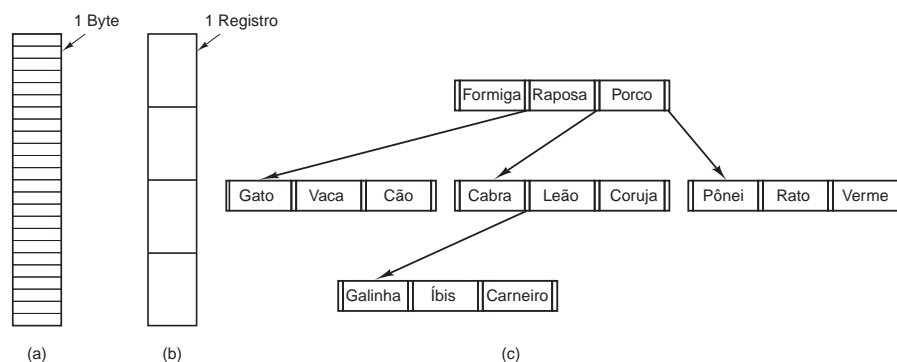


Nomeação de arquivos

Estrutura de arquivos (1/2)

Estrutura de arquivos (2/2)

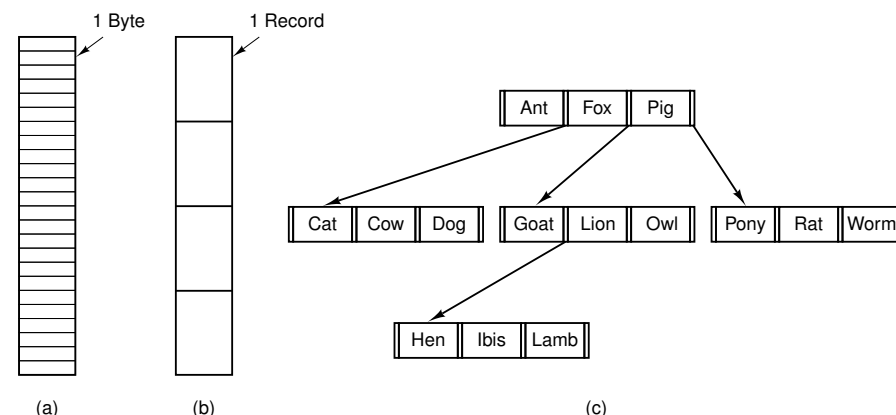


(a) Sequência de bytes (b) Registros (c) Árvore

Tipos de arquivos

- **Arquivos regulares:** contêm dados de usuários
- **Diretórios:** arquivos do sistema que contêm a estrutura do sistema de arquivos
- **Arquivos especiais de caracteres:** usados para acessar dispositivos orientados a caracter
- **Arquivos especiais de bloco:** usados para acessar dispositivos de E/S orientados a bloco
 - filosofia do UNIX: arquivos são a interface de acesso aos dispositivos de E/S para os usuários
- Variam de acordo com o SO

Estrutura de arquivos (2/2)



(a) Sequência de bytes (b) Registros (c) Árvore

Acesso a arquivos

- **Sequencial:** arquivo precisa ser lido/escrito em sequência
 - para ler o byte 1 milhão, é preciso ler os 999.999 anteriores
 - único modo de acesso compatível com fita magnética, p. ex.
- **Direto (aleatório):** bytes/registros podem ser acessados em qualquer ordem
 - a operação de E/S pode especificar a posição ou pode haver uma operação específica de posicionamento
 - modo mais comum de acesso a disco

Atributos de arquivos

Operações com arquivos

Atributo	Significado
Proteção	Quem pode ter acesso ao arquivo e de que maneira
Senha	Senha necessária para ter acesso ao arquivo
Criador	ID da pessoa que criou o arquivo
Proprietário	Atual proprietário
Flag de apenas para leitura	0 para leitura/escrita; 1 se apenas para leitura
Flag de oculto	0 para normal; 1 para não exibir nas listagens
Flag de sistema	0 para arquivos normais; 1 para arquivos do sistema
Flag de repositório (archive)	0 se foi feita cópia de segurança; 1 se precisar fazer cópia de segurança
Flag ASCII/binário	0 para arquivo ASCII; 1 para arquivo binário
Flag de acesso aleatório	0 se apenas para acesso sequencial; 1 para acesso aleatório
Flag de temporário	0 para normal; 1 para remover o arquivo na saída do processo
Flag de impedimento	0 para desimpedido; diferente de zero para impedido
Tamanho do registro	Número de bytes em um registro
Posição da chave	Deslocamento da chave dentro de cada registro
Tamanho da chave	Número de bytes no campo-chave
Momento da criação	Data e horário em que o arquivo foi criado
Momento do último acesso	Data e horário do último acesso ao arquivo
Momento da última mudança	Data e horário da última mudança ocorrida no arquivo
Tamanho atual	Número de bytes no arquivo
Tamanho máximo	Número de bytes que o arquivo pode vir a ter

- Create
- Delete
- Open
- Close
- Read
- Write
- Append
- Seek
- Get Attributes
- Set Attributes
- Rename

Arquivos mapeados em memória

Sumário

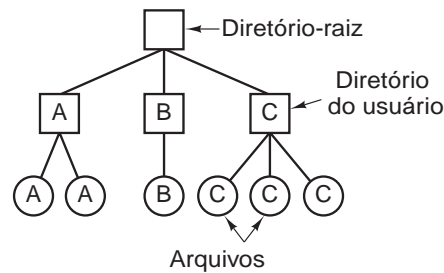
- Operações de acesso a arquivos são razoavelmente mais complicadas do que operações de acesso à memória
- Muitos sistemas suportam a noção de **arquivos mapeados em memória**
 - acessos a uma determinada região de memória são equivalentes a acessos ao arquivo
 - muitas vezes é possível mapear apenas partes de um arquivo
 - evita problemas com arquivos maiores que o espaço de endereçamento do processo
 - o que ocorre quando um processo A mapeia um arquivo em memória e outro processo B realiza uma leitura do modo tradicional?
 - modificações feitas por A só serão refletidas no arquivo quando a página correspondente for salva no disco

- Introdução
- Arquivos
- Diretórios
- Implementação de sistemas de arquivos
- Tópicos adicionais
- Sistemas de arquivos no Linux

Estruturas de diretórios

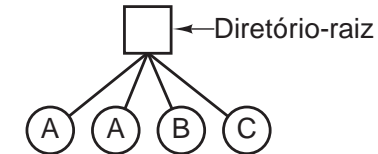
- Diretório em nível único
 - simplicidade
 - rapidez na busca
 - problemas na colisão de nomes
- Diretório em dois níveis
 - cada usuário possui seu diretório privado
- Diretórios hierárquicos
 - árvore de diretórios
 - facilita organização dos arquivos

Diretório em dois níveis

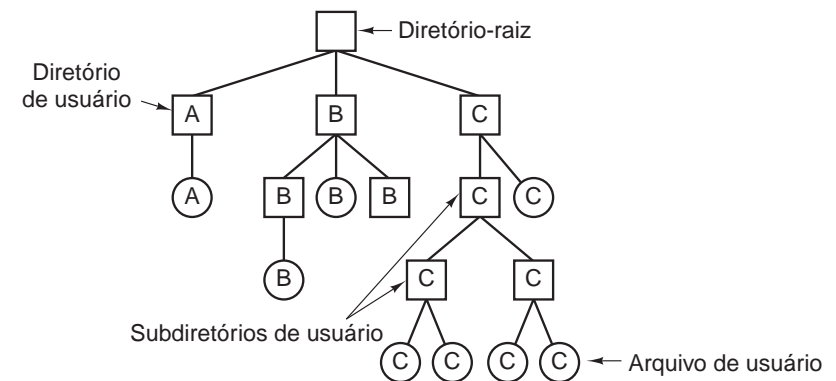


a letra indica o usuário dono do diretório/arquivo

Diretório em nível único



Diretórios hierárquicos



Lista encadeada com tabela na memória (2/2)

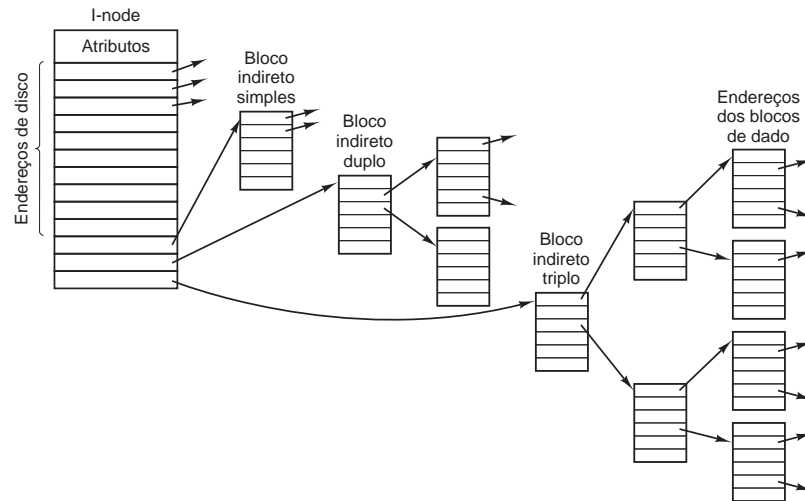
- | Bloco físico | | |
|--------------|----|---------------------------|
| 0 | | |
| 1 | | |
| 2 | 10 | |
| 3 | 11 | |
| 4 | 7 | ← O arquivo A começa aqui |
| 5 | | |
| 6 | 3 | ← O arquivo B começa aqui |
| 7 | 2 | |
| 8 | | |
| 9 | | |
| 10 | 12 | |
| 11 | 14 | |
| 12 | -1 | |
| 13 | | |
| 14 | -1 | |
| 15 | | ← Bloco sem uso |

arquivo B: blocos 6, 3, 11, 14

I-nodes (2/4)

-
- | Atributos do arquivo | |
|--------------------------------|---|
| Endereço do bloco 0 do disco | → |
| Endereço do bloco 1 do disco | → |
| Endereço do bloco 2 do disco | → |
| Endereço do bloco 3 do disco | → |
| Endereço do bloco 4 do disco | → |
| Endereço do bloco 5 do disco | → |
| Endereço do bloco 6 do disco | → |
| Endereço do bloco 7 do disco | → |
| Endereço do bloco de ponteiros | → |
- Bloco de disco contendo endereços de disco adicionais

I-nodes (3/4)



Um i-node do UNIX V7 (1979)

I-nodes (4/4)

- Se existem 10 blocos diretos em um i-node, os blocos de dados são de 512 bytes e o endereço do bloco tem 32 bits, qual o tamanho máximo de um arquivo usando a estrutura de i-nodes do slide anterior?

Cálculo

$$TamArq_{max} = blocos_{diretos} + blocos_{simples} + blocos_{duplo} + blocos_{triplo}$$

$$blocos_{diretos} = 10 \times 512 = 5120 \text{ bytes} = 5 \text{ KB}$$

$$blocos_{simples} = (512/4) \times 512 = 65.536 \text{ bytes} = 64 \text{ KB}$$

$$blocos_{duplo} = (512/4) \times (512/4) \times 512 = 8.388.608 \text{ bytes} = 8.192 \text{ KB}$$

$$blocos_{triplo} = (512/4) \times (512/4) \times (512/4) \times 512 = 1.048.576 \text{ KB}$$

$$\text{Total} = 1.056.837 \text{ KB} = 1,008 \text{ GB}$$

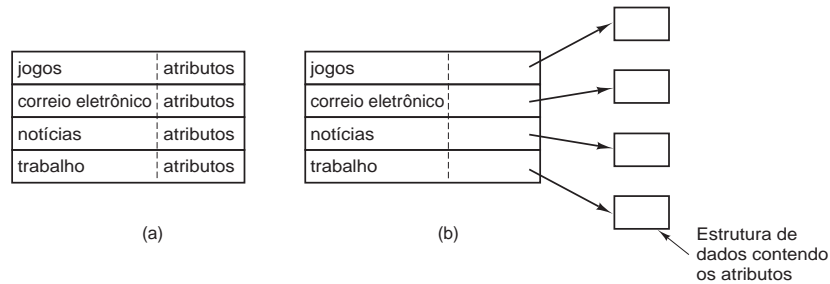
Implementação de diretórios (1/4)

- A entrada de um diretório fornece as informações para encontrar os blocos de disco correspondentes ao arquivo
 - endereço de todo o arquivo (alocação contígua)
 - número do primeiro bloco (lista encadeada)
 - número do i-node
- Função do diretório é mapear o nome do arquivo à informação necessária para localizá-lo

Implementação de diretórios (2/4)

- Projeto simples: lista de entradas de tamanho fixo contendo o nome do arquivo e seus atributos
 - uma entrada por arquivo
 - solução usada no MS-DOS
- Projeto com i-nodes: as entradas contêm o nome do arquivo e um ponteiro para o descritor com os atributos

Implementação de diretórios (3/4)



(a) Um diretório simples

- entradas de tamanho fixo
- endereços de disco e atributos na entrada de diretório

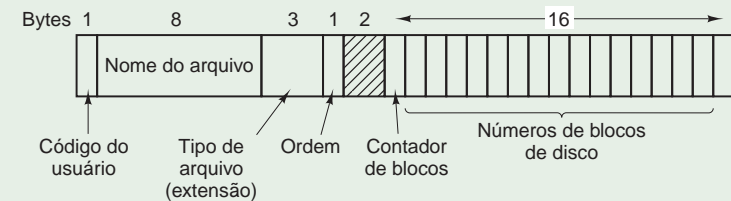
(b) Diretório no qual cada entrada se refere apenas a um i-nó

Gerenciamento de espaço em disco

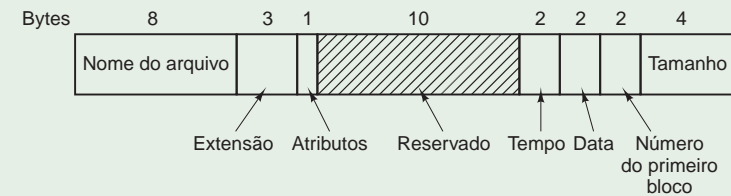
- Como alocar n bytes no disco?
 - usando uma área contígua de n bytes
 - dividindo os n bytes em blocos, não necessariamente contíguos
- Problema semelhante a segmentação vs paginação
 - soluções de compactação em disco são mais caras que a compactação de memória

Implementação de diretórios (4/4)

Implementação de diretórios no CP/M



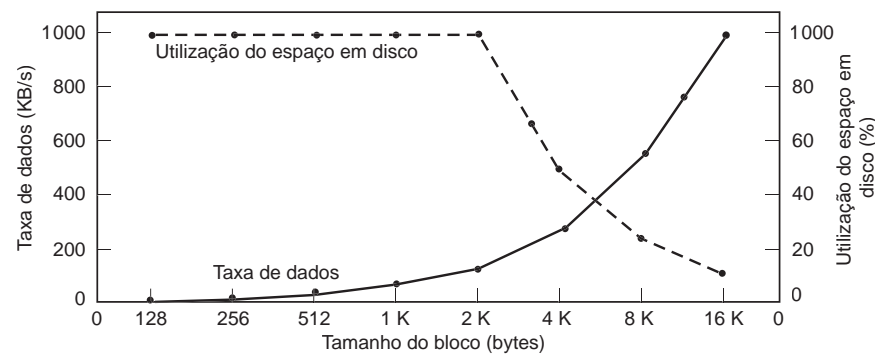
Implementação de diretórios MS/DOS (FAT)



Tamanho do bloco (1/2)

- Qual o tamanho de bloco mais adequado?
- Estudos mostram que o tamanho médio de um arquivo no UNIX é de 2 KB
- Blocos grandes têm melhor desempenho, mas desperdiçam mais espaço
- Blocos pequenos têm pior desempenho (muitos blocos a serem lidos/escritos), mas aproveitam melhor o espaço

Tamanho do bloco (2/2)

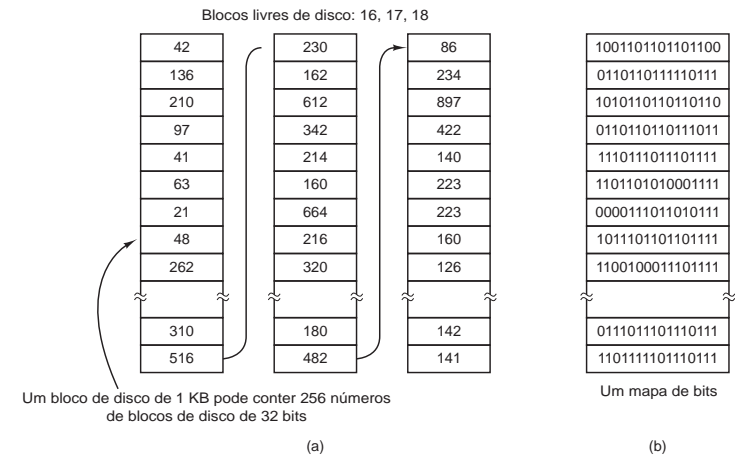


arquivos de 2 KB

Sumário

- 1 Introdução
- 2 Arquivos
- 3 Diretórios
- 4 Implementação de sistemas de arquivos
- 5 Tópicos adicionais
- 6 Sistemas de arquivos no Linux

Gerência de espaço livre



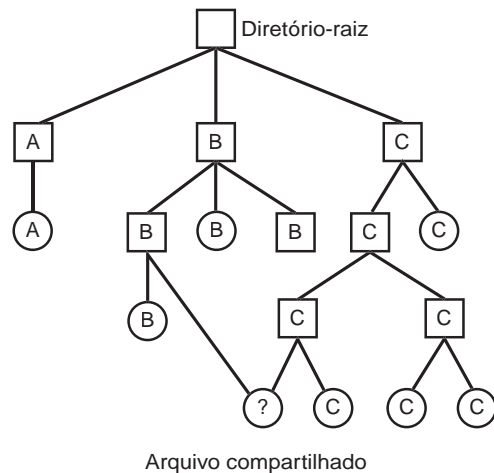
(a) lista encadeada

(b) mapa de bits

Arquivos compartilhados (1)

- Frequentemente é desejável que um mesmo arquivo esteja presente em vários diretórios diferentes
 - múltiplas cópias desperdiçam espaço em disco e introduzem problemas de consistência
- **Ligações (links)** permitem que um arquivo seja referenciado por múltiplos nomes
 - árvore de diretórios se torna um grafo acíclico dirigido

Arquivos compartilhados (2)

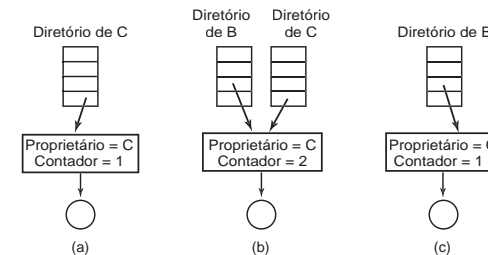


Ligações simbólicas

- Ligação simbólica: entrada no diretório aponta para o i-node de um arquivo do tipo ligação, que contém o nome de caminho do arquivo original
 - atalhos no Windows
- Pode ser usado com diretórios
- Problemas
 - overhead na tradução de nomes de caminho → acessos extras a disco
 - se o arquivo original for removido, as ligações se tornam inválidas
 - percurso recursivo de diretórios requer tratamento especial para evitar duplicação de conteúdo e laços
 - duplicação também se aplica a ligações estritas

Ligações estritas (*hard links*)

- Ligação estrita: entrada no diretório aponta para o mesmo i-node do arquivo original
 - i-node possui um contador de ligações
 - incrementado quando uma nova ligação é criada
 - decrementado quando uma ligação é removida
 - i-node e dados só são removidos quando o contador chega a zero
- Não pode ser usada com diretórios → ciclos no grafo
- Problemas na contabilização do espaço ocupado pelos usuários



Consistência do sistema de arquivos

- Escritas em disco não são operações atômicas, e falhas (falta de energia, p.ex.) podem deixar o SA em um estado inconsistente
 - se a falha for durante a escrita de metadados o problema pode ser agravado
- Geralmente há programas que analisam o SA para detectar e corrigir eventuais inconsistências
 - fsck (UNIX), Scandisk/CHKDSK (Windows)
- Esses programas se baseiam na redundância interna do SA para repará-lo

Funcionamento básico do fsck (2)

- Figure 1 illustrates the state of a memory pool with 16 blocks, indexed 0 to 15. The diagrams show the allocation and deallocation of blocks over time.

 - (a) Initial state: All 16 blocks (0-15) are free.
 - (b) After allocating 10 blocks: Blocks 0-9 are in use; blocks 10-15 are free.
 - (c) After allocating 2 more blocks: Blocks 0-12 are in use; blocks 13-15 are free.
 - (d) After deallocating 4 blocks: Blocks 0-11 are in use; blocks 12-15 are free.

Funcionamento básico do fsck (4)

- Verificação por arquivos
 - percorre recursivamente a árvore de diretórios, contando o número de referências a cada i-node
 - pode ser > 1 devido às ligações estritas
 - compara o número de referências observadas com a contagem de ligações presente em cada i-node
 - idealmente, são iguais
 - se a contagem do i-node for maior, o i-node e os blocos do arquivo não serão liberados depois que a última ligação for removida
 - se a contagem do i-node for menor, o i-node e os blocos serão liberados quando ainda estiverem em uso
 - em ambos os casos, a solução é atualizar a contagem do i-node com o número observado
- fsck ainda realiza outras verificações
 - compara informações resumidas no superbloco/i-nodes com resultados da análise estrutural
 - pode analisar permissões dos arquivos em busca de permissões sem sentido ou inseguras

Sistemas de arquivos *journaling* (2)

- Journaling exige que as operações registradas no log sejam **idempotentes**
 - podem ser repetidas várias vezes sem afetar o resultado final
 - `x=171` (idempotente) vs. `x++` (não idempotente)
- Exemplos de operações idempotentes
 - marcar o bloco *k* como livre no mapa de bits
 - remover a entrada `bar` no diretório `foo`
- Exemplos de operações não idempotentes
 - colocar o bloco *k* na lista de blocos livres (ele pode já estar lá)
 - decrementar o contador de ligações do *i*-node *i*
- Muitos sistemas de arquivos atuais usam journaling
 - `ext3/ext4`, ReiserFS, NTFS

Sistemas de arquivos virtuais (2)

-
- ```
graph TD
 subgraph User_Space [User Space]
 UP((User process))
 end
 subgraph Kernel_Space [Kernel Space]
 subgraph VFS_Layer [Virtual File System]
 VFS([Virtual file system])
 end
 subgraph FS_Layer [File Systems]
 FS1([FS 1])
 FS2([FS 2])
 FS3([FS 3])
 end
 BC[Buffer cache]
 end
 UP -- POSIX --> VFS
 VFS -- VFS interface --> FS1
 VFS -- VFS interface --> FS2
 VFS -- VFS interface --> FS3
 FS1 <--> BC
 FS2 <--> BC
 FS3 <--> BC
```

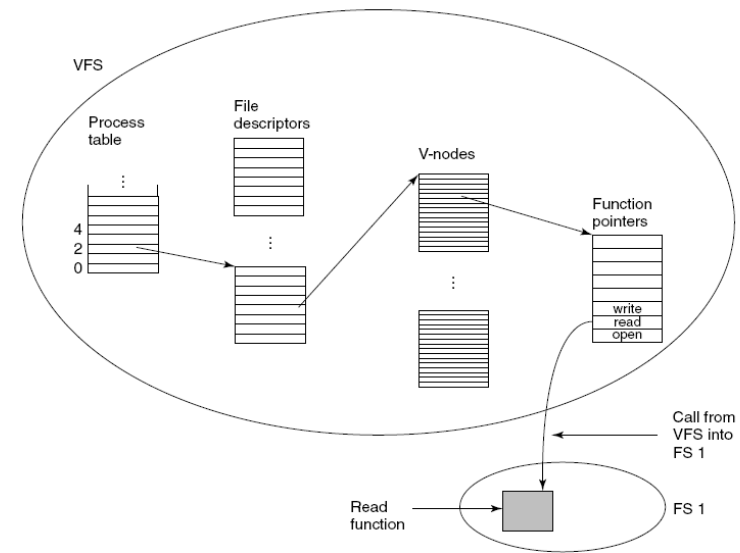
## Implementação do VFS (1)

- VFS mantém várias estruturas de dados
  - tabelas de descritores de arquivos em uso
    - para cada arquivo aberto é mantido um v-node em memória
  - v-nodes (equivalente a i-nodes)
    - contém dados que permitem acessar o SA subjacente
  - tabela de ponteiros de função para as operações de cada SA
    - funções que implementam `open()`, `read()`, `write()`, `close()`, `link()`, `unlink()`, etc.
    - projeto OO implementado em C

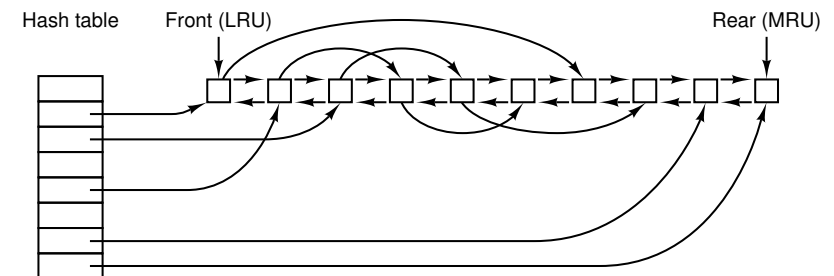
## Cache de blocos (1)

- Para reduzir o tempo das operações de disco, o SO mantém uma cache dos blocos mais acessados em memória
  - cache de blocos ou *buffer cache*
- Estrutura típica usa uma tabela hash com lista de colisão, indexada pelo dispositivo e pelo número do bloco
  - blocos são mantidos em uma lista duplamente encadeada ordenada pelo algoritmo MRU (menos recentemente usado)
    - LRU (*least recently used*)
    - bloco acessado por último vai para o final da fila
    - blocos acessados há mais tempo migram para o início da fila
    - quando for necessário substituir um bloco na cache (cache cheia), pega-se o primeiro da fila, que é gravado se foi modificado (→ sujo)

## Implementação do VFS (2)



## Cache de blocos (2)



## Cache de blocos (3)

- Heurísticas podem ser usadas para determinar se um bloco é inserido no início ou no final da fila
  - blocos com baixa probabilidade de reuso vão no início
    - são rapidamente escritos no disco e “recicladados”
  - blocos com maior probabilidade de reuso vão no final
    - blocos que estão sendo preenchidos, p.ex.
- Esperar até que um bloco de metadados chegue ao início da fila para ser gravado pode deixar o SA inconsistente
  - blocos de metadados são escritos rapidamente, em muitos casos de forma síncrona
- Podem ser usadas chamadas de sistema que forçam a escrita de todos os blocos sujos
  - `sync` (UNIX), `FlushFileBuffers` (Windows)
  - versões anteriores do Windows usavam cache de escrita direta (*write-through*)
    - modificações na cache eram imediatamente gravadas em disco
- Alguns sistemas integram cache de blocos e cache de páginas

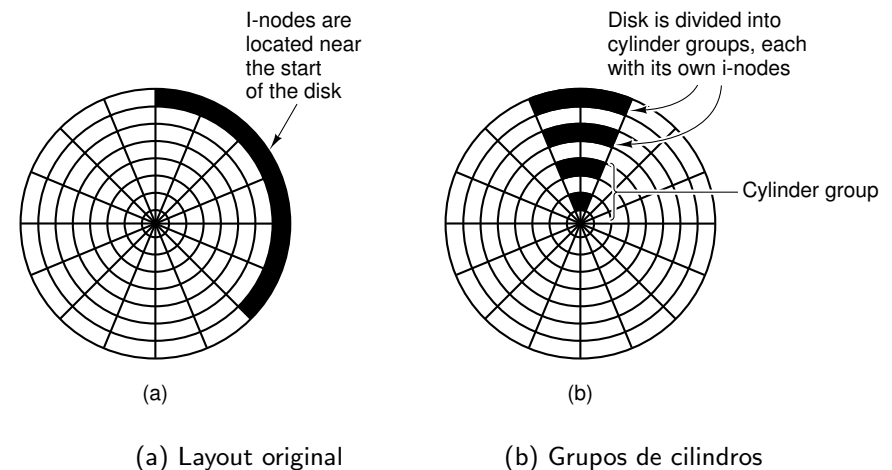
### Redução do movimento do braço do disco (1)

- Usar estratégias de alocação que minimizem deslocamentos (*seeks*)
- Originalmente, o sistema de arquivos do UNIX mantinha todos os i-nodes no início do disco, seguidos pelos blocos de dados
- Grupos de cilindros: espalha i-nodes pelo disco, mantendo-os próximos aos seus blocos de dados
  - cada grupo contém uma cópia do superbloco, i-nodes, um mapa de bits e blocos de dados
  - o SA tenta alocar todos os i-nodes de um diretório dentro do grupo
  - quando um bloco é requisitado para um i-node, o SA tenta alocar um bloco dentro do grupo
  - se não houver blocos disponíveis no mesmo grupo, recorre-se aos grupos vizinhos
  - blocos para arquivos grandes são alocados em diversos grupos, usando uma heurística

## Leitura antecipada de blocos

- Tenta trazer para a cache blocos que provavelmente serão necessárias em breve
- Exemplo: se uma aplicação está lendo um arquivo sequencialmente e solicita o bloco  $k$ , o SA já escalona a leitura do bloco  $k + 1$
- Em acesso aleatório, degrada o desempenho
- Pode ser usada uma heurística para analisar o comportamento da aplicação e decidir se vale a pena fazer leitura antecipada
  - inicialmente considera-se que o acesso é sequencial, e se faz leitura antecipada
  - se a aplicação reposiciona o ponteiro de arquivo, considera-se que o acesso é aleatório, e não há antecipação
  - se o acesso voltar a ser sequencial, pode-se voltar à leitura antecipada
  - tipo corrente de acesso pode ser indicado por um bit na entrada correspondente na tabela de arquivos abertos

## Redução do movimento do braço do disco (2)





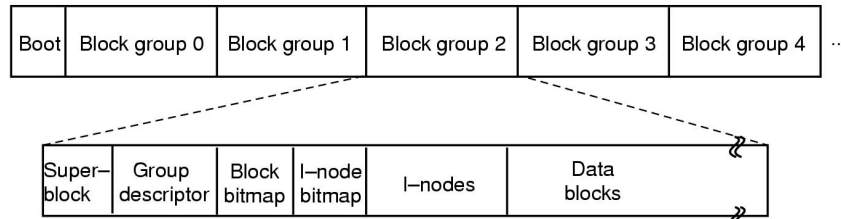
## Sumário

- 1 Introdução
- 2 Arquivos
- 3 Diretórios
- 4 Implementação de sistemas de arquivos
- 5 Tópicos adicionais
- 6 Sistemas de arquivos no Linux

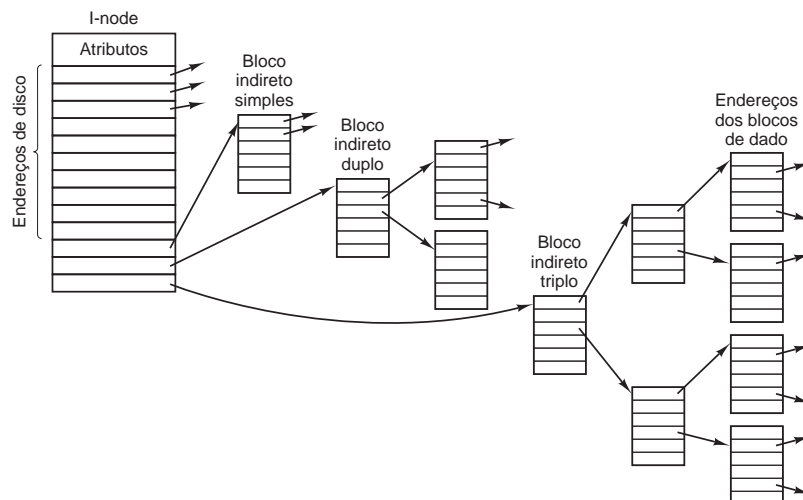
## Layout de disco no ext2 (1/2)

- O ext2 suporta partições de até 4 TB ( $2^{42}$  bytes)
- Cada partição é dividida em grupos de blocos
  - análogos aos grupos de cilindros do FFS (BSD UNIX)
  - no. de grupos depende do tamanho do SA
- Cada grupo de blocos é subdividido em
  - **superbloco**: contém os parâmetros e informações de controle do SA como um todo
    - cada grupo tem uma cópia do superbloco, que é único para o SA
  - **descriptor do grupo**: contém parâmetros e informações de controle do grupo
    - localização dos mapas de bits, no. inodos/blocos livres, no. diretórios
  - **mapas de bits**: controlam alocação de inodos e blocos de dados
    - cada mapa ocupa um bloco → limitador do tamanho do grupo
  - **inodos**: descritores de arquivo
  - **blocos de dados**

## Layout de disco no ext2 (2/2)



## Inodos no ext2 (2/2)



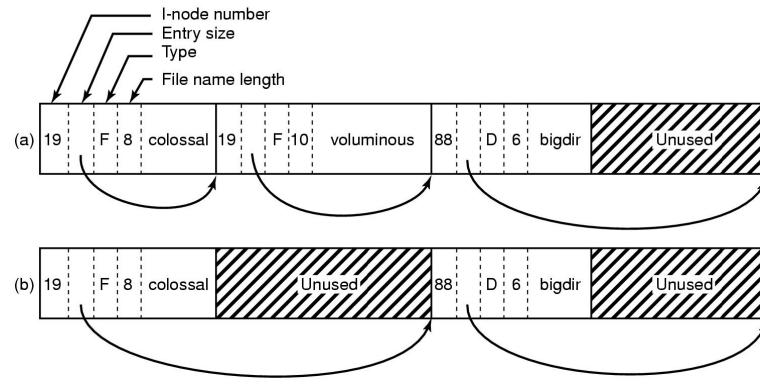
## Inodos no ext2 (1/2)

- O ext2 segue a estrutura clássica do UNIX
  - 12 ponteiros diretos
  - indireção simples/dupla/tripla
  - ponteiros de 32 bits (4 bytes)
  - blocos de 1, 2 ou 4 KB
    - o tamanho de bloco é definido na formatação do SA, por escolha do usuário ou em função do tamanho da partição
- cada inodo ocupa 128 bytes

## Diretórios no ext2 (1/2)

- Um diretório é uma lista encadeada de entradas de tamanho variável
- Cada entrada possui
  - inodo
  - tamanho da entrada
  - tamanho do nome
  - tipo de arquivo
  - nome
- O encadeamento é implementado usando o tamanho da entrada
  - quando uma entrada é removida, o tamanho da entrada anterior aumenta de modo a apontar para a próxima
    - o inodo também é zerado

## Diretórios no ext2 (2/2)



(a) diretório com 3 entradas

(b) após a remoção da entrada voluminous

## Bibliografia Básica



[Andrew S. Tanenbaum.](#)

*Sistemas Operacionais Modernos*, 3ª Edição. Capítulo 4.

[Pearson Prentice-Hall](#), 2010.



[Abraham Silberchatz, Greg Gagne e Peter Baer Galvin.](#)

*Fundamentos de Sistemas Operacionais*, 6ª Edição. Capítulos 11-12.

[LTC - Livros Técnicos e Científicos Editora S.A.](#), 2004.