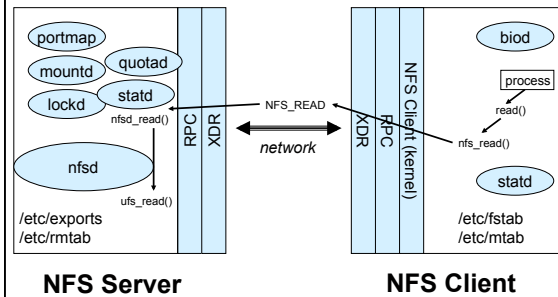


## CSE-595: Storage Systems

### NFS: Protocols, Programming, and Implementation

Erez Zadok  
ezk@cs.sunysb.edu

## The BIG Picture



Spring 2015

CSE595: Storage Systems

## NFS Overview

- using RPC: Remote Procedure Calls
  - ◆ which use XDR: eXternal Data Representation
- stateless server
  - ◆ crash recovery
- client side caching (data and attributes)
  - ◆ request retransmission
- file handles: 32 bytes opaque to client
  - ◆ server encodes: fsid, inum, igen, possibly more

Spring 2015

CSE595: Storage Systems

## XDR: eXternal Data Representation

- de/serializes data into network-order bytes
- repeated calls encode/decode more "XDR" bytes

```
bool_t xdr_long(XDR *xdrs, long *lp);
struct foo {
    int i;
    char *buf;
};
bool_t xdr_foo(XDR *xdrs, struct foo *foop) {
    if (!xdr_int(xdrs, &foop->i))
        return FALSE;
    if (!xdr_wrapstring(xdrs, &foop->buf))
        return FALSE;
    return TRUE;
}
```

Spring 2015

CSE595: Storage Systems

## RPC: Remote Procedure Call

- server does:
 

```
registerrpc(prognum, versum, procnum, s_inproc, in, s_outproc, out);
svc_run();
```
  - client issues:
 

```
callrpc(char *host, rpcprog_t prognum, rpcvers_t versum,
         rpcproc_t procnum, xdrproc_t inproc, char *in,
         xdrproc_t outproc, char *out);
```
  - which contacts server's portmapper, then RPC server w/ procnum.
- when client request comes
- ◆ find procnum
  - ◆ call s\_inproc to decode client args
  - ◆ call s\_outproc to encode output to client
  - ◆ return => client returns (or times out)
- rpcgen produces headers and .c stubs from .x files

Spring 2015

CSE595: Storage Systems

## Additional NFS Components

- on server:
  - ◆ mountd:
    - + listen for mount requests
    - + authenticate requests
    - + return root handles
- on client:
  - ◆ biod: dirty page clustering, simulate async writes
- on both:
  - ◆ lockd: coordinates local/remote record locks
    - + flock() uses lockd; lockf() only local locks;fcntl() can use both
  - ◆ statd: synchronizes lock information
    - + client reboot: tell server to release locks
    - + server reboot: tell all clients to reclaim locks
  - ◆ portmapper: the mother of all RPC servers

Spring 2015

CSE595: Storage Systems

### Example: mounting a remote server

- get handle (via MOUNTPROC\_MNT rpc to mountd)
- fill in struct nfs\_args
  - ◆ struct nfs\_args na
- call mount(2) syscall
  - ◆ mount("/mnt", flags, "nfs", &na, sizeof(na))

Spring 2015

CSE595: Storage Systems

### Contents of struct nfs\_args

```

NA->addr {sockaddr_in} (len=16) =      NA->version = 3
"02000801803b14640000000000000000"  NA->flags = 0x0
NA->addr.sin_family = "2"             NA->rsize = 4096
NA->addr.sin_port = "264"            NA->wsize = 4096
NA->addr.sin_addr = "803b1464"       NA->bsize = 0
NA->hostname = "opus"                NA->timeo = 7
NA->namlen = 255                     NA->retrms = 3
NA->filehandle =                     NA->acregmin = 3
"008000f400000002000a000000000026e NA->acregmax = 60
065b6c000a0000000000026e065b6c"   NA->acdirmin = 30
                                       NA->acdirmax = 60
    
```

Spring 2015

CSE595: Storage Systems

### NFS V.2 (1984)

- Built on top of UDP
  - 17 calls
- |                       |   |                    |    |
|-----------------------|---|--------------------|----|
| <b>NFS_NULL</b>       | 0 | <b>NFS_CREATE</b>  | 9  |
| <b>NFS_GETATTR</b>    | 1 | <b>NFS_REMOVE</b>  | 10 |
| <b>NFS_SETATTR</b>    | 2 | <b>NFS_RENAME</b>  | 11 |
| <b>NFS_ROOT</b>       | 3 | <b>NFS_LINK</b>    | 12 |
| <b>NFS_LOOKUP</b>     | 4 | <b>NFS_SYMLINK</b> | 13 |
| <b>NFS_READLINK</b>   | 5 | <b>NFS_MKDIR</b>   | 14 |
| <b>NFS_READ</b>       | 6 | <b>NFS_RMDIR</b>   | 15 |
| <b>NFS_WRITECACHE</b> | 7 | <b>NFS_READDIR</b> | 16 |
| <b>NFS_WRITE</b>      | 8 | <b>NFS_STATFS</b>  | 17 |
- (why no lseek?)*

Spring 2015

CSE595: Storage Systems

### Ex: NFS\_READ Call

```

struct readargs {
    fhandle file;
    unsigned offset;
    unsigned count;
    unsigned totalcount;
};

union readres switch (stat status) {
    case NFS_OK:
        fattr attributes;
        nfsdata data;
    default:
        void;
};
    
```

Spring 2015

CSE595: Storage Systems

### NFS V.3 (1994)

- TCP and UDP
- 64 byte file handles
- files > 2GB
- ACLs supported
- Kerberos authentication type
- All ops return old/new attributes
  - ◆ saves on most popular call, getattr (update client caches faster)

Spring 2015

CSE595: Storage Systems

### NFS V.3 Protocol

- Removed: ROOT and WRITECACHE
- Added:
  - ◆ **REaddirPLUS:** 17
    - ✦ also returns file handles
    - ✦ saves on NFS\_LOOKUPS
  - ◆ **FSSTAT:** 18
  - ◆ **FSINFO:** 19
  - ◆ **PATHCONF:** 20
  - ◆ **COMMIT:** 21
    - ✦ Saves cached data to disk

Spring 2015

CSE595: Storage Systems

## NFS V.4 (2004)

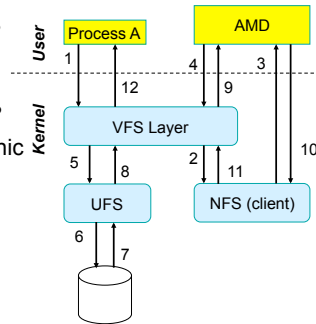
- IETF design, not Sun
- Integrated file locking and mount protocol
- Stronger security w/ negotiation
  - ◆ Public file handles
  - ◆ Works with firewalls & proxies
- Compound operations
- Internationalization
- Better suited for Internet (i.e., WAN)
- Migration and replication
- Extensible protocol

Spring 2015

CSE595: Storage Systems

## User Level NFS-Based File Servers

- Context switches
- extra communication
- Amd dead/hung?
- CFS: cryptographic file server
- FUSE like



Spring 2015

CSE595: Storage Systems

## Resources

- RFC 1094/1813
  - ◆ Usenix papers [Sandberg 84] and [Pawlowski 94]
- NFS V.2/3/4 specs and drafts
  - ◆ ietf.org

Spring 2015

CSE595: Storage Systems