

## 1 Lezione del 29-04-25

Riprendiamo il discorso del DMA nella prospettiva di un esempio concreto.

### 1.1 Hard disk e DMA

Fra i dispositivi visti finora solo l'hard disk è quello capace di fare DMA nel kernel. Dentro la macchina virtuale QEMU è disponibile un'emulazione dell'hard disk del PC AT (l'HD ATA visto in 4.1). Questo non era capace di fare DMA in autonomia, ma era bensì collegato ad un controllore DMA.

Fra i comandi disponibili per comunicare con l'hard disk ci sono quindi comandi dedicati a letture e scritture in DMA. Quando tali comandi vengono inviati all'hard disk, questi si occupa di coinvolgere il controllore DMA.

Questa non è più la situazione odierna: l'hard disk ATA con cui comunica la macchina emulata è situato sul bus ATA, che si collega al bus PCI con un ponte PCI-ATA. E' quindi il ponte a comportarsi come il controllore DMA, lato bus ATA.

Considerazioni storiche a parte, vediamo la struttura del controllore DMA dell'hard disk ATA, come descritto nella specifica reperibile a <https://calcolatori.iet.unipi.it/deep/idems100.pdf>. Abbiamo che questo può gestire due dischi separati, denominati *primario* e *secondario*, con relativi registri:

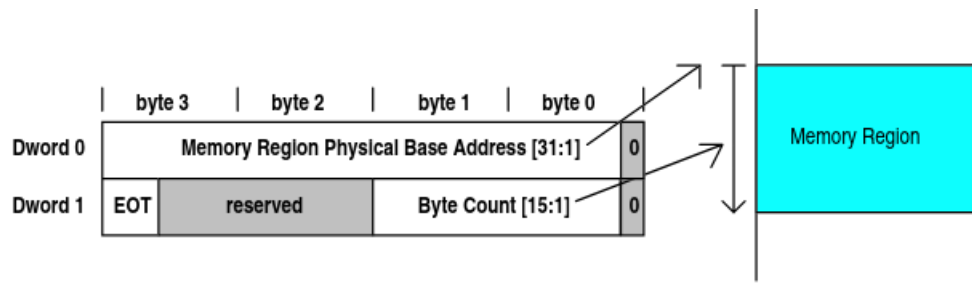
	<i>Primario</i>
0x?? + 0	<b>BMCMD</b> , <i>Bus Master Command</i>
0x?? + 1	Specifico al dispositivo
0x?? + 2	<b>BMSTR</b> , <i>Bus Master Status Register</i>
0x?? + 3	Specifico al dispositivo
0x?? + 4-7	<b>BMDTPR</b> , <i>Bus Master Descriptor Table Pointer</i>
	<i>Secondario</i>
0x?? + 8	<b>BMCMD</b> , <i>Bus Master Command</i>
0x?? + 9	Specifico al dispositivo
0x?? + a	<b>BMSTR</b> , <i>Bus Master Status Register</i>
0x?? + b	Specifico al dispositivo
0x?? + c-f	<b>BMDTPR</b> , <i>Bus Master Descriptor Table Pointer</i>

Riguardo a ogni registro avremo:

- **BMCMD**, *Bus Master Command*: questo specifica il tipo di operazione che vogliamo eseguire (lettura o scrittura), e ne specifica l'inizio. Per lanciare un'operazione, infatti, il software dovrà impostare il bit di *Read or Write Control* (bit 3), e successivamente alzare il bit *Start/Stop Bus Master* (bit 0);
- **BMSTR**, *Bus Master Status Register*: indica lo stato corrente del dispositivo a cui corrisponde. In particolare ci sono di interesse i primi 3 bit meno significativi (gli altri danno principalmente informazioni rispetto alle funzioni supportate dai dispositivi). Questi saranno:
  - Bit 2: rappresenta l'**interruzione**, viene alzato quando la trasmissione di dati in DMA è stata completata;
  - Bit 1: rappresenta uno stato di **errore**;
  - Bit 0: indica se il bus mastering è attivo o meno, cioè viene alzato quando il software scrive 1 sul bit start/stop bus master del BMCMD.

Abbiamo poi che i bit 1 e 2 possono essere resettati scrivendovi 1 (ed è questo passo che termina l'handshake col controllore DMA).

- **BMDTPR**, *Bus Master Descriptor Table Pointer*: questo punta alla prima entrata della cosiddetta tabella **PRD**, *Physical Region Descriptor Table*. Questa è una tabella di entrate da 8 byte, allineate ai 4 byte, che indicano l'indirizzo base della regione da trasferire, il numero di byte da trasferire e se l'entrata corrente è l'ultima della tabella (il controllore DMA continua a scorrere le entrate finché non raggiunge l'ultima). La struttura delle entrate PRD è la seguente:



Notiamo che le regioni indicate dall'indirizzo base dell'entrata PRD può essere al massimo di 64 KiB. Per questo lato hardware si può usare un sommatore a sole 16 cifre. In ogni caso, questo non sarà un problema in quanto vorremo trasferire buffer in memoria virtuale una pagina (4 KiB) alla volta.

A questo punto basterà definire i passaggi di un operazione di trasferimento:

1. Si prepara una tabella PRD in memoria;
2. Si carica l'indirizzo base della tabella PRD nel registro BMDTPR, quindi si ripuliscono i bit di interruzione ed errore del registro di stato BMSTR;
3. Si fornisce il comando appropriato sul registro BMCMD;
4. Si attiva il bit 0 del registro BMCMD per attivare il bus mastering;
5. Il controllore DMA trasferisce i dati secondo quanto disposto finora;
6. Alla fine della trasmissione il controllore segnala la fine dell'operazione su una linea di interruzione;
7. In risposta all'interruzione, si resetta il bit 0 del registro BMCMD, e si legge lo stato dal controllore e dal disco per capire se l'operazione è andata a buon fine.

### 1.1.1 Controller IDE su bus PCI

Per l'inserzione di un controllore di questo tipo in un bus PCI dobbiamo renderci conto di alcuni dettagli: Nei registri dello spazio di configurazione del dispositivo si devono attivare dei flag particolari per segnalare la possibilità che questo lavori in bus mastering.

### 1.1.2 Controller IDE nel kernel

Vediamo infine come il controllore DMA dell'hard disk ATA viene gestito nel kernel. La libreria `libce` definisce i registri del controllore:

```
1 namespace bm {
2     extern ioaddr iBMCMD; // Bus Master Command
3     extern ioaddr iBMSTR; // Bus Master Status Register
4     extern ioaddr iBMDTPR; // Bus Master Descriptor Table Pointer
5 }
```

e le relative funzioni per l'inizializzazione, l'acknowledge, ecc...

L'unica interfaccia ATA montata nel sistema è quindi descritta dal descrittore:

```
1 // descrittore di interfaccia ATA
2 struct des_ata {
3     // Ultimo comando inviato all'interfaccia
4     natb comando;
5     // Indice di un semaforo di mutua esclusione
6     natl mutex;
7     // Indice di un semaforo di sincronizzazione
8     natl sincr;
9     // Quanti settori resta da leggere o scrivere
10    natb cont;
11    // Da dove leggere/dove scrivere il prossimo settore
12    natb* punt;
13    // Array dei descrittori per il Bus Mastering
14    natl* prd;
15 };
```

che tiene conto dell'operazione corrente.

A questo punto il processo esterno dedicato all'hard disk dovrà limitarsi ad inviare i comandi corretti seguendo la scaletta appena riportata. Unica parte di interesse è quella della preparazione della tabella PRD, per cui bisogna tenere conto che il controllore DMA necessita di indirizzi fisici, e che legge sequenzialmente a partire da tali indirizzi fisici (perciò non si possono superare i 4 KiB della dimensione di pagina). Per fare questo, e tenere conto di buffer in memoria che iniziano potenzialmente a metà pagina, si sfrutta la funzione `prepare_prd()`:

```
1 bool prepare_prd(des_ata *d, natb* vett, natb quanti)
2 {
3     natq n = quanti * DIM_BLOCK;
4     int i = 0;
5
6     while (n && i < MAX_PRD) {
7         paddr p = trasforma(vett);
8         natq r = DIM_PAGINA - (p % DIM_PAGINA);
9         if (r > n)
10            r = n;
11         d->prd[i] = p;
12         d->prd[i + 1] = r;
13
14         n -= r;
15         vett += r;
16         i += 2;
17     }
18     if (n)
19         return false;
20     // il bit end of table
21     d->prd[i - 1] |= 0x80000000;
22     return true;
23 }
```

23 }

A questo punto si possono fornire all'utente primitive per l'accesso all'hard disk sia a controllo interruzione (come avevamo già visto, implementato in `libce`) sia in DMA. Queste saranno:

- **Controllo interruzione:** vediamo ad esempio l'operazione di ingresso.

```

1 // fondamentale un wrapper per hd::start_cmd di libce, che
  aggiorna il descrittore
2 void starthd_in(des_ata* d, natb vetti[], natl primo, natb quanti)
3 {
4     d->cont = quanti;
5     d->punt = vetti;
6     d->comando = hd::READ_SECT;
7     hd::start_cmd(primo, quanti, hd::READ_SECT);
8 }
9
10 // la primitiva vera e propria
11 extern "C" void c_readhd_n(natb vetti[], natl primo, natb quanti)
12 {
13     des_ata* d = &hard_disk;
14
15     // controlli (c_access)
16
17     sem_wait(d->mutex);
18     starthd_in(d, vetti, primo, quanti);
19     sem_wait(d->sincr);
20     sem_signal(d->mutex);
21 }
```

- **DMA:** vediamo sempre l'operazione di ingresso:

```

1 void dmastarthd_in(des_ata* d, natb vetti[], natl primo, natb quanti)
2 {
3     // passo 1 della scaletta
4     if (!prepare_prd(d, vetti, quanti)) {
5         flog(LOG_ERR, "dmastarthd_in: numero di PRD insufficiente");
6         sem_signal(d->sincr);
7         return;
8     }
9
10    d->comando = hd::READ_DMA;
11    d->cont = 1;
12
13    // passo 2
14    paddr prd = trasforma(d->prd);
15    bm::prepare(prd, false);
16
17    // passo 3
18    hd::start_cmd(primo, quanti, hd::READ_DMA);
19    bm::start();
20 }
```

A operazioni terminate, il processo esterno dovrà chiaramente interpretare correttamente le interruzioni che riceve in base al tipo di comando dato:

```

1 void estern_hd(natq)
2 {
3     des_ata* d = &hard_disk;
4     for(;;) {
```

```
5     d->cont--;
6     hd::ack();
7     switch (d->comando) {
8         // questi sono i casi gia visti
9         case hd::READ_SECT:
10            hd::input_sect(d->punt);
11            d->punt += DIM_BLOCK;
12            break;
13         case hd::WRITE_SECT:
14            if (d->cont != 0) {
15                hd::output_sect(d->punt);
16                d->punt += DIM_BLOCK;
17            }
18            break;
19         case hd::READ_DMA:
20         case hd::WRITE_DMA:
21            // qui si fa l'acknowledge, passo 7 della scaletta
22            bm::ack();
23            break;
24     }
25     if (d->cont == 0)
26         sem_signal(d->sincr);
27     wfi();
28 }
29 }
```