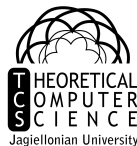


Systemy Operacyjne.

Jakub Kozik

semestr zimowy 2016/2017

Informatyka Analityczna
tcs@jagiellonian



Zasady

"This is not Nam. [...] There are rules."

Zaliczenie i egzamin

- 2 duże zadania + obrony projektów (2 x 30 punktów),
- 2 małe zadania (2 x 10 punktów),
- aktywność na ćwiczeniach (po 1 punkcie na każde ćwiczenia)
- egzamin pisemny (20 punktów).

Za każde z dużych zadań oraz za egzamin trzeba uzyskać przynajmniej 50% możliwych punktów.

Przeliczenie punktów na ocenę:

50-60 3.0; 60-70 3.5; 70-80 4.0; 80-90 4.5; 90-100 5.0

Zaliczenie w II-gim terminie

Po terminie oddania maksymalna liczba punktów za każde z zadań spada liniowo do 50% w ciągu dwóch tygodni.

Program Wykładu

- 1 POSIX - strona użytkownika.
- 2 MINIX - strona systemu.

Zagadnienia:

- Procesy.
- Wejście/Wyjście.
- Pamięć.
- System plików.

THE MINIX BOOK



Andrew S Tanenbaum, Albert S Woodhull,
Operating Systems Design and Implementation,
3rd Edition, Pearson Prentice Hall 2009

- Andrew S. Tanenbaum, **Systemy operacyjne**
 - Abraham Silberschatz, James L. Peterson, Peter B. Galvin, **Podstawy systemów operacyjnych**
- 1 <http://www.minix3.org/>
 - 2 POSIX.1-2008 - IEEE Std 1003.1™-2008 - The Open Group Technical Standard Base Specifications, Issue 7.

Outline

1 Zasady

2 POSIX

• Wstęp

- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

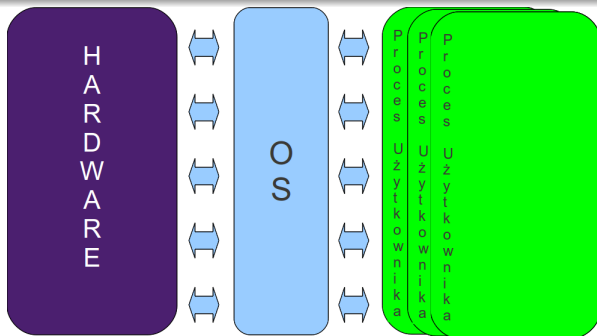
- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

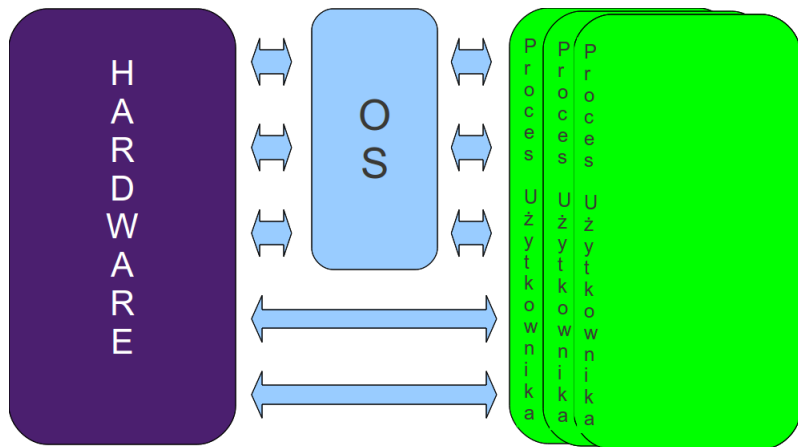
- Batch systems

Główne funkcje systemu.

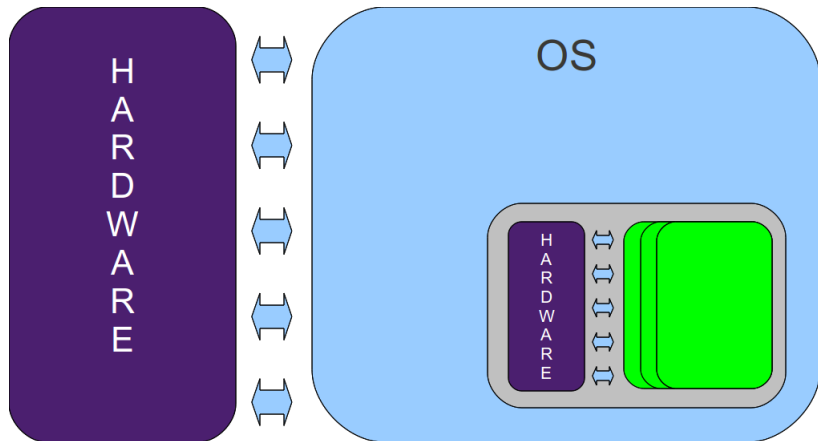
- Extended Machine
- Resource Management



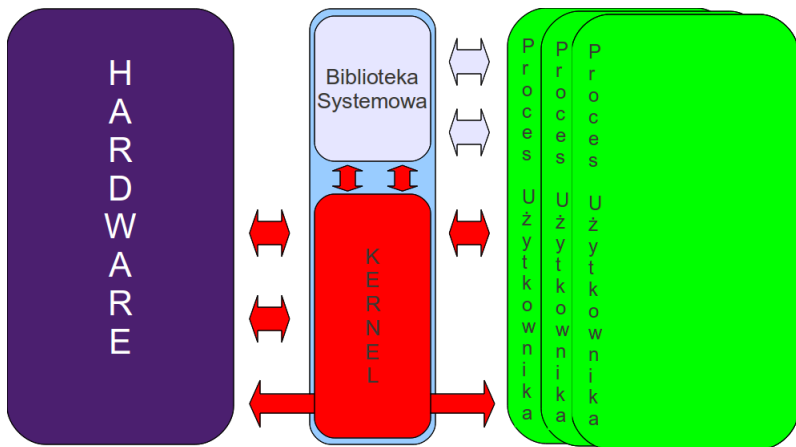
Bez zarządzania zasobami.



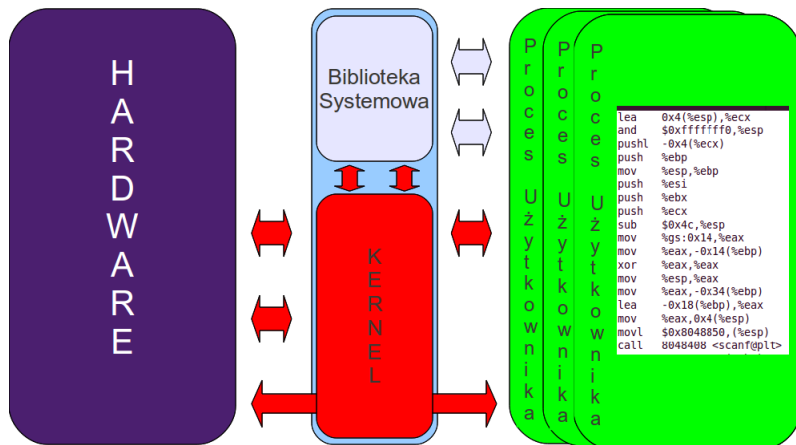
Wirtualna maszyna.



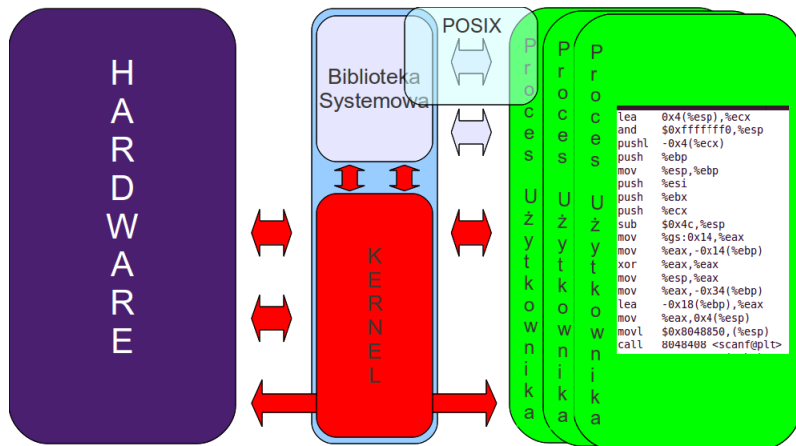
"Złoty środek."



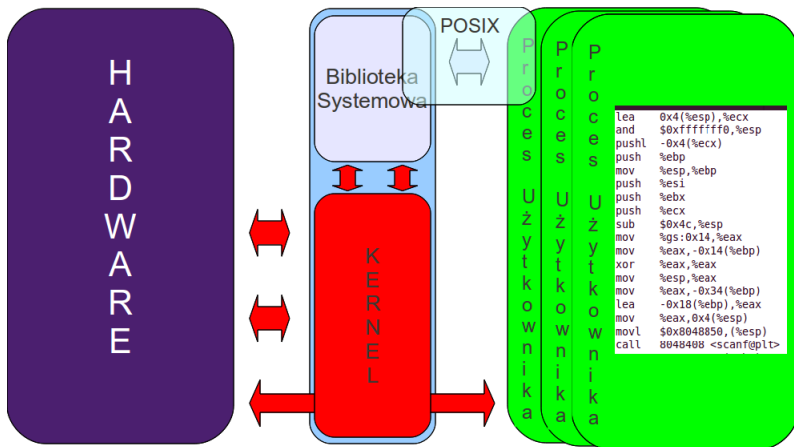
System calls - wywołania systemowe.



POSIX



POSIX programming.



Outline

1 Zasady

2 POSIX

- Wstęp
- **POSIX - standard**
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

POSIX

Portable Operating System Interface

”**POSIX.1-2008** is simultaneously IEEE Std 1003.1TM-2008 and The Open Group Technical Standard Base Specifications, Issue 7.”

POSIX principles:

- Application-Oriented
- Interface, Not Implementation
- Source, Not Object, Portability
- The C Language (ISO C)
- No Superuser, No System Administration
- Minimal Interface, Minimally Defined
- Broadly Implementable
- Minimal Changes to Historical Implementations
- Minimal Changes to Existing Application Code

IEEE Std. 1003.1-1990 Standard for Information Technology –
Portable Operating System Interface (POSIX) –
ART 1. System Application Programming Interface (API)
[C Language].

Donald Lewine, POSIX Programmers Guide, O'Reilly Media 1991

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- **POSIX - procesy**
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Proces

Program w trakcie wykonywania.

Procesy - system calls

```
1 #include <stdio.h>
2
3 int
4 main( int argc , char *argv [])
5 {
6     printf( "Hey, _you _sass _that _hoopy _Ford _Prefect? \n" );
7 }
```

Procesy - system calls

```
1 #include <stdio.h>
2
3 int
4 main(int argc , char *argv[])
5 {
6     printf("Hey, _you _sass _that _hoopy _Ford _Prefect?\n");
7 }
```

exit

```
#include <unistd.h>
void _exit(int status);

#include <stdlib.h>
void exit(int status);
```

Procesy - system calls

Linux - x86

08048080 <_start>:

8048080: b8 04 00 00 00	mov	\$0x4,%eax
8048085: bb 01 00 00 00	mov	\$0x1,%ebx
804808a: b9 a0 90 04 08	mov	\$0x80490a0,%ecx
804808f: ba 06 00 00 00	mov	\$0x6,%edx
8048094: cd 80	int	\$0x80
8048096: b8 01 00 00 00	mov	\$0x1,%eax
804809b: cd 80	int	\$0x80

exit

```
#include <unistd.h>
void _exit(int status);
```

```
#include <stdlib.h>
void exit(int status);
```

fork

```
#include <unistd.h>
pid_t fork(void);
```

```
int
main(int argc, char *argv[])
{
    int k;

    printf("%d,%d\n", \
        getpid(), getppid());
    k= fork();
    printf("%d,%d,%d\n", \
        k, getpid(), getppid());
}
```

- unique process ID.
- different parent process ID
- own copy of the parent's descriptors.
- no pending signals, inactive alarm timer

Fork bomb.

```
int main(){
    while (1) fork();
}
```

execve

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

```
int execve(const char *path, char *const argv[], char *const envp[]);
```

```
extern char **environ;
```

```
int execl(const char *path, const char *arg0, ... /*, (char *)0 */);
```

```
int execlp(const char *path, const char *arg0, ... /*,  
           (char *)0, char *const envp[] */);
```

```
int execlp(const char *file, const char *arg0, ... /*, (char *)0 */);
```

```
int execlp(const char *file, const char *arg0, ... /*, (char *)0 */);
```

```
int execlp(const char *file, const char *arg0, ... /*, (char *)0 */);
```

```
int execlp(const char *file, const char *arg0, ... /*, (char *)0 */);
```

```
int execlp(const char *file, const char *arg0, ... /*, (char *)0 */);
```

Deskryptory procesu wywołującego exec pozostają otwarte (domyślnie).

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

```
int execve(const char *path, char *const argv[], char *const envp[]);
```

tic.c

tictac.c

```
1 #include <stdio.h>
2
3 int
4 main(int argc, char *argv[])
5 {
6     int i;
7
8     for (i=0; i<10; i++) {
9         printf("%s\n", argv[1]);
10        sleep(1);
11    }
12 }
```

```
1 #include <unistd.h>
2
3 int
4 main(int argc, char *argv[])
5 {
6     char* str;
7
8     if (fork()) str = "tic";
9     else str = "tac";
10
11    execl("tic", "tic", str, NULL);
12 }
```

waitpid

```
#include <sys/wait.h>
```

```
pid_t wait(int *stat_loc);
```

```
pid_t waitpid(pid_t pid, int *stat_loc, int options);
```

```
wait(stat_loc)  $\equiv$  waitpid(-1, stat_loc, 0)
```


waitpid

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/types.h>
5 #define BSIZE 100
6
7 int main(){
8     char str[BSIZE];
9     pid_t chld_pid;
10
11     while (fgets(str,BSIZE,stdin)){
12         chld_pid = fork();
13         if (!chld_pid){
14             execlp("echo","echo",str,NULL);
15             exit(1);
16         } else
17             waitpid(chld_pid, NULL, 0);
18     }
19 }
```

waitpid – exit(1) ???

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <sys/types.h>
5 #define BSIZE 100
6
7 int main(){
8     char str[BSIZE];
9     pid_t chld_pid;
10
11     while (fgets(str,BSIZE,stdin)){
12         chld_pid = fork();
13         if (!chld_pid){
14             execlp("echo","echo",str,NULL);
15             exit(1);
16         } else
17             waitpid(chld_pid,NULL,0);
18     }
19 }
```

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- **POSIX - pliki**
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

File Descriptor

“A **per-process unique, non-negative integer** used to identify an open file for the purpose of file access.

The value of a file descriptor is from zero to OPEN_MAX.”

limits.h

```
#define _POSIX_OPEN_MAX    16 /* a process may have 16 files open */  
...  
#define OPEN_MAX          20 /* # open files a process may have */
```

Open File Description

“A record of how a process or group of processes is accessing a file. Each file descriptor refers to exactly one open file description, but an open file description can be referred to by more than one file descriptor. The file offset, file status, and file access modes are attributes of an open file description.”

Deskryptory plików

Domyślnie otwarte deskryptory.

0 - stdin

1 - stdout

2 - stderr

open (zwraca deskryptor dla otwartego pliku)

```
#include <sys/types.h>
```

```
#include <fcntl.h>
```

```
int open(const char *path, int flags [, mode_t mode]);
```

O_RDONLY open for reading only

O_WRONLY open for writing only

O_RDWR open for reading and writing

O_NONBLOCK do not block on open

O_APPEND append on each write

O_CREAT create file if it does not exist

O_TRUNC truncate size to 0

O_EXCL error if create and file exists

Semafor na plikach.

Atomic lock.

`(O_CREAT | O_EXCL)` - `open()` shall fail if the file exists.

```
#define LOCKFILE "/etc/ptmp"
...
int pfd; /* Integer for file descriptor returned by open() call. */
...
if ((pfd = open(LOCKFILE, O_WRONLY | O_CREAT | O_EXCL,
    S_IRUSR | S_IWUSR | S_IRGRP | S_IROTH)) == -1)
{
    fprintf(stderr, "Cannot open /etc/ptmp. Try again later.\n");
    exit(1);
}
...
```

creat & close

close

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

```
int close(int d);
```

creat

```
#include <sys/types.h>
```

```
#include <fcntl.h>
```

```
int creat(const char *name, mode_t mode)
```

```
creat(path, mode)  $\equiv$  open(path, O_WRONLY|O_CREAT|O_TRUNC, mode)
```

read

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

```
ssize_t read(int d, void *buf, size_t nbytes);
```

(zwraca liczbę przeczytanych byte'ów)

(0 \rightarrow EOF)

If a `read()` is interrupted by a signal before it reads any data, it shall return `-1` with `errno` set to `[EINTR]`.

If a `read()` is interrupted by a signal after it has successfully read some data, it shall return the number of bytes read.

write

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

```
ssize_t write(int d, const void *buf, size_t nbytes);
```

(zwraca liczbę zapisanych byte'ów)

If `write()` is interrupted by a signal before it writes any data, it shall return `-1` with `errno` set to `[EINTR]`.

If `write()` is interrupted by a signal after it successfully writes some data, it shall return the number of bytes written.

lseek

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

#define SEEK_SET 0    /* offset is absolute */
#define SEEK_CUR 1    /* relative to current position */
#define SEEK_END 2    /* relative to end of file */

off_t lseek(int d, off_t offset, int whence)
```

```
1#include <sys/stat.h>
2#include <fcntl.h>
3#include <unistd.h>
4
5int main(int argc, char* argv[]){
6    int fd=open("foo", O_RDWR|O_CREAT|O_TRUNC, S_IRUSR|S_IWUSR);
7
8    lseek(fd,10000000000L, SEEK_CUR); /*~10GB*/
9    write(fd, "a", 1);
10   close(fd);
11 }
```

pipe

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

```
int pipe(int fildes[2])
```

```
1#include <stdio.h>
2
3int main(int argc, char* argv[]){
4    int fd[2];
5
6    if (pipe(fd) != 0) return 1;
7    if (fork()){
8        write(fd[1], "say_something", 13);
9    } else {
10        char buf[21];
11        int n;
12        if (n = read(fd[0], buf, 20) >= 4){
13            buf[n] = 0;
14            printf("%s\n", buf+4);
15        }
16    }
17    return 0;
18}
```

Named pipe - FIFO

mkfifo & mknod

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

```
int mknod(const char *path, mode_t mode, dev_t dev)
int mkfifo(const char *path, mode_t mode)
```

pipe

Mknod may be invoked only by the super-user, unless it is being used to create a fifo.

The call `mkfifo(path, mode)` is equivalent to

```
mknod(path, (mode & 0777) | S_IFIFO, 0)
```

Pipe r/w rules.

Bad news

The behavior of multiple concurrent **reads** on the same pipe, FIFO, or terminal device is **unspecified**.

From read() - rationale

I/O is intended to be atomic to ordinary files and pipes and FIFOs. Atomic means that all the bytes from a single operation that started out together end up together, without interleaving from other I/O operations. It is a known attribute of terminals that this is not honored, and terminals are explicitly (and implicitly permanently) excepted, making the behavior unspecified. The behavior for other device types is also left unspecified, but the wording is intended to imply that future standards might choose to specify atomicity (**or not**).

Pipe r/w rules.

Good news

Write requests of `PIPE_BUF` bytes or less shall not be interleaved with data from other processes doing writes on the same pipe.

Writes of greater than `PIPE_BUF` bytes may have data interleaved, on arbitrary boundaries, with writes by other processes

```
1#include <sys/stat.h>
2#include <string.h>
3#include <stdio.h>
4int main(){
5    int fd[2],n;
6    char buf[4];
7    pipe(fd);
8    if (!fork()) {
9        while ((n = read(fd[0], buf, 3))>0){
10            buf[n] = 0;
11            printf("%s\n", buf);
12        }
13    } else {
14        sleep(1);
15        if (fork()) strcpy(buf, "tic");
16        else strcpy(buf, "tac");
17
18        for (n=0; n<10; n++){
19            write(fd[1], buf, 3);
20            sleep(1);
21        }
22    }
23 }
```

fcntl - file descriptor control functions

```
#include <fcntl.h>
```

```
int fcntl(int fd, int cmd, [data])
```


`fcntl(fd, F_DUPFD, int fd2)`

```
1 #include <fcntl.h>
2 #include <unistd.h>
3
4 int main(){
5     int fd[2];
6     pipe(fd);
7     if (!fork()) {
8         close(0);
9         close(fd[1]);
10        fcntl(fd[0], F_DUPFD, 0);
11        execlp("cat", "cat", NULL);
12    } else {
13        write(fd[1], "say_hello\n", 10);
14        close(fd[1]);
15        wait(NULL);
16    }
17 }
```

fcntl(fd, F_GETFD, int fd2) - fd flags

```
1#include <fcntl.h>
2#include <unistd.h>
3
4int main(){
5    int fd[2];
6    pipe(fd);
7
8    int flags = fcntl(fd[1], F_GETFD);
9    flags |= FD_CLOEXEC;
10   fcntl(fd[1], F_SETFD, flags);
11
12   if (!fork()) {
13       close(0);
14       /* close(fd[1]); */
15       fcntl(fd[0], F_DUPFD, 0);
16       execlp("cat", "cat", NULL);
17   } else {
18       write(fd[1], "say_hello\n", 10);
19       close(fd[1]);
20   }
21 }
```

`fcntl(fd, F_GETFL, int fd2)` - file status flags

`fcntl(fd, F_GETFL)`

Return the file status flags and file access modes associated with the file associated with file descriptor `fd`.

`fcntl(fd, F_SETFL, int flags)`

Set the file status flags of the file referenced by `fd` to `flags`. Only `O_NONBLOCK` and `O_APPEND` may be changed. Access mode flags are ignored.

```

1#include <fcntl.h>
2#include <errno.h>
3#include <unistd.h>
4int main(){
5    int fd[2],n;
6    pipe(fd);
7    int flags=fcntl(fd[0], F_GETFL);
8    fcntl(fd[0], F_SETFL, flags | O_NONBLOCK);
9    if (!fork()) {
10        char buf[20];
11        close(fd[1]);
12        while ((n=read(fd[0], buf, 20))!=0){
13            if (n>0) write(0, buf, n);
14            else if (errno!=EAGAIN) return 1;
15            else write(0, " still _nothing\n", 14);
16            sleep(1);
17        }
18    } else
19    for (n=0; n<5; n++){
20        sleep(3);
21        write(fd[1], " I _am _a _walrus.\n", 16);
22    }
23 }

```

O_NONBLOCK for open

```
1#include <fcntl.h>
2#include <errno.h>
3#include <unistd.h>
4#include <sys/stat.h>
5#include <stdio.h>
6
7int main(){
8    int fd,n; char buf[20];
9    mkfifo("mfifo", S_IWUSR | S_IRUSR);
10
11    if (!fork()) {
12        fd=open("mfifo", O_RDONLY|O_NONBLOCK);
13        write(1, "opened\n", 7);
14        sleep(10);
15        while ((n=read(fd, buf, 20))>0)
16            write(1, buf, n);
17    } else{
18        sleep(5);
19        fd=open("mfifo", O_WRONLY);
20        write(fd, "hello\n", 6);
21        write(1, "done\n", 5);
22    }
23 }
```

Advisory record locking.

```
fcntl(fd, F_GETLK, struct flock *lkp)
```

Find out if some other process has a lock on a segment of the file associated by file descriptor `fd` that overlaps with the segment described by the `flock` structure pointed to by `lkp`. [...]

```
fcntl(fd, F_SETLK, struct flock *lkp)
```

Register a lock on a segment of the file associated with file descriptor `fd`. [...] This call returns an error if any part of the segment is already locked.

```
fcntl(fd, F_SETLKW, struct flock *lkp)
```

Register a lock on a segment of the file associated with file descriptor `fd`. [...] This call blocks waiting for the lock to be released if any part of the segment is already locked.

```
struct flock {
    short    l_type;        /* F_RDLCK, F_WRLCK, or F_UNLCK */
    short    l_whence;      /* SEEK_SET, SEEK_CUR, or SEEK_END */
    off_t    l_start;       /* byte offset to start of segment */
    off_t    l_len;         /* length of segment */
    pid_t    l_pid;         /* process id of the locks' owner */
};
```

```

1#include <fcntl.h>
2#include <unistd.h>
3#include <sys/stat.h>
4int main(int argc, char * argv[]){
5    int fd;
6    struct flock fl;
7
8    fd= open("lock", O_CREAT | O_RDWR, S_IWUSR | S_IRUSR );
9    /* ... */
10    fl.l_type = F_WRLCK;
11    fl.l_whence = SEEK_SET;
12    fl.l_start = 0;
13    fl.l_len = 3;
14
15    fcntl(fd, F_SETLKW, &fl);
16    lseek(fd, 0, SEEK_SET);
17    write(fd, argv[1], 3);
18    sleep(30);
19    fl.l_type = F_UNLCK;
20    fcntl(fd, F_SETLK, &fl);
21    /* ... */
22    close(fd);
23}

```


Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- **POSIX - sygnały**
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Sygnały

Sygnał

Informacja o **asynchronicznym** zdarzeniu/błędzie.

Ctrl-c

Ctrl-c powoduje wysłanie sygnału SIGINT do wszystkich procesów z *foreground process group*.

Dzielenie przez 0

Dzielenie liczby (int) przez (int) 0 powoduje wysłanie sygnału SIGFPE do procesu.

Źródła sygnałów.

Terminal Ctrl-C SIGINT,
Ctrl-\SIGQUIT

Hardware dzielenie przez 0 SIGFPE,
niewłaściwe odwołanie do pamięci SIGSEGV,...

Proces syscall kill, domyślny sygnał SIGTERM

System - Software conditions SIGALARM,
SIGPIPE (broken pipe)

Sygnały które nie docierają do adresata

SIGKILL

SIGSTOP

Wysyłanie sygnałów.

```
#include <signal.h>
```

```
int kill(pid_t pid, int sig);
```

Permission

...the real or effective user ID of the sending process shall match the real or saved set-user-ID of the receiving process.

Adresaci - pod warunkiem że można do nich wysłać

`pid > 0` proces, którego ID jest równe `pid`

`pid = 0` procesy z tej samej grupy

`pid = -1` wszystkie procesy

`pid < -1` wszystkie procesy z grupy o ID równym `|pid|`

```
int raise(int sig);
```

Obsługa sygnałów - ISO C

```
#include <signal.h>

void (*signal(int signo, void (*func)(int)))(int);

typedef void Sigfunc(int);

Sigfunc *signal(int, Sigfunc *);
```

Obsługa sygnałów

SIG_DFL domyślna obsługa sygnału

SIG_IGN sygnał jest ignorowany

wskaźnik do funkcji która ma obsłużyć sygnał

Znikające i nieobsłużone sygnały.

```
1#include <stdio.h>
2#include <signal.h>
3
4void handler(int sig_nb){
5    write(1,"If everything seems under control , \
6you're just not going fast enough.\n",70);
7    sleep(1);
8    signal(SIGINT, handler);
9}
10
11int main(){
12    signal(SIGINT, handler);
13
14    while (1)
15        pause();
16}
```

Obsługa sygnałów - POSIX

```
#include <signal.h>
```

```
int sigaction(int sig, const struct sigaction *restrict act,  
              struct sigaction *restrict oact);
```

void	(*sa_handler)(int)	Pointer to a signal-catching function or one of the SIG_IGN or SIG_DFL.
sigset_t	sa_mask	Set of signals to be blocked during execution of the signal handling function.
int	sa_flags	Special flags.
void	(*sa_sigaction)(int, siginfo_t *, void *)	Pointer to a signal-catching function.

```
#include <signal.h>

int sigemptyset(sigset_t *set);
int sigfillset(sigset_t *set);
int sigaddset(sigset_t *set, int signo);
int sigdelset(sigset_t *set, int signo);
int sigismember(const sigset_t *set, int signo);
```

Signal mask for the duration of the signal-catching function

This mask is formed by taking the union of the current signal mask and the value of the `sa_mask` for the signal being delivered, and unless `SA_NODEFER` or `SA_RESETHAND` is set, then including the signal being delivered.


```

1 #include <unistd.h>
2 #include <sys/types.h>
3 #include <signal.h>
4
5 volatile int ready = 0;
6 void handler(int sig_nb){
7     ready = 1;
8 }
9
10 int main(){
11     pid_t other;
12     char* str="tic\n";
13     struct sigaction act;
14
15     act.sa_handler= handler;
16     act.sa_flags = 0;
17     sigemptyset(&act.sa_mask);
18     sigaction(SIGUSR1,&act,NULL);
19
20     if (!(other=fork())){
21         str = "tac\n";
22         other= getpid();
23     } else ready = 1;
24
25     while (1) {
26         if (ready){
27             ready = 0;
28             sleep(1);
29             write(1,str,4);
30             kill(other,SIGUSR1);
31         }
32         pause();
33     }
34 }

```

Flaga SA_SIGINFO

If SA_SIGINFO is set and the signal is caught, the signal-catching function shall be entered as:

```
void func(int signo, siginfo_t *info, void *context);
```

info the reason why the signal was generated;

context the receiving thread's context that was interrupted when the signal was delivered.

Syscallle przerwane sygnałami.

read - przypomnienie

If a `read()` is interrupted by a signal before it reads any data, it shall return `-1` with `errno` set to `[EINTR]`.

If a `read()` is interrupted by a signal after it has successfully read some data, it shall return the number of bytes read.

write - przypomnienie

If `write()` is interrupted by a signal before it writes any data, it shall return `-1` with `errno` set to `[EINTR]`.

If `write()` is interrupted by a signal after it successfully writes some data, it shall return the number of bytes written.

```

1 #include <unistd.h>
2 #include <sys/types.h>
3 #include <signal.h>
4 #include <errno.h>
5
6 #define BSIZE 100
7 void handler(int sig_nb){}
8
9 int main(){
10     int n,k,w;
11     char buf[BSIZE];
12     pid_t parent, child;
13     struct sigaction act;
14
15     act.sa_handler = handler;
16     act.sa_flags = 0;
17     sigemptyset(&act.sa_mask);
18     sigaction(SIGUSR1, &act, NULL);
19
20     if (!(child = fork())){
21         parent = getppid();
22         while (1) kill(parent, SIGUSR1);
23         exit(1);
24     }
25
26     while (n = read(0, buf, BSIZE)){
27         if ((n<0) && (errno!=EINTR)) break;
28         k = 0;
29         while(k<n){
30             w = write(1, buf+k, n-k);
31             if ((w<0) && (errno!=EINTR)) goto end;
32             if (w>0) k+=w;
33         }
34     }
35 end:    kill(child, SIGTERM);
36 }

```

SA_RESTART

If set, and a function specified as interruptible is interrupted by this signal, the function shall restart and shall not fail with [EINTR] unless otherwise specified.

Przykłady

read, write, open, waitpid, fcntl (F_SETLKW)

```

1 #include <unistd.h>
2 #include <sys/types.h>
3 #include <signal.h>
4 #include <errno.h>
5
6 #define BSIZE 100
7 void handler(int sig_nb){}
8
9 int main(){
10     int n,k,w;
11     char buf[BSIZE];
12     pid_t parent, child;
13     struct sigaction act;
14
15     act.sa_handler = handler;
16     act.sa_flags = SA_RESTART;
17     sigemptyset(&act.sa_mask);
18     sigaction(SIGUSR1, &act, NULL);
19
20     if (!(child = fork())){
21         parent = getppid();
22         while (1) kill(parent, SIGUSR1);
23         exit(1);
24     }
25
26     while ((n = read(0, buf, BSIZE))>0){
27         k = 0;
28         while(k<n){
29             w = write(1, buf+k, n-k);
30             if (w<0) goto end;
31             k += w;
32         }
33     }
34 end:    kill(child, SIGTERM);
35 }

```

UWAGA na errno!

```
1 #include <unistd.h>
2 #include <sys/types.h>
3 #include <signal.h>
4 #include <errno.h>
5
6 #define BSIZE 100
7 void handler(int sig_nb){ errno = 0;}
8
9 int main(){
10     int n,k,w;
11     char buf[BSIZE];
12     pid_t parent, child;
13     struct sigaction act;
14
15     act.sa_handler= handler;
16     act.sa_flags=0;
17     sigemptyset(&act.sa_mask);
18     sigaction(SIGUSR1, &act, NULL);
19
20     if (!(child = fork())){
21         parent = getppid();
22         while (1) kill(parent, SIGUSR1);
23         exit(1);
24     }
25
26     while (n = read(0,buf, BSIZE)){
27         if ((n<0) && (errno != EINTR)) break; else n = 0;
28         write(1, buf, n);    // nie dbamy o przerwane write'y
29     }
30 end:  kill(child, SIGTERM);
31 }
```

Flaga SA_NOCLDSTOP

SIGCHLD

Child process terminated, stopped, or continued.

SA_NOCLDSTOP

Do not generate SIGCHLD when children stop or stopped children continue.

Uwaga

If a process sets the action for the SIGCHLD signal to SIG_IGN, the behavior is unspecified.

Normalne sygnały NIE są kolejkowane!

```
1 #include <unistd.h>
2 #include <stdlib.h>
3 #include <stdio.h>
4 #include <sys/types.h>
5 #include <signal.h>
6
7 volatile int s=0;
8 void handler(int sig_nb){ s++; }
9
10 int main(){
11     int n,k,w;
12     pid_t parent;
13     struct sigaction act;
14
15     act.sa_handler = handler;
16     act.sa_flags = 0;
17     sigemptyset(&act.sa_mask);
18     sigaction(SIGUSR1, &act, NULL);
19
20     if (!fork()){
21         parent= getppid();
22         for (n=0; n<10; n++) kill(parent, SIGUSR1);
23         exit(0);
24     }
25
26     pause();
27     while (s-- >0) {
28         printf("received\n");
29         sleep(3);
30     }
31 }
```

```

1 #include <unistd.h>
2 #include <stdio.h>
3 #include <sys/types.h>
4 #include <sys/wait.h>
5 #include <signal.h>
6
7 #define CHILDREN 10
8 volatile int z=CHILDREN;
9 void handler(int sig_nb){
10     pid_t child;
11     do{
12         child = waitpid(-1,NULL,WNOHANG);
13         if (child>0) z--;
14     } while (child >0);
15     sleep(1);
16 }
17
18 int main(){
19     int n;
20     struct sigaction act;
21     act.sa_handler = handler;
22     act.sa_flags = 0;
23     sigemptyset(&act.sa_mask);
24     sigaction(SIGCHLD,&act,NULL);
25
26     for (n=0;n<CHILDREN;n++){
27         if (!fork()) {
28             sleep(n);
29             return 0;
30         }
31
32         while (z>0) {
33             printf("%d\children/zombies_left.\n",z);
34             sleep(1);
35         }
36         printf("No_more_zombies.\n");
37 }

```

SIG_IGN dla SIGCHLD

```
1 #include <unistd.h>
2 #include <signal.h>
3
4 #define CHILDREN 10
5
6 int main(){
7     int n;
8     struct sigaction act;
9     act.sa_handler = SIG_IGN;
10    act.sa_flags = 0;
11    sigemptyset(&act.sa_mask);
12    sigaction(SIGCHLD,&act,NULL);
13
14    for (n=0;n<CHILDREN;n++)
15        if (!fork()) {
16            return 0;
17        }
18    sleep(10);
19 }
```

LINUX

Ignoring SIGCHLD can be used to prevent the creation of zombies.

async-signal-safe functions

safe

_exit, close, kill, read, write, ...

unsafe

malloc, exit, printf ...

alarm() function → SIGALRM signal

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

```
unsigned alarm(unsigned seconds);
```

Co może pójść źle w poniższym programie?

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <signal.h>
4
5 #define TIME 5
6 void handler(int sig_nb){ }
7
8 int main(){
9     int n;
10    struct sigaction act;
11
12    act.sa_handler = handler;
13    act.sa_flags = 0;
14    sigemptyset(&act.sa_mask);
15    sigaction(SIGALRM, &act, NULL);
16
17    alarm(TIME);
18    pause();
19    printf("No_time_..._no_time_to_lose.\n");
20 }
```

```
#include <signal.h>

int sigprocmask(int how, const sigset_t *restrict set,
                sigset_t *restrict oset);
```

how values

SIG_BLOCK The resulting set shall be the union of the current set and the signal set pointed to by set.

SIG_SETMASK The resulting set shall be the signal set pointed to by set.

SIG_UNBLOCK The resulting set shall be the intersection of the current set and the complement of the signal set pointed to by set.

alarm()

Trochę lepsze rozwiązanie.

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <signal.h>
4
5 #define TIME 5
6 void handler(int sig_nb){ }
7
8 int main(){
9     int n;
10    struct sigaction act;
11    sigset_t mask;
12
13
14    act.sa_handler = handler;
15    act.sa_flags = 0;
16    sigemptyset(&act.sa_mask);
17    sigaction(SIGALRM, &act, NULL);
18    sigaction(SIGINT, &act, NULL);
19
20    sigfillset(&mask);
21    sigdelset(&mask, SIGALRM);
22
23    sigprocmask(SIG_SETMASK, &mask, NULL);
24    alarm(TIME);
25    pause();
26    printf("No_time..._no_time_to_lose.\n");
27 }
```

sigsuspend

```
#include <signal.h>
```

```
int sigsuspend(const sigset_t *sigmask);
```

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <signal.h>
4 #define TIME 5
5 void handler(int sig_nb){ }
6
7 int main(){
8     struct sigaction act;
9     sigset_t mask;
10
11     act.sa_handler = handler;
12     act.sa_flags = 0;
13     sigemptyset(&act.sa_mask);
14     sigaction(SIGALRM, &act, NULL);
15     sigaction(SIGINT, &act, NULL);
16
17     sigemptyset(&mask);
18     sigaddset(&mask, SIGALRM);
19     sigprocmask(SIG_BLOCK, &mask, NULL);
20
21     sigfillset(&mask);
22     sigdelset(&mask, SIGALRM);
23     alarm(TIME);
24     sigsuspend(&mask);
25     printf("No time to lose.\n");
26 }
```



```
#include <signal.h>

int sigpending(sigset_t *set);
```

The `sigpending()` function shall store, in the location referenced by the `set` argument, the set of signals that are blocked from delivery to the calling thread and that are pending on the process or the calling thread.

read with timeout - prawie poprawne

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <signal.h>
4
5 volatile int time_is_up;
6 void handler(int sig_nb){ time_is_up = 1;}
7
8 int tread(char * buf, int n, int timeout){
9     int r;
10    struct sigaction act, oact;
11    sigset_t mask, omask;
12
13    act.sa_handler = handler;
14    act.sa_flags = 0;
15    sigemptyset(&act.sa_mask);
16    sigaction(SIGALRM,&act,&oact);
17
18    time_is_up = 0;
19    alarm(timeout);
20    do r = read(0,buf,n);
21    while ((!time_is_up) && (r<0));
22    alarm(0);
23
24    sigaction(SIGALRM,&oact,NULL);
25    if (r<0) return 0;
26    return r;
27 }
28 int main(){
29     char buf[100];
30     int n = tread(buf,100,5);
31     write(1,buf,n);
32 }
```

select

```
#include <sys/select.h>

int select(int nfds, fd_set *restrict readfds,
           fd_set *restrict writefds, fd_set *restrict errorfds,
           struct timeval *restrict timeout);

void FD_CLR(int fd, fd_set *fdset);
int FD_ISSET(int fd, fd_set *fdset);
void FD_SET(int fd, fd_set *fdset);
void FD_ZERO(fd_set *fdset);
```

Upon successful completion, the `pselect()` or `select()` function shall modify the objects pointed to by the `readfds`, `writefds`, and `errorfds` arguments to indicate which file descriptors are ready for reading, ready for writing, or have an error condition pending, respectively, and shall return the total number of ready descriptors in all the output sets.

```
struct timeval {
    long    tv_sec;          /* seconds */
    long    tv_usec;        /* microseconds */
};
```

read with timeout - poprawne

```
1 #include <unistd.h>
2 #include <stdio.h>
3 #include <errno.h>
4 #include <sys/select.h>
5 #include <sys/time.h>
6
7 int tread(char * buf, int n, int timeout){
8     int r;
9     fd_set rfd;
10    struct timeval timeout_s;
11
12    FD_ZERO(&rfd);
13    FD_SET(0,&rfd);
14
15    timeout_s.tv_sec = timeout;
16    timeout_s.tv_usec = 0;
17
18    do r= select(1,&rfd,NULL,NULL,&timeout_s);
19    while ((r<0) && (errno==EINTR));
20
21    if (r<=0) return 0;
22
23    return (read(0,buf,n));
24 }
25
26 int main(){
27     char buf[100];
28     int n = tread(buf,100,5);
29     write(1,buf,n);
30 }
```

```
#include <sys/select.h>

int pselect(int nfd, fd_set *restrict readfds,
            fd_set *restrict writefds, fd_set *restrict errorfds,
            const struct timespec *restrict timeout,
            const sigset_t *restrict sigmask);
```

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- **POSIX - remanent**

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Procesy

`getpriority, setpriority` - get and set scheduling priority
`setsid, getpgrp` - create process group, get process group id
`setuid, setgid` - set user or group ID's
`brk, sbrk` - change data segment size

File System

```
access  - determine accessibility of file
chmod   - change mode of file
chown   - change owner and group of a file
link    - make a hard link to a file
mkdir   - make a directory file
mount, umount - mount or umount a file system
rename  - change the name of a file
rmdir   - remove a directory file
stat, lstat, fstat - get file status
sync, fsync - update dirty buffers and super-block
unlink  - remove directory entry
umask   - set file creation mode mask
utime   - set file times
```


Info

```
gettimeofday - get date and time  
getuid, geteuid - get user identity  
time, stime - get/set date and time  
times - get process times  
uname - get system info
```

Inne

- `chroot` - change root directory
- `ptrace` - process trace
- `reboot` - close down the system or reboot
- `mmap` - request memory mapping

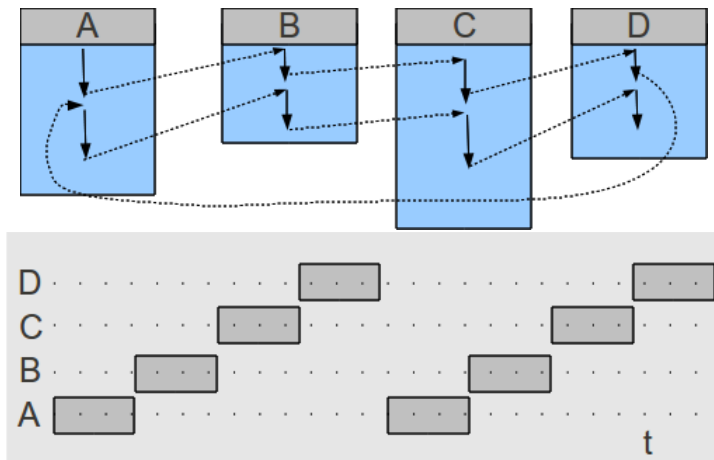
Outline

- 1 Zasady
- 2 POSIX
 - Wstęp
 - POSIX - standard
 - POSIX - procesy
 - POSIX - pliki
 - POSIX - sygnały
 - POSIX - remanent
- 3 Sequential Processes
 - Multiprogramming
 - IPC - InterProcess Communication
- 4 Architektura
 - Monolithic kernel
 - Micro kernel
 - Tanenbaum - Torvalds debate
- 5 Scheduling
 - Batch systems

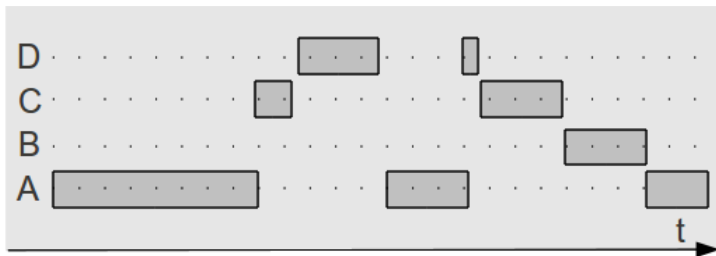
Multiprogramming (multitasking)

“Wirtualne procesory”

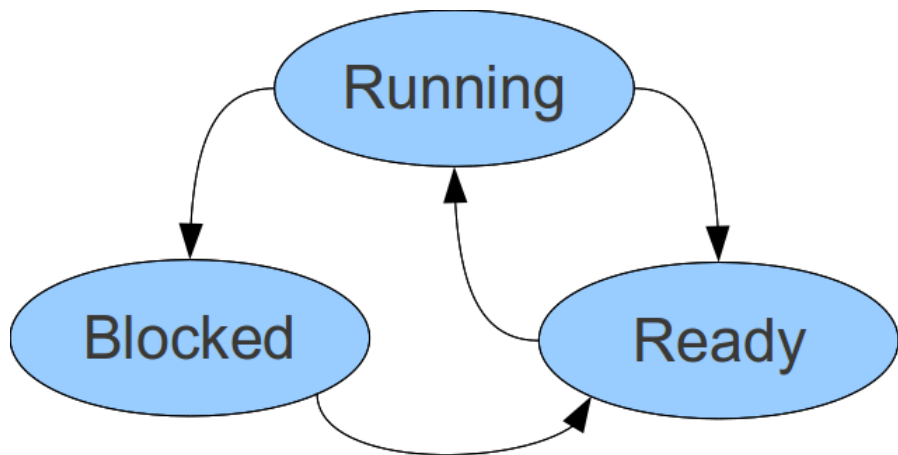
Każdy proces pracuje jak gdyby miał kopię procesora (z ograniczoną funkcjonalnością) dla siebie.



Multiprogramming (multitasking)



nierównomierny postęp w czasie → nie wolno robić założeń dotyczących czasu rzeczywistego.



Opis procesu

Kernel

- Registers
- Program counter
- Program status word
- Stack Pointer
- Process state
- Current scheduling priority
- Maximum scheduling priority
- Scheduling ticks left
- Quantum size
- CPU time used
- Message queue pointers
- Pending signal bits
- Various flag bits
- Process name

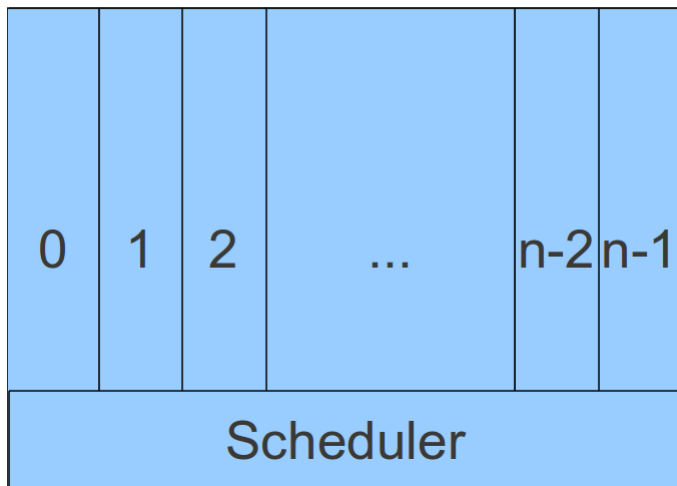
Process management

- Pointer to text segment
- Pointer to data segment
- Pointer to bss segment
- Exit status
- Signal status
- **Process ID**
- Parent Process
- Process group
- Children's CPU time
- Real UID
- Effective UID
- Real GID
- Effective GID
- File info for sharing text
- Bitmaps for signals
- Various flag bits
- Process name

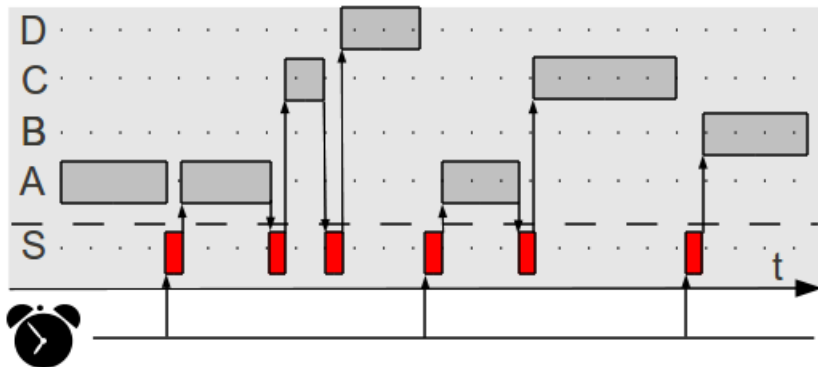
File management

- UMASK mask
- Root directory
- Working directory
- File descriptors
- Real UID
- Effective UID
- Real GID
- Effective GID
- Controlling tty
- Save area for read/write
- System call parameters
- Various flag bits

Procesy



Przerwania zegarowe.



MINIX - 60 przerwań na sekundę.

Pamięć - zmiany kontekstu

Kernel

- **Registers**
- **Program counter**
- **Program status word**
- **Stack Pointer**
- **Process state**
- Current scheduling priority
- Maximum scheduling priority
- Scheduling ticks left
- Quantum size
- CPU time used
- Message queue pointers
- Pending signal bits
- Various flag bits
- Process name

Process management

- **Pointer to text segment**
- **Pointer to data segment**
- **Pointer to bss segment**
- Exit status
- Signal status
- Process ID
- Parent Process
- Process group
- Children's CPU time
- Real UID
- Effective UID
- Real GID
- Effective GID
- File info for sharing text
- Bitmaps for signals
- Various flag bits
- Process name

File management

- UMASK mask
- Root directory
- Working directory
- File descriptors
- Real UID
- Effective UID
- Real GID
- Effective GID
- Controlling tty
- Save area for read/write
- System call parameters
- Various flag bits

Wiele “lekkich procesów” w tej samej przestrzeni adresowej.

Zmiana na inny wątek w tym samym procesie

- Registers
- Program counter
- Stack pointer
- State

wspólna pamięć → konieczna współpraca (synchronizacja).

systemowe / użytkownika

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

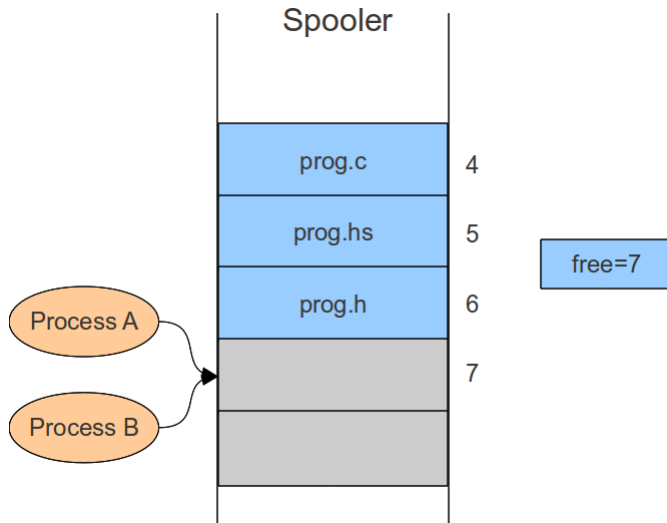
- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Współdzielona pamięć.

Race condition



Sekcje krytyczne

Sekcja krytyczna

Fragment programu z odwołaniami do współdzielonej pamięci.

Cel

Zawsze co najwyżej jeden proces/wątek w sekcji krytycznej.

Dodatkowo wymagamy:

- żaden proces działający poza sekcją krytyczną nie może blokować innego procesu,
- żaden proces nie powinien czekać w nieskończoność na wejście do sekcji krytycznej,
- mechanizm powinien działać niezależnie od szybkości i liczby procesorów.

Wyłączanie przerw.

CLI - processor instruction

Clear Interrupt Flag - Clears the interrupt flag (IF) in the rFLAGS register to zero, thereby masking external interrupts received on the INTR input.

Wady

- konieczne uprawnienia do blokowania przerw
- nieskuteczne w systemach wieloprocessorowych

Busy waiting - spin lock

```
while (TRUE) {  
    while (lock!=0);  
    lock = 1;  
    critical_region();  
    lock = 0;  
    noncritical_region();  
}
```

```
while (TRUE) {  
    while (lock!=0);  
    lock = 1;  
    critical_region();  
    lock = 0;  
    noncritical_region();  
}
```

Busy waiting - spin lock + turns

```
while (TRUE) {  
    while (turn!=0);  
    critical_region();  
    turn = 1;  
    noncritical_region();  
}
```

```
while (TRUE) {  
    while (turn!=1);  
    critical_region();  
    turn = 0;  
    noncritical_region();  
}
```

*Dude, [...] this determines who enters the next round robin. Am I wrong?
Am I wrong?*

Rozwiązanie Peterson'a

```
int turn;    //shared
int interested[2]; //shared

void enter_region(int process){
    int other;

    other = 1 - process;
    interested[process] = TRUE;
    turn = process;
    while ((turn==process) && (interested[other]==TRUE));
}

void leave_region(int process){
    interested[process] = FALSE;
}
```

Test and Set Lock instruction

TSL

TSL reg, lock - wczytuje zawartość pamięci lock do rejestru reg oraz zapisuje niezerową wartość pod adresem lock.

```
enter_region:
```

```
    TSL eax, lock
```

```
    CMP eax, 0
```

```
    JNE enter_region
```

```
    RET
```

```
leave_region:
```

```
    MOV lock, 0
```

```
    RET
```

x64

CMPXCHG reg/mem32, reg32

Compare EAX register with a 32-bit register or memory location. If equal, copy the second operand to the first operand. Otherwise, copy the first operand to EAX.

Busy Waiting \rightarrow Priority Inversion Problem

Sleep/Wakeup

```
#define N 100
int count=0;
void producer(void){
    int item;
    while (TRUE){
        item = produce_item();
        if (count==N) sleep();
        insert_item(item);
        count++;
        if (count==1) wakeup(consumer);
    }
}
void consumer(void){
    int item;
    while (TRUE){
        if (count==0) sleep();
        item = remove_item();
        count--;
        if (count==N-1) wakeup(producer);
        consume_item(item);
    }
}
```

Semafor

```
#define N 100  
semaphore mutex = 1;  
semaphore empty = N;  
semaphore full = 0;
```

```
void producer(void){  
    int item;  
    while (TRUE){  
        item = produce_item();  
        down(&empty);  
        down(&mutex);  
        insert_item(item);  
        up(&mutex);  
        up(&full);  
    }  
}
```

```
void consumer(void){  
    int item;  
    while (TRUE){  
        down(&full);  
        down(&mutex);  
        item = remove_item();  
        up(&mutex);  
        up(&empty);  
        consume_item(item);  
    }  
}
```

Semafory - UWAGA!

```
#define N 100  
semaphore mutex = 1;  
semaphore empty = N;  
semaphore full = 0;
```

```
void producer(void){  
    int item;  
    while (TRUE){  
        item = produce_item();  
        down(&mutex); //zmiana  
        down(&empty); //kolejności  
        insert_item(item);  
        up(&mutex);  
        up(&full);  
    }  
}
```

```
void consumer(void){  
    int item;  
    while (TRUE){  
        down(&full);  
        down(&mutex);  
        item = remove_item();  
        up(&mutex);  
        up(&empty);  
        consume_item(item);  
    }  
}
```


Monitory

```
monitor ProducerConsumer
    condition full, empty;
    int count;

    void insert(int item)
        if count==N then wait(full);
        insert_item(item);
        count++;
        if count==1 then signal(empty);

    int remove()
        if count==0 then wait(empty);
        remove = remove_item();
        count--;
        if (count= N-1) then signal(full);
```

Bez współdzielonej pamięci.

Message Passing (buffered)

```
#define N 100
```

```
void producer(void){  
    int item;  
    message m;  
  
    while (TRUE){  
        item = produce_item();  
        receive(consumer, &m);  
        build_message(&m, item);  
        send(consumer, &m);  
    }  
}
```

```
void consumer(void){  
    int item,i;  
    message m;  
  
    for (i=0; i<N; i++)  
        send(producer,&m);  
    while (TRUE){  
        receive(producer,&m);  
        item= extract_item(&m);  
        send(producer,&m);  
        consume_item(item);  
    }  
}
```

Message Passing - rendezvous

```
#define N 100
```

```
void producer(void){
    int item;
    message m;

    while (TRUE){
        item = produce_item();
        receive(consumer, &m);
        build_message(&m, item);
        send(consumer, &m);
    }
}
```

```
void consumer(void){
    int item,i;
    message m;

    send(producer,&m);
    while (TRUE){
        receive(producer,&m);
        item= extract_item(&m);
        send(producer,&m);
        consume_item(item);
    }
}
```

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Monolithic kernel - Linux

Procesy użytkownika pracują w dwóch trybach

- tryb użytkownika - user mode
- tryb jądra - kernel mode

Przejście do trybu jądra (system calls) poprzez specjalne bramki - w x86 przerwania programowe (`int 80h`) i specjalne instrukcje (`sysenter`, `syscall`).

Skutek: wiele procesów (użytkownika) może pracować równocześnie w trybie jądra.

Obsługa przerwań

- Scheduling.
- Drivery urządzeń w trybie jądra!

Do tego wątki systemu odpowiedzialne za takie rzeczy jak

- zapisywanie buforów (`pdflush`),
- kolejkę (drobnych) zadań jądra (`keventd`),
- zwalnianie/swapowanie stron pamięci (`kswapd`),

Outline

- 1 Zasady
- 2 POSIX
 - Wstęp
 - POSIX - standard
 - POSIX - procesy
 - POSIX - pliki
 - POSIX - sygnały
 - POSIX - remanent
- 3 Sequential Processes
 - Multiprogramming
 - IPC - InterProcess Communication
- 4 Architektura
 - Monolithic kernel
 - **Micro kernel**
 - Tanenbaum - Torvalds debate
- 5 Scheduling
 - Batch systems

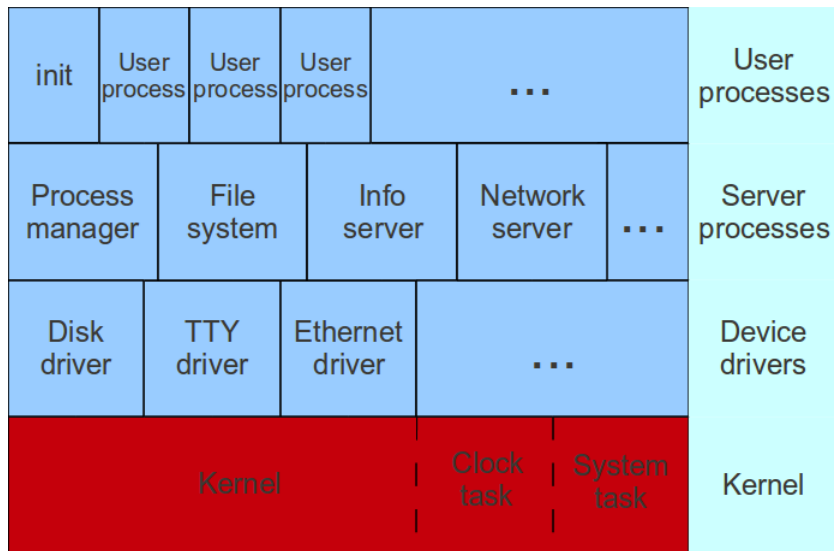
- Minimalizacja kodu pracującego w trybie uprzywilejowanym (kernel mode).
- Funkcjonalności systemu są zaimplementowane jako serwery (process manager, file system, network server).
- Komunikacja z serwerami - message passing.

Odpowiedzialność jądra

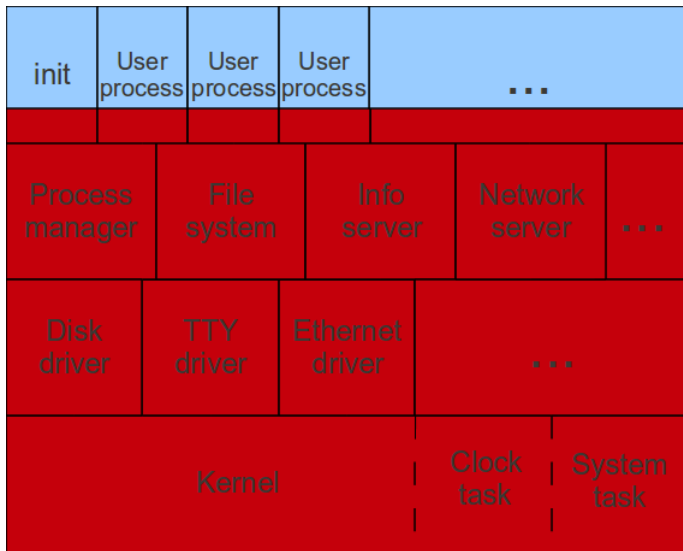
- multiprogramming,
- message - passing,
- mechanizmy dla obsługi urządzeń - drivery mogą pracować w trybie użytkownika.

Procesy użytkownika NIGDY nie pracują w trybie jądra.

Micro kernel - Minix



Monolithic kernel analogue



Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Tanenbaum - Torvalds debate

The Tanenbaum-Torvalds Debate, Open Sources: Voices from the Open Source Revolution, O'Reilly 1999

ast: ... LINUX is a monolithic style system. This is a giant step back into the 1970s. That is like taking an existing, working C program and rewriting it in BASIC. To me, writing a monolithic system in 1991 is a truly poor idea...

torvalds: ... True, linux is monolithic, and I agree that microkernels are nicer. With a less argumentative subject, I'd probably have agreed with most of what you said. From a theoretical (and aesthetical) standpoint linux loses. If the GNU kernel had been ready last spring, I'd not have bothered to even start my project: the fact is that it wasn't and still isn't. Linux wins heavily on points of being available now...

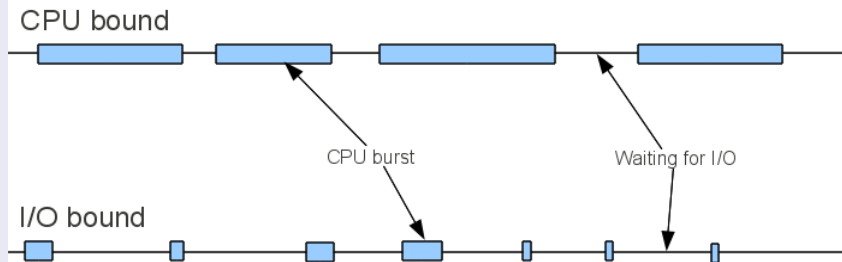
ast: ... Of course 5 years from now that will be different, but 5 years from now everyone will be running free GNU on their 200 MIPS, 64M SPARCstation-5...

PART II

Tanenbaum-Torvalds Debate: Part II, Tanenbaum's web page

Zachowanie procesów.

CPU-bound / I/O-bound



Kategorie procesów/schedulerów

- interaktywne
- wsadowe
- procesy czasu rzeczywistego

Kiedy scheduler wkracza do akcji?

Konieczna decyzja

- działający proces się blokuje
- działający proces się kończy

Możliwa zmiana procesu

- utworzenie nowego procesu
- przerwanie I/O
- przerwanie zegarowe (nonpreemptive / preemptive)

Wszystkie rodzaje

- Fairness - podobne procesy są podobnie traktowane
- Policy enforcement - uwzględnienie specyfiki procesów (priorytety itp.)
- Balance - wszystkie części systemu równo obciążone

Systemy wsadowe

- Throughput - maksymalizacja zadań na godzinę
- Turnaround time - minimalizacja czasu pomiędzy przyjściem a zakończeniem zadania (najlepiej znormalizowana wielkością zadania)
- CPU utilization - maksymalizacja obciążenia procesora

Systemy interaktywne

- Response time - szybka reakcja na bodźce
- Proportionality - opóźnienie reakcji systemu proporcjonalne do czasu trwania zadania

Systemy czasu rzeczywistego

- Meeting deadlines
- Predictability

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Batch systems: First-Come First-Served

- Procesy obsługiwane w kolejności przyścia.
- Zablokowane procesy trafiają na koniec kolejki kiedy przechodzą w stan gotowości.

Wady

- Słabe wskaźniki: Turnaround time, CPU utilization, Throughput.
- Wiele procesów może pracować współbieżnie - duże obciążenie pamięci.
- Jeśli wszystkie procesy są zablokowane na I/O to dodawane są nowe - czas oczekiwania na I/O się wydłuża.

Batch systems: Shortest Job First

- Zakładamy, że potrafimy oszacować czas wykonania zadania (totalny czas, nie tylko CPU).
- Wybieramy proces o najkrótszym czasie.

Zalety

- Off line: Optymalny z punktu widzenia Turnaround time

Wady

- On line: nie optymalny. Zadania A,B,C,D,E o wagach 2,4,1,1,1 przychodzą w chwilach 0,0,3,3,3. Średnia 4.6. Jeśli wykonać B,C,D,E,A to średnia 4.4.
- Nie ma mechanizmów równoważenia procesów (CPU-bound/I/O-bound)

Batch systems: Shortest Remaining Time Next

- Wybieramy proces o najkrótszym czasie pozostałym do zakończenia.
- Jeśli pojawia się nowy proces to dopuszczalne są wywłaszczenia.

Zalety

- Krótkie zadania mają krótki czas realizacji.
- Można uwzględnić blokowanie.
- Optymalny ze względu na Turnaround time w wersji on-line. Po zaniedbaniu kosztu zmiany kontekstu, schedulera, i ewentualnej korzyści z równoważenia.

Wady

- Nie ma mechanizmów równoważenia procesów (CPU-bound/I/O-bound).
- Pesymistyczny Turnaround-time może być bardzo duży.

Batch systems: Three-Level Scheduling

Admission scheduler

- Wybiera zadania które zaczną być wykonywane współbieżnie.
- Równoważenie zadań CPU-bound / I/O-bound.
- Minimalizacja normalizowanego Turnaround time.

Memory scheduler

- Działa jeśli procesy nie mieszczą się w pamięci.
- Wybiera procesy, które zostaną swapowane na dysk.
- Równoważenie zadań CPU-bound / I/O-bound.

CPU scheduler

- Wybiera następne zadanie do wykonania.
- Równoważenie zadań CPU-bound / I/O-bound.

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

Interactive systems: Round robin

- Każdy proces dostaje kwant czasu (quantum).
- Procesy ustawione są w kolejce.
- Proces może obliczać tak długo jak trwa jego kwant.
- Proces który zużyje kwant, jest przerywany, dostaje nowe kwant i trafia na koniec kolejki.
- Proces, który się blokuje przed wykorzystaniem quantu, po przejściu w stan Ready, trafia na początek kolejki z pozostałym czasem.

Ważny parametr: długość kwantu

- kwant za duży - system może mieć długi czas reakcji.
- kwant za mały - duży odsetek pracy procesora tracony na zmiany kontekstu.
- w praktyce 20-50 msec.

Wady

- Wszystkie procesy są równo traktowane.

Uogólnienie Round-robin. Wiele możliwości:

- Kwanty mogą zależeć od priorytetu.
- Procesy o wyższym priorytecie mają pierwszeństwo.
- Dla każdego priorytetu lista procesów - Round - robin.
- Dynamiczne/statyczne przydzielanie priorytetu (np. jeśli proces zużył $1/f$ swojego kwantu to dostaje priorytet f).
- Narastające kwanty w niższych priorytetach.

Ważny parametr: długość kwantu

- kwant za duży - system może mieć długi czas reakcji.
- kwant za mały - duży odsetek pracy procesora tracony na zmiany kontekstu.
- w praktyce 20-50 msec.

Guaranteed scheduling / Fair Share scheduling

- Każdy proces/użytkownik ma zagwarantowane $x\%$ procesora.
- Scheduler wybiera procesy tak aby wypełnić zobowiązanie.

Lottery scheduling

- Każdy proces dostaje kilka żetonów z puli.
- Scheduler wybiera losowo żeton i przydziela procesor właścicielowi.
- Współpracujące procesy mogą sobie przekazywać żetony.

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

- Batch systems

- hard real time / soft real time.
- periodyczne, aperiodyczne zdarzenia.
- szeregowanie statyczne/dynamiczne.

Schedulable system

C_i - czas obsługi zdarzenia i

P_i - częstotliwość zdarzenia i

$$\sum_i \frac{C_i}{P_i} \leq 1$$

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

5 Scheduling

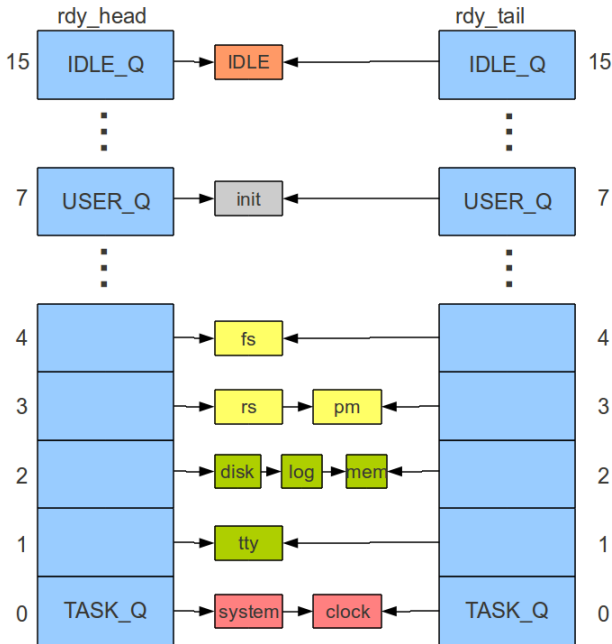
- Batch systems

Multilevel/Priority scheduling

- Dla każdego priorytetu kolejka.
- W obrębie kolejki round-robin.
- W kolejkach tylko procesy w stanie ready.
- Jeśli proces przechodzi w stan ready to:
 - Trafia na początek kolejki, jeśli nie wykorzystał całego czasu (kwantu) przed zablokowaniem,
 - wpp - trafia na koniec kolejki z nowym kwantem.

Wybór kolejnego procesu:

Weź pierwszy proces z niepustej kolejki o największym priorytecie.



+1 Proces zużył cały kwant i dwa razy z rzędu został mu przydzielony procesor.

-1 wpp.

- Nie dotyczy zadań jądra (kernel tasks) oraz procesu IDLE.
- Priorytet nie może być większy niż $IDLE_Q - 1$.
- Priorytet nie może być mniejszy niż maksymalny priorytet zależny od rodzaju procesu.


```

PUBLIC int _syscall(int who, int syscallnr, message * msgptr){
    int status;                                /*lib/other/syscall.c*/

    msgptr->m_type = syscallnr;
    status = _sendrec(who, msgptr);
    if (status != 0) {
        /* 'sendrec' itself failed. */
        msgptr->m_type = status;
    }
    if (msgptr->m_type < 0) {
        errno = -msgptr->m_type;
        return(-1);
    }
    return(msgptr->m_type);
}

__sendrec:                                     ! lib/i386/rts/_ipc.s
    push ebp
    mov ebp, esp
    push ebx
    mov eax, SRC_DST(ebp)      ! eax = dest-src
    mov ebx, MESSAGE(ebp)     ! ebx = message pointer
    mov ecx, SENDREC          ! _sendrec(srcdest, ptr)
    int SYSVEC                 ! trap to the kernel
    pop ebx
    pop ebp
    ret

```

```

_s_call:                ! kernel/mpx386.s
_p_s_call:
    cld                  ! set direction flag to a known value
    sub esp, 6*4         ! skip RETADR, eax, ecx, edx, ebx, est
    push ebp             ! stack already points into proc table
    push esi
    push edi
    o16 push ds
    o16 push es
    o16 push fs
    o16 push gs
    mov dx, ss
    mov ds, dx
    mov es, dx
    incb (_k_reenter)
    mov esi, esp         ! assumes P_STACKBASE == 0
    mov esp, k_stktop
    xor ebp, ebp         ! for stacktrace
                          ! end of inline save
                          ! now set up parameters for sys_call()
    push ebx             ! pointer to user message
    push eax             ! src/dest
    push ecx             ! SEND/RECEIVE/BOTH
    call _sys_call       ! sys_call(function, src_dest, m_ptr)
                          ! caller is now explicitly in proc_ptr
    mov AXREG(esi), eax  ! sys_call MUST PRESERVE si

```

```

! Fall into code to restart proc/task running.
_restart:                                ! kernel/mpx386.s

! Restart the current process or the next process if it is set.
    cmp (_next_ptr), 0                  ! see if another process is scheduled
    jz 0f
    mov  eax, (_next_ptr)
    mov  (_proc_ptr), eax              ! schedule new process
    mov  (_next_ptr), 0
0:  mov  esp, (_proc_ptr)              ! will assume P_STACKBASE == 0
    lldt P_LDT_SEL(esp)                ! enable process' segment descriptors
    lea  eax, P_STACKTOP(esp)          ! arrange for next interrupt
    mov  (_tss+TSS3_S_SP0), eax        ! to save state in process table
restart1:
    decb (_k_reenter)
    o16 pop gs
    o16 pop fs
    o16 pop es
    o16 pop ds
    popad
    add  esp, 4                        ! skip return adr
    iretd                             ! continue process

```

```

1 PUBLIC int sys_call(call_nr, src_dst, m_ptr)
2 int call_nr;          /* system call number and flags */
3 int src_dst;          /* src to receive from or dst to send to */
4 message *m_ptr;       /* pointer to message in the caller's space */
5 {
6     /.../
7     int function = call_nr & SYSCALL_FUNC; /* get system call function
8     /...// check several conditions
9     switch(function) {
10    case SENDREC:
11        /* A flag is set so that notifications cannot interrupt SENDREC
12        priv(caller_ptr)->s_flags |= SENDREC_BUSY;
13        /* fall through */
14    case SEND:
15        result = mini_send(caller_ptr, src_dst, m_ptr, flags);
16        if (function == SEND || result != OK) {
17            break; /* done, or SEND failed */
18        } /* fall through for SENDREC */
19    case RECEIVE:
20        if (function == RECEIVE)
21            priv(caller_ptr)->s_flags &= ~SENDREC_BUSY;
22        result = mini_receive(caller_ptr, src_dst, m_ptr, flags);
23        break;
24    /...// NOTIFY, ECHO, default
25    }
26    return(result);
27 }

```

```

1 PRIVATE int mini_send(caller_ptr, dst, m_ptr, flags)
2 register struct proc *caller_ptr; /* who is trying to send a message? */
3 int dst; /* to whom is message being sent? */
4 message *m_ptr; /* pointer to message buffer */
5 unsigned flags; /* system call flags */
6 {
7 /* Send a message from 'caller_ptr' to 'dst'. */
8 register struct proc *dst_ptr = proc_addr(dst);
9 register struct proc **xpp;
10 register struct proc *xp;
11 /*...*/ Check for deadlock by 'caller_ptr' and 'dst' sending to each other. */
12
13 /* Check if 'dst' is blocked waiting for this message. The destination's
14 * SENDING flag may be set when its SENDREC call blocked while sending.
15 */
16 if ( (dst_ptr->p_rts_flags & (RECEIVING | SENDING)) == RECEIVING &&
17 (dst_ptr->p_getfrom == ANY || dst_ptr->p_getfrom == caller_ptr->p_nr)) {
18 /* Destination is indeed waiting for this message. */
19 CopyMess(caller_ptr->p_nr, caller_ptr, m_ptr, dst_ptr, dst_ptr->p_messbuf);
20 if ((dst_ptr->p_rts_flags & ~RECEIVING) == 0) enqueue(dst_ptr);
21 } else if ( ! (flags & NON_BLOCKING)) {
22 /* Destination is not waiting. Block and dequeue caller. */
23 caller_ptr->p_messbuf = m_ptr;
24 if (caller_ptr->p_rts_flags == 0) dequeue(caller_ptr);
25 caller_ptr->p_rts_flags |= SENDING;
26 caller_ptr->p_sendto = dst;
27 /* Process is now blocked. Put in on the destination's queue. */
28 xpp = &dst_ptr->p_caller_q; /* find end of list */
29 while (*xpp != NIL_PROC) xpp = &(*xpp)->p_q_link;
30 *xpp = caller_ptr; /* add caller to end */
31 caller_ptr->p_q_link = NIL_PROC; /* mark new end of list */
32 } else {
33 return(ENOTREADY);
34 }
35 return(OK);
36 }

```

```

1 PRIVATE void enqueue(rp)
2 register struct proc *rp; /* this process is now runnable */
3 {
4     int q;          /* scheduling queue to use */
5     int front;       /* add to front or back */
6
7     /* Determine where to insert to process. */
8     sched(rp, &q, &front);
9
10    /* Now add the process to the queue. */
11    if (rdy_head[q] == NIL_PROC) { /* add to empty queue */
12        rdy_head[q] = rdy_tail[q] = rp; /* create a new queue */
13        rp->p_nextready = NIL_PROC; /* mark new end */
14    }
15    else if (front) { /* add to head of queue */
16        rp->p_nextready = rdy_head[q]; /* chain head of queue */
17        rdy_head[q] = rp; /* set new queue head */
18    }
19    else { /* add to tail of queue */
20        rdy_tail[q]->p_nextready = rp; /* chain tail of queue */
21        rdy_tail[q] = rp; /* set new queue tail */
22        rp->p_nextready = NIL_PROC; /* mark new end */
23    }
24
25    /* Now select the next process to run. */
26    pick_proc();
27 }

```

```

1 /*=====*/
2 *      pick_proc      *
3 /*=====*/
4 PRIVATE void pick_proc()
5 {
6 /* Decide who to run now. A new process is selected by setting 'next_ptr'.
7 * When a billable process is selected, record it in 'bill_ptr', so that the
8 * clock task can tell who to bill for system time.
9 */
10 register struct proc *rp; /* process to run */
11 int q; /* iterate over queues */
12
13 /* Check each of the scheduling queues for ready processes. The number of
14 * queues is defined in proc.h, and priorities are set in the image table.
15 * The lowest queue contains IDLE, which is always ready.
16 */
17 for (q=0; q < NR_SCHED_QUEUES; q++) {
18     if ( (rp = rdy_head[q]) != NIL_PROC) {
19         next_ptr = rp; /* run process 'rp' next */
20         if (priv(rp)->s_flags & BILLABLE)
21             bill_ptr = rp; /* bill for system time */
22         return;
23     }
24 }
25 }

```

```

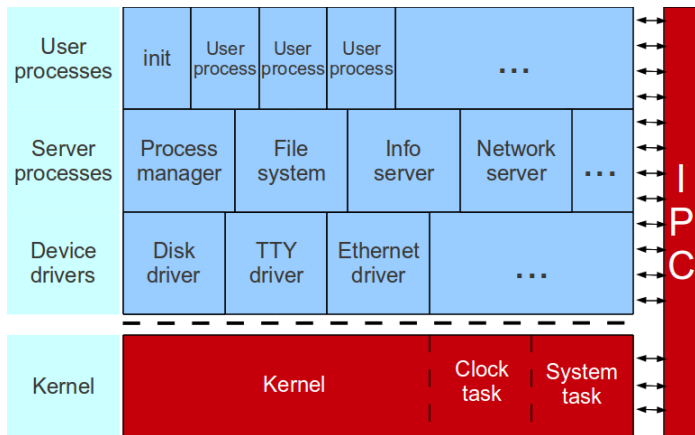
1 PRIVATE void sched(rp, queue, front)
2 register struct proc *rp;      /* process to be scheduled */
3 int *queue;                    /* return: queue to use */
4 int *front;                    /* return: front or back */
5 {
6     static struct proc *prev_ptr = NIL_PROC; /* previous without time */
7     int time_left = (rp->p_ticks_left > 0); /* quantum fully consumed */
8     int penalty = 0;            /* change in priority */
9
10    if ( ! time_left) {          /* quantum consumed ? */
11        rp->p_ticks_left = rp->p_quantum_size; /* give new quantum */
12        if (prev_ptr == rp) penalty ++; /* catch infinite loops */
13        else penalty --; /* give slow way back */
14        prev_ptr = rp; /* store ptr for next */
15    }
16
17    /* Determine the new priority of this process. The bounds are determined
18     * by IDLE's queue and the maximum priority of this process. Kernel task
19     * and the idle process are never changed in priority.
20     */
21    if (penalty != 0 && ! iskernelp(rp)) {
22        rp->p_priority += penalty; /* update with penalty */
23        if (rp->p_priority < rp->p_max_priority) /* check upper bound */
24            rp->p_priority = rp->p_max_priority;
25        else if (rp->p_priority > IDLE_Q-1) /* check lower bound */
26            rp->p_priority = IDLE_Q-1;
27    }
28
29    /* If there is time left, the process is added to the front of its queue,
30     * so that it can immediately run. The queue to use simply is always the
31     * process' current priority.
32     */
33    *queue = rp->p_priority;
34    *front = time_left;
35 }

```


Outline

- 1 Zasady
- 2 POSIX
 - Wstęp
 - POSIX - standard
 - POSIX - procesy
 - POSIX - pliki
 - POSIX - sygnały
 - POSIX - remanent
- 3 Sequential Processes
 - Multiprogramming
 - IPC - InterProcess Communication
- 4 Architektura
 - Monolithic kernel
 - Micro kernel
 - Tanenbaum - Torvalds debate
- 5 Scheduling
 - Batch systems

IPC - message passing



Message Passing Primitives

- SEND=1
- RECEIVE=2
- SENDREC=3
- NOTIFY=4

Kernel trap

```
push    ebp
mov     ebp, esp
push    ebx
mov     eax, SRC_DST(ebp) ! eax = dest-src
mov     ebx, MESSAGE(ebp) ! ebx = message pointer
mov     ecx, {SEND|RECEIVE|SENDREC|NOTIFY}
int     SYSVEC          ! 33
pop     ebx
pop     ebp
ret
```

```

1 PUBLIC int sys_call(call_nr, src_dst, m_ptr)
2 int call_nr;          /* system call number and flags */
3 int src_dst;          /* src to receive from or dst to send to */
4 message *m_ptr;       /* pointer to message in the caller's space */
5 {
6     /.../
7     int function = call_nr & SYSCALL_FUNC; /* get system call function
8     /...// check several conditions
9     switch(function) {
10    case SENDREC:
11        /* A flag is set so that notifications cannot interrupt SENDREC
12        priv(caller_ptr)->s_flags |= SENDREC_BUSY;
13        /* fall through */
14    case SEND:
15        result = mini_send(caller_ptr, src_dst, m_ptr, flags);
16        if (function == SEND || result != OK) {
17            break; /* done, or SEND failed */
18        } /* fall through for SENDREC */
19    case RECEIVE:
20        if (function == RECEIVE)
21            priv(caller_ptr)->s_flags &= ~SENDREC_BUSY;
22        result = mini_receive(caller_ptr, src_dst, m_ptr, flags);
23        break;
24    case NOTIFY:
25        result = mini_notify(caller_ptr, src_dst);
26        break;
27    /...// ECHO, default

```

Message Passing Primitives

- SEND, RECEIVE, SENDREC - są blokujące (rendez-vous). Uwaga na deadlock!
- NOTIFY - nie blokujące, nie buforowane
- Wiadomości wychodzące z jądra (kernel tasks) nie przechodzą przez sys_call.

Testy w sys_call

- Wiadomości mogą być wysyłane do zadań jądra wyłącznie przez SENDREC.
- Każdy proces ma maskę operacji IPC które może wykonywać (w strukturze priv) - procesy użytkownika mogą używać jedynie SENDREC.
- Argument src_dst musi być prawidłowym identyfikatorem procesu/zadania lub ANY dla RECEIVE.
- Wskaźnik do wiadomości musi się znajdować w pamięci procesu wołającego.
- Każdy proces ma maskę procesów do których może wysyłać (w strukturze priv) - procesy użytkownika mogą wysyłać do PM, FS, RS.
- Proces wskazywany przez src_dst musi być żywy lub ANY.

```

1 PRIVATE int mini_notify(caller_ptr, dst)
2 register struct proc *caller_ptr; /* sender of the notification */
3 int dst; /* which process to notify */
4 {
5     register struct proc *dst_ptr = proc_addr(dst);
6     int src_id; /* source id for late delivery */
7     message m; /* the notification message */
8
9     /* Check to see if target is blocked waiting for this message. A process
10      * can be both sending and receiving during a SENDREC system call.
11      */
12     if ((dst_ptr->p_rts_flags & (RECEIVING|SENDING)) == RECEIVING &&
13         ! (priv(dst_ptr)->s_flags & SENDREC.BUSY) &&
14         (dst_ptr->p_getfrom == ANY || dst_ptr->p_getfrom == caller_ptr->p_nr)) {
15
16         /* Destination is indeed waiting for a message. Assemble a notification
17          * message and deliver it. Copy from pseudo-source HARDWARE, since the
18          * message is in the kernel's address space.
19          */
20         BuildMess(&m, proc_nr(caller_ptr), dst_ptr);
21         CopyMess(proc_nr(caller_ptr), proc_addr(HARDWARE), &m,
22                 dst_ptr, dst_ptr->p_messbuf);
23         dst_ptr->p_rts_flags &= ~RECEIVING; /* deblock destination */
24         if (dst_ptr->p_rts_flags == 0) enqueue(dst_ptr);
25         return(OK);
26     }
27
28     /* Destination is not ready to receive the notification. Add it to the
29      * bit map with pending notifications. Note the indirectness: the system id
30      * instead of the process number is used in the pending bit map.
31      */
32     src_id = priv(caller_ptr)->s_id;
33     set_sys_bit(priv(dst_ptr)->s_notify_pending, src_id);
34     return(OK);
35 }

```

```

1 PRIVATE int mini_receive(caller_ptr, src, m_ptr, flags)
2 register struct proc *caller_ptr; /* process trying to get message */
3 int src; /* which message source is wanted */
4 message *m_ptr; /* pointer to message buffer */
5 unsigned flags; /* system call flags */
6 {
7     /*.../
8     if (!(caller_ptr->p_rts_flags & SENDING)) {
9
10    /* Check if there are pending notifications, except for SENDREC.
11    if (! (priv(caller_ptr)->s_flags & SENDREC_BUSY)) {
12        map = &priv(caller_ptr)->s_notify_pending;
13        for (chunk=&map->chunk[0]; \
14             chunk<&map->chunk[NR_SYS_CHUNKS]; \
15             chunk++) {
16            /* Find a pending notification from the requested source.
17            /*.../
18            /* Found a suitable source,
19            * deliver the notification message. */
20            BuildMess(&m, src_proc_nr, caller_ptr); /* assemble message
21            CopyMess(src_proc_nr, proc_addr(HARDWARE), \
22                    &m, caller_ptr, m_ptr);
23            return(OK); /* report success */
24        }
25    }
26    ...

```

```

1 ...
2  /* Check caller queue. Use pointer pointers to keep code simple. */
3  xpp = &caller_ptr->p_caller_q;
4  while (*xpp != NIL_PROC) {
5      if (src == ANY || src == proc_nr(*xpp)) {
6          /* Found acceptable message. Copy it and update status. */
7          CopyMess((*xpp)->p_nr,*xpp,(*xpp)->p_messbuf,caller_ptr,m
8          if (((*xpp)->p_rts_flags&= ~SENDING) == 0) enqueue(*xpp);
9          *xpp = (*xpp)->p_q_link;    /* remove from queue */
10         return(OK);                /* report success */
11     }
12     xpp = &(*xpp)->p_q_link;    /* proceed to next */
13 }
14 }
15 ...

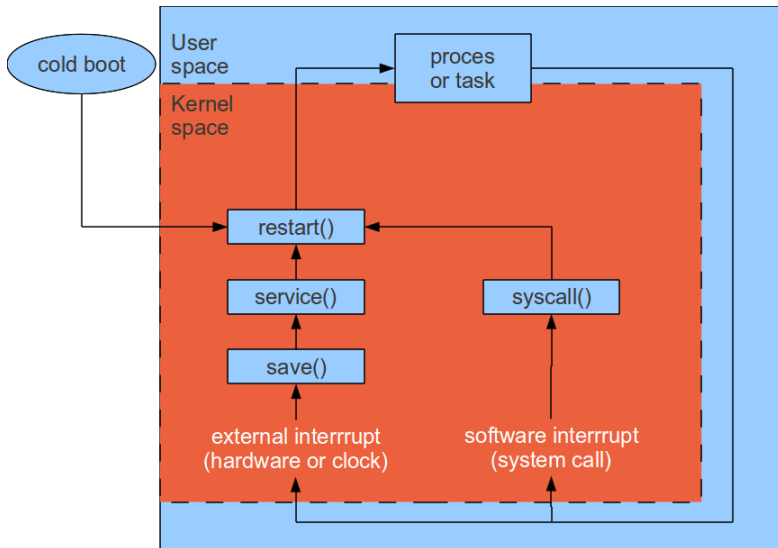
```



```

1 ...
2 /* No suitable message is available or the caller couldn't
3    * send in SENDREC.
4    * Block the process trying to receive ,
5    * unless the flags tell otherwise.
6    */
7 if ( ! (flags & NON_BLOCKING)) {
8     caller_ptr->p_getfrom = src;
9     caller_ptr->p_messbuf = m_ptr;
10    if (caller_ptr->p_rts_flags == 0) dequeue(caller_ptr);
11    caller_ptr->p_rts_flags |= RECEIVING;
12    return(OK);
13 } else {
14     return(ENOTREADY);
15 }
16 }

```



```

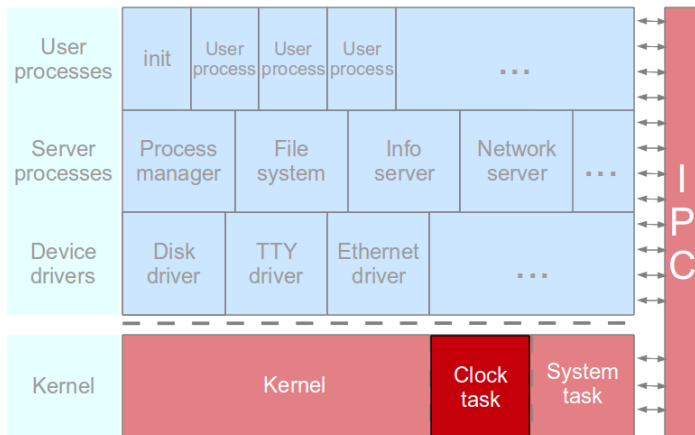
1 PRIVATE int generic_handler(hook)
2 irq_hook_t *hook;
3 {
4 /* This function handles hardware interrupt in a simple and generic
5 * way. All interrupts are transformed into messages to a driver.
6 * The IRQ line will be reenabled if the policy says so.
7 */
8
9 /* As a side-effect, the interrupt handler gathers random
10 * information by timestamping the interrupt events.
11 * This is used for /dev/random.
12 */
13 get_randomness(hook->irq);
14
15 /* Add a bit for this interrupt to the process' pending interrupts.
16 * When sending the notification message, this bit map will be
17 * magically set as an argument.
18 */
19 priv(proc_addr(hook->proc_nr))->s_int_pending |= (1<<hook->notify_id);
20
21 /* Build notification message and return. */
22 lock_notify(HARDWARE, hook->proc_nr);
23 return(hook->policy & IRQ_REENABLE);
24 }

```

Outline

- 1 Zasady
- 2 POSIX
 - Wstęp
 - POSIX - standard
 - POSIX - procesy
 - POSIX - pliki
 - POSIX - sygnały
 - POSIX - remanent
- 3 Sequential Processes
 - Multiprogramming
 - IPC - InterProcess Communication
- 4 Architektura
 - Monolithic kernel
 - Micro kernel
 - Tanenbaum - Torvalds debate
- 5 Scheduling
 - Batch systems

Clock task



`clock_handler()` Obsługa przerwań zegara - aktualizacja czasu, “księgowość”, sprawdzanie alarmów i potrzeby wywołania.

`clock_task()` Wywłaszczanie procesów których czas się wyczerpał, uruchamianie alarmów.

inne Funkcje dotyczące zegara:

- `get_uptime()`
- `set_timer()`
- `reset_timer()`
- `read_clock()`
- `clock_stop()`

```

1 PRIVATE int clock_handler(hook)
2 irq_hook_t *hook;
3 {
4     register unsigned ticks;
5
6     /* Acknowledge the PS/2 clock interrupt. */
7     if (machine.ps_mca) outb(PORT_B, inb(PORT_B) | CLOCK_ACK_BIT);
8
9     /* Get number of ticks and update realtime. */
10    ticks = lost_ticks + 1;
11    lost_ticks = 0;
12    realtime += ticks;
13
14    proc_ptr->p_user_time += ticks;
15    if (priv(proc_ptr)->s_flags & PREEMPTIBLE) {
16        proc_ptr->p_ticks_left -= ticks;
17    }
18    if (! (priv(proc_ptr)->s_flags & BILLABLE)) {
19        bill_ptr->p_sys_time += ticks;
20        bill_ptr->p_ticks_left -= ticks;
21    }
22
23    /* Check if do_clocktick() must be called. Done for alarms and scheduling.
24     * Some processes, such as the kernel tasks, cannot be preempted.
25     */
26    if ((next_timeout <= realtime) || (proc_ptr->p_ticks_left <= 0)) {
27        prev_ptr = proc_ptr; /* store running process */
28        lock_notify(HARDWARE, CLOCK); /* send notification */
29    }
30    return(1); /* reenale interrupts */
31 }

```

```

1 PUBLIC void clock_task()
2 {
3 /* Main program of clock task. If the call is not HARD_INT it is an error */
4 */
5 message m;          /* message buffer for both input and output */
6 int result;         /* result returned by the handler */
7
8 init_clock();        /* initialize clock task */
9
10 /* Main loop of the clock task. Get work, process it. Never reply. */
11 while (TRUE) {
12
13     /* Go get a message. */
14     receive(ANY, &m);
15
16     /* Handle the request. Only clock ticks are expected. */
17     switch (m.m_type) {
18     case HARD_INT:
19         result = do_clocktick(&m); /* handle clock tick */
20         break;
21     default:
22         /* illegal request type */
23         kprintf("CLOCK: illegal request %d from %d.\n", m.m_type, m.m_source);
24     }
25 }

```



```

1 PRIVATE int do_clocktick(m_ptr)
2 message *m_ptr;          /* pointer to request message */
3 {
4     if (prev_ptr->p_ticks_left <= 0 &&
5         priv(prev_ptr)->s_flags & PREEMPTIBLE) {
6         lock_dequeue(prev_ptr);    /* take it off the queues */
7         lock_enqueue(prev_ptr);    /* and reinsert it again */
8     }
9
10    /* Check if a clock timer expired and run its watchdog function. */
11    if (next_timeout <= realtime) {
12        tmrs_exptimers(&clock_timers, realtime, NULL);
13        next_timeout = clock_timers == NULL ?
14            TMR_NEVER : clock_timers->tmr_exp_time;
15    }
16
17    /* Inhibit sending a reply. */
18    return(EDONTREPLY);
19 }

```

Outline

1 Zasady

2 POSIX

- Wstęp
- POSIX - standard
- POSIX - procesy
- POSIX - pliki
- POSIX - sygnały
- POSIX - remanent

3 Sequential Processes

- Multiprogramming
- IPC - InterProcess Communication

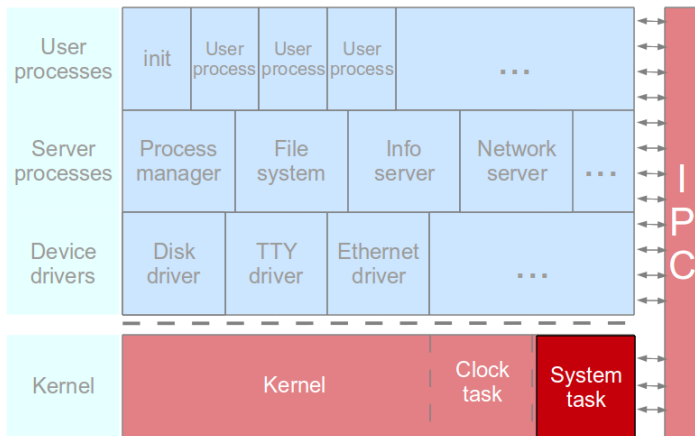
4 Architektura

- Monolithic kernel
- Micro kernel
- Tanenbaum - Torvalds debate

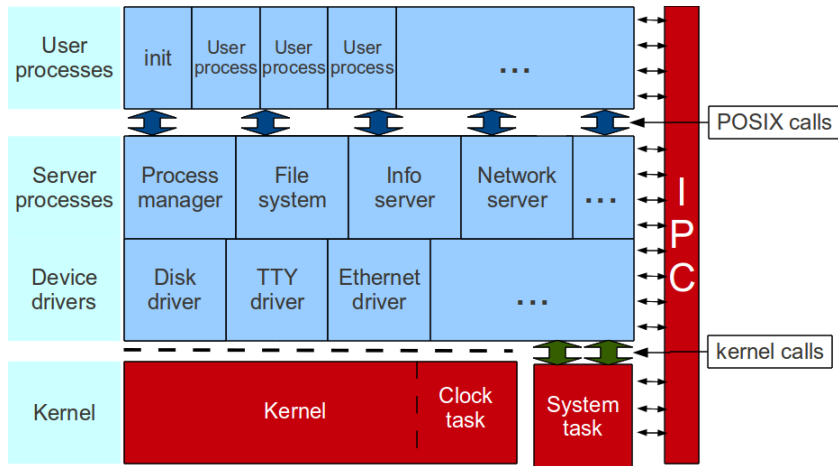
5 Scheduling

- Batch systems

System task



System task



Kernel calls

Message type	From	Meaning
sys_fork	PM	A process has forked.
sys_exec	PM	Set stack pointer after EXEC call
sys_exit	PM	A process has exited.
sys_nice	PM	Set scheduling priority.
sys_trace	PM	Carry out on operation of PTRACE call.
sys_kill	PM, FS, TTY	Send signal to a process after KILL call. (notify)
sys_getksig	PM	PM is checking for pending signals.
sys_endksig	PM	PM has finished processing signals.
sys_sigsend	PM	Send a signal to a process.
sys_sigreturn	PM	Cleanup after completion of a signal.
sys_newmap	PM	Set up a process memory map.
sys_memset	PM	Write char to memory area.
sys_times	PM	Get uptime and process times.
sys_setalarm	PM, FS, Drivers	Schedule synchronous alarm.
sys_abort	PM, TTY	Panic: MINIX is unable to continue

Kernel calls

Message type	From	Meaning
sys_privctl	RS	Set or change privileges.
sys_irqctl	Drivers	Enable, disable or configure interrupt.
sys_devio	Drivers	Read from or write to I/O port.
sys_sdevio	Drivers	Read or write string from/to I/O port.
sys_vdevio	Drivers	Carry out vector of I/O requests.
sys_int86	Drivers	Do a real-mode BIOS call
sys_segctl	Drivers	Add segment and get selector (far data access).
sys_umap	Drivers	Convert virtual address to physical addr.
sys_vircopy	FS, Drivers	Copy using pure virt. addressing.
sys_physcopy	Drivers	Copy using physical addressing.
sys_virvcopy	Any	Vector of VCOPY requests.
sys_physvcopy	Any	Vector of PHYSCOPY requests.
sys_getinfo	Any	Requests system information.

```

1 PUBLIC void sys_task()
2 {
3  /* Main entry point of sys_task. Get the message and dispatch on type. */
4  /* .../
5
6  initialize();
7
8  while (TRUE) {
9      /* Get work. Block and wait until a request message arrives. */
10     receive(ANY, &m);
11     call_nr = (unsigned) m.m_type - KERNEL_CALL;
12     caller_ptr = proc_addr(m.m_source);
13
14     /* See if the caller made a valid request and try to handle it. */
15     if (! (priv(caller_ptr) -> s.call_mask & (1<<call_nr))) {
16 kprintf("SYSTEM:_request_%d_from_%d_denied.\n", call_nr, m.m_source);
17 result = ECALLDENIED; /* illegal message type */
18 } else if (call_nr >= NR_SYS_CALLS) { /* check call number */
19 kprintf("SYSTEM:_illegal_request_%d_from_%d.\n", call_nr, m.m_source);
20 result = EBADREQUEST; /* illegal message type */
21 }
22 else {
23     result = (*call_vec[call_nr])(&m); /* handle the kernel call */
24 }
25
26 /* Send a reply, unless inhibited by a handler function. Use the kernel
27  * function lock_send() to prevent a system call trap. The destination
28  * is known to be blocked waiting for a message.
29  */
30 if (result != EDONTREPLY) {
31     m.m_type = result; /* report status of call */
32     if (OK != (s=lock_send(m.m_source, &m))) {
33         kprintf("SYSTEM,_reply_to_%d_failed:_%d\n", m.m_source, s);
34     }
35 }
36 }
37

```

```

1 PRIVATE void initialize(void)
2 {
3     register struct priv *sp;
4     int i;
5
6     /* Initialize IRQ handler hooks. Mark all hooks available. */
7     for (i=0; i<NR_IRQ_HOOKS; i++) {
8         irq_hooks[i].proc_nr = NONE;
9     }
10
11    /* Initialize all alarm timers for all processes. */
12    for (sp=BEG_PRIV_ADDR; sp < END_PRIV_ADDR; sp++) {
13        tmr_inittimer(&(sp->s_alarm_timer));
14    }
15
16    /* Initialize the call vector to a safe default handler. */
17    for (i=0; i<NR_SYS_CALLS; i++) {
18        call_vec[i] = do_unused;
19    }
20
21    /* Process management. */
22    map(SYS_FORK, do_fork);          /* a process forked a new process */
23    map(SYS_EXEC, do_exec);         /* update process after execute */
24    /* .../
25 }

```



```

1 PUBLIC int do_setalarm(m_ptr)
2 message *m_ptr;      /* pointer to request message */
3 {
4 /* A process requests a synchronous alarm, or wants to cancel its alarm. */
5 /* .../
6 /* Extract shared parameters from the request message. */
7 exp_time = m_ptr->ALRM_EXP_TIME; /* alarm's expiration time */
8 use_abs_time = m_ptr->ALRM_ABS_TIME; /* flag for absolute time */
9 proc_nr = m_ptr->m_source; /* process to interrupt later */
10 rp = proc_addr(proc_nr);
11 if (! (priv(rp)->s.flags & SYS_PROC)) return(EPERM);
12
13 /* Get the timer structure and set the parameters for this alarm. */
14 tp = &(priv(rp)->s.alarm_timer);
15 tmr_arg(tp)->ta_int = proc_nr;
16 tp->tmr_func = cause_alarm;
17
18 /* Return the ticks left on the previous alarm. */
19 uptime = get_uptime();
20 if ((tp->tmr_exp_time != TMR_NEVER) && (uptime < tp->tmr_exp_time) ) {
21     m_ptr->ALRM_TIME_LEFT = (tp->tmr_exp_time - uptime);
22 } else {
23     m_ptr->ALRM_TIME_LEFT = 0;
24 }
25
26 /* Finally, (re)set the timer depending on the expiration time. */
27 if (exp_time == 0) {
28     reset_timer(tp);
29 } else {
30     tp->tmr_exp_time = (use_abs_time) ? exp_time : exp_time + get_uptime();
31     set_timer(tp, tp->tmr_exp_time, tp->tmr_func);
32 }
33 return(OK);
34 }

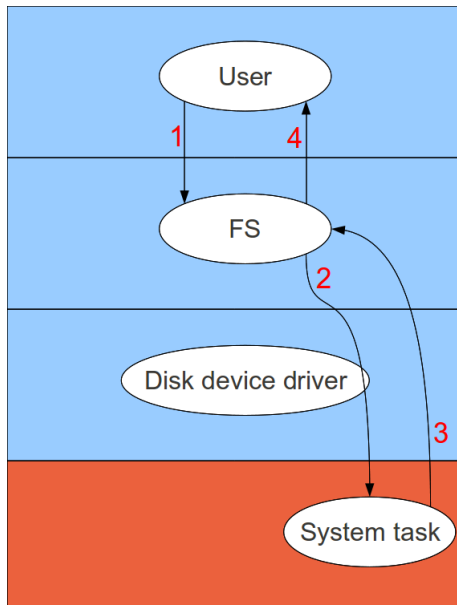
```

```

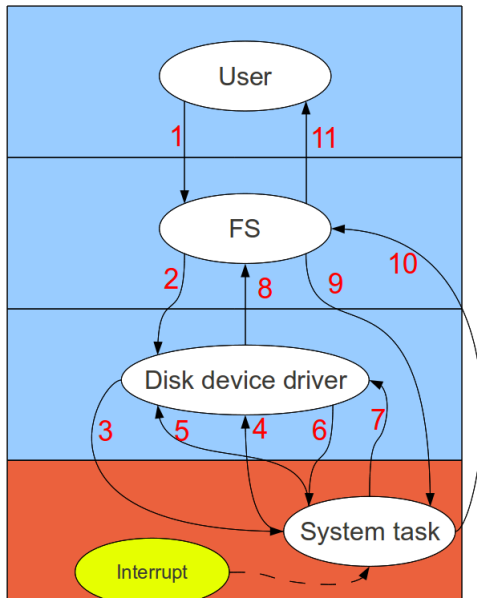
1 PRIVATE void cause_alarm(tp)
2 timer_t *tp;
3 {
4 /* Routine called if a timer goes off and the process requested
5  * a synchronous alarm.
6  * The process number is stored in timer argument 'ta_int'.
7  * Notify that process with a notification message from CLOCK.
8  */
9  int proc_nr = tmr_arg(tp)->ta_int;    /* get process number */
10  lock_notify(CLOCK, proc_nr);          /* notify process */
11 }

```

read()



read()



Minix kernel

