



# EESSI hackathon

## Using EESSI + Best Practices for using HPC @ UGent

20 Nov 2025

*Kenneth Hoste + Lara Peeters (HPC-UGent)*

*Pieter Asselman (UGent Dept. of Biology)*

# HPC-UGent in a nutshell



- Part of central UGent Functional Domain ICT (formerly DICT)
- Our mission:  
*HPC-UGent provides centralised **scientific computing** services, training, and support for researchers from Ghent University, industry, and other knowledge institutes.*
- Our core values:

Empowerment - Centralisation - Automation - Collaboration

# Centralised hardware in UGent datacenter (S10 @ Sterre)

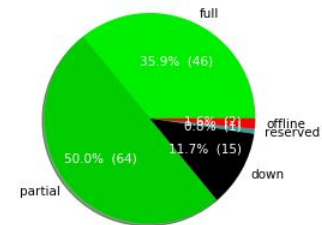


# Terminology: (worker)nodes

- Example cluster from the HPC-UGent  
Tier-2 infrastructure: doduo (*current default cluster*)
- 128 **(worker)nodes**, also referred to as “servers”
- 1 (worker)node is the equivalent of 1 computer  
(but with more cores, memory, faster network, ...)
- Check other HPC-UGent Tier-2 clusters

doduo

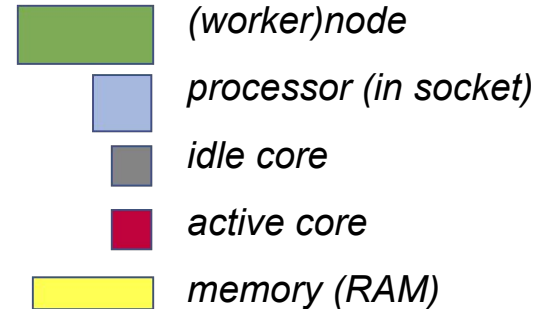
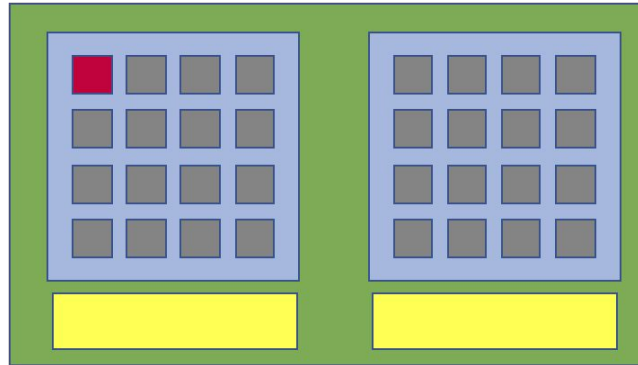
3501	3502	3503	3504	3505	3506	3507	3508
3509	3510	3511	3512	3513	3514	3515	3516
3517	3518	3519	3520	3521	3522	3523	3524
3525	3526	3527	3528	3529	3530	3531	3532
3533	3534	3535	3536	3537	3538	3539	3540
3541	3542	3543	3544	3545	3546	3547	3548
3549	3550	3551	3552	3553	3554	3555	3556
3557	3558	3559	3560	3561	3562	3563	3564
3565	3566	3567	3568	3569	3570	3571	3572
3573	3574	3575	3576	3577	3578	3579	3580
3581	3582	3583	3584	3585	3586	3587	3588
3589	3590	3591	3592	3593	3594	3595	3596
3597	3598	3599	3600	3601	3602	3603	3604
3605	3606	3607	3608	3609	3610	3611	3612
3613	3614	3615	3616	3617	3618	3619	3620
3621	3622	3623	3624	3625	3626	3627	3628



# Terminology: cores, CPUs, processors

Modern servers, also referred to as **(worker)nodes** in the context of HPC, include one or more *sockets*, each housing a **multi-core processor** (next to memory, disk(s), network cards, ...). A modern (micro)processor consists of **multiple cores** that are used to execute computations.

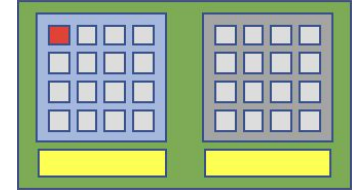
*Example:  
a single workernode  
with two 16-core  
processors running  
a single core job*



*Not shown here: local disk, network cards, GPUs, ...*

# Parallel vs sequential software (single-core)

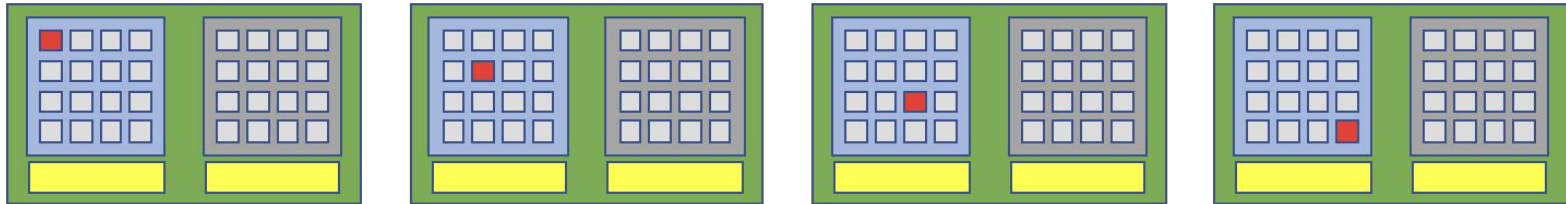
**Sequential** (a.k.a. serial) software does not do calculations in parallel, it only uses one **single core** of a single workernode.



**This type of software does not run faster by just throwing cores (or nodes) at it...**

But, you can run multiple instances of the same program at the same time!

*Example: running a Python script 100 times, each on 1 core, to quickly analyse 100 datasets*



# Parallel vs sequential software (single-node or multi-node)

In **parallel** software, many calculations are carried out simultaneously.

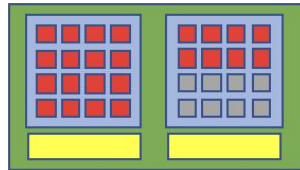
This is based on the principle that large problems can often be divided into smaller tasks, which are then solved concurrently (“in parallel”).

*Example: OpenFOAM can easily use 160 cores at the same time to solve a CFD problem.*

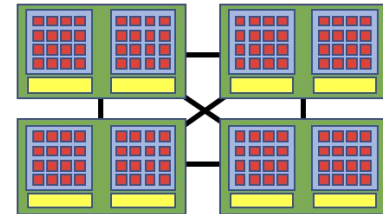
There are two common parallel programming paradigms (among others):

- **OpenMP** for shared memory systems (multi-threading) → using cores of a *single* node
- **MPI** for distributed memory systems (multi-processing) → using cores of *multiple* nodes

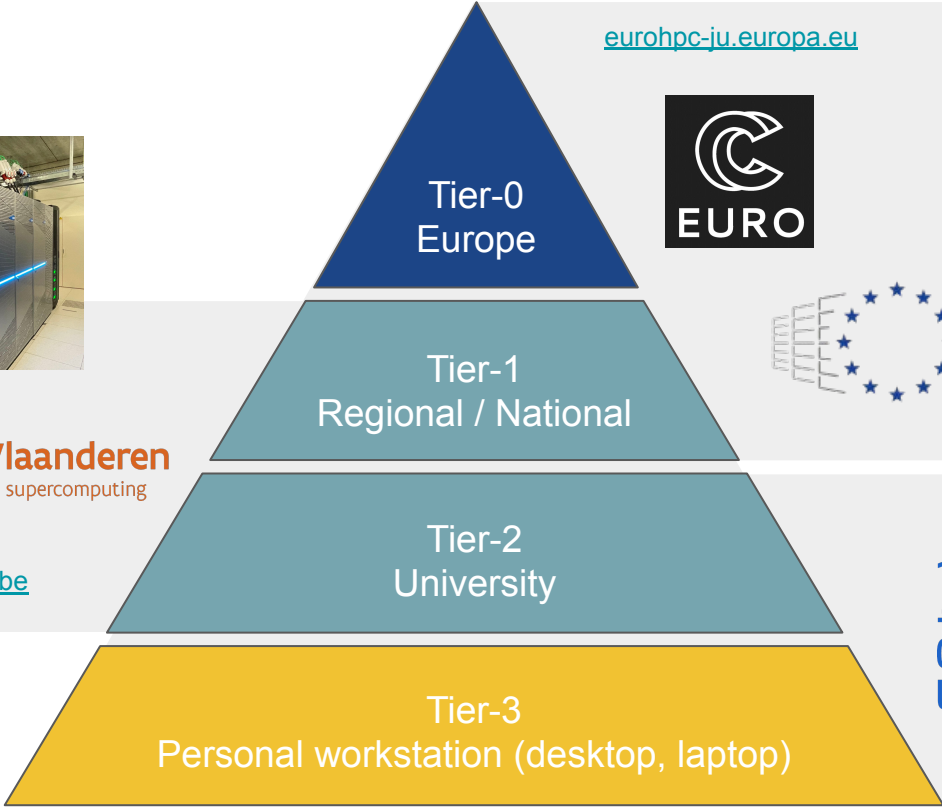
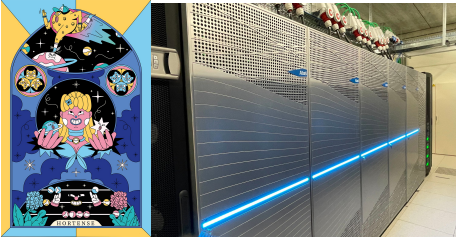
*OpenMP software can use multiple or all cores in a single node*



*MPI software can use (all) cores in multiple nodes*



# Different “tiers” of supercomputers



[eurohpc-ju.europa.eu](http://eurohpc-ju.europa.eu)



**EuroHPC**  
Joint Undertaking

VLAAMS  
SUPERCOMPUTER  
CENTRUM



**Vlaanderen**  
is supercomputing

[vscentrum.be](http://vscentrum.be)



**GHEENT  
UNIVERSITY**

[ugent.be/hpc](http://ugent.be/hpc)



# HPC-UGent Tier-2 infrastructure

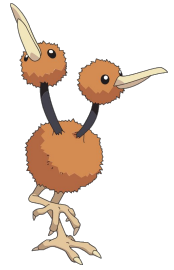
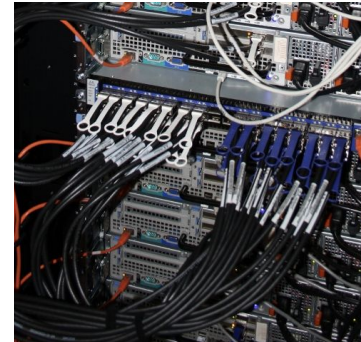


- HPC-UGent Tier-2 infrastructure currently consists of **7 clusters** (+ login nodes, shared storage, ...)
- Different types of clusters:
  - 3 CPU-only compute clusters
  - 3 GPU clusters
  - 1 CPU-only interactive + debug cluster (oversubscribed resources + strict user limits)
- **Available for academic researchers free of charge**, funding through [FWO](#); usage by industry via a pay-as-you-use contract (after free exploratory period)
- All Tier-2 clusters are running Red Hat Enterprise Linux 9 (RHEL9) as operating system



# HPC-UGent Tier-2 compute clusters

- `doduo`: 128 nodes, each with 96 cores (AMD Zen 2Rome) + ~250GiB of memory  
**(default)**
- `gallade`: 16 nodes, each with 128 cores (AMD Zen3 Milan) + ~940GiB of memory
- `shinx`: 48 nodes, each with 192 cores (AMD Zen4 Genoa) + ~370GiB of memory
- All with:
  - High-speed Infiniband network between nodes
  - Fast access to shared file systems
  - Fast local disk (SSD or NVMe)



# HPC-UGent Tier-2 GPU clusters



- `joltik`: 10 nodes,  
each with 32 CPU cores (Intel Cascade Lake),  
**4 NVIDIA V100 GPUs (32GB of GPU memory)**,  
~250GB of system memory
- `accelgor`: 9 nodes,  
each with 48 CPU cores (AMD Zen3 Milan),  
**4 NVIDIA A100 GPUs (80GB of GPU memory)**,  
~500GB of system memory
- `litleo`: 8 nodes,  
each with 48 CPU cores (AMD Zen4 Genoa),  
**2 NVIDIA H100 GPUs (96GB of GPU memory)**,  
~315GB of system memory
- All with high-speed network, fast access to shared filesystems, fast local disk (SSD)



# HPC-UGent Tier-2 interactive + debug cluster: donphan

- 16 nodes, each with 36 CPU cores (Intel Cascade Lake) + ~738GB of memory  
1 *shared* NVIDIA Ampere A2 GPU (16GB of GPU memory)
- Incl. high-speed network, fast access to shared storage, local disk (NVMe)
- Recycled hardware from old `kirlia` cluster (retired in May 2023)
- **Heavily oversubscribed!** More running jobs => All jobs run slower (due to CPU sharing)
- **Strict user limits:**
  - Max. 3 jobs running, 5 jobs in queue
  - Max. 8 cores + 27GB of memory in use (in total)
- => **No waiting time for jobs to start!** Perfect for debug jobs, or interactive use (web portal)
- See also [dedicated section in HPC-UGent documentation](#)



# VSC Tier-2 infrastructure

- You can use your VSC account to access HPC infrastructure provided by other VSC hubs
- Your `$VSC_HOME` and `$VSC_DATA` directories are available on each of these systems



# VSC Tier-1 compute cluster “Hortense”

(a.k.a. dodrio)

VLAAMS  
SUPERCOMPUTER  
CENTRUM



Vlaanderen  
is supercomputing

[compute@vscentrum.be](mailto:compute@vscentrum.be)

- Hosted, operated, and supported by HPC-UGent team since 2021
- 2x 384 CPU-only nodes (128-core AMD Rome or Milan CPUs) + 40 GPU nodes (4x NVIDIA A100)
- **Over 100,000 CPU cores in total!**
- High-speed Infiniband network (HDR-100) + 6PB of dedicated scratch storage



- **Project-based access** (free of charge, funded by FWO)
- 3 cut-off dates per year for submitting project proposals
- Project duration is typically 8 months
- 500k - 5M core hours (CPU-only) or 1k - 25k GPU hours

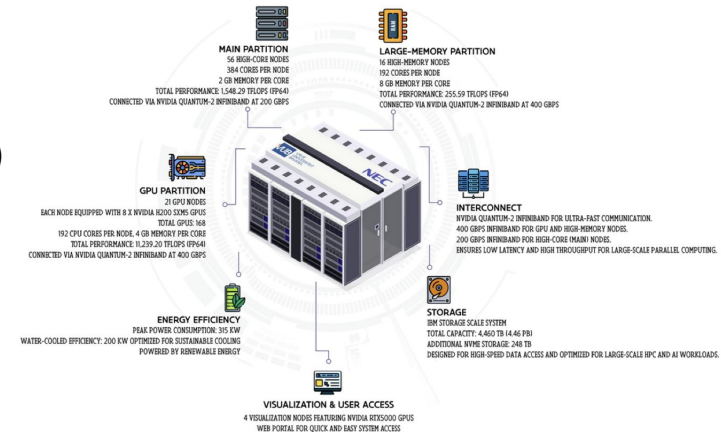
[vscentrum.be/compute](https://vscentrum.be/compute)

[docs.vscentrum.be/en/latest/gent/tier1\\_hortense.html](https://docs.vscentrum.be/en/latest/gent/tier1_hortense.html)

# Next VSC Tier-1 compute cluster

[compute@vscentrum.be](mailto:compute@vscentrum.be)

- Hosted at Green Energy Park in Zellik
- Operated by Vrije Universiteit Brussel (VUB)
- Expected to be ready to use by summer 2026
- Investment of 8.6 million Euro
- **~25,000 CPU cores (with 2GB or 8 GB RAM per core)**
- **168 NVIDIA H200 GPUs**
- **NVIDIA Quantum-2 InfiniBand network at 200 Gbps**
- 4 visualisation nodes
- 4460 TB of scratch storage + 24 TB NVME
- Peak power consumption of 315 kW, 200 kW water-cooled for improved energy efficiency

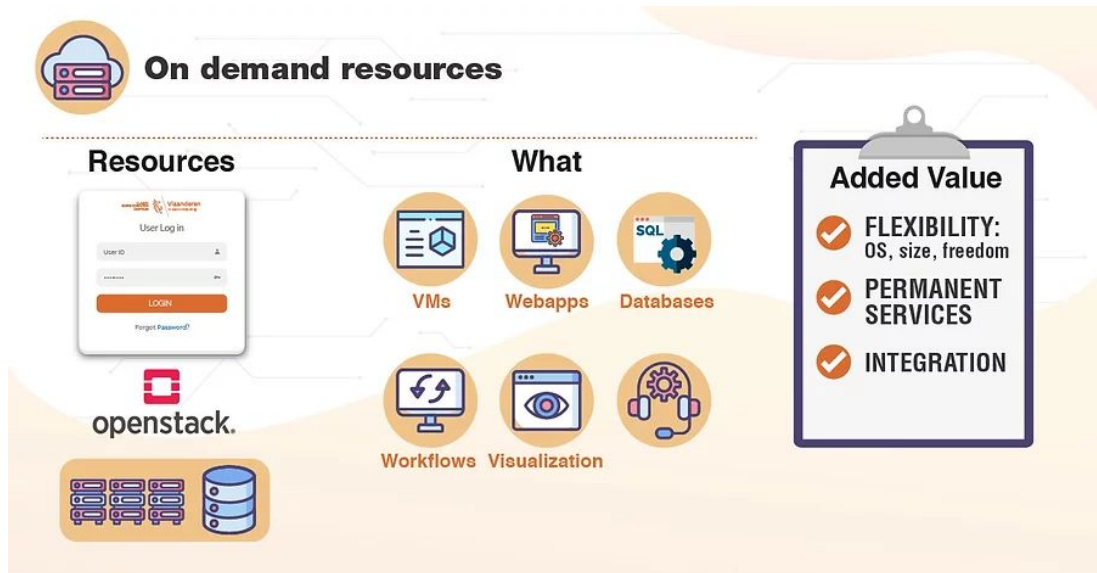


[vscentrum.be/compute](https://vscentrum.be/compute)

[vscentrum.be/post/flanders-invests-in-a-new-supercomputer-to-accelerate-research-and-innovation](https://vscentrum.be/post/flanders-invests-in-a-new-supercomputer-to-accelerate-research-and-innovation)

# VSC Tier-1 cloud

- **Project-based access**
- Free of charge
- **Self-managed virtual machines**
- For use cases that are not a good fit for compute clusters
- More info: [vscentrum.be/cloud](https://vscentrum.be/cloud)
- Contact: [cloud@vscentrum.be](mailto:cloud@vscentrum.be)



# Getting a VSC account



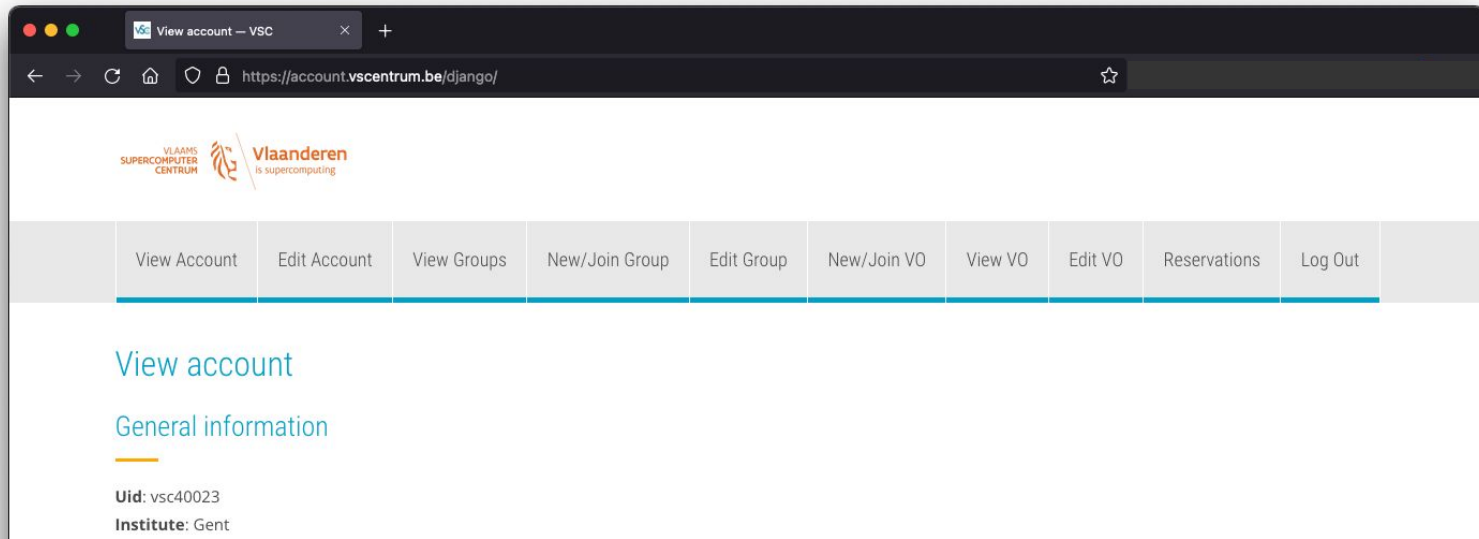
- All members of UGent association can request a VSC account
  - Researchers & staff
  - Master/Bachelor students
- **VSC account can be used to access HPC infrastructure on all VSC sites**
- Subscribed to `hpc-announce` and `hpc-users` mailing lists
- Beware of using HPC for teaching/exam purposes!
  - No guarantee on HPC availability (due unexpected power outage, maintenance, ...)
  - Have a backup plan at hand
  - Advisable teaching/exam formula: project work
- See also [HPC-UGent documentation](#)

# Managing your VSC account

You can manage your VSC account via the VSC account page

[account.vscentrum.be](https://account.vscentrum.be)

Can be used to join/leave user groups, consult storage usage, request more storage quota, ...  
manage your Virtual Organisation (VO), ...



The screenshot shows a web browser window with the address bar displaying `https://account.vscentrum.be/django/`. The page header includes the VSC logo and the text "Vlaanderen is supercomputing". Below the header is a navigation menu with the following items: View Account, Edit Account, View Groups, New/Join Group, Edit Group, New/Join VO, View VO, Edit VO, Reservations, and Log Out. The main content area displays the heading "View account" and "General information". Under "General information", the following details are shown: **Uid:** vsc40023 and **Institute:** Gent.

# High-level overview of HPC-UGent infrastructure

Getting access

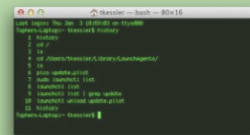
<https://login.hpc.ugent.be>

Web Portal (No SSH key required)



OR

ssh login.hpc.ugent.be  
Terminal (requires SSH key)

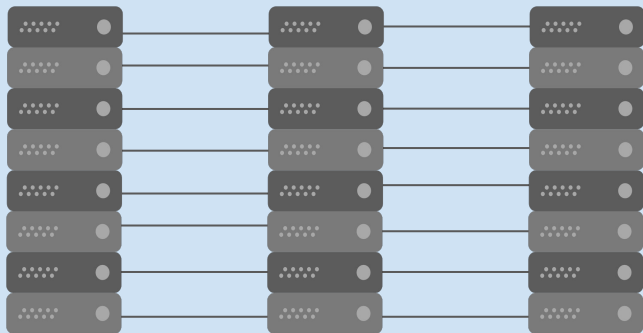


Login nodes



gligar09 or gligar10

Compute nodes (where your jobs will run)



File system

Data



Home



Data



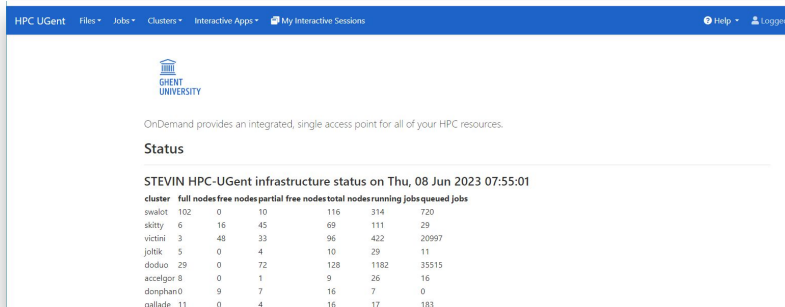
Scratch



Arcanine Scratch  
(SSD)

# Connecting to the HPC-UGent login nodes with web portal

Recommended!



The screenshot shows the HPC-UGent web portal interface. At the top, there is a navigation bar with 'HPC UGent' and several menu items: 'Files', 'Jobs', 'Clusters', 'Interactive Apps', and 'My Interactive Sessions'. Below the navigation bar is the Ghent University logo and a message: 'OnDemand provides an integrated, single access point for all of your HPC resources.' The main content area is titled 'Status' and displays the 'STEVIN HPC-UGent infrastructure status on Thu, 08 Jun 2023 07:55:01'. Below this title is a table with columns for 'cluster', 'full nodes', 'free nodes', 'partial free nodes', 'total nodes', 'running jobs', 'queued jobs', and 'jobs'.

cluster	full nodes	free nodes	partial free nodes	total nodes	running jobs	queued jobs	jobs
swalot	102	0	10	116	314	720	
skitty	6	16	45	69	111	29	
victini	3	48	33	96	422	20997	
jolke	5	0	4	10	29	11	
diadoa	29	0	72	128	1182	35515	
accelgor	8	0	1	9	26	16	
donphan0	9	7		16	7	0	
gallade	11	0	4	16	17	183	

- See [dedicated section of HPC-UGent docs](#)
- <https://login.hpc.ugent.be>
- Powered by [Open OnDemand](#)

- **Works with a standard internet browser** (Firefox, Chrome, ...)
- **Does not require SSH key pair** (only login via UGent account)
- Provides file browser, shell session, desktop environment, interactive apps, ...

# UGent web portal: interactive apps (Jupyter notebook)

1

HPG UGent Files Jobs Clusters Interactive Apps My Interactive Sessions

Desktops

- Bioimage AAnalysis Desktop
- Cluster Desktop
- Neurodesk

Servers

- Jupyter Lab
- Jupyter Notebook**
- RStudio server
- Shell (tmux)
- VS Code Tunnel

Testing

- Cluster desktop v2
- Code Server

On THU, 19 Sep 2024 12:4

cluster	full nodes free	running jobs	queued jobs
skitty	2	74	2
joltik	4	9	76
doduo	39	996	5376
accelgon	6	21	5082
donphan	0	14	14
gallade	3	9	39
shinx	3	43	351

2

Home / My Interactive Sessions / Jupyter Notebook

Interactive Apps

Desktops

- Bioimage AAnalysis Desktop
- Cluster Desktop
- Neurodesk
- Jupyter Lab
- Jupyter Notebook**
- RStudio server
- Shell (tmux)
- VS Code Tunnel
- Code Server

Jupyter Notebook

This app will launch a Jupyter Notebook server on one or more nodes.

Cluster: donphan (interactive/debug)

Time (hours): 1 hour

Number of nodes: 1 node

Number of cores (and default memory) per node: 1 core (3.3 GiB mem)

Jupyter Notebook version: 7.2.0 GCCore 13.2.0

Launch

3

Jupyter IPython Notebook (cluster/donphan-20004994) Queued

Created at: 2023-06-07 14:51:43 CEST

Time Requested: 12 hours

Session ID: dd4a26c3-a09e-44f3-b8a4-a558a99bdfc9

Please be patient... job currently sits in queue. The wait time depends on the number of cores as well as time requested.

Jupyter IPython Notebook (cluster/donphan-20004994) 1 node | 4 cores | Starting

Created at: 2023-06-07 14:51:43 CEST

Time Remaining: 11 hours and 59 minutes

Session ID: dd4a26c3-a09e-44f3-b8a4-a558a99bdfc9

Your session is currently starting... Please be patient as this process can take a few minutes.

4

Jupyter IPython Notebook (cluster/donphan-20004994) 1 node | 4 cores | Running

Host: >\_node0116.donphan.os

Created at: 2023-06-07 14:51:43 CEST

Time Remaining: 11 hours and 59 minutes

Session ID: dd4a26c3-a09e-44f3-b8a4-a558a99bdfc9

Connect to Jupyter Notebook

5

jupyter

Files Running Cluster

Select items to perform actions on

0 /

- examples
- ondemand

Quit Logout

Upload New

Name Last Modified File size

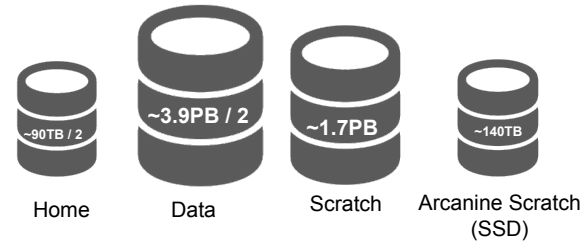
6 days ago

16 minutes ago

**Jupyter Notebook versions**

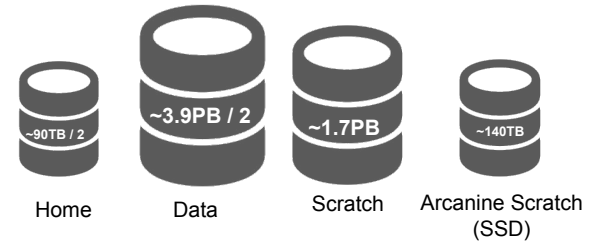
Make sure that the toolchain of the notebook matches the toolchain of the modules that you load, and the kernels and/or virtual environments that you make!

# Input/output data and shared filesystems



- See [dedicated section in HPC-UGent documentation](#)
- Think about input/output:
  - How and where will you *stage in* your data and input files?
  - How and where will you *stage out* your output and result files?
- Manually (on login nodes) vs automatically (as a part of job script)
- **Home filesystem** (`$VSC_HOME`): only for limited number of small files & scripts
- **Data filesystem** (`$VSC_DATA*`): 'long-term' storage, large files
- **Scratch filesystems** (`$VSC_SCRATCH*`): for 'live' input/output data in jobs

# Storage quota (disk space)



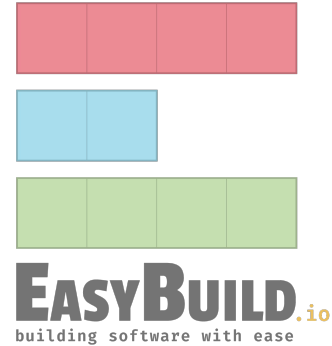
- Home directory (`$VSC_HOME`): 3GB (fixed!)
- Personal data directory (`$VSC_DATA`): 25GB (fixed!)
- Personal scratch directory (`$VSC_SCRATCH`): 25GB (fixed!)
- Current quota usage can be consulted on [VSC accountpage](#)
- **More storage quota (100s of GBs, even TBs) available for *virtual organisations (VOs)*; see [dedicated section on VOs in HPC-UGent documentation](#)**
- Additional quota can be requested via [VSC accountpage \(“Edit” tab\)](#)
- Shared directories with VO members: `$VSC_DATA_VO`, `$VSC_SCRATCH_VO`
- Personal VO subdirectories: `$VSC_DATA_VO_USER`, `$VSC_SCRATCH_VO_USER`

# Job scheduling

- All HPC-UGent clusters use a **fair-share scheduling** policy.
- No guarantees on when job will start (and impossible to predict), so **plan ahead!**
- Job priority is determined by various factors:
  - Historical usage
    - Aim is to balance usage over users
    - Infrequent/frequent users => higher/lower priority
  - Requested resources (# nodes/cores, walltime, memory, ...)
    - Larger resource request => lower priority
  - Time waiting in queue
    - Queued jobs get higher priority over time
  - User limits
    - Avoid that a single user fills up an entire cluster

# Software installations

- To submit a request for software installation, use the request form:  
[docs.hpc.ugent.be/software\\_installation\\_requests](https://docs.hpc.ugent.be/software_installation_requests)
- Requests may take a while to process (especially for new software), so **be patient...**
- Make the request sooner rather than later!
- All software installations are done using EasyBuild
- Originally developed by HPC-UGent,  
now a worldwide community of experts!
- See also [easybuild.io](https://easybuild.io)



# Questions, problems, getting help

**Don't hesitate to contact the HPC-UGent Tier-2 support team via [hpc@ugent.be](mailto:hpc@ugent.be)**

- Always include:
  - VSC login id
  - Clear description of the problem or question, include error messages, ...
  - Location of job script and output/error files in your VSC account
  - Preferably don't send files in attachment, we prefer to look at it "in context"
  - Also mention job IDs, which cluster was used, how job was submitted, etc.
- Preferably use your UGent email address
- Alternatives:
  - Short (Teams) meeting (for complex problems, big projects)
  - `hpc-users` mailing list

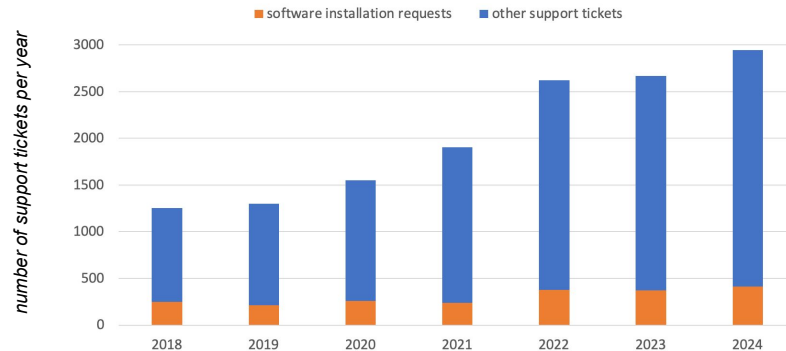
# Questions, problems, getting help (be patient...)

Don't hesitate to contact the HPC-UGent support team via [hpc@ugent.be](mailto:hpc@ugent.be),

but be patient...



- We're doing what we can to keep up with incoming support questions + software installation req.
- We have been getting **50-100 new tickets per week** recently, and have a backlog of ~100 tickets
- **Help us help you:** read the docs, provide sufficient details (like job IDs, output files, etc.), ...
- Feel free to send a reminder in the same ticket or mail thread, especially if your work is blocked



*What if you no longer have to install  
a **broad range of scientific software**  
from scratch on every laptop, HPC cluster,  
or cloud instance you use or maintain,  
**without compromising on performance?***



# European Environment for Scientific Software Installations (EESSI)

- **Shared repository of (optimized!) scientific software installations**
- Avoid duplicate work across by collaborating on a shared software stack
- Uniform way of providing software to users, regardless of the system they use!
- Should work on any Linux OS and system architecture
  - From laptops and personal workstations to HPC clusters and cloud
  - Support for different CPUs, interconnects, GPUs, etc.
- Focus on **performance, automation, testing, collaboration**
- Development effort funded through **MultiXscale** EuroHPC Centre-of-Excellence



**E E S S I**

EUROPEAN ENVIRONMENT FOR  
SCIENTIFIC SOFTWARE INSTALLATIONS

<https://eessi.io>

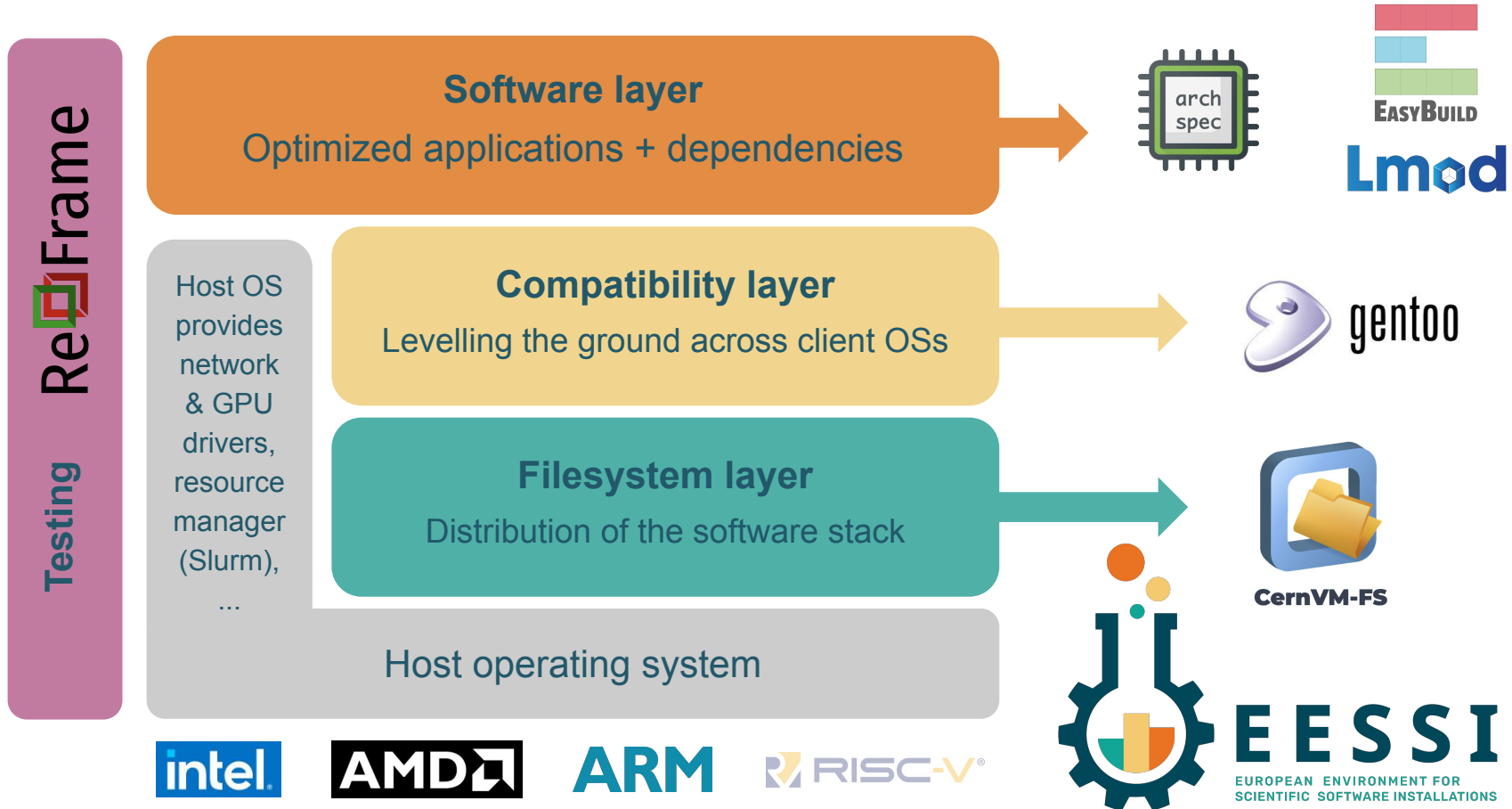
**MultiXscale**

# Major goals of EESSI



- Providing a truly **uniform software stack**
  - Use the (exact) same software environment everywhere
  - **Without sacrificing performance** for “mobility of compute” (like is typically done with containers/conda)
- **Avoid duplicate work** (for researchers, HPC support teams, sysadmins, ...)
  - Tools that automate software installation process (EasyBuild, Spack) are not sufficient anymore
  - Go beyond sharing build recipes => work towards a shared software stack
- Facilitate HPC training, development of (scientific) software, ...

# High-level overview of EESSI



# EESSI ingredients



gentoo linux™

## Compatibility layer

Abstraction from the host operating system



## Filesystem Layer

Global distribution of software installations via CernVM-FS



# E E S S I

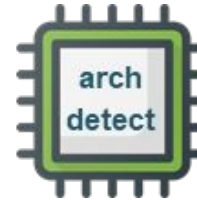
EUROPEAN ENVIRONMENT FOR  
SCIENTIFIC SOFTWARE INSTALLATIONS

## Software Layer



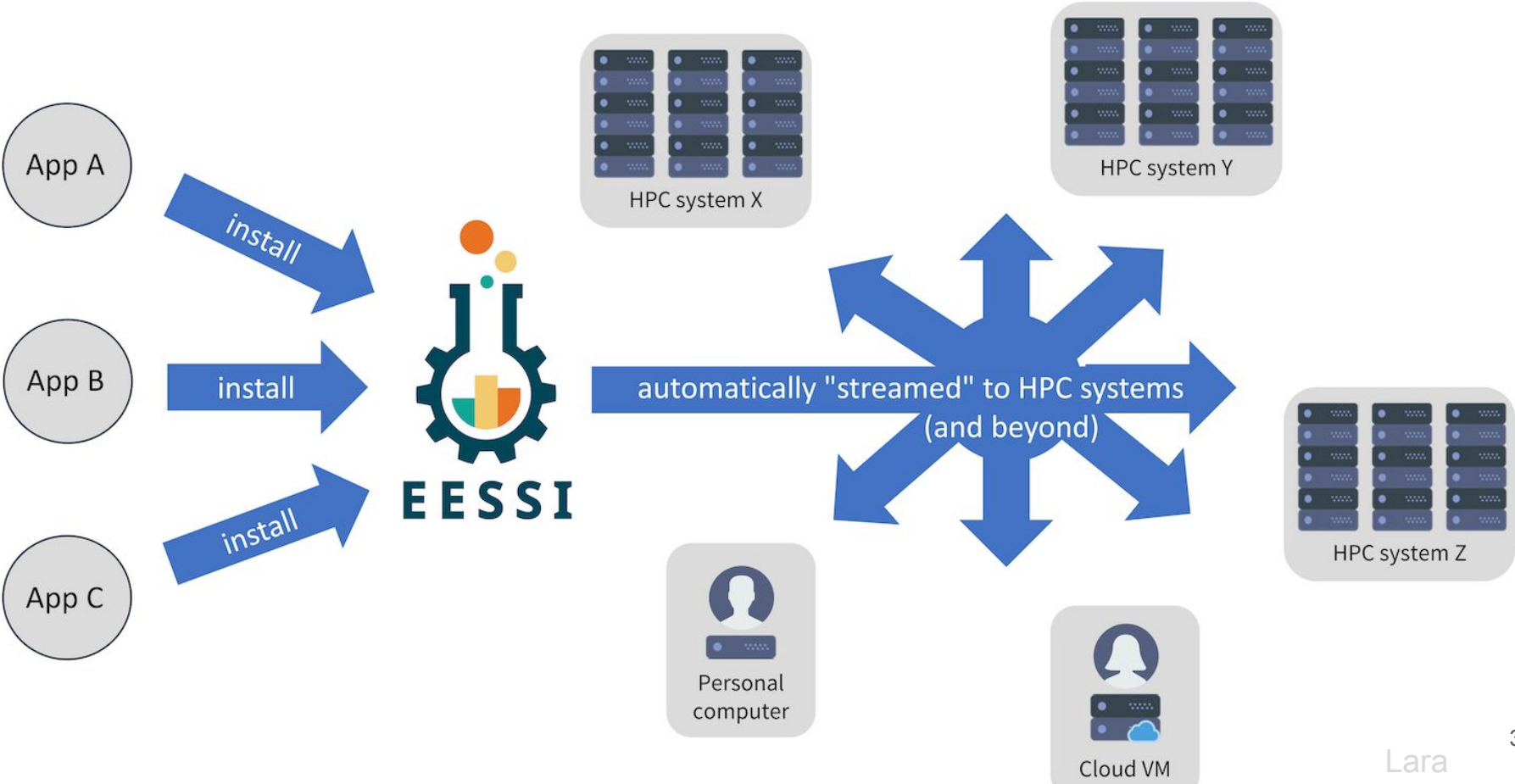
**Optimized** software installations for specific CPU microarchitectures

Intuitive user interface:  
module avail,  
module load, ...



Automatic selection of best suited part of software stack for CPU microarchitectures

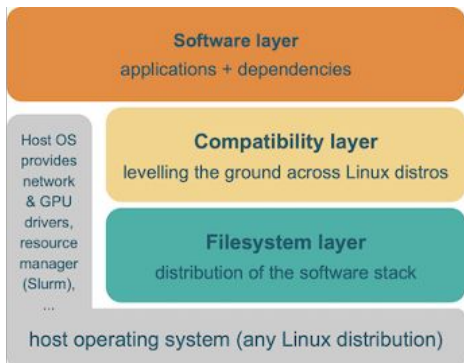
# EESSI as a shared software stack



# How does EESSI work?

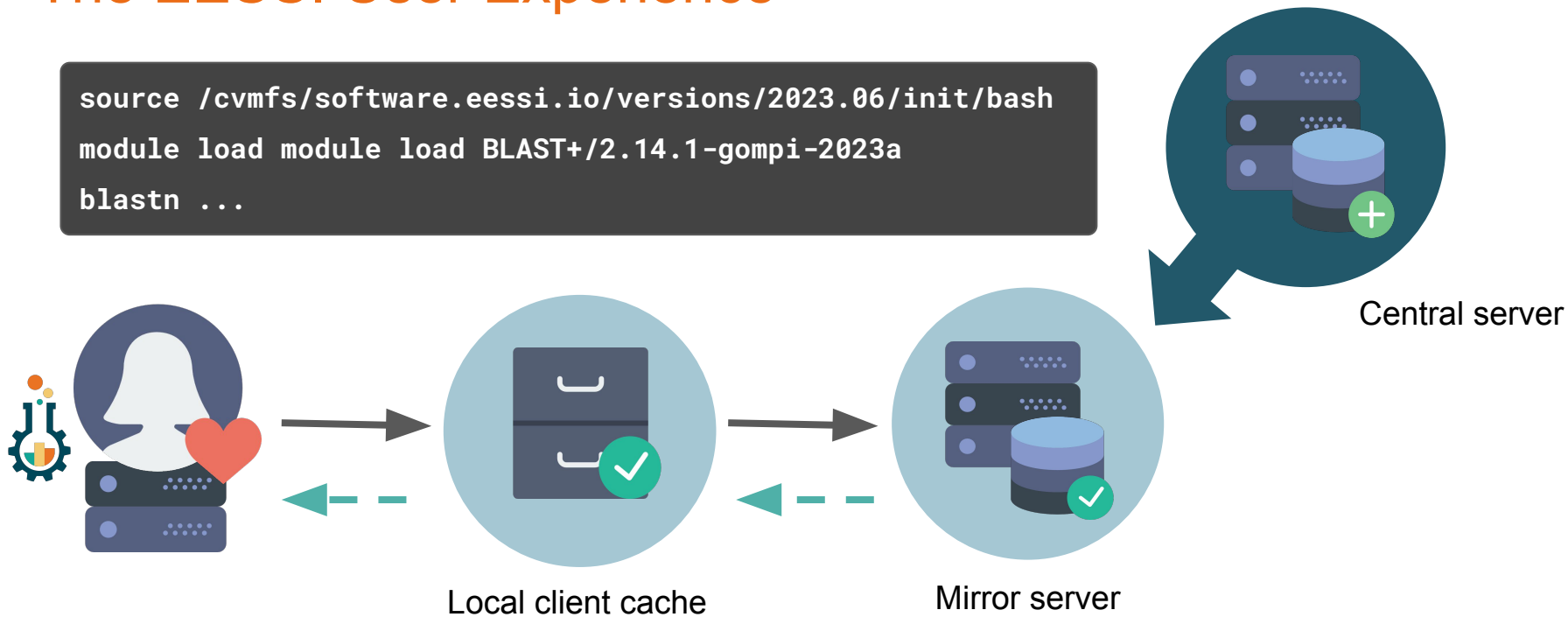


- Software installations included in EESSI are:
  - Automatically **“streamed in” on demand** (via CernVM-FS)
  - Built to be **independent of the host operating system**  
*“Containers without the containing”*
  - **Optimized** for specific CPU generations + specific GPU types
- Initialization script **auto-detects** CPU + GPU of the system



# The EESSI User Experience

```
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
module load module load BLAST+/2.14.1-gompi-2023a
blastn ...
```



EESSI provides **on-demand streaming**  
of (scientific) software (like music, TV-series, ...)

# Demo: Using EESSI on HPC-UGent Tier-2

```
# Access via HPC-UGent web portal (https://login.hpc.ugent.be)
# Start interactive shell session on donphan cluster
module --force purge # this is not required, but it is recommended to start afresh
source /cvmfs/software.eessi.io/versions/2023.06/init/bash

Found EESSI repo @ /cvmfs/software.eessi.io/versions/2023.06!
archdetect says x86_64/intel/cascadelake
archdetect found supported accelerator for CPU target x86_64/intel/cascadelake: accel/nvidia/cc80
Using x86_64/intel/cascadelake as software subdirectory.
Using /cvmfs/software.eessi.io/versions/2023.06/software/linux/x86_64/amd/cascadelake/modules/all
Using ...
...
Environment set up EESSI (2023.06), have fun!
module avail BLAST
module load BLAST+/2.14.1-gompi-2023a
```

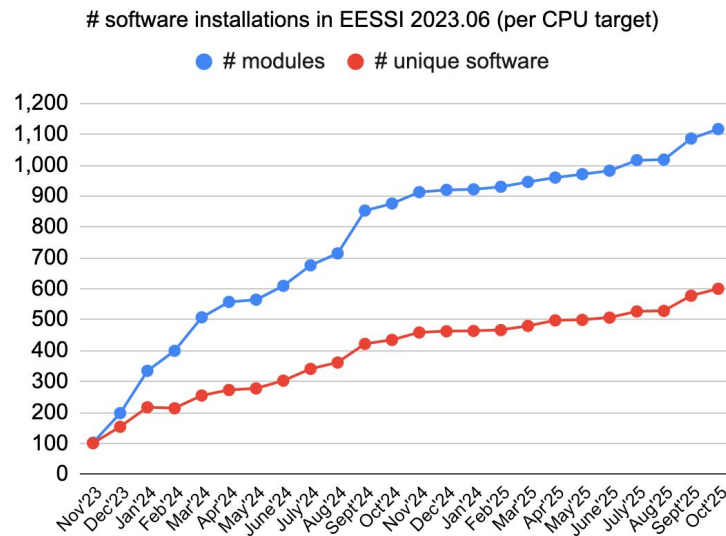
Alternative ways of accessing EESSI are available, via a container image, via cvmfsexec, ...  
[eessi.io/docs/getting\\_access/native\\_installation](https://eessi.io/docs/getting_access/native_installation) - [eessi.io/docs/getting\\_access/eessi\\_container](https://eessi.io/docs/getting_access/eessi_container)

# Overview of available software in EESSI

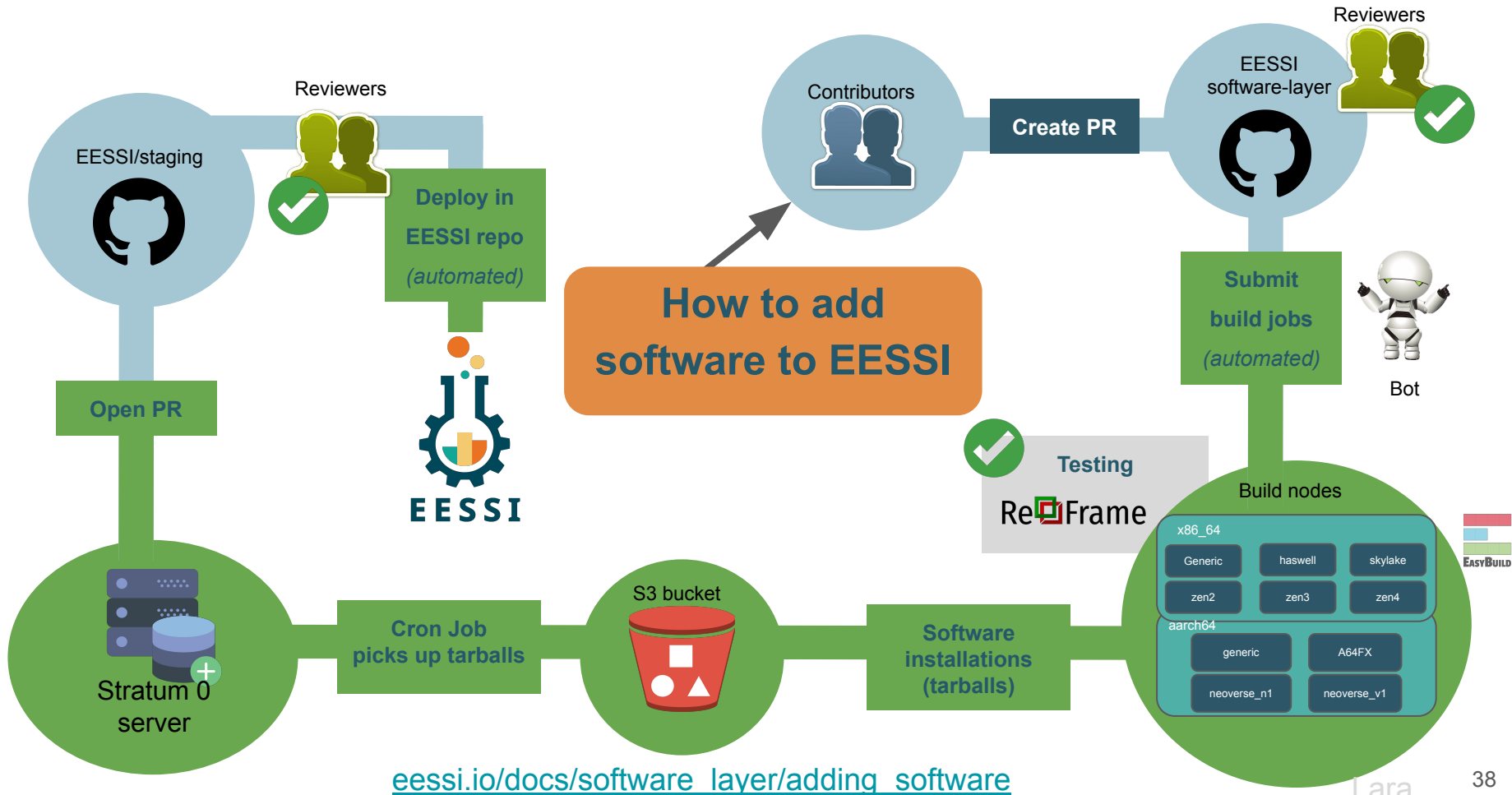


Currently more than 1,100 software installations available per supported CPU target via `software.eessi.io` CernVM-FS repository, increasing every week

- **14 fully supported CPU targets** (x86\_64 + Arm), see [https://eessi.io/docs/software\\_layer/cpu\\_targets](https://eessi.io/docs/software_layer/cpu_targets)
- **~600 different software packages**, excl. extensions: Python packages, R libraries
- **Almost 16,000 software installations in total**
- Including **BLAST+**, ESPResSo, **Flye**, GROMACS, LAMMPS, **MrBayes**, OpenFOAM, PyTorch, R, QuantumESPRESSO, TensorFlow, waLBerla, WRF, ...
- [https://eessi.io/docs/available\\_software/overview](https://eessi.io/docs/available_software/overview)
- Using `foss/2023a` and `foss/2023b` toolchains in EESSI 2023.06
- Using `foss/2024a` and `foss/2025a` toolchains in EESSI 2025.06



# Semi-automated workflow for adding software to EESSI



# Building software on top of EESSI



- You can use the software available in EESSI to **resolve dependencies** of the software you need
- To install software with EasyBuild on top of EESSI, first load the **E E S S I - e x t e n d** module to correctly configure EasyBuild
- To manually install software on top of EESSI, first load the **b u i l d e n v** module
- See also [https://eessi.io/docs/using\\_eessi/building\\_on\\_eessi](https://eessi.io/docs/using_eessi/building_on_eessi)

# On which systems is EESSI available?



- On VSC systems:
  - **Tier-2 infrastructure at UGent + VUB + Tier-1 Hortense (+ new Tier-1 soon)**
  - **Tier-1 cloud**: can deploy it yourself and have access to a full software stack in minutes
- EESSI is already available on various other European systems (and beyond)
  - EuroHPC JU systems incl. Vega, Karolina, MareNostrum 5, Deucalion, Discoverer, ...
  - Snellius @ SURF, EMBL, Univ. of Stuttgart, Sigma2 in Norway, etc.
- EESSI can be used in virtual machine in European Open Science Cloud (EOSC), see also <https://www.eessi.io/docs/blog/2025/10/22/eosc>
- **Overview of (known) systems that have EESSI available at <https://eessi.io/docs/systems>**

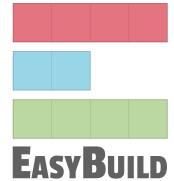
# Webinar series: Different aspects of EESSI

## Series of presentations (May-June 2025)

<https://eessi.io/docs/training/2025/webinar-series-2025Q2>

- Introduction to EESSI
- Introduction to CernVM-FS
- Introduction to EasyBuild
- EESSI for CI/CD
- Using EESSI as the base for a system stack

**Slides + recordings available**



# Best practices for running jobs @ HPC-UGent

- **Read the documentation** (HPC-UGent, software, ...)
- **Check resource usage** (CPUs, GPUs, memory, ...)
- Take into account **I/O overhead**: try different shared filesystems
- Use **more resources** (cores/nodes/memory), but check if that helps! (**scaling**)
- Consider using **different clusters** (newer is usually faster)
- Try using **recent software versions**: newer is usually better (but take care)
- Check details of **how the software was installed**:  
conda vs central software stack vs EESSI vs containers
- **Don't waste time over-optimizing...**

# Resources?

## Different types of resources

- CPU **cores** (# cores, a.k.a. ppn), also # **nodes** and # **GPUs** (if relevant)
- **Memory** per node (vmem)
- **Walltime**: amount of times that passes between start and end of width
- Also: (derived) memory per *core*, (local) storage space, network bandwidth, ...

# Some tips w.r.t. requesting resources



- **Small(er) resource request** in terms of cores/nodes/memory/walltime  
=> Usually **less waiting time** for jobs in queue  
(since a smaller “gap” is required)
- **Only request what you need**, and can **use efficiently**
  - Don't blindly request more cores, check if it results in (significant) **speedup!**
  - Evaluate impact of requesting more resources on job walltime
- **Total wait time** (for you) = **queue time + (wall)time required to complete job**
- Rule of thumb: if possible, use **quarter of a node** (or less) in terms of cores, helps to minimize waiting time in queue

# Checking resource usage: `slurm_jobinfo`

```
$ slurm_jobinfo 12345678
Name           : example
User           : vsc40000
Partition      : shinx
Nodes          : node4201.shinx.os
Cores          : 8
State          : COMPLETED
Submit        : 2025-11-10T10:58:47
Start         : 2025-11-10T11:01:58
End           : 2025-11-10T18:49:44
Reserved walltime : 3-00:00:0.0
Used walltime  : 07:47:46.0
Used CPU time  : 07:38:54.0
% User (Computation) : 82.03
% System (I/O)       : 17.97
Mem reserved      : 15G
Max Mem used      : 14.97G (node4201.shinx.os)
```

Specify job ID  
(make sure correct cluster  
module is loaded)

Used CPU time should be  
roughly equal to # cores  
times used walltime

Large ratio of system CPU  
is not good

# Scaling up: use more resources, but...



- More cores should result in **significant speedup** (but maybe less than you expect), if the software supports multi-threading or multi-processing
- Check **used walltime vs used CPU time** for your jobs
  - CPU time should be close to `walltime`  $\times$  `number-of-cores`
- Take into account **efficiency** (speedup divided by # cores), aim for min. ~50-60%
  - Avoid wasting resources, and longer waiting times in queue
- Requesting **multiple nodes** only makes sense if your software supports **multi-processing** (typically via MPI)
- Make sure the software you are using can make **efficient use of GPUs** (if at all)

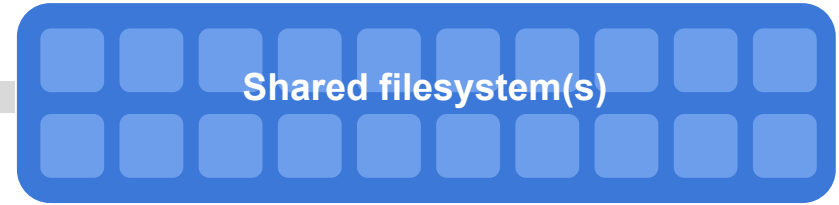
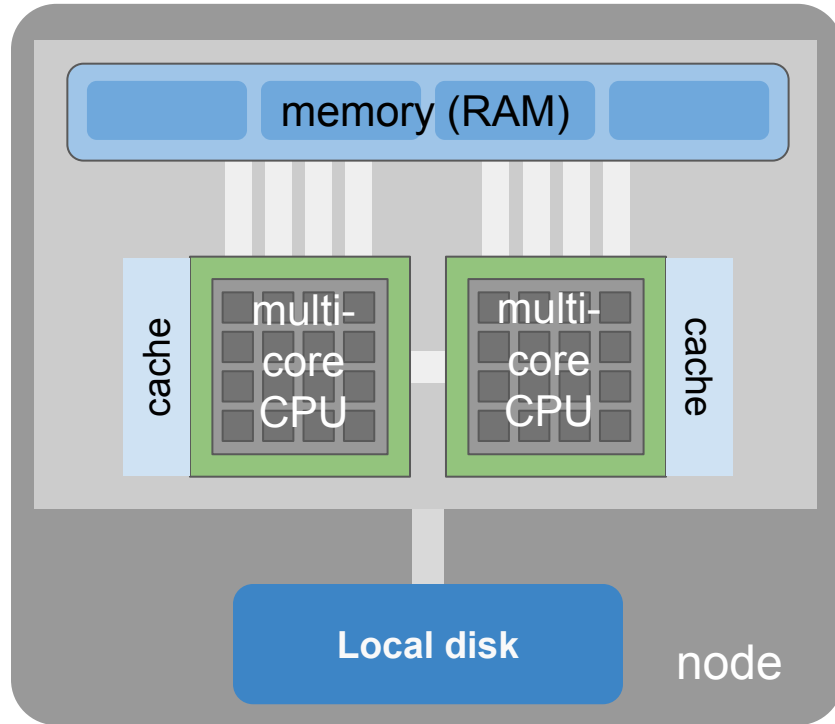
# Checking resource usage: live job inspection

- Sometimes it can be interesting to inspect the job **while it is running**...
- Check with `qstat -ant` on which node(s) your job is running
- Log into the node with `ssh node1234`
- A couple of useful commands:
  - `top, htop, btop` → total resource usage on node (quit with `q` or `Ctrl-C`)
  - `ps xfu` → list of running processes (for your VSC account), incl. CPU usage
  - `taskset -c -p 12345` → core pinning for process with ID 12345
  - `pidstat -u -p 12345 1` → CPU usage for with ID 12345, refresh every 1 sec.

# The \$OMP\_PROC\_BIND “curse”

- \$OMP\_PROC\_BIND is automatically set to TRUE in all jobs on HPC-UGent Tier-2 (and VSC Tier-1 @ UGent)
- Significant performance boost for *some* workloads that use OpenMP threading
- However, some other workloads (R, Python, BLAST+) get “locked” to a single core during startup of the software, and can not escape that jail anymore at runtime...
  - We reported a bug on this for the OpenMP runtime in GCC, but got told this is by design... see [https://gcc.gnu.org/bugzilla/show\\_bug.cgi?id=114765](https://gcc.gnu.org/bugzilla/show_bug.cgi?id=114765)
- **You may need to use “unset OMP\_PROC\_BIND” in your job script** to avoid that all threads are pinned to a single core... **Check scaling when using multiple cores!**
- If pinning of threads to cores is important, you should use the Likwid tool, see <https://github.com/RRZE-HPC/likwid>

# High-level overview of a (worker)node



Filesystems @ HPC-UGent Tier-2:

`$VSC_HOME`

“slow” & small

`$VSC_DATA`

“slow” & large (+ safe)

`$VSC_SCRATCH`

fast & large

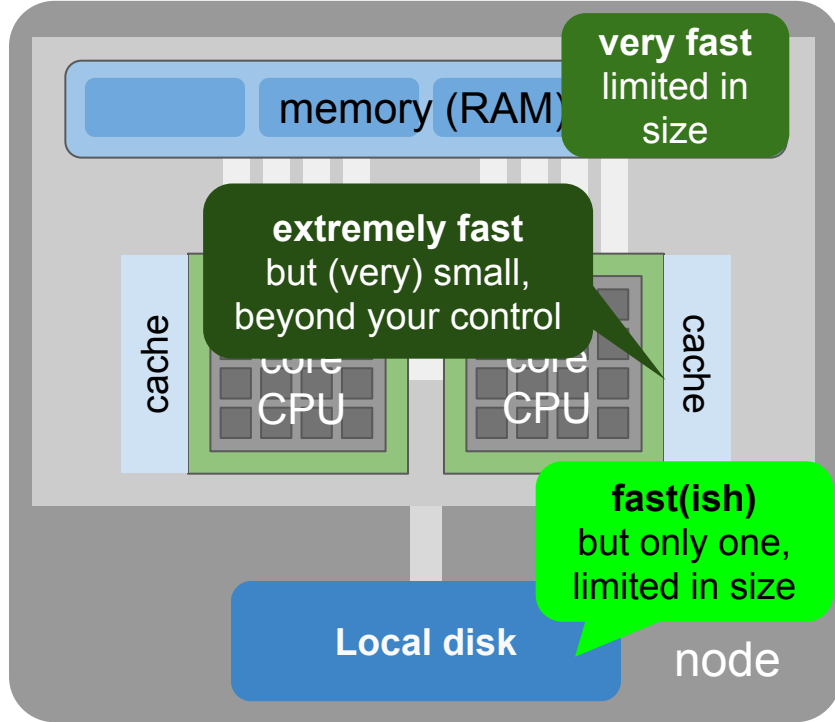
`$VSC_SCRATCH_ARCANINE`

faster & medium size

Types of disk:  
HDD | SSD | NVME

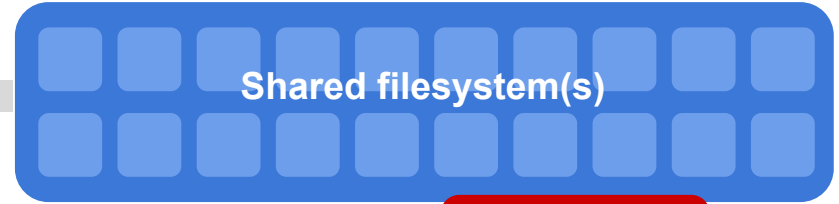
slow → fast

# Speed in which data can be read (or written)



Types of disk:  
HDD | SSD | NVME

slow → fast



Filesystems @ HPC-UGent

`$VSC_HOME`

**Don't use this for data**

"slow" & small

`$VSC_DATA`

**slow(ish)**

"slow" & large (+ safe)

`$VSC_SCRATCH`

**fast**

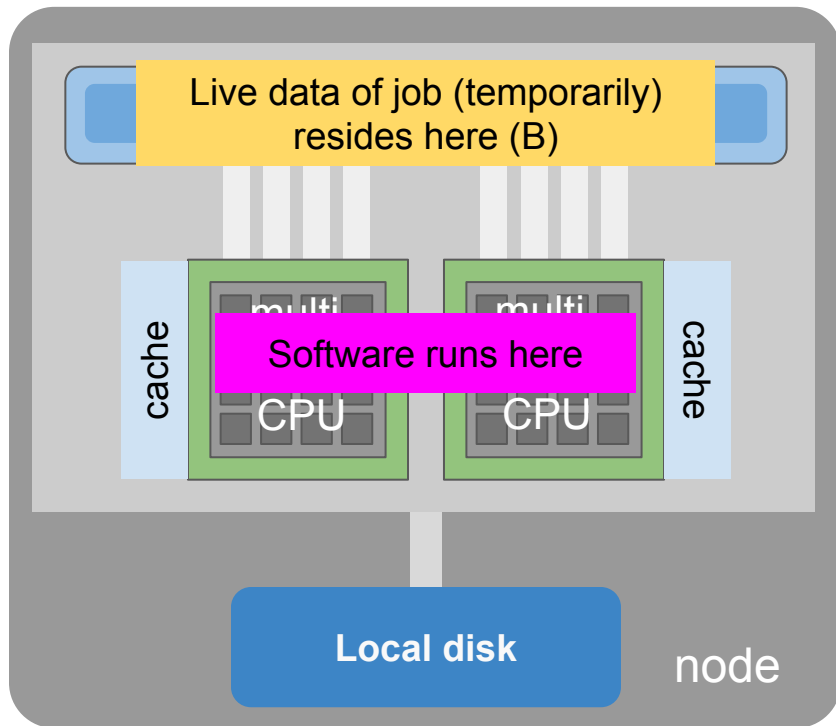
(& large) fast & large

`$VSC_SCRATCH_ARCANINE`

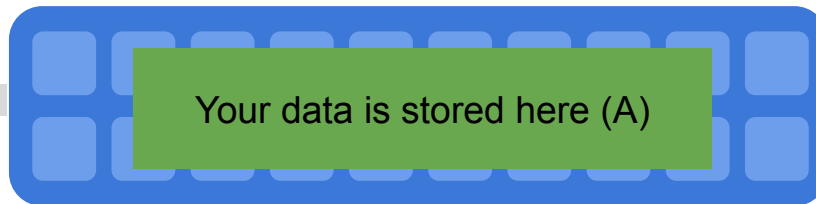
**faster**

medium size

# Where does your data come from?



**How the (input) data for your jobs gets from A to B can matter (a lot), especially when it's large volume (GBs), lots of (small) files, or is accessed "randomly".**



Filesystems @ HPC-UGent Tier-2:

<code>\$VSC_HOME</code>	"slow" & small
<code>\$VSC_DATA</code>	"slow" & large (+ safe)
<code>\$VSC_SCRATCH</code>	fast & large
<code>\$VSC_SCRATCH_ARCANINE</code>	faster & medium size

Types of disk:  
HDD | SSD | NVME

slow → fast

# I/O overhead: time spent on reading/writing data

- **Reading data from disk (shared filesystems) can be “slow”**
  - CPUs may not be able to get useful work done because of that...
- Correlates with “% System (I/O)” in `slurm_jobinfo` output
  - Should be  $\ll 5\%$ , or it will significantly slow down your job
  - Can be difficult to avoid for some type of workloads, but worth paying attention to
- Live job inspection with `htop`, `pidstat` can help you understand this better
- If used CPU time is not roughly equal to number of cores x used walltime, a lot of time is being lost to... something
- **Consider accessing the data through a different/faster shared filesystem:**  
`$VSC_SCRATCH`, `$VSC_SCRATCH_ARCANINE`

# Escaping I/O overhead by using local or RAMdisk



- If consuming data from shared filesystems is too slow (large I/O overhead), you can try using **local (SSD) disk or RAMdisk** (mostly interesting for single-node jobs)

```
export INPUT_DATA=$VSC_DATA_VO/large_dataset
export FAST_TMPDIR=/dev/shm/$USER # or $TMPDIR (local disk)
rsync -avP $INPUT_DATA $FAST_TMPDIR
export INPUT_DATA=$FAST_TMPDIR/large_dataset
do_work_with $INPUT_DATA
```

- Streaming copy from shared filesystem is usually fast (vs random access during job)
- Take into account that:
  - Time needed to copy the data is part of job walltime, so only interesting for long jobs
  - Size of the dataset may make this impossible (but may be possible on gallade)



- **How** the software was compiled/installed can impact **performance** a lot
- Software should be **optimized** for the specific CPUs on which it will be run
  - This is usually not the case when using container images or Conda environments...
  - Generic binaries are used that can run anywhere
  - Performance is sacrificed for mobility of compute...
- Vast majority of software installed in central software stack on HPC-UGent infrastructure is optimized for the specific CPUs of each cluster
  - Likewise for EESSI: broad range of CPUs supported, specific CPU is auto-detected
- Additional caveats with Conda: messes up your home directory very easily...



- You are working on a shared system, with **shared resources** => chaos...
- **Some things are outside of your control**
  - Load on shared filesystems varies a lot, you may see peaks in I/O overhead
  - Placement of jobs w.r.t. cores & sockets, memory channels, ...
  - Job(s) by other users running on same node may indirectly impact your job
- **CPU clock speed varies over time** based on load on CPU cores (throttling/boosting)
- Make sure you're running the **same workload** when comparing job walltimes!
  - Watch out with “converging” workloads...
  - Some software is “smart”, will detect or change things at runtime...

3. Do you have experience with installing software on a linux machine?Required to answer?

[More Details](#)



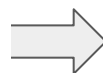
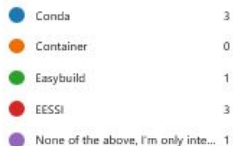
4. If you answered "Yes" to question 3, which method did you use?Required to answer.

[More Details](#)



5. Are you interested to see a quick demo on how to install software on HPC yourself via:

[More Details](#)



**“Do you use HPC for you BioInfo analyses?”**

**YES!**

*No, too many limitations.*

*Yes, GPU power!*

*Sometimes, I prefer long interactive sessions.*

*Yes, runs faster than my local laptop.*

*Yes, local servers are running out of storage.*

Getting Scientific Software Installed 2.0?

<https://www.vscentrum.be/>

see <https://github.com/vscentrum/gssi-training>

(recording available soon)

# MrBayes: An overview



**Bayes' theorem: what is the probability that a model parameter is true, given data and prior beliefs?**

**MrBayes** is a software program used for **Bayesian phylogenetic inference**.

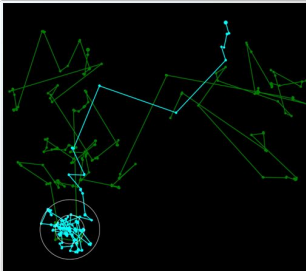
**In simple terms:**

- It builds evolutionary trees using Bayesian statistics and Markov Chain Monte Carlo (MCMC) sampling to estimate the most likely tree and the support for each branch.
- $\Theta$  includes parameters such as topologies, branch lengths;  $X$  is your sequence alignment

$$\text{posterior distribution } \leftarrow f(\theta | X) = \frac{f(\theta) f(X | \theta)}{\int f(\theta) f(X | \theta) d\theta}$$

prior distribution ← likelihood →  
normalizing constant ←

$\theta$  = model parameter  
 $X$  = data



- 2-dimensional image of posterior probability space
- 2 chains
- poor mixing
- no convergence

- Compute bound, so shinx (recent) should be faster than doduo (old)...
- Can use multiple cores (via Beagle dependency)
  - Via **multi-threading** (native threading, not OpenMP)
  - Via **multi-procesing** (MPI)
  - How well does it scale with # cores?
  - Is it worth using multiple *nodes*?
- Can also use **GPU**, that sounds promising...
- Only a small input file (~6MB), so no problems w.r.t. I/O
- Settings in \*.nex file:

```
mcmc nruns=2 nchains=4 temp=0.1 ngen=5000000  
samplefreq=1000 printfreq=1000 diagnfreq=5000;
```

# Running MrBayes jobs @ HPC-UGent

- Naive starting point: **single-core job**
- **Takes ~65h on shinx**

```
#!/bin/bash  
#PBS -l nodes=1:ppn=1  
#PBS -l walltime=72:00:00
```

```
source /cvmfs/software.eessi.io/versions/2023.06/init/bash  
module load MrBayes/3.2.7-gompi-2023a
```

```
cd $PBS_O_WORKDIR
```

```
mb mbtest.nex
```

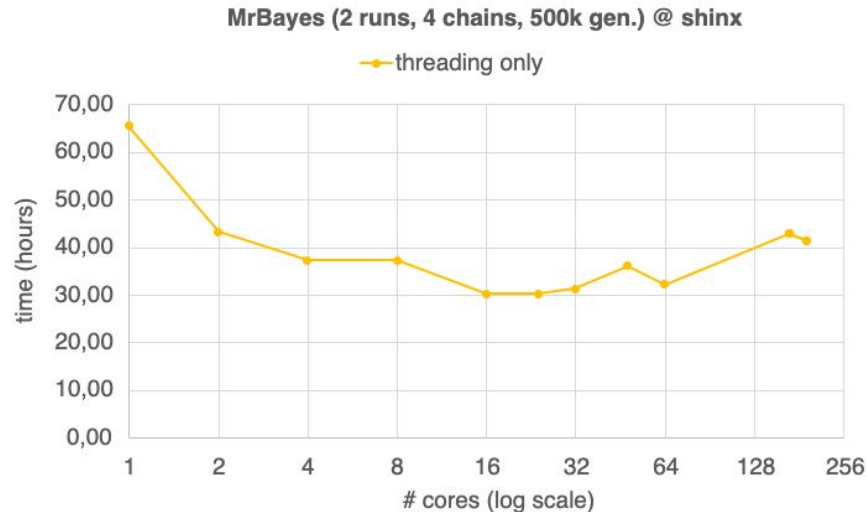
# Running MrBayes jobs @ HPC-UGent

- **Multiple cores via (only) threading** => limit: # cores in a single node
- In job script: only change `ppn=1` to desired number of cores/threads (example: `ppn=16`)
- In `*.nex` file, also indicate # threads to use via `beaglethreadcount` setting (to avoid that too many threads are being used based on auto-detection of # cores)

```
begin mrbayes;  
  set usebeagle=yes;  
  set beagledevice=cpu;  
  set beagleprecision=single;  
  set beaglescaling=dynamic;  
  set beaglesse=no;  
  set beaglethreadcount=16;
```

# Running MrBayes jobs @ HPC-UGent

- **Multiple cores via (only) threading** => limit: # cores in a single node
- **Some speedup**, but only 2.17x with 16 cores vs single core...
  - Only ~14% efficiency with 16 cores
  - Flatlines (or even gets worse) when trying to use more cores/threads



# Running MrBayes jobs @ HPC-UGent



- **Multiple cores via MPI** => limit: # runs x # chains == # cores (one MPI rank per core)
- In \*.nex file: use beaglethreadcount=1
- In job script: **use mympirun** tool, **change ppn=N** with N = 1, 2, 4, 8

```
#!/bin/bash
```

```
#PBS -l nodes=1:ppn=8
```

```
#PBS -l walltime=72:00:00
```

```
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
```

```
module load MrBayes/3.2.7-gompi-2023a
```

```
module load vsc-mypirun
```

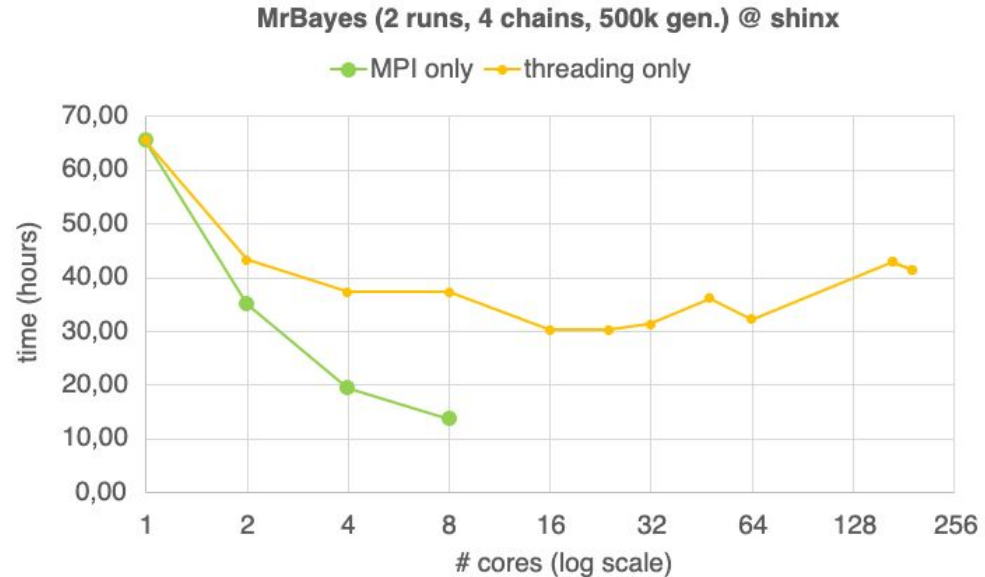
```
cd $PBS_O_WORKDIR
```

```
mympirun mb mbtest.nex
```

# Running MrBayes jobs @ HPC-UGent



- **Multiple cores via MPI** => limit: # runs x # chains == # cores (one MPI rank per core)
- **Pretty good scaling!**
  - With 8 cores (8 MPI ranks, each with 1 thread): **4.8x speedup** (~60% efficiency)
  - **Still takes ~13.5h...**



# Running MrBayes jobs @ HPC-UGent

- **Multiple cores via MPI + threading** (a.k.a. “hybrid”)
  - Job script:
    - Use **ppn=C** with C number of cores (per node) to use
    - Use **--hybrid N** to run N MPI ranks per node (instead of 1 per core)
  - In \*.nex file: use `beaglethreadcount=T` with T=2,4,8,...

```
#!/bin/bash
```

```
#PBS -l nodes=1:ppn=32
```

```
#PBS -l walltime=72:00:00
```

```
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
```

```
module load MrBayes/3.2.7-gompi-2023a
```

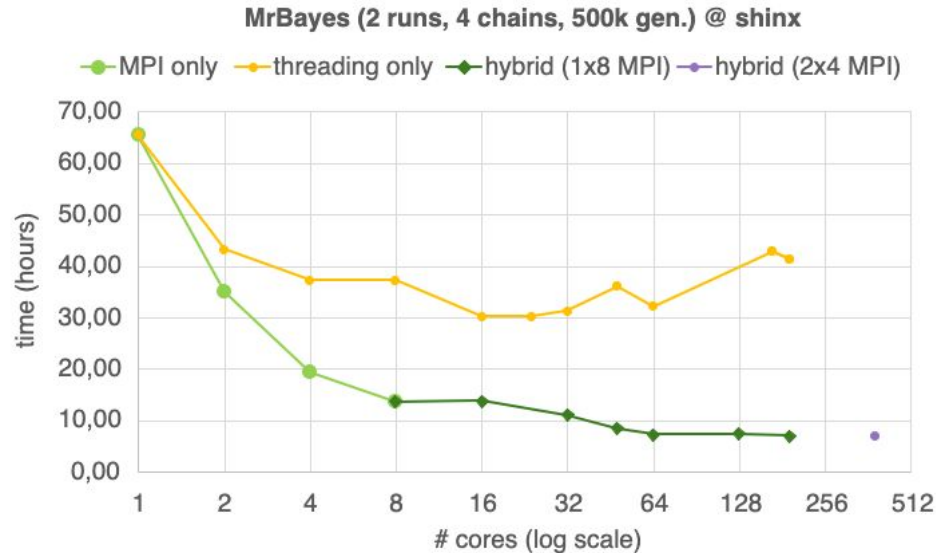
```
module load vsc-mypirun
```

```
cd $PBS_O_WORKDIR
```

```
mympirun --hybrid 8 mb mbtest.nex
```

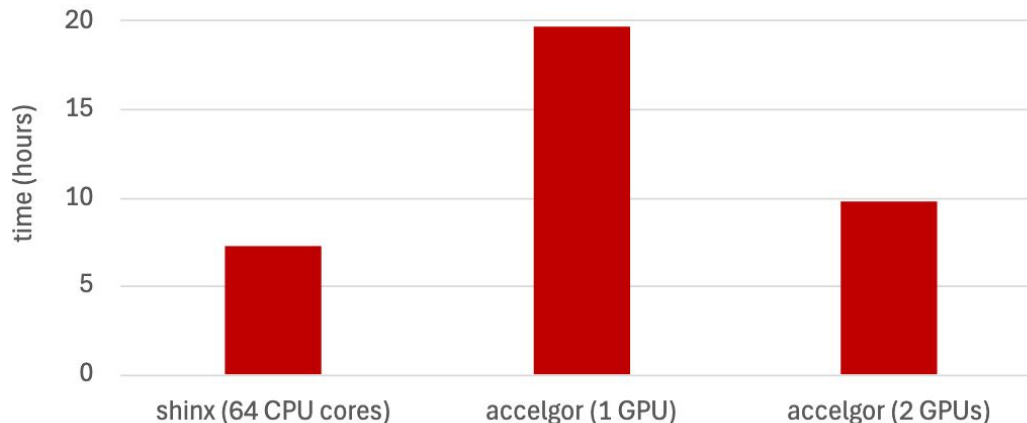
# Running MrBayes jobs @ HPC-UGent

- **Multiple cores via MPI + threading (a.k.a. “hybrid”)**
  - **More speedup beyond 8 cores!**
  - Almost 2x speedup with 8 MPI ranks with each 8 threads (64 cores): ~7h
  - No point to use multiple nodes: 2 nodes, 4x96 threads => no more speedup



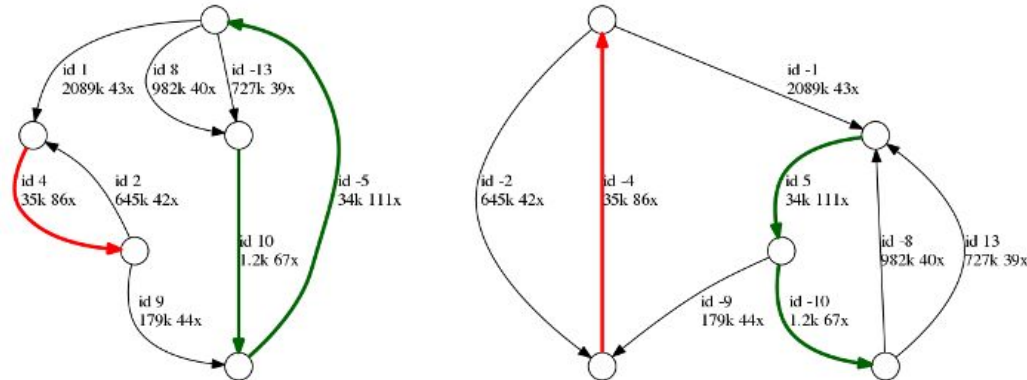
# Running MrBayes jobs @ HPC-UGent

- Can we get more speedup using GPUs? **No.**
  - 1 GPU in accelgor: ~20h, 2 GPUs: ~10h
  - Good scaling, but **makes no sense** compared to just 64 cores in shinx...
    - Takes **longer** on 2 GPUs vs 64 CPU cores
    - Queue wait times for GPUs are quite long



# Flye: An overview <https://github.com/mikolmogorov/Flye>

- de novo assembler
- PacBio & ONT
- Small to larger assemblies
- Complete pipeline (but assemblies are collapsed => HapDup alternative)
- Meta Flye => meta genome modus available



## 1. Assembly (disjointigs + k-mer counting + extending)

Reads analysis, detect repeats & computes frequency/coverage, initial graphs

## 2. Consensus computation

Remap reads to initial disjointigs - compute contig

## 3. Repeat graph

Builds graph - simplify paths

## 4. Contigger

Convert repeat graphs in contigs (scaffolding if possible)

## 5. Polishing

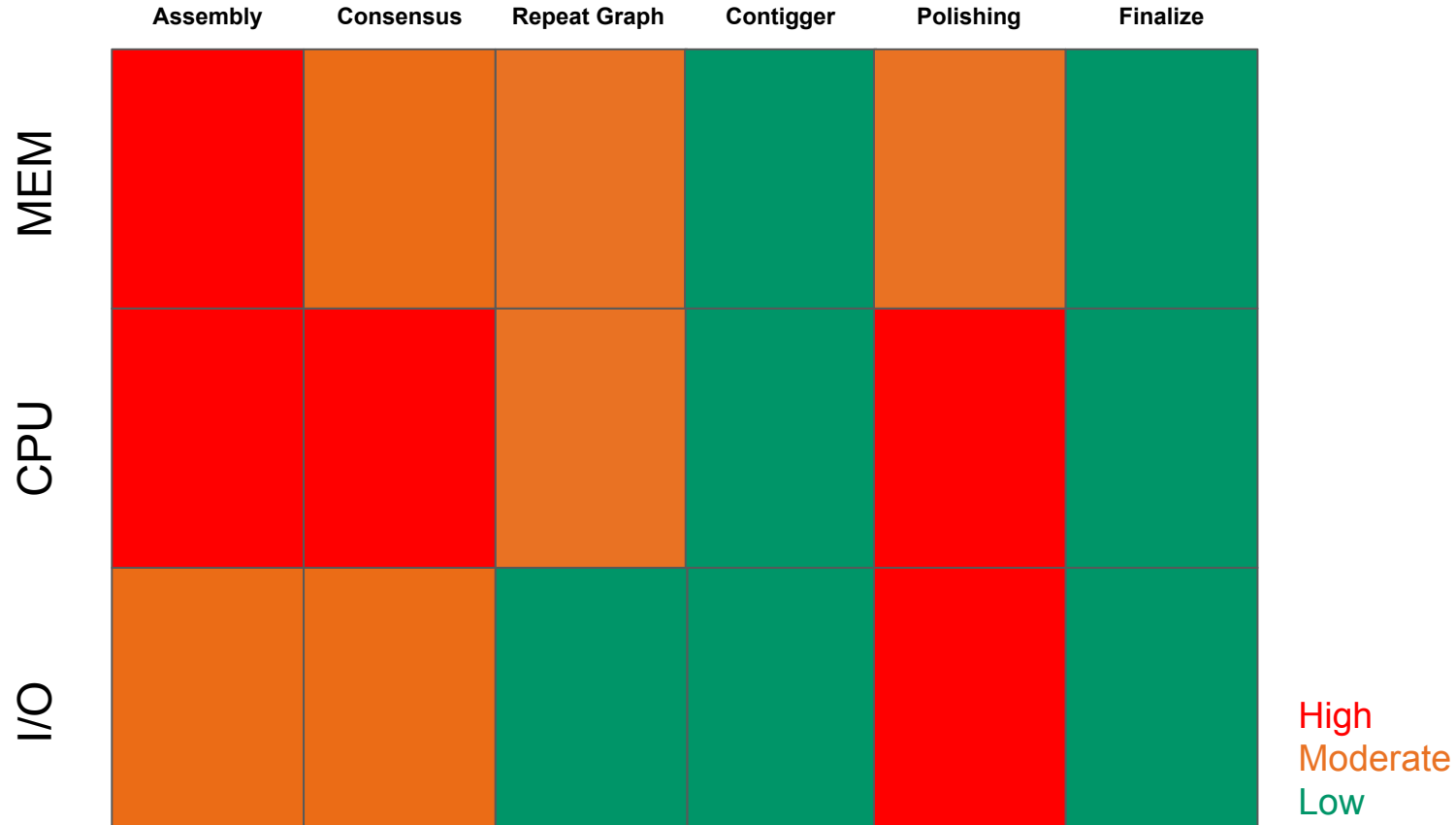
Align read back to contigs (resolve indels and mismatches)

## 6. Finalize

### flye.log

```
[ --dec={ 0 | 1 | 2 } ] [ --human ]  
[2025-11-12 11:33:37] INFO: Starting Flye 2.9.4-b1799  
[2025-11-12 11:33:37] INFO: >>>STAGE: configure  
[2025-11-12 11:33:37] INFO: Configuring run  
[2025-11-12 11:33:39] INFO: Total read length: 207507950  
[2025-11-12 11:33:39] INFO: Reads N50/N90: 14145 / 4332  
[2025-11-12 11:33:39] INFO: Minimum overlap set to 4000  
[2025-11-12 11:33:39] INFO: >>>STAGE: assembly  
[2025-11-12 11:33:39] INFO: Assembling disjointigs  
[2025-11-12 11:33:39] INFO: Reading sequences  
[2025-11-12 11:33:45] INFO: Counting k-mers:  
0% 10% 20% 30% 40% 50% 60% 70% 80% 90% 100%  
[2025-11-12 11:35:19] INFO: Filling index table (1/2)  
0% 10% 20% 30% 40% 50% 60% 70% 80% 90% 100%
```

# Flye: Computational overview



# Flye: Hardware & Software



	<b>CPU</b>	<b>RAM</b>	<b>Arch</b>	<b>Model</b>	<b>Inter connection</b>	<b>Platform</b>
<b>LOCAL</b>	4 cores - 8 threads	32Gib	x86_64	Intel(R) Core(TM) i7-8665U CPU @ 1.90GHz		Conda/EESSI/Container
<b>GROVER Tier-1 cloud</b>	20 cores (single threaded)	360Gib	x86_64	Intel Core Processor (Broadwell) 2.2GHz		Conda/EESSI/Container
<b>HPC (gallade)</b>	2x64 cores	940Gib	x86_64	AMD EPYC 7773X (Milan-X @ 2.2 GHz)	HDR-100 InfiniBand	Conda/EESSI(+generic)/Container

# Flye: Input data

Genome	Data	Asm.Size	NG50	CPU time	RAM
<a href="#">E.coli</a>	PB 50x	4.6 Mb	4.6 Mb	2 h	2 Gb
<a href="#">Zymo_Log</a>	ONT meta 16 Gb	29 Mb	N/A	100 h	76 Gb

**E.Coli:** This dataset includes **one SMRT Cell** of data gathered with a PacBio RS II System and P6-C4 chemistry on a size selected **20kb library of E. coli K12 substrain MG1655**. Assembling this SMRT Cell results in a **single contig at >Q50 accuracy using SMRT Analysis 2.0 or higher**.

**Zymo.Log:** 10 species (5 Gram-positive, 3 Gram-negative, 2 yeast) ranging from  $10^2$  -  $10^8$  genomic DNA abundance (total input  $5 \times 10^8$  cells)

Data available from:

- GridION (Zymo-GridION-LOG-BB-SN)
- PromethION (Zymo-PromethION-LOG-BB-SN)

## 1. Conda

- a. Create env (naive)
- b. Run flye

## 2. Container

- a. Build container
- b. Run flye

## 3. EESSI

- a. Activate eessi
- b. Run flye



## CONDA

- Package manager
- Multi -platform
- Easy access to programming tools
- Online repositories

### LINKS

<https://www.anaconda.com/docs/main>

<https://github.com/conda-forge/miniforge>

(aarch64 support and more lightweight!)

## APPTAINER

- Container platform
- Available on HPC systems (unlike Docker)
- Many open source tools
- OS incl.
- Online repositories

### LINKS

<https://github.com/BioContainers/containers>

<https://hub.docker.com/>

<https://biocontainers.pro/>

# Flye: Local-conda

```
# Create and activate a conda env
conda create -n flye_env python=3.12 sysstat
conda activate flye_env
conda install -c bioconda flye #will install latest version on bioconda channel

# Stage DATA
INPUTDATA="./E.coli_PacBio_40x.fasta"
PIDLOG="pidstat_flye"
OUTPUT="out_pacbio-local-conda"

# Start pidstat in the background for resource monitoring
pidstat -r -u -d -h -p ALL 5 > "${PIDLOG}.log" &
    #-r Report page faults and memory utilization
    #-u Report CPU utilization
    #-d Report I/O statistics
    #-h Display all activities horizontally on a single line, make parsing easier
    #-p Monitor all ongoing processes
    #5 Monitor every 5 seconds

# PIDSTAT_PID variable assign to later cleankill the reoccurring pidstat command
PIDSTAT_PID=!

# run flye
flye --pacbio-raw $INPUTDATA --out-dir "$OUTPUT" --threads 4

# Kill pidstat after Flye finishes
kill $PIDSTAT_PID

# cleanup PIDSTAT logfile
awk 'NR==1 || /(^# Time/ || tolower($NF) ~ /^(flye|flye-minimap2|flye-modules|flye-samtools|flye294|python3(\.10)?|starter-suid|bash|sh)$/'
"${PIDLOG}.log" > "${PIDLOG}.filtered.log"
```

```
#INSTALL CUSTOM VERSION

conda search -c bioconda flye
.
.
.
conda install -c bioconda flye=2.9.4
```

# Flye: Local-container-definition file

flye-ubuntu22.04.def



```
Bootstrap: docker
From: ubuntu:22.04
%post
# Update and install dependencies
apt-get update && apt-get install -y \
    build-essential python3 python3-pip python3-dev git zlib1g-dev libbz2-dev liblzma-dev libncurses5-dev \
    Libncursesw5-dev libcurl4-openssl-dev libssl-dev wget cmake make

apt install -y squashfuse fuse-overlayfs fuse2fs gocryptfs
# Clone Flye repository
cd /opt
git clone https://github.com/mikolmogorov/Flye.git
cd Flye
# Build Flye
make
# Optional: install Flye system-wide
ln -s /opt/Flye/bin/flye /usr/local/bin/flye

%environment
export PATH=/opt/Flye/bin:$PATH

%labels
Author Pieter Asselman
Version v1.0

%runscript
echo "Flye is ready to use. Try: flye --help"
```

**apt-get install** : traditional Debian/ubuntu  
**apt install**: newer frontend, nice progress bars

Will install  
latest  
version!



# Flye: Local-container-definition file

flye294-ubuntu22.04.def



```
Bootstrap: docker
From: ubuntu:22.04
```

```
%post
```

```
# Update and install dependencies
```

```
apt-get update && apt-get install -y \
```

```
Build-essential python3 python3-pip python3-dev git zlib1g-dev libbz2-dev liblzma-dev
```

```
Libcursesw5-dev libcurl4-openssl-dev libssl-dev wget cmake pip sysstat make
```

```
apt install -y squashfuse fuse-overlayfs fuse2fs gocryptfs
```

```
pip install "git+https://github.com/mikolmogorov/Flye@2.9.4"
```


```
%labels
```

```
Author Pieter Asselman
```

```
Version v1.0
```

```
%runscript
```

```
echo "Flye is ready to use. Try: flye --help"
```



Installs  
version  
2.9.4!



APPTAINER

# Flye: Local-container - Build&Run

```
# BUILD A CONTAINER FROM DEFINITION FILE
```

```
apptainer build containername.sif containername.def
```

```
# RUN flye from A CONTAINER FROM SIF FILE
```

```
apptainer exec ~/containers/flye-ubuntu22.04.sif \  
flye --pacbio-raw E.coli_PacBio_40x.fasta --out-dir flye-container --threads 4
```

```
# RRUN flye from within the container (or add applications in the container)
```

```
apptainer shell --writable flye-ubuntu22.04.sif  
#--writable by default container is read only, make file system available as read/write
```

**DON'T FORGET !**

- 1- PUT .SIF FILE IN \$VSC\_SCRATCH
- 2- PREP ENV for container runs

[https://hpcugent.github.io/vsc\\_user\\_docs/HPC/GentLinux/apptainer/?h=apptainer](https://hpcugent.github.io/vsc_user_docs/HPC/GentLinux/apptainer/?h=apptainer)  
D&E

# Flye: Local-EESSI



```
#!/bin/bash

# ACTIVATE EESSI
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
ml avail Flye
ml Flye/2.9.4-GCC-13.2.0

# Optional: Start pidstat in the background for resource monitoring
pidstat -r \ #-r Report page faults and memory utilization
-u \ #-u Report CPU utilization
-d \ #-d Report I/O statistics
-h \ #-h Display all activities horizontally on a single line, make parsing easier
-p ALL \ #-p Monitor all ongoing processes
5 \ #5 Monitor every 5 seconds
> pidstat_flye-grover-container.log &

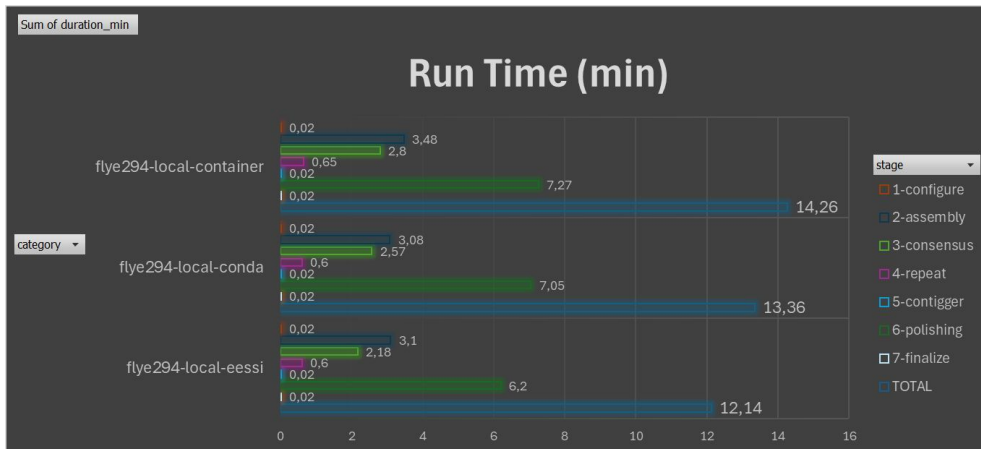
PIDSTAT_PID=$! # PIDSTAT_PID variable assign to later cleankill the reoccurring pidstat command

#Run flye command
flye --pacbio-raw E.coli_PacBio_40x.fasta --out-dir out_pacbio-local-eessi --threads 4

#Optional: Kill pidstat after Flye finishes
kill $PIDSTAT_PID
```



# Flye: Results - LOCAL



## Reminder

Container: pip install

Conda: conda install

EESSI: compiled based on (micro)arch detected

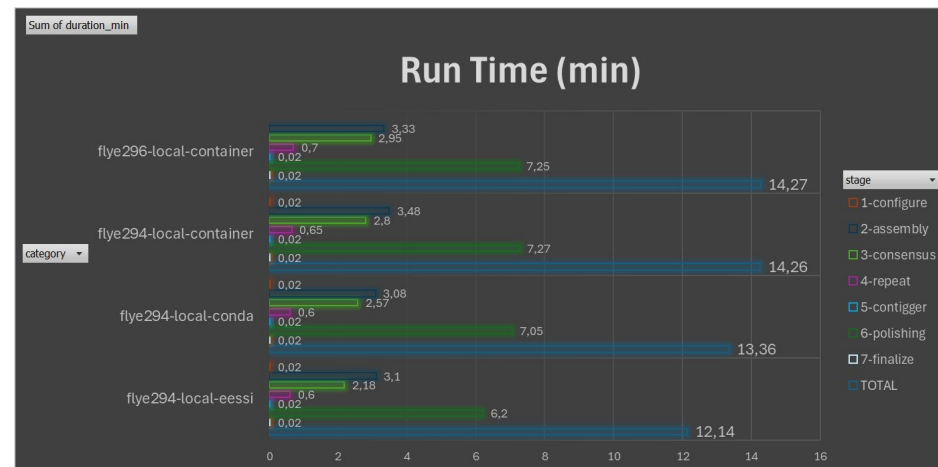
## FYI

conda/container => architecture compatibility **yes**

=> microarchitecture optimization **no**

## Key insights

- EESSI runs remarkably faster
- Flye 2.9.6 consensus building faster, but polishing longer



# Flye: Results - GROVER



## Key insights

- Same story here
- **BUT** Local EESSI(12.14m) runs faster than Grover EESSI (12.61m)???



# Flye: Results - GROVER

	CPU	RAM	Arch (CPU family)	Model (micro architecture)
<b>LOCAL (2019)</b>	4 cores - 8 threads	32Gib	x86_64	Intel(R) Core(TM) i7-8665U CPU @ <b>1.90GHz</b>
<b>GROVER Tier-1 cloud (2014)</b>	20 cores (single threaded)	360Gib	x86_64	Intel Core Processor (Broadwell) <b>2.2GHz</b>

## OLD vs NEW cores

- **IPC Advantage (Instructions per cycle):** Whiskey Lake (2019) has much higher instructions-per-cycle than Broadwell (2014).
- **Turbo Boost:** i7-8665U can hit ~4.8 GHz on 1–2 cores, while Broadwell likely tops out near 3 GHz.
- **Instruction Set:** Whiskey Lake supports newer optimizations (AVX2, better branch prediction).





# Flye: HPC - Can we Speed things up?

shinx : 48 nodes - 2\*96 cores - 370 GiB/node

=> Let's be modest, request half a node



BAD OR GOOD USAGE OF RESOURCES?

```
Name           : flye-conda
User            : vscxxxx
Partition      : shinx
Nodes          : node4232.shinx.os
Cores         : 96
State          : COMPLETED
Submit         : 2025-11-07T11:47:07
Start          : 2025-11-07T22:29:38
End            : 2025-11-07T23:12:33
Reserved walltime : 10:00:0.0
Used walltime  : 00:42:55.0
Used CPU time  : 14:29:18.0
% User (Computation): 99.36
% System (I/O)   : 0.64
Mem reserved   : 185G
Max Mem used    : 62.90G (node4232.shinx.os)
Max Disk Write  : 92.70G (node4232.shinx.os)
Max Disk Read   : 204.52G (node4232.shinx.os)
Working directory : /user/gent/433/vsc4xxxx
```

## Let's do the math

Walltime:  $42\text{min}/60=0.7$

CPU-time:  $14+(29/60)= 14.29$

#Cores Used= CPU-time/Walltime

=>  $14,29/0.7 \sim 20$  cores

**Requested too many  
cores!**

# Flye: HPC - Can we Speed things up?

shinx : 48 nodes - 2\*96 cores - 370 GiB/node

=> be wise, adapt, request 8 cores and 80GiB RAM

BAD OR GOOD USAGE OF RESOURCES?

```
Name       : flye-conda
User        : vsc4xxxx
Partition   : shinx
Nodes       : node4103.shinx.os
Cores      : 8
State       : COMPLETED
Submit      : 2025-11-14T15:06:03
Start       : 2025-11-14T15:06:03
End         : 2025-11-14T16:50:18
Reserved walltime : 20:00:0.0
Used walltime   : 01:44:15.0
Used CPU time : 11:32:40.0
% User (Computation): 99.26
% System (I/O) : 0.74
Mem reserved   : 80G
Max Mem used    : 62.36G (node4103.shinx.os)
Max Disk Write : 103.72G (node4103.shinx.os)
Max Disk Read  : 209.50G (node4103.shinx.os)
Working directory : /user/gent/433/vsc4xxxx
```

## Let's do the math

Walltime:  $1 + (44/60) = 1.7h$

CPU time:  $11 + (32/60) = 11.5$

CPU-time/Walltime =>  $11.5/1.7 \sim 6.76$  cores

**Requested a good amount  
of cores!**

**But walltime is double...**



# Flye: HPC - Can we Speed things up?



## Re(search)quest Resources

- Run on 8 cores
- 16
- 32
- 48
- ...
- Add MEM
- Look for I/O
- $\#cores = CPUtime / Walltime$



# Flye: HPC - Can we Speed things up?

Genome	Data	Asm.Size	NG50	CPU time	RAM
<a href="#">E.coli</a>	PB 50x	4.6 Mb	4.6 Mb	2 h	2 Gb
<a href="#">Zymo_Log</a>	ONT meta 16 Gb	29 Mb	N/A	100 h	76 Gb

GitHub: <https://github.com/mikolmogorov/Flye>

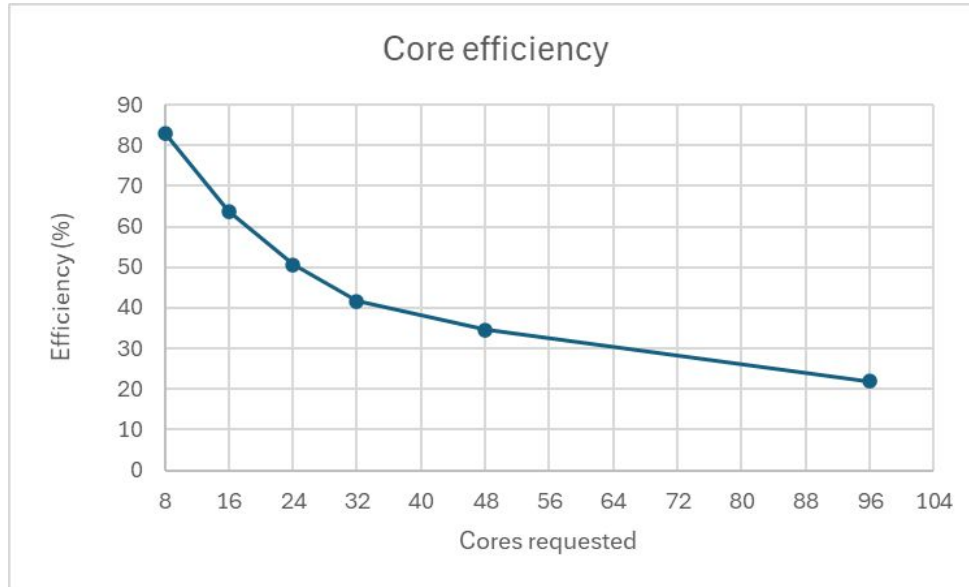
**E.Coli:** This dataset includes **one SMRT Cell** of data gathered with a PacBio RS II System and P6-C4 chemistry on a size selected **20kb library of E. coli K12 substrain MG1655**. Assembling this SMRT Cell results in a **single contig at >Q50 accuracy using SMRT Analysis 2.0 or higher**.

**Zymo.Log:** 10 species (5 Gram-positive, 3 Gram-negative, 2 yeast) ranging from  $10^2$  -  $10^8$  genomic DNA abundance (total input  $5 \times 10^8$  cells)

Data available from:

- GridION (Zymo-GridION-LOG-BB-SN)
- PromethION (Zymo-PromethION-LOG-BB-SN)

# Flye: HPC - Scaling # cores



Trade off...

A few minutes runtime gain

...

A whole lot more queue time

Calculations:

Used cores      => CPU-time / Walltime

Core efficiency   => (Used cores / Requested cores) \*100

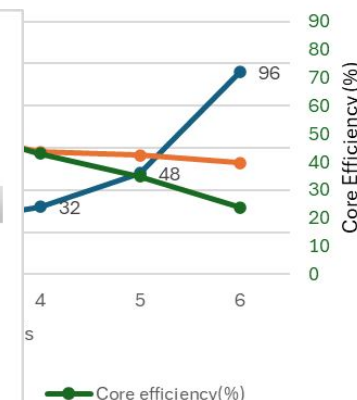
# Flye: HPC - Scaling - Memory



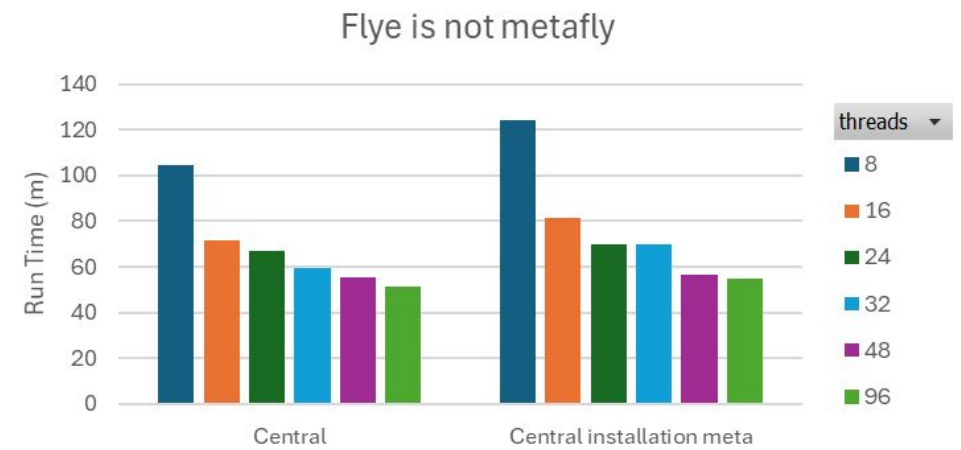
80 GiB - Central installation



160 GiB - Central installation



Sum of minutes



threads

- 8
- 16
- 24
- 32
- 48
- 96

Group

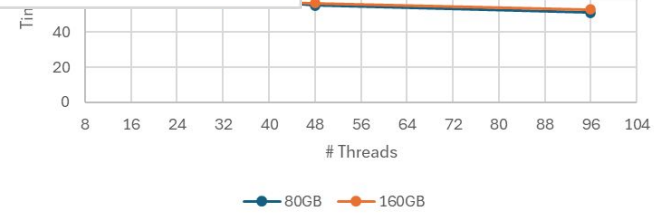
## Key insights

- Best balance
- >32 threads small gains
- 96 threads fastest runtime but 22% core efficiency
- Memory used tops off at ~62GB (assembly size ~24Mb)

16: Runti  
24: 66 m

mem please?

We have enough.



# Flye: HPC - Scaling - EESSI

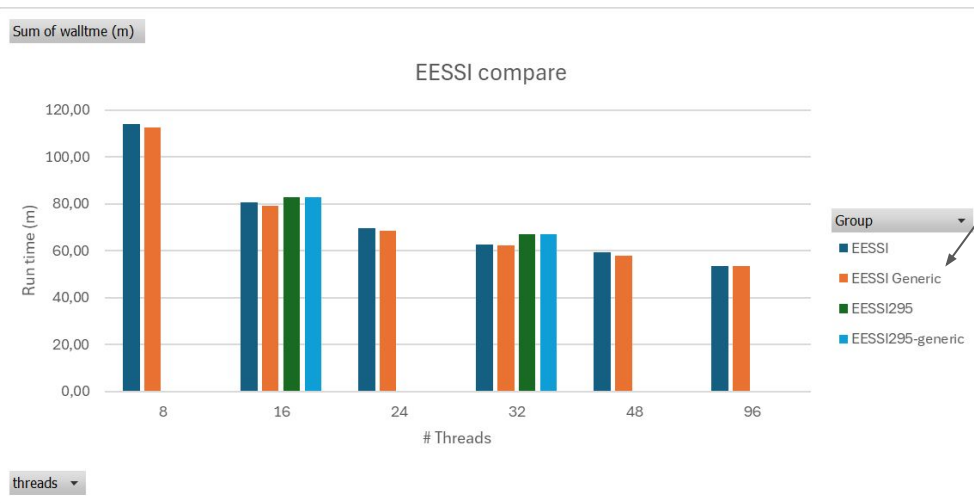


## Optimization level

EESSI (full option): arch + micro arch  
Generic - family level only

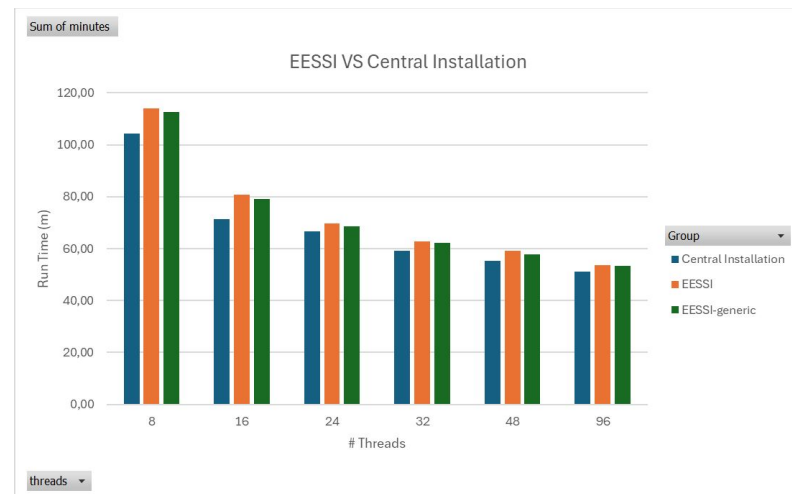
entries

x86_64		amd					intel		
generic	zen2	zen3	zen4	cascadelake	haswell	icelake			
x	x	x	x	x	x	x	x	x	x



## Key insights

- EESSI optimal VS generic => same result
- Central installation runs faster than EESSI on HPC => Different causes ... Kenneth Help!

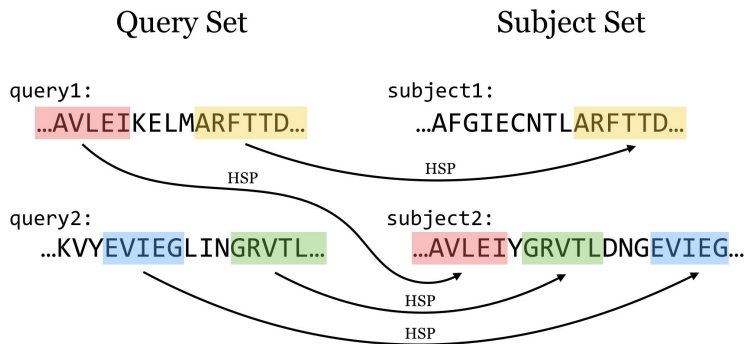


## Take home

- How did you get your software on your computer?
- Don't like queueing too much? Optimize your resource usage first.
- Do you like interactive sessions?  
=> Go play with donphan (interactive/DEBUG)!
- Central software installations are heavily optimized to run on the clusters.  
=> Yes, it takes time to get them installed centrally...
- Use EESSI locally, it runs faster!
- If you still see limitations, ask HPC-team if there is a solution.  
=> We can't all be computer scientists.

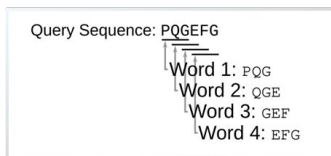
# BLAST+: An overview

[https://github.com/ncbi/blast\\_plus\\_docs](https://github.com/ncbi/blast_plus_docs)



## BLAST Workflow (condensed)

- Clean the query (remove low-complexity/repeats).
- Break the query into k-letter words.
- Score these words and keep only high-scoring ones.
- Build a fast lookup index for the selected words.
- Search the database for exact word matches.
- Extend matches into **high-scoring segment pairs (HSPs)**.
- Score HSPs and evaluate significance (E-values).
- Merge HSPs into longer alignments.
- Report significant hits.



## What does it do?

Basic Local Alignment Search Tool ([BLAST](#))

finds regions of local similarity between sequences.

Program	Query Type	Subject Type	Computation
blastn	N	N	~ 1X
blastp	P	P	~ 1X
blastx	N	P	~ 6X
tblastn	P	N	~ 6X
tblastx	N	N	~36X

(other BLAST types not listed: psiblast, deltablast, rpsblast)

## Query

## Subject

Blastn - Nucleotide	=>	Nucleotide
Blastp - Protein	=>	Protein
Blastx - Nucleotide (translated)	=>	protein
tBlastn - Protein	=>	Nucleotide (translated)
tBlastx - Nuc (transl)	=>	Prot (transl)

# BLAST+: Workload info

- `blastn` (nucleotide - nucleotide)
- Supports using multiple cores via multi-threading (OpenMP)
- Large input file (`*.fasta`): ~115,000 sequences
  - Easy to reduce to smaller input for testing/benchmarking: first N lines
- **Very** I/O-intensive: search in large databases
  - `core_cnt` NCBI database: ~260GB
  - `nt` NCBI database: ~700GB

# Running BLAST+ jobs @ HPC-UGent

Starting point:

- **Single-core, database on `$VSC_DATA` filesystem**
- **Reduced input file** (5,000 sequences), ~4% of full input of 115k sequences
- **Smaller database** (`core_nt`), ~260GB

```
#!/bin/bash
#PBS -l walltime=72:00:00
#PBS -l nodes=1:ppn=1

source /cvmfs/software.eessi.io/versions/2023.06/init/bash
module load BLAST+/2.14.1-gompi-2023a
export BLASTDB="$VSC_DATA_VO/blast/core_nt"

cd $PBS_O_WORKDIR
blastn -task megablast -query small-5k-seq.fasta -db $BLASTDB/core_nt
-outfmt '6 qseqid staxids bitscore std sscinames sskingdoms stitle' -evaluate
1e-10 -max_hsp 1 -max_target_seqs 5
```

# Running BLAST+ jobs @ HPC-UGent

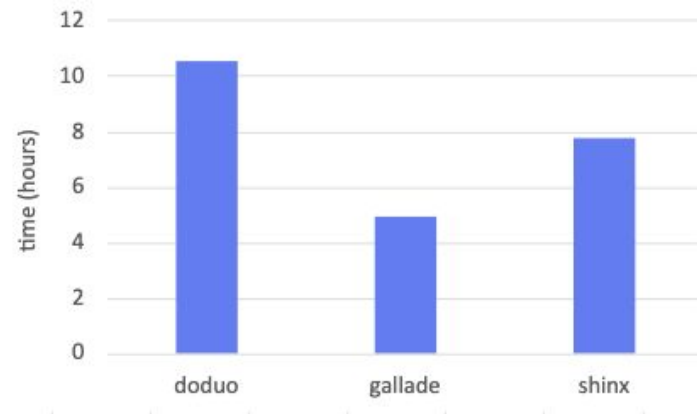
Starting point:

- **Single-core, database on `$VSC_DATA` filesystem**
- Reduced input file (5,000 sequences), ~4% of full input of 115k sequences
- Smaller database (core\_nt), ~260GB

**Results: 5-11 hours of required walltime**

Smells like trouble for full input file + nt database, given 72h walltime limit...

Naive extrapolation:  $x20 \times x3 \Rightarrow$  over 300h of walltime!



# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **multiple cores** via threading (OpenMP)?

- `core_cnt @ $VSC_DATA` + reduced input file (5k sequences)
- **First attempt: close to 2x speedup with 2 threads, then no further speedup?!**

```
#!/bin/bash
#PBS -l walltime=72:00:00
#PBS -l nodes=1:ppn=8

source /cvmfs/software.eessi.io/versions/2023.06/init/bash
module load BLAST+/2.14.1-gompi-2023a
export BLASTDB="$VSC_DATA_VO/blast/core_nt"

cd $PBS_O_WORKDIR
blastn -task megablast -query small-5k-seq.fasta -db $BLASTDB/core_nt
-num_threads $PBS_NUM_PPN -outfmt '6 qseqid staxids bitscore std sscinames
sstringdomains stitle' -evaluate 1e-10 -max_hsps 1 -max_target_seqs 5
```

# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **multiple cores** via threading (OpenMP)?

- `core_cnt @ $VSC_DATA` + reduced input file (5k sequences)
- **Live job inspection: blastn process with all threads stuck on a single core!**

```
$ qstat -ant 12345          # determine workernode on which job is running
$ ssh node1234            # SSH into that workernode
$ ps xfu                  # show list of running processes - look for blastn
USER          PID %CPU %MEM    VSZ   RSS TTY  STAT  START   TIME COMMAND
...
vsc40023    345678 99.9  4.2   123456 56789 ?    R    12:34   34:56  \_ blastn ...

$ taskset -c p 345678      # check on which cores the blastn process is running
pid 345678 current affinity list: 0
```

**Max. 100% CPU  
=> only using a single core**

**This should be a list of core IDs  
(not one single one...)**

# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **multiple cores** via threading (OpenMP)?

- `core_cnt @ $VSC_DATA` + reduced input file (5k sequences)
- Fixed job script (**no more pinning to a single core**)

```
#!/bin/bash
#PBS -l walltime=72:00:00
#PBS -l nodes=1:ppn=8
```

→ `unset OMP_PROC_BIND` # avoid being stuck to a single core...

```
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
module load BLAST+/2.14.1-gompi-2023a
export BLASTDB="$VSC_DATA_VO/blast/core_nt"
```

```
cd $PBS_O_WORKDIR
blastn -task megablast -query small-5k-seq.fasta -db -db $BLASTDB/core_nt
-num_threads $PBS_NUM_PPN -outfmt '6 qseqid staxids bitscore std sscinames
sskingdoms stitle' -evaluate 1e-10 -max_hsps 1 -max_target_seqs 5
```

# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **multiple cores** via threading (OpenMP)?

- `core_cnt @ $VSC_DATA` + reduced input file (5k sequences)
- Fixed job script (no more pinning to a single core)

```
$ qstat -ant 12346          # determine workernode on which job is running
$ ssh node1234            # SSH into that workernode
$ ps xfu                  # show list of running processes - look for blastn
USER          PID %CPU %MEM    VSZ   RSS TTY  STAT  START    TIME COMMAND
...
vsc40023    456789 798  32.5  123456 56789 ?    R    12:34   734:56  \_ blastn ...

$ taskset -c p 456789      # check on which cores the blastn process is running
pid 345678 current affinity list: 0-7
```

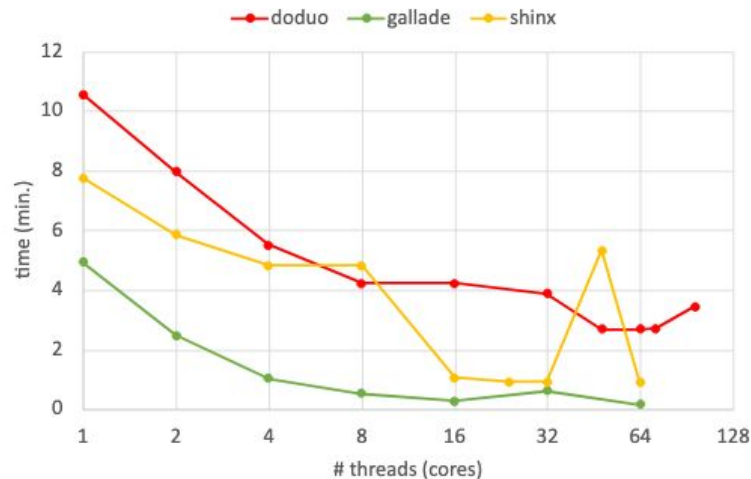
Close to 800% CPU  
=> using 8 cores 👍

List of core IDs 👍

# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **multiple cores** via threading (OpenMP)?

- `core_cnt @ $VSC_DATA` + reduced input file (5k sequences)
- Fixed job script (no more pinning to a single core)
- **Significant speedup** with multiple cores/threads (but **flatlines** at 8-16 cores)
- gallade cluster is *remarkably* better, well under 1h when using 8+ cores!
- Some quirky results...  
Using shared resources can lead to chaos



# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **different shared filesystem**?

- `core_nt @ $VSC_DATA , $VSC_SCRATCH, $VSC_SCRATCH_ARCANINE`
- Reduced input file (5k sequences)

```
#!/bin/bash
#PBS -l walltime=72:00:00
#PBS -l nodes=1:ppn=32

unset OMP_PROC_BIND    # avoid being stuck to a single core...

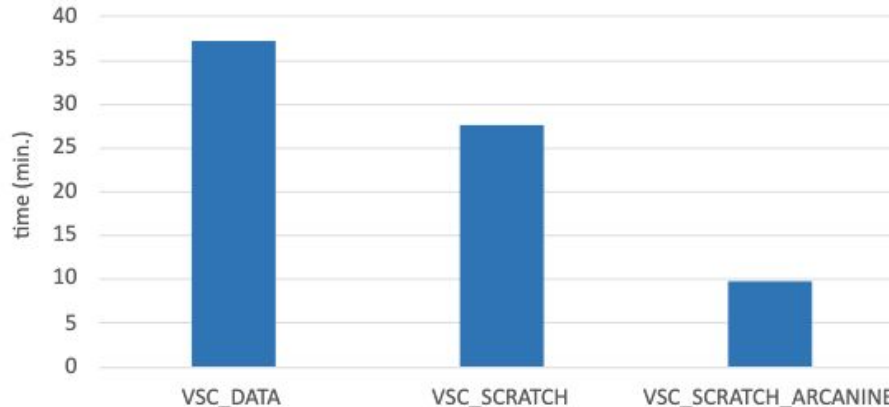
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
module load BLAST+/2.14.1-gompi-2023a
export BLASTDB="$VSC_SCRATCH_ARCANINE_VO/blast/core_nt"

cd $PBS_O_WORKDIR
blastn -task megablast -query small-5k-seq.fasta -db -db $BLASTDB/core_nt
-num_threads $PBS_NUM_PPN -outfmt '6 qseqid staxids bitscore std sscinames
sskingdoms stitle' -evaluate 1e-10 -max_hsps 1 -max_target_seqs 5
```

# Running BLAST+ jobs @ HPC-UGent

Can we make it faster using **different shared filesystem**?

- `core_nt @ $VSC_DATA , $VSC_SCRATCH, $VSC_SCRATCH_ARCANINE`
- Reduced input file (5k sequences)
- gallade cluster, 32 cores, required time in *minutes*
- ~4x faster with `$VSC_SCRATCH_ARCANINE` vs `$VSC_DATA`
- Can vary a lot due to load on shared filesystem



# Running BLAST+ jobs @ HPC-UGent



Can we run **full input** (115k seq.) using **large nt database** (~700GB) in < 72h?

# Running BLAST+ jobs @ HPC-UGent



Can we run **full input** (115k seq.) using **large nt database** (~700GB) in < 72h?

Database @ \$VSC\_SCRATCH\_ARCANINE, **multi-threaded** (quarter of cores per node)

```
#!/bin/bash
#PBS -l walltime=72:00:00
#PBS -l nodes=1:ppn=quarter

unset OMP_PROC_BIND # avoid being stuck to a single core...

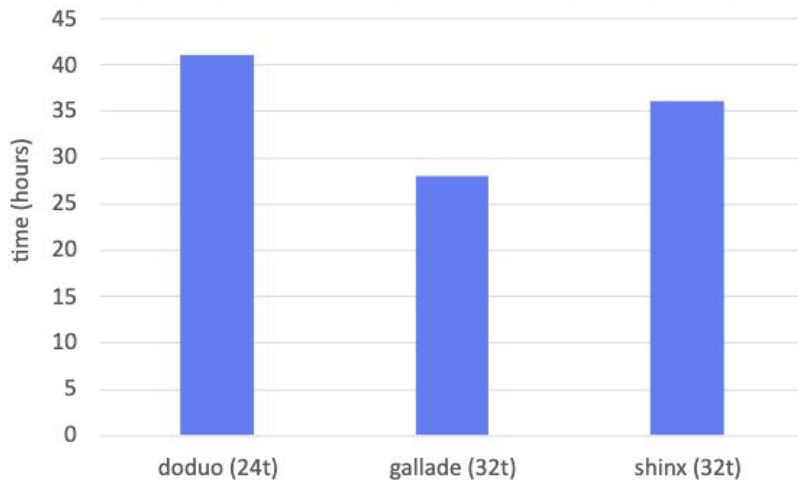
source /cvmfs/software.eessi.io/versions/2023.06/init/bash
module load BLAST+/2.14.1-gompi-2023a
export BLASTDB="$VSC_SCRATCH_ARCANINE_VO/blast/nt"

cd $PBS_O_WORKDIR
blastn -task megablast -query fastas/full-115k.fasta -db $BLASTDB/nt
-num_threads $PBS_NUM_PPN -outfmt '6 qseqid staxids bitscore std sscinames
sskingdoms stitle' -evaluate 1e-10 -max_hsp 1 -max_target_seqs 5
```

# Running BLAST+ jobs @ HPC-UGent

Can we run full input (115k seq.) using large nt database (~700GB) in < 72h? **Yes!**

- Database on \$VSC\_SCRATCH\_ARCANINE filesystem
- **Using reasonable # cores on each CPU-only cluster**
  - doduo (24 cores): ~41h
  - **gallade (32 cores): ~28h**
  - shinx (32 cores): ~36h

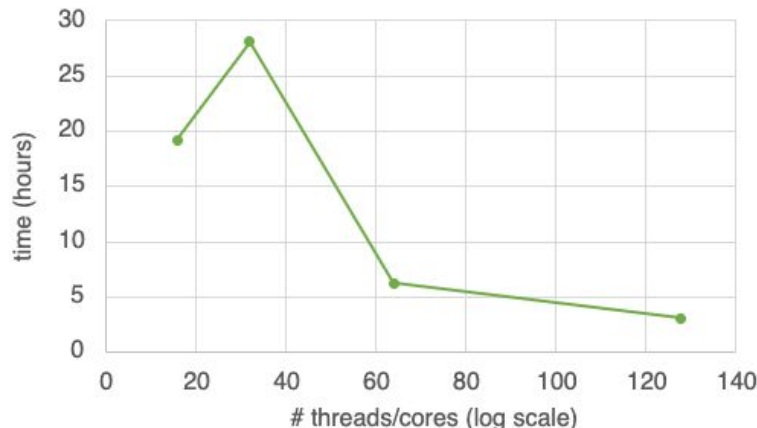


# Running BLAST+ jobs @ HPC-UGent

Can we run full input (115k seq.) using large nt database (~700GB) in < 72h? **Yes!**

- Database on \$VSC\_SCRATCH\_ARCANINE filesystem
- **Scaling with # cores on a single gallade node**
- **Done in ~3h15min with 128 cores (and ~900GB of memory)**

(and if you're lucky that load on \$VSC\_SCRATCH\_ARCANINE is sufficiently low)



# Time to get your hands dirty!

- Reservations for this hackathon:
  - On doduo cluster: 4 nodes (each w/ 96 cores and ~250GB memory)
  - On gallade cluster: 1 node (128 cores and ~940GB memory)
  - On shinx cluster: 4 nodes (each w/ 192 cores and 370GB memory)
  - On accelgor cluster: 1 node (4x A100 GPUs)

- Submit into reservation with:

```
module swap cluster/example # <== change this to name of right cluster
qsub --reservation=hackathon job-script.sh
```

- Or start interactive session in reservation via web portal
- Reservations will be active until end of today (7pm)