

I/O nel nucleo

G. Lettieri

9 Maggio 2017

Le operazioni di I/O offrono un'ottima opportunità per sfruttare appieno l'ambiente multiprogrammato. Tipicamente, un processo che inizia una operazione di I/O non può proseguire finché l'operazione non è terminata, e d'altro canto l'operazione stessa è normalmente lenta e non ha bisogno della CPU per essere portata avanti, o ne ha bisogno in minima parte. L'idea è quindi di *bloccare* i processi che iniziano una operazione di I/O e sbloccarli (riportarli in coda pronti) quando l'operazione è terminata. In questo modo altri processi potranno andare in esecuzione mentre è in corso l'operazione di I/O e la CPU sarà sfruttata più efficientemente. Si noti che il sistema deve sapere che l'operazione è completata mentre è in esecuzione un processo che, in generale, è completamente scorrelato da essa. Il completamento dell'operazione dovrà essere segnalato da una interruzione, in modo che il sistema possa riacquistare il controllo della CPU e svolgere le operazioni necessarie, tra cui sbloccare il processo che aveva richiesto l'operazione.

Si noti che quanto appena descritto è un esempio di *sincronizzazione*: vogliamo che il processo che ha iniziato l'I/O possa proseguire solo dopo che l'operazione di I/O è stata completata.

In ambiente multiprocesso il sistema deve anche preoccuparsi di coordinare l'accesso alle periferiche. Tipicamente, mentre una periferica è impegnata in una operazione di I/O, non può essere usata per altre operazioni. Se un processo vuole usare una periferica mentre questa è occupata, deve aspettare che la precedente operazione termini e la periferica ritorni libera. Si noti che, a differenza del precedente, questo non è un problema di sincronizzazione, ma di *mutua esclusione*: se due processi, P_1 e P_2 , vogliono entrambi usare la stessa periferica, è per noi indifferente se la usa prima P_1 e poi P_2 o viceversa; l'unica cosa che ci interessa è che non la usino contemporaneamente. Si noti infine che la mutua esclusione possiamo garantirla separatamente per ogni periferica: se P_1 e P_2 devono usare due periferiche distinte non hanno bisogno di coordinarsi.

Notiamo ora che i processi di cui stiamo parlando sono processi *utente* e, come al solito, sono da considerarsi non fidati. Non possiamo dunque aspettarci che gli utenti si coordinino tra di loro per usare le periferiche uno alla volta, o che sospendano volontariamente i propri processi per far andare avanti quelli degli altri. Per imporre la mutua esclusione e la sincronizzazione di cui sopra procediamo come al solito:

- impediamo agli utenti di parlare direttamente con le periferiche e
- forniamo delle primitive per svolgere le operazioni di I/O sotto il controllo del sistema.

Per realizzare il primo punto facciamo in modo che il campo IOPL del registro RFLAGS dei processi utente specifichi il valore “sistema”. In questo modo, mentre il processore si trova a livello utente, genera una eccezione ogni volta che si prova ad eseguire una istruzione di `in` o `out` ¹ Notare che, indipendentemente da questo problema, siamo sostanzialmente obbligati a settare il campo IOPL a “sistema”, in quanto questo è anche il modo in cui vietiamo agli utenti l’utilizzo delle istruzioni `sti` e `cli`, che abilitano e disabilitano le interruzioni esterne mascherabili. Per vietare l’accesso alle periferiche che hanno i registri mappati in memoria ricorriamo invece alla MMU, in uno dei due modi possibili: o non inserendo traduzioni che portino ai registri delle periferiche, o proteggendo le traduzioni con i bit “S/U” nei descrittori.

Occupiamoci ora della struttura generale delle primitive di I/O che il sistema deve fornire. Una tipica primitiva di lettura avrà una interfaccia simile a questa:

```
extern "C" void read_n(natl id, natb *buf, natl quanti);
```

La primitiva riceve un parametro `id`, che serve ad identificare la periferica da cui il processo vuole leggere, e un indirizzo `buf` che punta ad un buffer in cui l’utente vuole ricevere i dati. Faremo l’ipotesi che i dati siano sempre una sequenza di byte. Il parametro `quanti` specifica il numero di byte che l’utente vuole leggere. È l’utente che deve preoccuparsi di dichiarare un buffer grande a sufficienza per contenere i byte richiesti. Vedremo che il sistema può imporre altre limitazioni su questo buffer.

Analogamente, la tipica primitiva di scrittura avrà questa interfaccia verso gli utenti:

```
extern "C" void write_n(natl id, const natb *buf, natl quanti);
```

I parametri hanno lo stesso significato del caso precedente, ma ora il buffer puntato da `buf` deve contenere i dati che l’utente vuole scrivere (la primitiva si limita a leggerli, da cui il `const`).

In Figura 1 vediamo un esempio di programma utente che utilizza una primitiva di ingresso e una di uscita. In questo caso si tratta delle primitive che il sistema mette a disposizione per leggere una linea dalla tastiera (`readconsole`) e scrivere una linea sul video (`writeconsole`). Le primitive sono dichiarate in `sys.h` nel seguente modo:

```
extern "C" void readconsole(char* buff, natl& quanti);
extern "C" void writeconsole(const char* buff);
```

¹Per vietare `in` e `out` sono necessari anche altri accorgimenti, che tralasciamo per semplicità; maggiori dettagli si trovano nel codice.

```

1 #include <sys.h>
2 #include <lib.h>
3
4 process hello body hello_body(0), 20, LIV_UTENTE;
5
6 const int NAME_SIZE = 80;
7
8 char nome[NAME_SIZE];
9 natl lun = NAME_SIZE;
10
11 process_body hello_body(int a)
12 {
13     char buf[NAME_SIZE + 100], *ptr;
14     writeconsole("Ciao, come ti chiami?");
15     readconsole(nome, lun);
16     ptr = copy("Ciao_", buf);
17     ptr = copy(nome, ptr);
18     ptr = copy(", piacere di conoscerti", ptr);
19     writeconsole(buf);
20     pause();
21 }

```

Figura 1: Un programma utente con operazioni di I/O.

Entrambe non hanno bisogno di un parametro `id`, in quanto il sistema dispone di un'unica tastiera e un unico video, e usano il più comodo `char` invece di `natb`. La `writeconsole`, inoltre, è particolare in quanto non ha bisogno di un parametro che indichi il numero dei byte da scrivere (usa il terminatore di stringa per capire quando fermarsi) e, soprattutto, non blocca il processo che la invoca, in quanto la scrittura sul video richiede l'uso della CPU per copiare i caratteri in memoria video, e dunque non possiamo usare la CPU per eseguire un altro processo.

Nella `readconsole` il parametro `quanti` è passato per riferimento in quanto l'operazione legge una linea, terminata da *a capo*, la cui lunghezza non è nota *a priori*. L'utente scrive inizialmente in `quanti` la dimensione del buffer che lui ha allocato, in modo che la primitiva non tenti comunque di leggere più byte di quelli. Al ritorno, la primitiva scrive in `quanti` il numero di byte effettivamente letti.

La primitiva `readconsole` blocca il processo fino a quando l'operazione non è conclusa. Questo ha perfettamente senso: l'operazione si svolge alla velocità con cui l'utente umano batte i caratteri sulla tastiera; nel frattempo il processore ha tutto il tempo di portare avanti altri processi. Dal punto di vista del processo che ha chiamato la `readconsole` il blocco è completamente trasparente: alla riga 15 il programma chiede di leggere un massimo di `NAME_SIZE` caratteri nel buffer `nome` (allocato alla riga 8); alla riga 17 copia la stringa contenuta in `nome` all'indirizzo `ptr`. Come si vede, il programma assume che, una volta

che la `readconsole` è ritornata, l'operazione sia conclusa. Inoltre, dal programma non è possibile dedurre che il processo si blocca durante l'esecuzione della `readconsole`. Il programma non si preoccupa nemmeno di garantire che nessun altro stia cercando di usare la tastiera: si limita a chiamare la `readconsole` quando ne ha bisogno (e lo stesso per la `writeconsole`). Sincronizzazione e mutua esclusione, dunque sono nascoste all'interno delle primitive, e da queste garantite.

1 Realizzazione con primitiva e driver

Torniamo a considerare una generica operazione di lettura, con interfaccia utente

```
extern "C" void read_n(natl id, natb *buf, natl quanti)
```

Assumiamo che nel sistema siano installate diverse periferiche simili, identificate dunque dal parametro `id`. Assumiamo inoltre che queste periferiche siano in grado di trasferire un byte alla volta, inviando una richiesta di interruzione ogni volta che è disponibile un nuovo byte. L'interfaccia di ogni periferica avrà un registro di controllo per abilitare e disabilitare le interruzioni (CTL) e un registro di ingresso da cui leggere il nuovo byte (RBR). La lettura da RBR funziona da risposta alla richiesta di interruzione: l'interfaccia non genera una nuova richiesta di interruzione fino a quando non ha avuto una risposta a quella precedente.

L'operazione sarà svolta in parte dalla primitiva e in parte da un *driver*, che andrà in esecuzione ad ogni richiesta di interruzione da parte dell'interfaccia:

- la primitiva ha lo scopo di avviare l'operazione di I/O e bloccare il processo, garantendo anche la mutua esclusione;
- il driver ha il compito di trasferire effettivamente i byte e sbloccare il processo quando l'operazione si è conclusa.

Consideriamo ora un processo P_1 che invoca la primitiva `read_n`. Questa, come per tutte le primitive, è in realtà solo una piccola funzione scritta in Assembler nel file `utente.s`. La funzione invoca la primitiva vera e propria tramite una istruzione `int`, che permette l'innalzamento del livello di privilegio:

```
1     .global read_n
2     read_n:
3         int $IO_TIPO_RN
4         ret
```

All'offset `IO_TIPO_RN` della tabella IDT dovrà essere installato un gate con `GP=1` e `IND=a_read_n`. Questa sarà una funzione scritta in assembler nel file `sistema.s`:

```
1     .extern c_read_n
2     a_read_n:
3         cavallo_di_troia %rsi
```

```

4         cavallo_di_troia2 %rsi %rdx
5         call c_read_n
6         iretq

```

Ci sono alcune cose da notare:

- anche se P_1 verrà bloccato durante l'esecuzione della primitiva, la funzione `a_read_n` non chiama `salva_stato` e `carica_stato`;
- è necessario controllare il problema del Cavallo di Troia.

Il primo punto è molto importante e lo riprenderemo tra poco. Per il secondo punto, si osservi che l'utente potrebbe passare non l'indirizzo di un buffer che ha correttamente allocato, ma (usando dei cast o scrivendo direttamente in assembler) l'indirizzo di qualunque cosa, anche di parti della memoria che appartengono al sistema o ad altri processi. I controlli di protezione eseguiti dalla CPU e dalla MMU sono inefficaci in questo caso: il passaggio dell'indirizzo "cattivo" alla primitiva comporta solo la scrittura di un numero in un registro e la CPU non controlla alcunché, in quanto non conosce lo scopo di questa operazione. L'apparentemente innocuo Cavallo di Troia attraversa dunque le mura del sistema. Quando la primitiva tenta di usare l'indirizzo per leggere o scrivere, ecco che il cavallo si apre, in quanto la primitiva gira a livello sistema e dunque anche la MMU non esegue alcun controllo. Per evitare questo, le righe 2 e 3 controllano che tutto il buffer (passato nei registri `%rsi`, base, e `%rdx`, lunghezza) si trovi nella zona utente (cosa che si può fare guardando semplicemente l'indirizzo, senza accedere alla memoria indirizzata), altrimenti abortiscono il processo con un errore.

Passiamo ora alla parte C++ della primitiva. Dato l'identificatore `id`, la primitiva ha bisogno di ottenere alcune informazioni sulla corrispondente periferica (l'indirizzo dei suoi registri, ma non solo). Per memorizzare tutte queste informazioni prevediamo un descrittore di operazione di I/O definito come segue:

```

1     struct des_io {
2         natw iRBR, iCTL;
3         natb *buf;
4         natl quanti;
5         natl mutex;
6         natl sync;
7     };

```

È sufficiente avere un array di tali descrittori e usare `id` come indice al suo interno. Ogni descrittore contiene tre tipi di informazioni: gli indirizzi dei registri (riga 2), le informazioni (ricevute dall'utente) su dove i byte vanno trasferiti (righe 3-4) e due identificatori di semafori (righe 5-6). Per realizzare la sincronizzazione e la mutua esclusione, infatti, la `c_read_n` utilizzerà i semafori, che servono proprio a questo. Il semaforo `mutex` è inizializzato a 1 all'avvio del sistema e serve a garantire la mutua esclusione tra i processi che vogliono usare la periferica `id`. Il semaforo `sync` è inizializzato a 0 all'avvio del sistema e serve

a realizzare la sincronizzazione tra il processo che ha richiesto l'operazione e la conclusione dell'operazione stessa.

La `c_read_n` è strutturata nel modo seguente:

```
1   extern "C" void c_read_n(natl id, natb *buf, natl quanti)
2   {
3       struct des_io *d = &array_des_io[id];
4
5       sem_wait(d->mutex);
6       d->buf = buf;
7       d->quanti = quanti;
8       outputb(1, d->iCTL);
9       sem_wait(d->sync);
10      sem_signal(d->mutex);
11  }
```

Alla riga 3 ottiene un puntatore al descrittore della interfaccia, per comodità di scrittura. Si noti che bisognerebbe controllare che `id` sia effettivamente l'indice di una interfaccia esistente, ma qui tralasciamo questo controllo per semplicità.

Le righe 5 e 6 garantiscono la mutua esclusione: solo un processo alla volta può eseguire le righe 6-9.

Le righe 6 e 7 trasferiscono le informazioni sul buffer dell'utente all'interno del descrittore. Da qui le leggerà il driver quando andrà in esecuzione.

La riga 8 abilita le interruzioni (supponiamo che sia sufficiente scrivere 1 nel registro CTL). Da questo momento l'interfaccia può inviare richieste di interruzione ogni volta che ha un nuovo byte da trasferire. Si noti che per il momento le interruzioni esterne sono mascherate, in quanto ci troviamo nel modulo sistema.

La riga 9 blocca P_1 sul semaforo di sincronizzazione. A questo punto il sistema può passare ad eseguire altri processi. Mentre sono in esecuzione questi altri processi le interruzioni sono abilitate. Si noti che P_1 è bloccato all'interno della `sem_wait` alla riga 9, e dunque ancora dentro la zona di mutua esclusione. Fino a quando P_1 non verrà sbloccato e non eseguirà la `sem_signal(d->mutex)` alla riga 10 l'interfaccia non potrà essere usata da altri processi, come volevamo. Se un altro processo (andato in esecuzione in seguito al blocco di P_1) invocherà la `read_n` sulla stessa interfaccia, arriverà alla riga 5 e si bloccherà.

Arriviamo ora spiegare perché la `a_read_n` non deve invocare `salva_stato` e `carica_stato`. Il motivo è che sono le primitive dei semafori a bloccare (se necessario) il processo, e dunque sono loro che salvano e caricano opportunamente lo stato. Salvare e ripristinare lo stato nella `a_read_n` non è però soltanto inutile: è un errore che causa un malfunzionamento del sistema. Occorre ricordare che ogni processo ha un solo descrittore in cui salvare il proprio stato: quanto salvato da una eventuale `salva_stato` posta erroneamente all'inizio della `a_read_n` verrebbe sovrascritto dalla `salva_stato` chiamata dalle primitive semaforiche. In generale, non si può chiamare due volte di seguito una `salva_stato` sullo stesso descrittore senza che ci sia stata in mezzo una `carica_stato` di quel descrittore.

Passiamo ora ad esaminare il driver. Il driver andrà in esecuzione per effetto di una richiesta di interruzione da parte dell'interfaccia (richiesta che arriverà alla CPU tramite l'APIC). Il driver gira a livello sistema e, per poter tornare a livello utente, deve terminare anch'esso con una `iretq`. Anche il driver, dunque, avrà una parte scritta in assembler:

```

1     .extern c_driver
2     a_driver_i:
3         call salva_stato
4         movq $i, %rdi
5         call c_driver
6         call apic_send_EOI
7         call carica_stato
8         iretq

```

Per fare in modo che `a_driver_i` vada in esecuzione ogni volta che l'interfaccia genera una interruzione, occorre conoscere il tipo dell'interruzione generata e preparare il corrispondente gate della IDT con `GP=1` e `IND=a_driver_i`.

Chiamiamo P_2 il processo in esecuzione all'arrivo dell'interruzione.

Il driver salva e ripristina lo stato di P_2 (riga 3) perché può dover cambiare il processo in esecuzione. Infatti, quanto l'ultimo byte è stato trasferito, il driver deve risvegliare P_1 (che ora è bloccato nella `sem_signal(d->sync)` dentro `c_read_n`). Se P_1 ha una priorità maggiore di P_2 , è P_1 che deve andare in esecuzione, mentre P_2 deve tornare in coda pronti. La `carica_stato` alla riga 7, quindi, carica lo stato di P_2 o di P_1 , a seconda dei casi.

Alle righe 4-5 si chiama il driver vero e proprio, `c_driver`, scritto in C++. Dal momento che stiamo assumendo di trattare interfacce simili, `c_driver` può essere una funzione generica e ricevere un parametro che gli indichi l'interfaccia da gestire (linea 4). Si noti che il parametro che `a_driver_i` passa a `c_driver` è una costante. Questo perché ogni interfaccia genera una richiesta di interruzione di tipo diverso, quindi è sufficiente associare ad ogni tipo di interruzione una diversa copia di `a_driver_i`, ciascuna con una costante diversa.

Alla riga 6 inviamo l'End Of Interrupt al controllore APIC, in modo che lasci passare le interruzioni a priorità minore o uguale di quella appena gestita dal driver.

Si noti che il driver usa le risorse di P_2 . In particolare usa la sua pila sistema, dal momento che, quando è arrivata l'interruzione, in TR c'era il puntatore al descrittore di processo di P_2 . Usa anche la memoria virtuale di P_2 , dal momento che non stiamo cambiando il valore di `CR3`.

Vediamo ora il codice di `c_driver`:

```

1     extern "C" void c_driver(natl id)
2     {
3         des_io *d = &array_des_io[id];
4         char c;
5
6         d->quanti--;

```

```

7     if (d->quanti == 0) {
8         outputb(0, d->iCTL);
9         c_sem_signal(d->sync);
10    }
11    inputb(d->iRBR, c);
12    *d->buf = c;
13    d->buf++;
14 }

```

Alla riga 3 il driver ottiene un puntatore al descrittore della interfaccia, per comodità di scrittura. Si noti che questa volta non c'è bisogno di controllare che `id` sia valido, in quanto è stato passato da `a_driver_id` che, facendo parte del modulo sistema, è fidata.

Lo scopo principale del driver è leggere il nuovo byte dall'interfaccia e copiarlo nel buffer dell'utente, cosa che viene fatta alle righe 11–12. Alla riga 13 il puntatore al buffer viene incrementato, in modo che il prossimo byte venga copiato nella locazione successiva.

Alla riga 6 il driver decrementa `d->quanti`, che così contiene sempre il numero di byte ancora da leggere. Se `d->quanti` è arrivato a zero il byte letto alla riga 11 è l'ultimo byte richiesto dall'utente. Si deve quindi disabilitare l'interfaccia a generare interruzioni (riga 8) e risvegliare il processo che aveva iniziato l'operazione (riga 9).

Ci sono alcune cose da notare:

1. la disabilitazione delle interruzioni (riga 8) è eseguita *prima* della lettura del byte (riga 11);
2. invece di chiamare `sem_signal()` chiamiamo direttamente `c_sem_signal()` (riga 9);
3. la scrittura in `d->buf` (riga 12) è eseguita mentre è attiva la memoria virtuale di P_2 , anche se il buffer era stato allocato da P_1 .

Il punto 1 è una conseguenza del fatto che la lettura del byte fa da risposta alla richiesta di interruzione da parte dell'interfaccia, che a quel punto può generarne un'altra se ha un nuovo byte disponibile. Quindi, se leggiamo l'ultimo byte mentre le interruzioni sono abilitate, è possibile che l'interfaccia generi una nuova interruzione, rimandando in esecuzione il driver, anche se nessun processo ha iniziato una operazione di lettura. Il driver leggerebbe questo byte e lo copierebbe dove punta `d->buf`, andando a sovrascrivere parti casuali della memoria.

Il punto 2 è una conseguenza del fatto che `sem_signal()` salva e ripristina lo stato, ma il driver non è un processo e non ha un suo descrittore di processo. In particolare, se `c_driver` chiamasse `sem_signal()` salverebbe lo stato del driver nel descrittore processo attivo, che è P_2 . Per di più sovrascriverebbe lo stato salvato alla riga 3 di `a_driver_i` (avremmo violato la regola di non avere due `salva_stato` senza una `carica_stato` in mezzo).

Si noti che, chiamando `c_sem_signal`, il driver manipola le code dei processi, e dunque deve essere eseguito con le interruzioni disabilitate. Questo comporta che anche le richieste di interruzione a precedenza maggiore dovranno aspettare che il driver termini, prima di poter essere gestite.

Il punto 3 è un problema se P_1 ha allocato il buffer nella sua sezione utente/privata, dal momento che gli indirizzi virtuali delle sezioni private hanno significati diversi per ogni processo. In questo caso `c_driver` scriverebbe nella sezione utente/privata di P_2 , non di P_1 , causando un doppio problema: P_1 non riceverebbe i suoi dati e P_2 si vedrebbe sovrascrivere parti casuali della sua memoria. Dobbiamo quindi richiedere che i buffer per l'I/O vengano sempre allocati nella parte utente/condivisa. Dal punto di vista del linguaggio C++, nel nostro caso, questo comporta che i buffer devono essere dichiarati globali o allocati nello heap, e mai dichiarati come variabili locali. Questo perché le variabili locali vengono allocate in pila e abbiamo deciso che le pile si trovano in parti private. Nei sistemi che bloccano i processi che causano page fault, il buffer deve anche essere residente (non rimpiazzabile), in quanto il driver, non essendo un processo, non può essere bloccato.

Esaminiamo infine come devono essere impostati i campi, oltre IND e GP, del gate che porta a `a_read_n` e di quello che porta a `a_driver_i`. Per il primo avremo:

- il campo GL (Gate Level) che indica che il gate può essere utilizzato da livello utente;
- il campo L (Level) che indica che, dopo il salto, il processore si deve trovare a livello sistema;
- il campo I/T che può indifferentemente indicare “Interrupt” o “Trap”, dal momento che la primitiva non manipola direttamente le code e i descrittori dei processi (assumeremo “Interrupt”, per uniformità).

Per il gate della `a_driver_i` avremo:

- il campo GL (Gate Level) che indica che il gate può essere utilizzato solo da livello sistema;
- il campo L (Level) che indica che, dopo il salto, il processore si deve trovare a livello sistema;
- il campo I/T che indica “Interrupt”, per quanto detto prima.

L'impostazione del campo GL deriva da questa constatazione: tutti i gate della IDT possono essere usati indifferentemente da tutti e tre i meccanismi di interruzione (interruzioni esterne, eccezioni e istruzione INT). Impostando GL a livello sistema impediamo all'utente di invocare il driver in modo spurio, tramite una INT.