

## Remote FS Performance/Robustness

- 14C Remote Data - Robustness
- 14D Remote Data - Performance
- 14E Remote Data - Consistency

## Robustness: Embracing Failure

- Failures are inevitable
  - more components have more failures
  - complex systems have more modes of failure
  - we cannot build perfect components or systems
- We must build robust systems
  - additional capacity to survive failures
  - automatic failure detection
  - dynamically adapt to the new reality
  - continue service, despite component failures

Distributed File Systems

2

## Reliability and Availability

- Reliability ... probability of not losing data
  - disk/server failures to not result in data loss
    - RAID (mirroring, parity, erasure coding)
    - copies on multiple servers
  - automatic recovery (of redundancy) after failure
- Availability ... fraction of time service available
  - disk/server failures do not impact data availability
    - backup servers with automatic fail-over
  - automatic recovery (back up to date) after rejoin

Distributed File Systems

3

## Problems and Solutions

- Network Errors – support client retries
  - RFS protocol uses idempotent requests
  - RFS protocol supports all-or-none transactions
- Client Failures – support server-side recovery
  - automatic back-out of uncommitted transactions
  - automatic expiration of timed out lock leases
- Server Failures – support server fail-over
  - replicated (parallel or back-up) servers
  - stateless RFS protocols
  - automatic client-server rebinding

Distributed File Systems

4

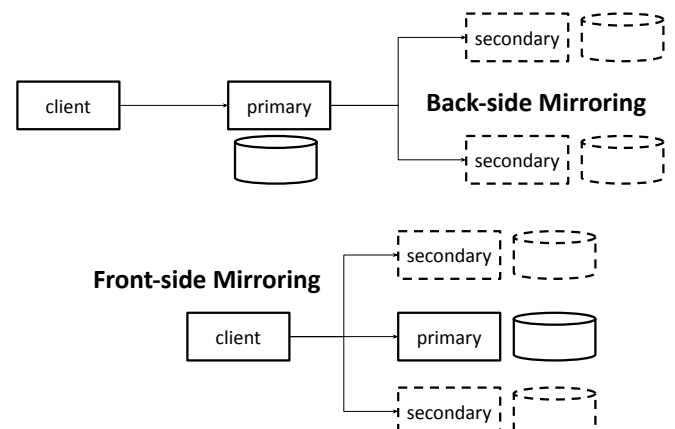
## Availability: Fail-Over

- data must be mirrored to secondary server
- failure of primary server must be detected
- client must be failed-over to secondary
- session state must be reestablished
  - client authentication/credentials
  - session parameters (e.g. working directory, offset)
- in-progress operations must be retransmitted
  - client must expect timeouts, retransmit requests
  - client responsible for writes until server ACKs

Distributed File Systems

5

## Reliability: Data Mirroring



Distributed File Systems

6

## Availability: Failure Detect/Rebind

- client driven recovery
  - client detects server failure (connection error)
  - client reconnects to (successor) server
  - client reestablishes session
- transparent failure recovery
  - system detects server failure (health monitoring)
  - successor assumes primary's IP address
  - state reestablishment
    - successor recovers last primary state check-point
    - stateless protocol

## Availability: Stateless Protocols

- a statefull protocol (e.g. TCP)
  - operations occur within a context
  - each operation depends on previous operations
  - successor server must remember session state
- a stateless protocol (e.g. HTTP)
  - client supplies necessary context w/each request
  - each operation is complete and unambiguous
  - successor server has no memory of past events
- stateless protocols make fail-over easy

## Availability: Idempotent Operations

- can be repeated many times with same effect
  - read block 100 of file X
  - write block 100 of file X with contents Y
  - delete file X version 3
  - non-idempotent operations
    - read next block of current file
    - append contents Y to end of file X
- if client gets no response, resend request
  - if server gets multiple requests, no harm done
  - works for server failure, lost request, lost response
    - but no ACK does not mean operation did not happen

## (nearly) Stateless Protocols

- client can maintain the session state
  - e.g. file handles and current offsets
- write operations can be made idempotent
  - e.g. associate a client XID with each write
- idempotence doesn't solve multi-writer races
  - competing writers must serialize their updates
  - clients cannot be trusted to maintain lock state
- we need a state-full Distributed Lock Manager
  - for whom failure recovery is extremely complex

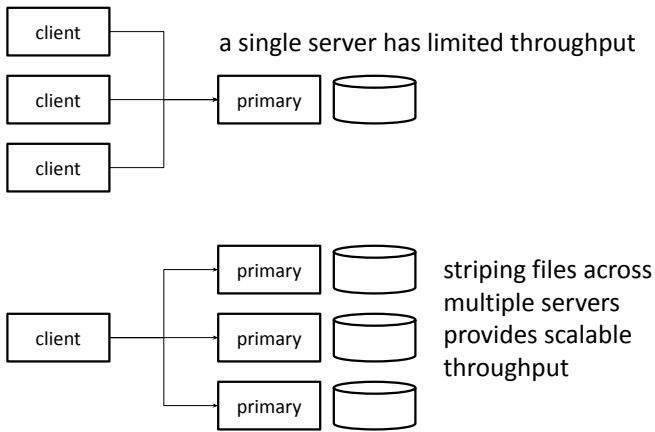
## Peter Deutsch Warned Us!

- POSIX semantics require coherent state
  - this complicates server fail-over
- there are numerous shared resources
  - must synchronize all updates to all of them
- location transparency
  - remote objects are much more expensive to use
- distributed is not really the same as local
  - performance is proportional to locality
  - adds more frequent and new modes of failure

## Performance Challenges

- single client response-time
  - remote requests involve messages and delays
  - error detection/recovery further reduces efficiency
- aggregate bandwidth
  - each client puts message processing load on server
  - each client puts disk throughput load on server
  - each message loads server NIC and network
- WAN scale operation
  - where bandwidth is limited and latency is high
- aggregate capacity
  - how to transparently grow existing file systems

## Performance: Bandwidth



## Performance: Minimize Messaging

- Protocol features
  - as few messages as possible
  - client-side caching to eliminate read requests
  - aggregation for fewer/larger write requests
- Work Partitioning
  - do as much as possible on the client
  - do as much as possible on a single server
  - eliminate multi-node coordination
  - eliminate multi-node request forwarding

## Performance: Read Requests

- client-side caching
  - eliminate waits for remote read requests
  - reduces network traffic
  - reduces per-client load on server
- whole file (vs. block) caching
  - higher network latency justifies whole file pulls
  - stored in local (cache-only) file system
  - satisfy early reads before entire file arrives
  - risk: may read data we won't actually use

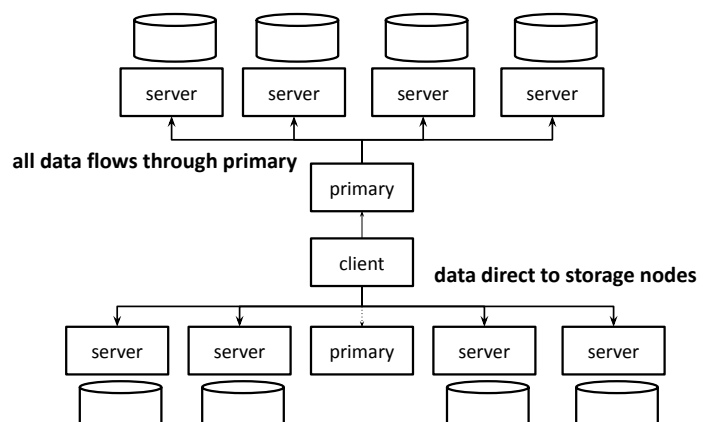
## Performance: Write Requests

- write-back cache
  - create the illusion of fast writes
  - combine small writes into larger writes
  - fewer, larger network and disk writes
  - enable local read-after-write consistency
- whole-file updates
  - wait until *close(2)* or *fsync(2)*
  - reduce many successive updates to final result
  - possible file will be deleted before it is persisted
  - enable atomic updates, close-to-open consistency

## Performance: Cost of Mirroring

- multi-host vs multi-disk mirroring
  - protects against host and disk failures
  - creates much additional network traffic
- mirroring by primary
  - primary becomes throughput bottleneck
  - replication traffic on back-side network
- mirroring by client
  - data flows directly from client to storage servers
  - replication traffic goes through client NIC
  - parity/erasure code computation on client CPU

## Performance: Direct Data Path



## (benefits of direct data path)

- architecture
  - primary tells clients where which data resides
  - client communicates directly w/storage servers
- throughput
  - data is striped across multiple storage servers
- latency
  - no intermediate relay through primary server
- scalability
  - fewer messages on network
  - much less data flowing through primary servers

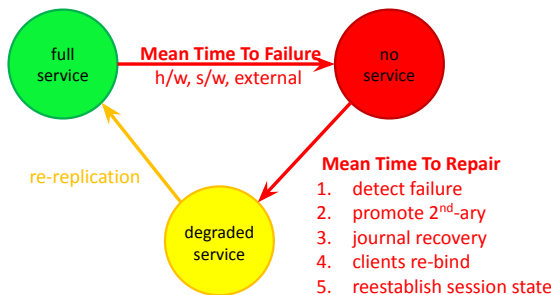
## Performance: Partitioning the Work

|                                   |                               |
|-----------------------------------|-------------------------------|
| open file instances, offsets      | <b>clearly on client side</b> |
| data packing and unpacking        |                               |
| -----                             |                               |
| authentication/authorization      | <b>either side (or both)</b>  |
| directory searching               |                               |
| block caching                     |                               |
| -----                             |                               |
| logical to physical block mapping | <b>clearly on server side</b> |
| on-disk data representation       |                               |
| device driver integration layer   |                               |
| device driver                     |                               |

## Performance: Recovery Time

Availability =  $\frac{MTTF}{MTTF + MTTR}$

we can try to maximize MTTF  
we can try to minimize MTTR



## (improving MTTR)

- MTTR (time before service can be restored)
  - primary failure detected (minimize)
  - secondary promoted to primary role (minimize)
  - recent/in-progress operations recovered
  - clients learn of change and re-bind
  - session state (if any) has been reestablished
- Degraded service may persist longer
  - restoring lost redundancy may take a while
  - heavily loading servers, disks, and network

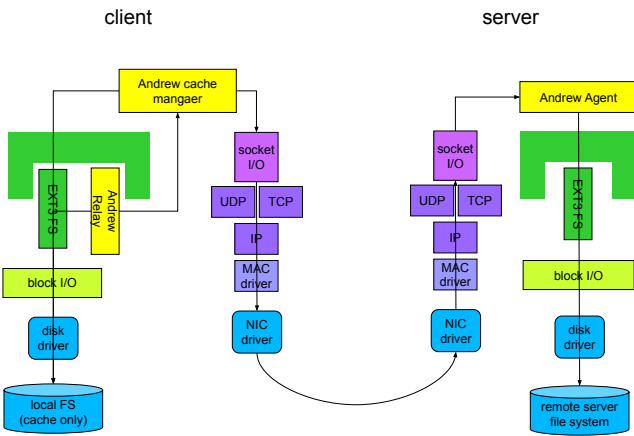
## Performance: Cost of Consistency

- caching is essential in distributed systems
  - for both performance and scalability
- caching is easy in a single-writer system
  - force all writes to go through the cache
- multi-writer distributed caching is hard
  - Time To Live is a cute idea that doesn't work
  - constant validity checks defeat the purpose
  - one-writer-at-a-time is too restrictive for most FS
  - change notifications are a reasonable alternative

## Andrew File System

- scalability, performance
  - large numbers of clients and very few servers
  - performance of local file systems
  - very low per-client load imposed on servers
  - no administration or back-up for client disks
- master files reside on a file server
  - local file system is used as a local cache
  - local reads satisfied from cache whenever possible
  - files are only read from server if not in cache
- simple synchronization of updates

# Andrew File System Architecture



Distributed File Systems



# (Andrew File System – Replication)

- check for local copies in cache at open time
  - if no local copy exists, fetch it from server
  - if local copy exists, see if it is still up-to-date
    - compare file size and modification time with server
  - optimizations reduce overhead of checking
    - subscribe/broadcast change notifications
    - time-to-live on cached file attributes and contents
- send updates to server when file is closed
  - wait for all changes to be completed
  - file may be deleted before it is closed



Distributed File Systems

26

# Andrew File System – Reconciliation

- updates sent to server when local copy closed
- server notifies all clients of change
  - warns them to invalidate their local copy
  - warns them of potential write conflicts
- server supports only advisory file locking
  - distributed file locking is extremely complex
- clients are expected to handle conflicts
  - noticing updates to files open for write access
  - dealing with conflicts: *not my job :-)*

Distributed File Systems

27

# Rating Andrew File System

- Performance and Scalability
  - all file access by user/applications is local
  - update checking (with time-to-live) is relatively cheap
  - both fetch and update propagation are very efficient
  - minimal per-client server load (once cache filled)
- Robustness
  - no server fail-over, but have local copies of most files
- Transparency
  - mostly perfect - all file access operations are local
  - pray that we don't have any update conflicts



Distributed File Systems

28

# Andrew File System vs. NFS

- design centers
  - both designed for continuous connection client/server
  - NFS supports diskless clients w/o local file systems
- performance
  - AFS generates much less network traffic, server load
  - they yield similar client response times
- ease of use
  - NFS provides for better transparency
  - NFS has enforced locking and limited fail-over
- NFS requires more support in operating system

Distributed File Systems

29

# Complication: Failure & Rejoin

- a file server goes down
  - no problem another server handles his clients
- then he comes back up and reports for work
  - he needs to get all the updates he missed
- How do we know what updates he missed?
  - we could compare all of his files with all of ours
    - that could take a very long time
  - we can keep a log of all recent updates
    - but we have to know which ones he already has
    - maybe files are versioned, or updates are numbered

Distributed File Systems

30

## Complication: Split-Brain

- suppose we had a network failure
  - that partitioned our file servers
  - and each half tried to take over for the other
  - and each half processed different write operations
- How could we reconcile the changes
  - we could merge updated versions of different files
  - what about files that were changed in both halves?
- Quorum rules can prevent “dueling servers”
  - servers that can’t make quorum are read-only

## Complication: Disconnected Operation

- Consider a notebook and a file server
  - I synchronize my notebook with the file server
  - I go away on a trip and update many files
  - others may change the same files on the server
- How can we identify all of the changes?
  - Intercept & log all changes (e.g. Windows Briefcase)
  - Differential Analysis vs. a baseline (e.g. rsync)
- How can we correctly reconcile conflicts?
  - perhaps some can be handled automatically
  - some may require manual (human) resolution

## Reading and Assignments

### Reading:

- Arpaci C10 ... Symmetric Multi-Processor scheduling
- Distributed Consensus
- Kampe: Multi-Processors
- Kampe: Clustering concepts

### Projects:

- we will help w/project 4C problems in the lab

## Supplementary Slides

## Network File System

- transparent, heterogenous file system sharing
  - local and remote files are indistinguishable
- peer-to-peer and client-server sharing
  - diskfull clients can export file systems to others
  - able to support diskless (or dataless) clients
  - minimal client-side administration
- high efficiency and high availability
  - read performance competitive with local disks
  - scalable to huge numbers of clients
  - seamless fail-over for all readers and some writers

## Network File System – Protocol

- idempotent operations and stateless server
  - built atop a Remote Procedure Call protocol
  - with eXternal Data Representation, server binding
  - versions of RPC over both TCP or UDP
  - optional encryption (may be provided at lower level)
- Scope – basic file operations only
  - lookup (open), read, write, read-directory, stat
  - supports client or server-side authentication
  - supports client-side caching of file contents
  - locking and auto-mounting done w/other protocol

# Network File System – Replication

- file systems can be replicated
  - improves read performance and availability
  - only one of these copies can be written to
- client-side agent (in OS) handles fail-over
  - detects server failure, rebinds to new server
- limited transparency for server failures
  - most readers will not notice failure (only brief delay)
  - users of changed files may get "stale handle" error
  - active locks may have to be re-obtained

# Network File System – Updates

- server does not prevent conflicting updates
  - as with local file systems, this is applications job
- auxiliary server/protocol for file and record locking
  - all leases are maintained on the lock server
  - all lock/unlock operations handed by lock server
- client/network failure handling
  - server can break locks if client dies or times out
  - "stale-handle" errors inform client of broken lock
  - client response to these errors are application specific
- lock server failure handling is very complex

# Rating NFS

- Transparency/Heterogeneity
  - local/remote transparency is excellent
  - NFS works with all major ISAs, OSs, and FSs
- Performance
  - read performance may be better than local disk
  - replication provides scalable read bandwidth
  - write performance slower than local disk
- Robustness
  - transparent fail-over capability for readers
  - recoverable fail-over capability for writers

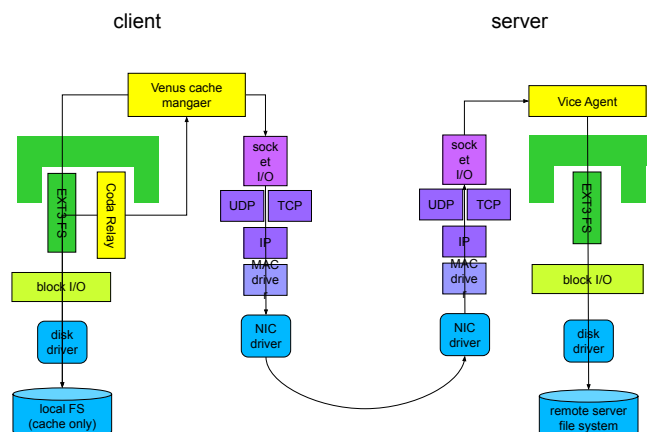
# CIFS vs. NFS

- functionality
  - NFS is much more portable (platforms, OS, FS)
  - CIFS provides much better write serialization
- performance and robustness
  - NFS provides much greater read scalability
  - NFS has much better fail-over characteristics
- security
  - NFS supports more security models
  - CIFS gives the server better authorization control

# CODA File System – Goals

- Start with the Andrew File System
  - scalability, performance, simple client admin
- Enhanced to support intermittent connection
  - clients can continue despite server outages
  - nomadic clients, with only occasional connection
- Enhanced to support server fail-over
  - for both read and write clients
- Enhanced conflict detection/resolution
  - automatic detection, and (when possible) resolution

# Coda File System Architecture



## Venus vs Andrew Cache manager

- Loss of contact with server
  - client attempts to switch to an alternate server
  - Client switches to disconnected mode
    - add all updates to Client Motivation Log (CML)
    - add files that need to be fetched to CML
  - upon reconnection, replay the CML
- File Hoarding
  - list of “important” files
    - can be generated automatically, or tuned manually
  - automatically keep local copies current
    - periodically check for changes to all of them

## CODA Conflict/Resolution

- clients do read one, write many
  - reads are satisfied from the nearest server
  - updates are multi-cast to all available servers
- clients detect version conflicts
  - at open time, fetch attributes from all servers
  - if multiple versions found, initiate resolution
- servers attempt automatic resolution
  - e.g. directories, e-mail, appointments, etc.
  - failing that, they request manual resolution

## Rating CODA File System

- Performance and Scalability
  - file access very similar to Andrew File System
  - “hoard walks” impose some additional server load
- Robustness
  - transparent fail-over for readers and many writers
  - disconnected operation (if all needed files are hoarded)
- Conflict Detection/Resolution
  - does not try to prevent all conflicts
    - requires more expensive distributed locking and lock manager
  - it does detect them, and attempts automatic resolution

## Approaches to Change Propagation

- Synchronous block mirroring
  - disk driver or controller replicates every block write
- Post-commit file replication
  - when a file change is committed, other servers are notified
- Automatic detection and periodic synchronization
  - intercept & log all change operations, replay to other servers
- Periodic differential comparison and copying
  - compare both sides and try to infer what has changed
- Copy what you remember when (and if) you remember
  - not a very good approach, but a surprisingly popular one