

The Abdus Salam International Centre for Theoretical Physics

Introduction to Storage for Climate Data

Stefano Cozzini(CNR/IOM Democritos), Gilberto Díaz(ULA)

May 16, 2011

イロト イヨト イヨト イヨト

Preliminaries

Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イヨト イヨト イヨト

Introduction

Many applications perform relatively simple operations on vast amounts of data. In such cases, the performance of a computer's data storage devices impact overall application performance more than processor performance.

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イヨト イヨト イヨト

Introduction

Data storage devices include, but are not limited to:

- processor registers
- caches memories
- main memory
- disk (hard, compact disk, etc.)
- and magnetic tape.

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

Memory Hierarchy

- Computers architectures try to keep data close to the processors in order to feed them continuosly.
- However, while the capacity of storage devices increases, the distance to the processors also increases.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

< ロ > < 同 > < 三 > < 三 >

Current Mass Storage Devices

- We are interested particularly in low cost storage devices with big capacity and high performance.
- Nowaday, Magnetic hard Disk Drives still are the technology which include all these features.

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イポト イヨト イヨト

Current Hard Disk Drives Technologies

We can find several magnetic hard disk tecnologies today

- Serial Advanced Technology Attachment (SATA)
- Small Computer System Interface (SCSI)
- Serial Attached SCSI (SAS)

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

Current Hard Disk Technologies

- Today disk space is cheap, a single (SATA) disk drives provides up to 3TB.
- However, performance is another story. Fastest hard disk drive (SAS) bandwith is around 3Gbps

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

A typical HDD includes a plurality of magnetics disks spun by a spindle motor.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

Read/Write heads supported bye slider suspension assembly which are moved by some actuators in radial direction.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

Read/Write heads supported bye slider suspension assembly which are moved by some actuators in radial direction.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

We can identify, on each plate, specific zones.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

If we see one specific track on all plates, we can see a cylinder.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

If we see one specific track on all plates, we can see a cylinder.



・ロト ・回ト ・ヨト ・ヨト

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

The policy used by the heads, to respond every read/write operation is similar to the policy used by a lifter.



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

HDD components

In this way, the most visited zone of the plate is the zone in the middle.



・ロト ・回ト ・ヨト ・ヨト

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

Redundant Array of Independent Disks

- One way to improve bandwidth is to define a logical device which consists of multiple disks.
- With this sort of approach a single I/O transaction can simultaneously move blocks of data to multiple disks.
- For example, if a logical device is created from eight disks, each of which is capable of sustaining 10 MB/sec, then this logical device is capable of delivering up to 80 MB/sec of I/O bandwidth.

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イポト イヨト イヨト

Redundant Array of Independent Disks

We can get some better performance using RAID.

- S: Hard disk drive size.
- N: Number of hard disk drives in the array.
- *P*: Average performance of a single hard disk drive (MB/seg).

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

・ロト ・四ト ・ヨト ・ヨト

臣

Linear RAID



 $\begin{aligned} & Performance = P \\ & \text{No redundancy} \\ & Capacity = N * S \end{aligned}$

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

RAID 0



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

RAID 1



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ●

æ

RAID 10



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イヨト イヨト イヨト

э

RAID 4



Parity Disk Bottleneck

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イヨト イヨト イヨト

э

RAID 5

One disk can fail Distributed parity



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

RAID 6

Two disk can fail Double parity code



イロン 不同 とくほど 不同 とう

э

Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

Logical Volume Manager

From phisical devices we can create:

- Volume groups (Mlvg)
- Logical volumes (logical partitions): /home, /opt, /tmp



Preliminaries Memory Hierarchy Basic Concepts of Hard Disk Drives Hardware Parallelism Software Parallelism

イロト イヨト イヨト イヨト

Software RAID

Several Operating Systems provide software RAID. For instance, Linux provides

- Linear RAID
- RAID 0
- RAID 1
- RAID 1+0
- RAID 5
- RAID 6

File Systems

File System

A file system is a set of methods and data structures used to organize, store, retrieve and manage information in a permanent storage medium, such as a hard disk. Its main purpose is to represent and organize resources storage.

< ロ > < 同 > < 三 > < 三 >

File System

The elements of a file system are:

- Name space: It is a way to assign names to the items stored and organize them hierarchically.
- API: Is a set of calls that allow the manipulation of stored items.
- Security Model: Is a scheme to protect, hide and share data.
- Implementation: Is the code that couples the logical model to the storage medium.

Basic Concepts

- Disk: A permanent storage medium of a certain size.
- Block: The smallest unit writable by a disk or file system. Everything a file system does is composed of operations done on blocks.
- Partition: A subset of all the blocks on a disk.
- Volume: The term is used to refer to a disk or partition that has been initialized with a file system.
- Superblock: The area of a volume where a file

<u>system stores its critical data.</u>

Stefano Cozzini(CNR/IOM Democritos), Gilberto Díaz(ULA) Introducti

Introduction to Storage for Climate Data

File Systems

Basic Concepts

- Metadata: A general term referring to information that is about something but not directly part of it. For example, the size of a file is very important information about a file, but it is not part of the data in the file.
- Journaling: A method of insuring the correctness of file system metadata even in the presence of power failures or unexpected reboots.
- Attribute: A name and value associated with the name. The value may have a defined type (string, integer, etc.).

Stefano Cozzini(CNR/IOM Democritos), Gilberto Díaz(ULA)

File Systems

イロン イヨン イヨン イヨン

File Systems

- FAT
- NTFS
- Ext2, Ext3, Ext4
- Reiserfs, xfs, jfs

File Systems

Some Performance Tests

(Write) record=512 KB, threads = 4



File Systems

Some Performance Tests



File Systems

臣

Some Performance Tests

4 threads, 16GB file, 512KB record size



Distributed File Systems Parallel File Systems

Distributed File Systems

A distributed file system makes available the data across the network: NFS, coda, AFS, etc.



Distributed File Systems Parallel File Systems

イロト イヨト イヨト イヨト

Parallel File Systems



Distributed File Systems Parallel File Systems

イロト イヨト イヨト イヨト

臣

Parallel File Systems



Distributed File Systems Parallel File Systems

・ロト ・回ト ・ヨト ・ヨト

Parallel File Systems

We can find several implementations of Parallel File Systems:

- Lustre from Cluster File System, then from Sun, now from Oracle (GPL)
- GPFS from IBM (Private)
- PVFS from ANL (GPL)
- GlusterFS from Gluster Inc. (GPL)
- PanFS from Panasas (Private)
- IBRIX from Ibrix (Private).
- Etc.

Distributed File Systems Parallel File Systems

イロト イヨト イヨト イヨト

Parallel File Systems



