



Spatial analysis with Python in CSC supercomputer Puhti

Zoom, 10.03.2022

Kylli Ek, Samantha Wittke





Timetable

12:30-13:30 Introduction spatial analysis with Python in
CSC computing environments

13:30-13.45 Break

13:45-15:15 Hands-on exercises in Puhti

Non-profit state organization with special tasks



Turn over
in 2021
60 M€



Headquarters in
Espoo,
datacenter in
Kajaani



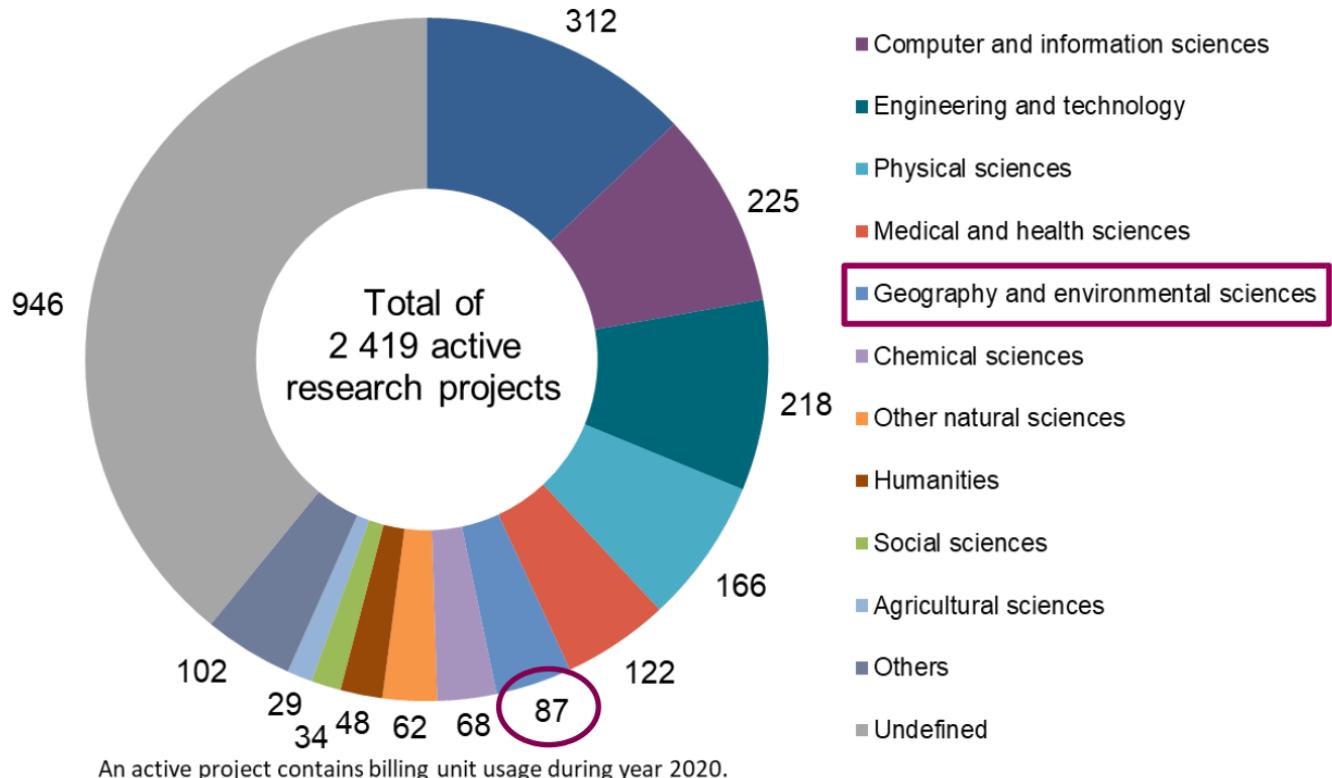
Owned by state (**70%**)
and all Finnish higher education
institutions (**30%**)



Circa
518
employees
in 2022

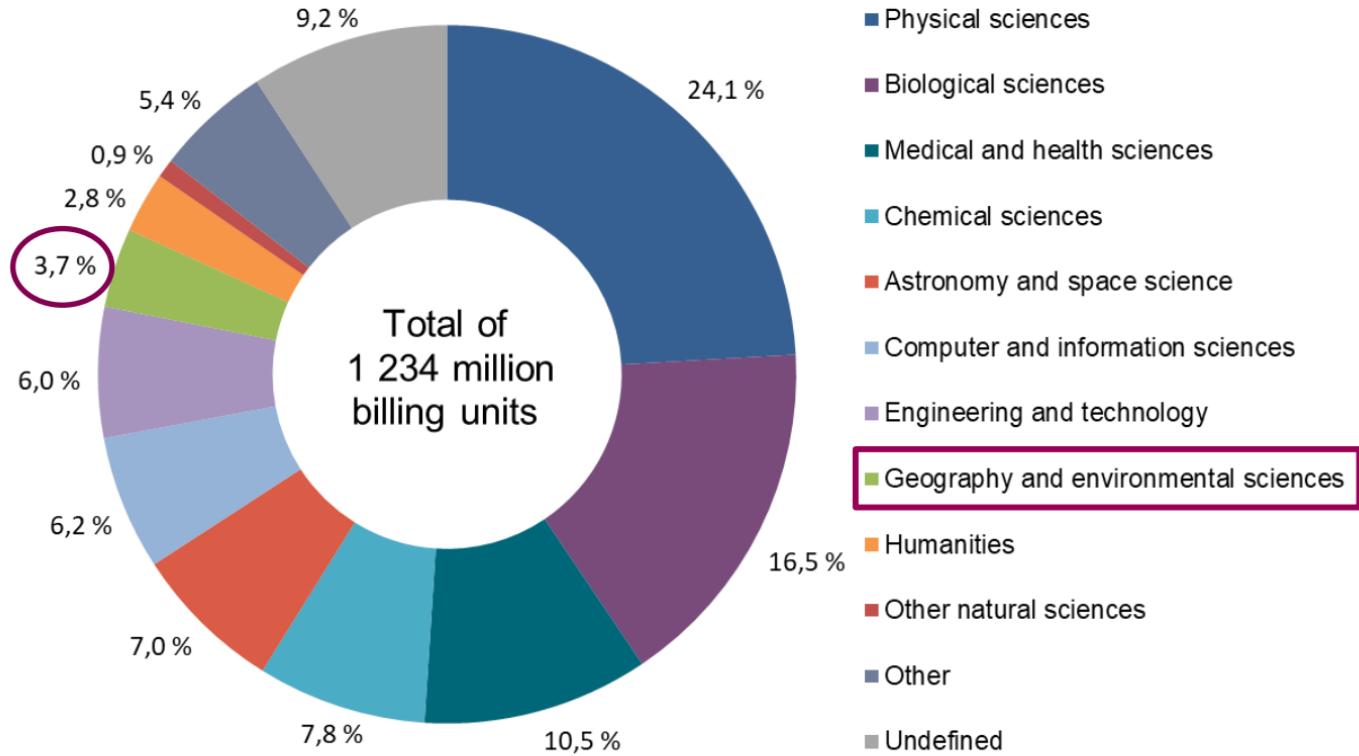
You are not alone!

Active research projects by science area in 2020 (all servers)



You are not alone!

Computing resource usage by science area in 2020 (incl. all computational and storage use)





Why using CSC resources?

- Long computing time (> 2 hours)
- High memory usage (> 8 GB)
- Large datasets (> 50 GB)
- Outsource computations, keep own computer free
- Software readily available
- Computing support
- Server needs -> cPouta
- Course computers (same setup) -> Notebooks
- It's free for open science and teaching for Finnish university and state research institute users

-> Puhti

PUHTI



	Puhti HPC	cPouta cloud
System	Supercomputer	Virtual machine cloud
Software	Pre-installed software + user-installed software	User-installed software
Data	Main Finnish datasets	-
Use cases	Run demanding analyses with numerous CPUs or GPUs	Setup your own virtual machine and environment
Max per job / VM / container	4000 CPUs / 80 GPUs 1500 GB memory	48 CPUs / 4 GPUs 240 GB memory



Average laptop:
4-8 CPUs /
8-32 GB memory

PUHTI SUPERCOMUPUTER





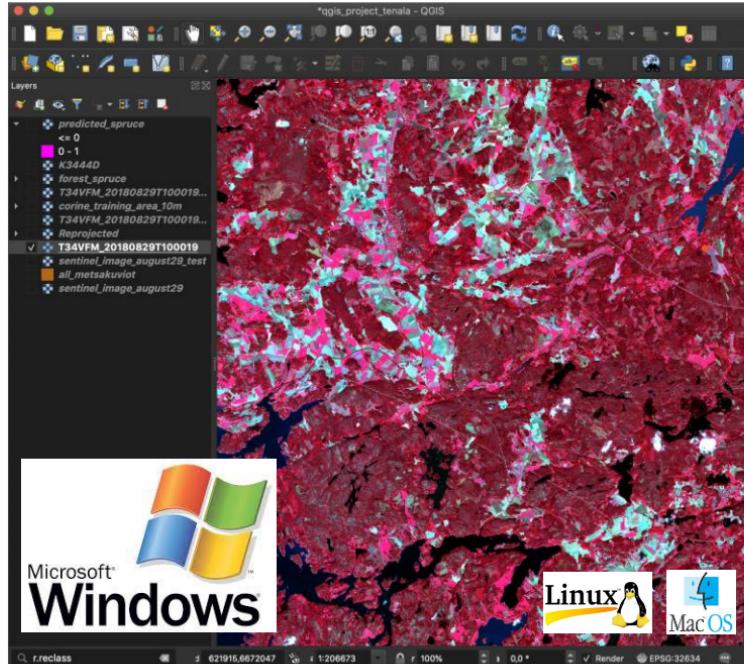
PUHTI

Supercomputers in Kajaani

- Puhti (national HPC, CSC):
 - 27 280 CPU cores and 320 GPUs
 - Wide stack of geospatial software
- Mahti (national HPC, CSC):
 - 182 784 CPU cores and 96 GPUs
 - No geospatial software
- LUMI (EuroHPC, international LUMI consortium)
 - 196 608 CPU cores and 10 240 GPUs
 - Software stack unclear
 - 20% for companies
 - CPU side in use, GPU coming later this year



Key to geocomputing in HPC: GUI → Scripts



A screenshot of a terminal window titled 'johannes@puhti-login1: ~\$'. The window shows a user guide for CSC (Tieteen tietotekniikan keskus - IT Center for Science) and a Python script named 'mergeOutFiles.py'. The script is used to merge predicted tiles and add coordinate system information. It includes comments explaining the purpose of each part of the code. On the right side of the terminal window, there is a large 'Linux' logo with Tux the Penguin.

```
puhti.csc.fi - Atos BullSequana X400 - 682 CPU nodes - 88 GPU nodes
Contact:
Servicedesk : 09-457 2821, servicedesk@csc.fi
User Guide
https://docs.csc.fi
Manage my account
https://my.csc.fi
Billing
Billing has changed significantly from
disk and lustre scratch space are all
https://docs.csc.fi/#accounts/billing/
Software
Main partitions
Small : 1-48 cores 3 days m
Large : 1:4000 cores 3 days m
hugemem : 1:168 cores 3 days m
longrun : 1:48 cores 14 days m
gpu : 1:88 GPUs 3 days m
See https://docs.csc.fi/#computing/run
Storage
In Puhti there are three main disk areas
home Personal home folder
projappj Per project Folder where app
scratch Per project Folder for runs
after 90 days
Run csc-workspaces to see your folders
See https://docs.csc.fi/#computing/disk
News
2019-08-30: Puhti opened for production
availability September 2.

[johannes@puhti-login1: ~]$
```

```
def mergeOutFiles(predicted_tiles_folder):
    """This function loops over the predicted tiles, adds the Coordinate system information to the tiles.

    Args:
        predicted_tiles_folder (str): Path to the folder containing the predicted tiles.

    Returns:
        None
    """
    list_of_rasters = []

    # Loop through all files in the predicted_tiles_folder and change their crs to match the reference tile
    for filename in os.listdir(predicted_tiles_folder):
        if filename.endswith('.tif'):
            full_filepath = os.path.join(predicted_tiles_folder, filename)
            reference_tile = rasterio.open(os.path.join(prediction_image_tile_subfolder, filename))
            with rasterio.open(full_filepath, 'r') as raster:
                raster.crs = reference_tile.crs
                raster.transform = reference_tile.transform

            raster = rasterio.open(full_filepath)
            list_of_rasters.append(raster)

    print("Successfully added CRS information to the predicted tiles")
    print(list_of_rasters)
    mosaic, out_trans = rasterio.merge.merge(list_of_rasters)
    print(mosaic.shape)
    out_metafile = raster.meta.copy()

    out_metafile.update({
        "driver": "GTiff",
        "height": mosaic.shape[0],
        "width": mosaic.shape[1],
        "transform": out_trans,
        "crs": "proj4text+zone=33N"
    })
    print(out_metafile)

    output_path = os.path.join(home_to_projects, 'mosaic.tif')
    with rasterio.open(output_path, 'w', **out_metafile) as dest:
        dest.write(mosaic)
```

Scripts: Python, R, shell, Matlab, ...



PUHTI

Keep it real

- Single core of Puhti = ~speed of laptop
- But: Puhti has **many cores**
- ... and **more memory** and **faster input-output (I/O)**
 - Running single core scripts on Puhti does not make them faster
 - Speedup: multi-core parallel processes
 - ... or optimize your script



Geospatial Software in Puhti

geo
portti

Finnish Geospatial
Research and
Education Hub

- ArcGIS Python API
- FORCE & SPLITS
- GDAL/OGR
- GRASS GIS
- LasTools
- MatLab / Octave
- Mapnik
- OpenDroneMap
- Orfeo Toolbox
- PCL
- PDAL
- Python geospatial packages: geoconda
- QGIS, for viewing point clouds
- R geospatial packages
- SagaGIS
- SNAP, Sen2cor
- WhiteboxTools
- Zonation
- Deep learning: Keras, PyTorch
- Something missing?
 - Ask us :)
servicedesk@csc.fi

Parallel computing with geospatial software

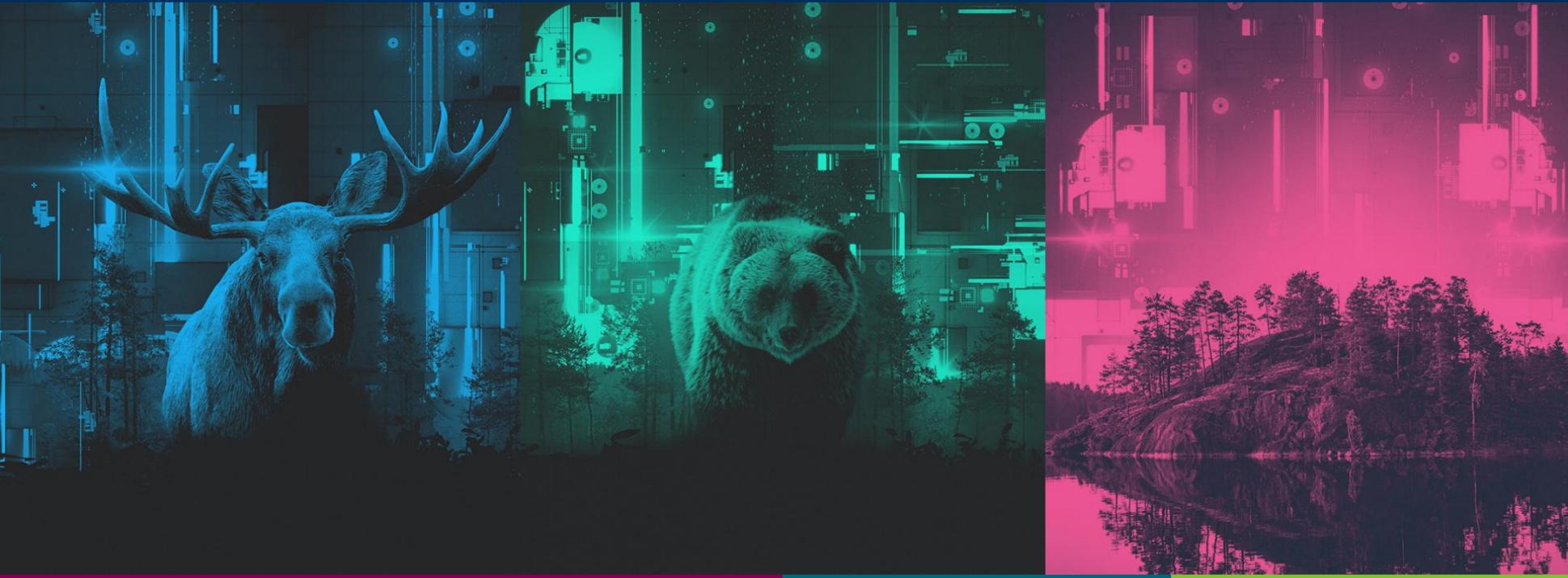
- Some tools have at least some support built-in
 - lidR, OpenDroneMap, Whiteboxtools (?)
 - SAGA GIS, GRASS GIS (some parts)
- Parallel libraries for Python
 - dask
 - multiprocessing
 - Joblib
 - ...
- Parallel libraries for R
 - future
 - snow
 - foreach
 - ...

Example scripts

<https://github.com/csc-training/geocomputing>

- Batch job files
- Parallel examples for **R** and **Python**
 - Different parallelization libraries
 - Array jobs as well as parallel jobs
- Examples for **Allas** data transfers with **R** and **Python**
- Sentinel image download example (python)
- **SNAP** array job example

Running Python in parallel — principles and practice



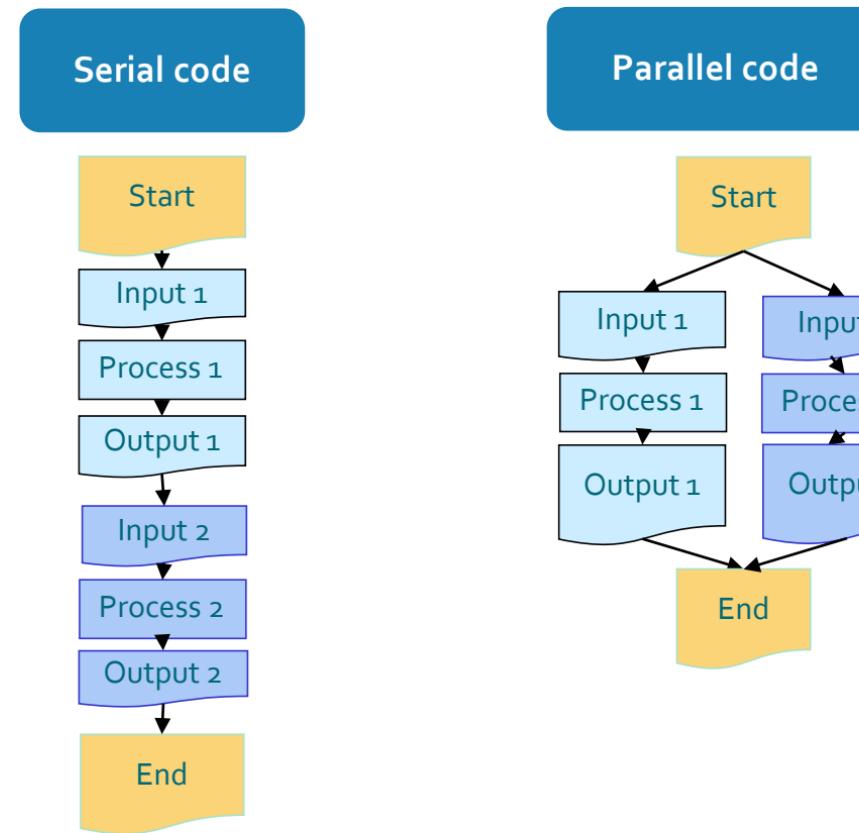


PUHTI

"My Python code is slow... what can be done?"

- Try to understand which part of the code takes time and why
 - Use `time` or `timeit` packages
- Always be suspicious of `for`-loops!
- Going parallel *may* help
- Unfortunately, increasing the number of cores will *never* decrease the time needed in the same proportion

PUHTI





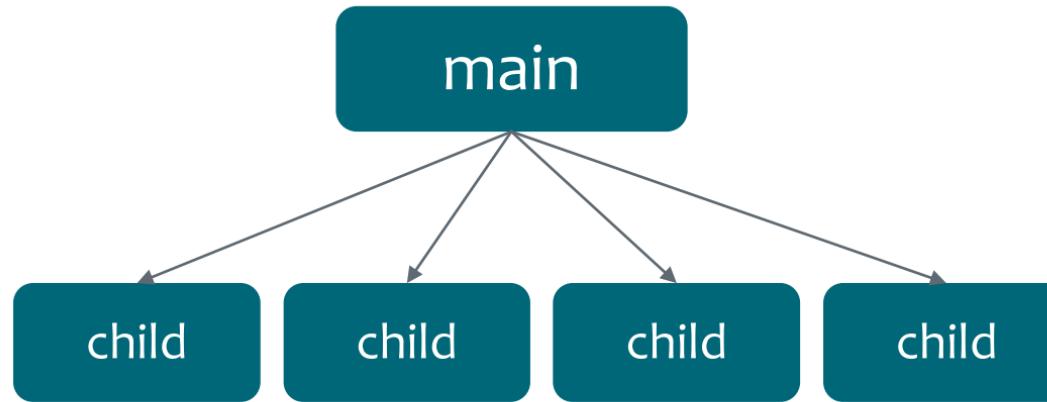
PUHTI

Parallelization in GIS

- Easiest parts for parallelization:
 - for loops
 - map
- Logical level
 - input files / map sheets / chunks of vector data,
 - scenarios, variables, time periods
- In many cases if using map sheets the borders need special care.

multiprocessing

- Parallelization package in Python
- Alternatives: joblib, dask



multiprocessing

function and input

```
def slow_function(i):
```

```
    time.sleep(3)
```

```
    return(i)
```

```
input = range(1, 7)
```

for loop

```
a = []
```

```
for i in input:
```

```
    a.append(slow_function(i))
```

map

```
a = list(map(slow_function, input))
```

SERIAL

multiprocessing

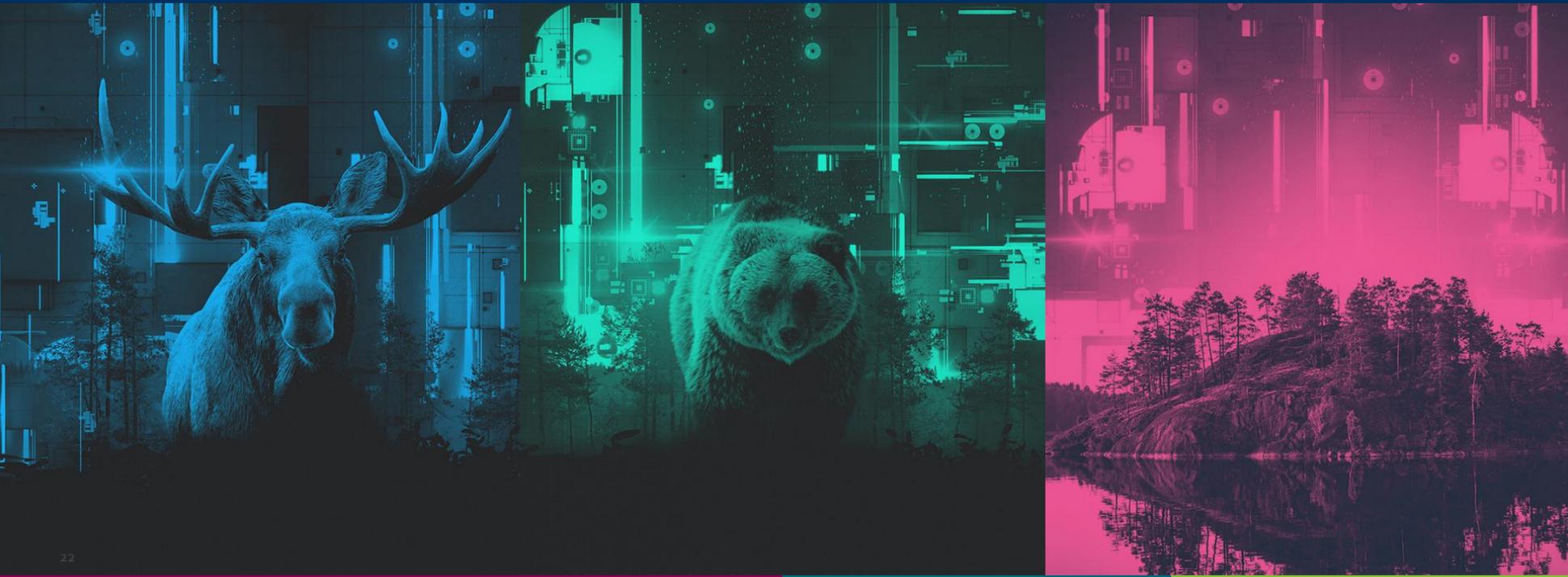
PARALLEL

```
from multiprocessing import Pool
```

```
pool = Pool(os.sched_getaffinity(0))
```

```
a = pool.map(slow_function, input)
```

DATA



Directories

Directory	Owner	Capacity	Num. files	Path	Cleaning
home (\$HOME)	Personal	10 GiB	100 000	/users/<user-name>	No
projappl	Project	50 GiB	100 000	/projappl/<project>	No
scratch	Project	1 TiB	1 000 000	/scratch/<project>	Yes - 90 days

- The quota of ***scratch*** and ***projappl*** directories can be increased
- **No back-up**
- More information: <https://docs.csc.fi/computing/disk/>



ALLAS

Allas object storage

- Central data service for all academic projects
- Up to 200 Tb of space for a project
- Access from Puhti, cPouta, personal computer or any other place
- GDAL has good support for reading directly from object storage
- Option to make data public

See also: CSC webinar - **Allas and Geospatial data**
https://www.youtube.com/watch?v=mnFXe2-dJ_g

Geospatial data in Puhti

- Hosts large commonly used datasets with open license
- Removes transfer bottleneck
- Located at: </appl/data/geo/>
- All Puhti users have **read** access

Datasets available (~ 13TB)

- Paituli data +
 - SYKE open datasets
 - LUKE Multi-source national forest inventory
 - NLS Virtual rasters for DEMs
 - Sentinel and Landsat mosaics

- Finnish Digital and Population Data Services Agency
- Finnish Food Agency
- Finnish Meteorological Institute (FMI)
- Finnish Transport Infrastructure Agency, Digiroad
- Institute for the Languages of Finland (KOTUS)
- Karelia University of Applied Sciences
- Latuviitta
- National Land Survey (MML)
- Natural resource institute Finland (LUKE)
- Statistics Finland
- University of Helsinki, Digital Geography Lab

<https://docs.csc.fi/data/datasets/spatial-data-in-csc-computing-env/>

Want more? - contact us at servicedesk@csc.fi



Virtual rasters in Puhti

- Ready made virtual rasters for 2m and 10m DEMs
- Multiple files usable as if they were a single file
- XML pointing to actual raster files
- External overviews and xml headers
- Possible to have all data in Allas and only virtual raster in Puhti

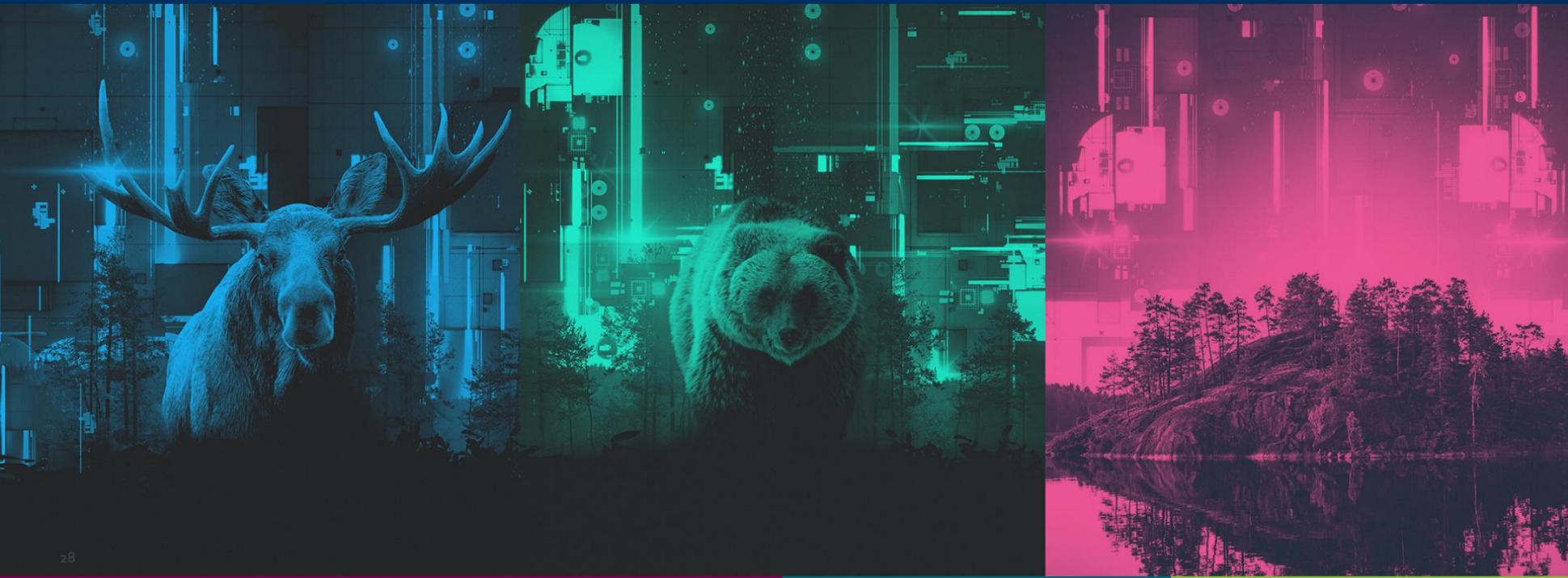
Python script to create your own virtual raster available

Data: tips for moving files

- When downloading from external services, download directly to CSC, not via your local PC
- When lot of/big files use :
 - command-line tools: rsync, wget, curl
 - Python/R packages for moving files

<https://docs.csc.fi/data/moving/scp/>

GETTING STARTED





PUHTI

Set up for Puhti, administrative

- Check your **eligibility** for using CSC resources:
<https://research.csc.fi/free-of-charge-use-cases>
- Get a **user account**: <https://docs.csc.fi/accounts/how-to-create-new-user-account/>
- Create or let your PI **create a new project**:
<https://docs.csc.fi/accounts/how-to-create-new-project/>
- **Add members** to your project: <https://docs.csc.fi/accounts/how-to-add-members-to-project/>
- **Add services** to your project: <https://docs.csc.fi/accounts/how-to-add-service-access-for-project/>

Set up for Puhti, Billing Units (BU)

- Each project has a certain amount of so-called Billing Units (BU)
- Using CSC resources consumes Bus, based on:
 - Computation time and number of resources (CPU, memory)
 - GPUs in Puhti (**a lot of BUs!**)
 - Increased storage quota in Puhti
 - Used storage in Allas
 - Cpouta uptime (flavor dependent)

You can apply for more BUs in my.csc.fi by providing a description what what you are doing



PUHTI

Set up for Puhti, tools

Minimum:

- web-browser for Puhti web interface

Recommended:

- **Moving data:**
 - FileZilla/WinSCP, rsync, wget, curl
- **Connecting:**
 - Terminal in Linux/Mac
 - PowerShell/Putty/MobaXterm in Windows



Puhti web interface

- Graphical interface to Puhti
- Used with web-browser
- Visual Studio Code, Jupyter Notebooks, Rstudio, TensorBoard
- QGIS, SNAP, GRASS, SagaGIS
- File browser
- Terminal

The screenshot shows a file browser window with the following details:

- Home Directory:** /projapp/project_2000745
- Current Path:** /scratch/project_2000745
- Buttons:** Open in Terminal, New File, New Directory, Upload, Download, Copy/Move, Delete.
- File List:**

Type	Name	Size	Modified at
Folder	ktegel	-	3/7/2022 8:23:19 PM
Folder	training081	-	3/8/2022 10:39:29 AM

The screenshot shows the RStudio Server interface with the following details:

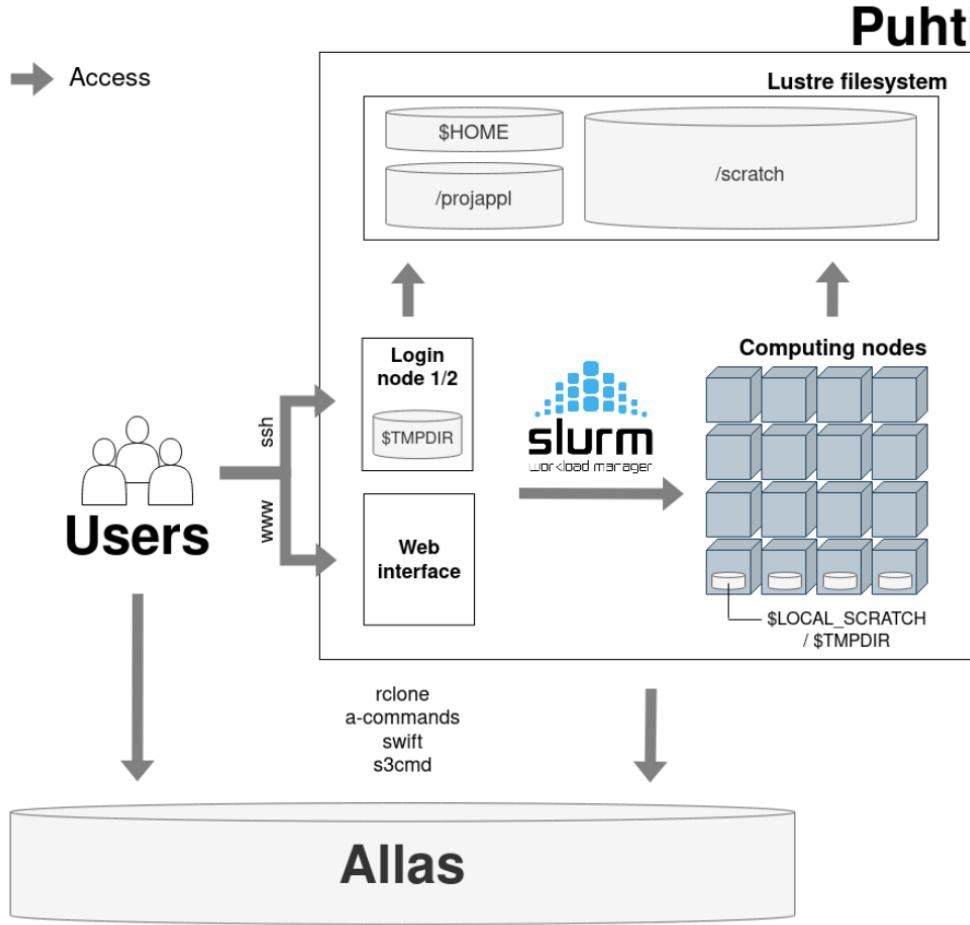
- RStudio Server:** puhti.csc.fi:port/10251/build/15973/
- R Script:** TestScript.R (netcdf_test.R)

```
library(raster)
raster_brick <- brick('/scratch/project_2000599/raster_raster')
raster_raster <- raster('/scratch/project_2000599/raster_rast')
library(ncdf4)
a <- nc_open('/scratch/project_2000599/tas_2001')
```

- Environment:** Global Environment
- Data:**
 - a
 - raster_brick
 - raster_raster
 - terra_raster
- Console:** R 4.1.1 ->

```
R version 4.1.1 (2021-08-10) -- "Kick Things"
Copyright (C) 2021 The R Foundation for Statistics
1 Computing
Platform: x86_64-pc-linux-gnu (64-bit)
```

- File Browser:** Shows files like .Renvironment, .Rhistory, and dash.worker-space.



The module system

- Puhti is a shared computing environment with hundreds of users
- Software is loaded with **modules**
 - Mutually incompatible software
 - One module: single program or group of similar programs
- **Example.** Loading module for geospatial Python tools

```
module load geoconda
```

- Check: <https://docs.csc.fi/apps/#geosciences> for module names



geoconda

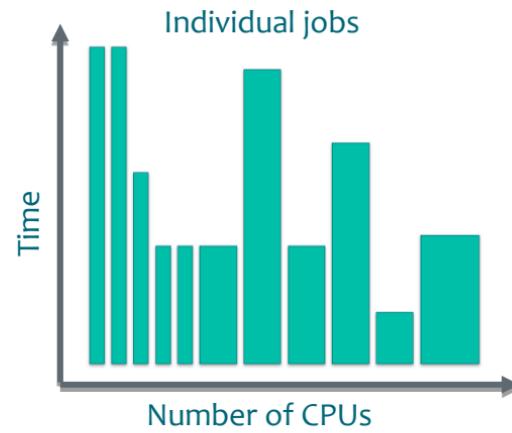
Module of pre-installed conda environment with:

- all main packages raster, vector and lidar data analysis
- scikit-learn for machine learning
- packages for accessing data in Allas
- possibility to add packages with pip

<https://docs.csc.fi/apps/geoconda/>

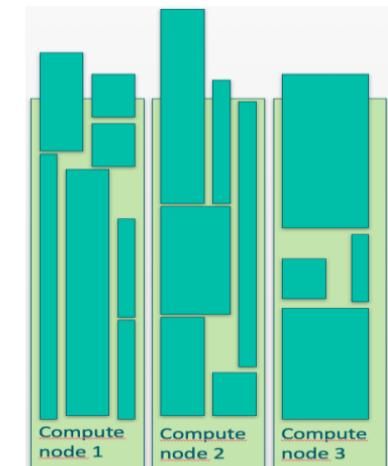
The batch job system (SLURM)

- Running jobs in Puhti requires you to use the batch job system
- You request resources for a batch job
 - CPU cores, memory, GPU etc.
 - Time
- Batch job optimization of the queue by SLURM:



SLURM places jobs
on computing nodes

In the most efficient
way resource wise





PUHTI

Batch job partitions

Partition	Time limit	Max tasks	Max memory
test	15 minutes	80	190 GiB
interactive	7 days	8	64 GiB
small	3 days	40	382 GiB
large	3 days	1040	382 Gib
longrun	14 days	40	382 GiB
hugemem	3 days	160	1534 GiB
hugemem_longrun	7 days	40	1534 GiB

Batch job scripts

- For requesting resources and submitting job description
- Specific bash scripts (.sh)

```
#!/bin/bash
#SBATCH --job-name=myTest
#SBATCH --account=<project>
#SBATCH --time=02:00:00
#SBATCH --cpus-per-task=4
#SBATCH --mem-per-cpu=2000
#SBATCH --partition=small

module load geoconda

srun python my_python_script.py
```

- Batch job is submitted with the command
sbatch <your-batch-job-script>
- Cancel the job with
scancel <your-job-id>
- See if your job has started running
queue -u <your-user-name>
- After the job, see how much resources it used
seff <your-job-id>

Geospatial parallel jobs

- With parallel jobs you submit one job but reserver for it plenty of CPU cores
- Your script has to divide the workload to different cores, otherwise 1 core is doing the work, others are just idling.
- In GIS, you often divide the dataset and give each worker (CPU core) their own subset. It could be vector features, raster subsets, or text inputs in a .CSV file
- How many workers should you utilize depends how long handling one subset takes. Communication always takes extra time so your workload for one worker should not take less than ~5 minutes



Steps for running a Python script in Puhti

1. (Set up: account, project, services, BU)
2. Log in to Puhti.
3. Move your data and scripts to Puhti.
4. Open VSCode
5. Check which Python packages do you need and if they are available in Puhti.
 - * If needed, install it yourself or ask CSC - servicedesk@csc.fi.
6. Fix the paths of your input/output files.
7. Test your script with some test data
8. Run your scripts with all data as batch job (or interactively).
9. (Make use of several cores in your Python code.)



Puhti, next steps

PUHTI

Read Puhti documentation

- Geocomputing: <https://research.csc.fi/geocomputing>
- GIS software specific pages: <https://docs.csc.fi/apps/#geosciences>
- Puhti GIS examples: <https://github.com/csc-training/geocomputing>
- CSC Linux tutorial for basic Linux commands:
<https://docs.csc.fi/support/tutorials/env-guide/overview/>

Come to a course

- Using CSC HPC environment efficiently: <https://ssl.eventilla.com/event/ge9Ey>

Contact CSC

- servicedesk@csc.fi

Join gis-hpc-mailing list

- <https://postit.csc.fi/sympa/info/gis-hpc>



PUHTI

Summary

Puhti is an excellent tool if you

- need more computing power
- don't want to run long analyses on your personal computer
- have **a lot** of data
- are using data already available in Puhti in large quantities
- are willing to use scripts for your work
- have really basic Linux skills
- are willing to learn to use Puhti