



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y R  
E N  
C E**

# DAFS Extensions for NFS

**Mark Wittle**

**Technical Director**

**Network Appliance**

**[mwittle@netapp.com](mailto:mwittle@netapp.com)**



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# Presentation Outline

- DAFS Overview
- Direct Access Networks
- DAFS Performance Operations
- DAFS Local File Sharing Enhancements
- Applicability for NFS



**N I C  
F N O  
S D N  
U S F  
T R E  
R E N  
C E**

# DAFS Overview



# DAFS Derived from NFSv4

- Many things are the same
  - Client-server, request-response messages
  - Basic file & directory operation set
  - File Attributes, ACLs
  - Locks, leases, delegations
  - Authentication types, internationalisation

**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS - NFS Differences

- **Some things are different**
  - Session-based (with authentication & attributes)
  - Credit-based flow-control (request throttling)
  - Transport channel aware
  - Different wire protocol format (XDR, Endian)
- **Two new areas of focus**
  - High performance
  - Local file sharing semantics



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# DAFS Performance Focus

- **RDMA Support**
  - Direct data placement
  - Pre-allocate transfer buffers
  - Split protocol header from application data
  - User-space I/O support
- **Sophisticated I/O Controls**
  - Batched and chained I/O operations
  - File access pattern hints for cache efficiency



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# DAFS Local File Sharing

- **Application Cluster Environment**
  - Access control for cooperating processes
  - Atomic file append operations
- **Failure Recovery**
  - Extended locking semantics
  - Fencing to control partial cluster failures
  - Server response cache management



N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E

# Direct Access Networks





**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C  
E**

# Direct Access Networks

- **Technology protocols**
  - InfiniBand, VIA, iWARP
- **Capabilities**
  - Multiple connections multiplexed in HW
  - Memory pre-registration
  - Send/receive messages
  - Remote DMA read/write
  - Asynchronous I/O completion



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# Direct Access Benefits

- No data packet fragmentation or reassembly
  - Benefits similar to IP Jumbo Frames, but with larger packets
  - Less transmission overhead with fewer HW interrupts
  - No ordering & space management issues
  - No data copying to recreate contiguous buffers



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# Direct Access Benefits (2)

- **No realignment of data copies**
  - Protocol headers and data buffers transmitted separately
  - Allows data alignment to be preserved
- **No user/kernel boundary crossing**
  - Less system call overhead
- **No user/kernel data copies**
  - Data transferred directly to application buffers



N I C  
F N O  
S D U  
I S T  
N D R  
D E  
Y E  
N C  
E

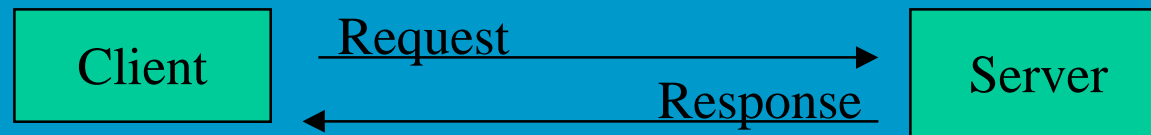
# DAFS Performance Operations



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# DAFS Inline I/O Operations

- **Inline Read and Write**
  - **Just like NFS**



- **Pre-allocated buffers, pre-registered with the transport**
- **Configurable message size limit**
- **Configurable data offset**
- **Low transport latency, simple model**

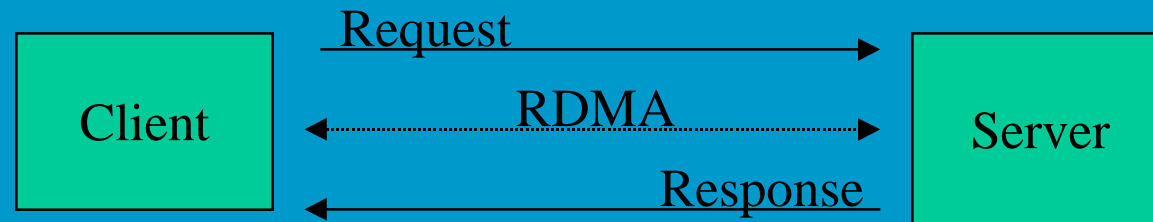


**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS Direct I/O Operations

- Direct Read and Write

- 3-part transfer



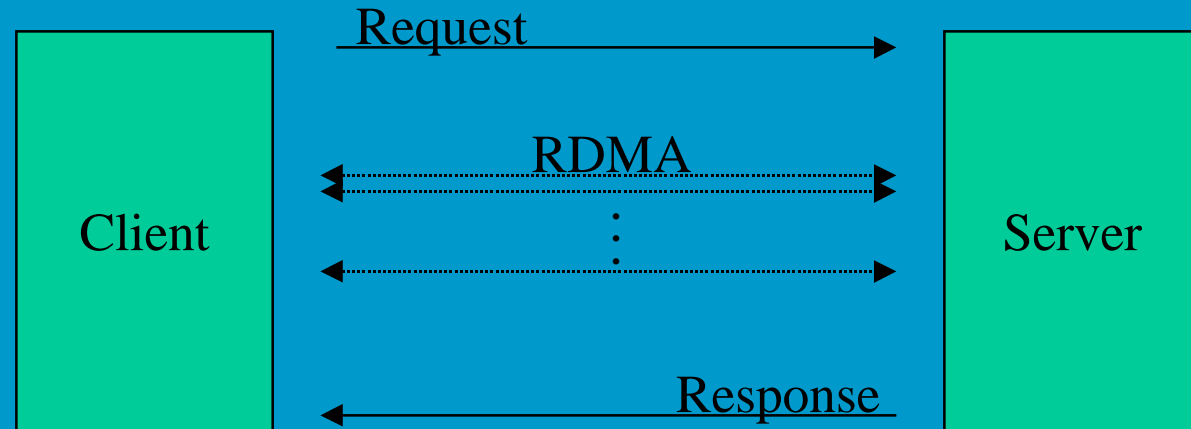
- Server initiates RDMA operation
- Receive buffer sized & placed at runtime
- Used for large messages
- Zero-copy, low CPU cost



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS Batch I/O Operations

- Synchronous Batch I/O
  - Multi-part transfer



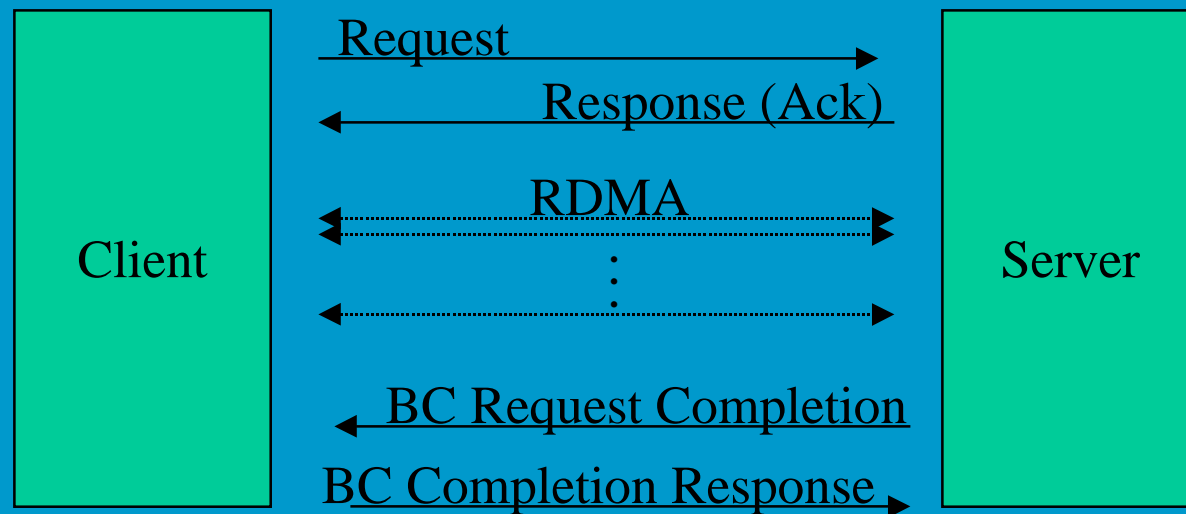
- Server initiates RDMA operation
- Supports multiple files, read/write, scatter-gather, sync/async, stable delay hint



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS Batch I/O Ops (2)

- Asynchronous Batch I/O
  - Multi-part transfer



- Server initiates RDMA operation





**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# DAFS Cache Hints

- **Similar to POSIX madvise**
  - Access prediction to inform caching policy
    - Random, Sequential, Unknown, Will Need, Don't Need
  - Buffer-specific future use hint to inform specific caching decisions
    - Two separate likelihood hints for future Read & Write access
- **Cache hint operation and hints included in I/O Operations**



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# Multi-component Lookup

- Multi-component path name lookup
  - An entire pathname is looked up in one operation
  - If an error is encountered partway through the lookup, partial lookup information is returned along with an error status.



**N I C  
F N O  
S D N  
U F  
S T R  
R E  
Y N  
C  
E**

# DAFS Local File Sharing



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# Shared Key Reservations

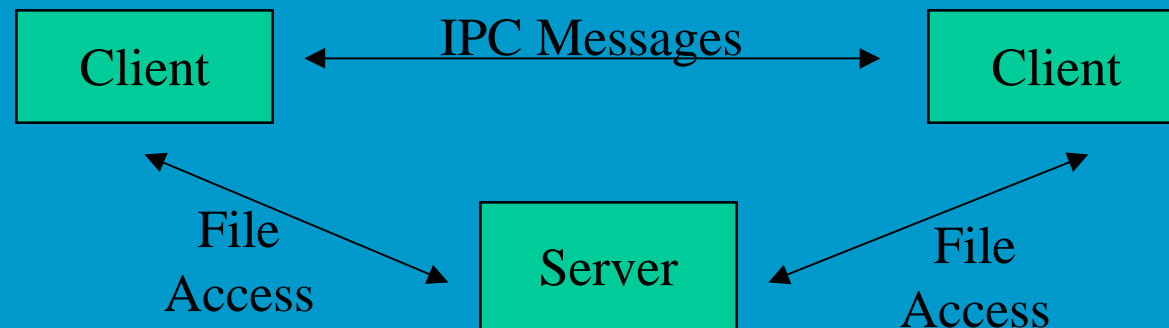
- Access control for cluster member instances
  - Addresses partial restarts
    - Same client, but different temporal instance
  - Reservation (lock) obtained at open with current "key"
  - Access control for cooperating clients
  - Key distribution left to application



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS Fencing

- Access control for cluster failures
  - Stable state

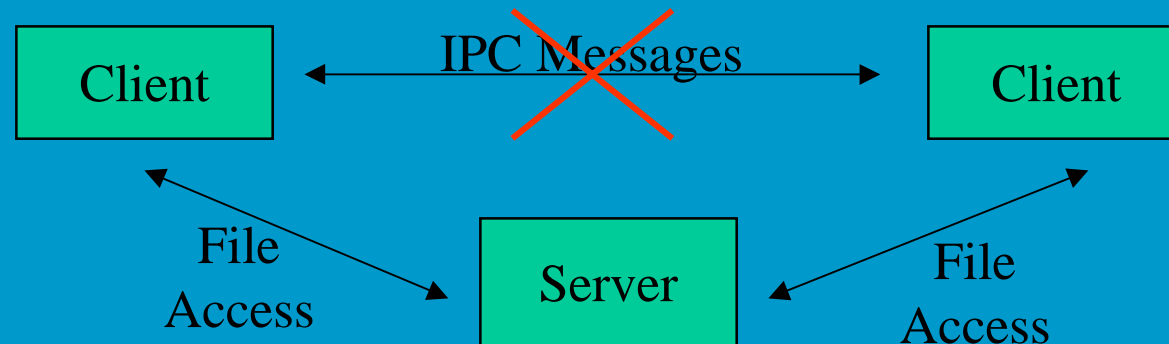




**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS Fencing (2)

- Access control for cluster failures
  - Partial failure

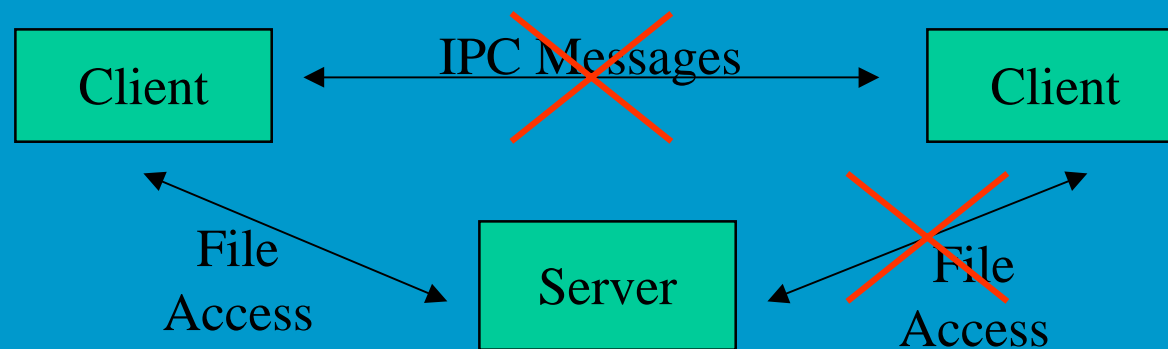




N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E

# DAFS Fencing (3)

- Access control for cluster failures



- Immediate access denial
- Drain operations on server to establish valid serialization point
- Manage a “fencing access list” of clients: allow, deny



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# DAFS Extended Locking

- **Persist locks**
  - take effect on lease expiration or on server reboot
  - remain locked until explicitly repaired
  - allow 3rd party recovery following failure
- **Auto-Recovery locks**
  - limited undo capability after failure
  - modifications performed under a rollback lock are undone following lease expiration or server reboot





**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
E**

# DAFS Atomic Append

- **Atomic relative to other writers**
  - Closes the gap between GET\_ATTR and WRITE
  - Server determines current file size and performs the write operation atomically
- **Atomic relative to the data buffer**
  - Server performs the I/O operation atomically, up to server file system attribute MAX\_APPEND\_SZ
- **Atomic relative to server failures**
  - Server performs synchronous I/O to stable storage



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# Other DAFS Operations

- **Set file attributes during creation**
  - Allow file attributes to be specified when a file is created
- **Open Unlinked**
  - Create a file without having a name in any directory. Subsequently the file can be named by linking it into the directory.
- **Delete On Last Close**
  - Addresses the .nfs file problem.



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
E**

# DAFS Response Cache

- **Session-based non-volatile Response Cache**
  - holds recent state modifying operations
  - number of entries bounded: transport message ordering properties and DAFS flow control constraints
  - Client queriable following a failure
- **Provides at-most-once semantics for state modifying file operations**



**N I C  
F N O  
S D N  
U S F  
T R E  
R Y N  
C E**

# Applicability to NFS (Speculation)



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C E**

# Performance Operations

- **Inline I/O operations**
  - Header padding for data for alignment
- **Batch I/O**
  - Similar to NFS Compound
  - Stable Hint useful for high throughput applications
- **Cache Hints**
  - Conduit for application madvise calls
  - Buffer caching hints for 2nd-level cache?



**N I C  
F N O  
S D N  
U S F  
T R E  
R E N  
Y N C  
E**

# Local File Sharing

- **Clustered application support**
  - Fencing, Shared key reservations
  - Are cluster applications becoming more prevalent?
  - No POSIX API, but good for user-space NFS
- **Atomic append**
  - Really a cluster application performance enhancement



**N I C  
F N O  
S D N  
U S F  
T R E  
R E N  
C E**

# Local File Sharing (2)

- **Multi-component Lookup**
  - Similar to NFS Compound?
- **Require persistent server state**
  - Open Unlinked (No POSIX API)
  - DOLC
  - Extended Locking (No POSIX API)
  - Persistent Response Cache



**N I C  
F N O  
S D N  
U F  
S E  
T R  
R E  
Y N  
C  
E**

# Conclusions

- **DAFS – NFS structural differences**
  - DAFS design (reliable transport assumption, session-orientation, flow-control) enables efficient buffer management, reliable response cache, etc
  - Not trivial to adapt to existing NFS
- **DAFS – NFS file sharing features**
  - Some features easy to introduce
  - May require more server state
  - POSIX API a limitation for some