# OpenSFS Project
# Lustre SMP Node Affinity

Liang.zhen@intel.com

Aug, 28  2012
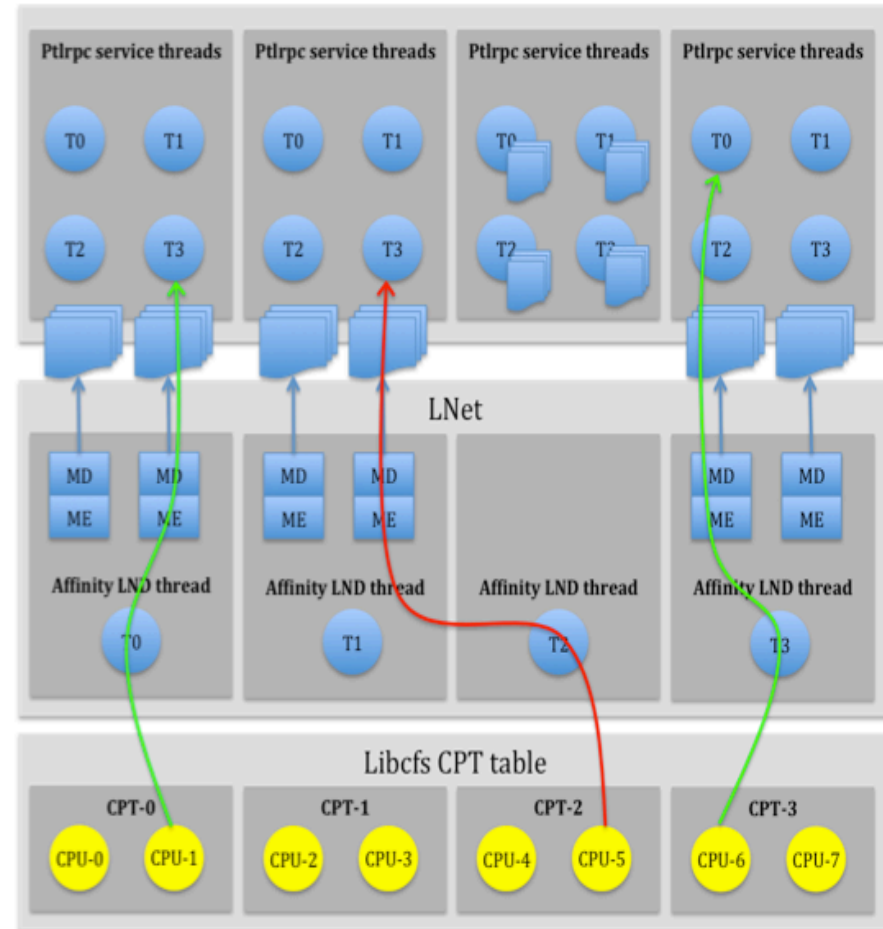
# Agenda

- Background

- Demonstration

- Tuning Lustre on SMP machine

**Intel Architecture Group** (intel)

# Background

- Goal of this project
  - Improve SMP scalability of LNet
  - Improve metadata performance for single MDS
  - Funded by OpenSFS

- Code landed to 2.3
  - 16K+ LOC

Intel Architecture Group (intel)

# Partitioned Lustre Server

- CPU Partition (CPT)
  - Similar to cpuset of linux
  - Can be easily used by kernel thread

- Partitioned LNet(LND)
  - LND thread-pool for each CPT
  - Core LNet has partition data

- Partitioned ptlrpc service
  - Ptlrpc service thread-pool for each CPT
  - Request-queue & wait-queue for each CPT
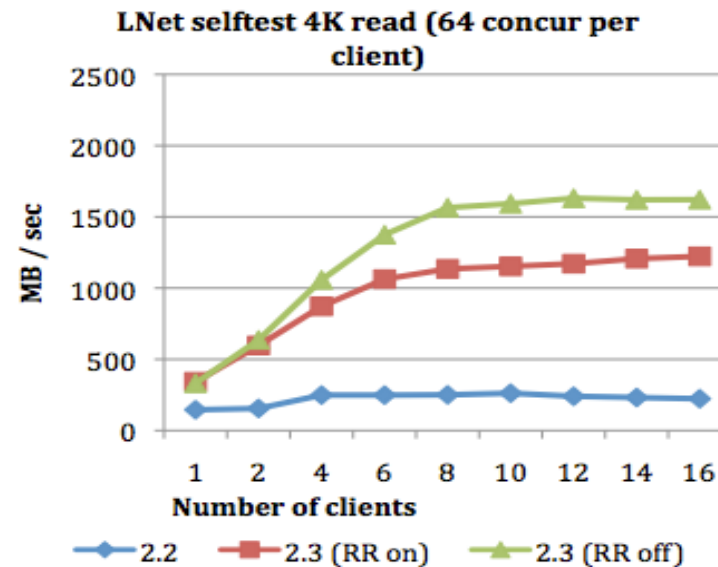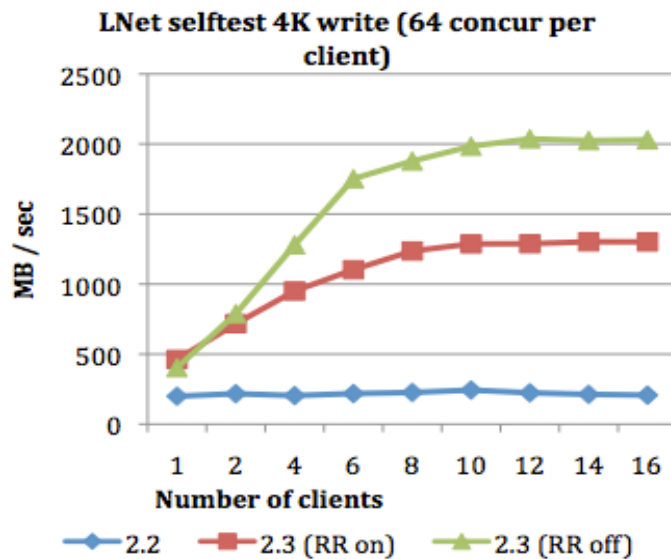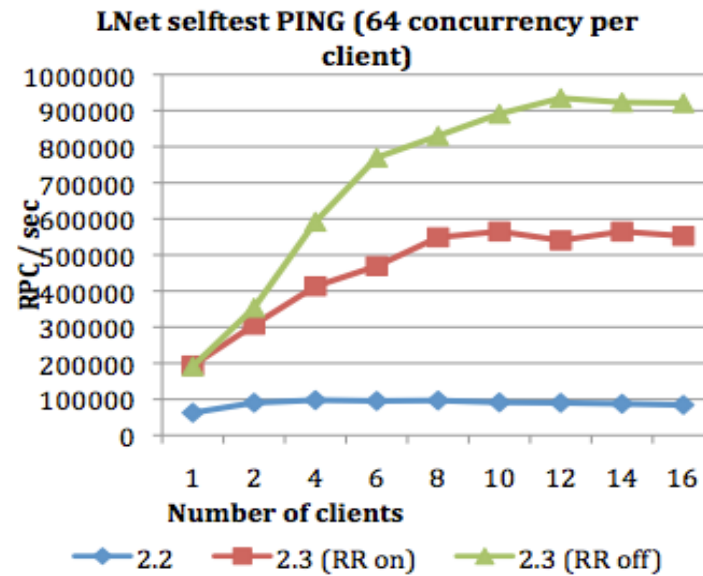
Intel Architecture Group (intel)

# LNet performance tests

- Hardware
  - Server: 6-core CPU (2-HT), 2 sockets
  - Client: 4-core, 1 socket
  - QDR infiniband

- LNet selftest
  - Selftest ping
  - Selftest 4K read/write
  - Concurrency

- Portal Round-Robin (Portal RR)
  - NID affinity in LNet (LND)
  - Enable/disable NID affinity of incoming message for upper layer (ptlrpc service, or LNet selftest)

# LNet performance

- 2.3 ping is 900% of 2.2 with Portal-RR OFF
- 2.3 ping is 600% of 2.2 with Portal-RR ON
- 2.3 4K-BRW is 600%-700% of 2.2 with Portal RR OFF
- 2.3 4K-BRW is 500% of 2.2 with Portal RR ON



LNet selftest PING (64 concurrency per client)



LNet selftest 4K write (64 concur per client)



LNet selftest 4K read (64 concur per client)

Intel Architecture Group   (intel)
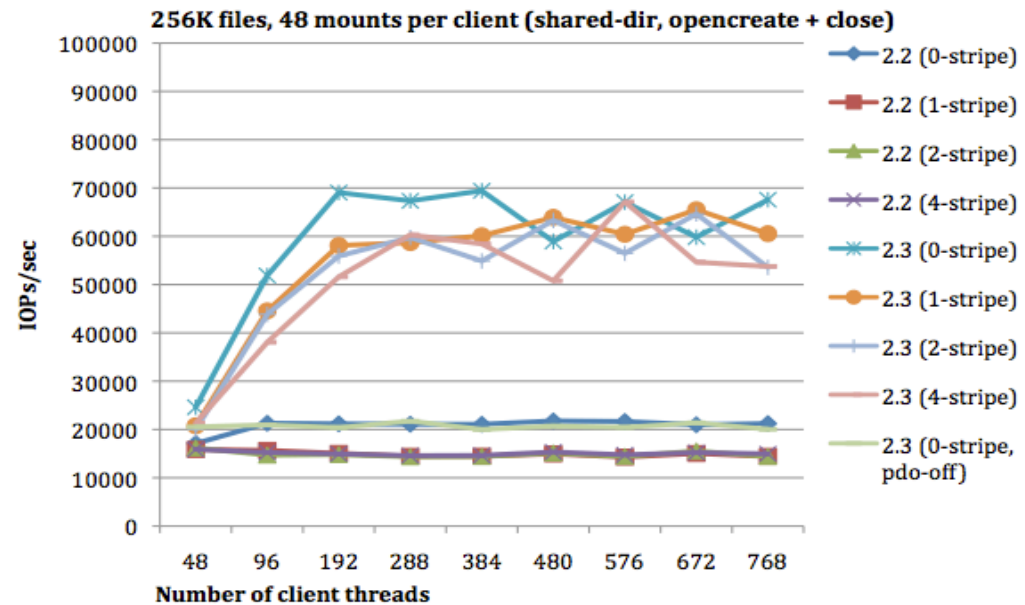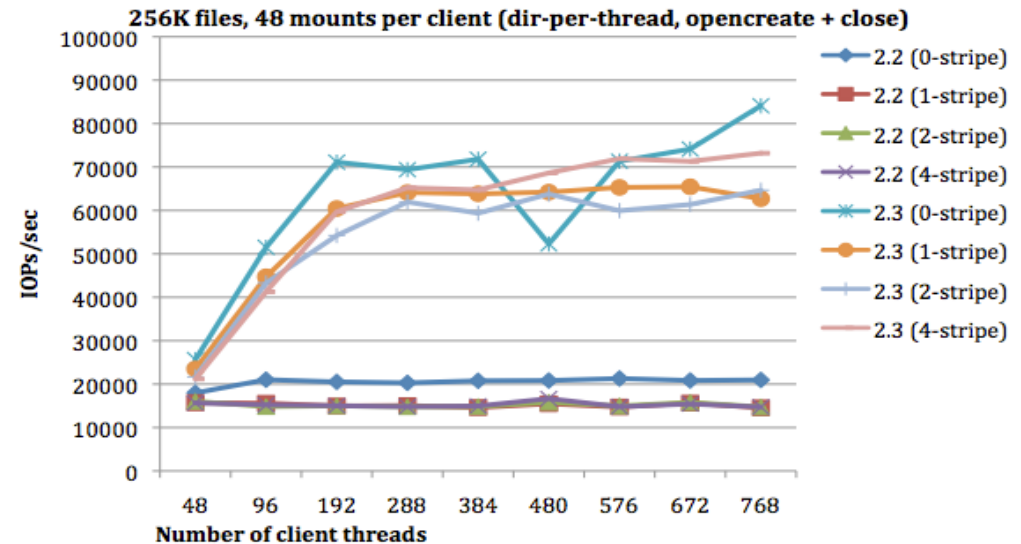
# mdtest

- Hardware
  - MDS
    - 6-core CPU (2-HT), 2 sockets
    - 8G SSD as MDT journal
  - OSS
    - 3 OSSs, 6 OSTs per OSS
  - Client: 4-core, 1 socket
  - QDR infiniband

- Mdtest patches
  - multi-mount
    - Simulate high work load with small number of clients
    - Disable mdc_rpc_lock can't help shared directory tests
  - 0-stripecount file
    - w/o OST object creation
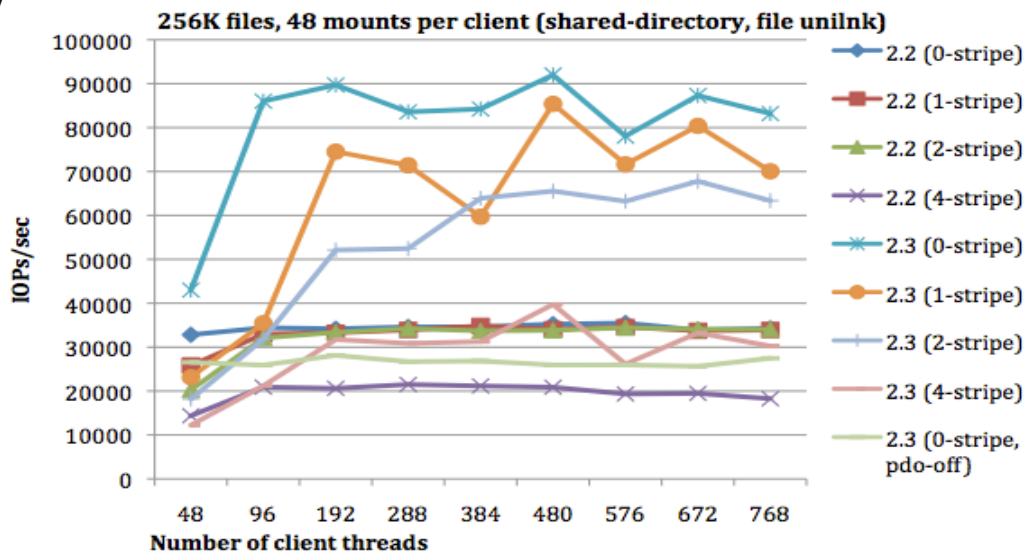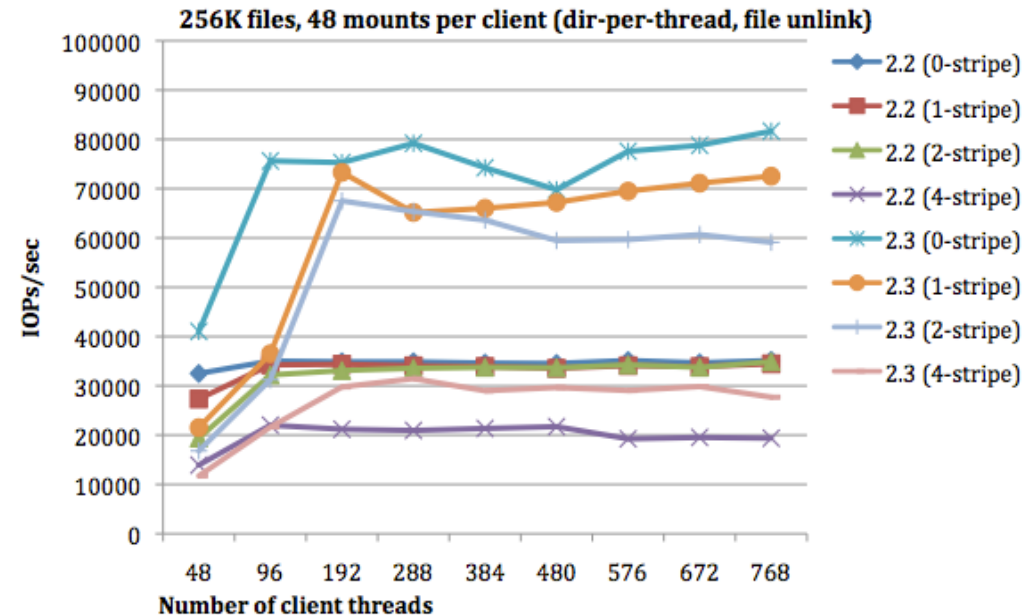
**Intel Architecture Group** (intel)

# File creation performance

- Iterate over 1,2, 6, 4, 8,10, 12,14, 16 clients
- Each client has 48 threads
- Each thread is running under a private mount
- 2.3 opencreate performance is 350%-400% of 2.2
- OST object pre-creation works pretty good
- Turning off PDO, shared directory opencreate performance of 2.3 is similar to 2.2



256K files, 48 mounts per client (dir-per-thread, opencreate + close)

- 2.2 (0-stripe)
- 2.2 (1-stripe)
- 2.2 (2-stripe)
- 2.2 (4-stripe)
- 2.3 (0-stripe)
- 2.3 (1-stripe)
- 2.3 (2-stripe)
- 2.3 (4-stripe)



256K files, 48 mounts per client (shared-dir, opencreate + close)

- 2.2 (0-stripe)
- 2.2 (1-stripe)
- 2.2 (2-stripe)
- 2.2 (4-stripe)
- 2.3 (0-stripe)
- 2.3 (1-stripe)
- 2.3 (2-stripe)
- 2.3 (4-stripe)
- 2.3 (0-stripe, pdo-off)
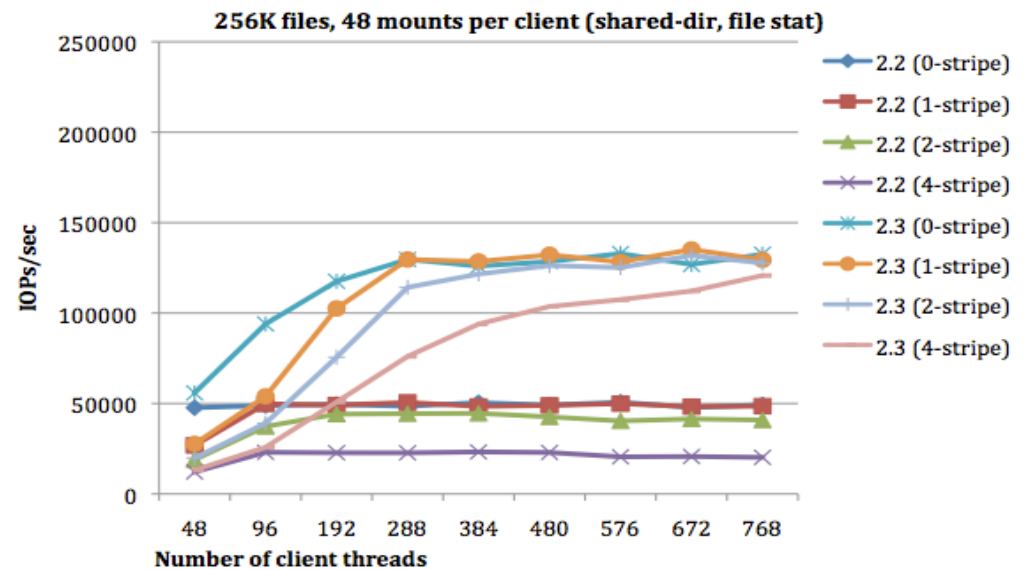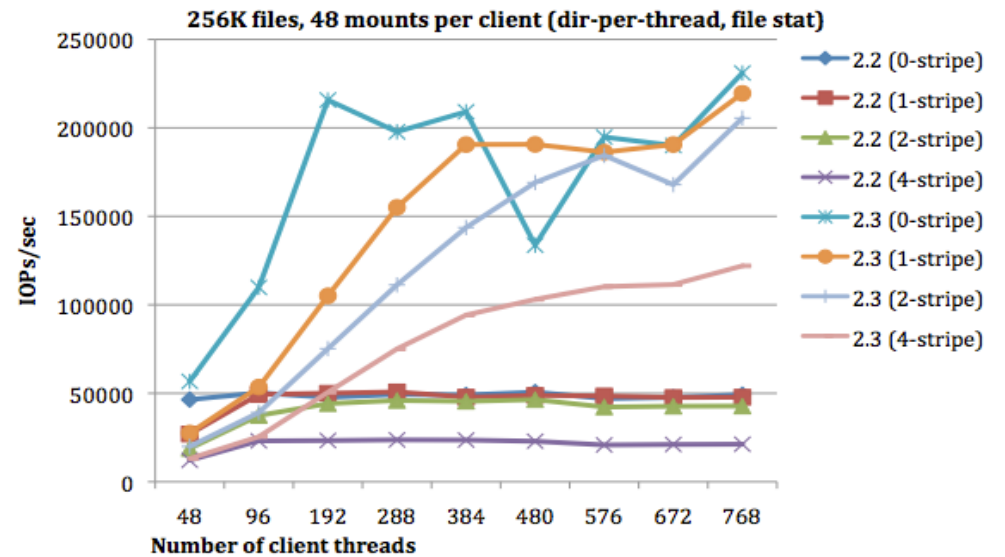
Intel Architecture Group (intel)
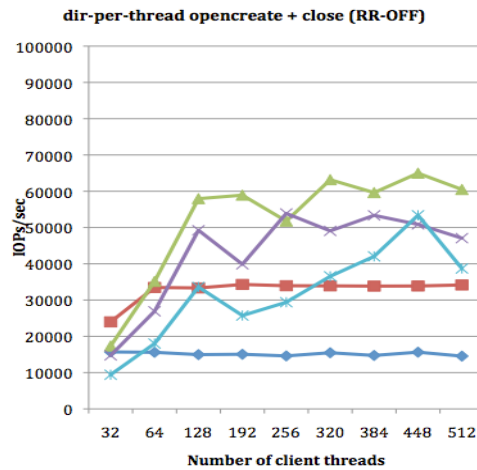
# File unlink performance

- Iterate over 1,2, 6, 4, 8,10, 12,14, 16 clients
- Each client has 48 threads
- Each thread is running under a private mount
- 2.3 unlink performance is 150%-300% of 2.2
- Client needs to send RPC to destroy each OST object
- Turning off PDO, shared directory opencreate performance of 2.3 is even worse than 2.2



256K files, 48 mounts per client (dir-per-thread, file unlink)



256K files, 48 mounts per client (shared-directory, file unilnk)

Intel Architecture Group (intel)

# File stat performance

- Iterate over 1,2, 6, 4, 8,10, 12,14, 16 clients
- Each client has 48 threads
- Each thread is running under a private mount
- 2.3 stat performance is 200%-400% of 2.2
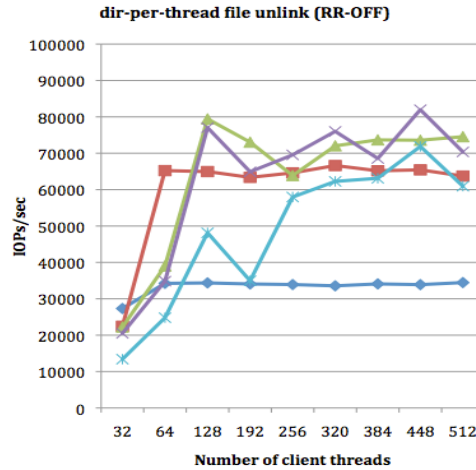- Client needs to send RPC to stat each OST object



256K files, 48 mounts per client (dir-per-thread, file stat)



256K files, 48 mounts per client (shared-dir, file stat)

# Performance of different CPT configurations

- MDS has 12 cores (24 HTs)

- 1 CPT

- 2 CPTs
  - Portal-RR ON

- 4 CPTs (default)
  - Portal-RR ON & OFF
  - 2 CPTs for LNet, 2 CPTs for ptlrpc service
  - 1 CPT for LNet, 3 CPTs for ptlrpc service

- 6 CPTs
  - Portal-RR ON & OFF
  - 2 CPTs for LNet, 4 CPTs for ptlrpc service
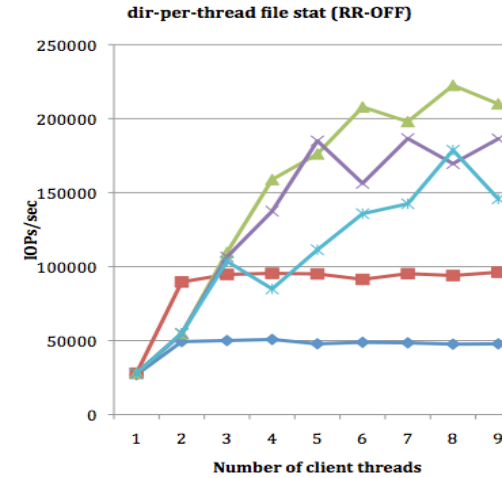
- 12 CPTs
  - Portal-RR ON & OFF

**Intel Architecture Group** (intel)

# Performance of different CPT configurations

Intel Architecture Group (intel)

# Lustre SMP configurations (libcfs)

- **Many chip types**
  - Server-1: Dual-core CPU, 8 sockets
  - Server-2: 50 cores, 1 socket
  - Server-3: 4 sockets, 2 NUMA nodes
  - Server-4: 2 sockets, 4 NUMA nodes

- **Default**
  - Preferred value "N"
    - $2 * (N - 1)^2 < NCPUS <= 2 * N^2$
  - Adjust "N" based on number of sockets or NUMA nodes

- **Configure CPU partitions for libcfs**
  - Libcfs cpu_npartitions=NUMBER
    - Prefer to put siblings in same CPT
  - Libcfs cpu_pattern=STRING_PATTERN
    - Example: libcfs cpu_pattern="0[0-6/2] 1[1-7/2]"
    - Example: libcfs cpu_parttern="N 0[0,2]  1[1,3]"

# Lustre SMP configurations (LNet)

- NID affinity
  - Hash NID by default
  - Bind NI on CPTs
    - O2ib0(ib0)[0, 1], tcp(eth0)[2, 3]

- Credits
  - NI credits
  - Router buffer credits

- Portal Round-Robin
  - /proc/sys/lnet/portal_rotor

- LND threads number
  - Decrease default threads number
  - Add extra threads for multiple interfaces

# Lustre SMP configurations (Lustre server & client)

- Bind service on CPTs
  - Both for MDS and OSS

- Use-cases
  - 32 cores machine, 4 sockets
  - Default
    - 4 partitions, LNet and ptlrpc services can run on all partitions
  - Config-1, one IB interface MDS
    - Lnet networks="o2ib0(ib0)[0]"
    - Mdt mdt_num_cpts="[1,2,3]"
  - Config-2, user only want to run Lustre client on one socket.
    - Libcfs cpu_pattern="0[0-31/4]"
    - Need some changes to set affinity for client threads

**Intel Architecture Group** (intel)

# Thank You