# Tools for Physicists: Boost your Analysis with High Performance Computing (HPC)

Hands on Trivial Parallelisation, Peter-Bernd Otte, 11.5.2022

#### Lecture Today

- Course webpage: <a href="https://indico.him.uni-mainz.de/event/136/">https://indico.him.uni-mainz.de/event/136/</a>
- Part of "Tools for Physicists" series: <u>https://www.hi-mainz.de/tfp22</u>

#### Talk (40')

- Motivation for High Performance Computing (HPC)
- Cluster building blocks and our HIMster2
- Trivial Parallelisation

#### Hands on (60')



#### **Trivial Parallelisation**

- todays course covers only trivial parallelisation and skips theory
   → see lecture next semester "Parallel Programming with OpenMP and MPI"
- Basic principle: run your existing analysis N times in parallel

time

#### **Trivial Parallelisation**

- todays course covers only trivial parallelisation and skips theory
   → see lecture next semester "Parallel Programming with OpenMP and MPI"
- Basic principle: run your existing analysis N times in parallel
- $\rightarrow$  How do we get there?

time

#### Worked out example

#### building of a house

- 1 worker = 1 year
- 3 workers = 4 months
- 9 workers = ?

#### $\rightarrow$ Scaling?



#### Running in parallel

Your analysis consists of 100 files to analyse

- On your desktop computer:
  - \$ ./myAnalysisExec InputFile1.dat OutputFile1.dat
- 8 cores:
  - ./myAnalysisExec InputFile1.dat OutputFile1.dat &
    ./myAnalysisExec InputFile2.dat OutputFile2.dat &
    ...

 $\rightarrow$  Your room mate has a computer, too, so why not use it?

### HPC out of distributed desktop computers?

- FLOPS / computer (floating-point operation per second):
  - FLOPS =  $f \times N_{cores} \times N_{instr per cycle}$
  - Intel E5-2670 (2,6 GHz, 8 cores): 2,6GHz × 8 × 8 = 166,4 GFLOPS
- N<sub>computers</sub>: 200 (=25 offices / floor, 4 floors, 2 people / office, 1 computer / person)
- 33TFLOPS cluster "for free" ⇔ Clover = 106TFLOPS, HIMster2/Mogon2: 2801TFLOPS

#### Drawbacks:

- OS: Windows (20%), MacOS (20%), Linux (50%) other (10%) all on a different version level
- Temperature in office rooms, closed window, 15th July: 0W = 29°C, with 400W = 50°C (simulated with: <u>www.thesim.at</u>)
- Network: 1GBit/s, Backbone 10GBit/s (HIMster2: 100GBit/s)
  - 10GBit/s / 200 computers / 8 cores = 780kByte/s
  - Compare bisection bandwidth (minimal accumulated bandwidth between any bisections of the network): fat tree binary tree
- Storage?
- No node checks, difficult to maintain, reduced availability



bisection bandwidth

# Why High Performance Computing (HPC)?

### Why HPC?

- Intense computational problem  $\rightarrow$  single desktop computer not capable enough
- Run on a "super computer"
  - 1. <2002: fast single core super computer
  - 2. Since 2002: parallel systems as super computers
    - $\rightarrow$  Why parallel systems?

### The Era of Moore's Law

- 1900-2000
- source: Wikipedia



#### Microprocessor transistor counts 1971-2011 & Moore's law The Era of Moore's Law

- Moore's law (1965) = observation number of transistors in a IC doubles every ~2a.
- Still valid, no natural law.

Cramming More Components onto IC (1965): <u>ftp://download.intel.com/sites/channel/museum/Moores\_Law/Articles-</u> Press\_Releases/Gordon\_Moore\_1965\_Article.pdf or

https://ieeexplore.ieee.org/document/658762?tp=&arnumber=658762



Date of introduction

### Single-Core Performance



<sup>© 2007</sup> Elsevier, Inc. All rights reserved.

### Why not further increase frequency?

- Core speeds topped out at 2-4 GHz
  - World record standard CPU: 8722.78 MHz with liquid nitrogen cooling (http://hwbot.org/benchmark/cpu\_frequency/)
- Problem #1: cooling the chip
- Finding: "Dennard scaling" (constant power density) no longer valid
  - No longer (since 2000's) true since 90nm gate sizes (leakage current!)
  - The two things that consume energy (CMOS gate):
    - 1. switching state (1  $\Leftrightarrow$  0) (10 $\mu$ W/MHz, prop with f^1.75)
    - 2. leakage current (10nW / CMOS-Gate, anti-prop with Vdd and gate size)
- Increasing f: increase in power on same area
  - $\rightarrow$  compensate this: shrinking gate sizes and lower Vdd
    - But: smaller gates have higher leakage current. → New innovations needed.
       → multi-cores at fixed f to gain performance



not gate



Power=P(Vdd, f)

#### Answer: multicores



#### Moore's Law scaling with cores



# Recap: Trivial Parallelisation and Multicore systems

- No drawback using a multi core machine
- We have single independent jobs
  - Assign single analysis runs to single cores



# Recap: Trivial Parallelisation and Multicore systems

- No drawback using a multi core machine
- We have single independent jobs
  - Assign single analysis runs to single cores



 $\rightarrow$  We are on the right path, so let's dive in.

## HPC building blocks

### What is High Performance Computing (HPC)

- Basic building blocks are:
  - 1. compute nodes (~1000)
  - 2. fast interconnect (1x)
  - 3. parallel file system (1x)



• Usage remotely, non interactively

#### Anatomy of a node

• cache coherent Non-Uniform Memory Access (ccNUMA, AMD: 2003, industry wide: 2011)



• ccNUMA uses inter-processor communication between cache controllers

#### Machine (32GB)

Socket P#0 (16GB)					Socket P#1 (	l6GB)								
NUMANode P#0 (8192MB)					NUMANode P#2 (8192MB)									
L3 (8192KB)					L3 (8192KB)									
L2 (2048KB)	L2 (2048KB)	L2 (2048KB)	L2 (2048KB)		L2 (2048KB)			L2 (2048KB)		L2 (2048KB)		L2 (2048KB)		
L1i (64KB)	L1i (64KB)	L1i (64KB)	L1i (64KB)		L1i (64KB)			L1i (64KB)		Lli (64KB)		Lli (64KB)		
Lld (16KB) Lld (16KB)	L1d (16KB) L1d (16KB)	L1d (16KB) L1d (16KB)	L1d (16KB) L1d (16KB)		L1d (16K	3) L1d (16KB)		L1d (16KB) L1d (	(16KB)	L1d (16KB)	Lld (16KB)	L1d (16KB)	L1d (16KB)	
Core P#0 Core P#1 PU P#0 PU P#1	Core P#2 Core P#3 PU P#2 PU P#3	Core P#4 Core P#5 PU P#4 PU P#5	Core P#6 Core P#7 PU P#6 PU P#7		Core P#0	Core P#1 6 PU P#17	]	Core P#2 Core PU P#18 PU	P#3 P#19	Core P#4 PU P#20	Core P#5 PU P#21	Core P#6 PU P#22	Core P#7 PU P#23	
NUMANode P#1 (8192MB)					NUMANode P#3 (8192MB)									
L3 (8192KB)					L3 (8192KB)									
L2 (2048KB)	L2 (2048KB)	L2 (2048KB)	L2 (2048KB)		L2 (2048KB)			L2 (2048KB)		L2 (2048KB)		L2 (2048KB)		
L1i (64KB)	L1i (64KB)	L1i (64KB)	L1i (64KB)		L1i (64KB)			L1i (64KB)		L1i (64KB)		L1i (64KB)		
Lld (16KB) Lld (16KB)	L1d (16KB) L1d (16KB)	L1d (16KB) L1d (16KB)	Lld (16KB) Lld (16KB)		L1d (16K	3) L1d (16KB)		L1d (16KB) L1d (	(16KB)	L1d (16KB)	Lld (16KB)	Lld (16KB)	L1d (16KB)	
Core P#0         Core P#1           PU P#8         PU P#9	Core P#2 Core P#3 PU P#10 PU P#11	Core P#4 Core P#5 PU P#12 PU P#13	Core P#6 Core P#7 PU P#14 PU P#15		Core P#0	4 Core P#1	]	Core P#2 Core PU P#26 PU	P#3 P#27	Core P#4 PU P#28	Core P#5 PU P#29	Core P#6 PU P#30	Core P#7 PU P#31	

- Output of: hwloc
- Topology of a ccNUMA Bulldozer server, 2 socket system

#### Anatomy of a cluster computer

• N \* ccNUMA = cluster:



- Fast lossless interconnect: OmniPath between ccNUMA nodes
- Inside a node: NUMA, ccNUMA
- Multiple nodes: Distributed memory parallelisation (DMP)

### Anatomy of a cluster computer

• Latencies:



(all numbers are platform dependent)

#### HIMster II Specs

- 320 Compute Nodes (256 theory, 64 experiment) in 8 racks
  - dual socket Intel 6130 @ 2.1GHz (à 16 cores)
  - 3GB RAM /core
  - OmniPath 100 Gbit/s interconnect
  - 400 GB local SSD scratch
  - <a href="https://mogonwiki.zdv.uni-mainz.de/dokuwiki/start:mogon\_cluster:nodes">https://mogonwiki.zdv.uni-mainz.de/dokuwiki/start:mogon\_cluster:nodes</a>
  - Parallel File System: 747TB Lustre volume
- Software
  - 1. organized in modules
    - eg:module avail; module module load lang/Python/3.6.6-foss-2018b
    - See: <u>https://mogonwiki.zdv.uni-mainz.de/dokuwiki/start:software:software\_usage</u>
  - 2. More via nfs mount: /cluster

#### HIMster II

- HIMster II, Mogon IIa and Mogon IIb form a compound state
  - share login nodes, maintenance servers
  - interconnect: OmniPath (100GBit/s)
- situated in the institute's basement computing room, 660kW
- 2PFlops Linpack (20% contributes HIMster II)
- account registration via PI of HIM or <u>it@him.uni-mainz.de</u>.
   University of Mainz account is mandatory (→ HIM Admin will contact you).
- ssh pbotte@miil01-miil03 (only ssh-key login possible, with 2<sup>nd</sup> factor)
  - <u>https://mogonwiki.zdv.uni-</u> mainz.de/dokuwiki/start:mogon\_cluster:access\_from\_outside\_unix
- home directory: quota 300 GB
- Rules apply: <u>https://www.en-zdv.uni-mainz.de/regulations-for-use-of-the-data-center/</u>

#### HIMster II: Info and do's

- Per core memory bandwidth: Clover, HIMster II = 5.6 GByte/sec
- HIMsterII has Skylake CPUs (eg AVX512 avail.)
- Storage / Parallel File system:
  - NO BACKUP of data
  - Try to use large files: Source code should be in /home/
  - Try not to put too many files into one directory (less than 1k)
  - Try to avoid too much metadata load:
    - DO NOT DO ls -1 unless you really need it
    - In your scripts avoid excessive tests of file existence (put in a sleep statement between two tests say 30 secs)
    - Use lfs find rather than GNU tools like find
    - Use O\_RDONLY | O\_NOATIME (readonly and no update of access time)

### Batch System: SLURM

- Batch system, introduces fair share
  - Accounts (e.g. m2\_himkurs, m2\_himexp, etc.)
  - Queues
  - Reservations
- Introduction and docu:
  - <u>https://mogonwiki.zdv.uni-mainz.de/</u> <u>dokuwiki/start:working\_on\_mogon:slurm\_submit</u>
  - <u>https://slurm.schedmd.com/tutorials.html</u>
- Today:
  - account to use: m2\_himkurs
  - Reservation: himkurs
  - Submit into partition: parallel
    - srun --pty -p parallel -A m2\_himkurs --reservation himkurs bash -i
- Check what is running: squeue -h | grep pbotte
  - 1184615\_79 parallel N203r001 **pbotte** R 1:00:40 52 z[0367-0386,0403-0413,0430-0450]
  - SSH login into your occupied nodes possible: eg ssh z0367
    - only for debugging, do not launch analysis tasks!



#### SLURM scheduler: Multifactor Priority

https://slurm.schedmd.com/priority\_multifactor.html



resources

### SLURM scheduler: Backfilling

Performed only when jobs with higher prio are not affected

https://slurm.schedmd.com/sched\_config.html



#### Batch System: SLURM



- Submit script for later execution (batch mode)
  - sbatch --partition=himster2\_exp
- Create job allocation and start a shell to use it (interactive mode)
  - salloc -p himster2\_exp -N 1 --time=02:00:00 -A m2\_him\_exp
- srun: Create a job allocation (if needed) and launch a job step (typically MPI job)
  - srun --pty -p himster2\_exp -N 1 --time=02:00:00 -A m2\_him\_exp bash -i
- sattach: Connect stdin/out/err for an existing job

#### Sample Submit Script

Examples only – not for today's hands on!

- Define and reserve resources (nodes with RAM)
- 2. Once allocated, run the executables as defined or interactively

More examples https://mogonwiki.zdv.unimainz.de/ dokuwiki/start:working\_on\_ mogon:slurm\_submit

#!/Dln/k	Jash
#SBATCH	-o /home/pbotte/test/myjob.%j.%N.out
#SBATCH	-D /home/pbotte/test/
#SBATCH	-J MyJobName
#SBATCH	-A m2_him_exp 🗧 🗲 account (NOT your account)
#SBATCH	-N 1 🗲 Request number of nodes
#SBATCH	<pre>partition=himster2_exp</pre>
#SBATCH	mem-per-cpu=1G
#SBATCH	mail-type=FAIL
#SBATCH	mail-user=pbotte@uni-mainz.de
#SBATCH	time=8:00:00

```
module load gcc/6.3.0
echo TEST...
srun myExecutable
```

Submit with: sbatch submitScript.sh

#### HIMster compute nodes 8 racks

<sup>111</sup>1 1 1

1 1 1

#### Cooling power for up to 750kW

**Power and OmniPath Interconnect** 



## Optimisation and usage

#### Amdahl's Law

• Given a program consisting of a non-parallelisable and a perfectly parallelisable part

• Fraction **s** of the non-parallelisable part:  $T(p)=T_{seq} + T_{par}(p) = T(1) * s + T(1) * (1-s)/p$ 



- Speed-up: S(p) = (1+(1-s)/p)<sup>-1</sup>
  - $p \rightarrow \inf: S(p) = 1/s$
  - If S(p) > 1/s → "super-scaler speedup", problem fit's into CPU cache.



### Order of optimisation

How to speed up your existing analysis:

• Apply trivial parallelisation (todays topic!)

Want to go further?

 $\rightarrow$  Identify bottlenecks (and only optimise them)

- 1. Optimise algorithm
- 2. Write algorithm on single core
- 3. Expand code to multicore, single node with OpenMP
- 4. Expand to multi node with MPI
- 5. Optimise multi node system

 $\rightarrow$  Not covered today, lecture in winter semester.

#### Parallel Programms: Worked out example

• Task: calculate sum of numbers distributed over N cores



### Trivial vs full usage of HPC

- Trivial parallelisation:
  - Run your analysis several times (with different parameters)
  - Out of the box with any non-interactively linux program
  - Outcome / speedup unclear, but works very good for 10-100 jobs in parallel Mainly disc access is limiting.
- Full usage (not covered today):
  - No automated process to convert a single-core to a multi-core program
  - Write parallel code or use existing.

## Preparing Hands on...

### Today's Setup

https://indico.him.uni-mainz.de/event/136

- Linux basics used:
  - Bash: launch a program with different parameters
  - SSH: generate a key and log into HIMster2
  - modules: list and load different modules
  - (versioning with git)
- (ideal world: work in groups of 2 on one computer.)
- (If not present, familiarise with it OR find the right team mate!) Ask your neighbour or tutor.

### Public-key cryptography

• Asymmetric Encryption





#### • connecting to Himster2:

- 1. run "ssh-keygen" if you have no keys yet
  - Private key: ~/.ssh/id\_rsa
  - Public key: ~/.ssh/id\_rsa.pub
- 2. Copy to ZDV via: <u>https://account.uni-mainz.de/my-account/add-ssh-key</u>
- 3. Authenticate to get 2nd factor via request to hpc@uni-mainz.de
- 4. ssh into mogon2 / HIMster2

Pictures from wikipedia

### Trivial Parallelisation (1)

- Submit a single core job multiple times
- Quick and often only solution for large software blobs (large packages used in collaborations)
  - No principal difference compared to running on your desktop computer
- limits:
  - required RAM (3GB/core)
  - licensees (Mathematica, max 10 concurrent usages in university for such uses cases)
  - shared scratch (under "/localscratch") in node (200GB-400GB)
  - parallel filesystem (loading at start, writing back results) max.  $\rightarrow$  10-100 starting jobs in parallel
- Hint: use job arrays
  - <a href="https://mogonwiki.zdv.uni-mainz.de/dokuwiki/start:working\_on\_mogon:workflow\_organization:job\_arrays">https://mogonwiki.zdv.uni-mainz.de/dokuwiki/start:working\_on\_mogon:workflow\_organization:job\_arrays</a>
  - Less work load for SLURM
- Disadvantage (for single and array jobs):
  - Single job on Mogon2 parallel partition always node exclusive: Single job blocks the complete node, independent on how many resources requested!
  - No control over speedup
  - Node health check (~1min) and batch system overhead (~1min) for every step
     → bundle them to larger blocks → use a workload manager!

Helper MPI-Script

https://gitlab.rlp.net/pbotte/workload-manager

- Occupy  $N_{cores}$  cores on  $|N_{cores}/20|$  (HIMster 2:  $|N_{cores}/32|$ ) different machines simultaneously
- Provide a directory with files to process (N<sub>files</sub>)
- Controlling instance on core 0
- Starts your analysis executable on workers (cores 1..N-1)
- Feedback or pull requests welcome

- Suits short and long running analysis (avoid node health checks)
- Occupies a complete node
- Does load distribution
- Takes care of in and output files



Helper MPI-Script

https://gitlab.rlp.net/pbotte/workload-manager

- Occupy N<sub>cores</sub> cores on |N<sub>cores</sub>/20| (HIMster 2: |N<sub>cores</sub>/32|) different machines simultaneously
- Provide a directory with files to process (N<sub>files</sub>)
- Controlling instance on core 0
- Starts your analysis executable on workers (cores 1..N-1)
- Feedback or pull requests welcome

- Suits short and long running analysis (avoid node health checks)
- Occupies a complete node
- Does load distribution
- Takes care of in and output files



Helper MPI-Script

https://gitlab.rlp.net/pbotte/workload-manager

- Occupy N<sub>cores</sub> cores on |N<sub>cores</sub>/20| (HIMster 2: |N<sub>cores</sub>/32|) different machines simultaneously
- Provide a directory with files to process (N<sub>files</sub>)
- Controlling instance on core 0
- Starts your analysis executable on workers (cores 1..N-1)
- Feedback or pull requests welcome

- Suits short and long running analysis (avoid node health checks)
- Occupies a complete node
- Does load distribution
- Takes care of in and output files



Helper MPI-Script

https://gitlab.rlp.net/pbotte/workload-manager

- Occupy  $N_{cores}$  cores on  $|N_{cores}/20|$  (HIMster 2:  $|N_{cores}/32|$ ) different machines simultaneously
- Provide a directory with files to process (N<sub>files</sub>)
- Controlling instance on core 0
- Starts your analysis executable on workers (cores 1..N-1)
- Feedback or pull requests welcome

- Suits short and long running analysis (avoid node health checks)
- Occupies a complete node
- Does load distribution
- Takes care of in and output files



#### Hint: Reservation for this problem class

\$ scontrol show reservation

ReservationName=himkurs StartTime=2022-05-11T14:00:00 EndTime=2022-05-11T17:00:00 Duration=03:00:00
Nodes=z[0500,0566,0751-0758,0811,0813] NodeCnt=12 CoreCnt=240 Features=(null) PartitionName=parallel Flags=
TRES=cpu=480
Users=(null) Groups=(null) Accounts=m2\_himkurs Licenses=(null) State=INACTIVE BurstBuffer=(null) Watts=n/a
MaxStartDelay=(null)

\$ salloc -p parallel --reservation=himkurs -A m2\_himkurs -N 1

#### Exercise 1: HIMster 2 log in

Learning objectives:

• Usage of ssh, asymmetric keys and 2<sup>nd</sup> factor

- 1. Go back to slide "Public-key cryptography" and create your keys
- 2. Login passwordless but with 2<sup>nd</sup> Factor into Mogon 2/ HIMster2

#### Exercise 2: Reserve Resources

Learning objectives:

- Reserve resources
- Check number of cores on node

- 1. Log into Himster 2 ssh miil01.zdv.uni-mainz.de
- Reserve a complete HIMSter node for 1h: salloc -p parallel --reservation=himkurs -A m2\_himkurs -N 1 -t 1:00:00 This step might take some minutes to complete. Wait until the prompt returns after "salloc: Nodes z0133 are ready for job"
- 3. Confirm that information with this cross check: squeue -u \$USER
- 4. Find out how many cores your node has with ssh [YOUR Node hostname] cat /proc/cpuinfo read the info and logout of that node with logout

#### Exercise 3: Single core test run

Learning objectives:

• Perform a test drive of your demo analysis

- 1. If not already done so, reserve first resources as described in exercise 2. Check with: squeue -u \$USER
- 2. Open 2 more ssh connections to run "top" two times: (1) on the head node (2) on the node
- 3. In your home directory, prepare the demo analysis with: git clone https://gitlab.rlp.net/pbotte/learnhpc/ cd /home/pbotte/learnhpc/openMP/exercise1 cc -o pi pi\_start.c
- 4. Make sure, you are working on the head node, run your program: ./pi Check, with your other SSH connections (see step 2), the binary runs on the head node.
- 5. Make sure, you are working on the head node, run your program: srun ./pi Check, with your other SSH connections (see step 2), the binary runs on the node.

#### Exercise 4: batch job

Learning objectives:

• Your first batch job

- 1. If not already done so, reserve first resources as described in exercise 2. Check with: squeue -u \$USER
- 2. Write a batch job file (name it "job.batch") as shown on the right. Replace "USERNAME" -> with your username!
- 3. Queue it: sbatch job.batch
- 4. Check what is does regularly: (a) squeue -u \$USER

  - (b) via top on the node (c) your email account.

#!/bin/bash
#SBATCH -o /home/USERNAME/test/myjob.%j.%N.out
#SBATCH -D /home/USERNAME/test/
#SBATCH -J MyJobName
#SBATCH -A m2_himkurs
#SBATCH -n 1
#SBATCHpartition=parallel
#SBATCHmem-per-cpu=1G
#SBATCHmail-type=FAIL
#SBATCHmail-user=USERNAME@uni-mainz.de
#SBATCHtime=0:05:00
echo Here comes some test
<pre>srun /home/USERNAME/learnhpc/openMP/exercise1/pi</pre>



Learning objectives:

• Run an executable on several cores in parallel

#### Steps:

- 1. If not already done so, reserve first resources as described in exercise 2. Check with: squeue -u \$USER
- 2. After you double checked that you are on the head node again (prompt start with "login2..."), load the modules needed for python3 (for the workload manager) and solver (Sundials) in exactly this order: module load math/SUNDIALS/2.7.0-intel-2018.03 module load lang/Python/3.6.6-foss-2018b lgnore the error messages (bonus: where do they come from?)
- Run the analysis with the help of the workload manager by typing on the headnode: srun -n 20 ~/workload-manager/wkmgr.py -v ~/workload-manager/examples/LGS/PulsedLGS ~/workloadmanager/examples/LGS/Run27\_LaPalma\_Profile\_I50
- 4. Output is send back to the head node delayed and not in order. Check with a second console any output written,
  - run "top" to see running processes, or
  - type ",tree", in the output directory ~/workload-manager/outputXXX
- 5. How much compared to a single core machine is roughly the speed up?

Hint: Like this, you can easily run on several nodes each having 32 cores.

### Exercise 6: mpi4py hello world (Bonus)

Learning objectives:

• Use MPI with Python the first time aka: make Python run its code on several cores and machines in parallel

Detailed description: https://gitlab.rlp.net/pbotte/learnhpc/-/tree/master/mpi4py/exercise1

#### Steps:

1. Download the starter files:

git clone https://gitlab.rlp.net/pbotte/learnhpc.git cd learnhpc/mpi4py/exercise1/

2. Copy the skeleton:

cp start.py ex1.py

- 3. Load environment: module load lang/Python/3.6.6-foss-2018b
- 4. Try with different number of ranks ("-n"), start with 3: mpirun -n 3 ./ex1.py