

Possibili domande aperte, prima parte

Prese da tutti compitini/appelli a disposizione

1_ Spiegare la differenza tra interruzioni multiple e interruzioni annidate, discutendo le differenti modalità di trattamento da loro richieste.

Per trattare le interruzioni si possono usare due approcci, il primo è disabilitare gli interrupt durante l'elaborazione di un interrupt, questi interrupt disabilitati rimarranno pendenti fino alla riabilitazione da parte della CPU e verranno gestiti in sequenza senza quindi tener conto di eventuali priorità da parte di qualcuno di essi. Il secondo invece consiste nel definire delle priorità e consentire quindi ad un interrupt di priorità maggiore di essere trattato subito, interrompendo quindi la gestione di un interrupt a priorità inferiore.

2_ Descrivere caratteristiche e le operazioni principali delle memorie a semiconduttore, considerando sia memorie RAM che ROM.

L'elemento base delle memorie a semiconduttore è la cella di memoria che presenta specifiche proprietà: due stati stabili, che rappresentano il bit 0 e 1, la possibilità di scrivere nella cella per impostare lo stato, la possibilità di leggere lo stato, inoltre tutti i tipi di memoria a semiconduttore sono ad accesso casuale (si accede alla parola tramite circuito dedicato). Il più comune è la memoria RAM la cui caratteristica principale è la possibilità di leggere e scrivere dati in e da memoria in modo semplice e rapido tramite segnali elettrici. Tale memoria è inoltre definita volatile per il fatto che necessita di alimentazione costante se si vuole preservare i dati ed è per questo quindi che la RAM viene utilizzata solo per una memorizzazione temporanea. Essa si divide in 2 tipologie SRAM e DRAM. La DRAM acronimo di Dynamic RAM è composta di celle che memorizzano i dati sotto forma di cariche nei condensatori. L'assenza o presenza di cariche determina lo stato 0 o 1. Poiché i condensatori tendono a scaricarsi naturalmente dopo un certo periodo di tempo le RAM necessitano di un refresh periodico e da ciò deriva l'aggettivo "dinamiche" per il fatto quindi che la carica scompare dinamicamente anche in presenza di alimentazione. Le SRAM utilizzano invece gli stessi elementi base di un processore, infatti i valori binari vengono memorizzati mediante porte logiche e diversamente dalle DRAM non necessitano di un refresh delle cariche. Un altro tipo di memoria ad accesso casuale è la ROM (read only memory) la quale presenta uno schema di dati predefinito e non modificabile. Tali memorie non sono volatili e mantengono quindi i dati anche in assenza di alimentazione, però questo tipo di dispositivi possono essere solo letti. Il vantaggio di queste memorie è che i dati e programmi di piccole dimensioni possono essere scritti in esse e non devono essere caricati da dispositivi esterni. In esse inoltre, i dati vengono inseriti durante il processo di fabbricazione processo costoso che richiede poi grande precisione poiché l'errata scrittura di un bit causa la perdita di tutto il lotto di ROM. Esse possono dividersi in PROM: una ROM programmabile che viene utilizzata qualora sia necessario un particolare contenuto, non è volatile ed è scrivibile solo una volta. Un'altra variante di queste memorie è la memoria principalmente di lettura utile quando le operazioni di lettura sono molto più frequenti di quelle di scrittura ed è richiesta una memoria non volatile. Si dividono in 3 tipi: EPROM, EEPROM, flash.

Le EPROM sono memorie cancellabili e programmabili otticamente e vengono lette e scritte elettricamente come le PROM. Prima della scrittura le celle di memoria devono essere cancellate e portate tutte allo stesso stadio mediante esposizione ai raggi UV. Le EEPROM sono memorie cancellabili e programmabili elettricamente, in esse è possibile scrivere in qualunque momento senza cancellare i dati, aggiornando solo i byte indirizzati (la scrittura richiede molto più tempo della lettura). Tali memorie sono più costose e meno dense delle EPROM. Le memorie flash invece sono chiamate così per la velocità con la quale possono essere riprogrammate. Così come le EEPROM adottano un sistema di

cancellazione elettrica e la sua memoria può essere eliminata molto più velocemente di una EPROM ed è inoltre possibile cancellare solo alcuni blocchi piuttosto che tutta la memoria.

3_ Nel contesto di una gerarchia di memoria, spiegare i possibili modi di realizzazione del mapping dei blocchi, discutendo vantaggi e svantaggi.

I possibili metodi di mapping dei blocchi in una gerarchia di memoria sono tre: diretto (direct mapping), associativo (n-way set) e set associativo (n-way set associative).

L'indirizzamento diretto è la tecnica più semplice in quanto assegna ad ogni blocco una sola possibile linea di cache. Per accedere alla cache poi, ogni indirizzo in memoria centrale viene diviso in tre campi: tag, linea, parola. Questo metodo di traduzione si distingue per la semplicità con cui quest'ultima avviene da indirizzo ILI (memoria) ad ILS (cache) e per la veloce determinazione di hit o miss. D'altro canto necessita di contraddistinguere il blocco presente in cache tramite un'etichetta (tag) e porta inoltre ad un numero elevato di swap per accedere ai dati di blocchi adiacenti. Il secondo metodo invece supera lo svantaggio dell'indirizzamento diretto poiché ogni blocco della memoria centrale può essere caricato su qualsiasi linea della cache. Così facendo l'indirizzo di memoria presenta solamente due campi: tag e parola e garantisce nel complesso una massima efficienza di allocazione. L'unico svantaggio è dato dalla complessità circuitale richiesta per esaminare in parallelo i tag di tutte le linee di cache. L'indirizzamento set-associativo invece è un compromesso che unisce i punti di forza dei precedenti riducendo i loro svantaggi. Qui la cache è suddivisa in insiemi (set) di k linee e il blocco può essere assegnato a qualunque linea dell'insieme, l'indirizzo in memoria è suddiviso nei campi tag, set e parola. Questo tipo di indirizzamento vanta una buona efficienza di allocazione a fronte di una discreta complessità di ricerca.

4_ Descrivere gestione dell'I/O tramite DMA.

Il DMA (direct memory access) è un modulo hardware aggiuntivo che sostituisce la CPU per la maggior parte delle attività di I/O. Esso si occupa del trasferimento dei dati da o verso la memoria senza passare attraverso il processore e invia a quest'ultimo un segnale di interrupt quando il trasferimento è completato. In questo modo la CPU è coinvolta solo all'inizio e alla fine del trasferimento e nel frattempo può eseguire altre attività. Prima di delegare una determinata operazione di I/O al DMA, il processore gli comunica le seguenti informazioni: lettura/scrittura, indirizzo dispositivo interessato, indirizzo iniziale in memoria del blocco dati coinvolto nell'operazione, quantità di dati da trasferire. Il DMA è connesso al bus di sistema e può accedere al canale dati in 2 modi. Una parola alla volta, sottraendo di tanto in tanto alla CPU il controllo del canale (cycle stealing) o per blocchi, prendendo in possesso il canale per una serie di trasferimenti (burst mode). Una configurazione ideale per consumare meno cicli di bus è integrare le funzioni di DMA e I/O. In questo modo il DMA utilizza il bus solo per scambiare dati con la memoria, mentre comunica in altri modi con i moduli I/O.

5_ Spiegare a cosa serve il bus di sistema, com'è strutturato e in che modo viene utilizzato in un calcolatore.

Il Bus di sistema è un mezzo di comunicazione condiviso che connette le componenti di un calcolatore e che permette lo scambio di dati fra di esse. Esso è composto da più linee (che va da 50 a 100) che trasmettono in parallelo un dato (sotto forma di segnali che indicano i bit 0 e 1) e tali linee vengono classificate in tre gruppi funzionali: dati, indirizzi e controllo. Le linee dati forniscono il percorso su cui viaggiano i dati tra i moduli del sistema e la loro ampiezza (ossia il numero di linee) varia da 32 a qualche centinaio ed è indispensabile per le prestazioni del sistema. Le linee indirizzi invece indicano la sorgente o la destinazione dei dati e per queste linee, l'ampiezza determina la quantità totale di memoria di un sistema. Infine le linee di controllo vengono utilizzate per controllare

l'accesso e l'uso delle linee dati e indirizzi. Il bus in un calcolatore opera come segue: se un modulo desidera inviare dati a un altro deve prima di tutto ottenere l'uso del bus e poi trasferire i dati mentre se desidera richiedere i dati deve, dopo aver ottenuto l'uso del bus, prima trasferire una richiesta all'altro modulo sulle appropriate linee di indirizzo e controllo e poi attendere l'invio dei dati.

6_ Come funziona il codice di correzione Hamming? Dare un esempio concreto di codifica nel caso di memorizzazione di un insieme di 4 bit.

Quando i dati devono venire memorizzati si esegue un calcolo su di essi per produrre un codice che verrà memorizzato assieme ai dati. Quando si deve leggere una parola tale codice verrà impiegato per rilevare eventuali errori negli M bit di dati. Un nuovo codice è generato a partire dagli M bit dati e confrontato con i K bit precedentemente prelevati. Il confronto ottiene 3 risultati: nessun errore viene rilevato, viene rilevato un errore ed è possibile correggerlo, perciò i bit dati e i bit per la correzione vengono inviati ad un correttore che produce un insieme corretto di M bit da emettere, viene rilevato un errore ma non è possibile correggerlo, ciò viene segnalato. Il codice a correzione d'errore più semplice è quello di hamming. In un insieme di parole di 4 bit utilizziamo il diagramma di eulero venn per dividere il piano in 3 cerchi che determinano 7 regioni, i 4 bit dati verranno assegnati agli scomparti interni. I restanti scomparti vengono riempiti con i cosiddetti bit di parità. Ciascun bit di parità è scelto in modo che il numero totale di 1 nel proprio cerchio sia pari. Ora se avvenisse una modifica dei bit dati l'errore sarebbe facilmente rilevato e potrebbe essere corretto cambiando il determinato bit che lo causa.

7_ Nel contesto di una gerarchia di memoria, discutere la modalità e la granularità del trasferimento delle informazioni fra i vari livelli della gerarchia.

E' possibile organizzare gerarchicamente le memorie. Scendendo lungo la gerarchia, troviamo un decrescente costo per bit, una capacità crescente, un tempo di accesso crescente e una decrescente frequenza di accesso da parte del processore. In questo modo, memorie più piccole, più costose e più veloci sono integrate da memorie più grandi, economiche e più lente. La CPU utilizza direttamente il livello più alto della gerarchia. Ogni livello inferiore, deve contenere tutti i dati presenti ai livelli superiori (più altri). La dimensione di un blocco e' la quantità minima indivisibile di dati che occorre prelevare (e copiare) dal livello inferiore. L'indirizzo di un dato diviene l'indirizzo del blocco che lo contiene, sommato alla posizione del dato all'interno del blocco.

8_ Discutere il modo in cui le informazioni sono organizzate in un CD-ROM.

Il CD è un dispositivo ottico non cancellabile che sfrutta una serie di pozzetti (pit) per memorizzare dati binari su una superficie di policarbonato. Tali informazioni, scritte da un laser ad alta intensità, vengono poi recuperate da un laser a bassa potenza. Le informazioni sono organizzate su una traccia a spirale che inizia vicino al centro e si svolge verso il bordo del disco. I settori in prossimità dell'estremità del disco sono della stessa lunghezza di quelli interni. Le informazioni sono quindi impacchettate uniformemente sul disco in segmenti della stessa dimensione e questi vengono letti ruotando il disco a velocità variabile. I pit sono poi letti dal laser a velocità lineare costante. Il disco ruota più lentamente per gli accessi sul lato esterno rispetto a quelli vicini al centro.

9_ Descrivere in dettaglio il ciclo di esecuzione con trattamento delle interruzioni.

All'inizio di ogni ciclo esecutivo il processore legge un'istruzione dalla memoria. Dopo il prelievo di tale istruzione il processore incrementa il valore del registro PC (Program

Counter) in modo che la prossima istruzione eseguita sia quella posizionata all'indirizzo di memoria immediatamente successivo. L'istruzione viene prima caricata in un registro del processore noto come IR(instruction register) e poi analizzata per determinare il tipo di operazione da eseguire l'operando da utilizzare. In seguito si determina l'indirizzo dell'operando e subito dopo avviene la sua lettura. Infine viene eseguita l'operazione indicata nell'istruzione e il risultato viene scritto nella memoria e viene inviato alla periferica. Dopodichè il processore controlla se è avvenuto qualche interrupt, se non vi sono interrupt pendenti, il processore procede al ciclo di prelievo e legge la successiva istruzione del programma altrimenti opera come segue: sospende l'esecuzione del programma in esecuzione salvandone il contesto (memorizzando quindi la l'indirizzo della prossima istruzione ed altri dati rilevanti per l'attività corrente) imposta il PC all'indirizzo di partenza di una routine per la gestione dell'interrupt. Il processore poi procede al ciclo di fetch e legge la prima istruzione del programma di gestione dell'interrupt. Quando questa routine termina, il programma riprende l'esecuzione del programma utente dal punto dell'interruzione.

10_ Spiegare in dettaglio le differenze tra un modulo di memoria DRAM ed un modulo di memoria SRAM, discuterne vantaggi e svantaggi.

La DRAM acronimo di Dynamic RAM è composta di celle che memorizzano i dati sotto forma di cariche nei condensatori. L'assenza o presenza di cariche determina lo stato 0 o 1. Poichè i condensatori tendono a scaricarsi naturalmente dopo un certo periodo di tempo le ram necessitano di un refresh periodico e da ciò deriva l'aggettivo "dinamiche" per il fatto quindi che la carica scompare dinamicamente anche in presenza di alimentazione. Le SRAM utilizzano invece gli stessi elementi base di un processore, infatti i valori binari vengono memorizzati mediante porte logiche e diversamente dalle DRAM non necessitano di un refresh delle cariche. Sia le ram dinamiche che quelle statiche sono volatili, una cella di memoria dinamica però è più piccola e semplice di una statica perciò le DRAM sono più dense e meno costose ma richiedono circuiti aggiuntivi per il refresh periodico. In sostanza le SRAM sono più veloci delle DRAM e vengono usate per la memoria cache mentre le DRAM per la memoria centrale.

11_ Nel contesto di una gerarchia di memoria, spiegare come i "miss" possono essere categorizzati in diversi tipi e dire quali strategie, per ogni tipo, si possono adottare per diminuirne il numero. Discutere criticamente tali strategie.

I miss in una gerarchia di memoria possono essere caratterizzati in 3 tipi diversi: Miss di primo accesso, inevitabile e non riducibile, Miss per capacità insufficiente, quando la cache non può contenere tutti i blocchi necessari all'esecuzione del programma e miss per conflitto, quando più blocchi possono essere mappati (con associazione diretta o a gruppi) su uno stesso gruppo. Le tecniche di risoluzione classiche per i miss per capacità insufficiente possono essere una maggiore dimensione del blocco la quale è una buona tecnica per fruire di località spaziale che però causa un aumento di miss per conflitto (a causa del numero ridotto di blocchi disponibili) mentre per i miss per conflitto una soluzione efficiente può essere la maggiore associatività che causa però un incremento del tempo di localizzazione in gruppo ed è soggetta alla regola del 2:1 (cache di N blocchi stessa probabilità di miss di cache a N/2 con ass a 2 vie).

Altre tecniche di risoluzione possono essere l'adozione di una cache multilivello, la separazione di cache dati e cache istruzioni e l'ottimizzazione dei dati mediante compilatori che permettono le seguenti operazioni ossia il posizionamento accurato delle procedure ripetitive, la fusione di vettori in strutture (località spaziale) e la trasformazioni di iterazioni annidate (località spaziale)

12_ Discutere le ragioni per cui è stato sviluppato il sistema RAID. Inoltre si descriva in dettaglio il livello 0,2,4 del RAID.

Il sistema RAID (redundant array of independent disk) venne elaborato primaditutto per aumentare le prestazioni di un calcolatore grazie all'utilizzo di più dischi in parallelo, in questo modo si garantiva una più veloce gestione delle diverse operazioni di input/output e, con l'aggiunta della ridondanza, una maggiore affidabilità poichè i dati vengono distribuiti su più dischi e sono recuperabili in caso di guasto di uno dei dischi. Il livello 0 non include la ridondanza a favore delle prestazioni, i dati sono distribuiti sui vari dischi quindi se vi sono due richieste di blocchi che si trovano in dischi diversi esse vengono servite in parallelo. I dati vengono distribuiti in strisce (strips) che possono essere blocchi fisici, settori o altro e sono distribuite a rotazione (round robin) sui dischi. Il raid 2 sfrutta l'accesso parallelo ossia le testine di ciascun disco si trovano nella stessa posizione in ogni momento. Si usa lo striping dei dati con strisce molto piccole (singoli byte o parola), vengono generati codici di correzione errori di Hamming memorizzati in più dischi, perciò se si riscontra un errore il controllore può riconoscere e correggere l'errore istantaneamente in modo che il tempo di accesso per la lettura non venga rallentato (tale configurazione non viene commercializzata). Nel raid 4 invece ogni disco opera indipendentemente, così facendo le richieste di I/O possono essere gestite contemporaneamente. Usa lo striping dei dati con strisce grandi e utilizza una striscia di parità bit a bit sulle corrispondenti strisce di ciascun disco dati, e i bit di parità sono memorizzati nella corrispondente striscia sul disco di parità.

13_ Descrivere in dettaglio il ciclo completo di fetch/execute delle istruzioni.

(**Fetch**) All'inizio di ogni ciclo esecutivo il processore legge un'istruzione dalla memoria. (**Calcolo indirizzo istruzione**) Dopo il prelievo di tale istruzione il processore incrementa il valore del registro PC (Program Counter) in modo che la prossima istruzione eseguita sia quella posizionata all'indirizzo di memoria immediatamente successivo. L'istruzione viene prima caricata in un registro del processore noto come IR (instruction register) e poi analizzata per determinare il tipo di operazione da eseguire l'operando da utilizzare (**Decodifica istruzione**). In seguito si determina l'indirizzo dell'operando e subito dopo avviene la sua lettura (**Calcolo indirizzo operando/lettura operando**). Infine viene eseguita l'operazione indicata nell'istruzione e il risultato viene scritto nella memoria e viene inviato alla periferica (**Operazione sui dati/memorizzazione risultato**).

14_ Descrivere in dettaglio le memorie DRAM.

La DRAM acronimo di Dynamic RAM è composta di celle che memorizzano i dati sotto forma di cariche nei condensatori. L'assenza o presenza di cariche determina lo stato 0 o 1. Poichè i condensatori tendono a scaricarsi naturalmente dopo un certo periodo di tempo le ram necessitano di un refresh periodico e da ciò deriva l'aggettivo "dinamiche" per il fatto quindi che la carica scompare dinamicamente anche in presenza di alimentazione. Per le operazioni di scrittura si applica tensione alla linea di bit (che a seconda dell'intensità determina se bit è 1 o 0) successivamente si applica segnale alla linea degli indirizzi trasferendo la carica al condensatore. Per quelle di lettura invece prima di tutto si seleziona la linea indirizzo poi la carica viene convogliata in una linea di bit collegata ad un amplificatore che confronta la tensione con un valore di riferimento e determina se la cella contiene il bit 1 o 0 e infine si applica un refresh per ripristinare la carica nel condensatore.

15_ Descrivere in dettaglio i livelli 1 e 4 del RAID.

In raid 1 la ridondanza, diversamente dagli altri raid, viene ottenuta duplicando tutti i dati, questo implica un numero doppio di dischi fissi. Si usa lo striping dei dati e ogni striscia si trova fisicamente su due dischi quindi una lettura può essere soddisfatta dal disco al quale si accede più velocemente e che presenta la minor latenza di rotazione, la scrittura

avviene su entrambi i dischi quindi le prestazioni sono condizionate dalla più lenta delle due scritture. In caso di guasto, i dati sono disponibili nell'altro disco, questo permette un'alta affidabilità ma un costo elevato.

Nel raid 4 invece ogni disco opera indipendentemente, così facendo le richieste di I/O possono essere gestite contemporaneamente. Usa lo striping dei dati con strisce grandi e utilizza una striscia di parità bit a bit sulle corrispondenti strisce di ciascun disco dati, e i bit di parità sono memorizzati nella corrispondente striscia sul disco di parità. RAID 4 risulta penalizzato dalle richieste di scrittura di piccola dimensione.

16_ Descrivere l'organizzazione e formattazione dei dati nei dischi rigidi.

Nei dischi magnetici la testina è in grado di leggere e scrivere su una porzione del disco rotante. Ciò origina la disposizione fisica dei dati in anelli concentrici, chiamati tracce (track). Le tracce hanno la stessa larghezza della testina, e ne esistono migliaia per ciascun piatto. Tracce adiacenti sono separate da spazi (gaps). Ciò minimizza gli errori dovuti al disallineamento della testina o all'interferenza tra i campi magnetici. Il trasferimento dati avviene per settori, che generalmente sono un centinaio per traccia, di lunghezza fissa o variabile (odiernamente ammonta a 512 byte), i settori adiacenti sono inoltre separati da spazi. I bit più vicini al centro del disco ruotano attorno al punto fisso, come la testina, più lentamente dei bit esterni, questa velocità viene quindi compensata per permettere alla testina di leggere tutti i bit alla stessa velocità facendo ruotare il disco a velocità angolare costante. La formattazione invece è l'operazione con la quale si prepara il disco per renderlo idoneo all'archiviazione e consiste ad esempio nell'inserimento di un criterio che determini la posizione di un dato settore il suo inizio e la sua fine e un punto di partenza sulla traccia. Questi requisiti sono rispettati tramite dati di controllo memorizzati sul disco con i quali esso viene inizializzato.

17_ Descrivere in dettaglio la gestione da programma I/O.

Nell'I/O da programma i dati vengono scambiati tra processore e modulo I/O. Quando il processore sta eseguendo un programma e incontra un'istruzione correlata con l'I/O, esso esegue l'istruzione inviando un comando al modulo di I/O appropriato. Nell'I/O da programma, il modulo eseguirà l'azione richiesta e imposterà i bit appropriati nel suo registro di stato. Il modulo non intraprende nessun'altra azione per allertare il processore. In particolare, esso non interrompe il processore, perciò il processore periodicamente controlla lo stato del modulo I/O finché non rileva il completamento dell'operazione. Ci sono 4 tipi di comandi che il processore può inviare al modulo I/O e sono Controllo (avvia una periferica e le dice cosa fare), test (testa le condizioni di stato dei moduli di I/O), lettura (ottiene i dati dalla periferica attraverso il modulo I/O) e scrittura (impone al modulo di trasmettere tramite bus i dati alla periferica).

18_ Descrivere in cosa consiste l'architettura Von Neumann.

La progettazione di quasi tutti i calcolatori odierni è basata su concetti sviluppati da John Von Neumann. Essa è nota come Architettura di Von Neumann ed è basata su 3 concetti chiave:

- 1: i dati e le istruzioni risiedono in un'unica memoria di lettura e scrittura
- 2: i contenuti di questa memoria sono accessibili per indirizzo, indipendentemente dal tipo di informazione rappresentata;
- 3: l'esecuzione avviene in modo sequenziale, da un'istruzione a quella immediatamente successiva.

Per eseguire un programma, possiamo costruire i componenti logici in modo che il risultato sia quello voluto. Questo è il modo di costruire un "programma cablato", cioè in forma hardware, che non può essere modificato. Esso è un sistema non flessibile, che può

eseguire solo istruzioni determinate. Tramite circuiti generici accetta segnali di controllo e produce risultati. Per nuovi programmi basta quindi dare i giusti segnali di controllo. La programmazione invece è molto più facile perché invece che ridefinire ogni volta l'hardware basta fornire una nuova sequenza di codici (ossia una nuova istruzione). Questo processo viene detto "programmazione software".

19_ Nel contesto di una gerarchia di memoria spiegare perché la memoria viene suddivisa in blocchi e, relativamente alle prestazioni della cache, discutere pregi e difetti dell'adozione di una dimensione di blocco elevata.

Per realizzare un'organizzazione gerarchica della memoria, che soddisfi i parametri di velocità, ampiezza e costo, conviene suddividere la memoria in blocchi. La dimensione di un blocco è la quantità minima indivisibile di dati che occorre prelevare (copiare) dal livello inferiore. L'indirizzo di un dato diviene l'indirizzo del blocco che lo contiene sommato alla posizione del dato all'interno del blocco.

Se si adotta un blocco con dimensione elevata si guadagna in termini di costi e capacità ma si perde in termini di velocità.

20_ Spiegare cosa sono gli errori soft, come ovviare a tali errori e fare eventualmente un esempio.

Le memorie primarie sono soggette ad errori. Questi possono essere classificati in guasti hardware, che sono permanenti, e guasti software, chiamati anche "soft error", che sono casuali e non distruttivi. Essi infatti alterano i contenuti di una o più celle di memoria, senza però danneggiarla fisicamente. Gli errori "soft" possono essere rilevati ed eventualmente corretti usando ad esempio il codice correttore di hamming. Esempio nel caso di memorizzazione di un insieme a 4 bit.

Esso funziona tramite bit di parità, ovvero, con insiemi da 4 bit, si hanno 3 cerchi che si intersecano e all'interno di ogni cerchio il numero di bit complessivo deve essere pari, se no viene rilevato l'errore il quale può essere eventualmente corretto cambiando il numero di bit.

21_ Descrivere le politiche di scrittura write through e write allocate evidenziando differenze, vantaggi e svantaggi di entrambe.

Write through:

Ogni dato modificato nella cache viene contemporaneamente modificato nella memoria centrale. In questo modo i dati sono sempre coerenti tra i vari livelli di memoria. Lo svantaggio principale però è che per frequenti scritture sul medesimo blocco si verifica un aumento di traffico nel bus con conseguente collo di bottiglia.

Write back:

La scrittura in memoria centrale avviene solo quando il corrispondente blocco in cache viene rimpiazzato. Ciò consente un'ottimizzazione del traffico tra livelli ma causa periodi di incoerenza tra di essi, inoltre occorre sempre ricordare se sono avvenute operazioni di scrittura nel blocco tramite "dirty bit". Il problema di questa tecnica è che parti della memoria centrale non sono aggiornate, e dunque gli accessi tramite moduli I/O possono essere consentiti solo attraverso la cache.

22_ Descrivere in che modo vengono gestite le interruzioni (sia per la componente hardware che per quella software) nel caso di I/O input driven.

Per evitare che il processore debba costantemente verificare lo stato del dispositivo di I/O esso invia al modulo un comando e quest'ultimo interromperà il processore una volta pronto allo scambio di dati, dopodiché il processore esegue lo scambio e riprende poi l'elaborazione interrotta. A livello hardware, quando una periferica completa un'azione di I/O avviene la seguente sequenza di eventi: il dispositivo invia un interrupt al processore, il

processore conclude l'operazione corrente e poi controlla l'eventuale presenza di interrupt e se c'è riscontro, invia un segnale di riconoscimento al dispositivo che ha inviato tale interrupt che rimuoverà quindi il proprio segnale. Il processore poi salva il contesto del programma ponendo PSW (ovvero lo stato del processore) e la locazione della prossima istruzione da eseguire contenuta nel PC in cima alla pila di sistema. Infine il processore scrive nel PC l'indirizzo della prima istruzione della routine di gestione dell'interrupt considerato. Dopodiché il controllo viene trasferito al programma di gestione dell'interrupt che procederà nel modo seguente: vengono salvate le restanti informazioni di stato del processo dai registri PSW e PC, successivamente avviene l'elaborazione dell'interrupt che può includere un esame dell'informazioni sullo stato delle operazioni di I/O oppure comporta un invio di ulteriori comandi al dispositivo di I/O. Quando l'elaborazione dell'interrupt è completa, i valori dei registri salvati vengono recuperati dalla pila e riportati nei registri ed infine i valori di PSW e PC vengono ripristinati.

23_ Descrivere la differenza tra architettura e organizzazione di un calcolatore.

L'architettura di un calcolatore si riferisce agli attributi di un sistema visibili al programmatore o, in altre parole, agli attributi che hanno un impatto diretto sull'esecuzione logica di un programma. L'organizzazione invece fa riferimento alle unità operative e alle loro interconnessioni che realizzano le specifiche architetturali. Attributi architetturali sono il repertorio delle istruzioni, il numero di bit usati per rappresentare i vari tipi di dato, i meccanismi di I/O e le tecniche di indirizzamento della memoria mentre quelli organizzativi includono aspetti hardware trasparenti al programmatore (come i segnali di controllo, le periferiche e la tecnologia delle memorie).

24_ Nel contesto di una gerarchia di memoria, spiegare come funziona la politica di scrittura write back. Discutere criticamente i problemi che possono insorgere nell'adottarla.

La politica di scrittura write back è una tecnica alternativa alla write through che evita il collo di bottiglia creato da quest'ultima minimizzando le scritture in memoria e applicando gli aggiornamenti solo nella cache. Il problema di questa tecnica è che parti della memoria centrale non sono aggiornate, e dunque gli accessi tramite moduli I/O possono essere consentiti solo attraverso la cache. Ciò richiede circuiti complicati e costituisce un potenziale collo di bottiglia. Oltre a ciò in un'organizzazione a bus in cui più di un dispositivo è connesso e dove la memoria centrale è condivisa un'alterazione dei dati in una cache può invalidare tutti i dati nella memoria centrale e anche quelli nelle altre cache eventualmente connesse al bus. Per ovviare al problema esistono 3 approcci: un monitoraggio del bus con write through dove i controllori della cache osservano le linee di indirizzi ed invalidano eventuali scritture in una locazione di memoria condivisa da parte di qualche gestore del bus qualora tale locazione di memoria sia già residente nella cache. Una trasparenza hardware che permetta mediante un hardware aggiuntivo un aggiornamento costante della memoria centrale e di tutte le cache presenti e una memoria *noncachable*, ovvero una porzione condivisa nella quale gli accessi sono dei cache miss dato che essa non viene mai copiata nella cache. Essa può essere identificata via hardware o tramite indirizzi riservati.

