# Parallel computation methods on CPU architectures

Claudia Blaas-Schenner

VSC Research Center, TU Wien

VIENNA SCIENTIFIC CLUSTER

06/2021

Univerza *v Ljubljani*

TU WIEN — TECHNISCHE UNIVERSITÄT WIEN

CINECA consorzio interuniversitario

VSB TECHNICAL UNIVERSITY OF OSTRAVA | IT4INNOVATIONS NATIONAL SUPERCOMPUTING CENTER

# HPC systems

- VSC – joint high performance computing (HPC) facility of Austrian universities



VSC-3



VSC-4

VSC-5

soon to come…

https://vsc.ac.at
https://vsc.ac.at/training

VIENNA SCIENTIFIC CLUSTER

# components of HPC clusters

login nodes: l40 l41 l42 . . . l49

SLURM

storage $HOME

compute nodes & interconnect:
n# n# n# . . . . . . . . .
n# n#
n#

storage $DATA

login nodes

compute nodes

shared (login, storage)

user exclusive (compute nodes)
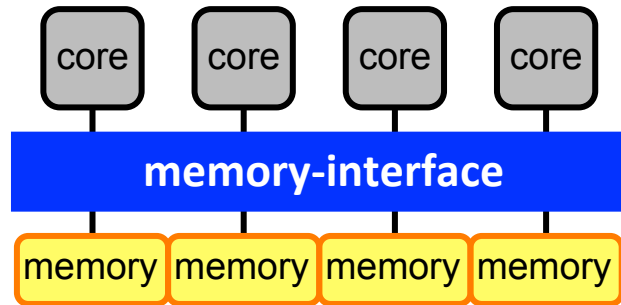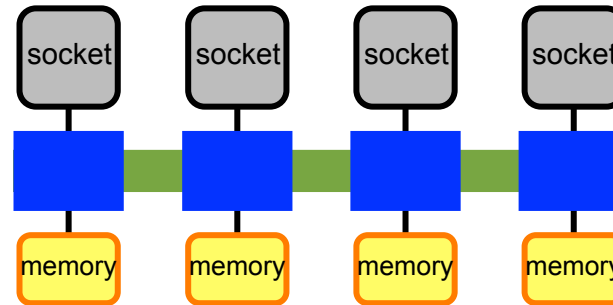
# parallel hardware

## shared memory

**socket: → memory-interface**
UMA (uniform memory access)
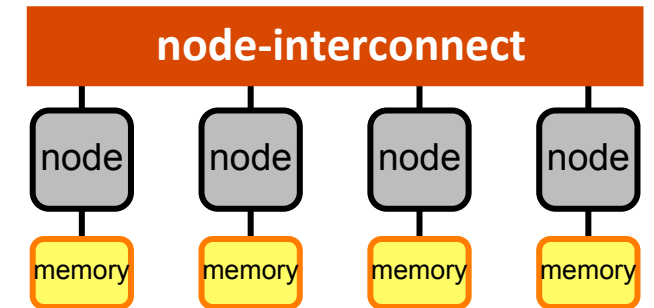SMP (symmetric multi-processing)

**socket / CPU**

**node: → hyper-transport**
ccNUMA (cache-coherent non-uniform…)
! first touch, pinning !

**node**

## distributed memory
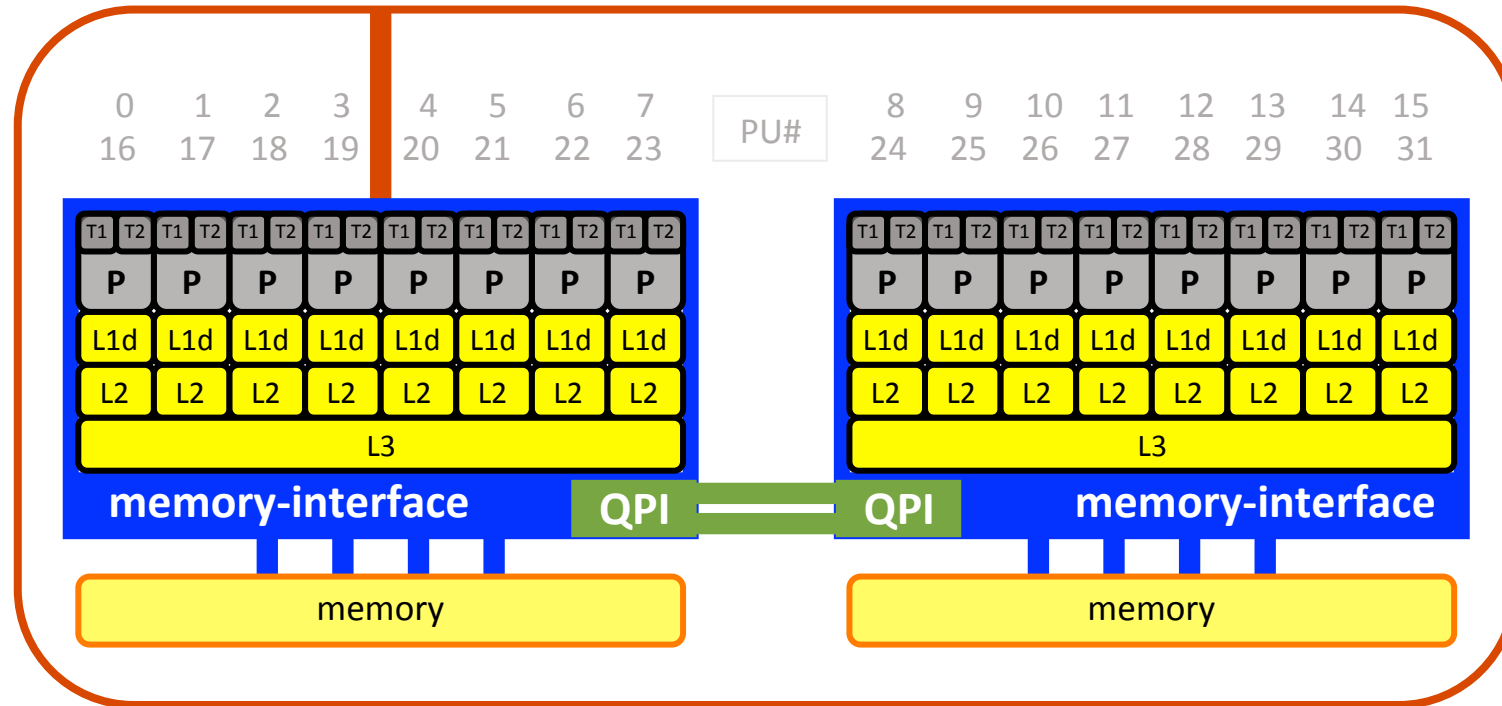
**node-interconnect**

**cluster: → node-interconnect**
NUMA (non-uniform memory access)
! fast access only to its own memory !

**cluster**

shared memory programming with **OpenMP**

**MPI** works everywhere

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | PU# | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |

T1 T2 | T1 T2 | T1 T2 | T1 T2 | T1 T2 | T1 T2 | T1 T2 | T1 T2

P P P P P P P P

L1d L1d L1d L1d L1d L1d L1d L1d

L2 L2 L2 L2 L2 L2 L2 L2

L3

**memory-interface**   QPI ══ QPI   **memory-interface**

memory   memory

**example:**

1 node

2 sockets (CPUs)

8 cores per socket (P)

2 threads per core (T1/T2)

1 HCA (host channel adapter)
        (node-interconnect)

**info about nodes:**

numactl --hardware     [Linux]

cpuinfo -A             [Intel]
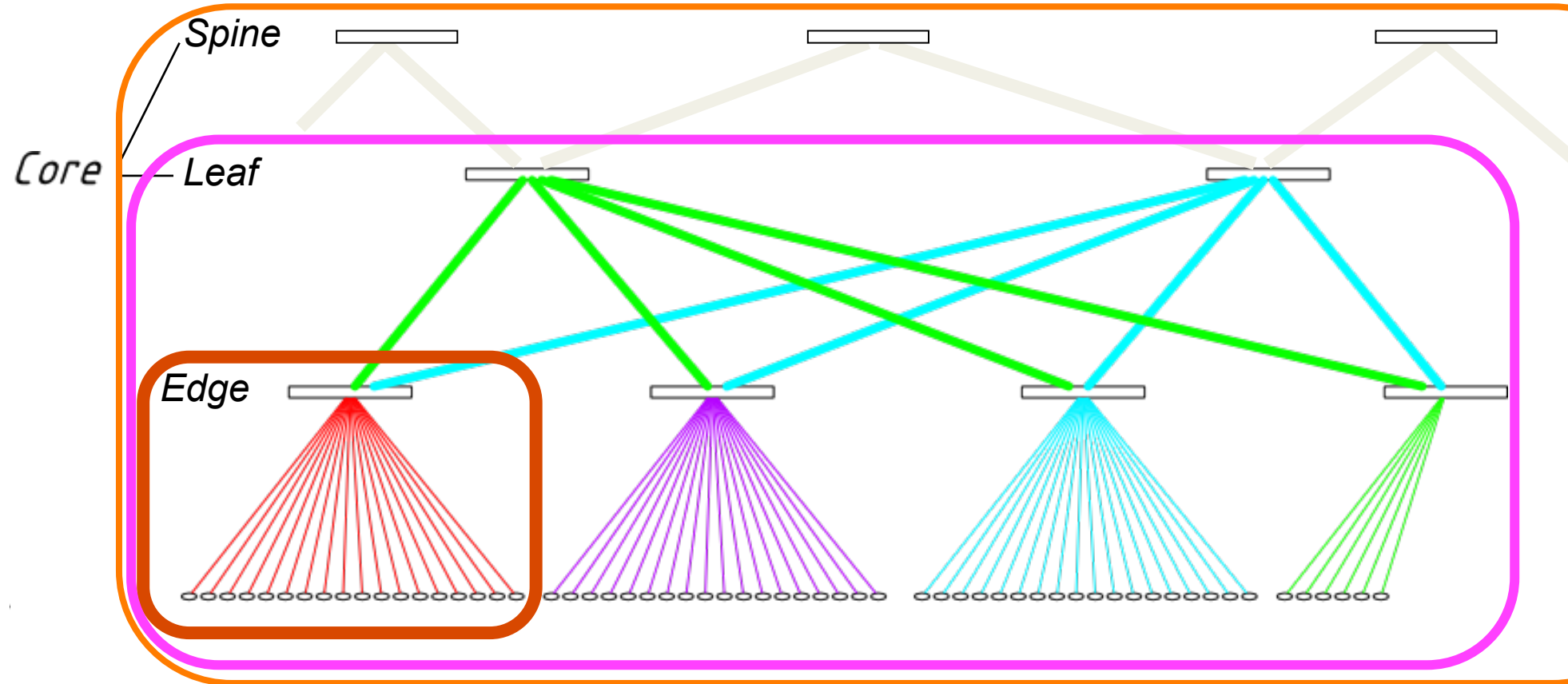
likwid-topology -c -g     [LIKWID]

# node-interconnect

**schematic figure:**

3-level fat tree

2-level fat-tree

1st level switches

compute nodes

attached to the

lowest level

# Amdahl's Law

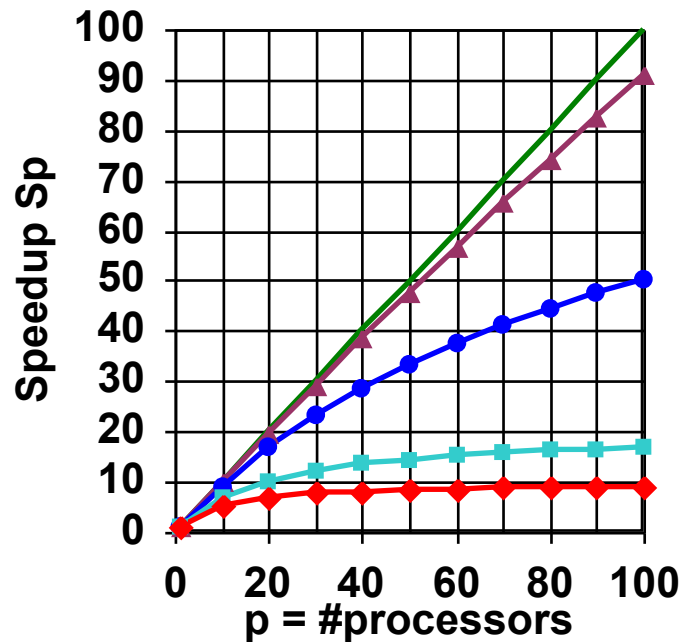$$T_{parallel,\,p} = f \cdot T_{serial} + (1-f) \cdot T_{serial} / p$$

f … sequential part of code     **neglecting time for communication**

$$S_p = T_{serial} / T_{parallel,\,p} = 1 / (f + (1-f) / p)$$

Speedup is limited: $S_p < 1 / f$     **neglecting load imbalance**

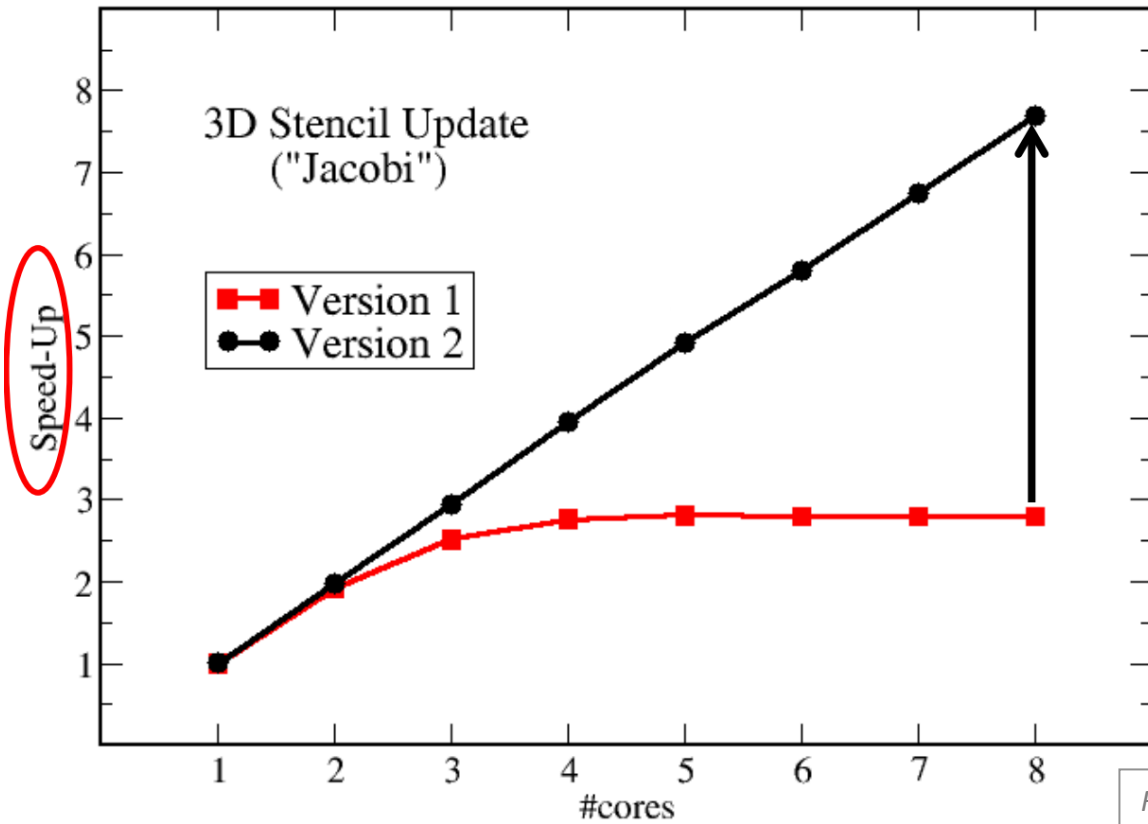

- Sp = p (ideal speedup)
- f=0.1%  =>  Sp < 1000
- f=  1%  =>  Sp < 100
- f=  5%  =>  Sp < 20
- f= 10%  =>  Sp < 10

*Figures courtesy of Rolf Rabenseifner.*

**Speedup = ratio – no absolute performance !**

7

# scalability vs. performance

3D Stencil Update ("Jacobi")

Single core/socket efficiency is key issue!

*Figures courtesy of Georg Hager.*

3D Stencil Update ("Jacobi"): `y(i,j,k) = b*(  x(i-1,j,k) + x(i+1,j,k) + x(i,j-1,k)`
                              `+ x(i,j+1,k) + x(i,j,k-1) + x(i,j,k+1))`

8

# pinning ?

**no pinning**

OpenMP
STREAM benchmark

*Benchmark & plots courtesy of Georg Hager.*

**MPI** will give the very same picture **!**



**pinning
(physical cores first,
first socket first)**

why should we care about **pinning** ?

- eliminating performance variations
- making use of architectural features
- avoiding resource contention

# HPC = computation – communication – I/O

HPC

| LATENCY | ← typical values → | BANDWIDTH | | |
|---|---|---|---|---|
| 1–2 ns | L1 cache | 100 GB/s | **computation** | *exclusive* |
| 3–10 ns | L2/L3 cache | 50 GB/s | **node / core** | |
| 100 ns | memory | 10 GB/s | **communication** | *exclusive* |
| 1–10 μs | HPC networks | 1–8 GB/s | **message passing** | *(BF)* |
| 50 μs | Gigabit Ethernet | 100 MB/s | | |
| 500 μs | Solid state disk | 100 MB/s | **I/O** | *shared* |
| 10 ms | Local hard disk | 50 MB/s | **parallel FS** | *with all users* |
| 50 ms | Internet | 10 MB/s | | |

Understand HW features!

Know your code!

Know the sys. environment!

→ Take control!

➔ **Avoiding slow data paths is the key to most performance optimizations!**

# login to a cluster

- **username and password** (ssh-keys)

- restricted IPs (firewall)

- two-factor authentication
- ➢ **graphical user interface (GUI)**

- terminal: xterm, terminal, PuTTY
- X-server, XQuartz, Xming

- ssh <username>@<cluster>
- ssh -X <username>@<cluster>

- ➢ Linux command-line access
- **NoMachine** (remote virtual desktop)

→ @viz.hpc.fs.uni-lj.si

# modules & compiling

- **module environment** [spack]

  ```
  module list | purge | load | avail [2>&1 | less]
  ```

- **compiling with GCC**

  ```
  module load foss/2019a
  ```

  ```
  cc --version
  mpicc --version
  ```

  ```
  cc [-fopenmp] program.c
  mpicc [-fopenmp] program.c
  ```

- **compiling with Intel**

  ```
  module load intel          → @viz.hpc.fs.uni-lj.si
  ```

  ```
  icc --version
  mpiicc --version
  ```

  ```
  icc [-qopenmp] program.c
  mpiicc [-qopenmp] program.c
  ```

# SLURM what kind of nodes

- partitions

  ```
  sinfo -o %P

  sinfo
  ```

- qos (quality of service)

  ```
  sacctmgr show qos
  ```

- @VSC → `sqos -acc` & `sqos`

- use other than default

  ```
  #SBATCH --qos=<qos>

  #SBATCH --partition=<partition>

  #SBATCH --account=<account>
  ```

- more detailed info

  ```
  scontrol show partition <part.>
  scontrol show node <node>
  scontrol show reservation
  ```

# job submission via SLURM

- **SLURM** job script

  `#!/bin/bash`            → has to be a shell script

  `#SBATCH`               → header lines for the job scheduler

  `do_my_work`            → whatever needs do be done by the job

- **SLURM** queuing system

- `sbatch job.sh`         → submit

- `squeue -u $USER`       → check

- `scancel JOB_ID`        → cancel

- `slurm-*.out`           → output

- **recommended** @ `~/.bashrc`

  **`alias sq='squeue -u $USER'`**
  `export LC_CTYPE=en_US.UTF-8`
  `export LC_ALL=en_US.UTF-8`
  `unset LD_PRELOAD`

  `source ~/.bashrc`

SCtrain training week provides a node reservation (daily, 9 am – 5 pm) to avoid queuing times:
`sbatch --reservation=sctrain job.sh`

```
#!/bin/bash                          #  → @viz.hpc.fs.uni-lj.si


#SBATCH -J test                      # SLURM_JOB_NAME
#SBATCH -N 2                         # SLURM_JOB_NUM_NODES
#SBATCH --tasks-per-node=24  # SLURM_NTASKS_PER_NODE   [1 mpi/core]


# <do_my_work>
module purge                         # recommended to be done in all jobs !!!!!
module load foss/2019a        # load only modules actually needed by job


mpirun -n $SLURM_NTASKS ./a.out
```

Recommendation: reduce queuing times by allowing SLURM to do backfilling:
#SBATCH –time=hh:mm:ss    |    sbatch --time=hh:mm:ss job.sh

# demo HPC literacy

- **make yourself familiar with: login, modules & compiling, job submission**

- `cp –a ~cblass/HPC .` → copy the HPC exercises

- `cd ~/HPC` → go to the folder

```
module load foss/2019a          → @viz.hpc.fs.uni-lj.si    → GCC 8.2.0 & OpenMPI/3.1.3


export MPI_PROCESSES=4          → co-co.c  demo / hello world with  conditional compilation
export OMP_NUM_THREADS=6        → only here (built into coco.c): –DUSE_MPI


cc co-co.c                      mpicc –DUSE_MPI co-co.c
./a.out | sort -n               mpirun -n $MPI_PROCESSES ./a.out | sort -n
sbatch job-serial.sh            sbatch job-mpi.sh


cc -fopenmp co-co.c             mpicc –DUSE_MPI -fopenmp co-co.c
./a.out | sort -n               mpirun -n $MPI_PROCESSES ./a.out | sort -n
sbatch job-openmp.sh            sbatch job-hybrid.sh
```

# Thank you for your attention!

## http://sctrain.eu/

Univerza *v Ljubljani*

TU WIEN | TECHNISCHE UNIVERSITÄT WIEN

CINECA
consorzio interuniversitario

VSB TECHNICAL UNIVERSITY OF OSTRAVA | IT4INNOVATIONS NATIONAL SUPERCOMPUTING CENTER