

*P*⁶ PROPER PINNING PREVENTS PRETTY POOR PERFORMANCE

January 26, 2021 | T. Hater | JSC



Member of the Helmholtz Association

Superlinear Speed-Up?





Superlinear Speed-Up?

No, just a bad baseline...

• Default process placement switched between two cases.



Superlinear Speed-Up?

No, just a bad baseline...

- Default process placement switched between two cases.
- Second configuration is better for this benchmark.



STREAM benchmark

Heavily Optimised for Target Architecture, ...

- 1GiB, only triad (3 double per element).
- De-activated bindings by MPI and OpenMP.
- 10 runs each averaged over 5 repetitions, pick top result.
- -Ofast -march=native -mtune=native
- -std=c++17 -fno-builtin -fno-rtti -fno-exceptions -fopenmp
- Cache line blocked and aligned, SIMD, single fork/join, first touch aware, RMW optimised.



STREAM benchmark



e, RMW optimised.



Slide 3124

STREAM benchmark



e, RMW optimised.



Also: Binding, Affinity, ...

• Force a process or thread to execute only on a given set of cores.



- Force a process or thread to execute only on a given set of cores.
- Increases performance predictability and absolute performance.



- Force a process or thread to execute only on a given set of cores.
- Increases performance predictability and absolute performance.
- Enforced by the OS, driven by user space tools.



- Force a process or thread to execute only on a given set of cores.
- Increases performance predictability and absolute performance.
- Enforced by the OS, driven by user space tools.
- In HPC this is (partially!) handled by the scheduler (SLURM) or MPI.



- Force a process or thread to execute only on a given set of cores.
- Increases performance predictability and absolute performance.
- Enforced by the OS, driven by user space tools.
- In HPC this is (partially!) handled by the scheduler (SLURM) or MPI.
- But you can (should?) take control.



A Cartoon CPU



Many cores, each with its own memory hierachy.



A Cartoon CPU



- Many cores, each with its own memory hierachy.
- Shared global memory, but...



A Cartoon CPU



- Many cores, each with its own memory hierachy.
- Shared global memory, but...
- ...*affinitiy* to memory partitions.



A Cartoon CPU



- Many cores, each with its own memory hierachy.
- Shared global memory, but...
- ...affinitiy to memory partitions.

OS manages allocation,...



A Cartoon CPU



- Many cores, each with its own memory hierachy.
- Shared global memory, but...
- ...*affinitiy* to memory partitions.

- OS manages allocation,...
- ...task placement, and...



A Cartoon CPU



- Many cores, each with its own memory hierachy.
- Shared global memory, but...
- ...affinitiy to memory partitions.

- OS manages allocation,...
- ...task placement, and...
- ...swaps tasks in and out.



Scenario 1: Task Migration





Scenario 1: Task Migration





Scenario 1: Task Migration







Scenario 1: Task Migration



Important

Swapping tasks in and out is basically free, but task *migration* leads to data migration. Granularity is a *cache line* (often 128 *B*); be aware of *false sharing*.



Slide 6124

Scenario 2: NUMA

NUMA: Non-Uniform Memory Access, ie memory performance depends on relative location.





Scenario 2: NUMA

NUMA: Non-Uniform Memory Access, ie memory performance depends on relative location.





Scenario 2: NUMA

NUMA: Non-Uniform Memory Access, ie memory performance depends on relative location.





Scenario 2: NUMA

NUMA: Non-Uniform Memory Access, ie memory performance depends on relative location.



Important

All modern CPUs are NUMA architectures; might even have more than one NUMA domain! Memory is actually allocated on initialisation, use same parallel configuration as consumer. There will be no automatic migration.





Scenario 3: Sharing Resources



In some instances resources might be shared

- Hardware Threads (HWT) on a core might share computational units.
- Cores on a socket might share memory bandwidth, caches, ...

This can lead to sub-optimal performance by leaving some parts idle and others saturated. The inverse *might also be true*, eg it might be beneficial to share caches for read-only data.



Scenario 4: Specialisation



- Accelerators/network interfaces might be attached to a specific socket.
- If tasks/threads have specialised jobs, like MPI communication, ...
- ...scheduling them close to the relevant hardware can improve performance.
- Again: Beware the context switch.





This Talk

- \checkmark Motivation: Suboptimial and/or unpredictable performance
- ✓ Definition: What is pinning?
- ✓ Mechanism: Why does it improve performance?
- □ Learn to know your hardware.
- □ How to pin your processes.
- □ How to bind your threads.



```
> ml hwloc
> hwloc-ls # IMPORTANT: Run this on the *compute node*, eq via srun!
Machine (754GB total)
  Package L#0
    NUMANode L#0 (P#0 376GB)
    L3 L#0 (28MB)
      L2 L#0 (1024KB) + L1d L#0 (32KB) + L1i L#0 (32KB) + Core L#0
        PU L#0 (P#0)
        PU 1#1 (P#40)
      L2 L#1 (1024KB) + L1d L#1 (32KB) + L1i L#1 (32KB) + Core L#1
        PU 1#2 (P#1)
        PU 1#3 (P#41)
    HostBridge
      PCIBridge
        PCI 3b:00.0 (InfiniBand)
          Net "ib0"
          OpenFabrics "mlx5_0"
  Package L#1
    NUMANode L#1 (P#1 378GB)
    L3 L#1 (28MB)
```

hwloc documentation



ASCII Art Edition

chine (504GB total) Package L#0 +													
							L3 (16MB)				+ L3 (16MB)		
							L2 (512KB) L1d (32KB)	++ L2 (512KB) L1d (32KB) ++	L2 (512KB) L1d (32KB)		L2 (512KB) L1d (32KB)	L2 (512KB) L1d (32KB)	+ L2 (512KB) L1d (32KB) +
Core L#0 ++ PU L#0 PU L#1 ++	++ Core L#1 ++ PU L#2 PU L#3 ++	++ Core L#2 ++ PU L#4 PU L#5 ++		++ Core L#21 ++ PU L#42 PU L#43 ++	++ Core L#22 ++ PU L#44 PU L#45 ++	++ Core L#23 ++ PU L#46 PU L#47 ++							



Accelerators and Network Devices

hwloc-ls --output-format=pdf > node.pdf

1. A. of Long and a contract for contract for and be come for come for come CoPres rol 38 GB



Accelerators and Network Devices

Groupe							
NUMANode L#1 P#1 (63GB)							
L3 (16MB)	L3 (16MB)						
L2 (512KB) L2 (512KB) L2 (512KB)	L2 (512KB) L2 (512KB) L2 (512KB)						
L1d (32KB) L1d (32KB) L1d (32KB)	L1d (32KB) L1d (32KB) L1d (32KB)						
L11 (32KB) L11 (32KB) L11 (32KB)	L11 (32KB) L11 (32KB) L11 (32KB)						
Core L#6 PU L#12 P0 L#13 P0 L#13 P0 L#13 P0 L#15 P0 L#15 P0 L#15 P0 L#15 P0 L#15 P0 L#15 P0 L#16 P0	Core L819 Core L811 PU L812 PU L821 PU L815 PU L821 PU L821 PU L821 PU L823 PU L823						
OpenFabrics mix5_0 32 16 PCI 44:00,0 GPU nvnl1 32 32 PCI 45:00,0							



Options for Binding

Usually, a hybrid model is used: MPI tasks \times threads (OpenMP/pthreads/...)

Processes

- Resource Managers: SLURM, ...
- MPI implementations: OpenMPI, PSMPI, ...
- Linux: taskset, numactl, ...
- HWLoc CLI tools

Threads

- OpenMP Environment variables (if used)
- Linux Kernel API
- OpenMP API (if used)
- HWLoc API



Bind

--bind=[options] Enable binding

verbose Print binding masks. cores threads Use preset masks. rank Bind tasks to CPU IDs matching to task rank. rank_ldom Like rank, but distribute across NUMA domains. mask_cpu=0x.. List of bit masks, can be generated by hwloc tools.

Note: binding a process with threads still allows migration between the available HWT.

Warning

SLURM might still generate bad distributions, see examples later on.





Distribute

N=block|cyclic|fcyclic HWT

where

block keep tasks as close together as possible cyclic round-robin distribution of requested tasks fcyclic round-robin distribution of requested CPUs slurm documentation



Examples: Single-node

System JUWELS GPU

Node 2 sockets \times 20 cores \times 2 HWT

Request 1 node with 8 tasks \times 3 CPUs

Goal: Optimise for using as much of the hardware as possible, assuming the application does not benefit from co-locating tasks.





Examples: Single-node

System JUWELS GPU

Node 2 sockets \times 20 cores \times 2 HWT

Request 1 node with 8 tasks \times 3 CPUs

Goal: Optimise for using as much of the hardware as possible, assuming the application does not benefit from co-locating tasks.



Examples: Single-node

System JUWELS GPU

Node 2 sockets \times 20 cores \times 2 HWT

Request 1 node with 8 tasks \times 3 CPUs

Goal: Optimise for using as much of the hardware as possible, assuming the application does not benefit from co-locating tasks.

X Tasks split over sockets.

Each requested CPU acquires a full physical core.





Examples: Single-node

System JUWELS GPU

Node 2 sockets \times 20 cores \times 2 HWT

Request 1 node with 8 tasks \times 3 CPUs

Goal: Optimise for using as much of the hardware as possible, assuming the application does not benefit from co-locating tasks.





Examples: Multi-node

System JUWELS GPUs

Node 2 sockets \times 20 cores \times 2 HWT

Request 2 nodes with 20 tasks \times 2 CPUs



bind=threads
distribution=block:cyclic:fcyclic



Examples: Multi-node

System JUWELS GPUs

Node 2 sockets \times 20 cores \times 2 HWT

Request 2 nodes with 20 tasks \times 2 CPUs



bind=threads
distribution=block:cyclic:fcyclic

Images: PinningWebtool



Examples: Advanced Usage

System JUWELS Booster: NIC/GPUs attached to NUMA domains 1, 3, 5, 7 Goal 4 dedicated tasks for driving accelerators and communication each.



Examples: Advanced Usage

System JUWELS Booster: NIC/GPUs attached to NUMA domains 1, 3, 5, 7

Goal 4 dedicated tasks for driving accelerators and communication each.

```
> # Compute masks for all HWT in the relevant NUMA domains
> numa=`huloc-calc numa:1 numa:3 numa:5 numa:7
    # Generate masks for the distribution of & tasks across these
> mask=`huloc-distrib & --single --taskset --restrict $numa | xargs | tr ' ' ','`
> # Run application
> srun --cou_bind=verbose.cpu_mask=$mask -N 1 -n 8 -c 1 -- app.exe
```



Examples: Advanced Usage

System JUWELS Booster: NIC/GPUs attached to NUMA domains 1, 3, 5, 7

Goal 4 dedicated tasks for driving accelerators and communication each.

```
> # Compute masks for all HWT in the relevant NUMA domains
> numa=`huloc-calc numa:1 numa:3 numa:5 numa:7`
    # Generate masks for the distribution of & tasks across these
> mask=`huloc-distrib 8 --single --taskset --restrict $numa | xargs | tr ' ' ','`
> # Run application
> srun --cou_bind=verbose.cpu_mask=$mask -N 1 -n 8 -c 1 -- app.exe
```

Warning

Masks can be computed by hand, but keeping track of the numbering and bitsets is tedious and errorprone. The numbering scheme may change by: vendor, CPU generation, OS, ...







JUWELS Booster Default

Just use the default if your application does not have special requirements.

srun -N 1 -n 4 --gpus=4 --cpu-bind=socket -- app.exe

This does the right thing and **also** restricts the tasks' visible GPUs to the closest one.





Threads

- When using threads within tasks, these can use affinitiy as well.
- Without, threads will be mobile within the task-level masks.
- Consequently, we need to add another level of bindings...
- ...and take care not to conflict with task-level masks.



Threads: OpenMP Environment Variables

OMP_PROC_BIND=[...] Inhibit migration, bind threads to true First location it runs on. spread Spread over allowable set. close Block threads together. OMP_PLACES=[...] Bind threads to a set of places threads Individual hardware threads cores All HWT of a core sockets All cores of a socket $\{1, ...\}$ List of HWT ids

Migration is still allowed within a place when binding is not enabled. Using threads cores sockets with task binding is safe.

OpenMP specification





- Be aware of your application, we cannot provide a general solution.
- Binding for more performance and more predictability.
- Tools like hwloc allow mapping node topologies.
- High-level settings for performance and portability. Example: SLURM and OpenMP.
- Low-level tools, eg hwloc-API, for ultimate control.











Summary

- Be aware of your application, we cannot provide a general solution.
- Binding for more performance and more predictability.
- Tools like hwloc allow mapping node topologies.
- High-level settings for performance and portability. Example: SLURM and OpenMP.
- Low-level tools, eg hwloc-API, for ultimate control.