

# Running Python code in CSC Taito supercluster

## Kylli Ek, CSC

Espoo, 13.11.2018



*CSC – Suomalainen tutkimuksen, koulutuksen, kulttuurin ja julkishallinnon ICT-osaamiskeskus*

Non-profit state organization with special tasks



Turnover in 2017  
**40,5** M€



Headquarters in Espoo, datacenter in Kajaani



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



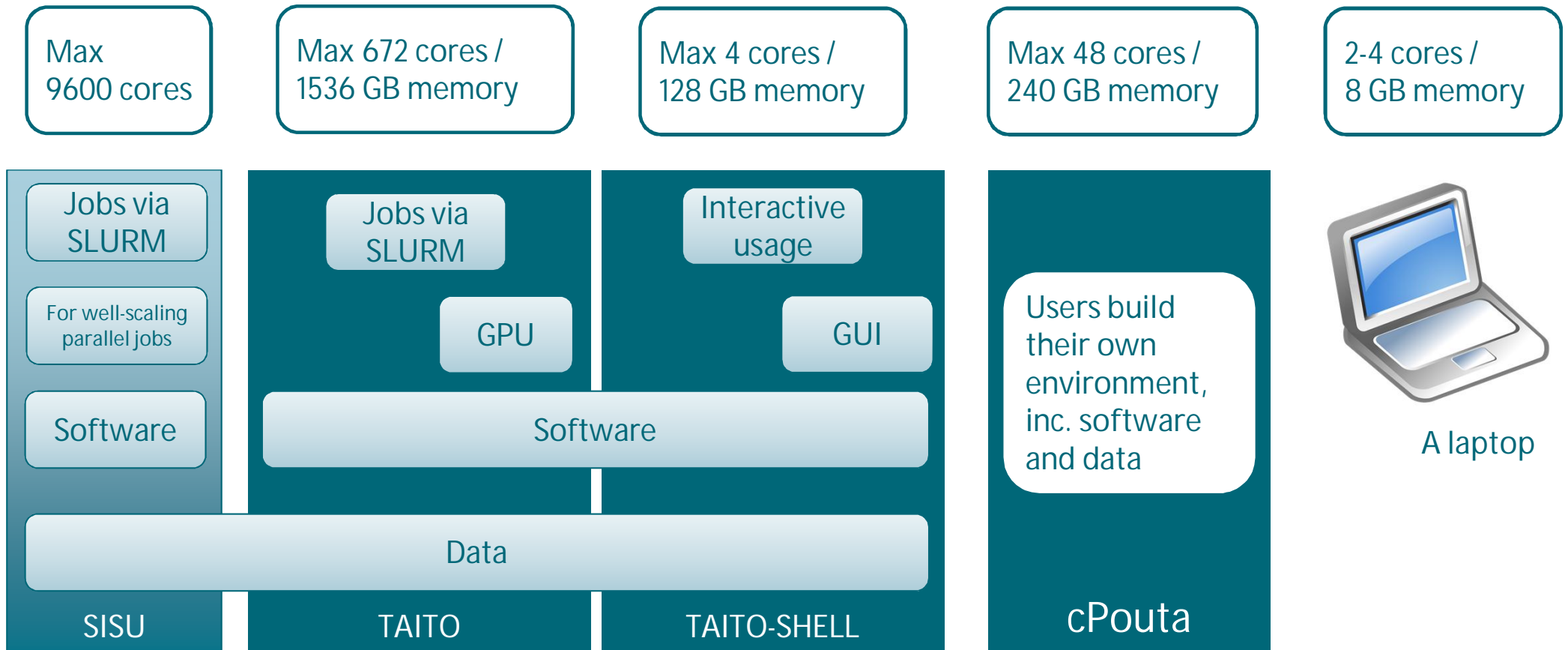
Circa  
**320**  
employees in 2017

## Reasons for using CSC computing resources

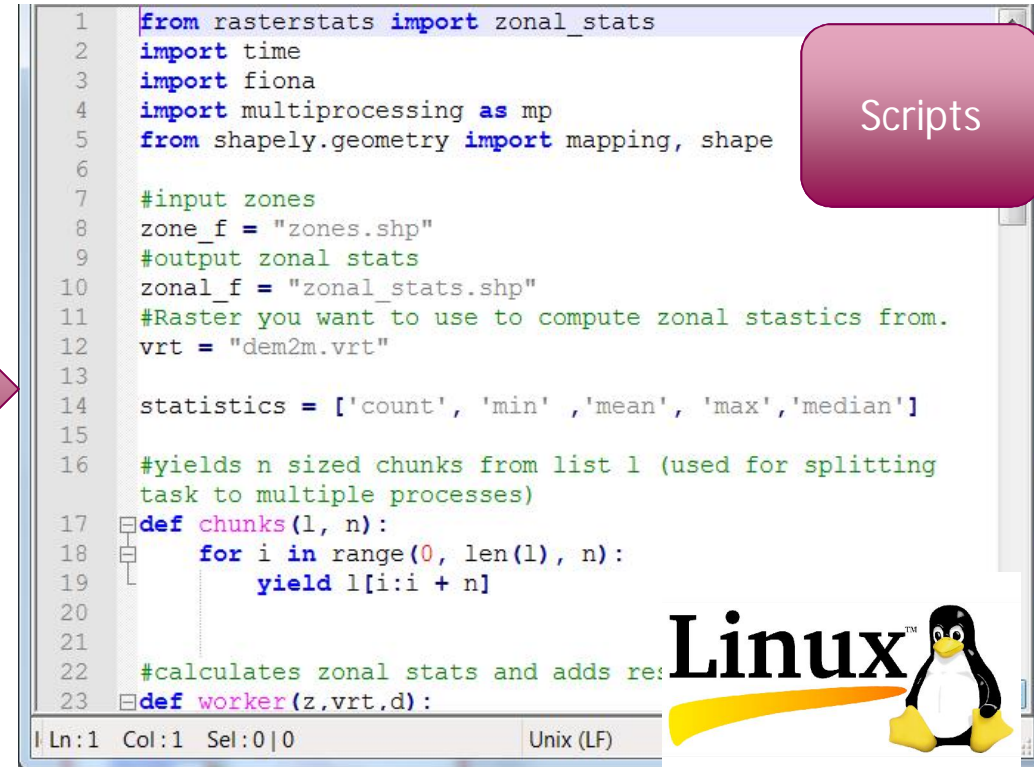
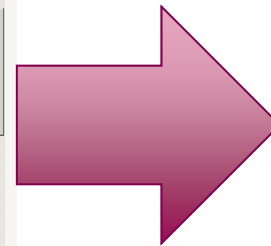
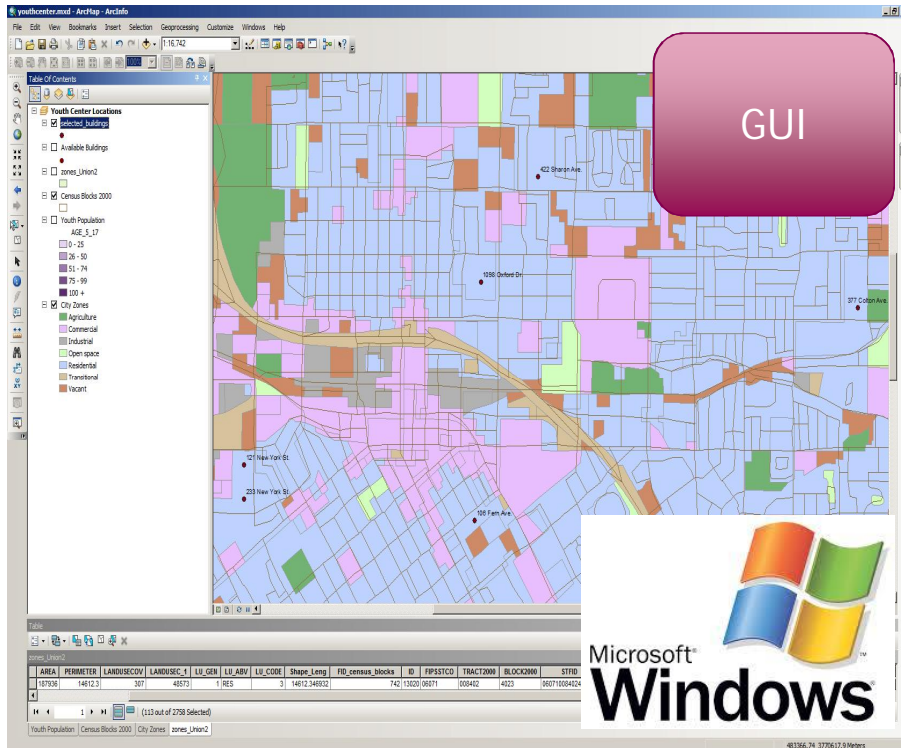
- Computing something takes more than 2-4 hours
- Need for more memory
- Very big datasets
- Keep your desktop computer for normal usage, do computation elsewhere
- Need for a server computer
- Need for a lot of computers with the same set-up (courses)
- Free for Finnish university users and for state research institutes



# CSC HPC resources



# The keys to geocomputing: Change in working style & Linux



ArcGIS, QGIS, ...

R, Python, shell scripts, Matlab, ...

## Taito / Taito-shell pre-installed software for GIS

- R
- Python
- MatLab / Octave
- GDAL/OGR
- GRASS GIS
- LasTools (some)
- PDAL
- Proj4
- QGIS
- SagaGIS
- SNAP, sen2cor
- Taudem
- Zonation

<https://research.csc.fi/software> -> Geosciences

# Geoconda

- NumPy, Scipy, Pandas etc.
- GIS Specific packages
  - Geopandas
  - Fiona
  - Shapely
  - Rasterio
  - Rasterstats
  - cartopy
  - GDAL/OGR
  - Networkx
  - Skimage
  - Pyproj
  - Pysal
  - Rtree
  - Descartes

```
module load geoconda
```

## Installing software

- Possibility to install software for own use
  - The software must be available for Linux
  - .. and installation must be possible without root access
- You can add also packages to Python with pip
- You can also make your own conda environments



## Realistic expectations

- A single core of a CSC machine is about as fast as one of a basic laptop.
- It has just a lot of them.
- .. and more memory and faster input-output.
  - Just running your single core script at CSC does not make it much faster.
  - For clear speed-ups you have to use several cores.
  - ... or optimize your script.

## Shared data area in Taito

- Hosts large commonly used datasets
- Reduces the need to transfer data to Taito
- Located at /proj/ogiir-csc/
- All Taito users have read access.
- Only CSC personnel have write access.
- For data with open license
  
- If you think some other dataset should be included here, ask from servicedesk@csc.fi

All Paituli open data

+

LUKE

Multi-source national forest inventory

NLS

Virtual rasters for DEMs

SYKE

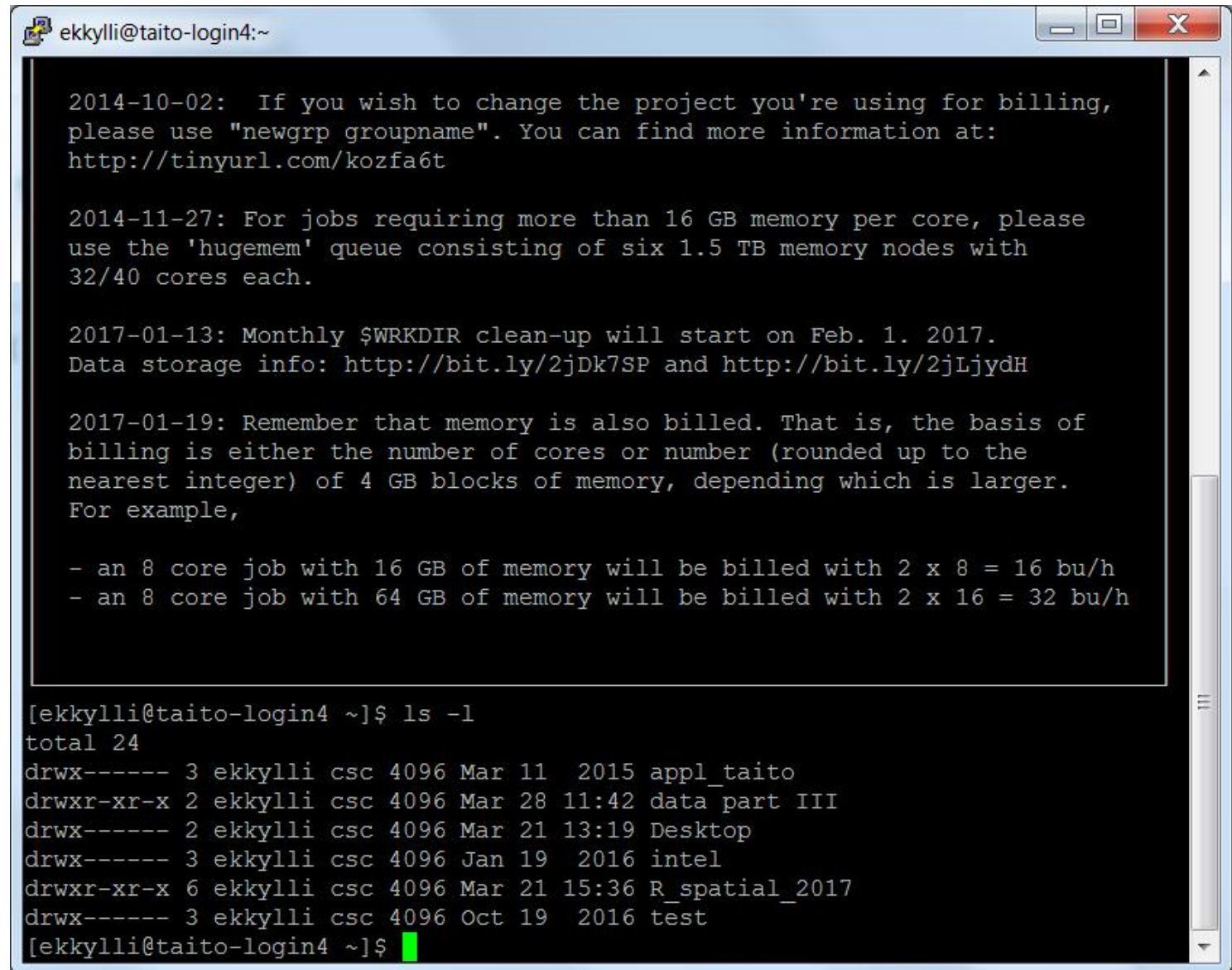
All open GIS data

[https://research.csc.fi/gis\\_data\\_in\\_taito](https://research.csc.fi/gis_data_in_taito)

## Access to Taito from Windows

- Putty for ssh connection
- FileZilla/WinSCP for moving data
- NoMachine for GUI
- Find about other access options and more information at:  
<https://research.csc.fi/taito-connecting>

# Putty

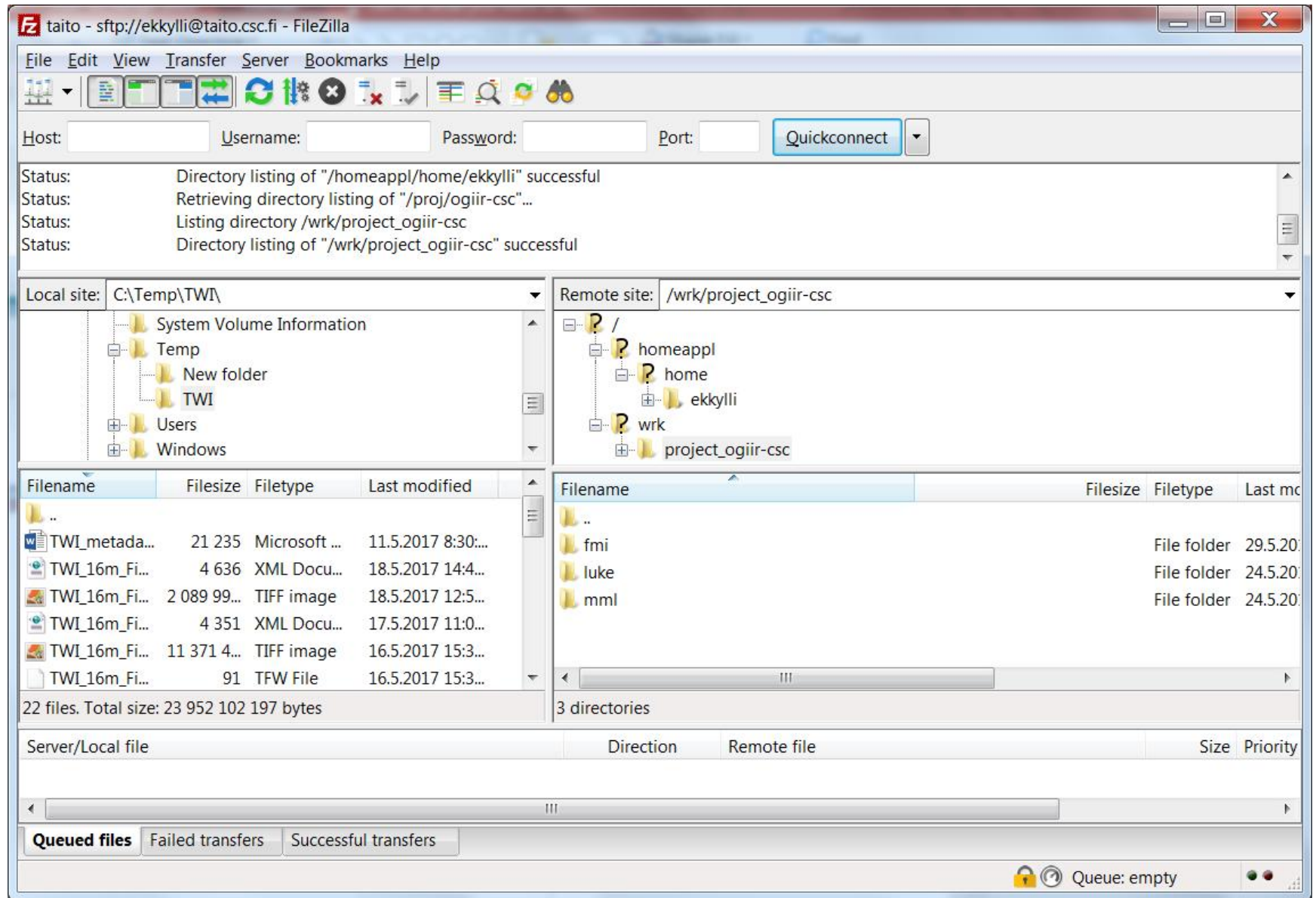


The screenshot shows a Putty terminal window titled "ekkylli@taito-login4:~". The terminal displays several system messages and a directory listing. The messages are dated and provide information about billing changes, memory requirements, and clean-up schedules. The directory listing shows the contents of the user's home directory.

```
ekkylli@taito-login4:~  
  
2014-10-02: If you wish to change the project you're using for billing,  
please use "newgrp groupname". You can find more information at:  
http://tinyurl.com/kozfa6t  
  
2014-11-27: For jobs requiring more than 16 GB memory per core, please  
use the 'hugemem' queue consisting of six 1.5 TB memory nodes with  
32/40 cores each.  
  
2017-01-13: Monthly $WRKDIR clean-up will start on Feb. 1. 2017.  
Data storage info: http://bit.ly/2jDk7SP and http://bit.ly/2jLjyDH  
  
2017-01-19: Remember that memory is also billed. That is, the basis of  
billing is either the number of cores or number (rounded up to the  
nearest integer) of 4 GB blocks of memory, depending which is larger.  
For example,  
  
- an 8 core job with 16 GB of memory will be billed with 2 x 8 = 16 bu/h  
- an 8 core job with 64 GB of memory will be billed with 2 x 16 = 32 bu/h  
  
[ekkylli@taito-login4 ~]$ ls -l  
total 24  
drwx----- 3 ekkylli csc 4096 Mar 11 2015 appl_taito  
drwxr-xr-x 2 ekkylli csc 4096 Mar 28 11:42 data part III  
drwx----- 2 ekkylli csc 4096 Mar 21 13:19 Desktop  
drwx----- 3 ekkylli csc 4096 Jan 19 2016 intel  
drwxr-xr-x 6 ekkylli csc 4096 Mar 21 15:36 R_spatial_2017  
drwx----- 3 ekkylli csc 4096 Oct 19 2016 test  
[ekkylli@taito-login4 ~]$
```



# FileZilla



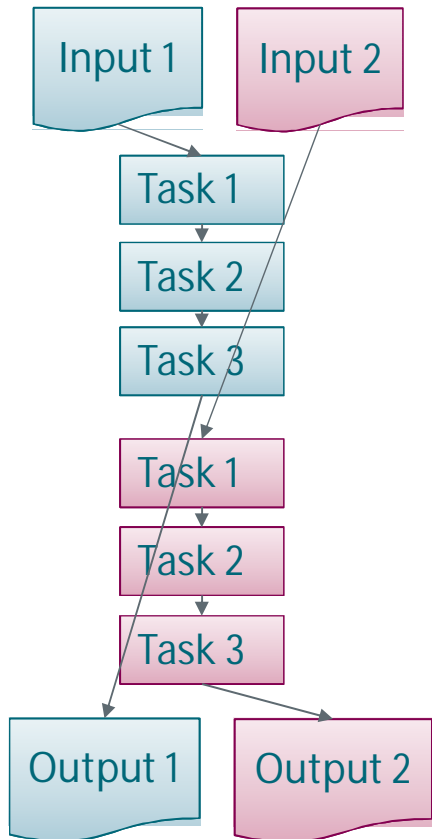
# NoMachine



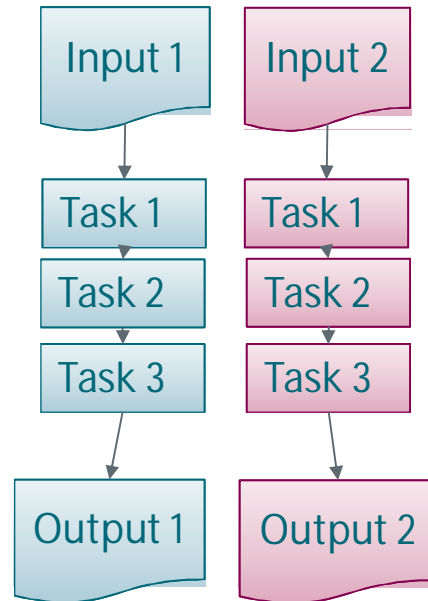
The screenshot displays a NoMachine session titled "NoMachine - Connection to nmkajaani.csc.fi". It features three main windows:

- Terminal Window (Left):** Shows shell commands and instructions for running R on Taito computing nodes. It includes instructions for using `rsun` and `Rmpi` for interactive work, and notes that starting R directly on a computing node will be killed. It also provides a URL for batch jobs: <https://research.csc.fi/~r/>.
- RStudio Editor (Center):** Displays the R script `Calc_contours.R`. The script defines command arguments, sets default output folders, and loads the `RSAGA` and `rgdal` libraries. It also includes code to create directories and download a DEM file from Kapsi.
- Console Window (Bottom):** Shows the execution of the R script, including directory creation and file downloading.
- Environment and Files Panels (Right):** The Environment panel shows variables like `args`, `gridFolder`, `imageFolder`, `inputfile`, `mainDir`, `mapsheet`, `shapeFolder`, `tiffFolder`, and `url`. The Files panel shows a directory listing for `~/R_spatial_2017` with various files and folders.

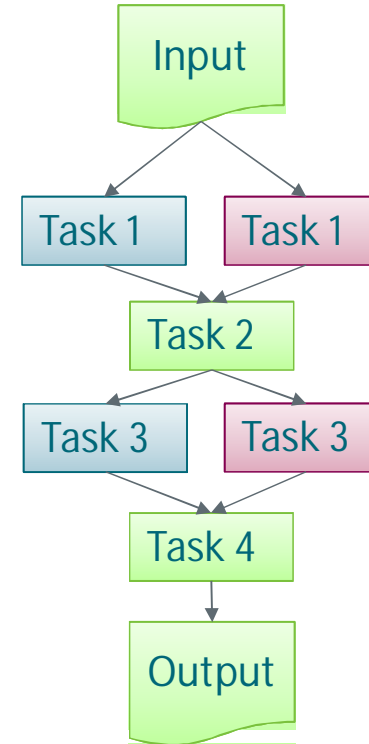
## Simple job



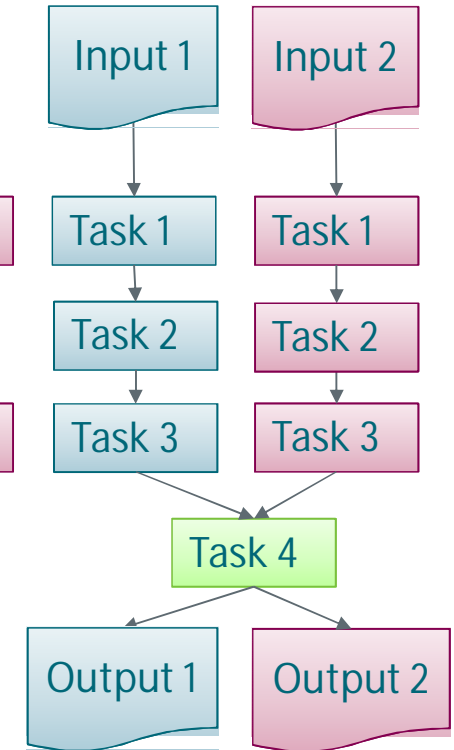
## Array job



## Parallel job 1



## Parallel job 2



## Example: steps for running your Python script in Taito

(0. Get yourself CSC user account)

1. Move your data and scripts to Taito (with FileZilla).

2. Log in to Taito (with Putty).

3. Open Spyder in Taito-shell with NoMachine.

4. Check which Python packages do you need and if they are available in Taito.

\* If needed, install it yourself or ask CSC - [servicedesk@csc.fi](mailto:servicedesk@csc.fi).

5. Fix the paths of your input/output files.

6. Test your script in Taito-shell with some test data.

7. Run your scripts with all data interactively on Taito-shell or in Taito as batch job.

(8. Make use of several cores using multiprocessing package in your Python code or with array jobs.)



# Directories at CSC Environment



<https://research.csc.fi/data-environment>

Directory or storage area	Intended use	Default quota/user	Storage time	Backup
\$HOME <sup>1</sup>	Initialization scripts, source codes, small data files. Not for running programs or research data.	50 GB	Permanent	Yes
\$USERAPPL <sup>1</sup>	Users' own application software.	50 GB	Permanent	Yes
\$WRKDIR <sup>1</sup>	Temporary data storage.	5 TB	90 days	No
\$WRKDIR/DONOTREMOVE	Temporary data storage.	Incl. in above	Permanent	No
\$TMPDIR <sup>3</sup>	Temporary users' files.	-	~2 days	No
Project <sup>1</sup>	Common storage for project members. A project can consist of one or more user accounts.	On request	Permanent	No
HPC Archive <sup>2</sup>	Long term storage.	2 TB	Permanent	Yes
IDA <sup>2</sup>	Storage and sharing of stable data.	On request	Permanent	No, multiple storage copies

<sup>1</sup>: Lustre parallel (<sup>3</sup>:local) file system in Kajaani    <sup>2</sup>: iRODS storage system in Espoo

## Batch system

- Has to be used on Taito (not in Taito-shell)
- Optimizes resource usage by filling the server with jobs
- You have to reserve time, cores and memory for your job
- Several queues: parallel, serial, longrun, test and hugemem
- You have to write a batch job script
- <https://research.csc.fi/taito-batch-jobs>

## Taito module system

- Tool to set up your environment
  - Load libraries, adjust path, set environment variables
  - Needed on a server with hundreds of applications and several compilers etc.
- Example: initialize Python with GIS packages
  - `$ module load geoconda`

## Accounts

- University users can create an account themselves in SUI:  
<https://research.csc.fi/accounts-and-projects>
- University users can start using Taito without project with the default quota.
- Research institute users have to ask for account from supportdesk.
- For serious work create a project and apply for resources.
- For cPouta you always need a project.

## Billing units

- Each project is given certain amount of so-called billing units (BU).
- On Taito, if you are using batch jobs, the billing is based on actual time used, but on the number of cores and memory reserved.
- If you need help with estimating your job resource needs, see the seff command from the end of [this page](#) or see the webinar about estimating needed memory: <https://www.youtube.com/watch?v=4ThGRZq1G8U>
- Changing billing project: <https://research.csc.fi/billing-and-monitoring>
- Project saldo, to see how much BUs you have used: <https://research.csc.fi/saldo>

## Example code in CSC training Github

- Examples for doing spatial analysis in CSC computing environment with:
  - Python
  - R
- Examples include also batch job scripts suitable for Taito.
- Some of the examples include samples for serial, array and parallel jobs.

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



## GIS training 2018/19

<https://www.csc.fi/training>

- Geospatial data analysis with R
- Google Earth Engine
- Lidar data analysis
- Webinar: Paituli, cPouta for GIS

## Further information

- <https://research.csc.fi/taito-user-guide>
- <https://research.csc.fi/geocomputing>
- CSC Python documentation: [research.csc.fi/-/python](https://research.csc.fi/-/python)
- Geoconda module: <https://research.csc.fi/-/geoconda>
- Geo-env module: <https://research.csc.fi/-/geo-env>
- Multiprocessing: <https://docs.python.org/3/library/multiprocessing.html>

Training archive:

<https://www.csc.fi/web/training/-/geocomputing-in-taito>

Support: [servicedesk@csc.fi](mailto:servicedesk@csc.fi)

GIS@CSC e-mail list: [gis-hpc](mailto:gis-hpc)

<http://research.csc.fi/gis-csc-news>





## Contact

<http://research.csc.fi/geosciences>

Kylli Ek, +358 50 38 12 838

[giscoord@csc.fi](mailto:giscoord@csc.fi)