# Client Discovery of NFSv4.1 Server Multipath Addresses

**NetApp**®

Andy Adamson

andros@netapp.com

IETF 96, Berlin, 2016

# Motivation

- Motivated by desire for session trunking
  - Linux prototype code

- pNFS data servers have the GETDEVICEINFO multipath4 to list potential session trunking addresses

- No such protocol feature for non-pNFS or MDS NFSv4.1+ servers

- How does the client discover a list of potential session trunking addresses?

- Meeting topic at June 2016 NFSv4 Bakeathon

**NetApp**

# Multipath Options for Session Trunking

o **IP Layer**
  o MPTCP RFC 6820
  o SCPT RFC 4960

o **RPC Layer**
  o Linux feature: multiple RPC transports per mount

o **All take advantage of multiple network paths between the client and the server to:**
  o Fully utilize network resource
  o Achieve better throughput
  o Failover for network failures (HA)

**NetApp**

# MPTCP

o  Client and server RPC layer sees a single TCP connection which is multiplexed

  o  NFS server can not assign resources per multiplex TCP, but is this really needed…

  o  Automatic 'session trunking'

  o  No support for other transports such as RDMA

o  Detects existence of multiple network interfaces on the hosts and creates the multiple TCP flows

  o No need for additional client discovery of addresses

  o New features in discussion to respond to topology

    o Load balancing, path priority to name a few

**NetApp**

# SCTP

- o **SCTP association has multiple IP addresses**
  - o Originally for HA, new load sharing for performance.

- o **No specification for SCTP in ONC RPC**
  - o Will the NFS server see one SCTP association or the potentially multiple connections within an association?
  - o No support for other transports such as RDMA

- o **Requires client discovery of server multipath**
  - o SCTP association is defined as [a set of IP addresses at A]+[Port-A]+[a set of IP addresses at Z]+[Port-Z]

- o **Still a new protocol that needs work**
  - o Performance only slightly better than TCP

**NetApp**

# Linux Multiple RPC transports

- New Linux RPC feature allows for multiple transports per mount
  - Supports a mix of transports (TCP, MPTCP, RDMA, SCTP, …)

- RPC layer controls use
  - Currently only round-robin algorithm implemented
  - Needs work: network topology, which path to prefer

- Server sees each connection as separate

- Requires client discovery of transport addresses

NetApp

# Client Discovery: Multipath Addresses

o Multiple host names on mount command

o Use a DNS A-record, or special trunk record

o Use fs_locations or fs_locations_info attribute on the pseudo file system

NetApp

# Multiple host names on mount

o ## Solaris supports this feature
  - o Disliked due to the need to address all clients when when network conditions change. (as reported at June 2016 Bakeathon)

o ## Discussion on Linux kernel list
  - o Linux RFC prototype received a mixed response
  - o Support issues influenced the decision to not implement multiple host names mount feature
    - o Response to change in network conditions
    - o Dislike another mount option to support
    - o Server should supply addresses

**NetApp**

# DNS for Multipath

o Use DNS A-records, or special trunk record

o Changing a DNS record due to an interface change on a server is problematic

   o Making a DNS change in many organizations is difficult and takes a lot of time as making a mistake that brings DNS down stops all computing services.

o DNS caching means that when a change occurs, need to wait the cache timeout, typically an hour.

NetApp

# fs_locations

o **fs_locations/_info useful replica definition:**

  o "When a set of servers have corresponding file systems at the same path within their namespaces, an array of server names may be provided."

  o Client could test each replica address for session trunking using the EXCHANGE_ID test.

    o Most addresses will be to different replica servers

o **fs_locations/_info is per file system**

  o session trunking is server wide.

NetApp

# fs_locations proposal

- Use an fs_locations or fs_locations_info replica attribute on the pseudo file system
  - pseudo file system does not migrate and is not replicated (I guess it could be - but what is the point?

- Define this to apply to the whole server

- Client test each address for session trunking with the EXCHANGE_ID test
  - All addresses will be from the same server

**NetApp**

# Summary

- Do we need to specify a means of client discovery of NFSv4.1 server multipath addresses?

- MPTCP gets us most of the way for TCP
  - Protocol development is active

- Is fs_locations/_info on pseudo file system a reasonable approach?

**NetApp**

# Thank you

# Multipath Trunk Options

o Testing has shown that multiple connections
  - o Improves performance
  - o Improves availability

o Testing:

o MPTCP
  - o http://multipath-tcp.org/pmwiki.php?n=Main.50Gbps

o Prototype Linux NFS client with multiple RPC transports

o SCTP
  - o https://www.researchgate.net/publication/220776215_Perform ance_Comparison_of_SCTP_and_TCP_over_Linux_Platform

o All trunk options require some sort of discovery of multipath addresses

NetApp