



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

NFS Version 4

Spencer Shepler

Sun Microsystems

spencer.shepler@sun.com



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

Today's Goals

- Describe the NFSv4 protocol
- Identify deployment considerations
- Discuss areas of potential future enhancement



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

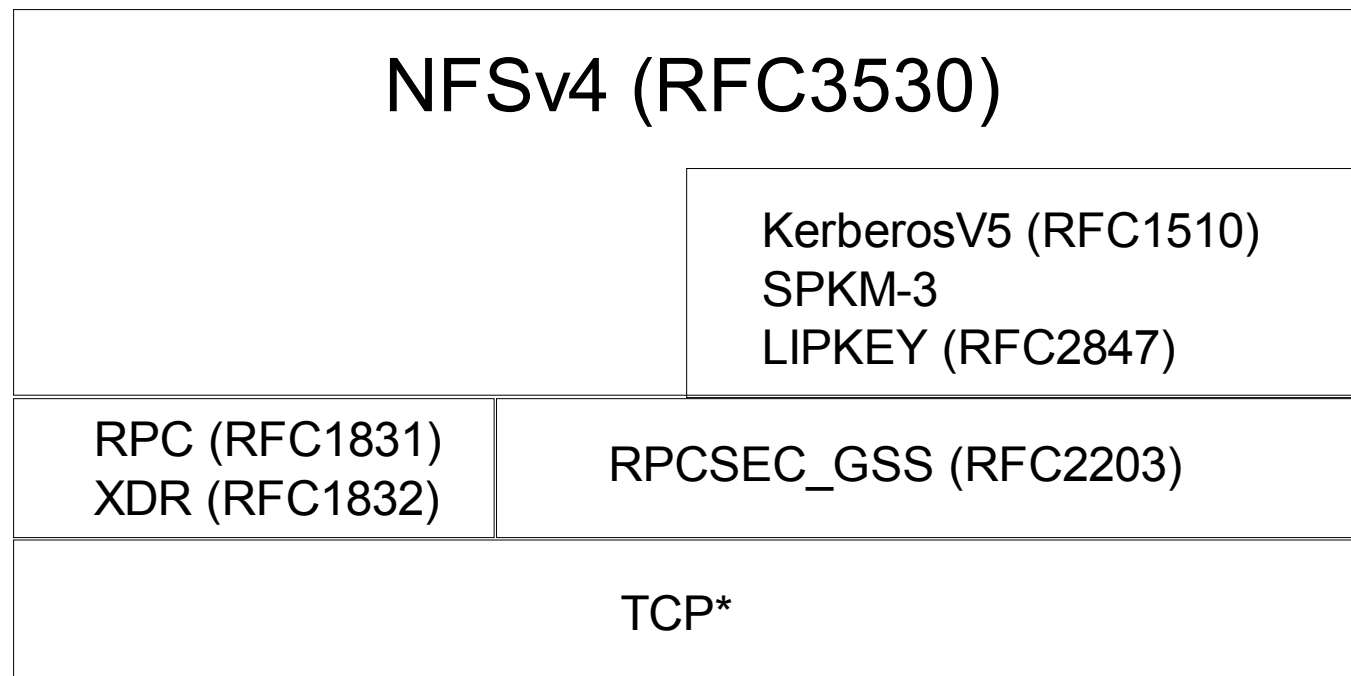
NFSv4 Design Goals

- Improved access and good performance on the Internet
- Strong security with negotiation built into the protocol
- Better cross-platform interoperability
- Designed for protocol extensions



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

NFSv4 Protocol “Stack”





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

Many To One

PORTMAP	→	Port 111
MOUNT	→	Dynamic
NFSv2/v3	→	Port 2049
LOCK/NLM	→	Dynamic
STATUS	→	Dynamic
ACL*	→	Dynamic

NFSv4	→	Port 2049 TCP
-------	---	------------------



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

Operation Count

- NFSv2 - 18 ops
- NFSv3 - 22 ops
- NFSv4 - 38 ops
- v2/v3 use “traditional” RPC
- v4 uses COMPOUND procedure to build sequence of operations



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

- ACCESS
- **CLOSE**
- COMMIT
- CREATE
- **DELEGPURGE**
- **DELEGRETURN**
- GETATTR
- **GETFH**
- LINK
- **LOCK**
- **LOCKT**
- **LOCKU**
- NULL
- **COMPOUND**
- LOOKUP
- **LOOKUPP**
- **NVERIFY**
- **OPEN**
- **OPENATTR**
- **OPEN_CONFIRM**
- **OPEN_DOWNGRADE**
- **PUTFH**
- **PUTPUBFH**
- **PUTROOTFH**
- READ
- READDIR
- READLINK
- REMOVE
- RENAME
- **RESTOREFH**
- **SAVEFH**
- **SECINFO**
- SETATTR
- **SETCLIENTID**
- **SETCLIENTID_CONFIRM**
- **VERIFY**
- WRITE
- **RELEASE_LOCKOWNER**



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

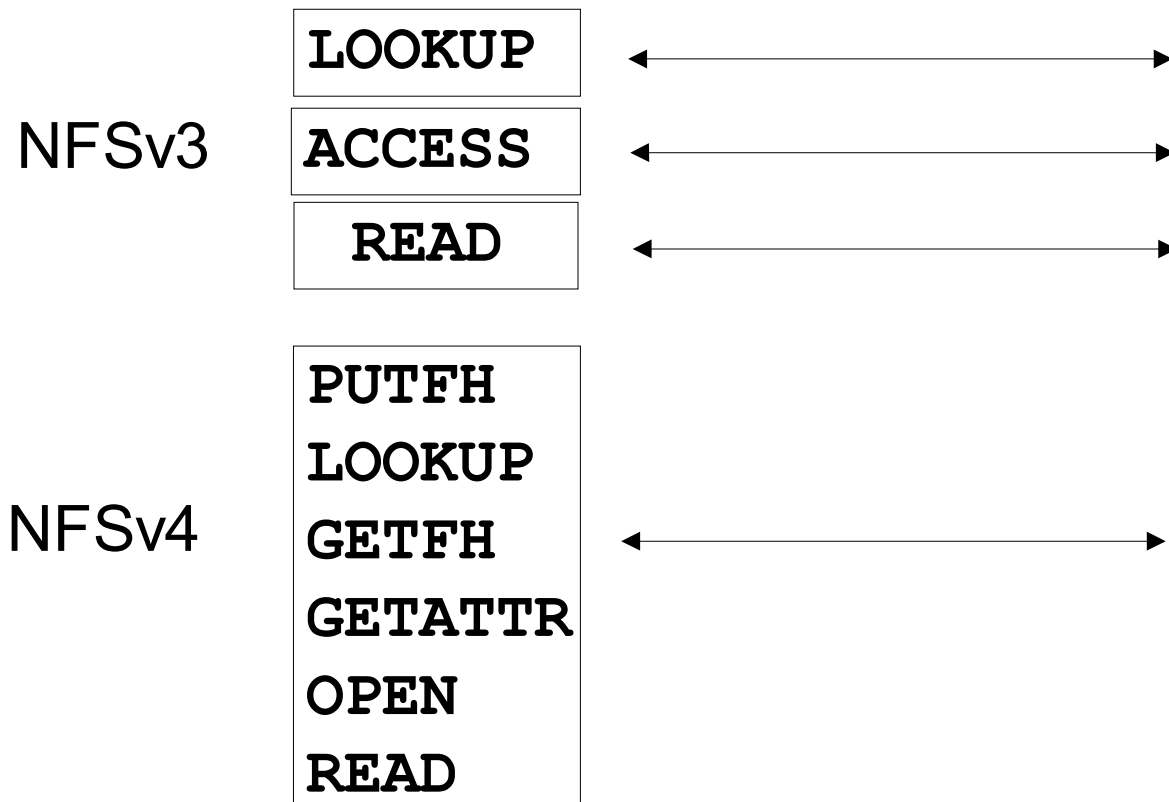
COMPOUND Procedure

- Group related operations in one RPC
- Evaluation stops at first error
- Reduce latency with fewer roundtrips
- Flexibility for various client environments



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

RPC vs. COMPOUND





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

Namespace

- Replaces use of MOUNT protocol
- Server provides access to filesystems from a “root filehandle”
- Server *pseudofs* joins exported subtrees with a read only virtual file system
- Client traverses *pseudofs* with LOOKUP

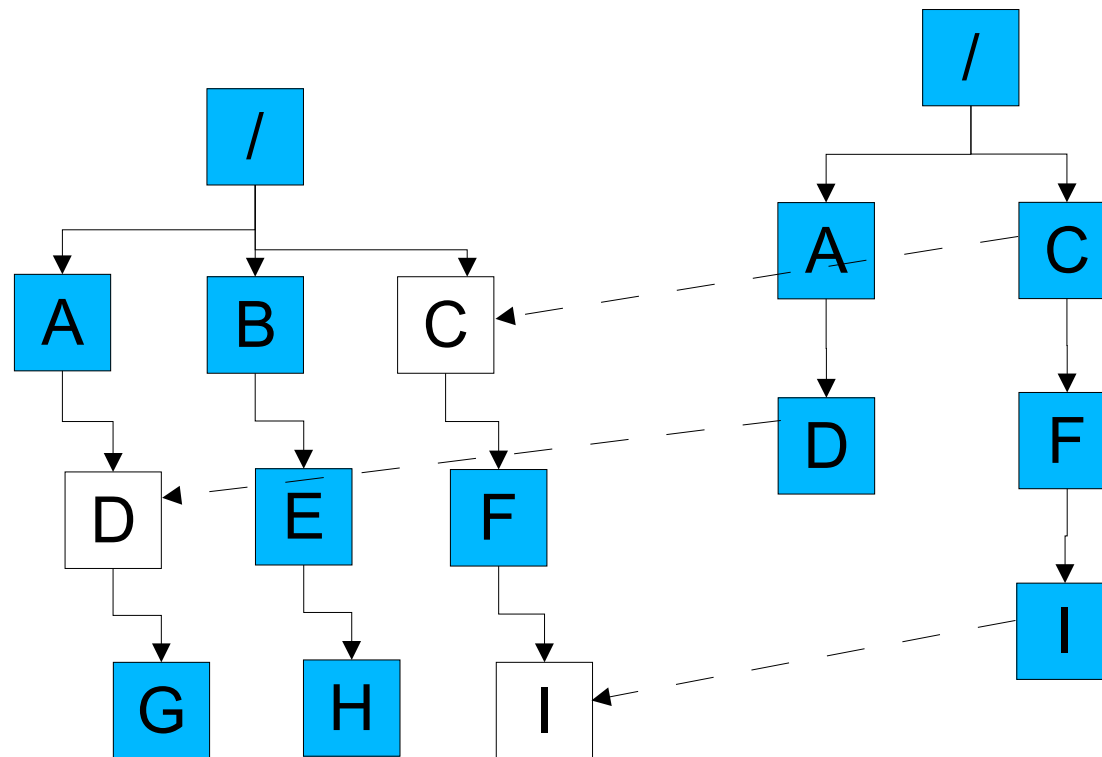


**N I C
F N O
S D N
U F
S E
T R
Y N
G E**

Building Pseudofs

Local FS

Pseudofs





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

Client Namespace

- Client starts at root filehandle of server
- Client inspects *fsid* attribute to determine when new filesystem is found
- At each new filesystem, client automatically mounts filesystem into its namespace

Security



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

- RPCSEC_GSS framework is basis for various security mechanisms
- Provides for Authentication, Integrity, and Privacy
- Kerberos V5 and SPKM/LIPKEY



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

Security Negotiation

- Policy set at each filesystem
- Access root filehandle with secure channel
- If client mismatches on security mechanism, server returns an NFS4ERR_WRONGSEC error.
- Client will use SECINFO operation to enumerate available mechanisms



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

Filehandle Types

- Persistent filehandles same as NFSv2/v3
- Volatile filehandles are new
 - Filesystem types like FAT or user-level server implementations may use
 - Filehandle may become invalid and *expire*
 - Different than traditional ESTALE
 - Expiration at server restart
 - Increases implementation burden on client



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

Volatile Filehandles

- Upon expiration, client attempts recovery
- At initial LOOKUP, client saves path of file
- Pathname used at recovery; client traverses pathname to find new filehandle
- Other expiration events may include RENAME or migration of filesystem



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

NFSv4 is Stateful



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

NFSv4's State

- Hierarchy of NFSv4 state begins with association between a single client and server
- This hierarchy is important because the LEASE period and recovery represents all client-server state
- NFSv4 supports many definitions of client



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

Lease Management

- Lease timeouts used to manage recovery
- Server determines lease period
- Period is for all state generated by client
- Lease renewal occurs at explicit RENEW or by any operation that uses *stateid*
 - CLOSE, DELEGRETURN, LOCK, LOCKU, OPEN, OPEN_CONFIRM, READ, RENEW, SETATTR, WRITE



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

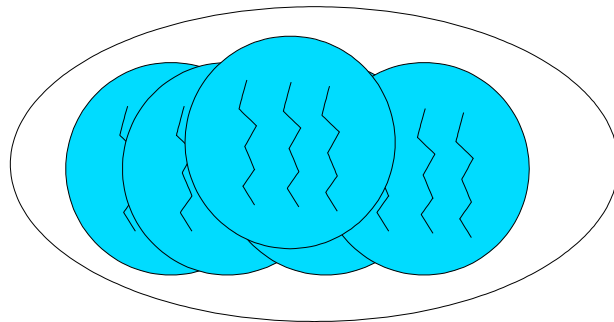
Lease Management (con't)

- Lease timeouts adds to complexity
- Client may track implicit lease renewal
- To ease implementation - use RENEW
- Server management of client state
 - Client state may be released at lease timeout
 - Must be released if lease expired and conflict with state (additional overhead)



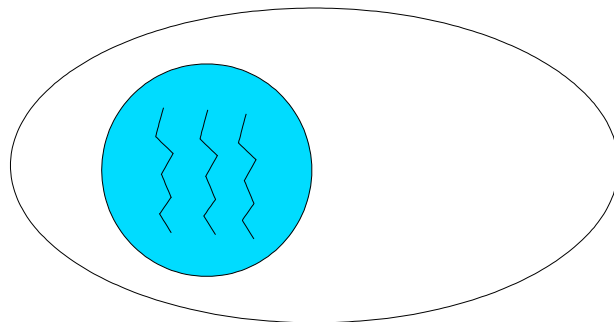
**N I C
F N O
S D N
U F
S E
T R
R E
Y N
C
E**

CLIENTID – may represent these entities at the NFS client



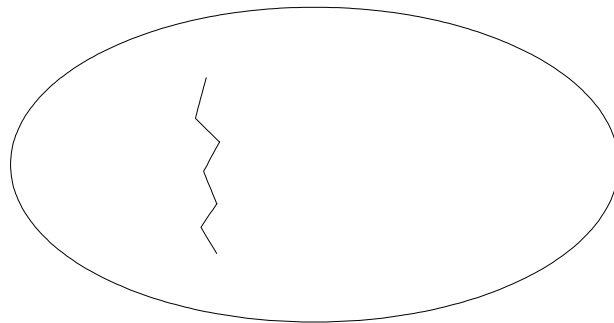
Traditional Unix NFS Client

(M-Processes / N-Threads)
(Y-Users/Principals)



User Level Client

(1-Processes / N-Threads)



Individual Thread



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

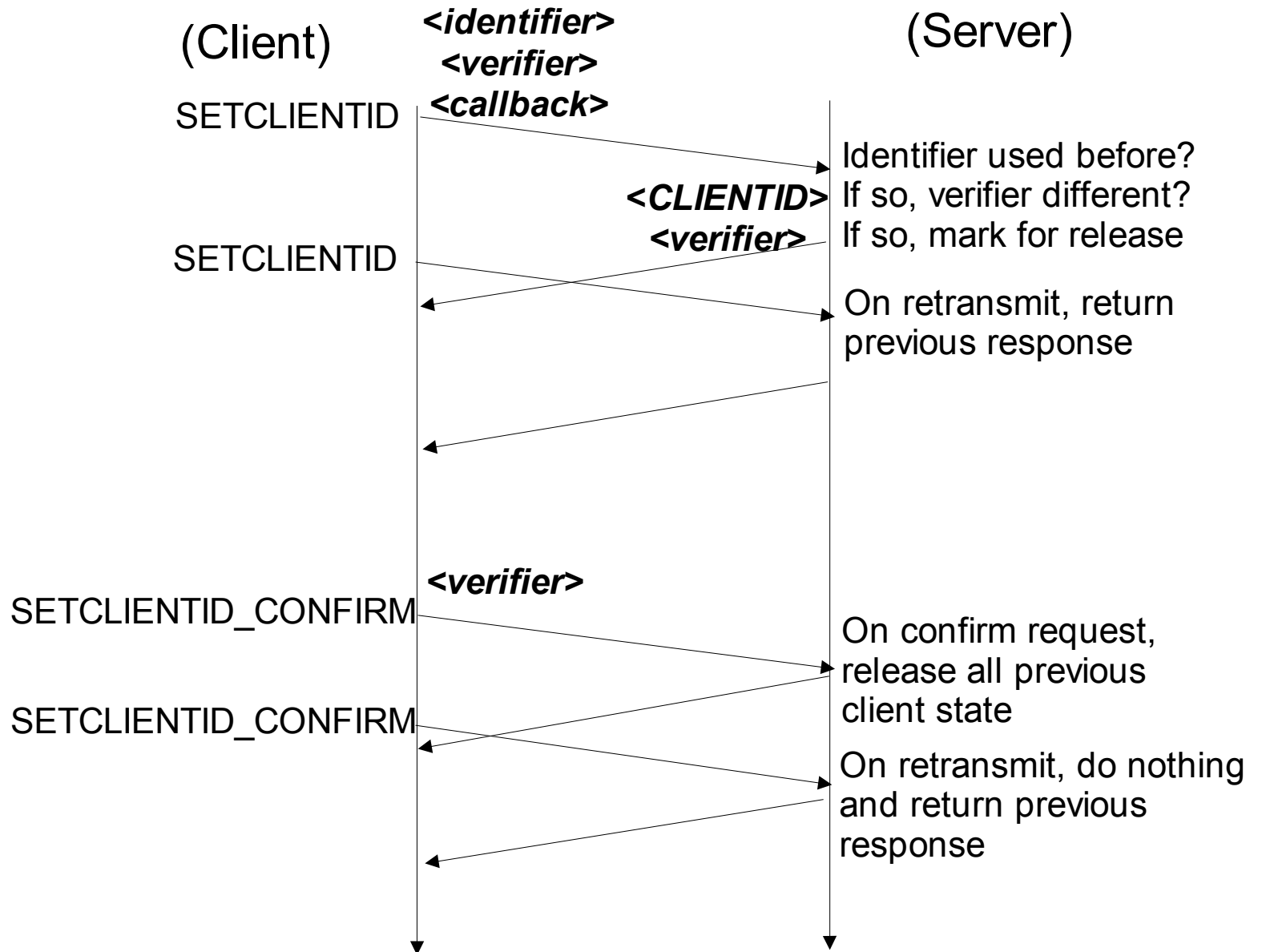
Creation of CLIENTID

- With SETCLIENTID, client chooses an opaque identifier and a verifier
- Uniquely identifies client and client instance
- Server assigns a *CLIENTID*
- Client confirms use with SETCLIENTID_CONFIRM
- Server uses RPC authentication to verify request (saves principal for future reference)

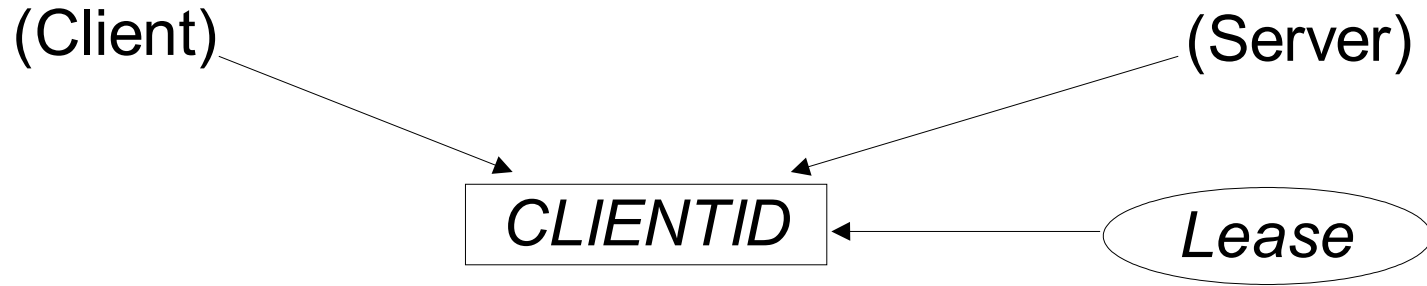


**N I C
F N O
S D N
U F
S E
T R
R E
Y N
C E**

CLIENTID creation



What do we have?



**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 F
S E
T R
R E
Y N
C
E**

OPEN and its state

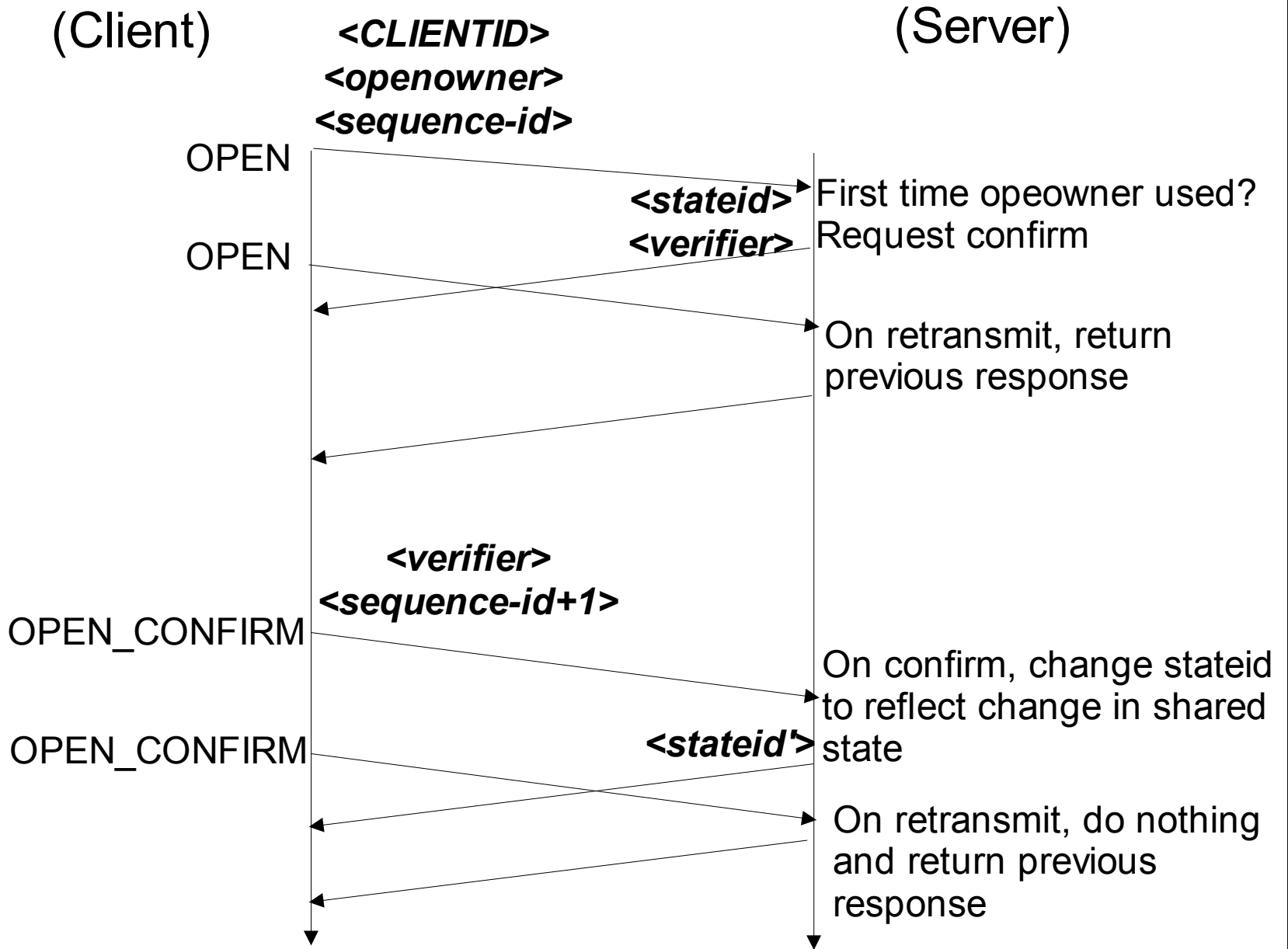
- Combines regular file CREATE, LOOKUP, and share reservations
- Requires name of regular file object
- Identifies owner for OPEN and LOCK state
- Server may provide delegation* in response



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

September 22-24

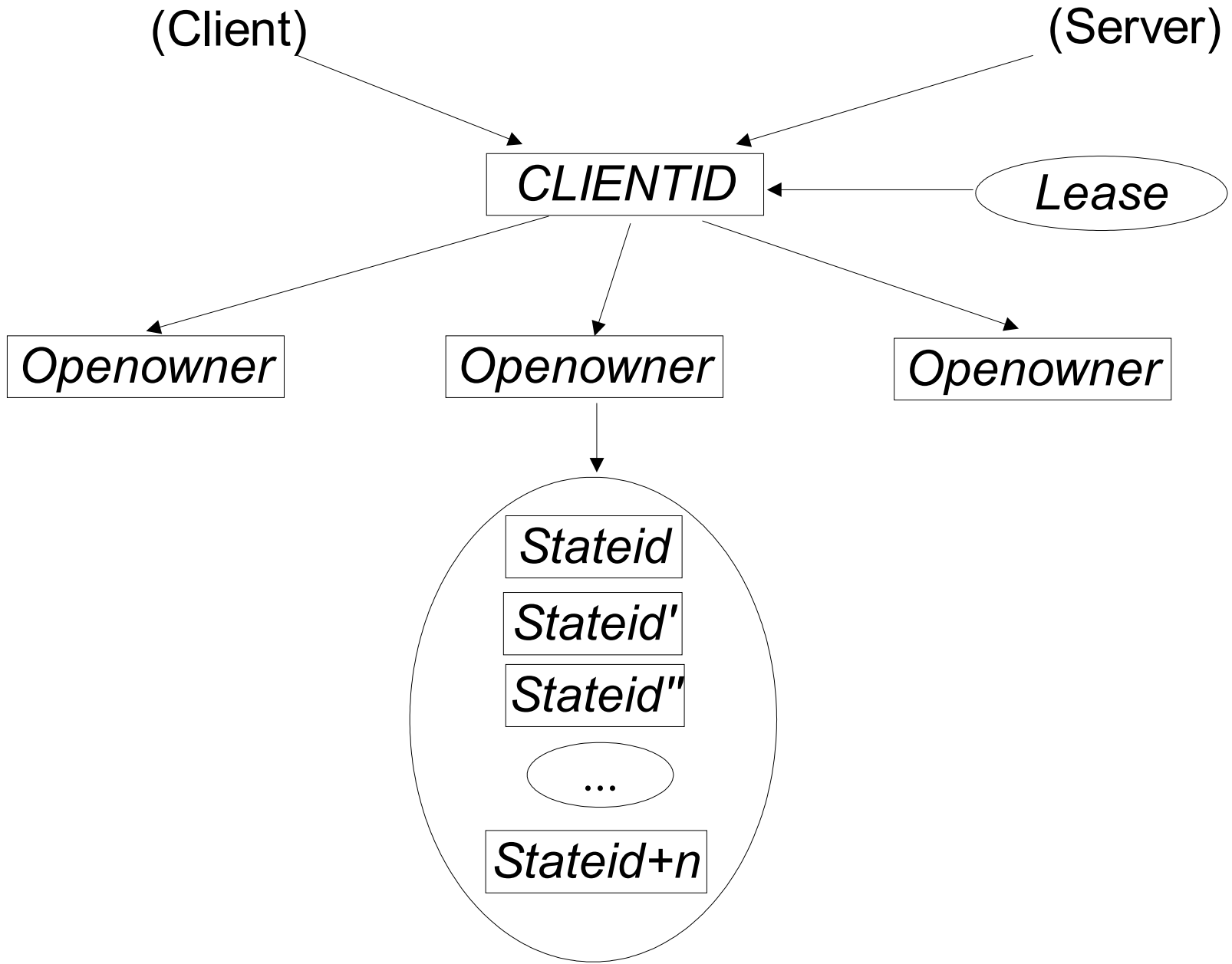
Openowner to openstateid





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

What do we have?





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

Delegation

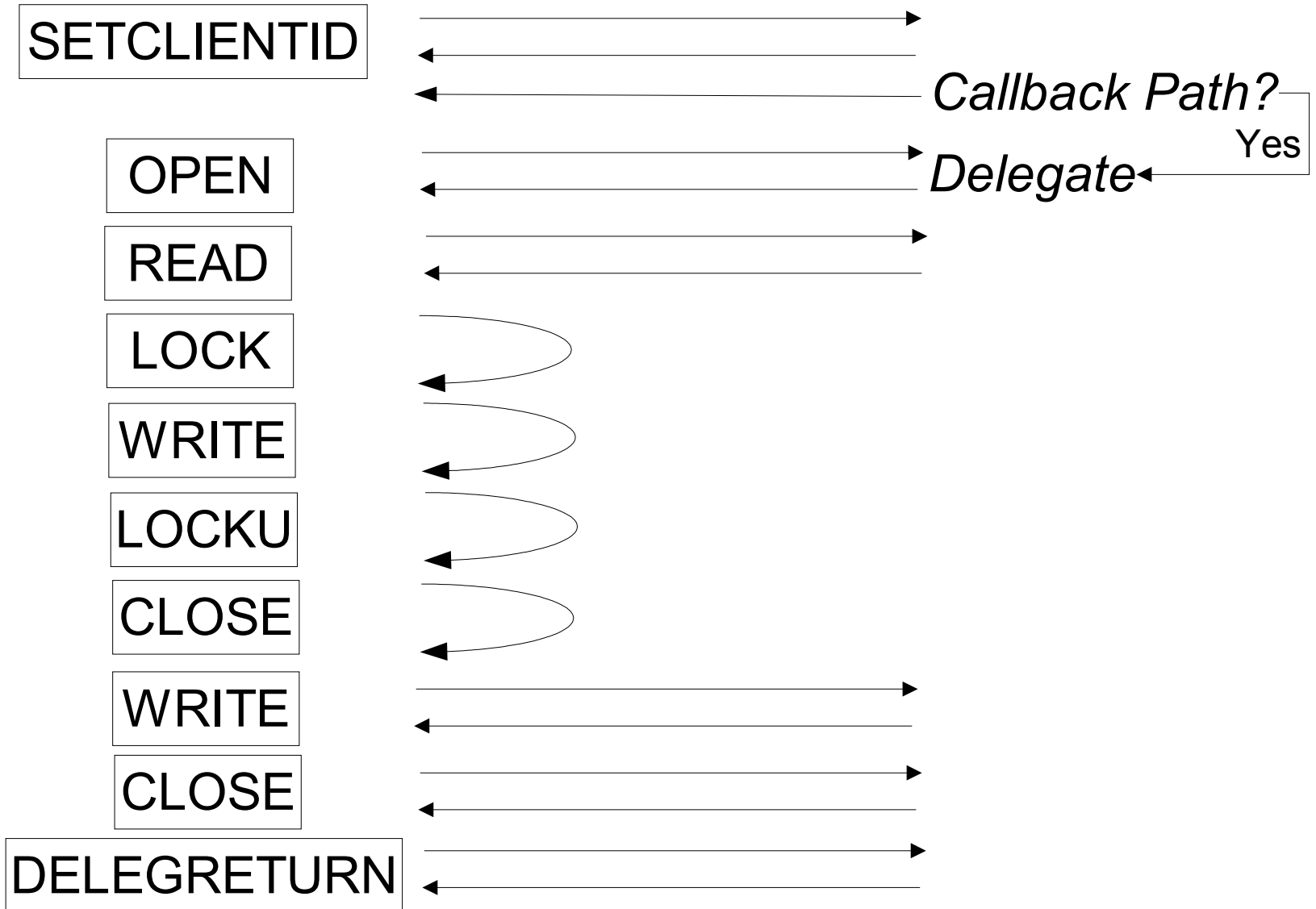
- Intended for minimal sharing environments
- Server decides when to provide delegation
- Not required for correct protocol operation
- Callback path to client must exist
- If delegation provided, client does not need to contact server for further OPEN, LOCK, READ, WRITE, CLOSE



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

Client

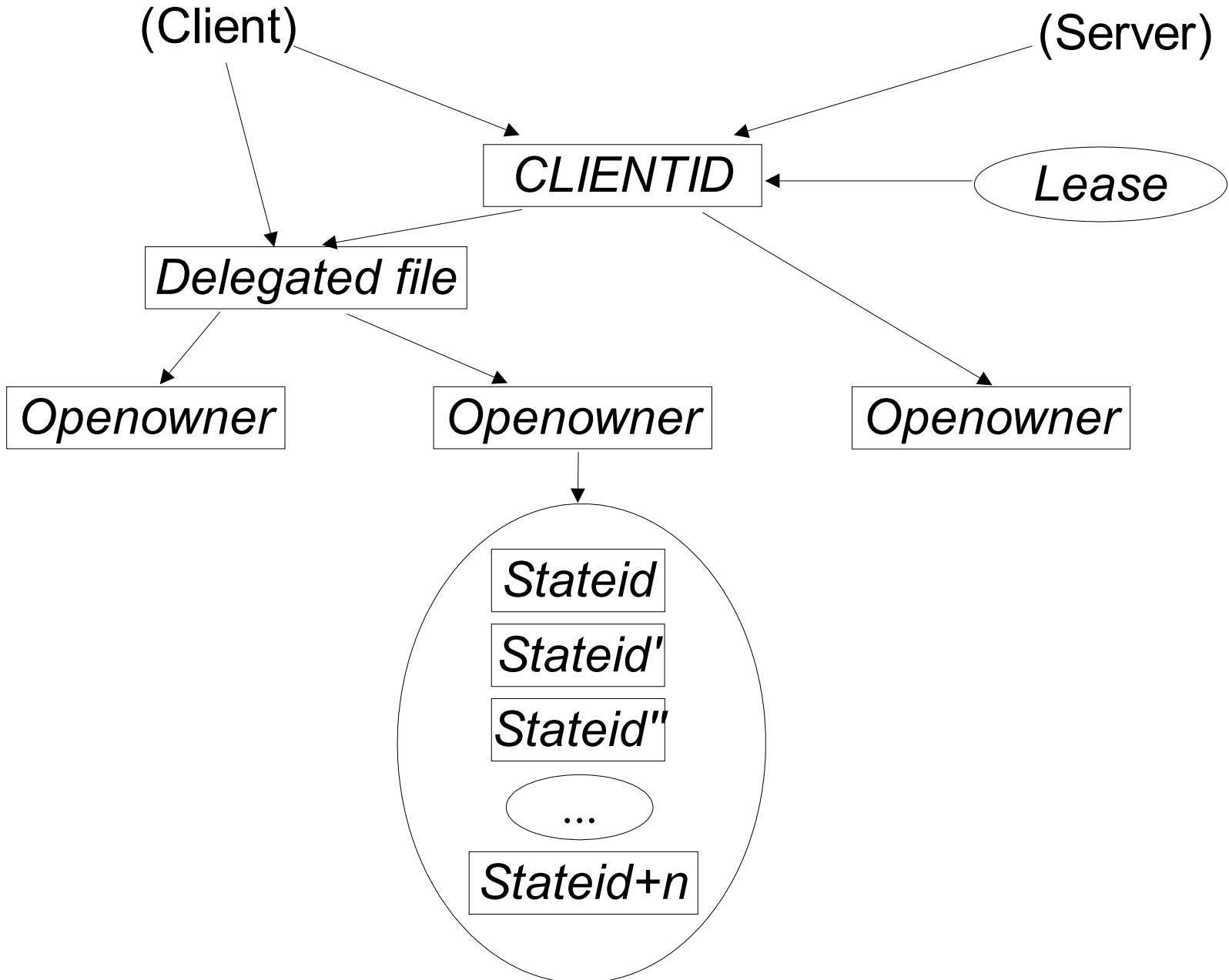
Server





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

What do we have?





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

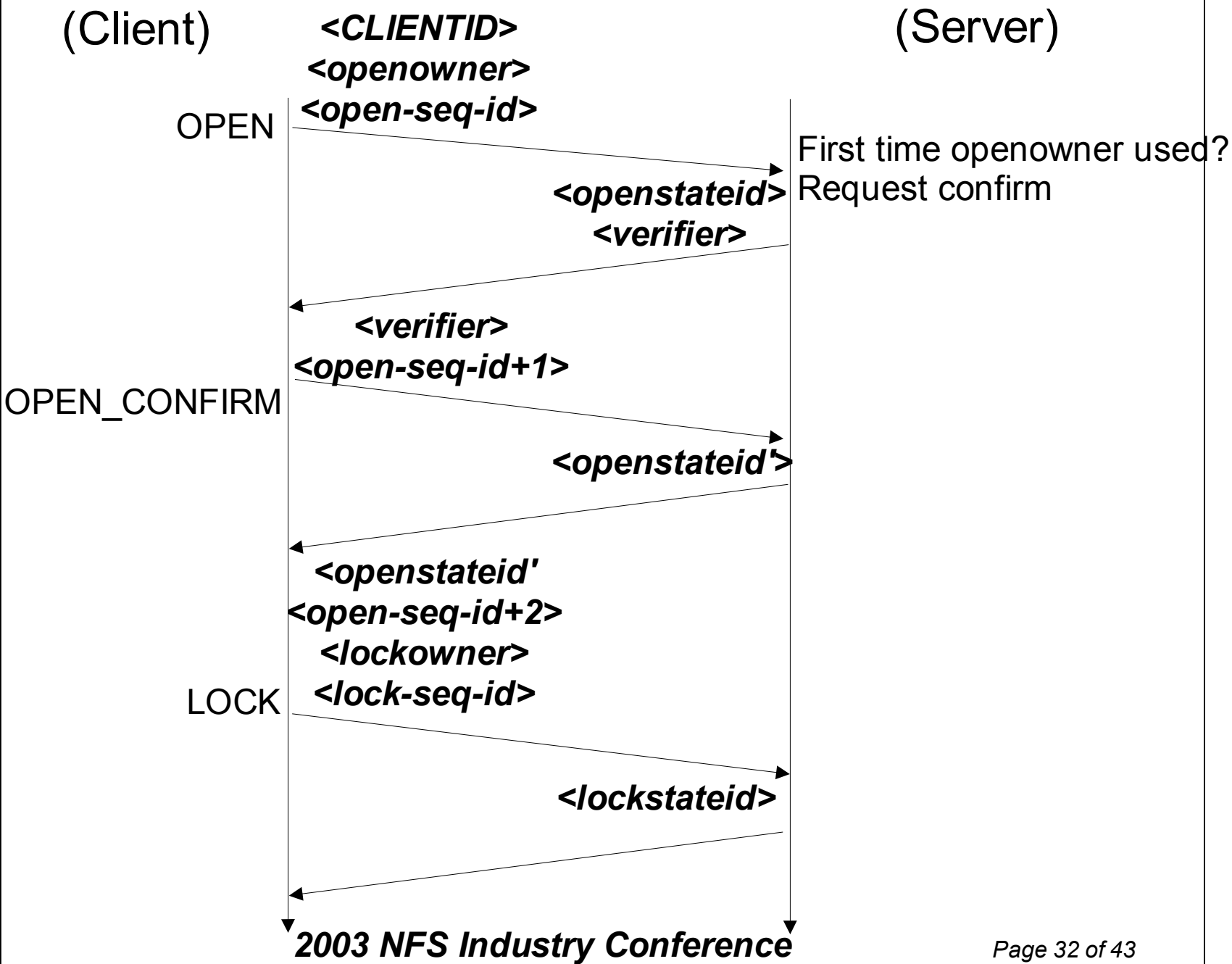
File Locking

- Better implementation a result of integration
- Byte-range locking - non-blocking, mandatory
- Callbacks not used*
- *Stateid* represents transitions of locking state
- *Sequence-id* preserves request ordering
- Locks recovered at lease expiration*
- Client does lock recovery at server restart



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

Openowner / Lockowner association



Now what do we have?

(Client)

(Server)

CLIENTID

Lease

Openowner

Openowner

Openowner

Lockowner

Lockowner

Lockowner

Stateid

Stateid'

...

Stateid+n

Stateid

Stateid'

...

Stateid+n



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

Attributes



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

- Mandatory, Recommended, and Named
- 52 mandatory/recommended attrs
- Extends beyond traditional Unix attrs
- Protocol encoding adds to overhead
- Encoding allows for extensibility
- Mandatory: type, expiration type, change, size, link?, symlink?, fsid, lease duration



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

Named or Extended Attrs

- Support is optional
- OPENATTR operation returns directory filehandle of named attributes
- Named attrs may be recursive
- Intended for application use
- Semantics are opaque to client and server



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

Access Control Lists

- Definition finally in the protocol
- Solves today's problem of 4+ ACL protocols
- ACLs may be manipulated by client
- Windows/NT ACL compatible
- Combined with RPC security provides for strong security mechanism
- POSIX stops at first denial; NT does full evaluation searching for allowance



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

Owner / Group Attributes

- String based identifiers
- Take place of numeric uid / gid
- Identifier looks like: user@domain
- External mapping outside of protocol is not defined



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

Filesystem Replication

- Intended for read-only filesystems
- Increases availability
- Client's policy directs when to switch to another replica
- *fs_locations* attribute enumerates replicas
- Client needs to reconstruct state at new server
- Server to server replication undefined



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

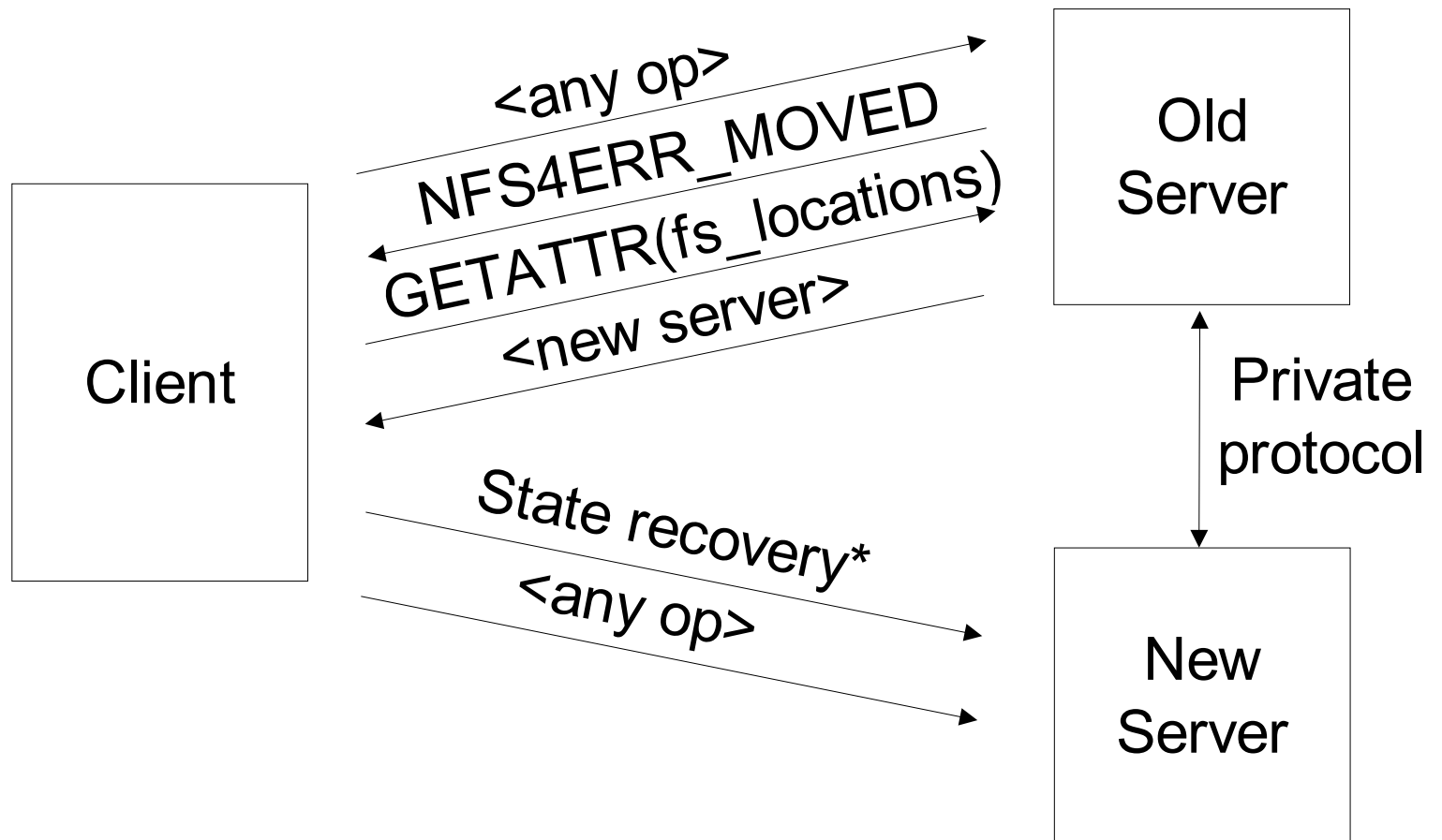
Filesystem Migration

- Enables load balancing or server reorganization
- Server to server transfer is undefined
- Client receives NFS4ERR_MOVED at migration event
- *fs_locations* attribute provides new location



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

Migration Event





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

Minor Versioning

- Difficult to revise previous versions (v2->v3)
- Realize that protocol is not perfect
- Protocol must be allowed to evolve
- Changes allowed in minor version
 - Operations (new argument types) added
 - Add new Attributes



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

Minor Version Examples

- SECINFO fixes
- Delegations for Directories
- Support for RDMA/RDDP



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

NFS Version 4

Spencer Shepler

Sun Microsystems

spencer.shepler@sun.com