

Long-term Information Storage

0. Must store large amounts of data
1. Information stored must survive the termination of the process using it
2. Multiple processes must be able to access the information concurrently

1

L'unità *convenzionale* di conservazione dei dati è il *file*

- insieme di dati
- insieme di attributi

Il concetto di "file" varia da S.O. a S.O.

Alcuni S.O. non hanno il concetto di file (allo scopo di meglio condividere informazioni fra le applicazioni)

Es. prototipi di MacOS (mai realizzati)

Es. PalmOS

2

Il concetto di file per MS-DOS

dati:
una sequenza di byte

attributi:

- *nome*
- bit read-only?
- bit hidden?
- bit system?
- bit archived?
- timestamps (creazione, ultima modifica, ultimo accesso)
- lunghezza

Nota: non c'è il proprietario (MS-DOS non ha il concetto di "utente")

3

Il concetto di file per Unix

dati:
una sequenza di byte

attributi:

- UID, GID (utente e gruppo proprietari del file)
- mode (permessi, es. 644 *rw-r--r--*)
- timestamps (creazione, ultima modifica, ultimo accesso)
- lunghezza

Nota: non c'è il nome; per Unix un file può avere tanti nomi diversi

4

Il concetto di file per NTFS

dati:
una o più sequenze di byte

- attributi:
- proprietario
 - ACL (Access Control List)
 - timestamps
 - lunghezza

5

Il concetto di file per MacOS

dati:
una sequenza di byte (data fork)

risorse:
un piccolo database strutturato (resource fork)

- icone
- font
- ... dati arbitrari

attributi:

- applicazione che lo gestisce
- *tipo* del file (es. TEXT)

Nota: non si basa sul *nome* del file per sapere che tipo di file è

Posso avere file di testo gestiti da una app, e altri file di testo gestiti da un'altra app *sulla stessa macchina*

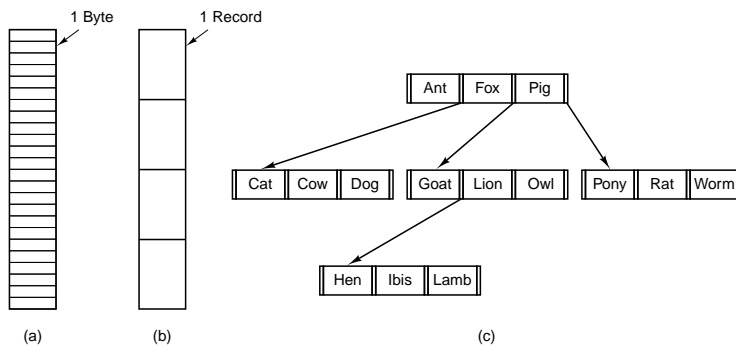
(impossibile in Windows! che usa l'associazione estensione–applicazione)

Posso associare a un semplice file di testo l'informazione di quale posizione stavo leggendo quando l'ho chiuso

Quando trasferisco un file da MacOS a un altro SO perdo la resource fork

6

Varie maniere di strutturare la parte "dati"



- (a) Unix, MS-DOS, NTFS, MacOS
- (b) VMS, AS/400, mainframe
- (c) AS/400, mainframe

7

Seq. di record versus seq. di byte

Vantaggi del modello Unix (seq. di byte)

- più semplice
- più flessibile

Vantaggi del modello a record

- tutte le app usano lo stesso meccanismo per conservare i dati
- molte app che potrebbero richiedere un DBMS possono farne a meno

Il modello a record è tuttora popolare nei mainframe

Il modello di Unix:

- la semplicità è la proprietà più importante
- la struttura a record viene vista dall'applicazione, non dal SO
- le funzioni di DB vengono fornite da librerie (es. Berkeley DB)

8

Il continuo FS—DB

Compiti di un File System

- gestire l'insieme dei blocchi di una partizione
- mantenere informazioni persistenti
- gestire accesso concorrente

Compiti di un Database Management System

- mantenere informazioni persistenti
- gestire accesso concorrente
- indicizzare le informazioni

⇒ c'è una grande sovrapposizione di compiti

Alcuni DBMS incorporano funzioni di FS

- es. Oracle lavora su partizioni nude

Alcuni FS incorporano funzioni di DBMS

- es. s/390 ha i B-tree nativi nel FS

9

Manipolazione di file in Unix

I file aperti vengono acceduti tramite file descriptor

file descriptor: intero non negativo (piccolo)

La shell di Unix mantiene la convenzione:

- fd 0 è lo *standard input*
- fd 1 è lo *standard output*
- fd 2 è lo *standard error*

il fd è usato come indice nella *file descriptor table* del processo

10

File Operations: aprire un file

Prima di usare un file occorre *aprirlo*

- input: file name
- input: flags
- output: file descriptor

```
int open(const char *pathname, int flags);
```

open(2) restituisce il più piccolo fd non usato

flags:

- O_RDONLY, O_WRONLY, O_RDWR
- O_APPEND
- O_CREAT
- O_TRUNC
- O_NONBLOCK
- ...

11

File Operations: creazione

esempio:

```
fd = open("pippo", O_CREAT | O_TRUNC | O_WRONLY, 0644);
```

Note:

- se O_CREAT è presente, si aggiunge un terzo argomento (*mode*)
- il mode in C è espresso come una costante ottale (inizia con 0!)

12

File operations: chiudere un file

```
int close(int fd);
```

Quando un processo termina, tutti i file aperti vengono chiusi dal kernel

13

File operations: accesso non sequenziale

Ogni file aperto ha un "current file offset"

Indica la posizione per la prossima operazione di lettura o scrittura

Può essere manipolato con la syscall `lseek(2)`

```
off_t lseek(int fd, off_t offset, int whence
```

I valori possibili per *whence* cambiano l'interpretazione di *offset*

- `SEEK_SET`: dall'inizio del file
- `SEEK_CUR`: dalla posizione corrente
- `SEEK_END`: dalla fine del file

Se ha successo restituisce il nuovo offset

Nota: *non tutti i file* supportano questa operazione

Se fallisce, restituisce -1

14

Esempio: testiamo standard input per vedere se supporta seek

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

int main() {
    if (-1 == lseek(0, 0, SEEK_CUR)) {
        printf("stdin non supporta seek\n");
    } else {
        printf("seek OK\n");
    }
    exit(0);
}
```

15

More seek fun

L'esecuzione di `lseek(2)` non causa mai I/O

Posso usare `lseek(2)` per spostare l'offset al di là della fine del file

Se a questo punto scrivo sul file, ottengo un file con un "buco" (hole)

I buchi non consumano blocchi sul disco

```
# ls -l file.hole
-rw-r--r--  1 matteo  users   1234774 Apr 14 21:55 file.hole
# du file.hole
4      file.hole
```

`ls(1)` riporta la lunghezza

`du(1)` riporta i blocchi occupati

16

File operations: leggere da un file

```
ssize_t read(int fd, void *buf, size_t nbytes);
```

Se ha successo, restituisce il numero di byte letti

Può essere meno di *nbytes*:

- in un file regolare, se raggiungo la fine del file
 - alla chiamata successiva, restituisce 0
- leggendo dal terminale, ottengo una riga alla volta
- leggendo dalla rete, ottengo solo i byte immediatamente disponibili

Se fallisce, restituisce -1

17

File operations: scrivere su un file

```
ssize_t write(int fd, void *buf, size_t nbytes);
```

Se ha successo, restituisce il numero di byte scritti

Di solito è uguale a *nbytes* a meno che

- esaurisco lo spazio su disco (errore)
- esaurisco la dimensione massima per file per processo (errore)
- scrivendo sulla rete, il processo ricevente è lento a leggere (non è un errore)

Se fallisce, restituisce -1

18

Esempio: copiare stdin su stdout

```
#define BUF_SIZE 8192
int main() {
    int n;
    char buf[BUF_SIZE];

    while ((n = read(0, buf, BUF_SIZE)) > 0) {
        if (n != write(1, buf, n)) {
            perror("write error");
            exit(EXIT_FAILURE);
        }
    }
    if (n < 0) {
        perror("read error");
        exit(EXIT_FAILURE);
    }
    exit(0);
}
```

19

Operazioni atomiche (i)

Per appendere dati in fondo a un file potrei usare

```
lseek(fd, 0, SEEK_END);
write(fd, buf, nbyter);
```

Ma questo codice ha una race condition

Per appendere dati in modo garantito atomico apro il file con `O_APPEND`

20

Operazioni atomiche (ii)

Due processi potrebbero sincronizzarsi in questo modo: entrambi cercano di creare un file; solo uno avrà successo. Il file funge da mutex.

```
open("/tmp/mymutex", O_CREAT | O_EXCL, 0600);
```

21

Testo di riferimento

Per imparare a usare *bene* Unix da *programmatore*:

W. Richard Stevens,
Advanced Programming in the Unix Environment
Addison-Wesley 1993

22

Strutture dati nel kernel: i processi e i file

La `task_struct` per ciascun processo

- contiene una file descriptor table
- ogni riga della fd table contiene
 - file descriptor flags
 - puntatore a un *file object*

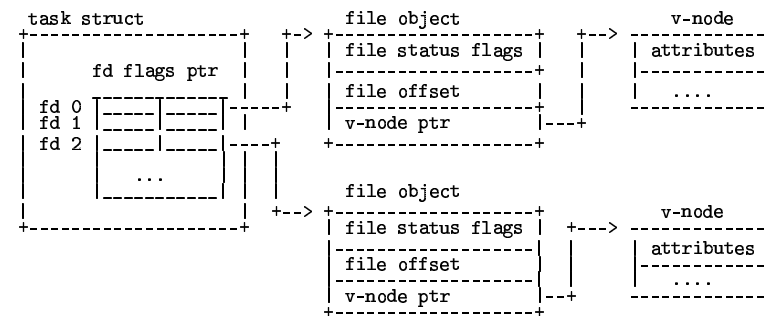
Il `file object` contiene

- file status flags (`O_RDONLY`, `O_WRONLY`, `O_NOBLOCK`,...)
- current file offset
- puntatore a un *v-node*

Il `v-node` è associato univocamente a un file (regolare o meno)

- attributi del file
- device
- lista dei blocchi

23



24

La struttura a tre livelli

Esiste fin dalle primissime versioni di Unix

Serve a consentire la condivisione dei file

0. fd table (una per processo)
1. file object (uno per ogni esecuzione di open(2))
2. v-node (uno per file)

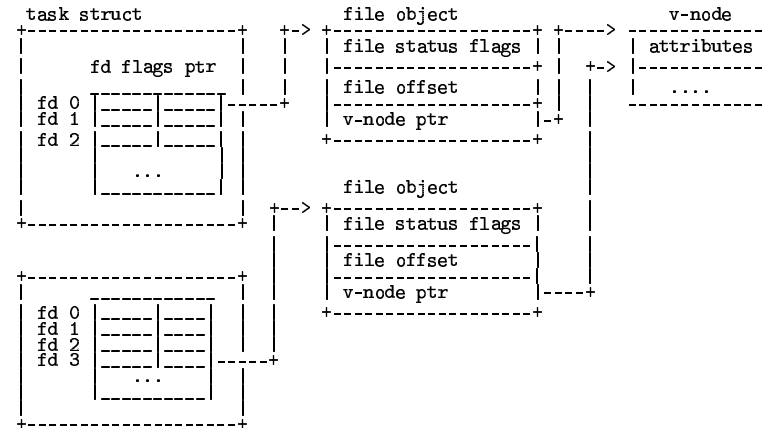
Operazioni compiute dalla open(2):

- percorre il pathname per trovare il file
- crea il v-node (se non è già aperto)
- crea il file object
- aggiorna la fd table

25

Condivisione, primo caso

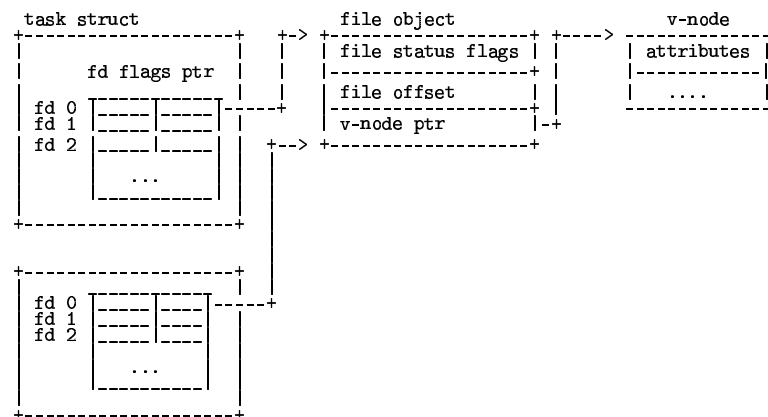
- due processi indipendenti aprono lo stesso file
- devono avere current file offset separati



26

Condivisione di file, secondo caso

Dopo una fork(2), i due processi hanno fd table identiche



27

Esempio: la redirectione della shell (i)

```
$ command > file
```

Ridirigi fd 1 sullo *file*

La shell lo implementa con

```
if (fork() > 0) {
    // siamo nel figlio
    close(1);

    // ora la open(2) restituirà 1
    open("file", O_CREAT | O_TRUNC | O_WRONLY, 0644);

    execve("command", ...);
} else {
    ...
}
```

28

Ora si spiega...

... perché in Unix usiamo fork + exec

Fra la fork e la exec, la shell può manipolare la fd table

29

Esempio: gli script di shell

Un altro caso in cui due processi condividono un file è dopo una fork(2)

Es. uno script di shell "prova.sh"

```
/bin/echo ciao
/bin/echo a tutti
```

eseguo

```
$ sh prova.sh > foo
```

Cosa contiene "foo" alla fine?

- a. ciao\na tutti
- b. a tutti
- c. (niente)

30

... il file descriptor 1 non viene mai chiuso durante l'esecuzione di questo script ... quindi l'offset cresce sempre!

31

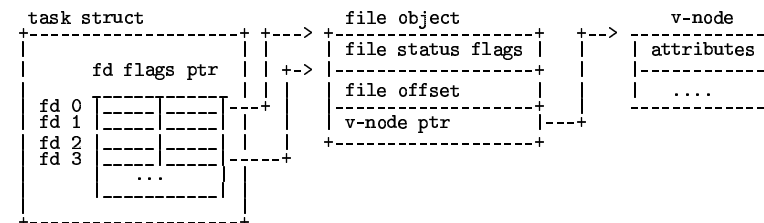
Condivisione di file, terzo caso

```
int dup(int fd);
```

"Duplica" il file descriptor *fd*

Se ha successo restituisce un nuovo file descriptor che è il più piccolo fra quelli non aperti

Il nuovo fd punta a una copia della riga del *fd* originale



32

Esempio: la redirectione della shell (ii)

```
$ command > file 2>&1
```

Ridirigi fd 2 sullo stesso file di fd 1

La shell la implementa con

```
if (fork() > 0) {  
    // siamo nel figlio  
    close(1);  
  
    // ora la open(2) restituirà 1  
    open("file", O_CREAT | O_TRUNC | O_WRONLY, 0644);  
  
    close(2);  
  
    // la dup restituirà 2  
    dup(1);  
  
    execve("command", ...);  
}
```

33

Esercizio...

Perché questi due frammenti di shell si comportano in modo diverso?

```
$ a.out > outfile 2>&1
```

```
$ a.out 2>&1 > outfile
```

Suggerimento: la shell interpreta i suoi argomenti da sinistra a destra

34

Strutture dati del kernel associate a un processo

- file descriptor table
- (ormai la conosciamo!)

- current working directory
 - serve a interpretare i pathname parziali

- current root
 - serve a “confinare” un processo in un sottoalbero del FS
 - si usa chroot(2)

35