

# 10

## Gestiunea fisierelor

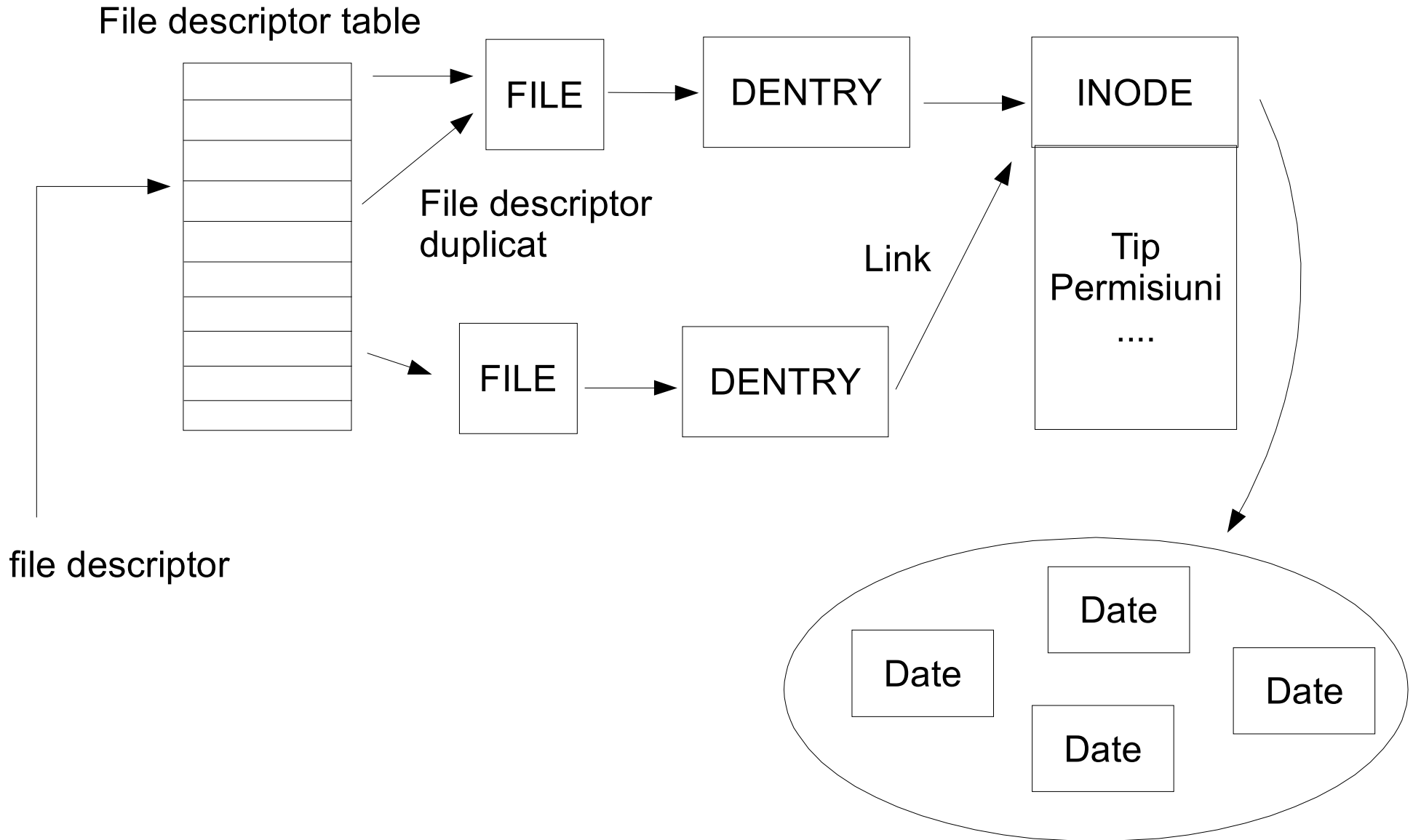
7 mai 2009

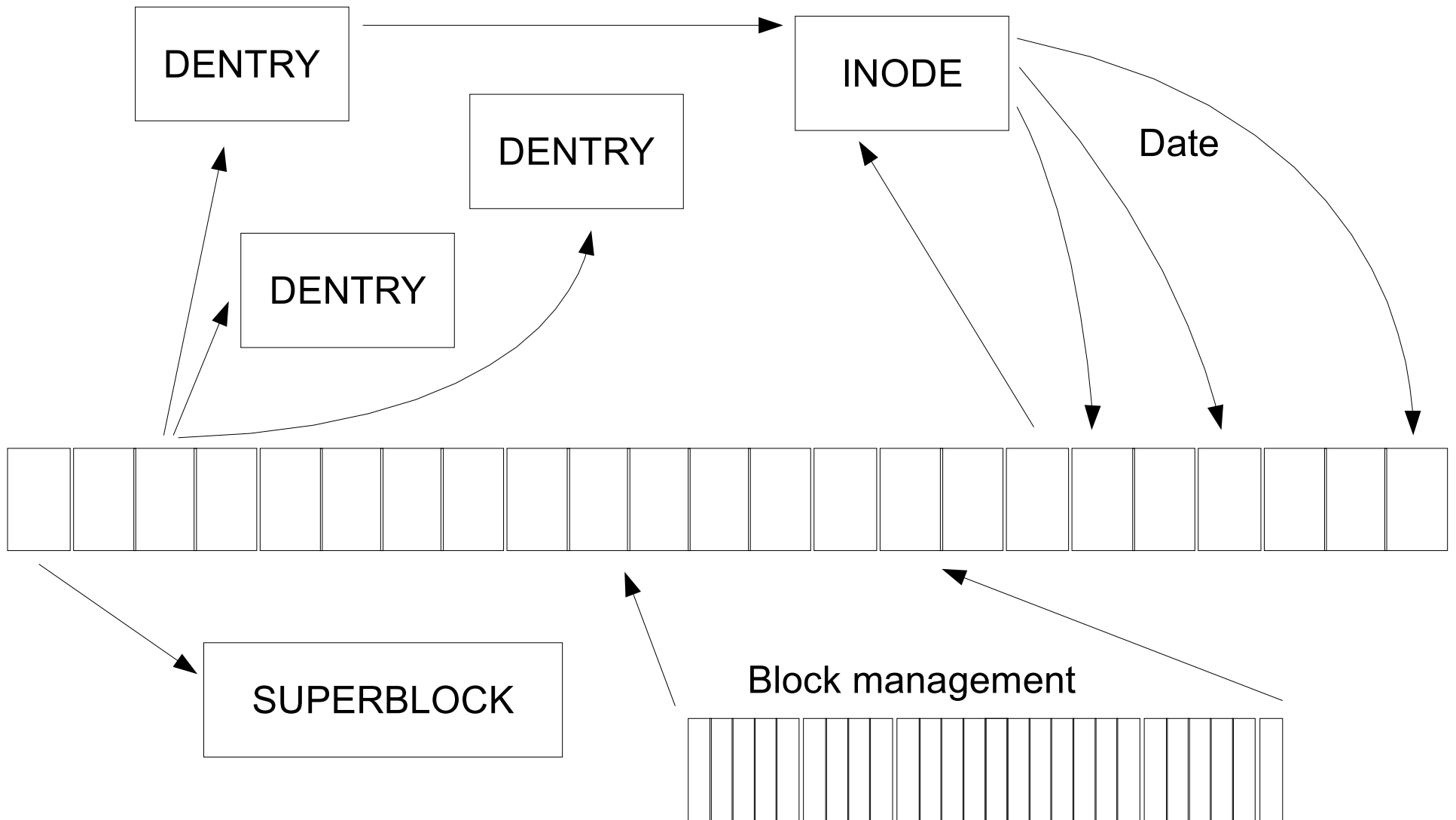
- Pe un sistem Linux cu split 3G/1G, la ce adresă fizică este mapată adresa virtuală 0xC0001000? Dar 0xF8001000?
- Care este complexitatea (e.g.  $O(1)$ ,  $O(n)$ , etc.) alocării unei pagini? Dar a unei structuri `task_struct`?
- Pentru folosirea unei zone de memorie dinamic alocată, în context întrerupere, ce primitive de alocare pot fi folosite în Linux? Dar în Windows?

- Gestiunea fișierelor – noțiuni generale
- Gestiunea fișierelor în Linux
  - VFS
  - Inode cache, dcache, page cache
  - Buffer cache
- Gestiunea sistemelor de fișiere în Windows
  - Cache manager
  - Fast I/O

- LKD: capitolele 12, 15
- WI: capitolele 11, 12

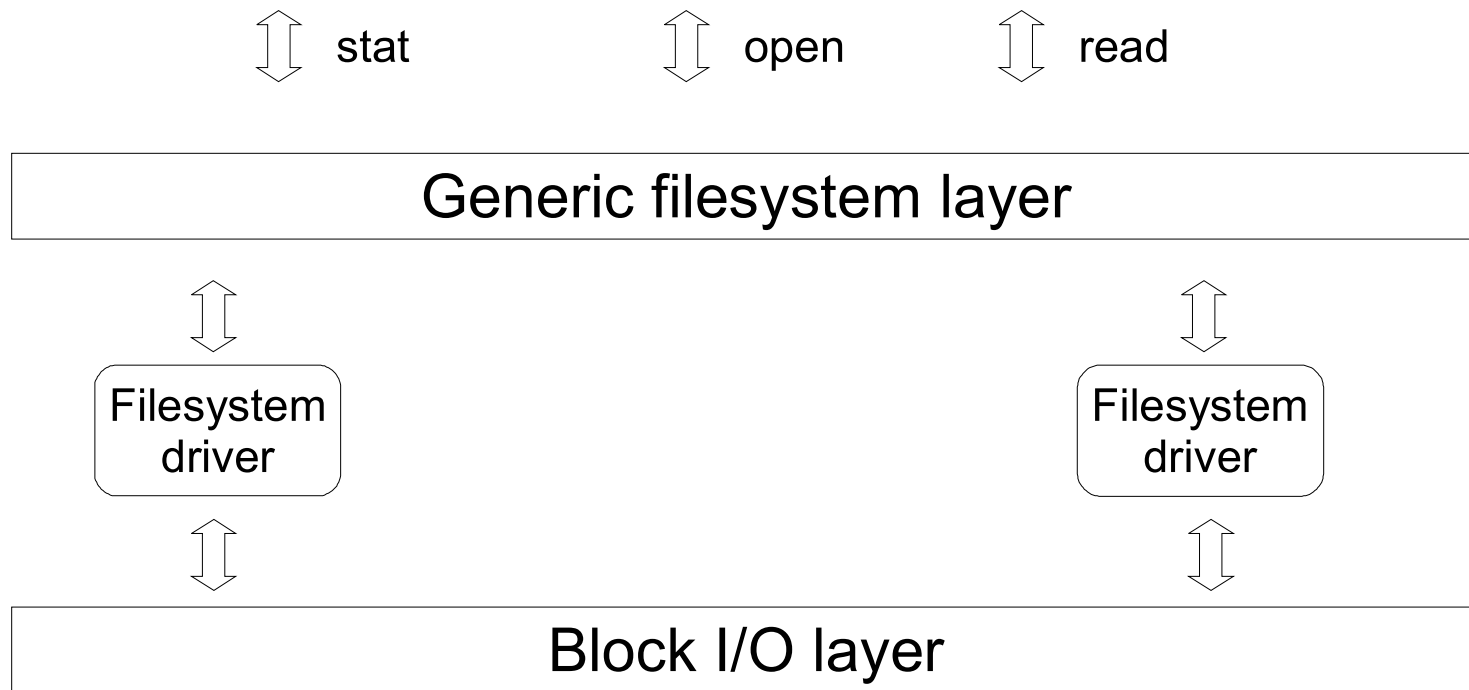
- Superblock – conține informații despre sistemul de fișiere (dimensiunea blocului, inode-ul rădăcină); există atât pe disc cât și în memorie
- File – structura internă SO ce descrie un fișier deschis; există doar în memorie
- Inode (FCB - file control block) – entitatea ce identifică un fișier în mod unic; există atât pe disc cât și în memorie
- Dentry- asociază un nume cu un fișier; există atât pe disc cât și în memorie





Superblock	IMAP	DMAP	IZONE	DZONE
------------	------	------	-------	-------





- Montarea
- Deschiderea unui fișier
- Determinarea atributelor fișierului
- Citirea de date din fișier
- Scrierea de date în fișier
- Închiderea unui fișier
- Crearea unui fișier
- Ștergerea unui fișier

- (sau înregistrarea lui în cadrul sistemului)
- Intrare: un disk (partiție)
- ieșire: un DENTRY către directorul rădăcină
- Operații: verificare partiție, determinare diverși parametri, determinare inode-ului rădăcină
- Exemplu PITIX:
  - Se verifică MAGIC-ul
  - Se determină dimensiunea blocului
  - Se citește inode-ul rădăcină și se crează dentry-ul asociat

- Intrare: o cale
- Ieșire: un file descriptor
- Operații:
  - Determinarea sistemului de fișiere
  - Pentru fiecare nume din cale, determinarea inode-ului asociat cu numele
  - Odată găsit inode-ul ultimului nume din cale se crează structura FILE și se alocă o intrare în tabela de descriptori de fișier

- Intrare: file descriptor
- ieșire: attributele fișierului
- Operații:
  - Se accesează inode-ul (file->dentry->inode)
  - Se citesc attributele din inode

- Intrare: file descriptor, offset, length
- ieșire: date
- Operații:
  - Se accesează inode-ul (file->dentry->inode)
  - Se determină blocurile de date
  - Se copiează datele de pe disk în memorie și sunt apoi trimise utilizatorului

- Intrare: file descriptor, offset, length, date
- Ieșire:
- Operații:
  - Se accează inode-ul
  - Se alocă unul sau mai multe noi blocuri pe disc – se caută blocuri libere, se marchează ca fiind ocupate
  - Se adaugă blocurile alocate la inode
  - Se copiează datele de la user în bufer interne și apoi se scriu pe disk

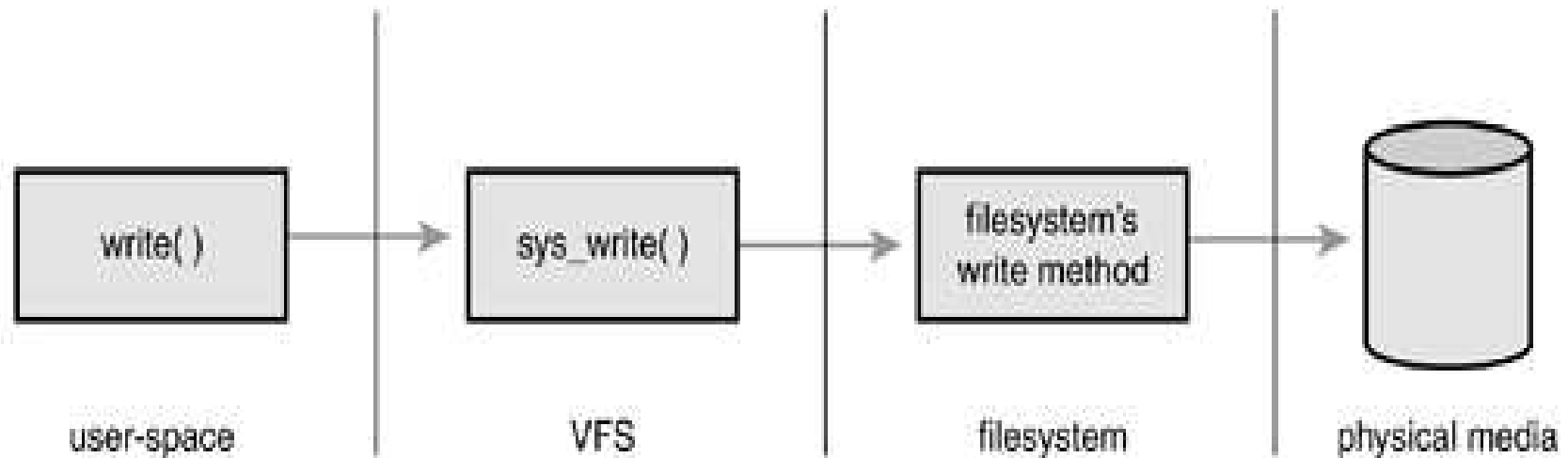
- Intrare: file descriptor
- ieșire:
- Operații:
  - Se decrementează reference counter-ul structurii FILE
  - Dacă referența counter-ului structurii FILE ajunge la 0, se șterge structura
  - Se setează pe NULL intrarea din tabela de descriptori de fișier



Directoarele sun fişiere ce conţin o înşiruire (vector, câteodată arbore) de DENTRY-uri.

- Intrare: o cale
- Ieșire:
- Operații:
  - Se determină inode-ul directorului în care trebuie adăugat fișierul
  - Se citesc blocurile de date
  - Se inserează un nou dentry
  - Se scriu pe disc blocurile de date modificate

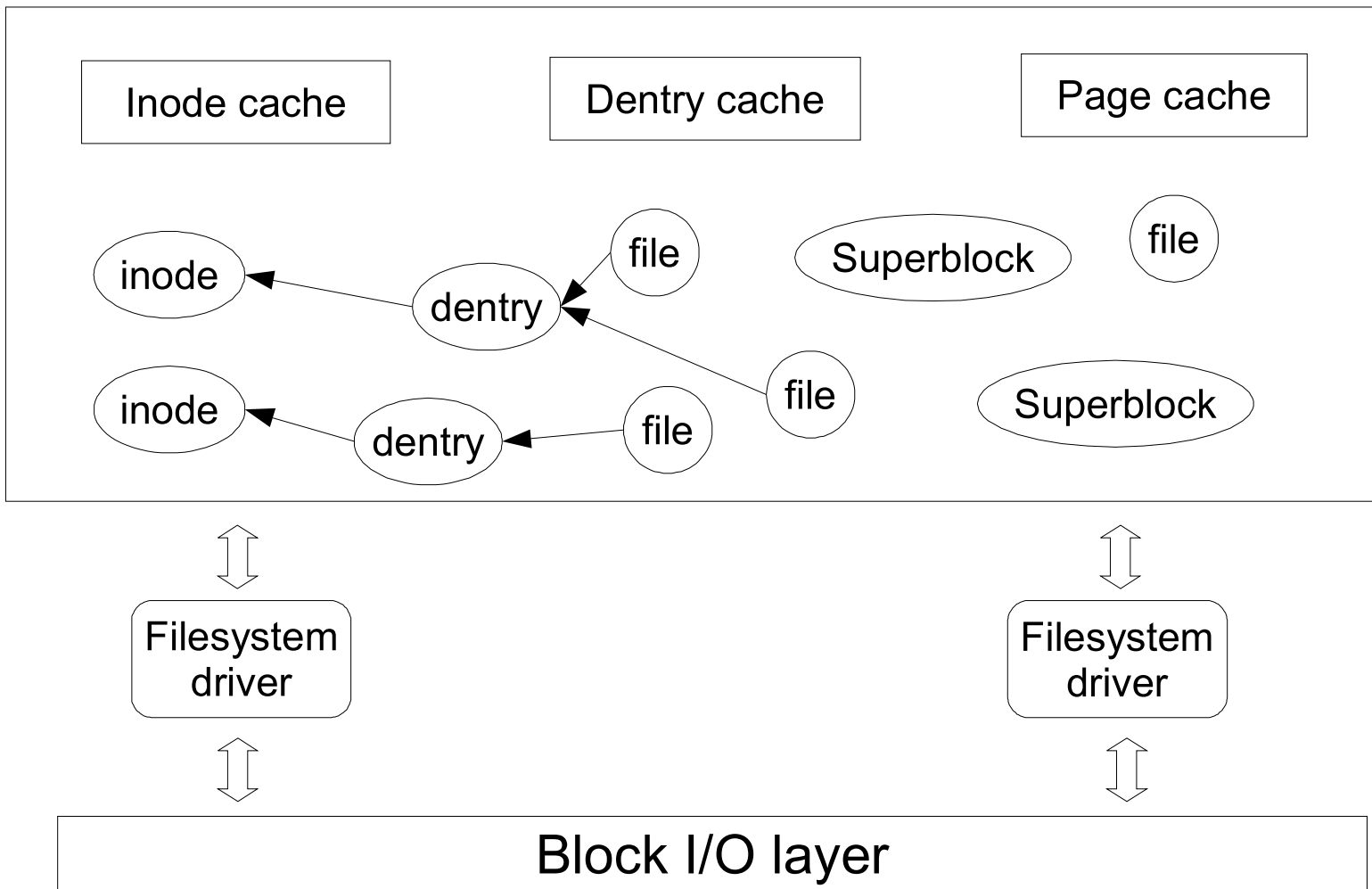
- Intrare: o cale
- Ieșire:
- Operații:
  - Se determină inode-ul directorului în care trebuie adăugat fișierul
  - Se citesc blocurile de date
  - Se caută și se șterge DENTRY-ul cu numele fișierului (caz special: link-uri)
  - Se scriu pe disc blocurile de date modificate



↕ stat

↕ open

↕ read



- fill\_super
- put\_super
- write\_super
- read\_inode
- write\_inode
- delete\_inode
- clear\_inode
- Statfs
- remount\_fs

- Create
- Lookup
- Link
- Unlink
- Symlink
- Mkdir
- rmdir
- Rename
- Readlink
- follow\_link
- put\_link
- Truncate
- ...

- Operațiile de citire a inode-urilor de pe disc sunt costisitoare
- Un inode citit se va menține în memorie (până când apar condiții de low memory)
- Căutarea inode-ului se face cu ajutorul unei tabele hash
- Funcția de hash (sb, inode\_number)



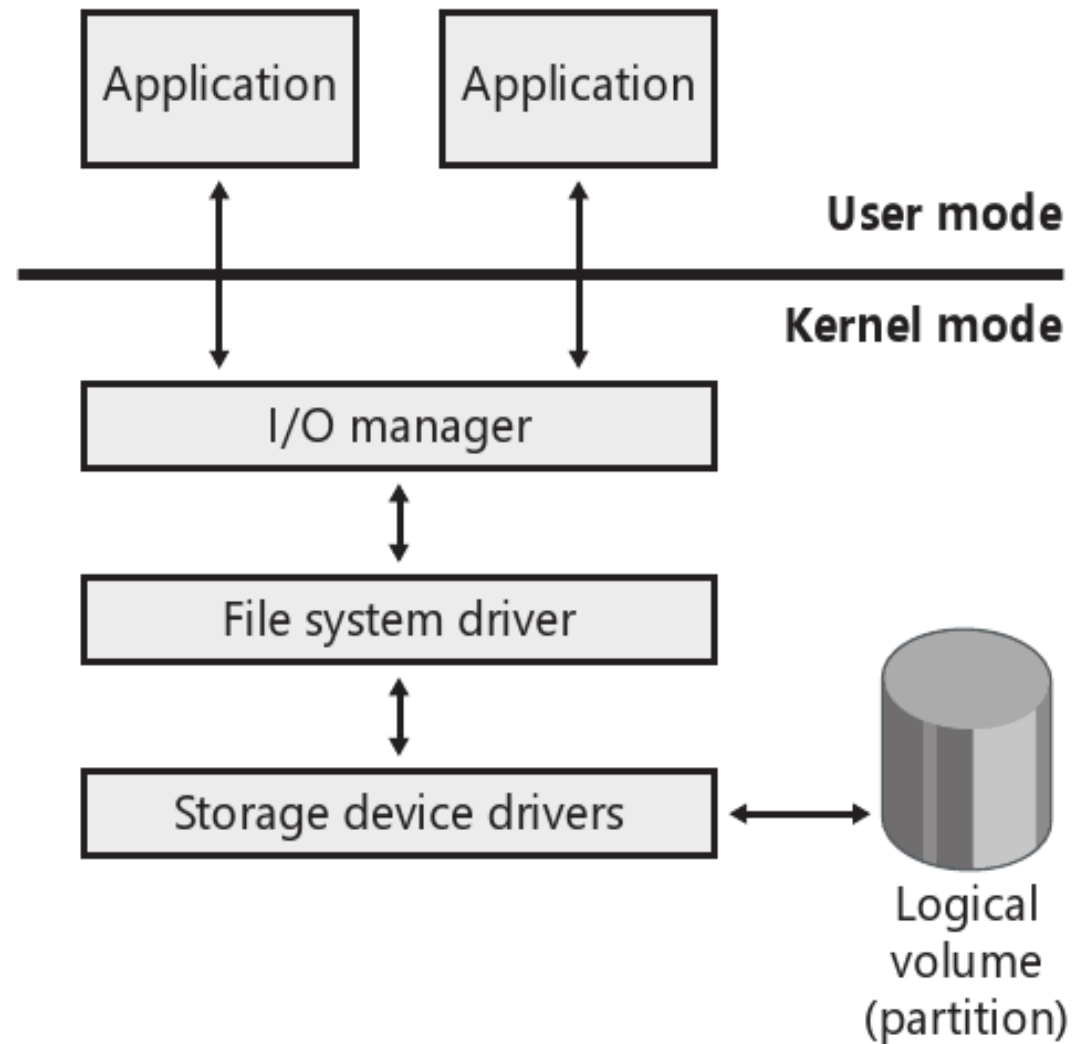
- Stări:
  - Used – d\_inode este valid și obiectul dentry este folosit
  - Unused – d\_inode este valid dar obiectul dentry nu este folosit
  - Negative – d\_inode nu este valid (NULL); fie inode-ul nu a fost încă încărcat fie a fost șters
- Dentry cache
  - Lista de dentry-uri folosite (dentry->d\_state == used)
  - Lista celor mai recent folosite dentry-uri (sortată după timpul de access)
  - Un hash table pentru a elimina parcurgerea arborelui

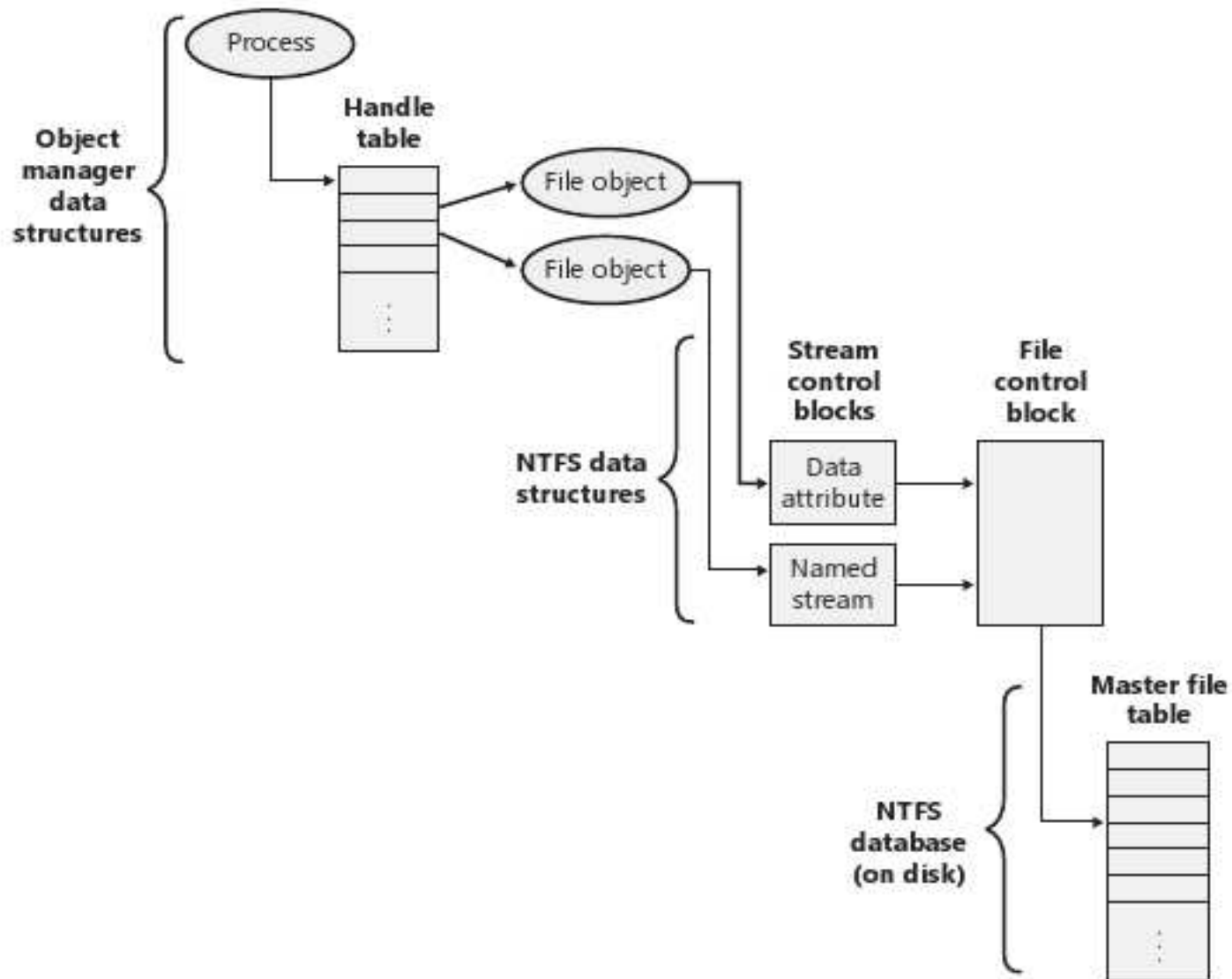
- Menține în memorie datele citite de pe disc
- Se folosește pentru:
  - read/write
  - Mmap
- Cache-ul se face la nivel de fișier și nu de block device -> este nevoie de o infrastructură pentru translatarea offset-urilor din fișiere în blocuri pe disc
- Cache-ul se bazează pe un radix tree

- Structura care este folosită de page-cache pentru a translata adrese
- Datorită acestei abordări se poate face caching nu doar pe mapări de fișiere
- Operații address\_space:
  - write\_page, read\_page
  - Bmap
  - direct\_IO
  - ...

- `generic_file_read` – implementare generică
  - Se verifică dacă pagina nu există în page cache
  - Dacă există se copiază datele din page cache în buffer-ul utilizatorului
  - Dacă nu există se apelază `a_ops->readpage()`
- `block_read_full_page` – implementare generică, folosește `a_ops->bmap()`

- De la versiunea 2.4 nu mai există un buffer cache separat
  - S-a eliminat problema double buffering
  - S-au eliminat problemele de sincronizare
- `buffer_head`-urile – structurile care fac maparea dintre un buffer în memorie și un block pe disc folosesc memorie alocată de page cache

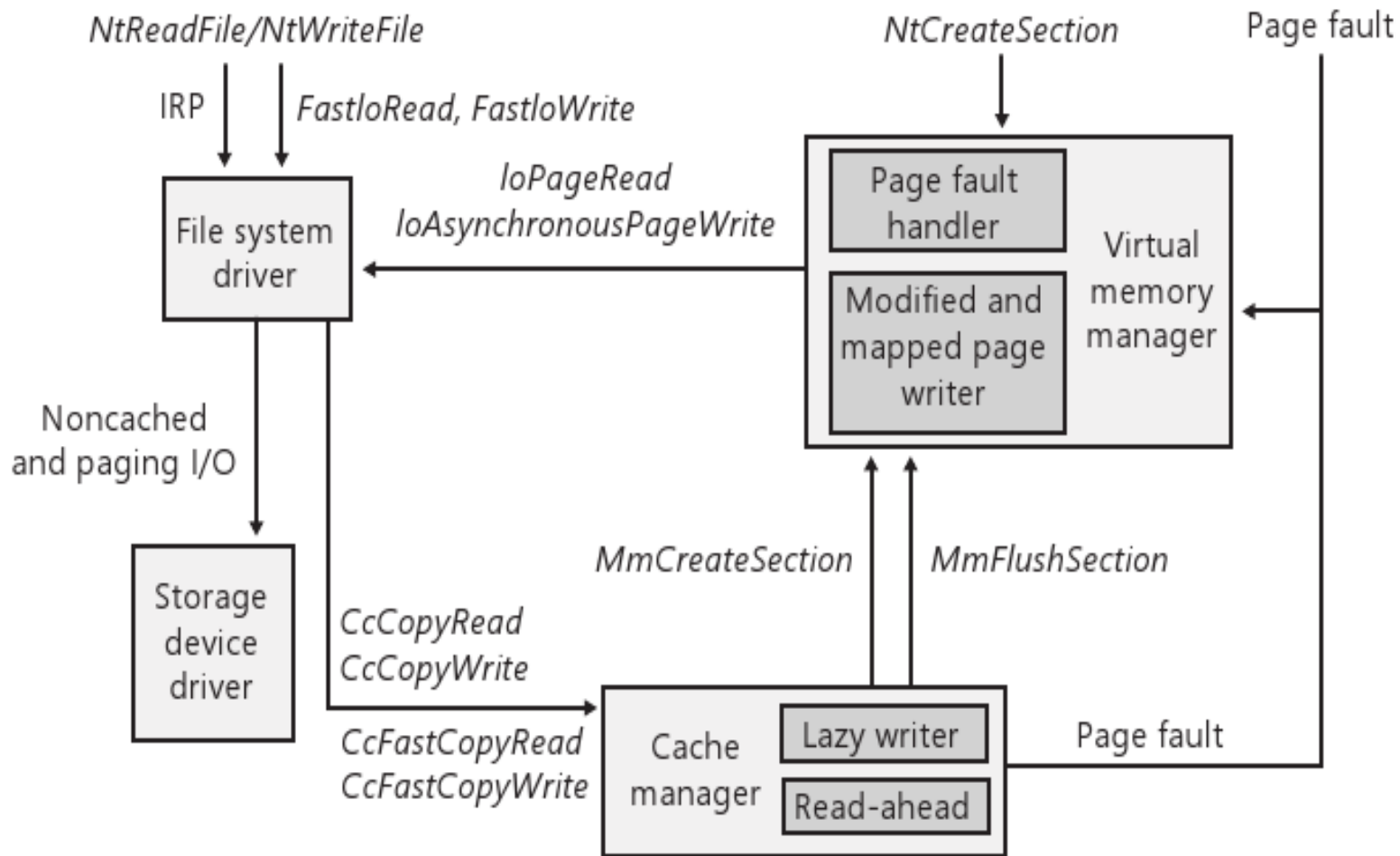




- Superblock – Volume Parameter Block
- Inode – File Control Block
- File – File Object
- Dentry -

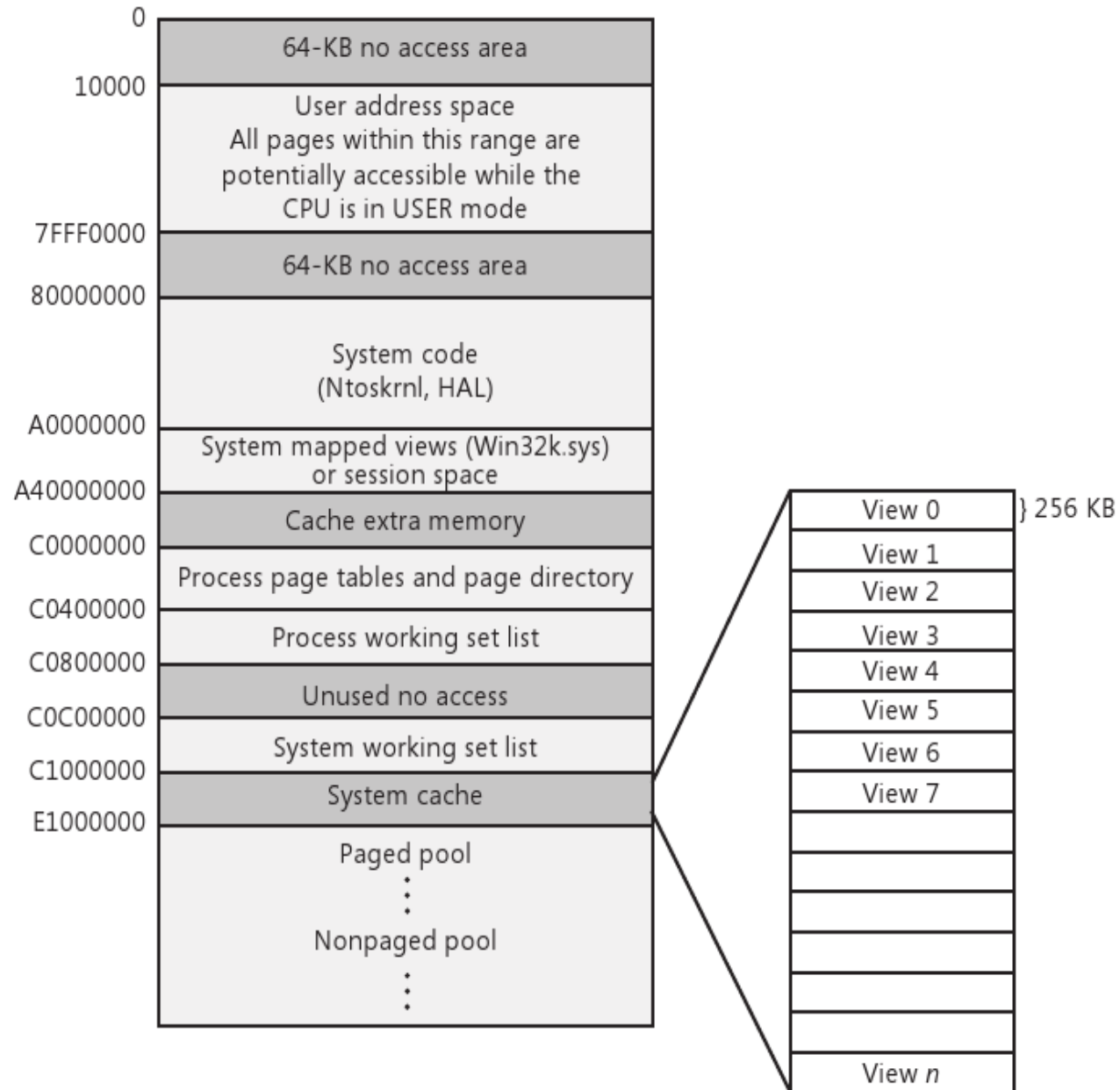


# Interacțiune FSD – MM – Cache manager

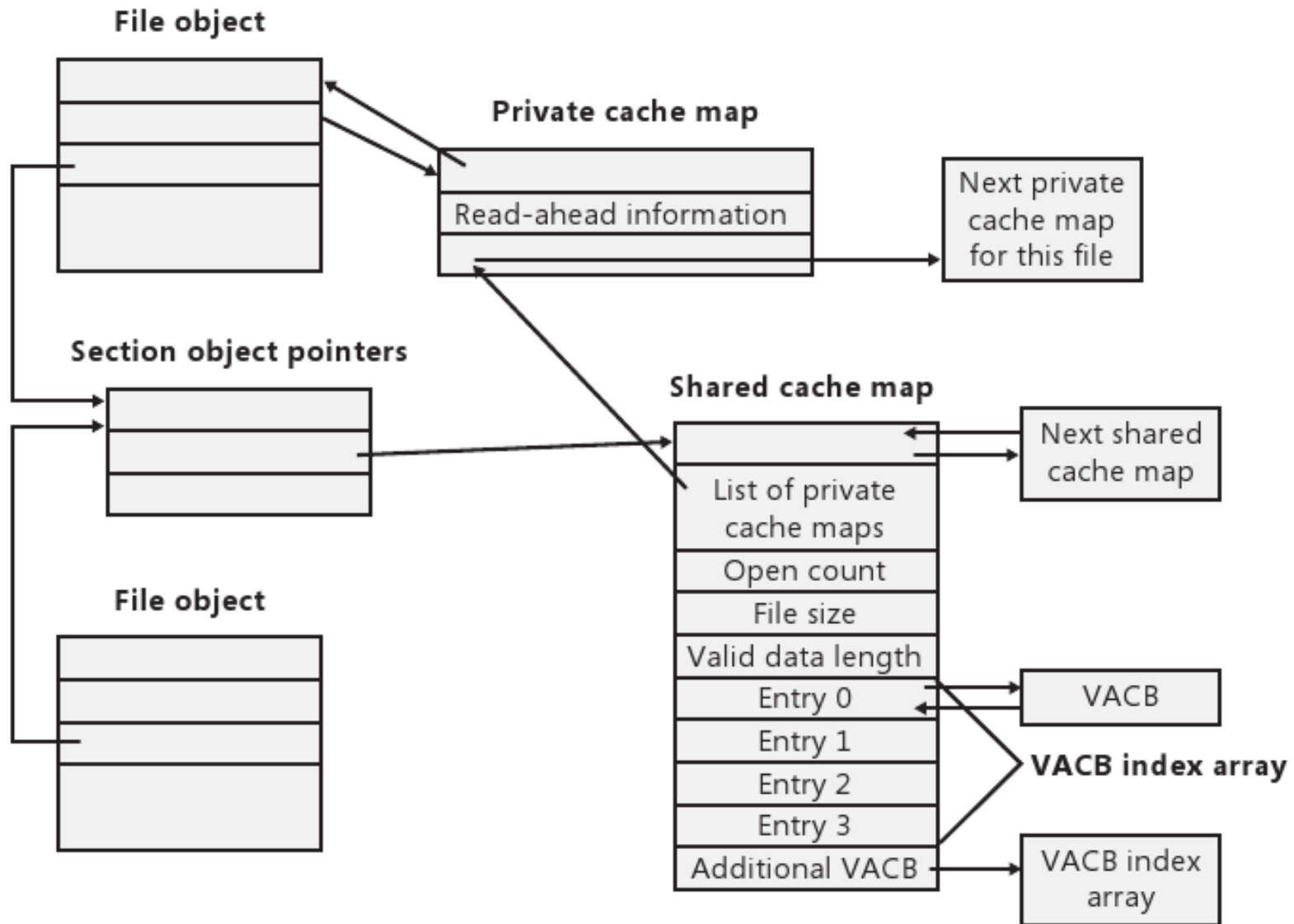


- Dacă `FILE_FLAG_NO_BUFFERING` nu a fost setat:
  - se copiează datele în buferul utilizatorului cu ajutorul cache managerului: `CcCopyRead`
- Altfel:
  - Se crează noi IRP-uri ce vor fi pasate storage device-ului

# System cache address space



# Maparea fișierelor în system cache



# Maparea fişierelor în system cache (2)

