

Storage Systems

NPTEL Course

Jan 2012

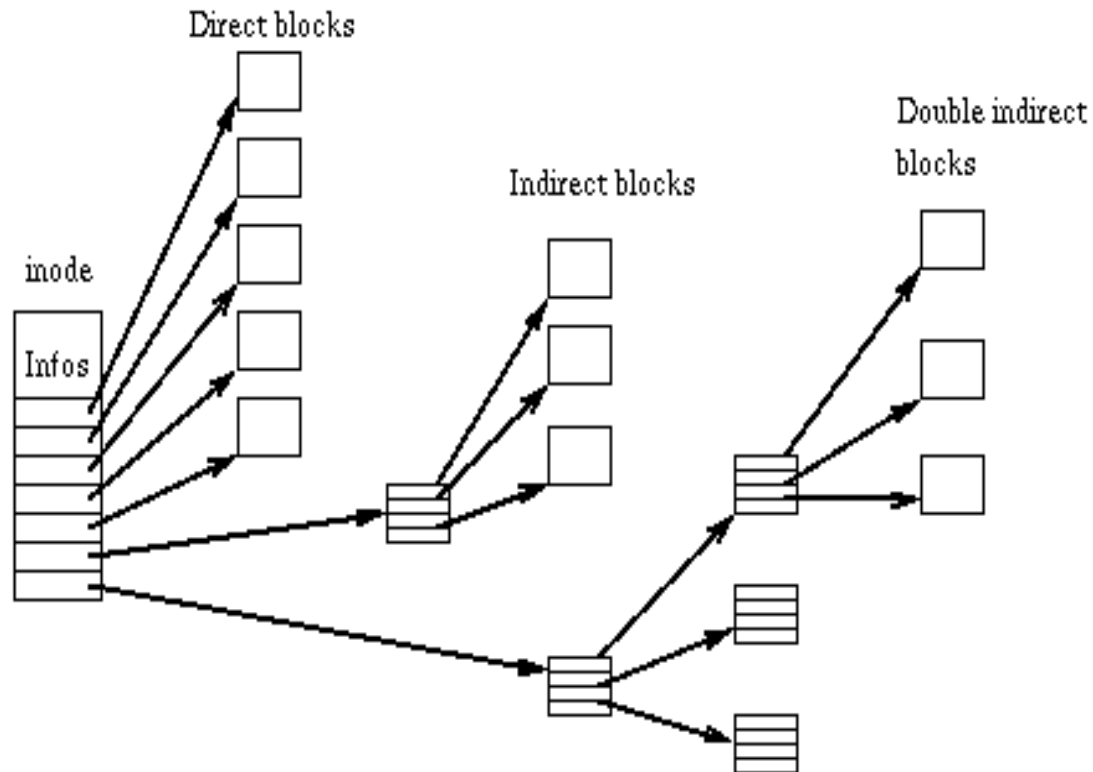
(Lecture 04)

K. Gopinath

Indian Institute of Science

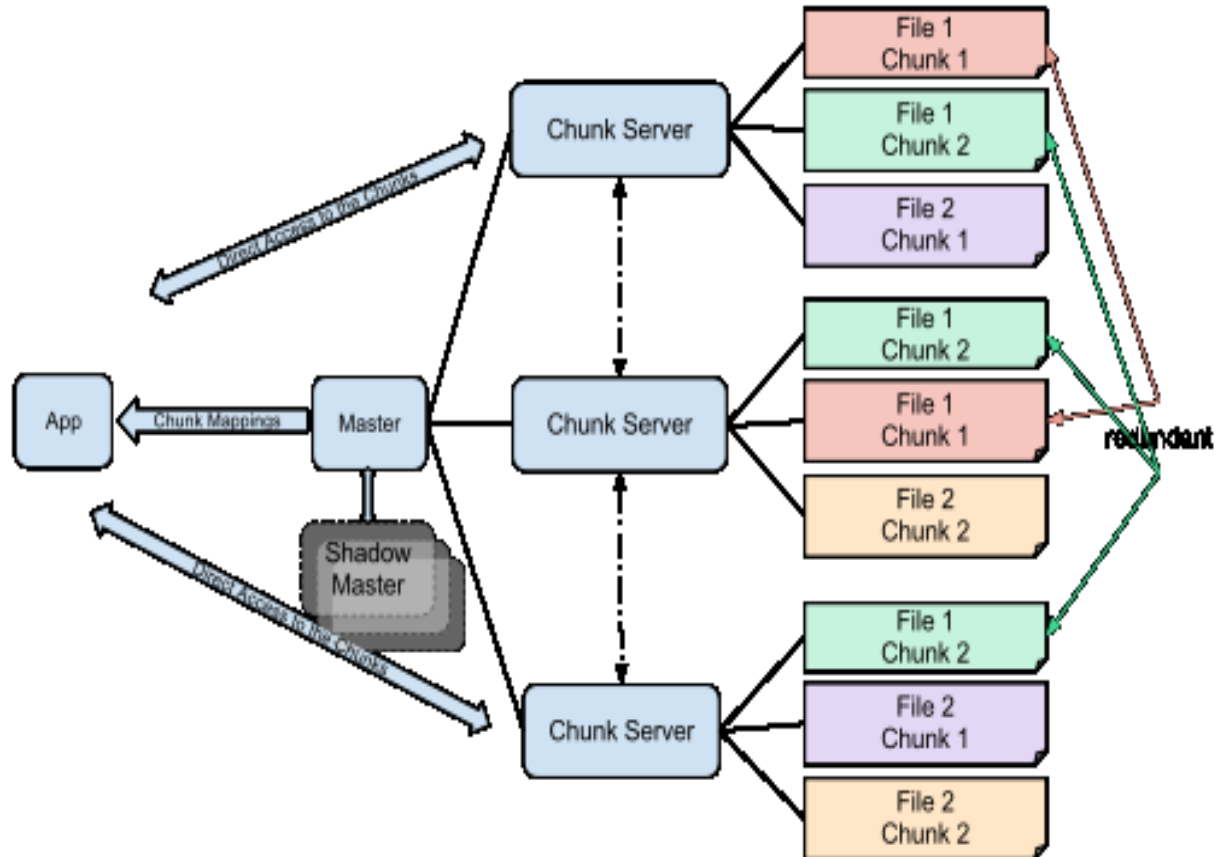
Ext2 FS

(from Wikipedia)



Google FS

(from Wikipedia)



Deep Storage Stack

- Various types of abstractions in stack
 - Device
 - Block
 - File
 - Application level (buffering in libc, for eg.)
- Finer sublayers in each layer
 - SCSI has upper (device-specific), mid (protocol specific) and lower (physical communication layer)
- For scalability, network stack part of storage stack
- A good part of stack in kernel
 - Increasingly, storage stack migrating out of kernel with separation of processing and storage (eg. GFS)

Storage Characteristics

- Concurrency arises naturally
 - Wide disparity in speeds of memory and storage
 - Have to mask slowness of storage
 - Make processing and storage go at their own rates and use interrupts (or sometimes polling) to signal completion of slow storage operations
 - Historic reason why operating systems developed
- Storage has to be typically persistent over time
 - Amount increases typically with time
 - But not all imp over time; keep imp part in fast storage?
- “Caching” and “tiering” arise naturally
 - Have to choose 2 of 3: speed, capacity or cost

Storage Performance

- Storage often slowest component
 - Cache!
- Within single device, efficiency by:
 - merging requests
 - scheduling requests in an order that is best wrt device (out-of-order execution commonplace)
 - Higher level software has to work around this aspect
 - If a particular order required, left to “user”
 - Semantically not much guaranteed
- Asynchronous processing often used: aio
- Parallelism across multiple devices/threads:
 - Multiple Heads (Disks)
 - Multiple chips (SSDs)

Optimization Framework

- Due to slowness of devices, optimization of accesses important
 - eg. what to cache, what to prefetch?
 - But usage patterns typically not known *a priori*
 - Big difference in performance whether sequential access or random
 - System slow if too many on demand migrations from slow to fast tier of storage (latency delays)
- Often, opts. critical and override “semantics”
 - Out-of-order processing typical
 - Complex higher-level software
- Learning on the job important
 - Simple and robust methods useful

Storage Protocols

- Interrupt driven rather than wait/poll
 - On completion, interrupt CPU or HBA
 - To avoid interrupt overhead, HBA or similar agents
 - Helps Segmentation and Reassembly (SAR)
- Split-phase transactions common
 - for eg: on completion of (a long) seek, slave takes bus
- Protocol endpoints preferably “virtualizable”
 - SCSI devices can be on an electrical bus, network or Internet if physical layer handled correctly
 - Protocols survive much longer
 - Devices can have arbitrary structure as long as they speak SCSI protocol
 - Even big servers!

Summary

- Storage systems design has many ramifications for the rest of the system
 - Provide abstractions based on application needs and devices
 - Design needs to be sensitive to cost, devices, manageability
 - Introduce newer abstractions with time
 - eg. key value stores
- Storage systems need to scale to support large scale computing systems