# Process placement

2020

# Process placement

▶ Crucial to get right
 – MPI processes must be pinned to get a good performance
 – If it is not correct, the program is often several times slower

▶ SLURM and Intel MPI have good defaults
 – Choosing the number of MPI and OpenMP correctly is often good enough

Atos

# What type of application?

► Pure MPI
- – Well-balanced code, scales fine.
- – Unbalanced code, usually multi-process codes (e.g. climate).

► Hybrid MPI-OpenMP
- – hyperthreading

► Options to verify process placement

AtoS

# Well-balanced MPI code

► Just use the default distribution, whether SLURM or Intel's mpirun
  – Each task on one core
  – Fill up the nodes core by core

Atos

# Unbalanced code

► These codes often have a lot of MPI time (close to 50%)

► MPI time is indicative of load imbalance. No real communication of data, but waiting for comunication to start from the other side.

► As MPI time is mostly waiting time, the redistribution of MPI processes has little impact on MPI time.

► Strategy to distribute the tasks differently across and within nodes could help to move load from one NUMA domain across multiple NUMA domains.

► „Trial & error"

Atos

# Unbalanced code
# E.g. a job on 3 nodes

► srun -m cyclic
- Rank 0 on node 0, rank 1 on node 1
- Rank 3 on node 0, …
- E.g. could benefit OCTOPUS

► srun -m plane=12 (try also 8,6,4,3,2)
- Ranks 0-11 on node 0, ranks 12-23 on node 1
- Ranks 36-47 on node 0

Atos

# Unbalanced code
# Process placement in a node

► srun --cpu-bind=map_cpu:0,1,2,3,24,25,26,27,48,49,50,51,72,73,74,75,4,5,…

  – Distribute the MPI tasks across NUMA domains within a node
  – Can be combined with the plane distribution from the last slide
  – E.g. could benefit fesom2

Atos

# Hybrid MPI-OpenMP code

▶ Use at least 4 MPI processes on a CLX-AP node, 1 on each NUMA domain.

▶ OpenMP threads inside a NUMA domain

▶ Try what nr. of OpenMP threads works best.

▶ 96 cores gives you a lot of possibilities.

| MPI/node | OpenMP | MPI*OpenMP |
|----------|--------|------------|
| 4        | 24     | 96         |
| 8        | 12     | 96         |
| 12       | 8      | 96         |
| 16       | 6      | 96         |
| 24       | 4      | 96         |
| 32       | 3      | 96         |
| 48       | 2      | 96         |

AtoS

# Hybrid MPI-OpenMP code hyperthreading

▶ „The proof of the pudding is in the eating"

▶ Try with hyperthreads, use export OMP_WAIT_POLICY=passive
  – srun/sbatch: -c flag is equal to OMP_NUM_THREADS
  – mpirun distributes processes evenly.

▶ Try without hyperthreads
  – srun/sbatch: -c flag is 2x OMP_NUM_THREADS
  – mpirun distributes processes evenly.

AtoS

# Options to verify placement

▶ **`export I_MPI_DEBUG=4`**

   – Intel MPI prints the affinity and node for each process

▶ **`export KMP_AFFINITY=verbose`**

   – Intel OpenMP runtime prints the affinity for each thread

▶ `srun` **`--cpu-bind=verbose`** `app`

   – print the affinity of all processes

▶ **`srun –l hwloc-bind --get app`**

   – it prints the affinity of all processes, independent from srun

Atos

# Terminology: bit masks

► --cpu-bind=verbose prints bit masks
  – hwloc-bind --get --pid as well
  – more tools use masks

► read from right to left
► each hexadecimal digit represents 4 logical CPUs, e.g.
  – 0x01 is a mask where the first logical CPU is on
  – 0x02 : second logical CPU is on
  – 0x03 : first and second logical CPU are on
  – 0xF0 : fifth to eighth logical CPUs are on

AtoS

# Options to verify placement (interactive)

▶ Login on compute node, then run `htop`

▶ All cores should be busy (green in htop)
  – Note that a core is shown as 2 logical cpus: core 1 is CPUs 0 & 96 in htop
  – Process should be busy on only one of the two logical CPUs.
  – However, if hyperthreading is used, both logical CPUs should be busy.

▶ Very little system time (red in htop)
  – A lot of system time usually points to a problem
  – Maybe I/O or task switching.

Atos

# htop

# Thanks for your attention

john.donners@atos.net

Atos