

Distributed File Systems

CS 256/456
Dept. of Computer Science, University
of Rochester

11/21/2013

CSC 2/456

1

Distributed System

- A collection of loosely coupled processors interconnected by a communication network
 - Do not share memory or a clock
 - Remote versus local resource access

11/21/2013

CSC 2/456

2

Network vs. Distributed Operating Systems

- Network operating system:
 - Users aware of multiple machines
 - Data/computation migration -> user's responsibility
- Distributed operating system:
 - Access to remote resources similar to local
 - Data/computation migration under control of OS

11/21/2013

CSC 2/456

3

Data Migration

- Whole files (automated ftp): e.g., AFS first version
- On-demand (demand paging): e.g., NFS, SMB, AFS newer versions

11/21/2013

CSC 2/456

4

Computation Migration

- RPC
- Messages (more concurrency)

11/21/2013

CSC 2/456

5

Process Migration

- Rationale
 - Load balancing
 - Computation speedup
 - Hardware preference
 - Software preference
 - Data access
- Either automated (transparent) or user specifies how the process should migrate

11/21/2013

CSC 2/456

6

Distributed File System

- One in which
 - clients,
 - servers (file servers/service software), and
 - storage devices (potentially multiple and independent)

are dispersed among the machines of a distributed system and requires networked/messaging activity

11/21/2013

CSC 2/456

7

Naming and Transparency

- Possible schemes
 - Host:local-name
 - Attach remote directories (mount: NFS)
 - Single global name structure (AFS)
- Location transparency vs. location independence?

11/21/2013

CSC 2/456

8

Caching and Coherence

- Remote-service mechanism
 - Stateful vs. stateless
- Caching and coherence
 - Caching granularity
 - Cache update policy (write through vs. delayed write)
 - Client-initiated vs. server-initiated

11/21/2013

CSC 2/456

9

Reliability and File Replication

- Reliability and file replication
 - Naming transparency
 - Availability vs. consistency

11/21/2013

CSC 2/456

10

NFS (Network File System)

- Deals with heterogeneity using RPC/XDR
- Stateless – no open and close of files
- Interface transparent to user via VFS

11/21/2013

CSC 2/456

11

The Andrew File System (AFS)

- Unified namespace (one /afs for everybody)
- One read-write copy
- Protection using access control lists (ACLs)
- Security using Kerberos 5 authentication

11/21/2013

CSC 2/456

12

AFS: Namespace

- Managed by dedicated servers called “Vice”
- Local (root file system) and shared namespace with mounting similar to NFS
- Files partitioned into volumes
 - typically files of a single client
 - Possible migration and replication of volumes
- Clients run “Virtue” protocol

11/21/2013

CSC 2/456

13

AFS: Security and Protection

- Security
 - Connection-based communication based on RPC
 - Encrypted
- Protection
 - Directories have access control lists
 - Allowed users or users not allowed
 - Regular UNIX bits for file protection

11/21/2013

CSC 2/456

14

AFS: Caching and Coherence

- Caching of entire files
 - Callback mechanism to eliminate cached copies
 - One read-write copy

11/21/2013

CSC 2/456

15

Distributed File Systems Issues

- Naming and transparency (location transparency versus location independence)
 - Host:local-name
 - Attach remote directories (mount)
 - Single global name structure
- Remote file access
 - Remote-service mechanism
 - Stateful vs. stateless
 - Caching and coherence
 - Cache update policy (write through vs. delayed write)
 - Client-initiated vs. server-initiated
- Reliability and file replication
 - Naming transparency
 - Availability vs. consistency

11/21/2013

CSC 2/456

16

Distributed File Systems: Issues of Scale

11/21/2013

CSC 2/456

17

Google File System - Motivation

- ▶ Component failures are normal
 - ▶ Fault tolerance and automatic recovery are needed
- ▶ Huge files are common (Multi GB)
 - ▶ Revisited I/O and block size assumptions
- ▶ Record appends are prevalent than random writes
 - ▶ Appending is the focus of performance optimization and atomicity guarantees
- ▶ Co-designing applications and the file-system API for flexibility
 - ▶ Relaxed consistency model
 - ▶ Atomic append

11/21/2013

CSC 2/456

18

What is GPFS (1/2)

- ▶ A *parallel, shared-disk* file system for *cluster* computers
- ▶ Developed by *IBM* since 1993
- ▶ Available on AIX, linux and MS Windows Server 2003.
- ▶ Used on 6 of the 10 most powerful supercomputers in the world
Extreme scalability!!

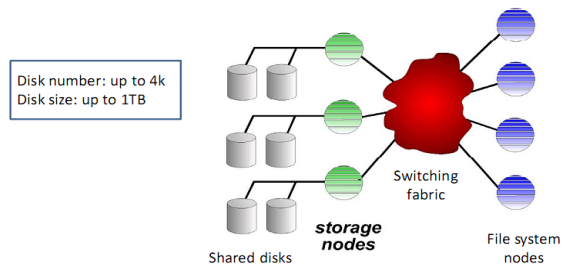
11/21/2013

CSC 2/456

19

What is GPFS (2/2)

▶ *Extreme scalability* from its shared-disk architecture



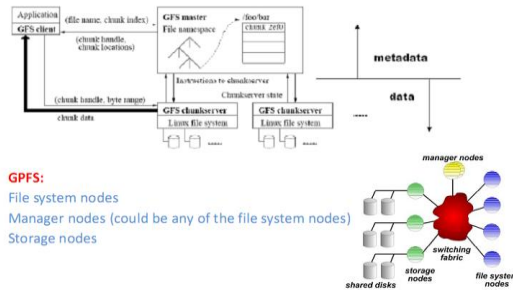
11/21/2013

CSC 2/456

20

GPFS vs. GFS

GFS: Single Master + Multiple trunk server + Multiple client



GPFS:

File system nodes

Manager nodes (could be any of the file system nodes)

Storage nodes

11/21/2013

CSC 2/456

21

	Google File System	GPFS
Application Type	Large, distributed data-intensive	Supercomputing other...
Data access assumption	Large streaming r/w Mainly record appends	None
File Size assumption	Usually huge	None
Consistency	Relaxed	-
Synchronization	Centralized Management	Distributed locking + Centralized Management
Caching	Not needed	Needed
Data Unit	Chunk (64MB)	Block(typically 256KB)
Fault Tolerance	Constant monitoring Fast recovery Replication	Logging & recovery Support of RAID Replication

11/21/2013

CSC 2/456

22