primer y único argumento, indicar si cada elemento es un archivo o un directorio, tabular la salida por nivel de anidamiento (**stat opendir readdir**)

Ejemplo de la salida (solo a modo ilustrativo):

```
$./tree .
d     dir1
f     file1
d     dir3
f          file1
d          dir1
f               file1
f     file2
```

# Ejercicio 2

1. Escribir un programa en C que lance 10 procesos que realicen alguna tarea que dure una cierta cantidad de tiempo considerable (por ejemplo, dormir una cantidad al azar de tiempo), el programa debe esperar a que sus procesos hijos terminen para terminar él mismo. (**fork waitpid**).

2. Hacer que cada hijo imprima su propio pid (**getpid**)

3. Separe el código de los hijos y del padre en dos unidades de compilación diferentes (**execve**)

# Ejercicio 3

Lance varios procesos hijos desde un padre, haga que el padre cuelgue en un loop infinito sin hacer `waitpid`, mientras que sus hijos terminan instantáneamente, salvo 1 que también quedará en un loop infinito.

Para observar el estado de los procesos se recomienda ejecutar

```
$ ps ax -O "%P"
```

1. ¿Qué imagina que va a pasar con los procesos hijos que terminaron? ¿Por qué cree que pasa?
2. Mate el proceso padre (`kill`) ¿Qué sucede?
3. ¿Cambió algo para el proceso hijo que estaba en un loop luego de matar al padre?
4. ¿Qué es el proceso `init`?
5. ¿Qué pasa en los sistemas `POSIX` con los procesos huérfanos?

# Ejercicio 4

Dados los programas dentro del repositorio ejemplos/producer-consumer y la resolución del inciso 3 del ejercicio 2:

1. Modifique la resolución del inciso 3 del ejercicio 2, de manera tal que el proceso padre ejecute *p* y *c* (**fork** y **exec**) conectando el stdout de p al stdin de c utilizando un pipe (**pipe**).
2. Modifique la resolución del inciso 3 del ejercicio 2, de manera tal que el proceso padre ejecute *p* en background. El programa *c* no participa en este ejercicio.
3. Modifique la resolución del inciso 3 del ejercicio 2, de manera tal que el proceso padre ejecute *p* y redireccione stdout a un archivo. (**close** y **open**). El programa *c* no participa en este ejercicio.
4. Modifique la resolución del inciso 3 del ejercicio 2, de manera tal que el proceso padre ejecute *c*, pero que en lugar de leer de stdin lea de un archivo. El programa *p* no participa en este ejercicio.

# Ejercicio 5

1.  Modifique la resolución del inciso 3 del ejercicio 2, de manera tal que el proceso padre ejecute 2 programas de forma secuencial considerando 3 escenarios posibles:
    a.  INCONDICIONAL, primero ejecuta un proceso y cuando termina ejecuta el otro.
    b.  AND, si el primer proceso falla (exit status != 0), NO se ejecuta el segundo
    c.  OR, si el primer proceso termina satisfactoriamente (exit status == 0), NO se ejecuta el segundo

    Puede utilizar los programas *true.c* y *false.c* para hacer pruebas. (hint: man **true**, man **false**)
    (**wait**)

# Ejercicio 6

1. Reflexión:
    a. ¿Fue necesario modificar p.c o c.c para resolver alguno de estos ejercicios?
    b. ¿Qué beneficios tiene que las syscalls **fork** y **execve** estén separadas y que **execve** preserve los archivos abiertos por el invocador?
    c. ¿Cómo cree que un intérprete de comandos como sh, bash, zsh, etc resuelve los siguientes comandos?
          i. `./p | ./c`
         ii. `./p &`
        iii. `./p > out`
         iv. `./c < in`
          v. `./true ; ./false`
         vi. `./false && ./true`
        vii. `./true || ./false`
       viii. `<( ) (process substitution)`
    d. Investigue cómo se resuelven estos problemas en windows.

# Glosario

- Pair programming
- Copy on write
- File descriptor
- fork bomb