

TDS: Cray XC (dom)

Lustre striping Lustre allows the user to have explicit control over how a file is striped over the OSTs: chunks are sent to the different OSTs to improve disk bandwidth.

- export `MPICH_MPIIO_STATS=1`
- `srun -n192 ./GNU.DOM`

```
lt -h out_1120x720x80.16x12.000*
-rw-r--r-- 1 piccinal csstaff 3.7G Mar 24 14:32 out_1120x720x80.16x12.0000.bin
-rw-r--r-- 1 piccinal csstaff 3.7G Mar 24 14:32 out_1120x720x80.16x12.0001.bin
-rw-r--r-- 1 piccinal csstaff 3.7G Mar 24 14:33 out_1120x720x80.16x12.0002.bin
-rw-r--r-- 1 piccinal csstaff 3.7G Mar 24 14:33 out_1120x720x80.16x12.0003.bin
```

- `lfs setstripe -c 1` and `lfs setstripe -c 2`

```
-----+
| MPIIO write access patterns for out_1120x720x80.16x12.0003.bin
| independent writes      = 0
| collective writes       = 1920
| independent writers     = 0
| aggregators            = 1
| stripe count            = 1
| stripe size             = 1048576
| system writes          = 3750
| stripe sized writes     = 3750
| total bytes for writes  = 3932160000 = 3750 MiB = 3 GiB
| ave system write size   = 1048576
| read-modify-write count = 0
| read-modify-write bytes = 0
| number of write gaps    = 0
| ave write gap size      = NA
| See "Optimizing MPI I/O on Cray XE Systems" S-0013-20 for explanations.
|-----+
Testing get_procmem... 7516160.000000 45846528.000000 38330368.000000
written grids of 80,80,80
written 4 iterations
MPI Elapsed time: 50.929825 sec
average 0.028762 Gbytes/sec
real 54.23
```

```
-----+
| MPIIO write access patterns for out_1120x720x80.16x12.0003.bin
| independent writes      = 0
| collective writes       = 1920
| independent writers     = 0
| aggregators            = 2
| stripe count            = 2
| stripe size             = 1048576
| system writes          = 3750
| stripe sized writes     = 3750
| total bytes for writes  = 3932160000 = 3750 MiB = 3 GiB
| ave system write size   = 1048576
| read-modify-write count = 0
| read-modify-write bytes = 0
| number of write gaps    = 0
| ave write gap size      = NA
| See "Optimizing MPI I/O on Cray XE Systems" S-0013-20 for explanations.
|-----+
Testing get_procmem... 7516160.000000 42070016.000000 34553856.000000
written grids of 80,80,80
written 4 iterations
MPI Elapsed time: 25.865780 sec
average 0.056632 Gbytes/sec
real 29.15
```