

# Systèmes 1 — Processus

December 30, 2012

## Contents

<b>1</b>	<b>Le concept de processus</b>	<b>2</b>
<b>2</b>	<b>Attributs d'un processus</b>	<b>2</b>
2.1	L'identité . . . . .	2
2.2	Propriétaire . . . . .	3
2.3	Répertoire de travail . . . . .	4
2.4	Session et groupe de processus . . . . .	4
2.4.1	Terminal de contrôle . . . . .	4
2.5	Date de création . . . . .	4
2.6	Temps CPU . . . . .	5
2.6.1	Le masque de création de fichiers . . . . .	5
2.7	Table de descripteurs . . . . .	5
<b>3</b>	<b>Espace d'adressage de processus</b>	<b>6</b>
<b>4</b>	<b>Création de processus</b>	<b>6</b>
<b>5</b>	<b>Primitives de recouvrement</b>	<b>8</b>
5.1	Attributs d'un processus après le recouvrement . . . . .	9
<b>6</b>	<b>Terminaison de processus</b>	<b>9</b>
<b>7</b>	<b>Processus zombie et synchronisation père fils</b>	<b>10</b>
7.1	wait . . . . .	10
7.2	waitpid . . . . .	10
<b>8</b>	<b>Job control</b>	<b>12</b>
<b>9</b>	<b>Duplication de descripteur avec dup() et dup2()</b>	<b>12</b>
<b>10</b>	<b>Communication par les tubes</b>	<b>15</b>
10.1	Tubes anonymes . . . . .	15
10.2	Création d'un tube anonyme . . . . .	16
10.3	Lecture d'un tube . . . . .	16
10.4	Écriture dans un tube . . . . .	17
10.5	Exemple - communication père vers fils par un tube anonyme . . . . .	17
10.6	Exemple - redirections . . . . .	18

<b>11 Tubes nommés</b>	<b>21</b>
11.1 Création et ouverture . . . . .	21
11.2 Lecture et écriture . . . . .	22
<b>12 Changement de mode : bloquant ↔ non bloquant pour les tubes nommés ou anonymes</b>	<b>22</b>
<b>13 setuid et setgid</b>	<b>22</b>

## 1 Le concept de processus

Un processus – l'exécution d'un programme binaire par la machine.

L'espace d'adressage: instructions, données (dont la pile d'exécution et le tas). Cette espace est dynamique et change avec le temps.

L'exécution de deux modes différents:

- **mode utilisateur** (user mode) – le processus exécute ses propres instructions et accède les données de son espace d'adressage. La tentative d'accéder à des données hors de son espace d'adressage provoque un échec et envoie d'un signal.
- Dans le **mode système, noyau, superviseur** (supervisor, system, kernel mode) le processus exécute des instruction n'appartenant pas au programme mais au noyau, cela permet d'accéder à l'ensemble de données de système (tables de système). Changement de mode via l'appel de système par une trappe (trap) ou par un événement (interruption). Dans le dernier cas le processus exécute le gérant de l'interruption (handler) correspondant et reprend au retour de handler.

L'ordonnancement (scheduling) de processus. L'utilisation de priorités. Les caractéristiques de processus dans le bloc de contrôle (BCP).

BCP composé de deux parties:

1. une partie conservée dans le noyau. Contient les information utiles pour le noyau même si le processus inactif: identité, les infos disponible uniquement en mode noyau,
2. les informations dans l'espace d'adressage de processus: descripteurs de fichiers, handlers de signaux.

## 2 Attributs d'un processus

Chaque processus possède un certain nombre d'attributs. Cette section donne la liste de fonctions qui permettent d'obtenir les différents attributs.

### 2.1 L'identité

L'identité de processus, **pid**, est un entier positif unique qui permet d'identifier le processus.

```
#include <sys/types.h>
#include <unistd.h>
pid_t getpid(void); //l'identite du processus
pid_t getppid(void); //l'identite du processus pere
```

`pid_t` c'est un type entier signé. `getpid()` donne le pid de processus qui ne change pas de la création jusqu'à la mort du processus. `getppid()` donne l'identité du processus père. Cette identité peut changer si le père meure.

## 2.2 Propriétaire

Chaque processus possède un propriétaire (identifié par une valeur de type `uid_t`, c'est une valeur entière) et un groupe de propriétaires (du type `gid_t`, aussi un type entier).

En fait il y a même deux propriétaires et deux groupes propriétaires, propriétaire réel et le propriétaire effectif et le groupe propriétaire réel et effectif.

On reviendra plus tard sur les différents propriétaires. Pour l'instant juste les fonctions qui permettent de connaître le propriétaire:

### Propriétaire réel

```
#include <sys/types.h>
#include <unistd.h>
uid_t getuid(void); /*retourne l'identite du proprietaire reel*/
```

Le propriétaire réel c'est celui qui a lancé le processus. Si vous lancez un processus depuis le terminale c'est le propriétaire de processus l'interprète de commande qui contrôle le terminal.

### Propriétaire effectif

```
uid_t geteuid(void); /* retourne id du proprietaire effectif*/
```

C'est le propriétaire effectif qui détermine les droits d'accès aux ressources (fichiers et processus).

```
int setuid(uid_t uid)
```

permet de changer le propriétaire réel et effectif. L'utilisateur non privilégié ne peut modifier que le propriétaire effectif en lui donnant la valeur du propriétaire réel.

### Groupe propriétaire réel

```
gid_t getgid(void)
```

### Groupe propriétaire effectif

```
gid_t getegid(void)
```

```
int setgid(gid_t)
```

Permet de changer le groupe propriétaire réel et effectif. L'utilisateur non privilégié peut modifier uniquement le groupe propriétaire effectif en lui attribuant la valeur de groupe propriétaire réel.

## 2.3 Répertoire de travail

Chaque processus possède le répertoire de travail utilisé pour résoudre les références relatives.

Le répertoire de travail est référencé par `.` (point).

On peut modifier le répertoire de travail avec :

```
#include <unistd.h>
int chdir(const char *reference)
```

où la référence donne le chemin vers le nouveau répertoire de travail. `chdir()` retourne 0 en cas de succès et `-1` en cas d'échec.

Pour obtenir la référence absolue (le chemin absolu) de répertoire de travail on utilise la fonction :

```
#include <unistd.h>
char *getcwd(char *tampon, size_t taille)
```

Si l'appel réussi le tampon contient le chemin absolu vers le répertoire de travail. `tampon` doit contenir une adresse valide du tampon où `getcwd()` doit mettre le résultat, `taille` donne la taille du tampon.

La valeur retournée par `getcwd()` est `tampon` si appel réussi et `NULL` dans le cas contraire. Les erreurs possibles: `EACCESS` problème d'accès, `ERANGE` le tampon trop petit. Si le tampon trop petit il faut agrandir et réessayer.

## 2.4 Session et groupe de processus

Une **session** est composée de processus lancés par utilisateur après la connexion. Les sessions sont utiles pour gérer les déconnexions de l'utilisateur pour provoquer la terminaison des tous les processus lancés par l'utilisateur.

Les **groupes de processus** sert au contrôle de tâches (job control). Permet de détacher les tâches en arrière plan et les ramener en avant plan.

### 2.4.1 Terminal de contrôle

Un terminal — un fichier spécial permettant d'effectuer les entrées sorties asynchrones. Fonctionne en mode *full duplex* c'est-à-dire traite les entrées et sorties en parallèle (clavier/écran).

Le terminal de contrôle est désigné par `/dev/tty`.

Même si on ferme les descripteurs 0, 1, 2 nous pouvons toujours (re)ouvrir le terminal de contrôle :

```
int d=open("/dev/tty", O_RDWR);
```

## 2.5 Date de création

La date de création d'un processus est donnée en nombre de secondes depuis EPOCH (1 janvier 1970).

## 2.6 Temps CPU

Pour connaître le temps CPU pour les deux modes: utilisateur et noyau, on utilise la fonction `times` en lui passant l'adresse de la structure `tms`:

```
#include <unistd.h>
clock_t times(struct tms *buf)
```

où `clock_t` un type entier. La fonction met à jour le contenu de la structure `struct tms`. La fonction retourne (`clock_t`) `-1` en cas d'échec et mets dans `errno` le code d'erreur.

Dans ce qui suit le temps `clock_t` est toujours mesuré en nombre de clics d'horloge. La structure `struct tms` définie dans

```
#include <sys/times.h>
struct tms{
    clock_t  tms_utime;  user CPU time
    clock_t  tms_stime;  system CPU time
    clock_t  tms_cutime; user CPU time of terminated child processes
    clock_t  tms_cstime; system CPU time of terminated child processes
}
```

Le temps `tms_utime` et `tms_stime` d'un fils sont ajoutés dans `tms_cutime` et `tms_cstime` pour chaque fils terminé au moment où `wait()` ou `waitpid()` retourne le *pid* de ce fils.

- `tms_utime` le temps que le processus a passé en mode utilisateur (en exécutant les instructions de l'utilisateur).
- `tms_stime` le temps que le processus a passé en mode système,
- `tms_cutime` — c'est la somme de `tms_utime` de processus courant et de `tms_cutime` de chaque processus fils terminé et attendu avec `wait()` ou `waitpid()`.
- `tms_cstime` — est la somme de `tms_stime` de processus courant et de `tms_stime` de chaque fils terminé et attendus avec `wait()` ou `waitpid()`.

Pour déterminer le nombre de clics d'horloge par seconde on utilise `sysconf(_SC_CLK_TCK)`. Ce nombre change d'une machine à l'autre.

### 2.6.1 Le masque de création de fichiers

Le masque de création de fichier est une propriété de processus. Pour changer de masque:

```
#include <unistd.h>
mode_t umask(mode_t masque)
```

Pour les détails voir les explications dans la partie FICHIERS.

## 2.7 Table de descripteurs

Le processus fils hérite la table de descripteurs de processus père.

### 3 Espace d'adressage de processus

L'espace d'adressage virtuel d'un processus est divisé en plusieurs segments :

- la région de code (text segment) contient les instructions que la machine doit exécutées, souvent ce segment est accessible uniquement en lecture pour que le programme ne modifie par erreur son propre code.
- le segment de données initialisées (initialized data segment) contient les variables déclarées en dehors de fonctions et initialisées explicitement dans le programme. Par exemple si le programme contient la déclaration

```
int toto=99;
```

à l'extérieur de toute fonction C alors la mémoire pour la variable `toto` se trouve dans ce segment est initialisée avec la valeur 99 au moment de chargement de programme pour exécution.

- le segment de données non-initialisées contient par les variables de niveau 0 (définies en dehors de fonctions) non-initialisées. Par exemple si

```
long tab[20];
```

est défini en dehors de toute fonction alors la mémoire pour `tab` se trouve dans ce segment et est initialisée à 0 par le noyau au moment de chargement de programme dans la mémoire.

- le tas (heap) pour les allocations dynamiques (les allocation par `malloc`)
- la pile (user stack) pour stocker les variables définies dans les fonctions.

La commande `size` permet d'afficher la taille de différents segments. Par exemple

```
size copier
text      data      bss      dec      hex filename
1172      272        8     1452     5ac copier
```

permet de voir que pour l'exécutable `copier` la taille de segment de code est 1172, segment de données 272, segment bss 8, total 1452 (tout en decimal).

### 4 Création de processus

Pour voir tous les processus la commande `ps axl`

Création d'un nouveau processus:

```
#include <sys/types.h>
#include <unistd.h>
pid_t fork(void)
```

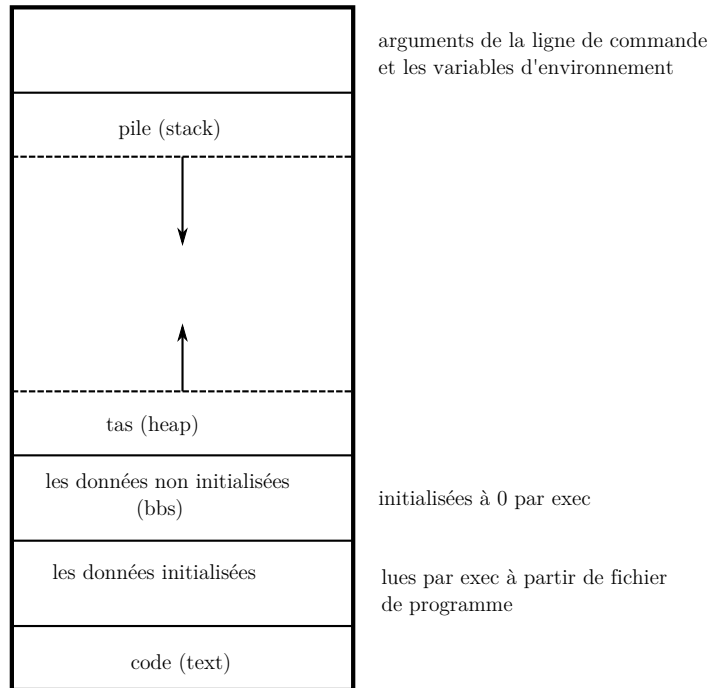


Figure 1: Les segments d'un programme C. Les flèches montrent comment grossissent la pile et le tas.

Tout processus (sauf 0) est créé avec `fork()`

Le processus fils: exécute le même programme que le processus père, avec les mêmes données (copie de données), la même pile d'appel (une copie de la pile du père), le même tas et il possède une copie de la table de descripteurs de fichiers hérité de son père.

A noter que les descripteurs de fichiers ouverts avant `fork` pointent toujours vers les mêmes entrées dans la table de fichiers ouverts, donc le changement de la position courante dans le fichier par un de deux processus change la position courante dans l'autre.

Les différences entre le processus fils et le processus père:

1. chez le fils `fork()` renvoie 0, chez le père il renvoie le *pid* du fils,
2. le processus enfant a son propre *pid* différents de tous les *pid* et *gid*,
3. le *ppid* de l'enfant initialisé avec le *pid* de père,
4. les mesures de temps consommé sont initialisées à 0 pour l'enfant,
5. les verrous posés par le père ne sont pas hérités,
6. si une minuterie est activé `unsigned int alarm(unsigned int secondes)` elle est désactivée chez le fils,
7. chez le fils l'ensemble de signaux pendants est initialisé à l'ensemble vide.

```
#include <sys/types.h>
#include <unistd.h>
pid_t vfork(void)
```

La fonction `vfork()` fait la même chose que `fork()` mais ne fait pas de copie de données, les données sont partagées entre le processus père et le processus fils. L'utilisation de `vfork()` appropriée uniquement pour faire immédiatement `exec()`

## 5 Primitives de recouvrement

Il y a six primitives de recouvrement. Différences entre ces primitives:

1. manière de faire passer les arguments de `main()`:
  - soit une liste `l` (`execl`, `execle`, `execlp`)
  - soit un vecteur `v` (`execv`, `execve`, `execvp`),
2. la manière utilisée pour trouver le fichier à charger:
  - soit en utilisant la variable `PATH` de l'environnement, ce que indique la lettre `p` dans le nom de la fonction, c'est le cas de primitives `execlp()` et `execvp()`,
  - soit relativement au répertoire de travail (ou une référence absolue),
3. l'environnement du processus après le recouvrement, dans le cas de `execve()` et `execle()` un nouvel environnement se substitue à l'ancien.

Les primitives de la famille `exec` ne retournent jamais en cas de succès. En cas d'échec la valeur de retour est `-1`.

```
#include <unistd.h>
int execl(const char *path, const char *arg0, ... /*, (char *)NULL*/);
int execv(const char *path, char *const argv[]);
int execle(const char *path,
           const char *arg0, ... /*, (char *)NULL, char *const envp[]*/);
int execve(const char *path, char *const argv[], char *const envp[]);
int execlp(const char *file, const char *arg0, ... /*, (char *)NULL */);
int execvp(const char *file, char *const argv[]);
```

**Exemples:** Pour remplacer le code actuellement exécuté par le code de `ls`:

```
execl("/bin/ls", "ls", "-l", "/", (char *)NULL);
```

`execlp()` se comporte de la même manière sauf que pour résoudre les les références relatives on utilise la variable `PATH`. Par exemple si `PATH` contient `./bin:/usr/bin` on pourra faire

```
execlp("ls", "ls", "-l", "/", (char *)NULL);
```

La même chose avec le vecteur d'arguments:

```
char *tab[]={"ls", "-l", "/", NULL};
execvp("ls", tab);
```

ou avec un vecteur alloué dynamiquement:

```
char **tab=malloc( 4 * sizeof(char *));
tab[0]="ls";
tab[1]="-l";
tab[2]="/";
tab[3]=NULL;
execvp("ls", tab);
```



## 5.1 Attributs d'un processus après le recouvrement

- (1) **propriétaire effectif** si le bit *set-uid* est positionné sur le fichier chargé le propriétaire de ce fichier devient propriétaire effectif du processus.
- (2) **groupe effectif** si *set-uid* est positionné alors le groupe propriétaire de fichier chargé devient groupe propriétaire effectifs du processus,
- (3) le *handler* par défaut est installé pour les signaux captés,
- (4) **descripteurs** si le bit `DF_CLOEXEC` de fermeture automatique en cas d'exéc est positionné par `fcntl()` alors ce descripteur sera fermé. Les autres descripteurs restent ouverts.

## 6 Terminaison de processus

Le processus se termine soit à sa demande soit à la réception d'un signal provoquant sa terminaison. Le processus qui se termine devient *zombie*, l'information sur la terminaison est enregistrée dans le bloc de contrôle et reste disponible jusqu'à la consultation par le processus père avec les primitives `wait` et `waitpid`

```
#include <unistd.h>
void exit(int valeur)
void _exit(int valeur)
```

provoque

1. `fflush()` sur tous les fichiers ouverts pour vider les tampons,
2. suppression de fichiers créés avec `tmpfile()`
3. la fonction appelle toutes les fonctions enregistrées avec `atexit()` dans l'ordre inverse d'enregistrement,
4. fait appel à `_exit(valeur)`

Seulement le bit de poids faible de `valeur` peut être récupéré. Deux constantes symboliques `EXIT_SUCCESS` et `EXIT_FAILURE` définies dans `stdlib.h` peuvent être utilisées.

L'appel à `_exit()` provoque

1. fermeture de fichiers et répertoires (sans `fflush()`)
2. envoi de signal `SIGHUP` à tous les processus de groupe si le processus qui se termine est le leader,
3. rattachement de fils au processus 1,
4. si le père est en attente sur `wait()` ou `waitpid()` alors il est réveillé sinon le fils devient zombie. Le signal `SIGCHLD` est envoyé au père.

L'instruction `return` exécutée dans `main()` provoque l'appel à `exit`

L'utilisation de `_exit()` est recommandée si on a fait `exec` qui n'a pas réussi après `vfork()`. Si on fait `exit` cela videra les tampons de fichiers utilisés aussi par le père.

## 7 Processus zombie et synchronisation père fils

Le processus qui termine devient *zombie*. Il est maintenu dans la table de processus mais ne consomme plus de ressources système. Le processus peut terminer pour deux raisons : (1) terminaison normale par `exit()`, `_exit()` ou `return` de la fonction `main`, dans ce cas le processus retourne le code de retour, (2) une terminaison à la suite de réception d'un signal.

La raison de garder le processus terminé dans la table de processus c'est pour pouvoir récupérer le code de retour ou l'information sur le signal qui a provoqué la terminaison.

Si on exécute sur la console la commande

```
ps -eo pid,user,args,s --sort user
```

pour voir tous les processus nous pouvons distinguer les processus zombies comme ceux désignés comme `defunct`. De plus l'état de processus est zombie est Z.

Les seules informations maintenues dans le bloc de contrôle après la terminaison de processus:

- son code de retour,
- le temps d'exécution dans les modes utilisateur et noyau,
- son identité et l'identité de son père.

Si le père termine le fils zombie sera attaché au processus 1 qui le supprimera.

Les fonctions `wait()` et `waitpid()` permettent d'obtenir les informations sur le statut d'un processus fils qui est soit terminé soit stoppé.

### 7.1 wait

```
#include <sys/types.h>
#include <sys/wait.h>
pid_t wait(int *p_status)
```

- (1) Si le processus qui fait appel à `wait()` n'a pas de fils la `wait` retourne `-1` et `errno=ECHILD`.
- (2) si le processus possède un fils zombie alors la fonction retourne l'identité d'un fils zombie, si `p_status` n'est pas NULL alors `p_status` reçoit les informations sur la terminaison. Le fils zombie est supprimé de la table de processus.
- (3) Si le processus possède de fils mais aucun fils zombie alors le processus est bloqué
  - en attendant qu'un fils devient zombie,
  - ou qu'il reçoit un signal. Dans ce deuxième cas la valeur de retour est `-1` et `errno` est `EINTR`.

### 7.2 waitpid

```
#include <sys/types.h>
#include <sys/wait.h>
pid_t waitpid(pid_t pid, int *p_status, int options)
```

permet de tester la terminaison d'un processus appartenant au même groupe de processus que le processus donné. Le paramètre `options` indique si l'attente est bloquante ou non.

Le paramètre `pid` permet de spécifier le processus attendu:

valeur <code>pid</code>	attente de la terminaison
<code>&lt; -1</code>	d'un processus quelconque dans le groupe de processus <code> pid </code>
<code>(pid_t)-1</code>	d'un processus fils quelconque (équivalent à <code>wait()</code> )
<code>0</code>	d'un processus fils quelconque du même groupe de processus que l'appelant
<code>&gt; 0</code>	du processus fils d'identité <code>pid</code>

Le paramètre `options` de `waitpid` est une combinaison bit à bit de valeurs suivantes :

<code>WNOHANG</code>	appel non bloquant
<code>WUNTRACED</code>	processus concerné est stoppé et cette information n'est pas encore transmise depuis que le processus est stoppé

Si `wait()` ou `waitpid()` retournent parce que l'état de processus fils est disponible ces deux fonctions retournent le `pid` du processus fils.

Dans le cas contraire les valeurs de retour de `wait()` et `waitpid()` :

- `-1` en cas d'erreur,
- `0` en cas d'échec, l'échec signifie que, en mode non bloquant, le processus fils concerné existe mais il n'est ni terminé ni arrêté.

Si l'appel à `wait()` ou `waitpid()` est réussi alors l'entier `*p_status` permet de connaître la raison de terminaison de processus fils : terminaison par `exit()` ou terminaison à la réception d'un signal.

**Rappel:** Dans le code de retour donné à l'appel de `exit()` juste l'octet de poids faible compte, donc les codes de retour valables sont entre 0 et 255.

Si le processus fils termine à cause d'un signal alors `*p_status` contient le numéro de signal augmenté de 128 si le fichier `core` a été créé.

Les macro-fonctions suivantes permettent d'interroger la valeur de `p_status` :

<code>WIFEXITED</code>	non nulle si le fils a terminé normalement
<code>WIFSIGNALED</code>	non nulle si le fils a terminé à cause d'un signal
<code>WIFSTOPPED</code>	non nulle si le fils est arrêté ( <code>waitpid()</code> avec option <code>WUNTRACED</code> )

Pour plus de précision on utilise les macro-fonctions suivantes :

<code>WEXITSTATUS</code>	donne le code de retour pour le processus terminé normalement
<code>WTERMSIG</code>	si le retour de la macro-fonction <code>WIFSIGNALED</code> est non nul alors cette macro-fonction donne le numéro du signal
<code>WSTOPSIG</code>	si la valeur de retour de la macro <code>WIFSTOPPED()</code> est non nulle <code>XSTOPSIG</code> donne le numéro du signal qui a stoppé le processus.

**Exemple** On attend la terminaison d'un fils et on récupère la valeur `exit` de ce fils (s'il a terminé par `exit`):

```
pid_t fils;
int status, code_ret;
.....
fils=wait(&status);
if( fils < 0 ){
    /* erreur ou signal */
}
```

```

}
else{
    /* fils zombie supprimer*/
    if( WIFEXITED(status) ){ /*le fils a termine par exit*/
        code_ret=WEXITSTATUS(status); /*recuperer le code
                                   * d'exit de fils*/
        ....
    }
}

```

## 8 Job control

Notion de session et de groupe de processus.

**ctr-z** arrête le processus qui est en avant plan dans le terminal en lui envoyant le signal **SIGTSTP** (plus exactement le signal est envoyé vers le groupe de processus en avant plan).

**ctr-q** start

En peut aussi envoyer un signal qui arrête le processus directement :

```
kill -SIGSTOP pid_processus
```

ou

```
kill -SIGTSTP pid_processus
```

Envoyer le signal “continuer exécution”

```
kill -SIGCONT processus
```

La fonction

```

#include <unistd.h>
pid_t getpgrp(void)

```

renvoie le pid du leader du groupe de processus auquel appartient l'appelant.

## 9 Duplication de descripteur avec dup() et dup2()

Le mécanisme de redirections des entrées et sorties standard est possible grâce à l'utilisation homogène de descripteurs de fichiers et s'applique non seulement aux fichiers réguliers mais aussi aux tubes permettant une communication entre les processus.

Pour cela nous avons besoin une opération de duplication de descripteur. Cette opération crée un nouveau descripteur dans la table de descripteurs d'un processus, le nouveau descripteur pointe vers le même entré dans la table de fichiers ouverts.

```

#include <unistd.h>
int dup(int descripteur)

```

L'argument de `dup()` un descripteur de fichier valable. Cette primitive cherche dans la table descripteurs le plus petit descripteurs qui n'est pas utilisé (qui reste disponible) et copie le contenu de descripteur passé en arguments vers le descripteur trouvé. La fonction retourne le nouveau descripteur ou `-1` en cas d'échec.

Le compteurs de descripteurs ouverts dans la table de fichiers ouverts est incrémenté de 1.

Le compteurs de descripteurs ouverts dans la table de fichiers ouverts est incrémenté de 1. Un exemple est présenté sur la figure 9.

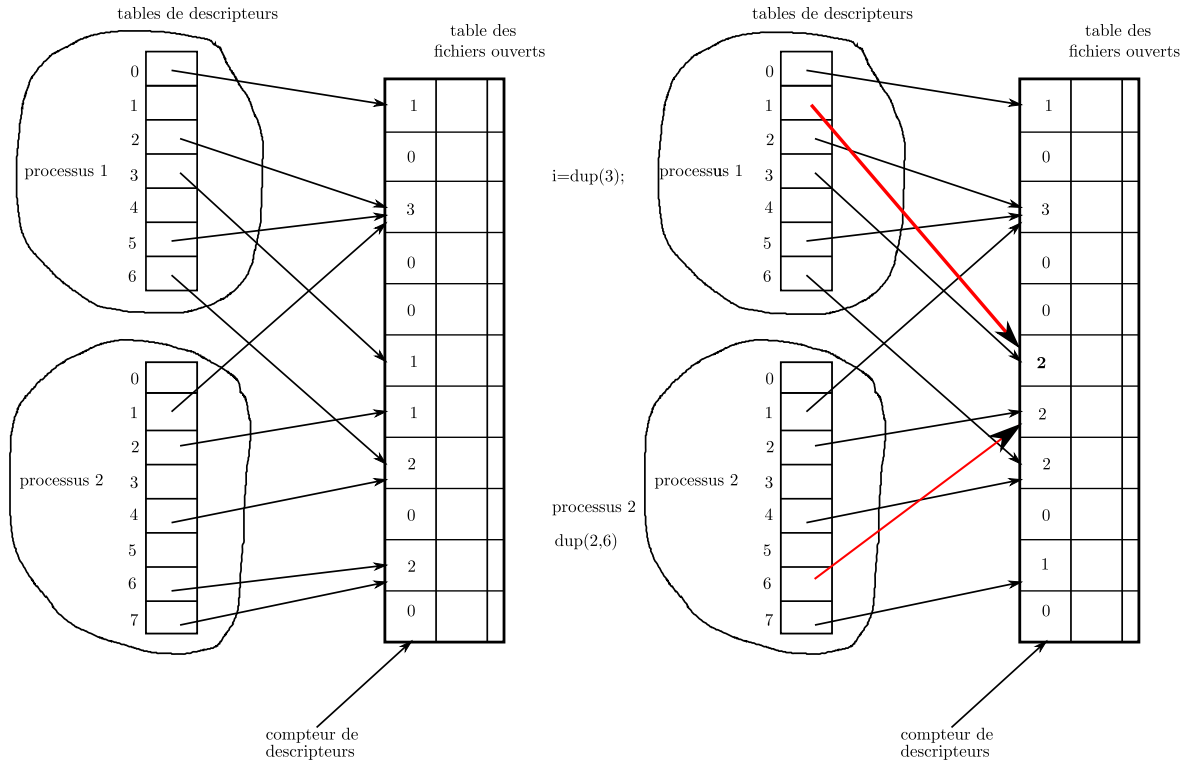


Figure 2: Processus 1 exécute  $i = \text{dup}(3)$ . Le premier descripteur libre pour ce processus est 1 donc  $\text{dup}$  retournera la valeur 1 et le descripteur 3 sera copié vers le descripteur 1. Les descripteurs 1 et 3 pointent maintenant vers le même fichier ouvert. Processus 2 exécute  $\text{dup}(2,6)$ . Cet appel ferme le descripteur 6 et ensuite copie le descripteur 2 vers 6. A gauche la configuration avant que les deux processus exécutent leur  $\text{dup}$ , à droite après.

```
#include <unistd.h>
int dup2(int d1, int d2)
```

La primitive  $\text{dup2}()$  ferme le descripteurs  $d2$  s'il est ouvert et en fait le synonyme de  $d1$ . La primitive retourne  $-1$  en cas de problèmes et  $d2$  sinon.

Après l'appel réussi les descripteurs  $d2$  et  $d1$  pointent vers la même entrée de la table de fichiers ouverts et le compteur de descripteurs pour ce fichier est augmenté de 1.

- Si  $d2 < 0$  ou  $d2 \geq \text{MAX\_OPEN}$   $\text{dup2}()$  retourne  $-1$  et  $\text{errno} = \text{EBADF}$ .
- Si  $d1$  est un descripteur valide et égale à  $d2$  alors la fonction  $\text{dup2}()$  retourne  $d2$  sans le fermer (donc  $\text{dup2}()$  n'a pas d'effet.)
- Si  $d1$  n'est pas un descripteur valide alors  $\text{dup2}()$  retourne  $-1$  mais ne ferme pas  $d2$ .

Le nouveau descripteur alloué par  $\text{dup}()$  et  $\text{dup2}()$  possède les mêmes verrous que l'ancien.

```
#include <fcntl.h>
int fcntl(int filedesc1, int cmd, ...)
```

Pour dupliquer les descripteurs on pourra utiliser aussi la primitive `fcntl()` de manière suivante :

```
int new_desc;
int a_partir;
```

```
new_desc=fcntl(source_desc,F_DUPFD,a_partir);
```

Cette appel cherche le premier descripteur libre  $\geq a\_partir$ , copie `source_desc` vers le descripteur trouvé et retourne le descripteur trouvé.

## 10 Communication par les tubes

Il existe deux types de tubes : anonymes et nommés. Les deux types de tubes ont plusieurs caractéristiques communes.

- (1) Chaque tube est représenté par un i-node, donc l'utilisation de tubes ressemble à l'utilisation de fichiers, en particulier dans les programmes les tubes sont accessibles au moyen de descripteurs.
- (2) Un tube possède deux extrémités, représentées par deux descripteurs différents : un descripteur utilisé pour la lecture et l'autre pour écriture. Donc les tubes réalisent une communication unidirectionnelle.
- (3) On peut utiliser le mécanisme de duplication de descripteurs sur les tubes. Donc il est possible d'avoir plusieurs descripteurs à l'entrée et à la sortie d'un tube.
- (4) Les fils hérite les descripteurs de tube de son père, c'est le mécanisme utilisé pour instaurer une communication entre père et fils.
- (5) Les tubes utilisent le mode fifo, l'ordre de lecture est le même que l'ordre d'écriture. L'information est toujours écrite à une extrémité et lue à l'autre. L'information lue est supprimée du tube. Il n'est pas possible de changer la position dans le tube avec `lseek()`
- (6) Le tube a une capacité limitée, donc si le tube est plein on ne peut plus écrire.
- (7) Si le nombre de lecteurs est nul (aucun descripteur ouvert pour la lecture) l'écrivain essayant écrire dans le tube reçoit le signal **SIGPIPE**.
- (8) Si le nombre d'écrivains est nul (le nombre de descripteurs ouverts pour l'écriture dans le tube est nul) alors le lecteur détecte la fin de fichier pendant la tentative de lecture.

Par défaut on utilise les tubes en mode bloquant:

- l'écrivain essayant écrire dans le tube plein est bloqué en attendant qu'un lecteur lit dans le tube en libérant la place.
- le lecteur essayant lire dans un tube vide attend qu'un écrivain y écrit.

Cependant on peut changer le mode de bloquant vers non bloquant.

### 10.1 Tubes anonymes

Le tube a le compteur de références nul (ces tubes n'apparaissent pas dans un répertoire). Par conséquent un tube anonyme existe (le i-node correspondant existe) s'il y a des processus avec les descripteurs ouverts sur ce tube.

Si le compteur de descripteurs sur le tube anonyme est nul alors le tube sera automatiquement supprimé.

Les premiers descripteurs d'un tube sont créés à la création du tube avec `pipe()`, ensuite d'autres descripteurs peuvent être créés avec le mécanisme de duplication de descripteurs.

Les descripteurs de tubes sont hérités par les processus fils comme d'autres descripteurs de fichiers.

Si le processus a fermé tous les descripteurs d'un tube anonyme soit pour la lecture soit pour l'écriture il lui est impossible d'en acquérir des nouveaux descripteurs pour le même mode (lecture ou écriture).

## 10.2 Création d'un tube anonyme

```
#include <unistd.h>
int pipe(int desc[2])
```

`desc` est un vecteur de deux entiers. Après l'appel `desc[0]` donne le descripteur ouvert en lecture tandis que `desc[1]` donne le descripteur ouvert en écriture.

Si l'appel réussit alors `pipe()` renvoie 0. En cas d'échec `pipe()` renvoie `-1` est

- `errno=EMFILE` – si la table de descripteurs du processus est pleine,
- `errno==ENFILE` – si la table des fichiers ouverts du système est pleine.

Exemple de création d'un tube :

```
int tube[2];

if( pipe(tube) == -1 ){
    perror("creation de tube");
    exit(EXIT_FAILURE);
}
/* maintenant tube[0] - le descripteur d'entree
    tube[1] - le descripteur de sortie
*/
```

## 10.3 Lecture d'un tube

La lecture s'effectue en utilisant le descripteur avec la primitive `read()` (la même qui est utilisée pour la lecture de fichiers ordinaires) :

```
char buf[TAILLE_BUFF];

int nb_lu = read(tube[0], buf, TAILLE_BUF);
```

Si le tube n'est pas vide et contient `n>0` caractères alors le primitif extrait `min(n, TAILLE_BUF)` caractères et les met dans `buf`. `read` renvoie le nombre de caractères lus.

Si le tube est vide

- et le nombre d'écrivains est nul, on détecte la fin du fichier, c'est-à-dire `read()` retourne 0.
- si le nombre d'écrivains n'est pas nul et la lecture est bloquante (ce qui est le cas par défaut) alors le processus est mis en attente jusqu'à ce que le tube ne soit pas vide ou qu'il n'y ait plus d'écrivains,
- si le nombre d'écrivains est non nul et la lecture est non bloquante alors le retour immédiat avec la valeur `-1` et `errno==EAGAIN`.



## 10.4 Écriture dans un tube

L'écriture est réalisée avec la primitive `write()`

```
nb_ecrit=write(tube[1], buf, n);
```

Si le nombre d'octets à écrire est `<PIPE_BUF` l'écriture est garantie atomique (les caractères écrits dans le tube d'un seul coup). Si le nombre de caractères à écrire est supérieur à `PIPE_BUF` l'opération peut-être décomposé par le systèmes.

Si le nombre de lecteur dans le tube est nul alors la tentative d'écriture provoque le signal `SIGPIPE` qui, s'il n'est capté, termine le processus écrivain.

Si le nombre de lecteurs est non nul alors

- si l'écriture est bloquante le retour de `write()` n'aura lieu que quand tous les octets sont écrits. Si `n<=PIPE_BUF` alors l'écriture sera atomique. Le processus peut être endormi s'il n'y a pas assez place libre dans le tube.
- si l'écriture est non bloquante
  - si `n>PIPE_BUF` le retour est un nombre de caractères écrits (inférieur à `n`) ou `-1` et `errno==EAGAIN`
  - si `n<=PIPE_BUF` et il y a au moins `n` places libres dans le tube alors une écriture atomique de `n` octets est réalisée,
  - si `n<=PIPE_BUF` et il n'y pas `n` places libres dans le tube alors rien n'est écrit, `write()` retourne `-1` et `errno==EAGAIN`.

Par défaut la lecture/écriture d'un tube anonyme s'effectue en mode bloquant. Pour changer le mode à non bloquant il faut procéder comme indique la section 12.

## 10.5 Exemple - communication père vers fils par un tube anonyme

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

int tube[2];
#define BSIZE = 100
char buf[BSIZE];
ssize_t nbytes;
int status;

status = pipe(tube);
if (status == -1 ) {
    /* traiter erreur de pipe() */
    ...
}

switch (fork()) {
case -1: /* traiter erreur de fork() */
```

```

        break;

case 0: /* fils - lit du tube */
        close(tube[1]);
                                /* fils n'ecrit pas,
                                * fermer le descripteur
                                * d'écriture */

        while( (nbytes = read(tube[0], buf, BSIZE) ) > 0){ /* lire dans le tube*/
                /* traitement */
        }

        close(tube[0]);
                                /* fils termine */
        exit(EXIT_SUCCESS);

default: /* processus pere - ecrit dans le tube */
        close(tube[0]);
                                /* le pere ne lit pas
                                * fermer le descripteur de lecture */

        write(tube[1], "Hello world\n", 12); /* pere ecrit */
        close(tube[1]);
                                /* fermer le tube pour que
                                le fils voie la fin de fichier */

        exit(EXIT_SUCCESS);
}

```

## 10.6 Exemple - redirections

Supposons que voulons établir une communication entre le processus parent et le processus enfant par deux tubes anonymes, la tube `parent2enfant` sert à envoyer les messages de parent vers l'enfant tandis que le tube `enfant2parent` permet d'envoyer le message de l'enfant vers le parent.

On suppose que le processus enfant lit après le recouvrement les messages en utilisant le descripteur `STDIN_FILENO` (descripteur 0) et écrit sur le descripteur `STDOUT_FILENO` (descripteur 1). Chez l'enfant la lecture sur le descripteur 0 doit correspondre à la lecture de tube `parent2enfant` et l'écriture sur le descripteur 1 doit correspondre à l'écriture dans le tube `enfant2parent`.

La création de tube, de processus enfant et les redirections sont réalisés par la fonction `redirection` suivante :

---

```

#define _POSIX_C_SOURCE 200112L
#include <unistd.h>
#include <sys/types.h>
#include "redirection.h"
#include "panic.h"
/*****
 * Cette fonction cree un processus enfant et
 * etablit une communication entre le parent et enfant par deux tubes anonymes.
 * Les parametres d'entree:
 * char *ref — chemin vers executable execute par le processus enfant
 * char **argv — le tableau d'arguments pour l'executable de l'enfant
 * Les parametres de sortie:

```

```

* int *in — *in le descripteur de lecture pour le parent
* int *out — *out le descripteur d'écriture pour le parent
* Valeur de retour : le pid de l'enfant.
*****/

pid_t redirection(char *ref, char **argv, int *in, int *out){
    /* les tableaux de descripteur :
     * parent2enfant pour le tube de parent vers l'enfant,
     * enfant2parent pour le tube de l'enfant vers le parent. */
    int parent2enfant[2], enfant2parent[2];
    pid_t pid;

    if( pipe(parent2enfant) < 0 ){
        PANIC(0);
        return -1;
    }
    if( pipe(enfant2parent) < 0 ){
        PANIC(0);
        return -1;
    }

    if( ( pid=fork() ) == (pid_t) -1 ){
        PANIC(0);
        return -2;
    }
    /** processus enfant **/
    else if( pid == 0 ){
        /*fermer les descripteurs inutiles pour l'enfant*/
        close(enfant2parent[0]); close(parent2enfant[1]);

        /*pour le processus enfant la lecture sur
         * le descripteur 0
         * devait donner la lecture dans le tube parent2enfant*/
        if( dup2( parent2enfant[0], STDIN_FILENO ) == -1 ){
            PANIC(0);
            return -3;
        }
        /*fermer le descripteur inutile apres la duplication*/
        close(parent2enfant[0]);

        /* pour l'enfant l'écriture sur le descripteur 1
         * devait en fait provoquer l'écriture dans
         * le tube enfant2parent*/
        if( dup2( enfant2parent[1], STDOUT_FILENO ) == -1 ){
            PANIC(0);
            return -4;
        }
        /* fermer le descripteur inutile apres la duplication*/
        close(enfant2parent[1]);

        /*recouvrement avec exec*/
        execvp(ref, argv);
        PANIC(0);
        return -5;
    }
}

```

```

    else{/*parent*/
        /* fermer les descripteurs inutiles pour le parent*/
        close(enfant2parent[1]); close(parent2enfant[0]);
        /*preparer les parametres de retour */
        *out = parent2enfant[1]; *in = enfant2parent[0];
        return pid;
    }
}

```

---

L'exemple d'utilisation de cette fonction :

---

```

#define _POSIX_C_SOURCE 200112L
#include <stdio.h>
#include <sys/types.h>
#include <stdlib.h>
#include <unistd.h>
#include <string.h>
#include <sys/wait.h>
#include "redirection.h"
#include "panic.h"
int main(void){
#define T 512
    char tampon[T];
    char *text[]={"text_a_envoyer_vers_le_processus_fils.\n",
                  "et_un_autre_text" };
    /*les enfants vont executer la commande tr*/
    char *ref="tr";
    char *arg[]={"tr", "a-z", "A-Z", (char *) NULL };

    int in, out;

    int dim=sizeof text/sizeof text[0];
    for(int i=0; i<dim; i++){
        int k;
        pid_t f = redirection(ref, arg, &in, &out);

        if(f < 0){
            fprintf(stderr, "echec_redirection\n");
            exit(1);
        }
        /* ecrire vers l'enfant*/
        if( (k = write(out, text[i], strlen(text[i])) ) < 0){
            PANIC(0);
            continue;
        }
        close(out);

        /*lire ce que l'enfant envoie*/
        k=read(in, tampon, sizeof tampon);
        if(k < 0){
            PANIC(0);
            continue;
        }
        close(in);
        tampon[k]='\0';
        printf("[%s]\n", tampon);
    }
}

```

```

    /* attendre la terminaison de l'enfant */
    wait(NULL);
}
exit(0);
}

```

---

## 11 Tubes nommés

### 11.1 Création et ouverture

Un tube nommé peut être créé depuis un interpréteur shell avec

```
mkfifo reference
```

Par exemple

```

mkfifo toto
ls -l|grep toto
prw-r--r-- 1 zielonka zielonka      0 2010-12-06 18:16 toto|

```

Nous pouvons voir que `mkfifo` a créé un fichier spécial `toto` et type de fichier est `p` (*pipe*). Faire `man mkfifo` pour voir les options de `mkfifo`.

Depuis un programme C la création d'un tube nommé s'effectue avec

```

#include <sys/types.h>
#include <sys/stat.h>
int mkfifo(const char *chemin, mode_t droits);

```

où `chemin` donne le chemin d'accès au tube et `droits` les droits d'accès (comme pour `open()`).

Une fois le tube créé on peut y accéder en l'ouvrant avec `open()`, de la même façon que pour les fichiers ordinaires. Les seuls modes possibles sont `O_RDONLY` et `O_WRONLY`. Par défaut `open()` est bloquant :

- une demande d'ouverture en lecture est bloquante en absence d'écrivain,
- une demande d'ouverture en écriture est bloquante s'il n'y pas de lecteurs.

Donc le processus peuvent se synchroniser en utilisant le tube, un en lecture et l'autre en écriture.

L'ouverture en mode non bloquant est réalisée avec le drapeau `O_NONBLOCK`

**Ouverture non bloquante en lecture** réussit immédiatement même en absence d'écrivains. Les opérations de lecture qui suivent restent non bloquantes (jusqu'à le changement de mode avec `fcntl()`).

**Ouverture non bloquante en écriture** sans lecteur échoue. La valeur de retour de `open()` est `-1` avec `errno==ENXIO`

## 11.2 Lecture et écriture

La lecture et écriture d'un tube nommé s'effectue de la même façon que pour le tube anonyme, voir les sections 10.3 et 10.4.

## 12 Changement de mode : bloquant ↔ non bloquant pour les tubes nommés ou anonymes

Le changement de mode bloquant/non bloquant pour les tubes s'effectue de la même façon que pour les fichiers ordinaires.

Pour un descripteur de tube (anonyme ou nommé) on peut passer de mode bloquant à non bloquant par un appel à `fcntl` :

```
#include <unistd.h>
#include <fcntl.h>
fcntl(descripteur, F_SETFL, O_NONBLOCK);
```

Cet appel retourne une valeur `int > 0` si l'appel réussit.

Dans le sens inverse il est possible de passer de mode non bloquant à bloquant par :

```
#include <unistd.h>
#include <fcntl.h>
fcntl(descripteur, F_SETFL, 0);
```

## 13 `setuid` et `setgid`

Propriétaire réel du processus – celui qui a lancé le processus, ou plus exactement c'est le propriétaire du processus qui a lancé notre processus avec `exec`.

Le propriétaire effectif – c'est

- soit le propriétaire réel (dans la plupart de cas),
- soit le propriétaire de fichier exécutable si `setuid` bit est positionné sur le fichier exécutable.

```
> ls -l |grep writer*

-rwxr-xr-x 1 zielonka zielonka 10161 2010-12-08 17:19 writer
-rw-r--r-- 1 zielonka zielonka   939 2010-12-07 01:13 writer.c

> chmod u+s writer*
> ls -l |grep writer*

-rwsr-xr-x 1 zielonka zielonka 10161 2010-12-08 17:19 writer
-rwSr--r-- 1 zielonka zielonka   939 2010-12-07 01:13 writer.c
```

La lettre `s` ou `S` indique que `setuid` bit est positionné, `s` indique que `setuid` est positionné sur un fichier exécutable, `S` indique que le fichier n'est pas exécutable mais `setuid` est positionné quand même.

Les bits `setuid` et `setgid` sont utilisés quand on charge un exécutable avec `exec` pour déterminer le propriétaire effectif de processus.