
ANNO ACCADEMICO 2022/2023

Sistemi Operativi

Teoria

Altair's Notes



DIPARTIMENTO DI INFORMATICA

| CAPITOLO 1 | INTRODUZIONE | PAGINA 5 |
|-------------------|---|-----------------|
| 1.1 | Idee base Eventi — 6 • Architetture — 6 • Caratteristiche di un SO — 7 | 6 |
| 1.2 | Gestione Generale del SO Programmi e Processi — 8 • Gestione del Sistema — 9 • SO come Ambiente di Lavoro — 10 | 8 |
| 1.3 | Gestione dei Processi Dettagli sui Processi — 13 • Commutazione e Scheduling — 14 • Algoritmi di Scheduling — 15 • Creazione e Terminazione — 17 | 12 |

| CAPITOLO 2 | SINCRONIZZAZIONE | PAGINA 19 |
|-------------------|--|------------------|
| 2.1 | Processi Cooperanti Introduzione — 19 • Scambi di Messaggi — 20 | 19 |
| 2.2 | Thread Heavyweight Process — 22 • Tipi di Thread — 23 • Lightweight Process — 24 | 22 |
| 2.3 | Implementare la Sincronizzazione Scheduling della CPU — 25 • Altri tipi di sincronizzazione: — 27 | 25 |

| CAPITOLO 3 | GESTIONE DELLA SINCRONIZZAZIONE | PAGINA 30 |
|-------------------|---|------------------|
| 3.1 | Semafori Introduzione — 30 • Problemi Classici — 32 | 30 |
| 3.2 | Monitor Introduzione — 35 | 35 |
| 3.3 | Transazioni e Lock Transazioni — 36 • Lock — 37 | 36 |
| 3.4 | Deadlock Introduzione — 38 • Prevenzione del Deadlock — 39 • Deadlock Avoidance — 40 • Rottura del Deadlock — 41 | 38 |

| CAPITOLO 4 | GESTIONE DELLA MEMORIA | PAGINA 43 |
|-------------------|---|------------------|
| 4.1 | Introduzione Indirizzi, Binding e Loading — 43 • Allocazione della RAM — 44 | 43 |
| 4.2 | Paginazione Struttura — 47 • Architettura di Paginazione — 49 • Segmentazione — 51 • Memoria Virtuale — 52 • Paginazione su Richiesta — 55 | 47 |
| 4.3 | Allocazione e Thrashing Allocazione — 58 • Thrashing — 59 • Pagine Attive e Working Set — 59 | 58 |

| | | |
|-----|---|----|
| 5.1 | Introduzione | 63 |
| | Cos'è il Filesystem? — 63 • Organizzazione del Filesystem — 64 | |
| 5.2 | Directory | 65 |
| | Che Cos'è una Directory? — 65 • Tipi — 66 | |
| 5.3 | Aprire un File | 69 |
| | Inode — 69 • Apertura in UNIX — 70 | |
| 5.4 | Allocazione | 70 |
| | Implementazione delle Directory — 70 • Spazio Disco ai Files — 71 | |
| 5.5 | Gestione e Implementazione | 73 |
| | Gestione dello Spazio Libero — 73 • Implementazione del Filesystem — 73 | |

Premessa

Licenza

Questi appunti sono rilasciati sotto licenza Creative Commons Attribuzione 4.0 Internazionale (per maggiori informazioni consultare il link: <https://creativecommons.org/version4/>).



Formato utilizzato

Box di "Concetto sbagliato":

Concetto sbagliato 0.1: Testo del concetto sbagliato

Testo contenente il concetto giusto.

Box di "Corollario":

Corollario 0.0.1 Nome del corollario

Testo del corollario. Per corollario si intende una definizione minore, legata a un'altra definizione.

Box di "Definizione":

Definizione 0.0.1: Nome delle definizioni

Testo della definizione.

Box di "Domanda":

Domanda 0.1

Testo della domanda. Le domande sono spesso utilizzate per far riflettere sulle definizioni o sui concetti.

Box di "Esempio":

Esempio 0.0.1 (Nome dell'esempio)

Testo dell'esempio. Gli esempi sono tratti dalle slides del corso.

Box di "Note":

Note:-

Testo della nota. Le note sono spesso utilizzate per chiarire concetti o per dare informazioni aggiuntive.

Box di "Osservazioni":

Osservazioni 0.0.1

Testo delle osservazioni. Le osservazioni sono spesso utilizzate per chiarire concetti o per dare informazioni aggiuntive. A differenza delle note le osservazioni sono più specifiche.

1

Introduzione

Domanda 1.1

Che cosa sono i sistemi operativi?

Definizione 1.0.1: Sistema Operativo

Permette di usare l'hardware attraverso programmi. È principalmente un gestore di risorse, ma si occupa anche di controllare l'esecuzione dei programmi applicativi.

Le risorse costituiscono il sistema e consentono la risoluzione di problemi. Possono essere di 2 tipi:

- *Hardware*: CPU, RAM, I/O...
- *Software*: code di messaggi, processi, thread...

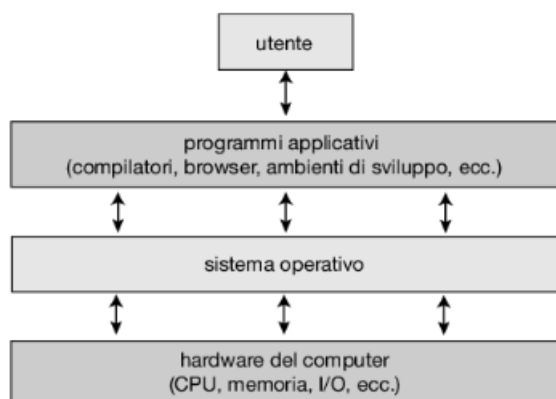


Figure 1.1: Schema generale di un computer.

1.1 Idee base

Definizione 1.1.1: Astrazioni

Il SO definisce delle “astrazioni”, ovvero delle rappresentazioni dei vari elementi del SO stesso, in modo che siano di facile comprensione da parte dell’utente. In generale il SO definisce delle astrazioni software di tutti i tipi di oggetti che occorre rappresentare e gestire per far funzionare un computer (compresi gli utenti), implementa opportuni algoritmi di controllo e contiene interfacce sia verso gli utenti sia verso i dispositivi fisici.

1.1.1 Eventi

Definizione 1.1.2: Program Event-Driven

Sono dei programmi che eseguono diverse routine basandosi su determinate situazioni chiamate eventi. A ogni evento corrisponde una routine di gestione predefinita detta “event-handler”.

Note:-

La gestione degli eventi deve essere molto efficiente in un SO.

Corollario 1.1.1 Dispatch

Gli eventi possono essere accumulati in una coda, e quando vengono gestiti si esegue un “dispatch”, che consiste nell’individuare ed eseguire la routine associata all’evento.

Note:-

È possibile implementare il dispatch attraverso un vettore di puntatori agli event handler: basta associare a ogni evento un numero (id dell’evento e al contempo indice nel vettore). Il vettore è detto vettore delle interruzioni (interrupt vector).

La gestione degli interrupt causa un context switch:

1. SO sospende il processo e salva le sue informazioni in RAM.
2. Carica nei registri CPU le informazioni necessarie per l’esecuzione dell’handler.

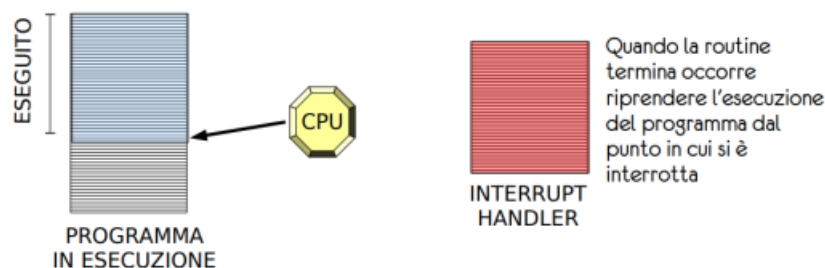


Figure 1.2: Context switch.

1.1.2 Architetture

Gli elaboratori possono essere divisi in 2 categorie:

- *Monoprocessore (Single Core)*: possiedono una singola CPU, ma possono presentare processori ausiliari dedicati ad attività specifiche (es: controller).

- **Multiprocessore (Multi Core):** posseggono due o più CPU e godono di maggiore capacità elaborativa, maggiore affidabilità e gestione della memoria più efficiente, ma possiedono un costo di risorse HW più elevato legato alla gestione e sincronizzazione dei processori.

Corollario 1.1.2 Cluster di Elaboratori

Sono sistemi multiprocessore costituiti da insiemi di elaboratori completi (chiamati nodi). Ogni nodo può essere single o multi-core. C'è bisogno di un software che gestisce il cluster, per controllare e sistemare i possibili crash e per dividere il carico sui nodi.

Definizione 1.1.3: Memoria Centrale (RAM)

E' l'unica memoria di grandi dimensioni direttamente accessibile dalla CPU. Queste interazioni avvengono tramite istruzioni load (RAM \rightarrow registro) e store (registro \rightarrow RAM).

Note:-

È molto più veloce di una memoria secondaria, ma è anche volatile, ovvero viene cancellata allo spegnimento della macchina, e ha una capacità limitata a causa dell'elevato costo di produzione.

Definizione 1.1.4: Memoria Secondaria

È una memoria permanente, con capacità elevata. In genere si tratta di un disco magnetico, ma gli SSD, nastri, CD/DVD/BD. Si differenziano in base al costo e la velocità di accesso.

I dispositivi di input/output vengono gestiti da un certo numero di controller, tutti connessi da un bus. Ogni controller:

- Gestisce uno o più dispositivi.
- Ha una memoria interna (buffer e registri).
- Ha un software chiamato driver che si interfaccia con il SO.
- Per operazioni più rapide utilizzano DMA (Direct Memory Access).

1.1.3 Caratteristiche di un SO

Definizione 1.1.5: Scheduling della CPU

- Multiprogrammazione: il SO gestisce contemporaneamente un insieme di processi (job, task) distribuendo l'utilizzo della CPU fra i vari processi.
- Multitasking: estensione della multiprogrammazione in cui si tiene conto dell'interazione con l'utente, facendo in modo che esso abbia la percezione che solo il suo job sia in esecuzione (parallelismo virtuale).

Note:-

La gestione di quali processi spostare in memoria centrale e viceversa è detta job scheduling.

Definizione 1.1.6: Dual Mode

Politica di protezione del SO che permette ai processi di utilizzare un Instruction Set "sicuro", che non li lascia eseguire modifiche drastiche alla memoria (i processi possono solo accedere all'area di memoria a loro assegnata).

In dual mode, tramite un **bit di validità**, è possibile decidere quale IS utilizzare:

- 0 → Modalità **kernel** (o supervisore): accesso all'intero IS.
- 1 → Modalità **utente**: non si ha accesso alle istruzioni di I/O.

Note:-

Al bootstrap, il bit di modalità è inizializzato a 0. Con la richiesta di esecuzione di processi utente, a seconda del processo, il bit di modalità cambierà valore.

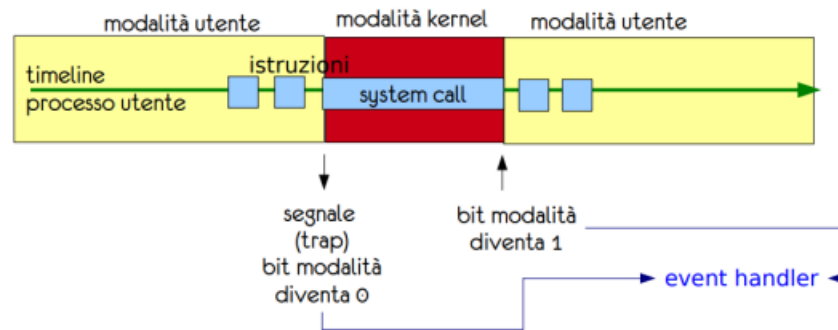


Figure 1.3: Dual Mode.

Definizione 1.1.7: System Call

Comandi che possono essere utilizzati dai processi per eseguire indirettamente (le op. non le esegue il processo ma il SO) operazioni di I/O. Durante l'esecuzione di una system call, il bit di modalità vale 0.

1.2 Gestione Generale del SO

1.2.1 Programmi e Processi

Definizione 1.2.1: Processo

Un programma in sé è un file, un'entità passiva che non fa nulla. La sua esecuzione produce un'entità attiva detta processo.

Osservazioni 1.2.1

- Un programma può avere diversi processi che possono portare ad eseguire task diversi, ognuno con una propria area di memoria che può essere anche condivisa.
- Di ogni processo viene mantenuta un'astrazione.

Il SO deve gestire i processi:

- Identificazione.
- Creazione e cancellazione.
- Sincronizzazione e comunicazione.
- Identificare e gestire il deadlock.
- Evitare starvation.

Corollario 1.2.1 Terminazione di un Processo

La terminazione dei processi è gestita da un system call, e la terminazione anormale può causare la copiatura dell'immagine del processo (dump o core). In questa maniera si può eseguire il codice passo per passo, per vedere cosa sia successo.

Un processo può anche fare delle azioni:

- Caricare un programma diverso da quello originario.
- Generare un altro processo.
- può essere sospeso per un lasso di tempo o in attesa di un evento (per esempio, un segnale)

Osservazioni 1.2.2

- Il SO “convive” con gli altri processi utente, perciò anche esso ha bisogno di utilizzare la CPU.
- Il SO può riprendere la CPU assegnata ad un processo utente tramite l'utilizzo di timer predefiniti, che causano un interrupt del processo.

1.2.2 Gestione del Sistema

Le operazioni principali per la gestione dei file che il SO esegue sono:

- Creazione e cancellazione.
- Apertura e chiusura.
- Lettura e scrittura.
- Riposizionamento e Spostamento.
- Copiatura.
- Lettura e modifica delle proprietà.

Ci sono varie operazioni che si possono fare sui dispositivi:

- Request o release.
- Lettura, scrittura e riposizionamento.

Note:-

Sono system call che trasferiscono informazioni come la data e l'ora, informazioni sui processi o informazioni sui file all'utente o a suoi programmi.

Ci sono 2 modelli principali per la comunicazione tra processi:

- *Scambio di messaggi:*
 - Identificare macchina e processo.
 - Mittente: apertura connessione.
 - Destinatario: accettazione connessione.
 - Scambio messaggio.
 - Chiusura connessione.
- *Memoria condivisa:*
 - Allocazione di un'area di memoria condivisa.
 - Aggancio (attach) di un'area di memoria condivisa.

Gestione delle memorie:

- Memoria principale:
 - Ricordarsi chi sta usando quale parte di memoria e quale parte è invece libera.
 - Assegnare/revocare lo spazio a seconda delle necessità.
 - Decidere quali processi vanno trasferiti in RAM e quali vanno rimossi.
- Memoria secondaria:
 - Assegnazione dello spazio.
 - Scheduling del disco.
 - Swapping.

Osservazioni 1.2.3 Memoria e File System

- Gli utenti vedono la memoria organizzata in *file*, unità logiche di archiviazione di solito organizzati in *directory*.
- Quando gli utenti utilizzano i nomi dei file e i loro “cammini” (path), il SO deve essere in grado di identificare i blocchi di memoria ad essi corrispondenti: perciò ogni file mantiene delle proprietà che aiutano al SO ad identificarli (es: tipo di file, permessi di accesso, data di creazione/modifica...).

1.2.3 SO come Ambiente di Lavoro

Il SO è in grado di offrire diversi servizi:

- *Interfaccia utente (UI, User Interface):*
 - Linea di comando (CLI, Command User Interface).
 - Interfacce grafiche (GUI, Graphical User Interface).
- *Esecuzione di programmi:*
 - I programmi possono richiedere l'accesso a file e dispositivi I/O.
 - Deve essere possibile rilevare lo stato di terminazione di un programma.
- *Comunicazione tra processi:*
 - Memoria condivisa: processi lavorano in un'area comune che deve essere gestita in modo da evitare inconsistenze.
 - Scambio di messaggi (o pacchetti).
- *Protezione:*
 - Gli utenti possono limitare l'accesso alle proprie informazioni.
 - Non tutti gli utenti possono eseguire tutti i comandi o accedere a ogni file.

Note:-

In molti SO l'interfaccia utente non è altro che un programma che viene avviato all'atto del login.

Esistono due tipi di interfacce utente:

- CLI (Command-line interface) o shell: esegue un ciclo infinito in cui attende un comando, lo esegue, e torna di nuovo in attesa.
- GUI (Graphical User Interface): si basa su un desktop, su cui sono posti oggetti selezionabili e attivabili tramite puntatore.

Implementazione di un SO:

1. *Scelta di un linguaggio che supporti la portabilità:* mentre in passato i SO venivano scritti in linguaggio macchina (assembly), attualmente le CPU moderne sono troppo complesse per fare ciò, quindi vengono scritti in linguaggi di alto livello (es: Unix, Linux e Windows sono scritti in C).
2. *Definizione delle politiche di alto livello:* serve definire criteri e meccanismi per le funzioni che l'utente vorrà utilizzare.
 - Criterio (o politica): specifica un comportamento da seguire in una certa circostanza.
 - Meccanismo: strumenti (anche software) neutri rispetto ai criteri, che vengono quindi utilizzati per seguire i criteri.
3. *Scelta di un'architettura per il SO:* per poter essere mantenuto e aggiornato, un SO deve essere ben strutturato, anche per evitare problemi dovuti alla crescita di complessità del programma e/o degli applicativi.

Definizione 1.2.2: Stratificazione

Per questioni di sicurezza e modularità, il sistema ha una struttura "a cipolla"^a: lo strato più interno corrisponde all'HW, quello più esterno all'interfaccia utente. Su ogni strato esiste un oggetto astratto, che incapsula dei dati. Le operazioni che poi elaborano questi dati possono scegliere se rendere visibili/accessibili all'esterno o no.

^aCit. Shrek.

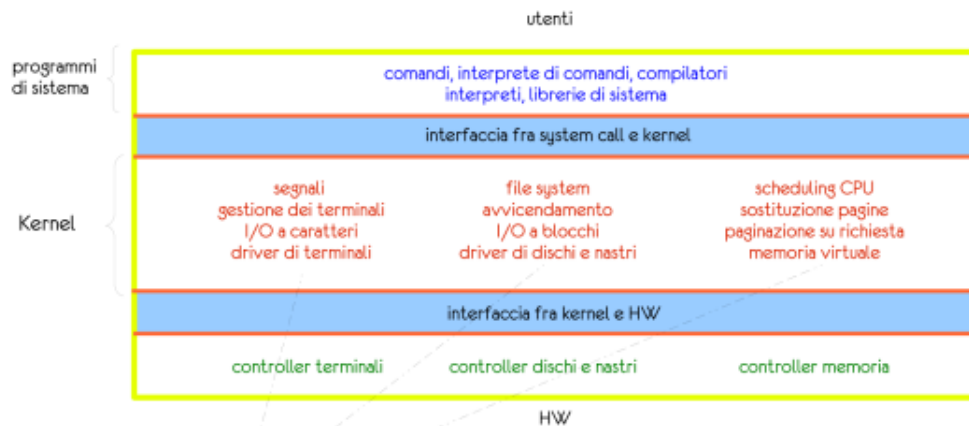


Figure 1.4: Struttura a strati.

Vantaggi e Svantaggi:

- Semplicità di progettazione: per realizzare uno strato basta sapere quali funzionalità ha a disposizione, non importa sapere come sono realizzate.
- Semplicità di debug e verifica del sistema: ogni strato usa solo funzionalità messe a disposizione dallo strato immediatamente inferiore, possiamo verificare gli strati uno per volta.
- Difficoltà: definire in modo opportuno gli strati, e quali funzionalità mettere in ognuno di loro.
- Minore efficienza in fase di esecuzione: ogni passaggio di stato comporta infatti la chiamata di una nuova funzione, da caricare, che può richiedere parametri, da caricare a loro volta...

Definizione 1.2.3: Macchine Virtuali

Le macchine virtuali sono un approccio alla stratificazione che implica la duplicazione del comportamento di ogni componente hardware del computer. Invece di installare un SO direttamente, installandolo su una macchina virtuale può lavorare su un computer che ha già un suo SO. I SO installati sulla macchina virtuale sono detti ospiti.

Note:-

Può essere utile per usare SO diversi sullo stesso computer (anche da utenti diversi) o verificare/debuggare applicativi su SO diversi dal proprio.

Definizione 1.2.4: Microkernel

Questo tipo di modularizzazione consiste nel rimuovere dal kernel tutto ciò che non è essenziale, spostandolo a livello di applicativo utente. I microkernel contengono i servizi minimi per la gestione di processi, comunicazione e memoria.

Vantaggi e svantaggi:

- Facilità di estensione: possibile aggiungere nuovi servizi senza modificare il kernel.
- Maggiore semplicità di adattamento a nuove architetture.
- Maggiore sicurezza: il kernel non viene direttamente affetto da cambiamenti.
- Possibile generare sovraccarichi quando processi utente vengono eseguiti con funzionalità di SO.
- Bottleneck: comunicazione indiretta.

Definizione 1.2.5: Tecnica a Moduli

L'approccio considerato il migliore adottato, dai SO moderni. A un kernel minimale possono essere aggiunti moduli nuovi, in fase di avvio ma anche di esecuzione. Aggiungere moduli significa aggiungere system call, ma anche capacità di gestire nuovo hardware.

Note:-

Come nei sistemi a microkernel, il kernel qui gestisce un nucleo essenziale, che viene espanso poi con moduli.

Vantaggi:

- Ogni modulo può usare qualsiasi altro modulo (*flessibilità*).
- La comunicazione tra moduli è diretta (*efficienza*).

1.3 Gestione dei Processi

Definizione 1.3.1: Processo

“Processo” è un’astrazione, una rappresentazione interna al SO, che consente di pensare e realizzare meccanismi quali multitasking, scheduling della CPU, protezione.

Osservazioni 1.3.1

- Il SO deve mantenere informazioni riguardo i diversi task: ogni task esegue un programma, elabora dei dati, ha un utente “proprietario”, può avere una priorità.
- Quando un task viene interrotto e ripreso occorre mantenere tutte le informazioni necessarie.

- Ogni processo ha una propria *sezione dati* composta dallo stack di esecuzione e dallo heap, ma in RAM viene anche memorizzato (per ogni processo) il *programma* in sé e il *program counter*, che indica la prossima istruzione da eseguire.

1.3.1 Dettagli sui Processi

Ogni processo può avere diversi stati che aiutano il SO ad eseguire scheduling e parallelismo:

- new: appena creato.
- running: in esecuzione.
- waiting: in attesa di un evento (es: ricevere un dato, completamento di un'operazione I/O, sospensione volontaria).
- ready: aspetta l'assegnazione della CPU (può passare running → ready con un interrupt).
- terminated: cessata esecuzione (exit, abort, kill).

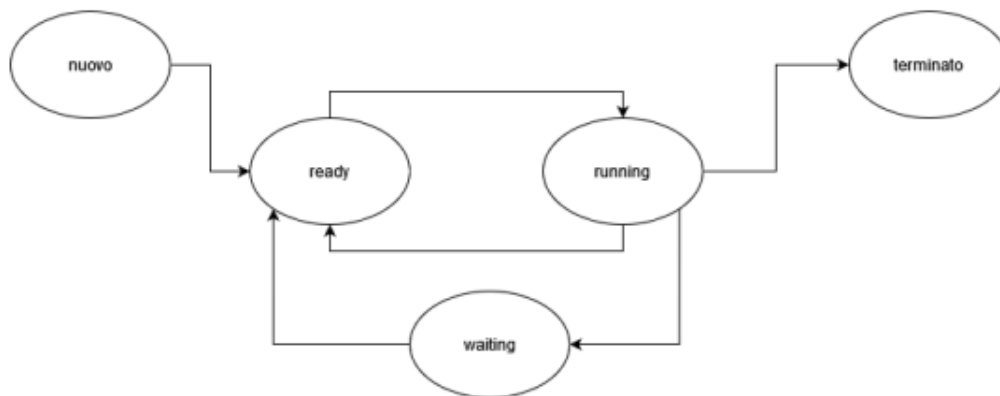


Figure 1.5: Diagramma di transizione degli stati di un processo.

Definizione 1.3.2: Process Control Block (PCB)

In un SO un processo è rappresentato dal PCB, Process Control Block, in cui sono rappresentate le sue informazioni:

- Stato del processo.
- PC (Program Counter).
- Copia dei registri di CPU (copiati nel PCB quando il processo passa da running a ready).
- Info sullo scheduling CPU (priorità e parametri di scheduling).
- Info sulla gestione della memoria (tabella delle pagine).
- Contabilizzazione risorse.
- Stato di I/O

Vengono Mantenuti in strutture dati dette *code*:

- Coda ready: contiene tutti i PCB dei processi caricati in memoria in stato ready.
- Coda di dispositivo: contiene tutti i PCB dei processi in stato waiting.

Definizione 1.3.3: Swapping

Quando un sistema ha in esecuzione più processi di quanto la RAM può sostenere, viene eseguito lo swapping: una parte dei processi, quindi i loro PCB, vengono trasferiti in memoria secondaria (swap out) fino a quando non sarà possibile eseguirli (swap in).

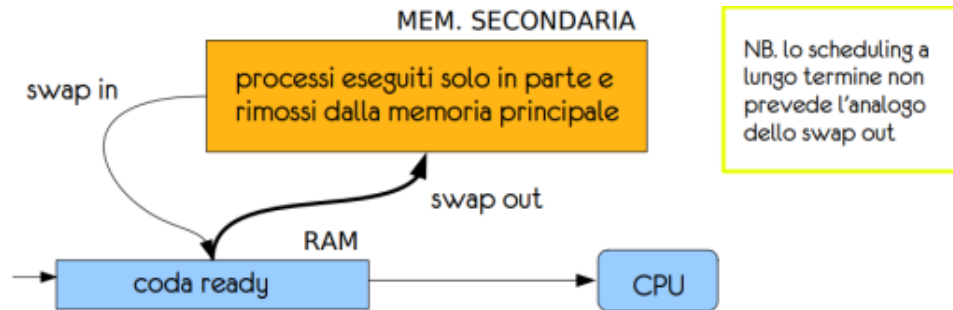


Figure 1.6: Swapping.

1.3.2 Commutazione e Scheduling

Definizione 1.3.4: Context Switch

Il passaggio di un processo da running a waiting, o anche da ready a running richiede quello che si chiama context switch.

Note:-

Questa operazione avviene esclusivamente in modalità kernel. Il tempo per un context switch dipende dall'architettura, in quanto si possono avere diversi set di registri.

Ci sono tre tipi di scheduling:

- *Scheduling a breve termine*: politica di avvicinamento alla CPU dei processi.
- *Scheduling a medio termine*: quando il grado di multiprogrammazione è troppo alto, e quindi non tutti i processi possono essere contenuti in RAM, a turno vengono spostati in memoria secondaria.
- *Scheduling a lungo termine*: presente in sistemi batch, in cui la coda dei processi era conservata in memoria secondaria. Si attiva quando un processo in esecuzione termina, scegliendone uno in memoria secondaria da caricare.

Definizione 1.3.5: Scheduling a Breve Termine

Seleziona dalla coda ready il processo a cui assegnare la CPU. Casi:

1. Running → Waiting.
2. Running → Ready tramite interrupt.
3. Waiting → Ready.
4. Running → Terminated.

Osservazioni 1.3.2

- Lo scheduling è *preemptive* se interviene in almeno uno dei casi 2 e 3, può occorrere anche in 1 o 4.
- Lo scheduling è non-preemptive se interviene solo nei casi 1 o 4.
- Il *dispatcher* effettua il cambio di contesto, effettua il posizionamento alla giusta istruzione, passa nella modalità di esecuzione giusta.

Criteri di scheduling:

- Mantenere la CPU il più attiva possibile.
- *Throughput*: numero di processi completati nell'unità di tempo.
- *Tempo di attesa*: tempo trascorso in coda ready.
- *Turnaround time*: tempo di completamento di un processo; somma del tempo di esecuzione e dei tempi di attesa.
- *Tempo di risposta*.

Definizione 1.3.6: Diagramma di Gantt

In un diagramma di Gantt si rappresenta l'occupazione della CPU con un rettangolo.

1.3.3 Algoritmi di Scheduling**Definizione 1.3.7: First Come First Served (FCFS)**

PCB organizzati in una coda FIFO, non è preemptive, il primo processo arrivato è il primo processo ad avere la CPU, e la mantiene per un intero CPU burst. Il tempo medio di attesa è abbastanza lungo.

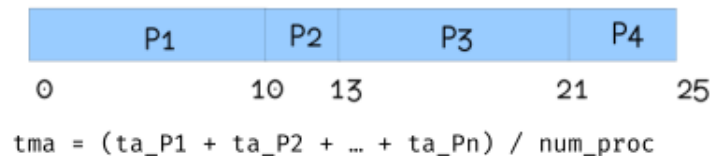


Figure 1.7: Diagramma di Gantt di FCFS.

Corollario 1.3.1 Prelazione

Meccanismo generale per il quale il SO può togliere una risorsa riservata per un processo (la CPU nello scheduling a breve termine) anche se il processo la sta utilizzando o la utilizzerà in futuro.

Osservazioni 1.3.3

- La prelazione richiede l'introduzione di meccanismi di sincronizzazione per evitare inconsistenze: se due processi condividono dati e uno dei due li sta aggiornando quando gli viene tolta la CPU, i dati possono risultare mancanti o corrotti.
- La prelazione non è sempre possibile: occorre avere un'architettura che supporti i timer (in caso contrario lo scheduling è solo preemptive).

Definizione 1.3.8: Shortest Job First

SJF è ottimale nel minimizzare il tempo medio di attesa, seleziona sempre il processo avente CPU burst successivo di durata minima. Se la durata dei CPU burst non è nota, si prevede la durata sulla base dei CPU burst precedenti, combinati con una media esponenziale.

Questo scheduler può essere:

- Preemptive (o shortest remaining time left): quando un nuovo processo diventa ready, lo scheduler controlla se il suo burst è inferiore a quanto rimane del burst del processo running:
 - Sì: il nuovo processo diventa running, context switching.
 - No: il nuovo processo va in coda ready.
- Non-preemptive: il processo viene messo in coda ready, eventualmente in prima posizione.

Definizione 1.3.9: Scheduling a Priorità

Ogni processo ha associata una priorità che decide l'ordine di assegnazione della CPU. Può essere con o senza prelazione.

La priorità può essere definita su due criteri:

- Internamente: sulla base di caratteristiche del processo. Per esempio, potremmo usare SJF e definire la priorità come l'inverso della durata stimata del CPU burst.
- Esternamente: non calcolabile, impostata in base a criteri esterni, dall'utente per esempio.

Corollario 1.3.2 Starvation

Tutti gli algoritmi a priorità sono soggetti a starvation: un processo potrebbe non ottenere mai la CPU perché continuano a passargli davanti nuovi processi a priorità maggiore. Per risolvere, la priorità dei processi viene alzata con il trascorrere del tempo (aging).

Definizione 1.3.10: Round Robin (RR)

Tipo di scheduling preemptive ideato per i sistemi time-sharing. La coda ready è FIFO circolare. A ogni processo viene assegnata la CPU per un quanto di tempo, se non è sufficiente a concludere, il processo viene inserito nuovamente in coda ready e la CPU riassegnata al successivo.

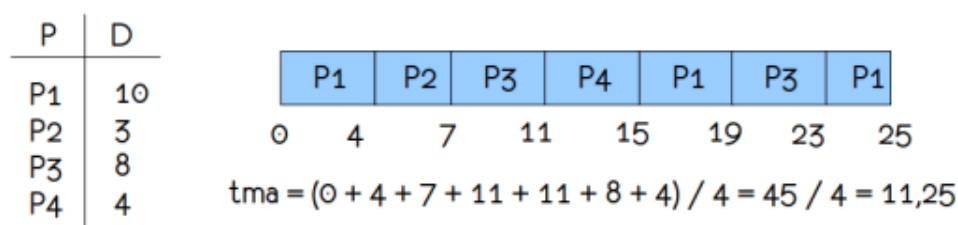


Figure 1.8: Diagramma di Gantt di RR.

Note:-

Il tma è abbastanza alto ma non si ha starvation: ogni processo viene servito ogni $(n_{proc} - 1) * \text{quanto msec.}$ La velocità dell'algoritmo dipende anche dal quanto di tempo scelto: se è troppo lungo RR tenta di diventare un FCFS, mentre se è troppo corto i tempi del context switch possono rallentare e appesantire il processo.

Definizione 1.3.11: Code a Multilivello con Feedback

Algoritmi utili quando si possono dividere i processi in categorie legate alla loro natura. La ready queue è divisa in tante code quante sono le categorie, che possono avere algoritmi di scheduling. Tra le code, poi, esiste una priorità.

Note:-

Se i processi possono cambiare code, si parla di multilivello con feedback.

1.3.4 Creazione e Terminazione**Definizione 1.3.12: Creazione di un Processo**

Un processo viene generato dall'unica entità attiva gestita dal SO: un altro processo. Con l'avvio del SO si genera un albero di processi che cresce e decresce dinamicamente, a seconda dell'evoluzione dell'elaborazione.

Osservazioni 1.3.4

- Ogni processo ha un PID.
- I processi figli, generati dal padre, in genere condividono delle risorse. In unix, per esempio, il figlio riceve una copia di tutte le variabili del padre.

Definizione 1.3.13: Terminazione

Un processo giunto alla sua ultima istruzione termina tramite la system call exit. Il SO libera le risorse associate al processo. A terminazione completata, il padre del processo viene notificato della terminazione. In certi SO, i dati sulla terminazione del figlio possono essere mantenuti finché non vengono ispezionati dal padre, che a sua volta si occupa di rimuoverli.

Un processo può esplicitamente terminare un altro processo (conoscendo il suo PID) tramite system call se:

- Il processo sta usando troppe risorse.
- Finisce di elaborare.
- Il padre è terminato (SO forza la terminazione dei figli).

2

Sincronizzazione

2.1 Processi Cooperanti

2.1.1 Introduzione

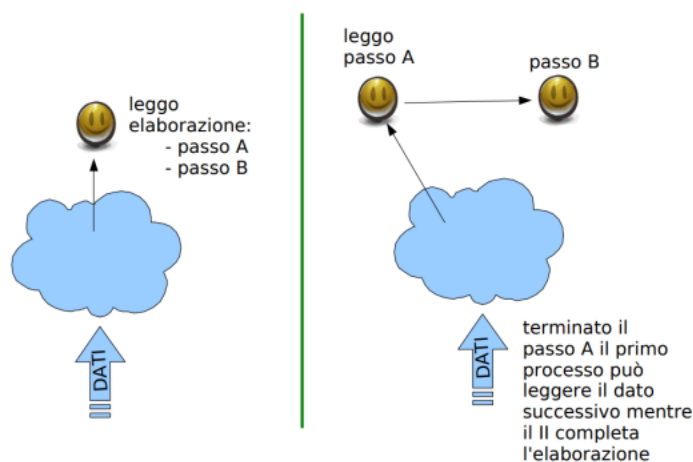


Figure 2.1: Processi che cooperano.

Classico esempio (produttore/consumatore):

- Il processo *produttore* produce dei dati.
- Il processo *Consumatore* consuma dei dati.

Definizione 2.1.1: Memoria Condivisa

Area di memoria fruibile da più processi. Richiesta tramite system call, i processi che la vogliono utilizzare devono agganciarla al proprio spazio di indirizzi. Può essere modificata a piacimento secondo qualsiasi tipo di dati utile ai programmi che la utilizzano.

```

PRODUTTORE
...
alloca b di tipo buffer_cond come memoria condivisa

while (1) {
    if (! pieno(b) ) {
        nuovo = ... produci ...;
        b.dato[b.inserisci] = nuovo;
        b.inserisci = (b.inserisci+1) % D;
    }
}

CONSUMATORE
...
aggancia b al proprio spazio indirizzi

while (1) {
    if (! vuoto(b) ) {
        nuovo = b.dato[b.preleva];
        b.preleva = (b.preleva+1) % D;
    }
}

```

Figure 2.2: Produttore e consumatore.

Definizione 2.1.2: Inconsistenza dei Dati

Per colpa dello scheduler della CPU su sistemi a singolo core, o su sistemi con 2 o più core, potrebbero verificarsi problemi di inconsistenza dei dati.



Figure 2.3: Caso di inconsistenza.

2.1.2 Scambi di Messaggi

Definizione 2.1.3: Scambio di Messaggi

Consente a due processi di comunicare senza condividere una stessa area di memoria.

Lo scambio di messaggi può essere:

- *Diretto* o *indiretto*: se diretto, ricevente deve comunicare il PID al mittente.
- *Sincrono* (bloccante) o *asincrono* (non bloccante):
 - Invio sincrono: mittente aspetta che ricevente riceve il messaggio.
 - Ricezione sincrona: ricevente aspetta il messaggio.
 - Rendez-vous: entrambi insieme.
- A gestione automatica.

Note:-

Questo è visto in dettaglio nel corso di "Modelli Concorrenti e Algoritmi Distribuiti".

Per comunicare, due processi utilizzano un canale, e le funzioni send e receive per scambiarsi messaggi:

- *Comunicazione diretta*:

- `send(P, msg)`: invia msg al processo P.
- `receive(P, msg)` / `receive(id, msg)`: attendi un messaggio da P, memorizza il messaggio in msg / ricevi un messaggio, memorizza in id il PID del mittente e in msg il messaggio.
- la reciproca conoscenza del PID finisce in un canale logico.

- *Comunicazione indiretta:*

- Si utilizza una mailbox come canale intermedio di comunicazione, distinte dall'identità del ricevente.
- `send(M, msg)`: invia msg alla mailbox M.
- `receive(M, msg)`: attendi un messaggio alla mailbox M.
- Più mittenti e/o riceventi possono usare la stessa mailbox.
- Una stessa coppia di processi può utilizzare diverse mailbox per comunicare.
- La mailbox ha una capacità N definita da un buffer:
 - * 0: il canale non ha memoria (meccanismo no buffering); il mittente rimane sospeso se il ricevente non ha ancora consumato il messaggio ricevente (gestione esplicita del buffer).
 - * $N > 0$: il mittente rimane in attesa solo se il buffer è pieno (automatic buffering).
 - * Illimitata: il mittente non attende mai (automatic buffering).

Definizione 2.1.4: Socket

Usato in sistemi client-server, socket è il nome dato ad un endpoint di un canale di comunicazione tra due processi. Due processi che vogliono comunicare usano una coppia di socket, formati da l'IP della macchina concatenato ad una porta

Note:-

Questo è visto in dettaglio nel corso di "Programmazione III".

Osservazioni 2.1.1

- Le porte < 1024 sono riservate. Se un processo utente richiede una porta, riceverà un numero maggiore di 1024.
- L'indirizzo di loopback (127.0.0.1) identifica la macchina stessa.

Definizione 2.1.5: Remote Procedure Call (RPC)

Meccanismo di scambio di messaggi utilizzato in sistemi client server per invocare l'esecuzione di una procedura che risiede su un'altra macchina connessa in rete. I messaggi sono ben strutturati: hanno identificatore della procedura e parametri.

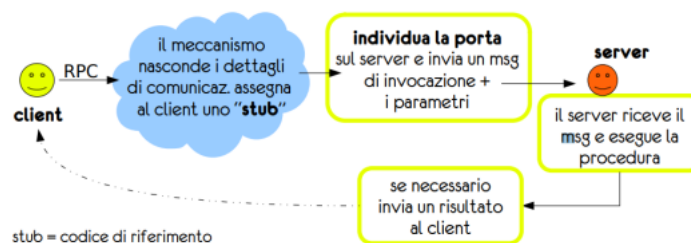


Figure 2.4: Remote Procedure Call.

Definizione 2.1.6: Pipe

Canale di comunicazione tra processi.

Esistono due tipi di pipe:

- *Pipe anonima:*
 - Canale simplex (unidirezionale), FIFO.
 - Uno produce, l'altro consuma.
- *Named pipe:*
 - FIFO.
 - In UNIX è simplex, in windows è full-duplex.
 - Può far comunicare più di due processi.
 - Spesso realizzata come file.

2.2 Thread

2.2.1 Heavyweight Process

Definizione 2.2.1: Thread

Sono le diverse parti di un processo che cooperano per eseguire tutte le sue funzioni. Molti SO li considerano come l'unità di base d'uso della CPU al posto del processo.

Un processo costituito solo da thread è detto heavyweight process.

Un thread è costituito da:

- Un identificatore.
- Un PC (Program Counter).
- Un insieme di valori di registri.

Note:-

Condivide con gli altri thread dello stesso processo il codice, la sezione dati, i file aperti, i segnali.

Vantaggi dell'utilizzo di thread:

- Usare thread al posto di processi cooperanti è più efficiente.
- Molti programmi contengono sezioni di codice che è possibile eseguire indipendentemente.
- I context switch avvengono più rapidamente.
- Si allocano meno risorse, in quanto molte vengono condivise.
- La comunicazione è più veloce, perché avviene tramite variabili condivise.
- È possibile assegnare thread diversi a core diversi in architetture multicore.

Osservazioni 2.2.1

- Mentre i thread sono definiti dal programmatore e gestiti dal programma, i processi sono generati in modo invisibile all'utente, non permettendo ai programmi di gestirli direttamente ma solo tramite system call.

- Alcuni linguaggi di programmazione includono specifiche istruzioni per la creazione e il controllo dei thread (es: Java, C#, Python), mentre per altri c'è bisogno di includere apposite librerie (es: C, C++; sono detti "a singolo flusso di controllo").

Operazioni sui thread:

- *Creazione*.
- *Terminazione*: più rapida di quella dei processi perché non richiede la gestione delle risorse.
- *Sospensione* (blocco).
- *Recupero* (risveglio).
- *Join*: comporta l'attesa da parte di un thread della terminazione di un altro.

2.2.2 Tipi di Thread

I thread si possono dividere in:

- *Thread Utente*.
- *Thread Kernel*.

Definizione 2.2.2: Thread Utente

Funzionamento a singolo contesto: ogni processo multithread deve gestire le info dei suoi thread, il loro scheduling, la comunicazione tra i thread. Sono creati da funzioni di libreria che non possono eseguire istruzioni privilegiate. Sono trasparenti al SO, che vede il processo come una sola entità.

Vantaggi:

- Possibile utilizzarli su SO senza multi-threading in quanto sono gestiti internamente al processo.
- Criteri di scheduling adattabili alle esigenze del programma.
- Esecuzione più rapida perché non usa nè interrupt nè system call.

Svantaggi:

- Non adattabili a sistemi multiprocessore.
- Operazioni I/O bloccano l'intero processo fino al termine dell'operazione.

Definizione 2.2.3: Thread Kernel

Funzionamento a molteplici contesti, uno per ogni thread creato. L'utente genera thread utente, e il SO li implementa a livello kernel associando a essi delle strutture proprie (thread kernel sono strutture separate).

Note:-

Sono gestiti dal SO, quindi esso gestisce anche il loro scheduling e sono più costosi perché per ogni thread deve mantenere gli appositi descrittori e le loro associazioni ad ogni processo.

Vantaggi:

- Possibile distribuire i thread su più processori.
- Le operazioni di I/O non bloccano l'intero processo (eseguite solo sul thread).
- Maggiore interattività con l'utente.
- I singoli processi sono più veloci.

Svantaggi:

- Non portabili su SO non-multithreaded.
- In caso di troppi thread, è possibile sovraccaricare il processore.

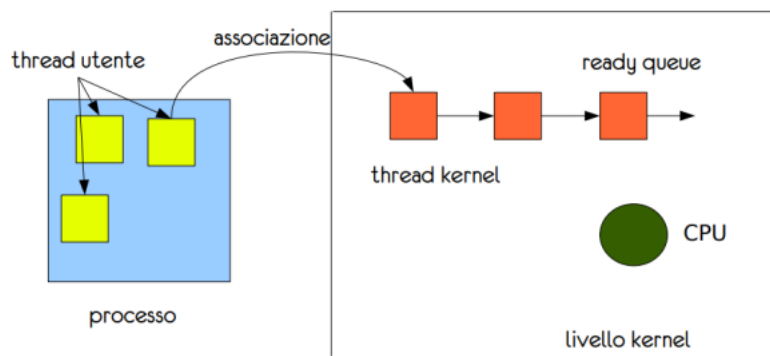


Figure 2.5: Tipi di Thread.

Associazione dei thread kernel:

- *Uno a uno*: un thread utente assegnato a un solo thread kernel.
- *Molti a uno*: un pool di thread utente assegnata a un solo thread kernel.
- *Molti a molti*: un pool di thread utente assegnata a molteplici thread kernel (di solito di numero inferiore).

2.2.3 Lightweight Process**Definizione 2.2.4: Lightweight Process**

Nei SO che usano il modello molti a molti o a due livelli, si ha la necessità di eseguire uno scheduling dei thread utente per l'accesso ai thread kernel. Questo viene fatto introducendo un layer tra thread utente e kernel, chiamato Lightweight Process (LWP). Questi LWP sono visti come processori virtuali, e corrispondono ad un thread kernel.

Note:-

Questa associazione viene fatta in modo esplicito dall'applicativo che deve usare delle "upcall" speciali.

Osservazioni 2.2.2

1. L'utente quindi effettua il proprio scheduling dei thread su un insieme di LWP messi a disposizione dal kernel. Un thread utente è in esecuzione se ha un LWP assegnato.
2. Se un thread esegue una syscall bloccante, il SO informa l'applicativo con una upcall.
3. L'applicativo quindi esegue un gestore della upcall, che salva lo stato del thread bloccante, e riassegna l'LWP ad un altro thread pronto per l'esecuzione.
4. Quando poi si verifica l'evento che sveglia il thread sospeso, il SO con un'altra upcall farà segnare all'applicativo il thread come pronto, a cui verrà assegnato un LWP.

2.3 Implementare la Sincronizzazione

2.3.1 Scheduling della CPU

In presenza di thread, lo scheduling della CPU è quindi a due livelli:

- *Process-Contention Scope (PCS)*: è lo scheduling effettuato all'interno di un processo per decidere a quali thread utente assegnare LWP. I thread kernel vengono gestiti come PCB.
- *System-Contention Scope (SCS)*: lo scheduling della CPU viene fatto tra tutti i thread kernel in queue, con una granularità più fine rispetto ai PCB, tutto in maniera indipendente dal processo di appartenenza.

Problema: interleaving delle istruzioni di diversi processi produttori eseguite in un'area di memoria condivisa. In assenza di controlli, diverse istruzioni sugli stessi dati possono non essere eseguite in modo atomico, producendo dati inconsistenti.

Note:-

Il problema è spesso presente in un SO con multitasking, in cui possono essere presenti allo stesso tempo diversi processi in modalità kernel.

Soluzione: per risolvere il problema, viene implementata nel programma la *sezione critica*.

Definizione 2.3.1: Sezione Critica

Una porzione di codice in cui un processo modifica variabili condivise secondo un certo criterio sicuro.

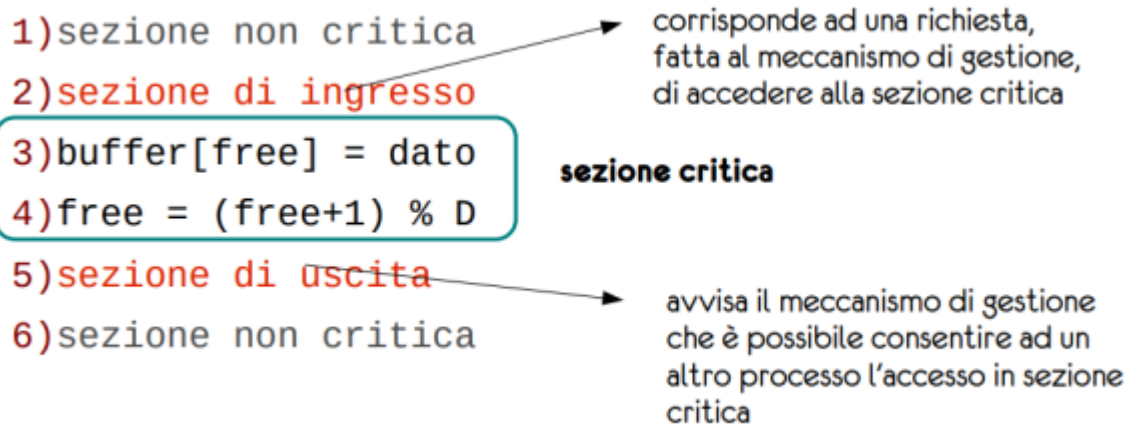


Figure 2.6: Sezione critica.

Note:-

Una sezione critica è determinata dalle variabili condivise utilizzate, e non deve essere eseguita con interleaving di istruzioni di altre sezioni critiche che usano le stesse variabili.

Criteri soddisfatti da una soluzione al problema della sezione critica:

- *Mutua Esclusione*: un solo processo per volta può eseguire la sua sezione critica.
- *Progresso*: nessun processo che non desideri utilizzare una variabile condivisa può impedirne l'accesso a processi che desiderano utilizzarla. Solo i processi che intendono entrare in sezione critica concorrono a determinare chi entrerà.
- *Attesa limitata*: esiste un limite superiore all'attesa di ingresso in sezione critica.

```

<sezione non critica>
while (turno != i) do no_op;
<sezione critica>
turno = j;
<sezione non critica>

```

Figure 2.7: Primo tentativo di soluzione.

Problema: se un programma non ha più bisogno di usare la sezione critica, non ci entrerà; se è il turno per un programma di utilizzare la sezione critica ma non ci entra, l'altro processo rimane in waiting (violazione progresso).

```

flag[i] = true;
while (flag[j]) do no_op;
<sezione critica>
flag[i] = false;
<sezione non critica>

```

Figure 2.8: Seconda soluzione.

Problema: quando entrambi i processi hanno il flag attivo, entrambi vanno in loop infinito (violazione attesa limitata).

```

flag[i] = true; turno = j;
while (flag[j] && turno == j) do no_op;
<sezione critica>
flag[i] = false;
<sezione non critica>

```

Figure 2.9: Terza soluzione: algoritmo di Peterson.

Note:-

L'algoritmo di Peterson garantisce tutti i criteri ma funziona solo per 2 processi (non è generalizzabile).

Note:-

L'algoritmo del fornaio garantisce tutti i criteri, funziona per N processi, ma il codice è complesso e i processi fanno busy waiting (anziché sospendersi, attengono il turno tenendo occupata la CPU).

Codice di P_i :

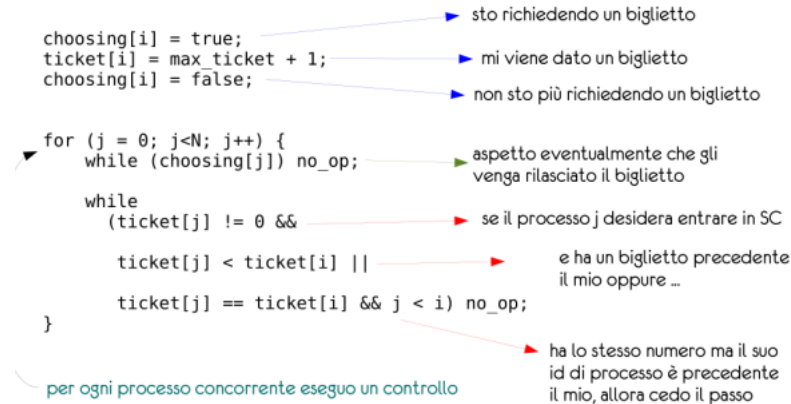


Figure 2.10: Quarta soluzione: algoritmo del fornaio.

2.3.2 Altri tipi di sincronizzazione:

Sincronizzazione hardware:

- Potremmo disabilitare gli interrupt quando entriamo in una sezione critica, questo eviterebbe la prelazione. Ma causa interferenze pesanti con lo scheduling della CPU, e gli interrupt non possono essere mantenuti disabilitati a lungo.
- Introdurre l'uso di lock: per entrare in una sezione critica, il processo deve avere ottenuto il giusto lock, che rilascerà al termine. Per questo motivo molte architetture forniscono operazioni di controllo e modifica del valore di una cella e istruzioni per lo scambio in modo atomico.

Definizione 2.3.2: TestAndSet

TestAndSet restituisce il valore precedente di lock che sarà falso se nessun altro processo è in una sezione critica controllata tramite lock. Solo in questo caso si esce dal while.

Note:-

L'esecuzione dell'intera routine è atomica.

```
bool TestAndSet (bool *variabile) {
    bool valore = *variabile;
    *variabile = true;
    return valore;
}
```

Figure 2.11: Implementazione di TestAndSet.

```
while (TestAndSet(&lock));

<sezione critica>

lock = false;
```

Figure 2.12: Esecuzione di TestAndSet.

Definizione 2.3.3: Swap

Scambia in modo atomico i valori dei suoi due parametri.

```
Swap(bool *a, bool *b) {
    bool temp = *a;
    *a = *b;
    *b = temp;
}
```

Figure 2.13: Implementazione di Swap.

Osservazioni 2.3.1

- La particolarità è nell'esecuzione della routine, che è atomica.
- Usiamo lo swap per realizzare l'accesso in mutua esclusione ad una sezione. Oltre al lock, variabile condivisa, abbiamo anche una variabile booleana locale "chiuso"
- Si esce dal ciclo while solo quando lock è false.
- All'uscita da Swap lock risulta automatica impostato a true.

```
chiuso = true;

while(chiuso) Swap(&lock, &chiuso);

<sezione critica>

lock = false;
```

Figure 2.14: Esecuzione di Swap.

Note:-

Entrambi i metodi visti garantiscono la mutua esclusione, ma non l'attesa limitata. Non c'è garanzia che un processo che vuole eseguire una delle due non venga sempre prevaricato da altri.

| <u>Sezione d'ingresso:</u> | <u>Sezione d'uscita:</u> |
|--|---|
| <pre>attesa[i] = true; chiave = true; while (attesa[i] && chiave) chiave = TestAndSet(&lock); attesa[i] = false;</pre> | <pre>j = (i+1) % n; while ((j != i) && !attesa[i]) j = (j+1) % n; if (j == i) lock = false; else attesa[j] = false;</pre> |

Figure 2.15: Algoritmo di attesa limitata.

3

Gestione della Sincronizzazione

3.1 Semafori

Strumento di sincronizzazione introdotti da Dijkstra per minimizzare il busy-waiting.

3.1.1 Introduzione

Definizione 3.1.1: Semafori

Il semaforo non è altro che una variabile a cui si può accedere solo con due operazioni atomiche:

- P (Proberen, verificare in olandese).
- V (Verhogen, incrementare, in olandese).

| Definizione di P | Definizione di V |
|--|--------------------------------------|
| <pre>P (S) { while (S <= 0) no_op; S--; }</pre> | <pre>V(S) { S++; }</pre> |

Figure 3.1: Definizioni.

Definizione 3.1.2: Spinlock

L'implementazione originale dei semafori, detta spinlock, aveva attesa attiva.

Note:-

Una possibile soluzione per evitare attesa attiva: ogni semaforo mantiene una lista di PCB dei processi sospesi su quel semaforo; quando un processo si sospende su quel semaforo, lo scheduler assegna la CPU a un altro processo, e quando il semaforo viene alzato, uno dei processi in attesa viene riattivato.

```

typedef struct {
    int valore; /* Valore del semaforo */
    processo *lista; /* Lista di attesa relativa al semaforo */
} semaforo;

P (semaforo *s) {
    S->valore - -; /* Decremento il valore del semaforo se il valore diventa negativo
                   occorre sospendere il processo */
    if (S->valore < 0) {
        <aggiungi il PCB di questo processo a S->lista>
        /* block è una system call che richiama lo scheduler della CPU, che deve riassegnare
           la medesima a un altro processo, e il dispatcher, che deve effettuare il context switch */
        block();
    }
}

V (semaforo *S) {
    S->valore++; /* Incremento il valore del semaforo */
                /* Se il valore è negativo, allora ci sono processi da
                   risvegliare */
    if (S->valore < 0) {
        <scegli un PCB P da S->lista>
        wakeup(P);
    }
}

```

Figure 3.2: Implementazione di un semaforo.

Note:-

Il valore del semaforo indica il numero di processi in attesa.

I valori dei semafori possono essere:

- 1: valore di inizializzazione, risorsa disponibile.
- 0, -1, -2, ..., -n: risorsa occupata.

Note:-

I semafori che possono assumere valori > 1 sono detti *semafori contatori*, il cui numero rappresenta una quantità di risorse disponibili.

I semafori possono realizzare molti tipi di sincronizzazione:

- *Mutua esclusione*: tutti i processi coinvolti separano le loro sezioni critiche con $P(\text{mutex})$ e $V(\text{mutex})$, dove mutex è un semaforo di mutua esclusione.
- *Accesso limitato*: tutti i processi coinvolti separano le loro sezioni critiche con $P(\text{nr})$ e $V(\text{nr})$, dove nr è un semaforo contatore che permette a molteplici processi di eseguire in parallelo una certa sezione critica.
- *Ordinamento*: l'ordine di esecuzione dei processi viene esplicitamente controllato.

Note:-

Le operazioni sui semafori devono essere eseguite in modo atomico, in quanto i semafori sono variabili condivise (le operazioni sono quindi dentro sezioni critiche).

Si può ottenere l'atomicità:

- Sui sistemi monoprocesso disabilitando gli interrupt (P e V sono brevi quindi non rallentano il sistema).
- Sui sistemi multiprocesso utilizzando gli spinlock (disabilitare gli interrupt fa calare le prestazioni).

3.1.2 Problemi Classici

Produttori e Consumatori

Usiamo un buffer circolare di SIZE posizioni in cui i produttori inseriscono i dati e i consumatori li prelevano.

```
typedef struct {...} item;
item buffer [SIZE];
semaphore full, empty, mutex;
item nextp, nextc;
int in = 0, out = 0;
full = 0 ; empty = SIZE; mutex = 1;
```

Figure 3.3: Dati condivisi e inizializzazione.

- full: conta il numero di posizioni piene del buffer.
- empty: conta il numero di posizioni vuote del buffer.
- mutex: semaforo binario, per l'accesso in mutua esclusione del buffer circolare e delle variabili in e out.
- in e out: servono per gestire il buffer circolare.

```
while (true) {
    ...
    produci un item in nextp
    ...
    wait(empty);
    wait(mutex);
    buffer[in] = nextp;    // inserisce nextp
    in = (in + 1) mod SIZE; // nel buffer
    signal(mutex);
    signal(full);
}
```

Figure 3.4: Codice di un PRODUTTORE.

```
while (true) {
    wait(full)
    wait(mutex);
    nextc = buffer[out]    // rimuove un
    out = out + 1 mod SIZE // elemento
    signal(mutex);
    signal(empty);
    ...
    consuma item in nextc
    ...
}
```

Figure 3.5: Codice di un CONSUMATORE.

Osservazioni 3.1.1

- Se ci sono più produttori, e buffer ha almeno due posizioni libere ($\text{empty} \geq 2$), ci possono essere

due (o più) processi che hanno superato la `wait(empty)` e cercano di inserire in buffer un item e di aggiornare la variabile `in`.

- Se ci sono più consumatori, e buffer ha almeno due posizioni piene ($\text{full} \geq 2$), ci possono essere due (o più) processi che hanno superato la `wait(full)` e cercano di prelevare un item dal buffer e di aggiornare la variabile `out`.

Lettori e Scrittori

Problema: condividere un file tra molti processi.

- Alcuni processi richiedono solo la lettura (*lettori*), altri possono voler modificare il file (*scrittori*).
- Due o più lettori possono accedere al file contemporaneamente.
- Un processo scrittore deve poter accedere al file in mutua esclusione con tutti gli altri processi.

```
wait(scrivi);
...
esegui la scrittura del file
...
signal(scrivi);
```

Figure 3.6: Codice di uno SCRITTORE.

```
wait(mutex); // mutua escl. per aggiornare numlettori
numlettori++;
if (numlettori == 1) wait(scrivi); // il primo lettore ferma
                                   // eventuali scrittori

signal(mutex);
... leggi il file ...
wait(mutex);
numlettori--;
if (numlettori == 0) signal(scrivi);
signal(mutex);
```

Figure 3.7: Codice di un LETTORE.

Cinque Filosofi

Problema: 5 filosofi passano il tempo seduti intorno a un tavolo pensando e mangiando a fasi alterne. Per mangiare un filosofo ha bisogno di due posate ma ci sono solo cinque posate in tutto.

Definizione 3.1.3: Deadlock

Lo stallo o deadlock indica una situazione in cui due o più processi o azioni si bloccano a vicenda, aspettando che uno esegua una certa azione (per esempio rilasciare il controllo su una risorsa come un file, una porta input/output ecc.) che serve all'altro e viceversa.

Definizione 3.1.4: Starvation

La starvation è l'impossibilità perpetua, da parte di un processo pronto all'esecuzione, di ottenere le risorse sia hardware sia software di cui necessita per essere eseguito.

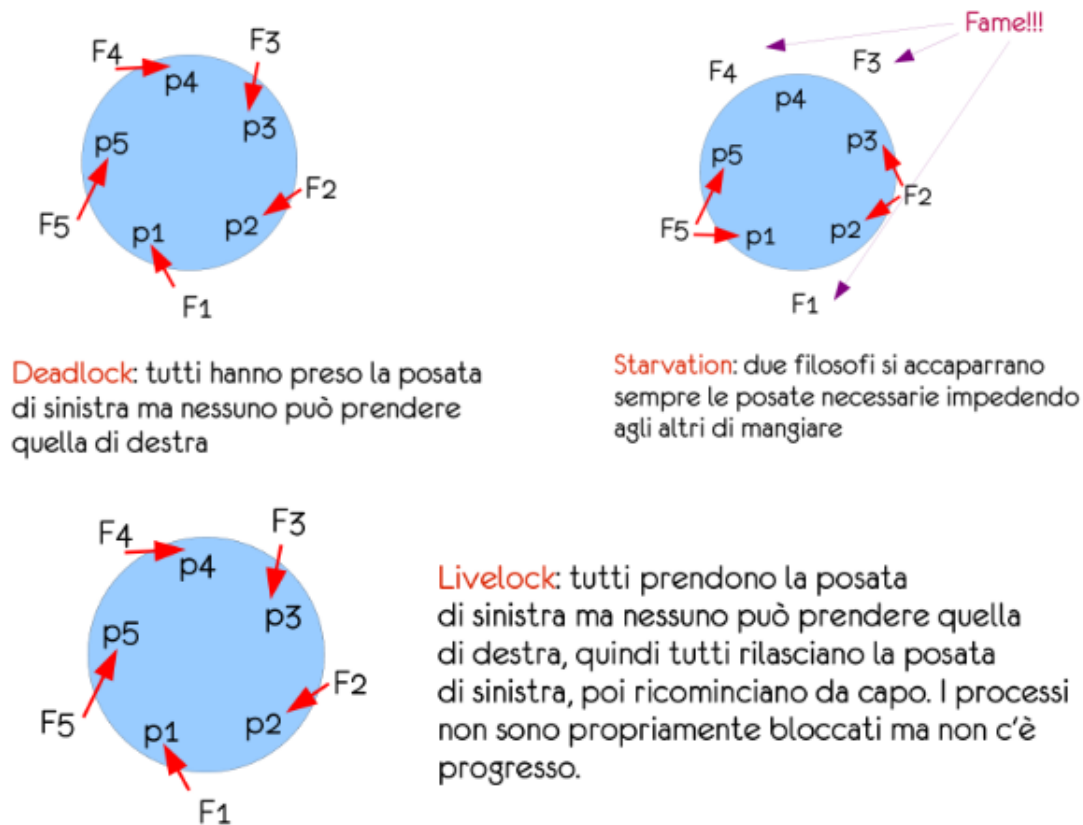


Figure 3.8: Problema dei 5 filosofi.

Note:-

Dijkstra propone una soluzione che consente di evitare il deadlock. Bisogna rompere la simmetria nell'accesso alle posate, introducendo la seguente regola: ogni filosofo deve prendere per prima la posata con indice minore.

Soluzione di Chandy-Misra: quando due processi sono in competizione l'algoritmo non favorisce sempre lo stesso (causando starvation). Viene introdotto un concetto di precedenza dinamica tra processi (cambia nel tempo), che funziona associando un concetto di stato alla risorsa.

1. Ogni forchetta (risorsa) può essere "sporca" o "pulita"; inizialmente sono tutte "sporche" e ogni filosofo ne ha una.
2. Quando un filosofo vuole mangiare, invia ai suoi vicini dei messaggi per ottenere le forchette che gli mancano.
3. Quando un filosofo, che ha una forchetta, riceve una richiesta: se sta pensando, la cede altrimenti se la forchetta è pulita ne mantiene il possesso, se è sporca la cede. Quando passa una forchetta ne pone lo stato a "pulita".
4. Dopo aver mangiato, tutte le forchette di un filosofo diventano "sporche". Se risultano richieste pendenti per qualche forchetta, il filosofo la pulisce e la passa.

3.2 Monitor

3.2.1 Introduzione

Definizione 3.2.1: Monitor

Sono dei costrutti di sincronizzazione contenenti i dati e le operazioni necessarie per allocare una risorsa condivisa usabile in modo seriale. Le variabili di un monitor sono condivise dai processi che usano quel monitor. Per accedere alle variabili condivise un processo deve eseguire una routine di accesso al monitor (possono essercene diverse). Un solo processo per volta può essere attivo all'interno di un monitor.

Note:-

È un abstract data type (ADT): in sé non è un tipo di dato, ma contiene dati e le istruzioni per utilizzarli.

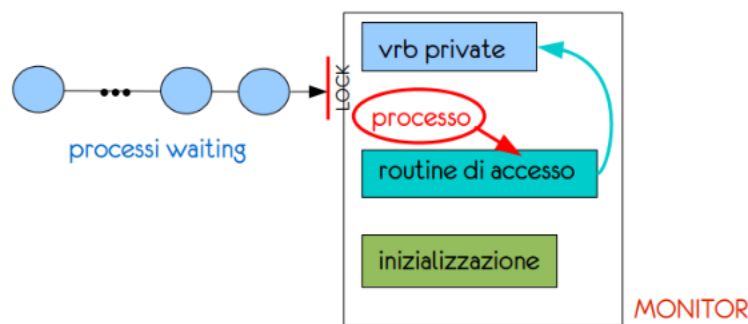


Figure 3.9: Monitor.

I monitor consentono anche di effettuare sincronizzazione tra processi attraverso variabili di tipo condition (boolean), su cui si possono eseguire solo le operazioni:

- wait(x): esecutore sospeso se $x == \text{false}$.
- signal(x): se un processo è sospeso sulla condition x, uno viene scelto e risvegliato.

monitor per un allocatore di risorse

```
boolean inUso = false;
```

```
condition disponibile;
```

```
monitorEntry void prendiRisorsa() {
    if (inUso) wait(disponibile);
    inUso = true;
}
```

```
monitorEntry void rilasciaRisorsa() {
    inUso = false;
    signal(disponibile);
}
```

se una qualche risorsa di interesse risulta usata da qualcun altro, allora mi sospendo sulla condition "disponibile". Quando mi risveglio continuo dalla linea di codice successiva: setto a true inUso.

Memento: il monitor è eseguito in ME quindi tutte le monitorEntry sono atomiche, a meno di sospensioni volontarie

quando rilascio una risorsa, faccio una notifica sulla condition "disponibile". Se qualcuno era in attesa si risveglia, altrimenti la notifica viene dimenticata

Figure 3.10: Esempio di monitor.

Quando un processo (thread) esegue signal rischiamo di avere due processi attivi nel monitor. Ci sono 2 soluzioni:

- Segnalare e attendere:
 - P attende.
 - Q riprende ed esegue.
- Segnalare e proseguire:
 - P continua.
 - Q aspetta che P finisca.

3.3 Transazioni e Lock

3.3.1 Transazioni

Definizione 3.3.1: Transazione

Un insieme di istruzioni che esegue una singola funzione logica, ovvero una sequenza di read e write che si conclude con un commit (successo) o con un abort (fallimento).

La transazione deve essere atomica, e l'atomicità dipende dai dispositivi di memoria utilizzati per mantenere i dati elaborati dalla macchina:

- *Memorie volatili* (RAM, cache, registri): cancellati a spegnimento.
- *Memorie non volatili* (dischi, EEPROM): persistenti, ma non sempre "eterni".
- *Memorie stabili* (RAID¹): supporti di memorizzazione che aggiungono politiche/strumenti di duplicazione, rendendo i dati "eterni".

Domanda 3.1

Cosa succede se una transazione viene abortita e la memoria volatile viene cancellata?

Note:-

Dobbiamo tenere traccia delle operazioni eseguite, un logfile mantenuto su memoria stabile.

Il log dovrà contenere quindi:

- L'inizio di una transazione T.
 - <T, start>
- Una sequenza di tuple, relativa a un'operazione di write da fare.
 - <ID transazione, ID dato modificato, valore precedente, nuovo valore>
- Il successo della transazione.
 - <T, commit>

A seguito di un crash di sistema, il SO controlla il log e per ogni transazione T registrata:

- se a <T, start> non corrisponde un <T, commit>, il SO esegue l'operazione undo(T) che ripristina tutti i valori modificati dalla transazione eseguita in modo parziale.
- Se a <T, start> corrisponde un <T, commit>, il SO verifica che le modifiche registrate siano state effettivamente eseguite. In caso contrario, il SO esegue un redo(T) attuando le modifiche.

¹Questo è visto in dettaglio nel corso di "Basi di Dati".

Note:-

Questa registrazione dei dati riduce leggermente l'efficienza dell'esecuzione, ma fornisce un'enorme affidabilità e stabilità al sistema.

Per evitare di scorrere tutto il log file ad ogni crash, introduciamo dei checkpoint. Questi checkpoint ci fanno sapere che:

- Tutte le transazioni prima di questo checkpoint sono state riportate in memoria stabile.
- Tutte le operazioni di scrittura registrate nel logfile sono state applicate con successo.

Note:-

In caso di crash, basta andare a trovare il primo checkpoint, cosicché si possa ignorare tutto quello che avviene prima, applicando le operazioni di undo/redo solo successive al checkpoint.

Definizione 3.3.2: Serializzabilità

La serializzabilità è la proprietà per cui la loro esecuzione concorrente è equivalente alla loro esecuzione in una sequenza arbitraria.

Operazioni conflittuali:

- Date due transazione T_1 e T_2 e le due operazioni O_1 e O_2 , se le operazioni appaiono in successione, accedono agli stessi dati e almeno una delle due operazioni è una write, allora O_1 e O_2 sono operazioni conflittuali.

3.3.2 Lock

Per garantire la serializzabilità si può usare un meccanismo di lock in cui:

- A ogni dato soggetto a transazione si associa un lock.
- Il lock di un dato può essere S (read only) o X (read & write).
- Una transazione che intenda usare un certo dato deve richiederne il lock appropriato, ed eventualmente attendere che un'altra transazione lo rilasci.

Definizione 3.3.3: Lock a 2 Fasi

Inizialmente le transazioni sono in fase di "crescita", ovvero possono ottenere nuovi lock per acquisire tutte le risorse necessarie, senza rilasciarne mai nessuno. Successivamente, in fase di "riduzione", le transazioni rilasciano via i lock per rimuovere le risorse inutilizzate, senza richiederne mai di nuovi.

Note:-

In questo caso è possibile però avere deadlock.

Definizione 3.3.4: Timestamp

Il timestamp è una rappresentazione univoca di un istante temporale.

Corollario 3.3.1 Protocollo basato su Timestamp

Il protocollo assegna a ogni transazione T_i un timestamp T prima della loro esecuzione. Se due transazioni sono tali che $T_5(T_1) < T_5(T_2)$, allora il sistema deve garantire l'esecuzione sequenziale con T_1 prima di T_2 .

Osservazioni 3.3.1

- L'algoritmo non impone un ordinamento corretto fra le transazioni, l'unico scopo è far sì che tutte le volte che una transazione legge un dato, abbia il dato originale oppure il dato modificato dalla

transazione stessa.

- Viene garantita l'atomicità funzionale, diversa dall'atomicità di esecuzione per il fatto che viene consentito l'interleaving delle istruzioni.
- Si hanno quindi transazioni che non interferiscono tra di loro perché usano dati diversi oppure non vengono eseguite in modo concorrente.
- L'algoritmo non è preventivo: identifica situazioni problematiche e le aggiusta.

3.4 Deadlock

3.4.1 Introduzione

Definizione 3.4.1: Deadlock

Un deadlock è una situazione per cui un insieme di processi sono fermi in attesa di un evento che solo uno dei processi appartenenti all'insieme stesso potrebbe causare.

Note:-

Se un processo vuole usare N istanze di una certa risorsa, dovrà chiederle al gestore delle risorse (SO).

Si hanno 3 fasi:

- Richiesta.
- Uso.
- Rilascio.

Definizione 3.4.2: Grafo di Assegnazione delle Risorsse

È una rappresentazione delle assegnazioni che permette di rilevare situazioni di deadlock, tramite un grafo $G = \langle V, E \rangle$.

Nel grafo:

- V è l'insieme dei vertici ed è partizionato in due sottoinsiemi P e R ($P \cap R = \emptyset$)
 - P è l'insieme di tutti i processi del sistema.
 - R è l'insieme di tutte le classi di risorse del sistema.
- E è l'insieme degli archi:
 - Un arco direzionato da R_i a P_j indica che una risorsa di classe R_i è stata assegnata al processo P_j .
 - Un arco direzionato da P_j a R_i indica che il processo P_j ha richiesto ed è in attesa di una risorsa di tipo R_i .

Note:-

Se il grafo non contiene ciclo, non c'è deadlock. La presenza di un ciclo è condizione necessaria ma non sufficiente per avere deadlock.

3.4.2 Prevenzione del Deadlock

Domanda 3.2

Cosa fare con il deadlock?

- Rilevarlo.
- Rompere il deadlock richiede la capacità di monitorare richieste e assegnazioni di risorse.
- Per prevenire il deadlock occorre definire opportuni protocolli di assegnazione delle risorse.
- Fare finta che il deadlock sia impossibile è la tecnica più usata e poco costosa perché non richiede politiche né risorse aggiuntive.

Per prevenire il deadlock è necessario rendere impossibile una delle 4 condizioni necessarie al deadlock:

1. **Mutua Esclusione:** la richiesta di usare le risorse in ME può essere rilasciata solo per alcuni tipi di risorse, altre sono intrinsecamente ME.
2. **Strategia di Havender 1** (possessione e attesa): se un processo ha bisogno di più risorse, le ottiene tutte insieme, oppure non ne ottiene nessuna.
3. **Strategia di Havender 2** (prelazione): quando un processo con N risorse ne richiede un'altra, la ottiene subito oppure rilascia tutte le altre.
4. **Strategia di Havender 3** (attesa circolare): imporre un ordinamento delle risorse e dei processi.

Definizione 3.4.3: Prima Strategia di Havender

Tutte le risorse necessarie ad un processo devono essere richieste insieme:

- Se sono tutte disponibili, il sistema le assegna e il processo prosegue.
- Se anche solo una non è disponibile, il processo non ne acquisisce e si mette in attesa.

Vantaggio: previene il deadlock.

Svantaggio: spreco di risorse (non vengono utilizzate in modo ottimale in quanto un processo può tenersi per più tempo risorse che non usa più).

Note:-

Questa strategia non funziona per processi heavyweight però se all'interno del processo riusciamo a distinguere più thread di esecuzione, ciascuno dei quali ha bisogno di un sottoinsieme delle risorse ed è generato solo quando occorre, la strategia può risultare efficace.

Definizione 3.4.4: Seconda Strategia di Havender

Quando un processo richiede una risorsa che gli viene negata, rilascia tutte le risorse accumulate fino a quel momento. Eventualmente, il processo richiederà tutte le risorse che ha perso, più quella che gli serviva.

Osservazioni 3.4.1

- È una tecnica costosa (perdere delle risorse può significare perdere un lavoro già compiuto in parte) e vale la pena solo se il sistema è tale per cui questa tecnica viene eseguita raramente.
- Il suo uso in congiunzione a un criterio di priorità che predilige l'assegnazione di risorse a processi che ne richiedono poche (può causare starvation).

Definizione 3.4.5: Terza Strategia di Havender

Ogni risorsa ha assegnato un numero utilizzato per quella risorsa soltanto, che le rende ordinabili in ordine crescente ($R_1 < R_2 < \dots < R_n$).

Osservazioni 3.4.2

- Un processo che ha bisogno di M risorse deve richiederle in ordine crescente.
- Non si può avere deadlock perché l'ordinamento delle richieste impedisce l'attesa circolare, ma non è molto flessibile.

3.4.3 Deadlock Avoidance

Definizione 3.4.6: Deadlock Avoidance

Non si può sempre evitare deadlock a priori. I metodi che consentono di fare ciò richiedono alcune informazioni, come per esempio il numero di risorse di cui hanno bisogno. L'algoritmo di deadlock avoidance esamina lo stato di allocazione delle risorse e garantisce che in futuro non si formeranno attese circolari.

Si introducono 2 nuove nozioni:

- Stato (del sistema) sicuro: si dice che un sistema è in un stato sicuro se il SO può garantire che ciascun processo finisca la propria esecuzione in un tempo finito.
- Sequenza sicura: una sequenza di processi viene detta sicura se le richieste che ogni processo deve ancora fare sono soddisfacenti usando le risorse attualmente libere più le risorse usate e liberate da altri processi.

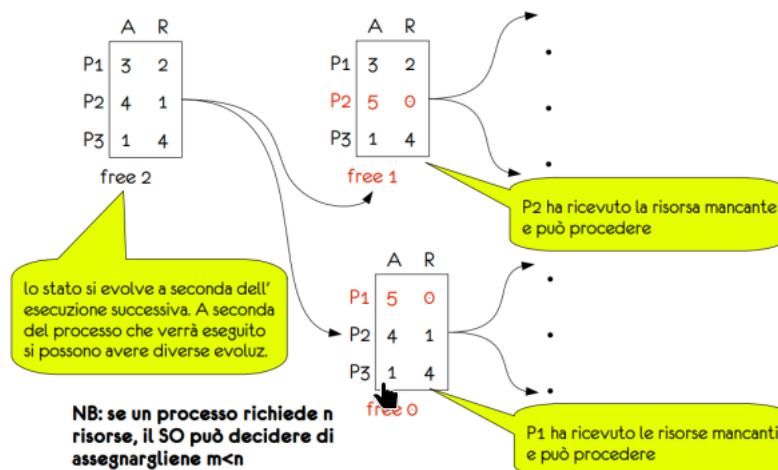


Figure 3.11: Stato di allocazione delle risorse.

Note:-

Uno stato è sicuro se da esso si dirama almeno una sequenza sicura, quindi se esiste almeno un ordinamento dei processi che è una sequenza sicura.
 Uno stato non sicuro non è necessariamente di deadlock ma può portare a esso.

Corollario 3.4.1 Algoritmo di Deadlock Avoidance

Variante del grafo di assegnazione che utilizza un terzo tipo di arco: l'arco di reclamo (claim edge).

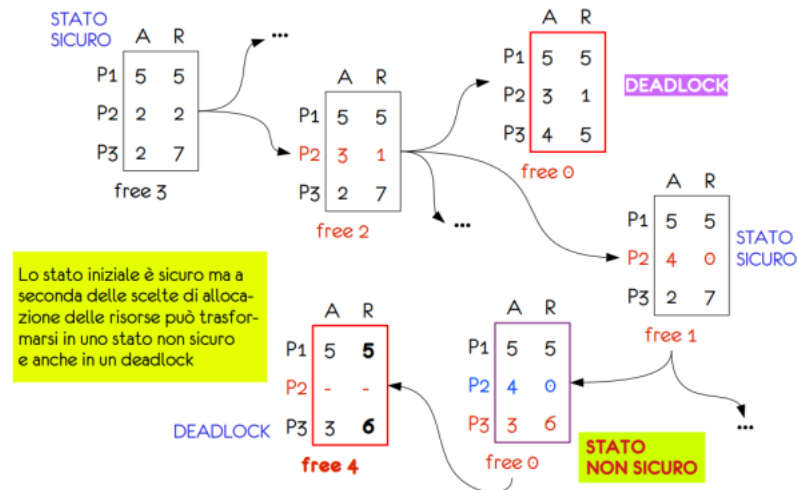


Figure 3.12: Esempio.

Funziona solo se ogni classe di risorsa ha una sola istanza. All'inizio tutti i processi inseriscono nel grafo di assegnazione un claim edge per ciascuna risorsa di cui avranno bisogno. Verranno trasformati in arco di richiesta solo se non si genera un ciclo (lo stato è quindi sicuro).

Alcuni algoritmi:

- Banchiere.
- Verifica della sicurezza.
- Gestione delle richieste.

3.4.4 Rottura del Deadlock

Tre possibili soluzioni:

- Terminare i processi coinvolti:
 - Terminare tutti i processi coinvolti.
 - Terminare un processo per volta fino alla risoluzione del deadlock, applicando l'algoritmo dopo l'abort di ciascun processo.
- Effettuare la prelazione delle risorse: sottrarre risorse ad altri processi per assegnarle ad altri. Anche qua bisogna identificare una vittima, su criteri economici.
- Riassegnare le risorse.

4

Gestione della Memoria

4.1 Introduzione

4.1.1 Indirizzi, Binding e Loading

Tipi di indirizzo:

- *Indirizzo logico*: indirizzo prodotto dalla CPU.
- *Indirizzo fisico*: indirizzo di ogni parola di memoria.

Definizione 4.1.1: Binding

Viene eseguito il collegamento tra lo spazio degli indirizzi logici (ind. prodotti dalla CPU) e lo spazio degli indirizzi fisici (ind. effettivi in RAM).

Osservazioni 4.1.1

- Quando il binding viene fatto a compile time, indirizzo logico = fisico.
- Quando il binding viene fatto a exec time, la corrispondenza deve essere calcolata (i due spazi di indirizzi non corrispondono).
- Il binding è a carico dell'MMU (Memory Management Unit), che in genere funziona con un meccanismo basato sul registro base (registro di rilocalizzazione).

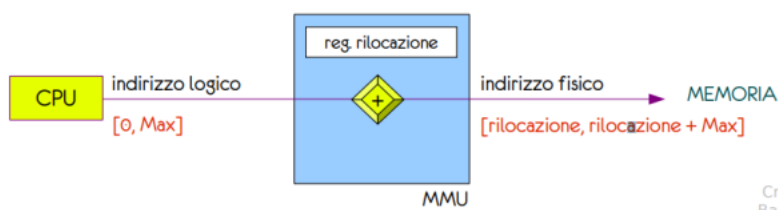


Figure 4.1: Binding.

Definizione 4.1.2: Linking

Il linking è il processo di composizione dei moduli che servono a un programma. Associa ai nomi di variabili e procedure usate da ciascun modulo (e non definiti in esso) le corrette definizioni.

Definizione 4.1.3: Loading

Il loading è copiare un programma eseguibile (o parte) nella RAM.

Note:-

Questi due processi si dicono dinamici quando sono svolti a tempo di esecuzione, statici quando precedono l'esecuzione.

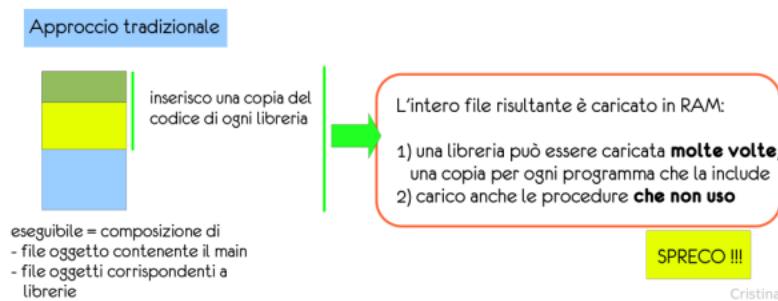


Figure 4.2: Linking statico.

In modo dinamico:

- Loading: una procedura è caricata in RAM alla sua prima invocazione.
 - Tutte le procedure risiedono in memoria secondaria sotto forma di codice rilocabile.
 - Molto vantaggioso rispetto a dover caricare l'intera libreria in RAM.
 - Funziona solo se il SO ci fornisce gli strumenti per realizzare librerie a caricamento dinamico.
- Linking: il collegamento del codice di una procedura al suo nome è effettuato alla sua prima invocazione (linker statico aggiunge solo uno stub della procedura).
 - Rimanda il collegamento reale di una libreria alla fase di esecuzione.
 - Dopo la compilazione, il linker statico arricchisce il programma aggiungendo gli stub relativi alle procedure appartenenti alle librerie dinamiche usate.
 - Durante l'esecuzione, lo stub verifica se il codice della procedura è già stato caricato in RAM:
 - * Se sì, sostituisce se stesso con l'indirizzo della procedura in questione.
 - * Se no, causa il caricamento del codice della procedura e poi si sostituisce con l'indirizzo del codice della procedura in questione.
 - Il vantaggio del linking dinamico risiede nella facilità di aggiornare le librerie condivise: se aggiorni una libreria dinamica, tutti i programmi che la usano faranno riferimento alla nuova versione senza dover ricompilare.

4.1.2 Allocazione della RAM

Esistono 3 approcci per la gestione della memoria principale:

- Allocazione *contigua*.
- *Paginazione*.

- *Segmentazione.*

Definizione 4.1.4: Allocazione Contigua

In questo modello si suddivide la RAM in due parti:

- Una parte riservata al SO, in genere allocata nella parte bassa.
- Una parte riservata ai processi utente.

Osservazioni 4.1.2

- Bisogna quindi proteggere la parte di memoria riservata al SO dalle letture/scritture dei processi utente, facilmente realizzabile usando un registro di rilocazione.
- Il codice del SO può essere, a sua volta, suddiviso in una parte sempre necessaria, e una parte che può essere utile o meno a seconda delle circostanze, chiamate *codice transiente*.
- Questa parte può essere rimossa dalla RAM quando non serve, e quindi ci occorre modificare la partizione riservata al SO.



Figure 4.3: Allocazione contigua.

Definizione 4.1.5: Allocazione a Partizioni Multiple

All'inizio, la RAM non utilizzata dal SO è libera. Il SO deve mantenere una lista di porzioni libere dette buchi. Un nuovo processo può essere caricato in RAM solo se esiste un buco abbastanza grande.

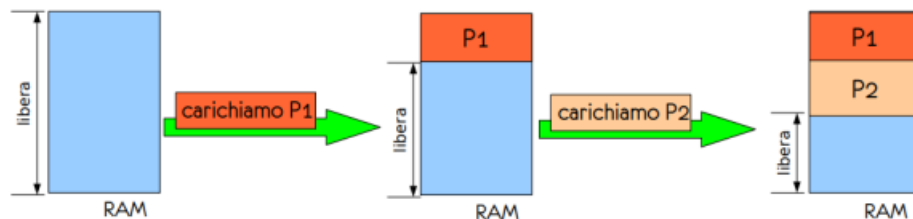


Figure 4.4: Allocazione a partizioni multiple.

Per caricare i processi in RAM, vengono utilizzati 3 criteri principali che aiutano a decidere quale porzione di memoria utilizzare:

- **Best-fit:** scelgo la porzione più piccola tra quelle adeguate a contenere l'immagine del processo.
- **First-fit:** scelgo la prima porzione abbastanza grande da contenere l'immagine del processo.
- **Worst-fit:** scelgo la porzione più grande tra quelle libere.

Note:-

Per capire qual è la migliore, bisogna analizzare la frammentazione della memoria.

Definizione 4.1.6: Frammentazione

La frammentazione è lo spezzettamento della memoria in tante parti; ne esistono due tipi:

- **Esterna:** se queste parti sono abbastanza grandi da poter essere utilizzabili.
- **Interna:** se sono troppo piccole per essere utilizzate, e di solito vengono unite alle partizioni precedenti.

Osservazioni 4.1.3

- La frammentazione è un problema, in quanto è molto facile avere ampie quantità di memoria inutilizzabile: secondo la regola del 50% usando il first-fit, per ogni N blocchi di memoria si hanno $\frac{N}{2}$ blocchi di memoria inutilizzabile ($\frac{1}{3}$ della memoria totale).
- In generale worst-fit è la strategia peggiore, first-fit e best-fit sono comparabili ma first-fit risulta computazionalmente meno costosa.
- Si può combattere la frammentazione attuando una politica di compattamento: spostare le immagini in memoria in modo che risultino contigue. Questo però, è applicabile solo se il binding tra indirizzi logici e fisici è effettuato a tempo di esecuzione.

Definizione 4.1.7: Swapping

In quanto la RAM ha dimensione limitata, è possibile che i processi in stato running/ready occupino più memoria di quanto ne sia disponibile. La soluzione consiste nel mantenere una parte dei processi ready in memoria secondaria e di tanto in tanto eseguire lo swapping, ovvero lo scambio dei processi tra le due memorie.

Swap:

- **Swap in:** si carica l'immagine di un processo ready dalla memoria secondaria (detta anche backing store) alla RAM.
- **Swap out:** si scarica l'immagine di un processo che non è in esecuzione dalla RAM alla memoria secondaria.

La collocazione del processo in RAM dipende da quando viene effettuato il binding delle variabili:

- Se il codice non è rilocabile, l'immagine del processo deve occupare la stessa sezione di RAM.
- Se il codice è rilocabile, è possibile inserirlo in qualsiasi posizione.

Corollario 4.1.1 Tempo di Swapping

Il tempo necessario al completamento dello swap è dato dal tempo di swap-out + tempo di swap in. Questo tempo è influenzato dalla dimensione della RAM usata dai singoli processi.

Note:-

Lo swapping non può essere effettuato se il processo sta effettuando operazione di I/O: queste operazioni non possono essere effettuate su variabili residenti in memoria secondaria.

4.2 Paginazione

4.2.1 Struttura

Definizione 4.2.1: Paginazione

La paginazione è un meccanismo di gestione della RAM alternativa alla allocazione contigua, e sua caratteristica fondamentale è che consente allo spazio degli indirizzi fisici di un processo di non essere contiguo. Questo consente di ridimensionare in modo dinamico lo spazio riservato ad un processo, aggiungendo o togliendo pagine su richiesta.

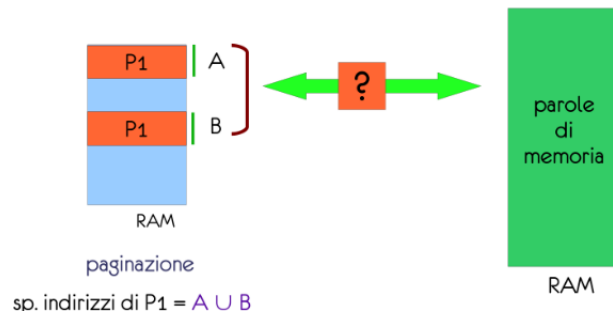


Figure 4.5: La paginazione.

Domanda 4.1

Come si associano le porzioni di processo a RAM? Come si organizzano le diverse porzioni in un tutt'uno?

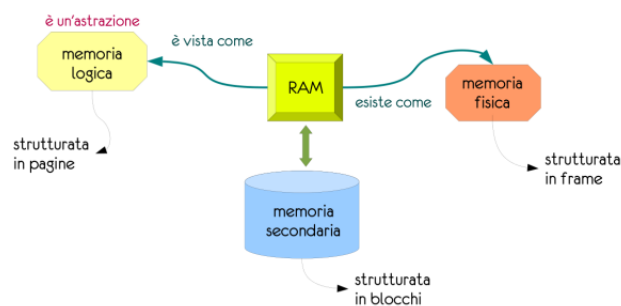


Figure 4.6: Struttura.

Note:-

Per fare il binding tra indirizzi logici e fisici, usiamo una tabella degli indirizzi.

Corollario 4.2.1 Dimensione di una Pagina

La dimensione è la stessa per tutte le pagine ed è definita dall'architettura (tra $2^9 = 512B$ e $2^{24} = 16MB$).

La memoria logica viene suddivisa secondo questa formula:

- Dimensione di una pagina = 2^n .
- Dimensione della memoria logica = 2^m .
- Numero di pagine = 2^{m-n} .
- Bit per rappresentare il numero di pagina = $m - n$.
- Bit per rappresentare lo scostamento (offset) all'interno di una pagina = n .

Note:-

L'indirizzo logico di una pagina è quindi una coppia $\langle p, o \rangle$ in cui p è il numero di pagina e o è l'offset.

Definizione 4.2.2: Rilocalizzazione

La rilocalizzazione nel contesto della paginazione significa consentire l'accesso ad una pagina indipendentemente dal frame in cui è caricata. Per effettuare la rilocalizzazione, il registro di rilocalizzazione è sostituito dalla entry nella tabella delle pagine, corrispondente a p . Il valore f individua l'indirizzo di inizio del frame.

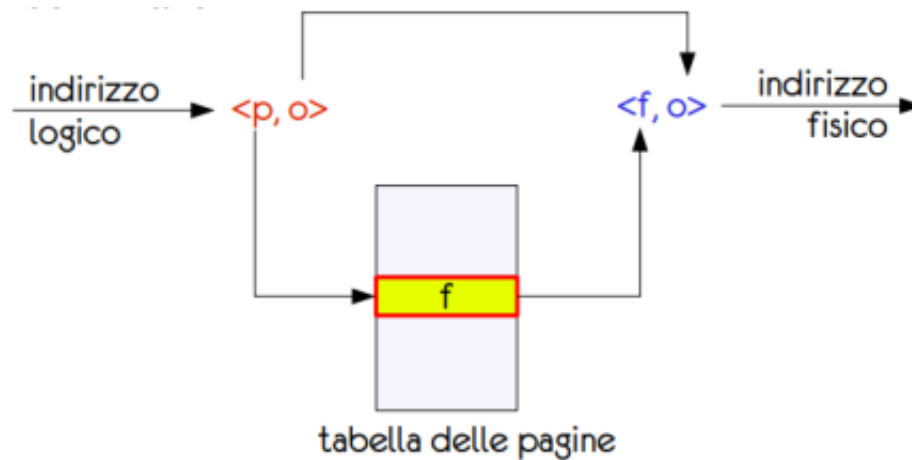


Figure 4.7: Rilocalizzazione.

Osservazioni 4.2.1

- La paginazione elimina il problema della frammentazione esterna, ma rimane il problema della *frammentazione interna*.
- Di ogni processo, solo l'ultima sua pagina allocata può presentare frammentazione interna, perché raramente la dimensione di un processo sarà un multiplo della grandezza di una pagina.
- Si può perciò dire che in media si ha mezza pagina inutilizzata per processo.

Occorre quindi trovare una dimensione ottimale per le pagine in modo da ridurre la frammentazione:

- Utilizzando pagine di piccole dimensione si riduce la frammentazione interna, ma aumenta il numero di pagine.
- Utilizzando pagine di grandi dimensioni si aumenta la memoria inutilizzata, ma il traferimento di dati dalla memoria secondaria è più veloce.

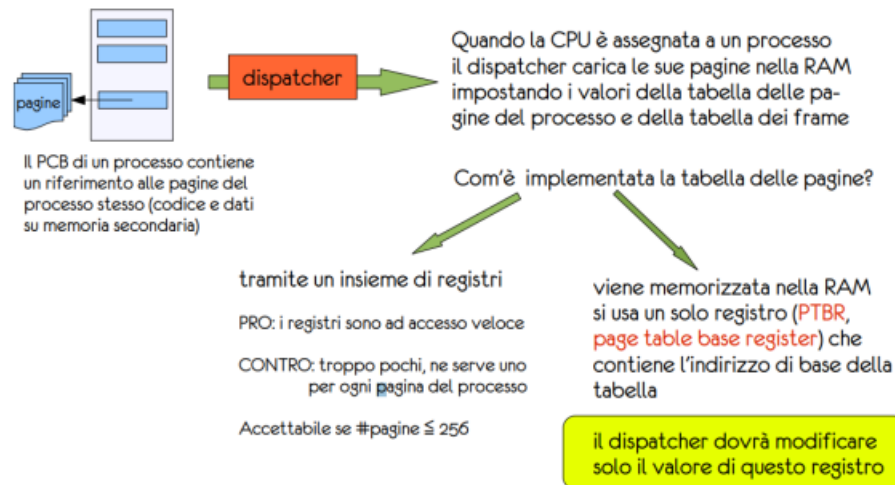


Figure 4.8: Architettura generale della paginazione.

4.2.2 Architettura di Paginazione

Tempi di accesso:

- Si accede alla RAM 2 volte, una per individuare la page table e un'altra per prelevare l'indirizzo fisico d'interesse.
- Questo doppio accesso rallenta troppo l'esecuzione, perciò viene implementata la *TLB* (*Translation Look-aside Buffer*), una cache molto veloce contenente delle coppie <chiave, valore>.
- Quando riceve un input lo confronta contemporaneamente con tutte le chiavi. Se trova una chiave corrispondente restituisce il valore associato.

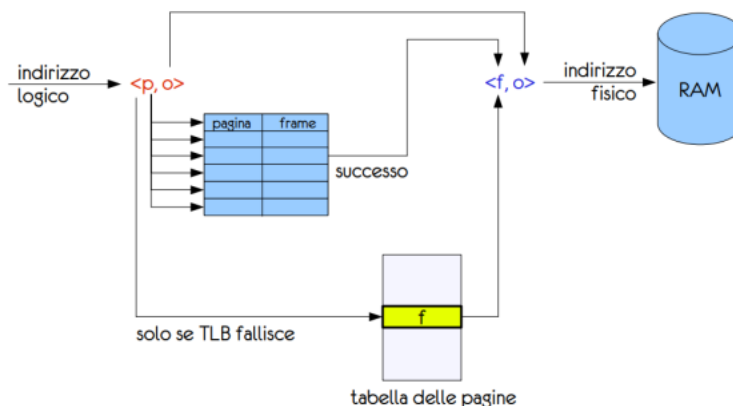


Figure 4.9: Accesso mediante TLB.

Definizione 4.2.3: Hit Ratio

È la percentuale di successo di reperimento di una pagina tramite TLB (TLB hit), e consente di calcolare il tempo medio effettivo di accesso a una pagina.

Esempio 4.2.1 (Hit Ratio)

Tempo di accesso al TLB = 20 nsec

Tempo di accesso alla RAM = 100 nsec

Hit Ratio = 0.82

Percentuale di accessi tramite page table = 0.18

 $T_{TLB} = \text{accesso TLB} + \text{accesso RAM} = (20 + 100) \text{ nsec} = 120 \text{ nsec}$ $T_{PT} = 2 * \text{accesso RAM} = 200 \text{ nsec}$ $T_M = (0.82 * T_{TLB}) + (0.18 * T_{PT}) = 134.4 \text{ nsec}$ **Quando non trovo una pagina di interesse (TLB miss):**

- Se c'è spazio nel TLB si inserisce la nuova coppia <pagina, frame>.
- Se non c'è spazio nel TLB, si sostituisce una coppia (in genere quella meno recente, LRU) con quella nuova.

Nella tabella delle pagine del processo viene anche memorizzata la modalità di accesso alle pagine, che può essere:

- Read-only.
- R/W.
- Exec.
- Invalid (non appartiene allo spazio degli indirizzi del processo).

Note:-

Si possono avere pagine non valide quando la tabella delle pagine di un processo è più piccola dello spazio riservato ad essa nella RAM. Se un processo tenta di accedere ad una pagina non valida, viene generato un interrupt.

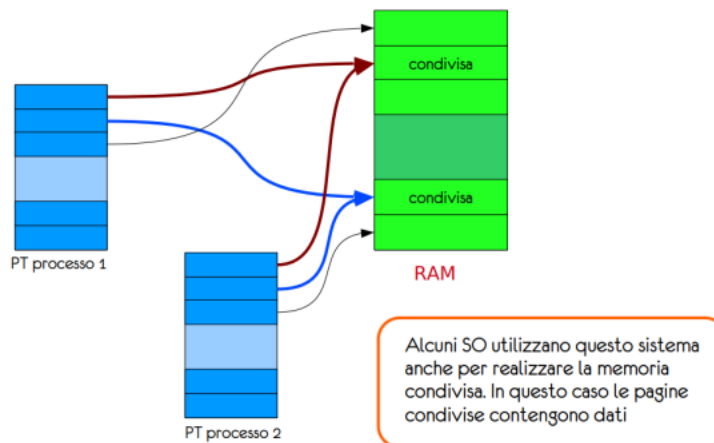


Figure 4.10: Condivisione delle pagine.

Definizione 4.2.4: Paginazione Multilivello

La dimensione della tabella delle pagine dipende dalla dimensione dello spazio degli indirizzi logici. In realtà, nessun processo usa l'intero spazio degli indirizzi. La maggior parte della tabella andrebbe sprecata. Per risolvere questo problema, possiamo suddividere la tabella in tante parti organizzate.

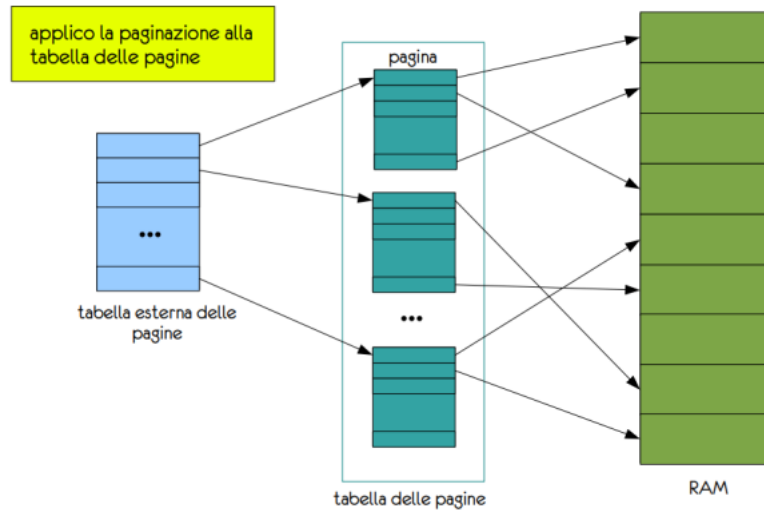


Figure 4.11: Paginazione a 2 livelli.

Vantaggi:

- Rappresentazione naturale se si adotta un approccio processo-centrico.
- Ordinata in modo tale che sia semplice per il SO identificare il punto in cui si trova l'indirizzo del frame di interesse.

Svantaggi:

- Può diventare troppo grande e invece di aiutare a gestire la RAM diventa essa stessa un problema da gestire.

4.2.3 Segmentazione**Definizione 4.2.5: Segmentazione**

La segmentazione è una modellazione della memoria in cui ogni processo ha una porzione di RAM organizzata come un insieme di segmenti di memoria di dimensione variabile.

Note:-

Per organizzare la memoria in questo modo, il compilatore del programma ha il compito di organizzare la memoria in tali segmenti. Spesso è più efficiente della paginazione.

Per esempio, per un programma in C serviranno i segmenti:

- Codice.
- variabili globali.
- Stack.
- Heap.
- Libreria standard.
- Librerie esterne.

Problemi:

- Occorre gestire dinamicamente la memoria libera.
- Si ha frammentazione esterna.

Esempio 4.2.2 (Pentium Intel)

Il Pentium Intel utilizza una tecnica mista tra paginazione e segmentazione.

Esempio: Pentium Intel

Ogni segmento è strutturato in pagine

dimensione max segmento 4GB
numero max di segmenti per processo 16K

Il processore ha 6 registri di segmento che permettono a un processo di fare riferimento a 6 segmenti contemporaneamente

8K riservati ∈ tabella locale dei descrittori
8K condivisi ∈ tabella globale dei descrittori

4.2.4 Memoria Virtuale**Definizione 4.2.6: Memoria Virtuale**

Separando la memoria logica dalla memoria fisica, nasce il concetto di memoria virtuale.

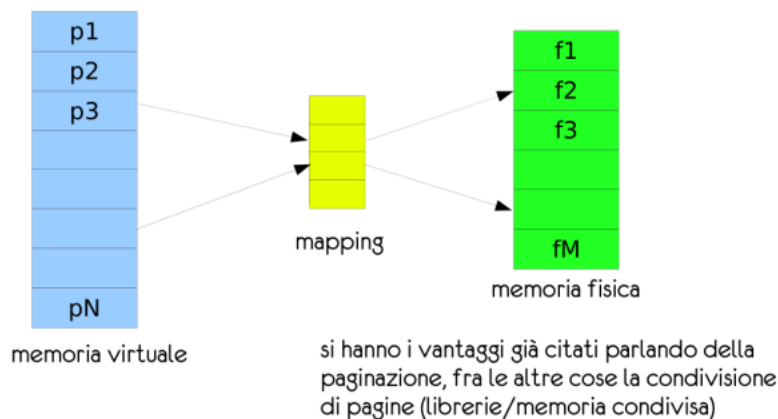


Figure 4.12: Memoria virtuale.

Definizione 4.2.7: Spazio degli Indirizzi

Si tratta di uno spazio predefinito, in genere molto più grande dello spazio occupato dai processi di riferimento. Con l'esecuzione, stack e heap cresceranno uno verso l'altro. Per contenere nuovi dati potranno essere allocati nuovi frame.

Note:-

Questo non cambia lo spazio degli indirizzi virtuali del processo, semplicemente una parte di indirizzi prima non corrispondenti ad alcun indirizzo assoluto saranno ora mappati su parole di memoria effettive.

Definizione 4.2.8: Lazy Swapping

I processi vengono caricati in modo parziale: quando una pagina diventa utile, la si carica in RAM al posto di una pagina già presente, che verrà scaricata in memoria secondaria.

Osservazioni 4.2.2

- Occorre definire un meccanismo per mantenere traccia della locazione delle pagine (RAM o mem. secondaria).
- La parte del SO che gestisce il caricamento delle pagine è detto *pager*.

Per ogni pagina nella page table di un processo è presente un bit che indica la validità della pagina:

- Una pagina è *valida* se è presente in RAM e appartiene allo spazio degli indirizzi del processo.
- Una pagina è *invalida* se è in memoria secondaria o non appartiene allo spazio degli indirizzi del processo.

Definizione 4.2.9: Page Fault

Quando un processo tenta di accedere ad una pagina non valida, viene generato un interrupt (page fault exception):

- Se questa pagina è conservata in memoria secondaria (ed è quindi valida), viene copiata da disco in un frame libero, si aggiorna la tabella delle pagine e il processo continua come se non fosse successo niente.
- Altrimenti, se la pagina è invalida, si termina il processo.

Gestione del page fault:

1. Si genera un'interruzione.
2. Salvataggio dei registri e dello stato del processo.
3. Verifica dell'interruzione: in questo caso si determina che si tratta di page fault.
4. Controllo della correttezza del riferimento e individuazione della locazione occupata dalla pagina mancante sul disco.
5. Lettura e copiatura della pagina da memoria secondaria in RAM (operazione di I/O).
6. Durante l'attesa per il completamento di questa operazione, allocazione della CPU a un altro processo.
7. Interrupt che segnala il completamento dell'operazione di caricamento della pagina.
8. Aggiornamento della tabella delle pagine.
9. Quando lo scheduling riavvia il processo sospeso, ripristino dello stato.
10. Riesecuzione dell'istruzione interrotta.

Osservazioni 4.2.3

Riassumendo:

- Servizio di interruzione per page fault.
- Lettura della pagina.
- Riavvio del processo.

Definizione 4.2.10: Area di Swap

Le pagine di un processo che non sono in RAM sono contenute in memoria secondaria, che viene vista proprio come un'estensione della RAM (nonostante i tempi di accesso maggiori).

Osservazioni 4.2.4

- Ogni SO gestisce l'area di swap in maniera molto diversa: sia per il dimensionamento, sia per i suoi contenuti. Per quanto riguarda la dimensione:
 - Linux suggerisce di allocare il doppio della quantità di RAM (outdated).
 - Solaris suggerisce di allocare una quantità pari alla differenza fra la dimensione dello spazio degli indirizzi logici e la dimensione della RAM.
- Per l'implementazione, ci sono due approcci:
 - Come file: l'area di swap diventa un file speciale del file system:
 - * Pro: si evita il problema di dimensionamento, perchè un file si può ridimensionare secondo le necessità.
 - * Contro: la gestione del file system rallenta ulteriormente l'accesso.
 - partizione a sé, non formattata, in cui si usa un gestore speciale per usare algoritmi ottimizzati per ridurre i tempi di accesso:
 - * Pro: maggiore velocità.
 - * Contro: per ridimensionarla bisogna ripartizionare il disco.

Note:-

Si preferisce mantenere solamente i dati nello swap: il codice tanto si può prelevare dal file system.

Definizione 4.2.11: Processi Padre e Figli

Quando un processo P esegue una fork, viene creato un processo “figlio”, che esegue lo stesso codice del “padre” (processo che lo ha generato) e ha una copia delle var. e dello stack del padre.

Note:-

Tramite codice è possibile (e comunemente usato) cambiare il codice che il figlio deve eseguire.

Con l'esecuzione della fork sorgono 2 problemi relativi alla gestione della memoria tra i processi:

- *Duplicazione delle pagine*: risolta condividendo le pagine tra processi padre e figlio.
- *Sovrascrittura delle pagine*: risolta usando la tecnica di “copiatura su scrittura”: quando uno dei due processi esegue una exec, si allocano delle nuove pagine per il processo chiamante e si carica il nuovo codice in esse.

Definizione 4.2.12: Sostituzione delle Pagine

Avviene quando, all'avvenimento di un page fault, non si trova spazio in RAM per la pagina richiesta. Il meccanismo richiede la copiatura di due pagine: vittima copiata in mem. secondaria, nuova pagina copiata in RAM.

Per evitare molteplici casi di duplice scrittura, si cerca di limitarla solo ai casi in cui serve:

- Serve copiare una volta per mettere la pagina nuova in RAM.
- Serve copiare un'altra volta per copiare la vittima in mem. secondaria, SOLO SE essa è stata modificata rispetto ad una sua copia già presente.

Note:-

Per controllare i casi di modifica viene mantenuto un bit per ogni pagina (chiamato *dirty bit*) per indicare se la pagina è stata modificata o meno.

Osservazioni 4.2.5

- Memoria logica e fisica sono completamente separate.
- Lo spazio di indirizzamento logico dei processi è maggiore dello spazio di indirizzamento fisico offerto dalla RAM.
- Serve realizzare un meccanismo dinamico tramite il quale caricare/sostituire pagine a seconda dell'esigenza.
- Serve introdurre meccanismi per aumentare l'efficienza, in quanto lettura e copiatura sono operazioni pesanti.
- Incremento del livello di multiprogrammazione.

4.2.5 Paginazione su Richiesta**Definizione 4.2.13: Paging on Demand**

Anche detta paginazione a richiesta, è un meccanismo per cui una pagina viene caricata in RAM solo se è stata richiesta (i processi vengono caricati senza pagine)

Note:-

Non servono ulteriori strutture per applicare questo meccanismo e vale il principio di località dei riferimenti, ma potrebbe essere più pesante di un meccanismo normale.

per implementare la paginazione su richiesta occorre sviluppare due algoritmi:

- *Allocazione del frame*: spartisce M frame liberi fra N processi.
- *Sostituzione delle pagine*: secondo un criterio predefinito e informazioni prese dalle pagine, sceglie le pagine da sostituire.

Alcune implementazioni:

- FIFO (First-In-First-Out).
- Ottimale.
- LRU (Last-Recently-Used).
- Approssimazioni di LRU:
 - Seconda chance.
 - Bit supplementare.
- Conteggio:
 - LFU (Least-Frequently-Used).
 - MFU (Most-Frequently-Used).

Definizione 4.2.14: Sostituzione FIFO

gni pagina si associa un marcatore temporale: quando è stata caricata in memoria. Si sceglie poi di sostituire la pagina più vecchia in memoria. Questo è possibile perché le pagine sono organizzate in una coda FIFO: andremo infatti a prendere sempre la prima pagina in coda.

Osservazioni 4.2.6

- Questa tecnica è molto semplice da realizzare, ma non si comporta sempre bene: il fatto che una pagina sia stata caricata da tempo non significa per forza che non sia più in uso.
- Questo algoritmo presenta anche l'*anomalia di Belady*: la frequenza delle assenze può aumentare con l'aumentare del numero di frame in RAM.

Definizione 4.2.15: Algoritmo Ottimale (OPT)

È l'algoritmo migliore: non presenta l'anomalia di Belady, e ha la frequenza di page fault più bassa. Si sostituisce la pagina che non verrà usata per il periodo di tempo più lungo.

Note:-

Non è implementabile, perché non sappiamo quando una pagina verrà richiesta.

Definizione 4.2.16: LRU

Simile all'algoritmo FIFO, basiamo la scelta sul tempo di utilizzo. Questo algoritmo è un'approssimazione dell'OPT, perché ci basiamo su quanto sia stata usata attualmente una pagina. Scegliamo sempre la pagina usata meno di recente.

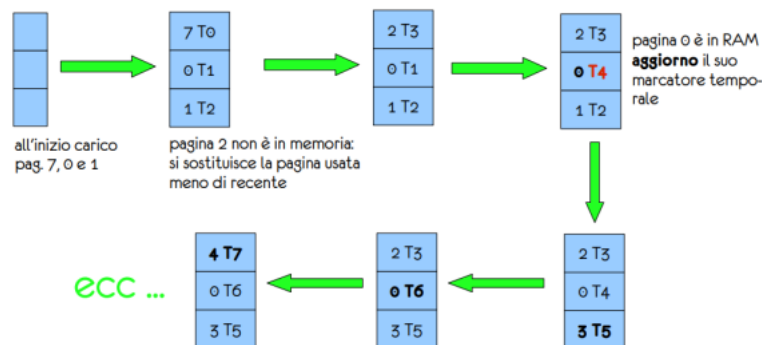


Figure 4.13: LRU.

Definizione 4.2.17: Algoritmo con Bit Supplementare

Si associa a ogni elemento nella tabella delle pagine una serie di bit che realizzano registri a scorrimento. A intervalli regolari un timer passa il controllo al SO, che:

- Sposta i bit nella sequenza traslandoli a destra di 1.
- Copia i bit di riferimento nel bit più significativo.
- Scarta il bit meno significativo.

Note:-

Quindi i bit di riferimento delle pagine vengono azzerati: una pagina che ha il suo registro a 00000000 non è stata usata negli ultimi 8 periodi di tempo.

Definizione 4.2.18: Algoritmo di Seconda Chance

Unisce il metodo LRU a FIFO: le pagine sono mantenute in una coda circolare FIFO e mantengono un bit di riferimento.

Funzionamento:

- Carico la pagina e le associo un bit di riferimento impostato ad 1.
- Quando diventerà necessario caricare una nuova pagina, la ricerca partirà dalla posizione successiva, in questo caso il bit di rif. ora vale 1.
- Quando il bit di riferimento è impostato a 1: la pagina viene saltata ma il suo bit di riferimento viene azzerato. Idem per la pagina successiva.
- A questo punto si incontra una pagina con bit di riferimento a 0 e la si sovrascrive mettendo il bit di riferimento a 1.

Osservazioni 4.2.7

- La struttura circolare della coda è un modo per concedere un po' di tempo in RAM a ciascuna pagina: infatti una pagina con bit che passa da 1 a 0 non viene rimossa subito ma solo quando si tornerà alla sua locazione dopo aver percorso tutta la lista.
- Se tutti i bit di rif. vengono impostati a 1, diventa FIFO.

Definizione 4.2.19: Algoritmo di Seconda Chance Migliorato

Ogni pagina ha associata una coppia di bit $\langle r, m \rangle$:

- Riferimento: indica se la pagina è stata usata di recente.
- Modificata: indica se la pagina è stata modificata di recente.

Si individuano 4 classi di pagine:

- $\langle 0, 0 \rangle$: non usata, non modificata.
- $\langle 0, 1 \rangle$: non usata, modificata.
- $\langle 1, 0 \rangle$: usata, non modificata.
- $\langle 1, 1 \rangle$: usata, modificata.

Note:-

Il criterio di scelta è basato sulla priorità delle 4 classi.

Sostituzioni su conteggio:

- LFU: sostituisce la pagina col minor numero di riferimenti. Si basa sull'idea che una pagina molto usata avrà un conteggio alto. Diventa però difficile distinguere fra una pagina che è stata molto usata in passato e una che è molto usata di recente.
- MFU: sostituisce la pagina col maggior numero di riferimenti. Si basa sull'idea che se una pagina ha un contatore basso è probabilmente stata usata di recente.

Definizione 4.2.20: Pool of Free Frames

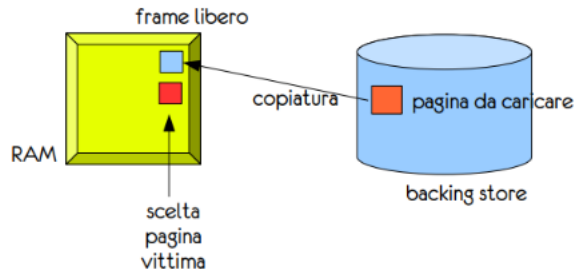
Spesso gli algoritmi di sostituzione delle pagine sono affiancati da altre procedure finalizzate a incrementare le prestazioni del sistema. L'idea è associare un piccolo pool di frame liberi.

Quando diventa necessario caricare una pagina nuova:

- La si copia in un frame libero associato al processo.
- Si sceglie una vittima e la si copia in mem. secondaria.
- Si libera quindi un frame che viene aggiunto al pool di frame liberi.

Note:-

In questo modo, non viene cancellata la vittima, che non dovrà essere ricopiata se nasce il bisogno di accedervi nuovamente.



La pagina vittima va ad arricchire il pool dei frame liberi assegnati al processo ma non viene cancellata o sovrascritta, al contrario rimane accessibile attraverso la tabella delle pagine

Figure 4.14: Pool of free frames.

4.3 Allocazione e Thrashing

4.3.1 Allocazione

Definizione 4.3.1: Allocazione dei Frame

L'algoritmo di allocazione dei frame si applica quando si hanno a disposizione N frame liberi, occorre caricare 1 o + processi e bisogna decidere come spartire i frame. In generale, il numero di page fault è inversamente proporzionale al numero di frame allocati per ciascun processo. L'idea è quindi di mantenere in memoria un numero minimo di frame per processo, per ridurre la probabilità che si generi page fault.

Il numero minimo di frame dipende dall'istruzione e dall'architettura:

- Se ha un solo operando occorrono ≥ 2 pagine (codice e dati).
- Se è un riferimento alla memoria occorrono ≥ 3 pagine (codice, dati, riferimento).
- Se è grande, può stare a cavallo tra due pagine.
- Può avere più livelli di indirizzamento indiretto, tutta la mem. virtuale potrebbe dover essere caricata.

Definizione 4.3.2: Allocazione Proporzionale

Indichiamo con VM^i la dimensione della memoria logica occupata dal nostro processo. Indichiamo poi con M il numero di frame disponibili, il numero di frame allocati al processo sarà:

$$m = \frac{VM^i}{V} * M$$

Note:-

Se abbiamo poi definito un numero minimo di frame necessari per caricare un'istruzione, potrebbe essere necessario incrementare tale valore per i processi più leggeri, o diminuirlo per quelli più pesanti.

Corollario 4.3.1 Dinamicità

Se il numero di processi in RAM aumenta bisognerà redistribuire il numero di frame ancora liberi ai nuovi processi, se diminuisce si può assegnare un numero maggiore di frame.

4.3.2 Thrashing**Definizione 4.3.3: Thrashing**

Il thrashing è un fenomeno che si verifica nella gestione della memoria virtuale: quando una pagina in realtà in uso da un processo viene sostituita, e subito dopo viene scatenato di nuovo un page fault, avviene questo fenomeno. Questo loop può continuare per molto tempo, e si spenderà più tempo nel sostituire le pagine che nell'eseguire il processo.

Note:-

Quando si ha thrashing perché i processi hanno a disposizione meno frame del numero di pagine attive, l'attività della CPU tende a diminuire: i processi tendono a rimanere in attesa del completamento di operazioni di I/O.

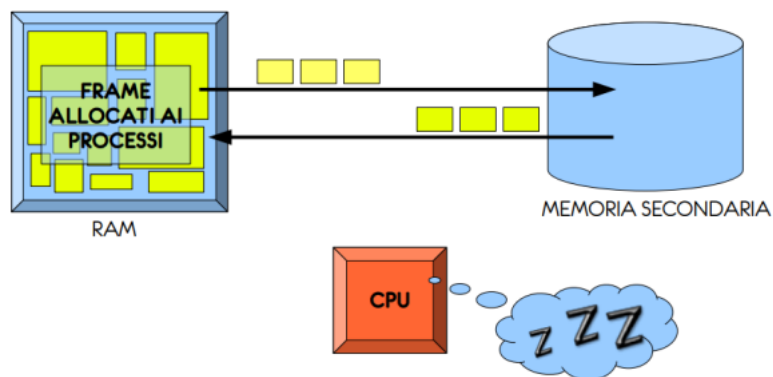


Figure 4.15: Effetti del thrashing.

Osservazioni 4.3.1

- Molti SO monitorano l'uso della CPU: quando questo cala, se ci sono processi conservati in mem. secondaria, portano in RAM qualche processo in più.
- Ai nuovi processi vengono allocati frame sottratti ad altri processi in RAM se la strategia di allocazione è globale, quindi aumenta la frequenza di page fault, diminuisce l'uso della CPU e si ripete fino al blocco totale.
- Si ha thrashing perché la politica di allocazione dei frame è globale: il SO può sottrarre frame ad un processo per assegnarli ad un altro, rendendo degeneri più processi alla volta.

4.3.3 Pagine Attive e Working Set

L'ideale è cercare di prevedere di quante pagine avrà bisogno un processo e assegnargli un numero sufficiente di frame. L'allocazione dipenderà dal numero di pagine attive per processo, e si distribuiranno frame liberi solo ai

processi che ne hanno bisogno. Poiché il numero di pagine attive varia nel tempo, se cresce si allocano nuovi frame e se decresce si liberano frame.

Definizione 4.3.4: Working Set

Si basa sul principio di località: per periodi di durata consistente, il programma avrà accesso solo ad un sottoinsieme del suo codice, variabili globali e locali (quindi un sottoinsieme di pagine, le pagine attive).

Note:-

Il WS è quindi l'insieme delle pagine attive di un processo in un certo periodo di tempo.

Definita la finestra σ è possibile calcolare il WS di ogni processo, e quindi calcolare la quantità di frame necessari ai processi in esecuzione:

$$D = \sum_{i=1}^n |WS^i|$$

- Quando $D >$ numero di frame liberi si genera thrashing: significa che qualche processo non ha a disposizione frame a sufficienza.
- Se $S <$ numero di frame liberi è possibile avviare nuovi processi (limitandosi ai frame disponibili).

Note:-

Quando il WS di un processo cresce, se non ci sono frame liberi, il SO sceglie uno dei processi, lo copia in memoria secondaria, lo sospende e assegna (parte dei) suoi frame al processo richiedente: così si evita il thrashing.

Definizione 4.3.5: Prepaginazione

La prepaginazione consiste nel memorizzare insieme al PCB anche il WS, in questo modo si caricano in memoria tutte le pagine utili subito.

Memoria virtuale e file:

- Quando i processi agiscono sui file tramite system call `read()` e `write()`, richiedono accessi sequenziali alla memoria secondaria, in quanto il file stesso non fa parte del loro codice.
- Eseguire operazioni direttamente sul disco comporta bassa efficienza: sarebbe meglio eseguire le operazioni sulla RAM (memoria più veloce) e riportare ogni tot le modifiche in mem. secondaria.

Definizione 4.3.6: Mappatura dei File

Procedimento nel quale si associano degli indirizzi virtuali di un processo ad una parte di un file a cui ha accesso.

Note:-

Scrivere su RAM renderà le operazioni di I/O più veloci, e se un file è usato da più processi, lo si carica in RAM una singola volta e si condividono le pagine.

È possibile utilizzare diverse strategie per riportare le modifiche dalle copie dei file in RAM ai file effettivi presenti in mem. secondaria:

- *Controllo periodico*: se durante il controllo la pagina risulta modificata, la si copia in mem. secondaria.
- *Chiusura del file*.
- Operazioni di *flush* esplicite nel codice del programma.

Definizione 4.3.7: Mappatura I/O

Oltre ai file, alcuni SO mappano anche i registri dei controller dei dispositivi di I/O. Così l'interazione con questi device avviene solamente scrivendo/leggendo le porzioni della RAM associate ai device.

Esempio 4.3.1 (Page Fault Malgestiti)

Supponiamo di dover inizializzare una matrice. Abbiamo due possibilità scorrere le colonne per riga o scorrere le righe per colonna.

Sebbene siano equivalenti dal punto di vista del risultato la seconda causa molti più page fault del necessario dato che le matrici sono salvate per righe.

Definizione 4.3.8: Allocazione dei Frame per il Kernel

Differente dall'allocazione per i processi utente a causa di caratteristiche differenti dei processi kernel:

- Necessità di strutturare dati di dimensioni variabili, spesso molto più piccole di una pagina.
- Alcuni dispositivi interagiscono direttamente con la RAM: c'è bisogno di avere aree di memoria contigua.
- Codice e dati del kernel non sono sottoposti a paginazione.

Corollario 4.3.2 Sistema Buddy

Il sistema buddy, anche noto come sistema gemellare, che consiste nell'utilizzo di pagine fisicamente contigue allocate in memoria in unità di dimensioni pari a potenze di 2 (4KB, 8KB, 16KB...), e arrotondando poi per eccesso le richieste (es: se servono 6KB, si allocano 8KB).

Note:-

Lo svantaggio principale di questo sistema è l'alta frammentazione interna: per allocare 33KB uso 64KB, per allocare 65KB ne uso 128 ecc.

Definizione 4.3.9: Allocazione a Slab

Una slab è una sequenza di pagine fisicamente contigue, e una cache un insieme di slab. Con questo sistema viene mantenuta una cache per ogni tipo di strutture dati usate dal SO, per esempio: i semafori, PCB... Ogni cache contiene delle istanze, e queste istanze possono essere libere o occupate. Quando viene richiesta una nuova istanza, il SO cerca la cache in questione, e se c'è un'istanza libera usa quella. Se non la trova, alloca una nuova slab.

Note:-

Questa tecnica è molto efficiente, perchè elimina il problema della frammentazione interna. Introdotta in Solaris, ora è usata anche da Linux (in precedenza usava il buddy system).

5

Filesystem

5.1 Introduzione

5.1.1 Cos'è il Filesystem?

Definizione 5.1.1: Filesystem

Programmi e dati sono memorizzati in memoria secondaria, organizzata in file e directory che insieme formano il file system.

Corollario 5.1.1 File

Il file è un'astrazione: è un insieme di informazioni correlate (a discrezione dell'autore) a cui è associato un nome.

Note:-

Per l'utente, i file sono gli elementi base in cui è organizzata la memoria, ma non sono uguali tra loro.

Ogni file può essere diviso in due parti:

- Contenuto del file.
- Metadati¹
 - Nome.
 - Tipo.
 - Dimensione.
 - Proprietario.
 - Protezione.
 - Data/Ora di modifica.

Note:-

In alcuni sistemi, per associare al file un suo tipo viene utilizzato un codice all'interno del file stesso, chiamato magic number.

¹Insieme di caratteristiche che descrivono il file stesso:

Corollario 5.1.2 Estensioni

Le estensioni dei file (.pdf, .java, .xml ecc.) servono per definire un particolare formato con cui i dati sono organizzati: i file immagine sono quindi organizzati diversamente dai file di testo, ma anche un file .docx (documento di Word) è organizzato diversamente da un file .html

Note:-

Curiosità: i docx files sono solo zip in disguise.

I file possono essere caratterizzati anche da tipologie più generali rispetto alla specifica applicazione che li gestirà, in particolare:

- *File alfanumerici*: contengono una sequenza di caratteri:
 - Testi.
 - Sorgenti.
- *File binari*: contengono byte organizzati secondo una struttura precisa ben diversa da un file di testo:
 - File oggetto.
 - Eseguibili.

5.1.2 Organizzazione del Filesystem

Esistono due visioni classiche:

- *Organizzazione piatta*: insieme di file tutti allo stesso livello (appunto “piatto”), non possono esserci file con lo stesso nome.
- *Organizzazione gerarchica*: memoria organizzata ad albero, i cui nodi interni fungono da contenitori e sono chiamati *directory*.

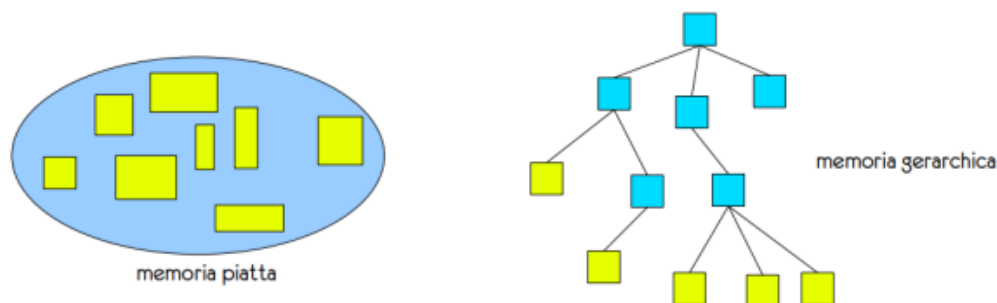


Figure 5.1: Tipi di organizzazione.

Un file aperto da un processo è inteso come una risorsa in uso da parte del processo:

- L'*apertura* di un file comporta l'allocazione di un insieme di risorse di file system che consentono l'accesso al file stesso.
- *Chiudere* un file significa rilasciare le risorse di file system allocate. Se il processo era l'unico utilizzatore del file, le pagine di RAM usate per consentire un'elaborazione efficiente dei dati possono essere liberate.

Osservazioni 5.1.1

- In genere, sia apertura che chiusura restituiscono/richiedono un handle come argomento, che definisce il file aperto/il file da chiudere.
- In C la funzione `fopen()` restituisce un `FILE*` (puntatore al file).
- La `fclose()` richiede come parametro il `FILE*` ottenuto con `fopen`.

Anche lettura e scrittura richiedono un handle, ma occorre anche definire il punto del file da cui leggere o scrivere:

- **Accesso sequenziale:** non si specifica esplicitamente la posizione, ma il SO mantiene un puntatore alla posizione corrente (mantenuta nella tabella dei file aperti del processo):
 - Quando si legge, si fa avanzare un puntatore che indica la posizione raggiunta all'interno del file.
 - Quando si scrive, il contenuto viene aggiunto al fondo del file.
- **Accesso diretto:** si può leggere/scrivere in punti specifici del file, che vanno indicati espressamente:
 - Il file è visto come una sequenza di record di uguale dimensione.
 - Conoscendo dimensione e posizione dei record è possibile accedervi direttamente.
- **Accesso a indice:** un file indicizzato è costituito da due file:
 - Il file dei contenuti, memorizzati in un formato specifico.
 - Il file indice, con i riferimenti ai record.

Definizione 5.1.2: Protezione dei File

Per gestire gli accessi ai file, alcuni SO richiedono di indicare la modalità di apertura di un file: solo lettura, lettura/scrittura, esecuzione. Certi SO consentono di associare ai file dei diritti di accesso che indicano le modalità di apertura del file consentite agli utenti.

Alcuni SO consentono di associare dei lock ai file:

- **Lock condiviso** (lettura): consente a n processi di effettuare determinate operazioni sullo stesso file, anche in parallelo.
- **Lock esclusivo** (scrittura): solo il processo che detiene il lock può usare il file.

Note:-

Questi lock possono essere consigliati od obbligatori.

5.2 Directory

5.2.1 Che Cos'è una Directory?

Definizione 5.2.1: Directory

È un entità che può contenere file o altre directory. Dal punto di vista astratto, la directory non è altro che una tabella che consente di accedere ai contenuti di un file a partire dal suo nome. Le operazioni possibili sono le solite: scrittura, lettura, ricerca, attraversamento.

Osservazioni 5.2.1

- Prima, molti filesystem erano piatti: cioè, i file sono tutti sullo stesso livello, senza alcuna organizzazione.

- Altri consentivano l'uso di un solo livello di directory.
- Il primo filesystem gerarchico è stato UNIX.

5.2.2 Tipi

Definizione 5.2.2: Directory a 1 Livello

Contengono i file tutti nella stessa directory. Non si può quindi avere due file con lo stesso nome, la multiutenza diventa difficile.

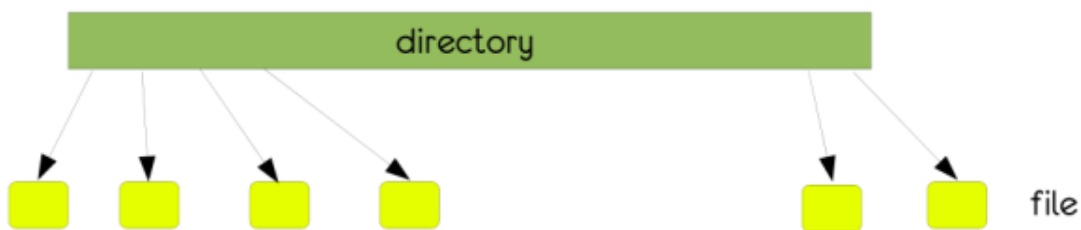


Figure 5.2: Directory a 1 livello.

Definizione 5.2.3: Directory a 2 Livelli

Ogni utente ha una propria cartella. Sistema il problema precedente, ma con molti file diventa di difficile uso.

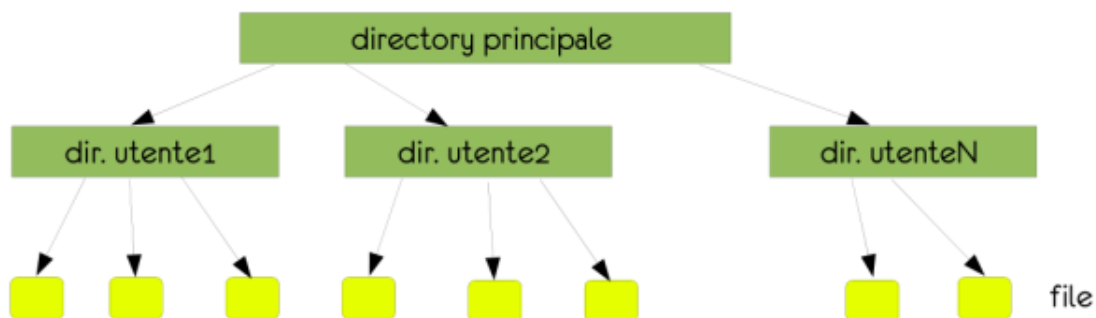


Figure 5.3: Directory a 2 livelli.

Definizione 5.2.4: Directory ad Albero

I file diventano foglie dell'albero, abbiamo infiniti livelli. Usare cammini assoluti diventa scomodo. Per questo l'utente viene posizionato virtualmente in una directory di lavoro.

Note:-

Ogni utente sarà proprietario di un sottoalbero, e un utente sarà proprietario della parte di sistema (amministratore).

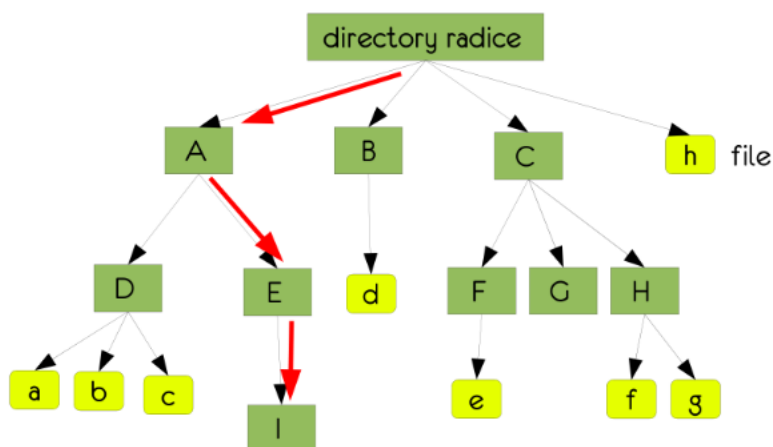


Figure 5.4: Directory ad albero.

Definizione 5.2.5: Directory con Grafo Aciclico

Un file/directory può essere accessibile da cammini diversi. Se si vuole eliminare un file che è puntato da due sottoalberi diversi, significa due cose:

- Il file non serve a nessuno dei due utenti: viene rimosso.
- Solo uno dei due utenti rinuncia ad usarlo: solo uno dei collegamenti viene rimosso.

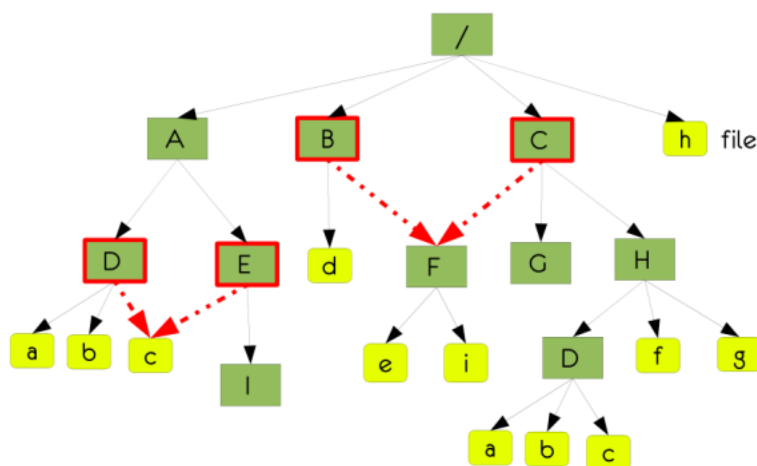


Figure 5.5: Directory a grafo.

Definizione 5.2.6: Link

I link sono riferimenti a file/directory.

Ci sono due possibili soluzioni:

- Soluzione demandata all'utente, l'utente ha due scelte:
 - *Link simbolico*: se il link creato viene eliminato, il file/directory a cui punta rimane intatto.
 - *Link fisico*: eliminare questo link vuol dire anche cancellare ciò a cui punta.

- Soluzione globale, mantenere il numero di riferimenti a quello che puntiamo:
 - Creare un link aumenta il valore.
 - Cancellare un link lo decrementa.
 - Quando il valore arriva a zero, il file/directory viene rimossa.

Definizione 5.2.7: Mount

Device possono essere aggiunti/rimossi; possiamo poi partizionare gli hdd, e montare suddette partizioni in punti diversi.

Note:-

Il mount su linux viene fatto nel file `/etc/fstab`.

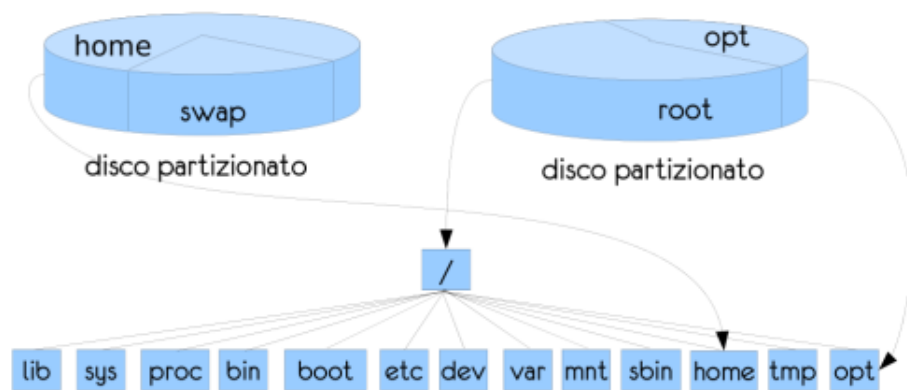


Figure 5.6: Mount.

Definizione 5.2.8: Filesystem Distribuiti

I file a cui si accede attraverso un FS possono anche risiedere su macchine diverse, collegate tra di loro in rete. Il modello client-server si basa sul concetto di una macchina chiamata server che contiene fisicamente i file condivisi. Le macchine che vogliono accedere ai file sono dette client.

Note:-

Questi FS distribuiti si possono montare esattamente come le partizioni locali.

Per proteggere i file da accessi impropri, si possono usare due soluzioni:

- *Access Control List* (ACL):
 - A ogni nodo del FS viene associata una lista che specifica gli utenti che possono accedere al FILE.
 - Questa soluzione però complica l'implementazione delle directory, perchè dovrebbe mantenere una ACL, che possono essere lunghe.
- Proprietario, Gruppo, Altri:
 - Ogni nodo ha un proprietario.
 - Tutti gli utenti sono divisi in gruppi di lavoro.



Figure 5.7: implementazione di FS distribuito.

5.3 Aprire un File

5.3.1 Inode

Definizione 5.3.1: Inode

Un FCB (o inode) viene conservato in memoria secondaria, e mantiene le informazioni relative ai files. Hanno dimensione fissa.

Note:-

Quando un file viene aperto all'inizio, viene arricchito da alcune informazioni aggiuntive.

In linux:

- Inode copiato da disco.
- Stato dell'inode dell'inode:
 - Inode locked (file non disponibile).
 - La copia dell'inode in RAM è diversa da quella su disco.
 - Il file è un punto di mount.
- Device number: identificatore del FS a cui appartiene il file.
- Inode number: identificatore dell'inode nella struttura dati su disco.
- Contatore del numero di riferimenti all'inode.
- Contatore del numero di riferimenti all'inode (numero di utilizzi del file attuale).

Supponiamo di volere aprire un file. Dovremo fare le seguenti operazioni:

- Il kernel mantiene una inode table di dimensione finita.
- Nell'eseguire la open: controlla se l'inode corrispondente al file è già stato caricato, se sì userà l'inode trovato altrimenti.
- Se c'è spazio, crea un nuovo inode, lo "locka" e va a copiarvi l'inode su disco corrispondente al file da usare.
- Se non c'è spazio l'operazione fallisce e viene restituito un errore (così il processo richiedente non rischia di rimanere sospeso per tempi lunghi).

Domanda 5.1

Gli inode, però, sono immagazzinati in una sequenza di blocchi. Come facciamo a trovare un file, a partire da un lungo percorso assoluto?

```
inode-number NAMEI(string cammino)
    if (la prima directory del cammino è /)
        current = root inode;
    else
        current = inode della working directory;
    repeat
        el = leggi prossimo elemento da input;
        if (el == null)
            return current;
        else if (el è contenuto in current)
            current = inode associato a el;
        else
            return (no inode);
    until (el == null)
return current;
```

Figure 5.8: Algoritmo di NAMEI.

5.3.2 Apertura in UNIX

Per aprire un file in UNIX, un processo deve:

- Eseguire `open(pathname, flags)`.
- Il SO verifica se il file è in uso da qualche processo.
- Se no:
 - Utilizza l'algoritmo namei per trovare il numero di inode del file.
 - Calcola l'indirizzo su disco che permette di accedere fisicamente all'inode.
 - Invia al controller del disco il comando di lettura, che risulterà nella copiatura dell'inode in una entry della tabella in RAM (in-core inode).
 - Aggiorna la tabella di sistema.
- Se sì:
 - Identifica l'in-core inode e lo rende accessibile al processo modificando la tabella di sistema.

Note:-

La system call `open` restituisce come handle del file un *file descriptor*, un numero intero. Ogni volta che viene chiamata la `open`, viene aggiunto il file nella tabella dei file, una tabella globale con i riferimenti a tutti i file aperti.

5.4 Allocazione

5.4.1 Implementazione delle Directory

Esistono possibili alternative:

- Lista lineare (UNIX): una directory è una sequenza di coppie `<nomeFile, inodeNumber>`.
- B-tree.
- Tabella hash.

Limiti della lista lineare:

- Per verificare se un file è in una directory, bisogna scorrerla interamente.
- La ricerca di tipo lineare è lenta.
- È difficile mantenere la lista ordinata senza appesantire la gestione.
- Sono necessarie strutture di appoggio e algoritmi per gestire i buchi creati dalla cancellazione dei file.

Definizione 5.4.1: B-tree

Un B-tree è un albero ordinato: ogni nodo può contenere da N a $2N$ elementi detti chiavi. Un nodo avrà poi K elementi, con $K+1$ puntatori a nodi del livello successivo. Il numero N è anche detto ordine dell'albero. Questo, è un albero di ordine 2.

Note:-

Per ricercare in un B-tree si parte sempre dalla radice. Il suo ordinamento consente tempi rapidi di ricerca. Il suo bilanciamento e fan-out (apertura dell'albero) sono caratteristiche sfruttate per velocizzare la ricerca.

5.4.2 Spazio Disco ai Files

Ci sono tre metodi per organizzare la memorizzazione dei dati su disco:

- *Contigua.*
- *Concatenata.*
- *Indicizzata.*

Definizione 5.4.2: Allocazione Contigua

Ogni file è allocato in una sequenza contigua di blocchi. In genere utilizzata negli HDD a disco magnetico.

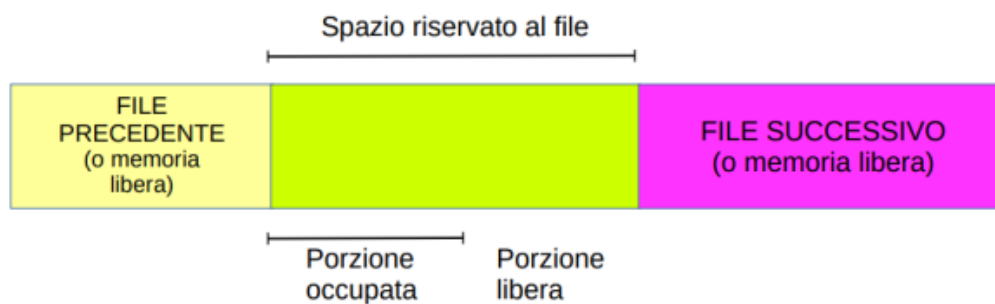


Figure 5.9: Allocazione contigua.

Vantaggi dell'allocazione contigua:

- Rapidità di accesso ai file: empo di seek (posizionamento della testina sulla traccia giusta) trascurabile.
- Quando si accede a un file si tiene traccia dell'ultimo blocco letto, quindi se abbiamo appena letto il blocco B , sia l'accesso sequenziale (al blocco $B+1$) sia l'accesso diretto (al blocco $B+k$) sono immediati.

Svantaggi della ricerca sequenziale:

- Difficile gestire la quantità di memoria da riservare per i file (può essere troppa o troppo poca, il file può crescere).
- Frammentazione esterna, c'è bisogno di deframmentazione ogni tanto.
- Bisogna gestire i buchi di memoria.

Definizione 5.4.3: Allocazione Concatenata

Consiste nello spezzare il file in parti che possono essere allocate in modo non contiguo (blocchi di dati). Ogni blocco contiene un puntatore a quello successivo.

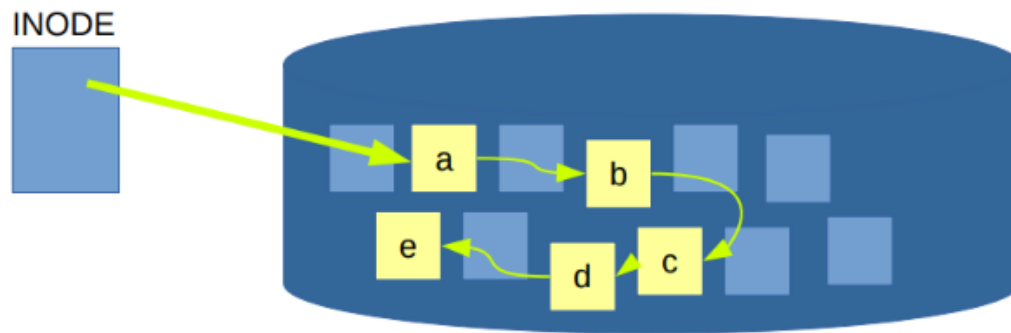


Figure 5.10: Allocazione concatenata.

Vantaggi dell'allocazione concatenata:

- Non è necessario preallocare memoria per i files.
- La lista concatenata è dinamica.
- Non è necessario deframmentare.

Svantaggi:

- Solo l'accesso sequenziale è efficiente, accesso diretto e indicizzato devono comunque scorrere la lista dei blocchi.
- I puntatori ai blocchi sono sparsi per il disco, e ogni salto ha latenza (spreco di tempo).
- Occorre spazio per mantenere i puntatori ai blocchi successivi.
- Se si corrompe un puntatore il resto del file va perso.

Corollario 5.4.1 Allocazione Concatenata FAT

Si riserva una sezione della partizione per mantenere una tabella che ha tanti elementi quanti blocchi. Se un blocco fa parte di un file, il contenuto della entry corrispondente in tabella è un riferimento al blocco successivo.

Note:-

Usata nei sistemi MS-DOS e successori.

Definizione 5.4.4: Allocazione Indicizzata

Questo tipo di allocazione risolve i problemi delle precedenti soluzioni introducendo un blocco indice. Posseduto da ogni file, il blocco indice è un array degli indirizzi dei blocchi che costituiscono il file. I riferimenti ai blocchi indice dei file sono mantenuti nelle directory.

Per l'accesso:

- Tramite la directory recupero il blocco indice.
- Tramite i riferimenti contenuti nel blocco indice posso accedere ai vari blocchi dati.

Osservazioni 5.4.1

- Più adeguata ad accessi diretti.
- L'allocazione indicizzata richiede di mantenere in RAM una parte dei blocchi indice, possono occorrere due o più accessi al disco se la RAM non è sufficiente:
 - Uno (o più) per accedere al blocco indice giusto.
 - Uno per raggiungere il dato di interesse.
- Alcuni sistemi combinano allocazione contigua e indicizzata: finché il file rimane di piccole dimensioni si usa l'allocazione contigua, oltre un certo limite si comincia ad usare un indice.

5.5 Gestione e Implementazione

5.5.1 Gestione dello Spazio Libero

Definizione 5.5.1: Gestione dello Spazio Libero

Per tenere traccia dei blocchi liberi, si utilizza una struttura basata su un array di bit, ognuno dei quali corrisponde a un blocco. Se il blocco è libero, il bit è impostato a 1

Tecniche di utilizzo:

- *Lista concatenata*: i blocchi liberi sono concatenati in una lista. Poiché i blocchi vengono allocati ai file uno per volta, basta pescare dalla testa della lista il primo blocco libero ed aggiornare il puntatore.
- *Raggruppamento*: concatenazione di blocchi indice. Si utilizza un blocco libero per mantenere N-1 puntatori ad altrettanti blocchi liberi, e l'ultimo puntatore per un eventuale ulteriore blocco simile.
- *Conteggio*: variante del precedente, in presenza di sequenze di blocchi liberi contigui si mantiene un riferimento al primo di tali blocchi e il numero di blocchi liberi ad esso consecutivi.

Quantità massima di memoria gestibile:

- Dipende dalle strutture del file system.
- Se devo gestire un disco di dimensioni superiori alla quantità di mem. massima gestibile, bisogna partizionare il disco in file system diversi.

5.5.2 Implementazione del Filesystem

Programmi e dati sono conservati in memoria secondaria. la memoria secondaria ha le seguenti caratteristiche:

- È organizzata in blocchi (un blocco può comprendere più settori).
- È possibile accedere direttamente a qualsiasi blocco.
- È possibile leggere un blocco, modificarlo e riscriverlo esattamente nella stessa posizione di memoria.
- Si accede alla memoria secondaria attraverso un FS.

Il FS è strutturato in una gerarchia di livelli:

- driver di dispositivo: si occupa del trasferimento dei dati da dispositivo di memoria secondaria a RAM e viceversa. È il gestore del dispositivo.
- FS di base: è proposto a passare comandi al driver di dispositivo.
- modulo di organizzazione dei file: è a conoscenza di come i file sono memorizzati su disco, è in grado di tradurre indirizzi logici in indirizzi fisici. Mantiene e gestisce anche l'informazione relativa ai blocchi liberi.
- FA logico: gestisce il FS a livello di metadati. Ogni file è rappresentato da un File Control Block (FCB).

Definizione 5.5.2: Struttura del Disco

Un disco può essere diviso in più partizioni. Ogni partizione può avere il suo FS, ed i diversi FS sono composti in una sola struttura.

Ogni FS non è altro che una sequenza di blocchi di memoria secondaria. Esistono poi cosiddetti blocchi speciali:

- *Boot Control Block*: blocco che in presenza di un SO contiene informazioni necessarie per la fase di bootstrap.
- *Volume Control Block*: descrive lo stato del FS, grandezza e altre informazioni (il numero di blocchi liberi, usati e la loro dimensione).
- Struttura delle directory: organizza i file, implementata in modo diverso, e contiene le coppie <nome file, FCB>

