

Operatii IO avansate (1)

Contents

- 1 Linux - multiplexarea I/O
 - ◆ 1.1 select
 - ◆ 1.2 poll
 - ◇ 1.2.1 Avantaje poll
 - ◇ 1.2.2 Dezavantaje poll
 - ◆ 1.3 epoll
 - ◇ 1.3.1 Crearea unui obiect epoll
 - ◇ 1.3.2 Adăugarea/eliminarea descriptorilor la/de la obiectul epoll
 - ◇ 1.3.3 Ateptarea unui eveniment I/O
 - ◇ 1.3.4 Edge-triggered sau level-triggered
 - ◇ 1.3.5 Exemplu folosire epoll
- 2 Linux - generalizarea multiplexarii
 - ◆ 2.1 eventfd
 - ◆ 2.2 signalfd
- 3 Windows - I/O asincron (overlapped)
 - ◆ 3.1 Apeluri pentru transfer asincron (Overlapped I/O)
 - ◆ 3.2 Interogarea/ateptarea operaiilor asincrone
 - ◆ 3.3 Notificarea încheierii operaiilor asincrone
 - ◆ 3.4 Operatii asincrone pe socketi
- 4 Exerciii
 - ◆ 4.1 Presentare
 - ◆ 4.2 Exerciii de laborator
 - ◇ 4.2.1 Linux
 - ◇ 4.2.2 Windows
- 5 Soluii
- 6 Resurse utile

Linux - multiplexarea I/O

Există situaii în care un program trebuie să trateze operaiile I/O de pe mai multe canale ori de câte ori acestea apar. Un astfel de exemplu este un program de tip server care folosește mecanisme precum pipe-uri sau socketi pentru comunicarea cu alte procese. Un program trebuie să citească practic simultan informații atât de la intrarea standard cât și de la un socket (sau mai mulți).

În aceste situaii nu pot fi folosite operaii obișnuite de citire sau scriere. Folosirea acestor operaii are drept consecință blocarea thread-ului curent până la încheierea operaiei. O posibilă soluție este folosirea de operaii non-blocante (spre exemplu folosirea flag-ului `O_NONBLOCK`) și interogarea succesivă a descriptorilor de fiier. Totuși, interogarea succesivă (polling) este o formă de așteptare ocupată (busy waiting) și este ineficientă.

Soluția este folosirea unor mecanisme care permit unui thread să aștepte producerea unui eveniment I/O pe un set de descriptori. Thread-ul se va bloca până când unul din descriptorii din set poate fi folosit pentru citire/scriere. Un server care folosește un mecanism de acest tip are, de obicei, o structură de forma:

```
set = setul de descriptori urmărit
```

```

while (true) {
    așteaptă producerea unui eveniment pe unul din descriptori
    pentru fiecare descriptor pe care s-a produs un eveniment I/O {
        tratează evenimentul I/O
    }
}

```

Detaliile variază de la o implementare la alta, dar secvența de pseudocod de mai sus reprezintă structura de bază pentru serverele care folosesc multiplexarea I/O.

select

O primă soluție este utilizarea funcțiilor `select` sau `pselect`. Folosirea acestor funcții conduce la blocarea thread-ului curent până la producerea unui eveniment I/O pe un set de descriptori de fiier, a unei erori pe set sau până la expirarea unui timer.

Funcțiile folosesc un set de descriptori de fiier pentru a preciza fiierile/socketii pe care thread-ul curent va aștepta producerea evenimentelor I/O. Tipul de date folosit pentru definirea acestui set este `fd_set`, care este, de obicei, o mască de bii.

Funcțiile `select` și `pselect` sunt definite conform POSIX.1-2001 în `sys/select.h`

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

```
int select(int nfd, fd_set *readfds, fd_set *writefds, fd_set *exceptfds, struct timeval *timeou
```

Deoarece și apelul **poll** este specificat în standardul POSIX (deci portabilitate mare), dar oferă performanțe mai bune, **nu** vom insista asupra apelului **select**!

Avantaje:

- simplitate;
- portabilitate: funcția `select` este disponibilă chiar și în API-ul Win32;

Dezavantaje:

- lungimea setului de descriptori este definită cu ajutorul lui `FD_SETSIZE`, și implicit are valoarea 64;
- este necesar ca seturile de descriptori să fie reconstruite la fiecare apel `select`;
- la apariția unui eveniment pe unul din descriptori, toi descriptorii pui în set înainte de `select` trebuie testați pentru a vedea pe care din ei a apărut evenimentul;
- la fiecare apel, același set de descriptori este transmis în kernel.

poll

Funcția `poll` consolidează argumentele funcției `select` și permite notificarea pentru o gamă mai largă de evenimente. Funcția se definește ca mai jos:

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

```
int poll(struct pollfd *ufds, unsigned int nfd, int timeout);
```

`select`

Timeout-ul este specificat în milisecunde. În caz de valoare negativă, semnificația este de așteptare pentru o perioadă nedefinită ("infinit").

Structura `struct pollfd` este definită în `sys/poll.h`:

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

struct pollfd {
    int fd;          /* file descriptor */
    short events;   /* evenimente solicitate */
    short revents;  /* evenimente apărute */
};
```

Funcția `poll` permite astfel așteptarea evenimentelor descrise de vectorul `ufds` de dimensiune `nfds`.

În cadrul structurii `struct pollfd`, câmpul `events` este o mască de bii în care se specifică evenimentele urmărite de `poll` pentru descriptorul `fd`. `revents` este, de asemenea, o mască de bii completată de kernel cu evenimentele apărute în momentul în care apelul se întoarce (`POLLIN`, `POLLOUT`) sau valori predefinite (`POLLERR`, `POLLHUP`, `POLLNVAL`) pentru situații speciale.

În caz de succes funcția returnează un număr diferit de zero reprezentând numărul de structuri pentru care `revents` nu e zero (cu alte cuvinte toți descriptorii cu evenimente sau erori). Se returnează 0 dacă a expirat timpul (timeout milisecunde) și nu a fost selectat nici un descriptor. În caz de eroare se returnează -1 și se setează `errno`. De asemenea, funcția `poll` poate fi întreruptă de semnale, caz în care va întoarce -1 și `errno` va fi setat la `EINTR`.

Un exemplu de utilizare a `poll` este prezentat în continuare:

```
#define MAX_PFDS      32

[...]
struct pollfd pfds[MAX_PFDS];
int nfds;
int listenfd, sockfd;      /* listener socket; connection socket */

nfds = 0;

/* read user data from standard input */
pfds[nfds].fd = STDIN_FILENO;
pfds[nfds].events = POLLIN;
nfds++;

/* TODO ... create server socket (listener) */

/* add listener socket */
pfds[nfds].fd = listenfd
pfds[nfds].events = POLLIN;
nfds++;

while (1) {          /* server loop */
    /* wait for readiness notification */
    poll(pfds, nfds, -1);

    if ((pfds[1].revents & POLLIN) != 0) {
        /* TODO ... handle new connection */
    }
    else if ((pfds[0].revents & POLLIN) != 0) {
        /* TODO ... read user data from standard input */
    }
}
```

```

    }
    else {
        /* TODO ... handle message on connection sockets */
    }
}
[...]
```

Avantaje poll

- transmiterea setului de descriptori este mai simplă decât în cazul funcției `select`;
- setul de descriptori nu trebuie reconstruit la fiecare apel `poll`;

Dezavantaje poll

- ineficiență - la apariția unui eveniment, trebuie parcurs tot setul de descriptori pentru a găsi descriptorul pe care a apărut evenimentul;
- la fiecare apel, acelai set de descriptori (care poate fi mare) este copiat în kernel și înapoi.

epoll

Funcțiile `select` și `poll` nu sunt scalabile la un număr mare de conexiuni pentru că la fiecare apel al lor trebuie transmisă toată lista de descriptori. În astfel de situații, la fiecare pas, trebuie construită lista de descriptori și apelat `poll` sau `select` care copiază tot setul în kernel. La apariția unui eveniment va fi marcat corespunzător descriptorul. Utilizatorul trebuie să parcurgă tot setul de descriptori pentru a-i da seama pe care dintre ei a apărut evenimentul. În acest fel se ajunge să se petreacă tot mai mult timp scanând după evenimente în setul de descriptori și tot mai puțin timp făcând I/O.

Din acest motiv, diverse sisteme au implementat interfee scalabile dar non-portabile:

- `/dev/poll` pe Solaris;
- `kqueue` pe FreeBSD;
- `epoll` pe Linux.

Aceste interfee rezolvă problemele asociate cu `select` și `poll` și rezolvă problemele de scalabilitate.

Pentru a folosi `epoll`, trebuie inclus `sys/epoll.h`. Funcțiile asociate sunt `epoll_create`, `epoll_ctl` și `epoll_wait`. Interfaa `epoll` oferă funcții pentru crearea unui obiect `epoll`, adăugarea sau eliminarea de descriptori de fiere/sockete la obiectul `epoll` și așteptarea unui eveniment pe unul dintre descriptori.

Crearea unui obiect epoll

Pentru crearea unui obiect `epoll` se folosește funcția `epoll_create`:

```
int epoll_create(int size);
```

Apelul `epoll_create` facilitează crearea unui descriptor de fiier ce va fi ulterior folosit pentru așteptarea de evenimente. Descriptorul întors va trebui închis folosind `close`.

Argumentul `size` este ignorat în versiunile recente ale nucleului, acesta ajustând dinamic dimensiunea setului de descriptori asociați obiectului `epoll`.

Adăugarea/eliminarea descriptorilor la/de la obiectul `epoll`

Operațiile de adăugare/eliminare de descriptori se realizează cu ajutorul funcției `epoll_ctl`:

```
int epoll_ctl(int epollfd, int op, int fd, struct epoll_event *event);
```

Apelul `epoll_ctl` permite specificarea evenimentelor care vor fi așteptate. Câmpul `event` descrie evenimentul asociat descriptorului `fd` care poate fi adăugat, ters sau modificat în funcție de valoarea argumentului `op`:

- `EPOLL_CTL_ADD`: pentru adăugare;
- `EPOLL_CTL_MOD`: pentru modificare;
- `EPOLL_CTL_DEL`: pentru tergere.

Primul argument al apelului `epoll_ctl` (`epollfd`) este descriptorul întors de `epoll_create`.

Structura `struct epoll_event` specifică evenimentele așteptate

```
typedef union epoll_data {
    void *ptr;
    int fd;
    __uint32_t u32;
    __uint64_t u64;
} epoll_data_t;

struct epoll_event {
    __uint32_t events;      /* Epoll events */
    epoll_data_t data;     /* User data variable */
};
```

Tipuri de evenimente care pot fi urmărite sau primite pe un descriptor:

- `EPOLLIN`: descriptorul asociat are date de citit;
- `EPOLLOUT`: descriptorul asociat are loc în buffer-ul asociat pentru a scrie date;
- `EPOLLERR`: a apărut o eroare pe descriptorul asociat.

Union-ul `epoll_data` poate fi folosit pentru a asocia o cheie cu descriptorul urmărit.

Așteptarea unui eveniment I/O

Thread-ul curent așteaptă producerea unui eveniment I/O la unul din descriptorii asociați obiectului `epoll` prin intermediul funcției `epoll_wait`:

```
int epoll_wait(int epollfd, struct epoll_event* events, int maxevents, int timeout);
```

Functia `epoll_wait` este echivalentul funcțiilor `select` și `poll`. Este folosită pentru așteptarea unui eveniment la unul din descriptorii asociați descriptorului `epollfd`.

La revenirea apelului programatorul nu va trebui să parcurgă toți descriptorii configurați ci numai cei care au evenimente produse. Argumentul `events` va marca o zonă de memorie unde vor fi plasate maxim `maxevents` evenimente de nucleu. Presupunând că valoarea câmpului `timeout` este `-1` (așteptare nedefinită), apelul se va întoarce imediat dacă există evenimente asociate, sau se va ploua până la apariția unui eveniment.

La fel ca și în cazul `select/pselect` și `poll/ppoll`, există apelul `epoll_pwait` care permite precizarea unei măști de semnale.

Edge-triggered sau level-triggered

Interfața `epoll` are două comportamente posibile: *edge-triggered* sau *level-triggered*. Se poate folosi unul sau altul, în funcție de prezența flag-ului `EPOLLET` la adăugarea unui descriptor în lista `epoll`.

Presupunem existența unui socket funcționând în mod non-blocant pe care sosesc 100 de octeți. În ambele moduri (edge sau level triggered) `epoll_wait` va raporta `EPOLLIN` pentru acel socket.

Vom presupune că se citesc 50 de octeți din cei 100 primii. Diferența între cele două moduri de funcționare apare la un nou apel `epoll_wait`. În modul level-triggered se va raporta imediat `EPOLLIN`. În modul edge-triggered nu se va mai raporta nimic, nici măcar la sosirea unor noi date pe socket. Se poate observa cum modul edge-triggered sesizează *schimbarea* stării descriptorului în relație cu evenimentul, iar level-triggered *prezenta* stării. Modul edge-triggered este implementat mai eficient în kernel, chiar dacă pare mai greu de folosit.

În continuare, sunt prezentate câteva reguli care trebuie urmărite cu o metodă sau alta. Pentru ambele metode este recomandată folosirea socketelor în modul non-blocant.

- *Level-triggered*
 - ◆ La apariția unui eveniment `EPOLLIN` se poate citi oricât, la următorul apel `epoll_wait` se va raporta din nou `EPOLLIN` dacă mai sunt date de citit.
 - ◆ `EPOLLOUT` nu trebuie configurat inițial pentru un socket pentru că astfel `epoll_wait` va raporta imediat că este loc de scris în buffer (inițial bufferul de scriere asociat cu socketul este gol). Acesta este o formă deghizată de busy waiting. Folosirea corectă implică scrierea normală pe socket și numai dacă la un moment dat funcțiile de scriere raportează că nu mai este loc de scriere în buffer (`EAGAIN`), se va activa `EPOLLOUT` pe descriptorul respectiv și salva ce mai este de scris. Când în sfârșit se face loc, se va raporta `EPOLLOUT` și atunci se poate încerca să se scrie datele păstrate. Dacă se reușește scrierea lor integrală, **trebuie** eliminat flagul `EPOLLOUT` pentru a nu intra într-un nou ciclu de busy-waiting. În concluzie, `EPOLLOUT` trebuie activat doar când nu se reușește scrierea integrală a datelor și scos imediat după ce acestea au fost scrise.
- *Edge-triggered*
 - ◆ La apariția unui eveniment `EPOLLIN` pe un descriptor, trebuie citit tot ce se poate citi înainte de reapelarea `poll_wait`, altfel nu va mai fi raportat `EPOLLIN` niciodată.
 - ◆ Pentru scrierea folosind edge-triggered se poate activa de la început `EPOLLOUT`. Aceasta va cauza apariția unui eveniment `EPOLLOUT` imediat după apelarea `epoll_wait` (pentru că

bufferul de scriere este gol) care ar trebui ignorat. La următorul apel `epoll_wait` nu se mai generează `EPOLLOUT` pentru că nu s-a schimbat starea de la ultimul apel. Dacă la un moment dat se încearcă scrierea unor date pe socket și acestea nu pot fi scrise integral, la următorul `epoll_wait` se generează `EPOLLOUT`, pentru că s-a schimbat starea socketului. Mai pe scurt, asta are ca efect faptul că nu mai trebuie activat/deactivat `EPOLLOUT` ca în cazul `level-triggered`.

Exemplu folosire `epoll`

Mai jos este prezentat un exemplu de utilizare a `epoll` echivalent cu exemplele pentru `select` și `poll` (server care multiplexează mai multe conexiuni pe socket și intrarea standard):

```
#define EPOLL_INIT_BACKSTORE      2

[...]
```

```
int listenfd, sockfd;           /* listener socket; connection socket */
struct epoll_event ev;

/* create epoll descriptor */
epfd = epoll_create(EPOLL_INIT_BACKSTORE);

/* read user data from standard input */
ev.data.fd = STDIN_FILENO;      /* key is file descriptor */
ev.events = EPOLLIN;
epoll_ctl(epfd, EPOLL_CTL_ADD, STDIN_FILENO, &ev);

/* TODO ... create server socket (listener) */

/* add listener socket */
ev.data.fd = listenfd;         /* key is file descriptor */
ev.events = EPOLLIN;
epoll_ctl(epfd, EPOLL_CTL_ADD, listenfd, &ev);

while (1) {                    /* server loop */
    struct epoll_event ret_ev;

    /* wait for readiness notification */
    epoll_wait(epfd, &ret_ev, 1, -1);

    if ((ret_ev.data.fd == listenfd && ((ret_ev.events & EPOLLIN) != 0)) {
        /* TODO ... handle new connection */
    }
    else if ((ret_ev.data.fd == STDIN_FILENO &&
              ((ret_ev.events & EPOLLIN) != 0)) {
        /* TODO ... read user data from standard input */
    }
    else {
        /* TODO ... handle message on connection sockets */
    }
}
[...]
```

Linux - generalizarea multiplexării

O problemă a funcțiilor de multiplexare de mai sus (`select`, `poll`, `epoll`) este aceea că sunt limitate la descriptori de fiier. Altfel spus, se pot aștepta doar evenimente asociate cu un fiier/socket: gata de citire, gata de scriere. De multe ori însă se dorește să existe un punct comun de așteptare a unui semnal, a unui semafor, a unui proces, a unei operații de intrare/ieșire, a unui timer. În Windows, acest lucru se poate realiza cu ajutorul funcției `WaitForMultipleObjects` și pe faptul că majoritatea mecanismelor din Windows sunt folosite cu ajutorul tipului de date `HANDLE`.

eventfd

Pentru a asigura în Linux posibilitatea așteptării de evenimente multiple s-a definit interfața **eventfd**. Cu ajutorul acestei interfețe și combinat cu interfețele de multiplexare I/O existente, kernel-ul poate notifica o aplicație utilizator de orice tip de eveniment.

Interfața **eventfd** este prezentă în nucleul Linux începând cu versiunea 2.6.22. În momentul de față (glibc 2.7), biblioteca standard C nu oferă suport pentru aceste mecanisme. Suportul va fi oferit în versiunea 2.8 a glibc. Pentru folosirea mecanismelor oferite de kernel se poate folosi interfața `syscall`:

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

static int eventfd(unsigned int initval, int flags)
{
    return syscall(__NR_eventfd, initval, 0);
}
```

Cele trei apeluri de bază pentru extinderea funcționalității multiplexării I/O sunt `eventfd`, `signalfd` și `timerfd_create`.

Toate cele trei apeluri întorc un descriptor de fiier pe care se vor putea primi notificări (semnale, timere etc.). Operațiile posibile pe descriptorul de fiier întors sunt:

- `write`: pentru transmiterea unui mesaj de notificare pe descriptor;
- `read`: pentru primirea unui mesaj care înseamnă primirea notificării;
- `select`, `poll`, `epoll`: pentru multiplexarea I/O;
- `close`: pentru închiderea descriptorului și eliberarea resurselor asociate.

În următorul exemplu, apelul `eventfd` este folosit pentru notificarea procesului părinte de către procesul fiu. Codul este cel prezent în pagina de manual (`man eventfd`).

```
[...]
int efd;
uint64_t u;

/* create eventfd file descriptor */
efd = eventfd(0, 0);

switch (fork()) {
case 0:
    /* notify parent process */
    s = write(efd, &u, sizeof(uint64_t));

    printf("Child completed write loop\n");
    exit(EXIT_SUCCESS);
}
```



```

default:
    printf("Parent about to read\n");

    /* wait for notification */
    s = read(efd, &u, sizeof(uint64_t));
    exit(EXIT_SUCCESS);
[...]
```

signalfd

Apelul `signalfd` este folosit în mod similar pentru recepționarea de semnale prin intermediul unui descriptor de fiier. Pentru a putea recepționa un semnal cu ajutorul interfeei `signalfd`, va trebui blocat în masca de semnale a procesului. La fel ca i exemplul de mai sus, codul de mai jos este cel prezent în pagina de manual (`man signalfd`).

```

/* at this point Linux-specific headers are required to use struct signalfd_siginfo */
#include <linux/types.h>
#include <linux/signalfd.h>

#define SIZEOF_SIG      (_NSIG / 8)
#define SIZEOF_SIGSET  (SIZEOF_SIG > sizeof(sigset_t)? \
                        sizeof(sigset_t): SIZEOF_SIG)

/* on par with the current glibc wrapper */
static int signalfd(int fd, const sigset_t *mask, int flags)
{
    return syscall(__NR_signalfd, fd, mask, SIZEOF_SIGSET);
}

[...]
```

```

sigset_t mask;
int sfd;
struct signalfd_siginfo fdsi;

sigemptyset(&mask);
sigaddset(&mask, SIGINT);           /* CTRL-C */
sigaddset(&mask, SIGQUIT);        /* CTRL-\ */

/*
 * Block signals so that they arent handled
 * according to their default dispositions
 */

sigprocmask(SIG_BLOCK, &mask, NULL);

/* create signalfd descriptor */
sfd = signalfd(-1, &mask);

for (;;) {
    /* wait for signals to be delivered by user */
    s = read(sfd, &fdsi, sizeof(struct signalfd_siginfo));

    if (fdsi.ssi_signo == SIGINT) {
        printf("Got SIGINT\n");
    } else if (fdsi.ssi_signo == SIGQUIT) {
        printf("Got SIGQUIT\n");
        exit(EXIT_SUCCESS);
    } else {
```

```

        printf("Read unexpected signal\n");
    }
}
[...]
```

Interfaa **eventfd** permite unificarea mecanismelor de notificare ale kernel-ului într-un descriptor de fiier care va fi folosit de utilizator.

În acest moment (kernel 2.6.29, glibc 2.9) a fost unificat subsistemul de operatii I/O al nucleului Linux cu interfaa **eventfd** si exista si suport in biblioteca standard C pentru acest apel (incepand cu 2.8).

ATENȚIE! Versiunea glibc instalata in laborator (2.7) nu are suport pentru eventfd, fiind necesara folosirea wrapper-ului de mai sus, ce foloseste syscall.

Windows - I/O asincron (overlapped)

În Windows se folosesc funcțiile `ReadFile`, `WriteFile` sau `ReadFileEx`, `WriteFileEx` pentru pornirea operatiilor asincrone. Pentru notificarea terminării operatiilor în Windows se pot folosi *evenimente*, *APC-uri*, *Thread Pooling* sau *Completion Port-uri*.

Apeluri pentru transfer asincron (Overlapped I/O)

În Windows, operatiile asincrone se numesc **overlapped I/O** (operatii suprapuse). Pentru folosirea overlapped I/O în cazul fiierelor se foloseste ultimul argument transmis funcției `ReadFile`, respectiv `WriteFile`, ca în exemplul de mai jos:

```

HANDLE hFile = CreateFile(...);
char buffer[BUFSIZ];
DWORD bytesRead;
OVERLAPPED ov;

/* init overlapped structure */
ov.Offset = 0;          /* read from beginning of file */

/* TODO ... ov.hEvent = ... - if using an event for notification */

if (ReadFile(
    hFile,
    buffer,
    BUFSIZ,
    &bytesRead,
    &ov) == FALSE) {
    fprintf(stderr, "ReadFile failed with %d\n", GetLastError());
    ExitProcess(-1);
}
[...]
```

Ultimul argument este un pointer la structura `OVERLAPPED` care descrie operaia asincronă de realizat. Pentru ca operatiile asincrone să poată decurge pe fiere, fiierul trebuie deschis cu flag-ul

`FILE_FLAG_OVERLAPPED`:

```

HANDLE hFile;
```

```

/* open file for reading using overlapped I/O */
hFile = CreateFile(
    files[i],
    GENERIC_READ,
    FILE_SHARE_READ,
    NULL,
    OPEN_EXISTING,
    FILE_ATTRIBUTE_NORMAL <PIPE> FILE_FLAG_OVERLAPPED,
    NULL);

```

Interogarea/ateptarea operaiilor asincrone

Pentru determinarea stării operaiei asincrone se poate folosi funcția `GetOverlappedResult`:

```

HANDLE hFile;
OVERLAPPED ov;
DWORD bytesTransferred;

/* TODO ... start overlapped I/O operation */

/* wait for completion */
GetOverlappedResult(hFile, &ov, &bytesTransferred, TRUE);

```

După cum se poate observa, ultimul argument al apelului `GetOverlappedResult` difereniază între interogarea operaiei asincrone (`FALSE`) sau așteptarea încheierii acesteia (`TRUE`).

Notificarea încheierii operaiilor asincrone

Pentru notificarea încheierii unei operaiei asincrone se pot folosi:

- *evenimente* - prin completarea corespunzătoare a câmpului `hEvent` al structurii `OVERLAPPED`;
- *APC* - prin folosirea apelurilor `ReadFileEx/WriteFileEx` care permit specificarea unei rutine apelată la încheierea operaiei;
- *I/O completion ports* - mecanismul scalabil de așteptare a operaiilor asincrone.

Operatii asincrone pe socketi

La fel ca la fisiere, folosirea operatii asincrone (overlapped I/O) pe socketi presupune folosirea unei structuri specializate (`WSAOVERLAPPED`). Pentru transmiterea/receptia de informatii in format asincron se folosesc functiile `WSARecv`, respectiv `WSASend`.

Un exemplu de folosire a functiei `WSARecv` pentru primirea asincrona a informatiilor este prezentat in continuare. Codul este inspirat din exemplul din [pagina MSDN asociata](#):

```

char buffer[BUFSIZ];
WSABUF DataBuf;
DWORD flags;
DWORD recvBytes;
WSAOVERLAPPED ov;
SOCKET s;

```

```

/* TODO ... connect socket */

DataBuf.buf = buffer;
DataBuf.len = BUFSIZ;
flags = 0;

/* start asynchronous I/O */
WSARecv(s, &DataBuf, 1, &recvBytes, &flags, &ov, NULL);

```

Exerciii

Prezentare

Pentru a urmări mai uor noiunile expuse la începutul laboratorului folosii [această prezentare \(pdf\)](#) ([odp](#)).

Exerciii de laborator

Folositi [arhiva de sarcini](#) a laboratorului.

Linux

1. **pollpipe (2p)** Creai folosind [fork\(2\)](#) o aplicaie de test pentru [poll\(2\)](#). Aplicaia foloseste un server (părintele) i `CLIENT_COUNT` clieni (copiii) ce comunică prin pipe-uri anonime.
 - ◆ server-ul:
 - ◇ construiete un vector de pipe-uri (în funcia `main`);
 - ◇ creează clienii;
 - ◇ blochează în așteptarea datelor pe acestea i tipărete datele primite;
 - ◇ termină execuia după ce a primit date de la fiecare client;
 - ◆ clienii:
 - ◇ ateaptă un număr aleator de secunde (mai mic decât 10);
 - ◇ scriu în pipe-ul corespunzător un ir de `MSG_SIZE` caractere de forma "`<pid>:<character random>`" (`'a' + random() % 30`);
 - ◇ scrierile i citirile în pipe-uri de până la `PIPE_BUF` octei (4096 pe Linux) sunt atomice.

Hints:

 - ◇ Consultai seciunea [poll](#) i [Pipe-uri în Linux](#).
2. **epollpipe (2p)** Modificai codul anterior pentru a folosi [epoll](#).
3. **eventpipe (2p)** Modificai codul anterior pentru a adăuga funcionaliata de deînregistrare a clienilor.
 - ◆ Se va folosi un descriptor [eventfd](#) prin intermediul căruia serverul este notificat de încheierea execuiei unui client.
 - ◆ server-ul:
 - ◇ creează `eventfd`-ul i îl adaugă la descriptor-ul `epoll`;
 - ◇ la primirea unui eveniment prin acesta, dacă primii 32 bii sunt `MAGIC_EXIT`, atunci scoate capătul pipe-ului corespunzător din `epoll` și îl închide;

- ◊ există 2 opțiuni pentru a determina pipe-ul care trebuie scos și închis, pe baza restului de 32 de bii din event:
 - îi folosești pentru PID-ul client-ului și suplimentar reții în server o asocierie client-pipe
 - va folosi de mostenirea descriptorilor la fork()
- ◊ face exit după ce a primit MAGIC_EXIT de la toți cei CLIENT_COUNT clienți
- ◆ clienții:
 - ◊ fac trimiterea datelor într-o buclă cu IT_COUNT iterații
 - ◊ la fiecare iterație, verifică dacă o valoare random() este multiplu de IT_COUNT, și dacă da trimite un eveniment MAGIC_EXIT și termină execuția.
- 4. **signalpipe (2p)** Înlocuiește în codul anterior notificarea (de terminare a clienților) bazată pe **eventfd** cu una bazată pe semnale.
 - ◆ server-ul:
 - ◊ creează un descriptor via signalfd pentru SIGCHLD și-l adaugă la epoll
 - ◊ la primirea unui semnal, prin read(2) pe descriptorul creat, determină PID-ul copilului defunct, afișează un mesaj și scoate pipe-ul din epoll.

Windows

1. **aio (2p)** Scrierea unor date aleatoare în fișiere folosind operații sincrone și asincrone. Implementați funcțiile `do_io_sync`, respectiv `do_io_async`.
 - ◆ Pentru operațiile sincrone puteți folosi funcția `xwrite` definită în fișier.
 - ◆ Va trebui să alocăți spațiu pentru structurile `OVERLAPPED` pentru toate cele 4 fișiere.
 - ◆ Pentru inițializarea structurilor `OVERLAPPED` se recomandă implementarea funcției `init_overlapped`.
 - ◆ Folosiți `GetOverlappedResult` pentru realizarea operațiilor asincrone pe cele 4 fișiere.
 - ◆ Funcțiile trebuie să scrie conținutul bufferului `g_buffer` în cele 4 fișiere cu numele dat de vectorul `files`.
 - ◆ Folosiți macro-ul `IO_OP_TYPE` pentru a determina comportamentul programului.

Soluii

Resurse utile

- `eventfd`
 - ◆ [Exemplu de cod integrare eventfd cu Linux AIO](#)
- `select/poll/epoll`
 - ◆ Linux System Programming - Chapter 2 - File I/O (Multiplexed I/O)
 - ◆ Linux System Programming - Chapter 4 - Advanced File I/O (The Event Poll Interface)
 - ◆ Beginning Linux Programming, 4th Edition - Chapter 15: Sockets (select)
 - ◆ Un [articol interesant](#) despre epoll, la un nivel mai înalt
- Windows I/O Completion Ports
 - ◆ [MSDN - I/O Completion Ports](#)
 - ◆ [Inside I/O Completion Ports](#)
 - ◆ [O introducere în Completion Ports](#)
 - ◆ [Tutorial IOCP și socketi](#)

- General
 - ◆ libevent - bibliotecă de operații asincrone
 - ◆ C10K - Problema celor 10.000 de clienți