

Big Data & Data Science

Infraestrutura Computacional **Parte 1: Linux e Shell**



Armazenamento de Dados

Hard Disk Drive



Solid State Drive



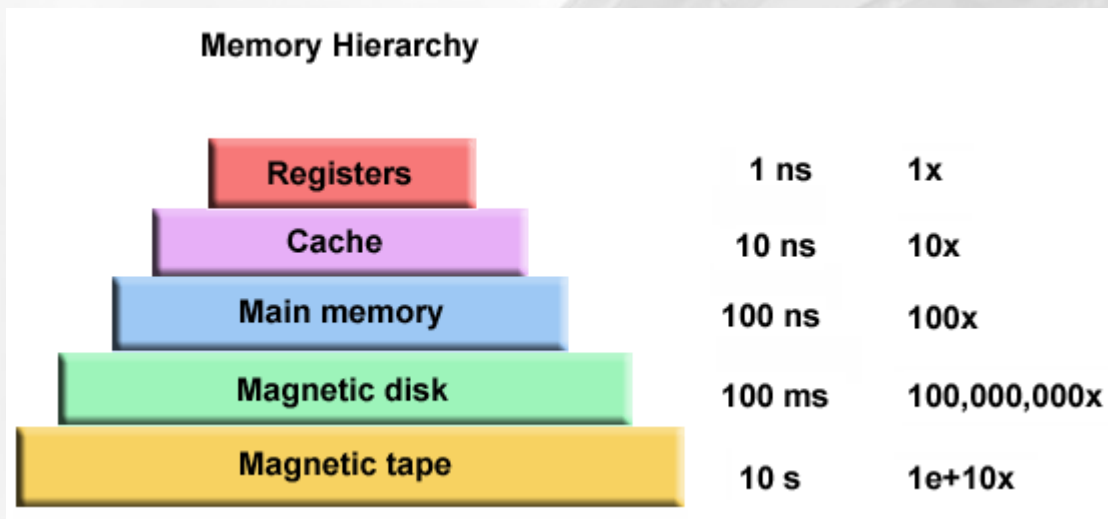
https://en.wikipedia.org/wiki/Solid-state_drive



Armazenamento de Dados

Volátil x não volátil

Acesso aleatório x sequencial



Redundant Array of Inexpensive Disks (RAID)

Combinação de dois ou mais discos

Padrões de organização, ou níveis:

- ▶ RAID 0: *stripping*
- ▶ RAID 1: espelhamento
- ▶ RAID 6: *stripping* de blocos com redundância de 2 discos
- ▶ RAID 1+0 ou 10: combinação

Implementação via software (md) ou hardware (controladora)

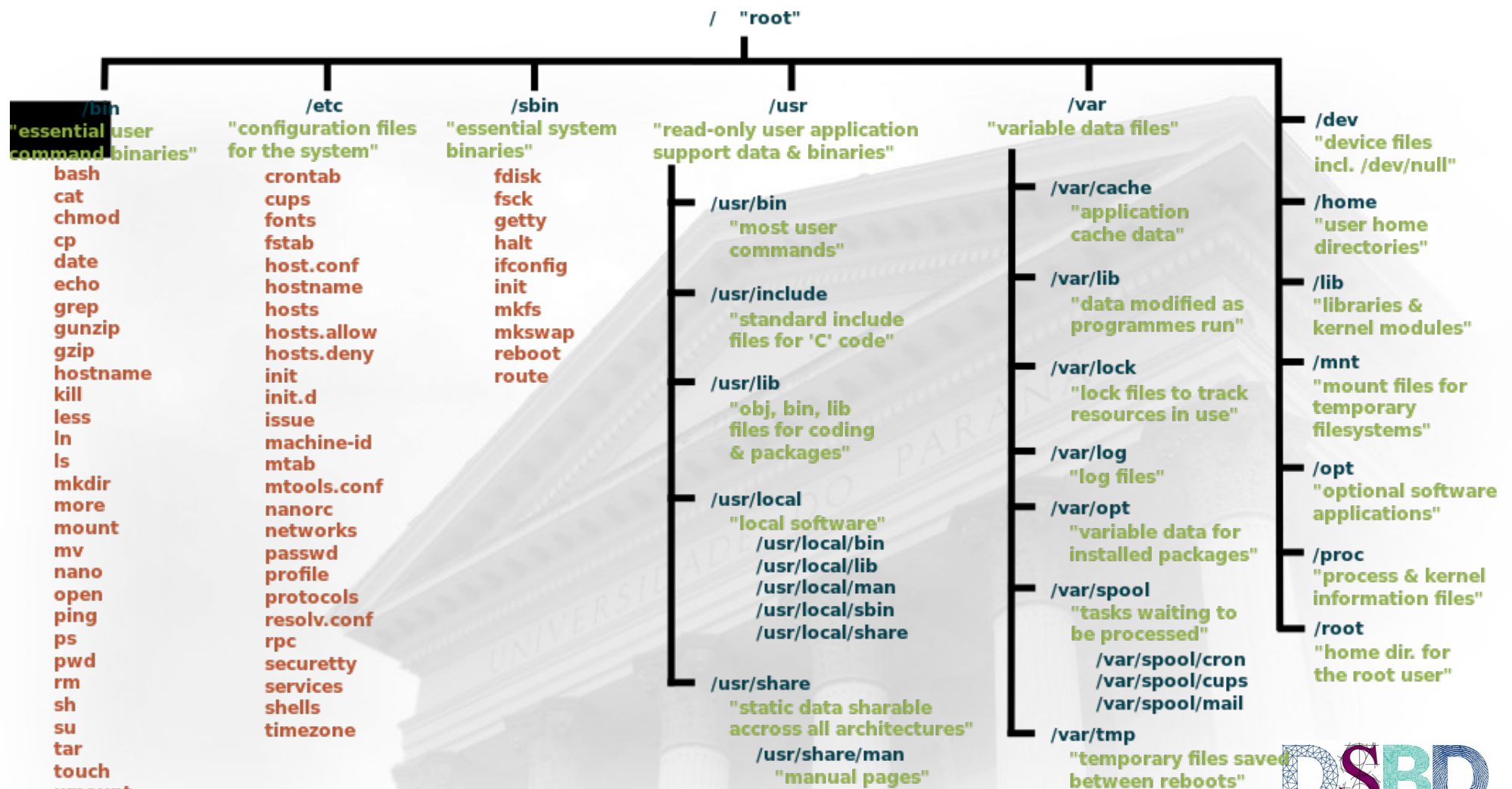


Sistemas de Arquivo

Sistemas de Arquivo

Virtual File System (VFS)

- ▶ Camada do Linux que permite acesso uniforme a diversos sistemas de arquivo
- ▶ Especifica uma interface (API) de acesso a arquivos. Padrão POSIX (open, close, read, write, seek, link, ...)
- ▶ 2 conceitos fundamentais: arquivos e diretórios
- ▶ Estrutura de diretórios em árvore, com diretório raiz “/”
 - ▶ . : referência ao próprio diretório
 - ▶ .. : referência ao diretório pai



Sistemas de aquivo de disco

Gerencia os blocos de dados (setores) do disco

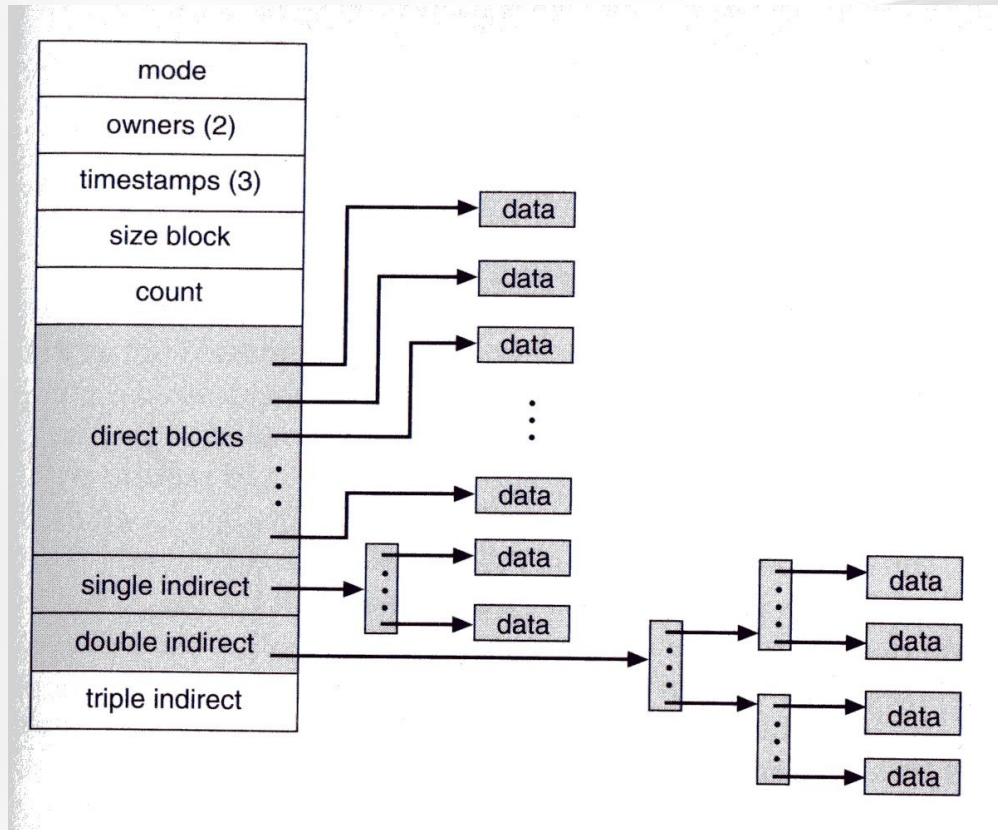
- ▶ Associa nomes de arquivos a blocos
- ▶ Mantém metadados
- ▶ Controla espaço livre (fragmentação), quota, permissões

Diversas implementações

- ▶ ext4, zfs, ntfs, fat
- ▶ Algumas implementações fazem versionamento ou *journaling*



Sistemas de Arquivo em disco



Partições

Partições proveem uma melhor separação dos dados em disco

- ▶ Cada partição tem seu próprio sistema de arquivos
- ▶ Comando **fdisk**, **parted**
- ▶ Partição de dados x partição de *swap*
- ▶ Segurança para falha no sistema de arquivos

Sistemas de Arquivo

Um SA precisa ser montado antes de ser acessado

- ▶ Montar significa indicar o diretório a partir do qual o SA será acessado neste computador
- ▶ A montagem sobrepõe qualquer dado existente (fica inacessível)
- ▶ **df -h**
- ▶ **mount**

Sistemas de Arquivo em Rede

Um SA pode ser exportado pela rede e montado por diversos clientes

- ▶ Arquitetura cliente-servidor
- ▶ Compartilhamento se dá através da montagem via protocolos específicos
- ▶ NFS, SAMBA, DNDB3

Clustered File System é um SA que distribui os dados em diversos servidores

- ▶ Redundância a falhas, acesso paralelo, escalabilidade
- ▶ Geralmente baseados em objetos (*object file system*):
 - ▶ Separação de Metadados e Dados
- ▶ GFS, Lustre, Hadoop, Gluster



Sistemas de Arquivo especiais

/dev

- ▶ Acesso direto aos dispositivos (*devices*) da máquina.

- ▶ `ls -l /dev/sd*`

/proc

- ▶ Acesso aos processos

- ▶ `cat proc/cpuinfo`

/sys

- ▶ Acesso aos dispositivos através do kernel. Muito utilizado para configurações

- ▶ `cat /sys/devices/system/cpu/cpu0/cpufreq/scaling_governor`



Arquivos

“Em um sistema UNIX, tudo é um arquivo; se algo não é um arquivo, é um processo”

Tipos de arquivo (opção **-l** do comando **ls**)

- ▶ **Arquivos regulares (-)**: armazenam dados. Não há divisão em nome.extensão
- ▶ **Diretórios (d)**: são arquivos que contém outros arquivos
- ▶ **Dispositivos (b, c)**: acessam dispositivos de hardware. Podem ser arquivos de bloco ou caracteres
- ▶ **Links (l)**: ponteiros para outros arquivos. Podem ser *soft* ou *hard links*.



Sistema de Arquivos Linux

Caminhos

- ▶ **\$PATH**
- ▶ caminhos relativos (. e ..) e absolutos (/)

Diretório HOME

- ▶ **quota -vs**
- ▶ **~**

Segurança de arquivos

GNU/Linux tem um sistema bastante rígido de permissões para arquivos

- ▶ Todo arquivo pertence a um usuário e um grupo
- ▶ As permissões de leitura, escrita, e execução devem ser definidas para o usuário, grupo e outros
- ▶ Comando `id` mostra usuário e grupos aos quais pertence
- ▶ Comando `ls -l` mostra permissões de forma posicional
 - ▶ usuário, grupo, outros



Modos de acesso para arquivos

Código	Significado
0 ou -	Acesso desta posição não concedido
1 ou x	Permissão de execução nesta posição
2 ou w	Permissão de escrita nesta posição
4 ou r	Permissão de leitura nesta posição
u	Permissão do usuário
g	Permissão do grupo
o	Permissão para outros
a	Permissão para todos

Cuidados no uso de arquivos

Acesso a arquivos é várias ordens de grandeza mais lento que o processador

Performance é dominada pelo número de acessos a disco

- ▶ ~10 ms por acesso

Custo do acesso é dominado pela latência

- ▶ *tempo de busca + latência de rotação + bytes / bandaDisco*
 - ▶ 1 setor: $5ms + 4ms + 2,5\mu s (\approx 512 B / 200 MB/s) \approx 9ms$
 - ▶ 100 setores: $5ms + 4ms + 0,25ms \approx 9,25ms$
 - ▶ 100 vezes mais dados com 3% de aumento no tempo



Perguntas?

Manipulando Arquivos

Manipulação de permissões

Comando	Significado
<code>chmod</code>	modifica as permissões de um arquivo
<code>chown</code>	modifica usuário ou grupo de um arquivo
<code>mkdir</code>	cria um diretório
<code>cp -R -avu</code>	copia um arquivo
<code>mv</code>	move um arquivo
<code>rm -f -r</code>	mostra o diretório corrente
<code>head</code> ou <code>tail</code>	mostra linhas iniciais ou finais do arquivo
<code>ln -s</code>	faz um link entre arquivos

Arquivos, usuários e permissões

Comando	Significado
<code>finger</code>	mostra informações sobre usuário
<code>who</code>	mostra usuários logados no sistema
<code>quota</code>	mostra a quota do usuário
<code>find</code>	procura arquivos



Manipulação de arquivos

Comando	Significado
<code>ls -a -l -R -F -t</code>	mostra a lista de arquivos de um diretório
<code>file <arq></code>	mostra o tipo do arquivo
<code>mkdir</code>	cria um diretório
<code>cp -R -avu</code>	copia um arquivo
<code>mv</code>	move um arquivo
<code>rm -f -r</code>	remove um arquivo
<code>head</code> ou <code>tail</code>	mostra linhas iniciais ou finais do arquivo
<code>ln -s</code>	faz um link entre arquivos
<code>touch</code>	muda a data de um arquivo



Wildcards

São caracteres especiais usados para criar padrões definindo um conjunto de arquivos ou diretórios

- ▶ * - representa zero ou mais caracteres
- ▶ ? - representa apenas um caractere
- ▶ [] - representa um intervalo de caracteres
- ▶ [^] - representa a negação de um intervalo de caracteres

Wildcards

Comando	Significado
<code>ls b*</code>	Arquivos iniciando com b
<code>ls -ld .g*</code>	Arquivos iniciando com .g
<code>ls ?i*</code>	Arquivos com um caractere seguido de i
<code>ls [sv]*</code>	Arquivos que iniciem com s ou v
<code>ls *[0-9]*</code>	Arquivos com um dígito de 0-9
<code>mv *.*g fotos/</code>	Move arquivos .png e .jpg para dir. fotos
<code>ls *[:upper:]*</code>	Arquivos com uma letra maiúscula

Referências

- ▶ Anatomy of the Linux kernel
- ▶ Linux OS Tutorial
- ▶ Introduction to UNIX
- ▶ Introduction to Linux

▪