



university of
groningen

center for
information technology

High Performance
Computing

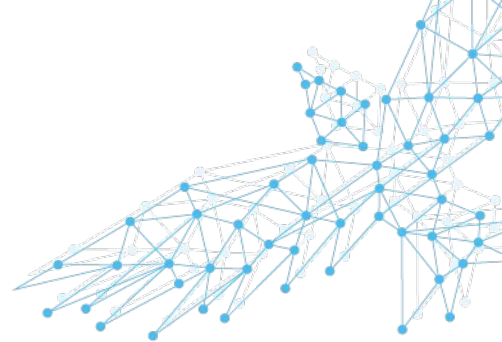
CI

From Peregrine to Hábrók

Fokke Dijkstra, Bob Dröge - 27-3-2023



Contents

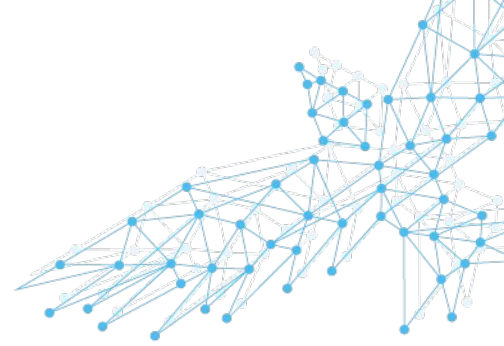


- Peregrine
- Design goals Habrok
- Habrok choices
- Migration
- Changes
 - Accounts, MFA
 - Logging in
 - Storage areas
 - Partitions
 - GPUs

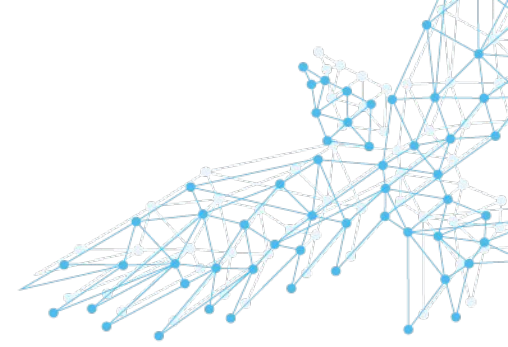


Almost 8 years of Peregrine

- Operational summer 2015
- Heavily used, increasing number of users
 - Genomics
 - Machine learning
- Addition of extra nodes (2016), GPUs (2019)



Peregrine: Main challenges

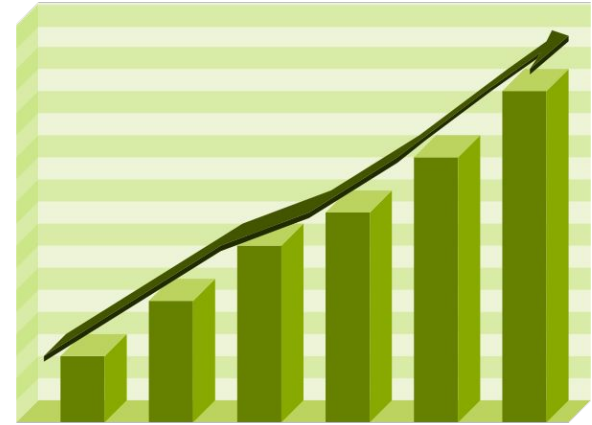
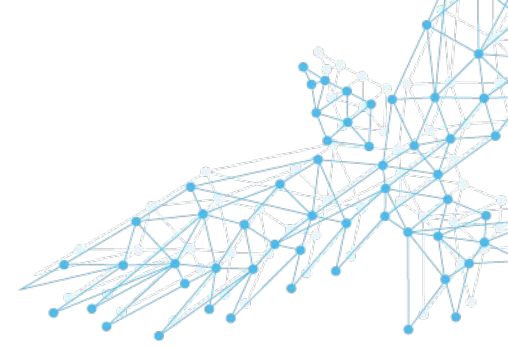


- Lack of capacity
- Storage performance
 - Huge number of files
 - Stability
- Security
 - Requests for working with sensitive data
- All hardware tied up in a single platform



Design goals Hábrók

- Increased capacity
- Improved security
- Better performing storage systems
- More flexibility

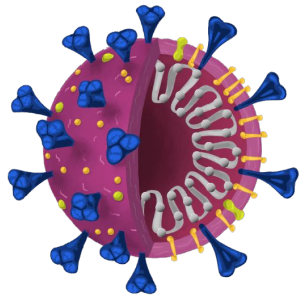


university of
 groningen

center for
 information technology

Delays & deadlines

- Covid pandemic
- Expired framework contract
- Understaffed team
- Big bang operation
- April 1st - End of contract with datacenter Peregrine



university of
 groningen

center for
 information technology

New cluster: Hábrók



The best of hawks

Yggdrasil's ash is
of all trees most excellent,
and of all ships, Skidbladnir,
of the Æsir, Odin,
and of horses, Sleipnir,
Bifröst of bridges,
and of skalds, Bragi,
Habrok of hawks,
and of dogs, Garm.

*Excerpt from Grímnismál - from
the poetic Edda - Iceland 13th
century*



university of
 groningen

center for
 information technology

New data center: Coenraad Bron

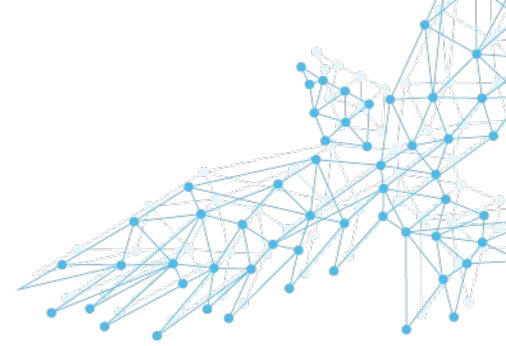


Increased capacity

- More capacity
 - Large number of CPU cores
 - 6024 -> 19184
 - 512 GB memory standard node
- GPUs
 - Newer GPUs, but only 6 of the 10 nodes
 - 4 Nvidia A100 cards per node
 - 36 existing GPU nodes
 - 1 Nvidia V100 card per node
- Big memory nodes
 - 4 x 4 TB memory
- Additional storage capacity
 - 2.7 PB of /scratch
- Redundant power



Improved security

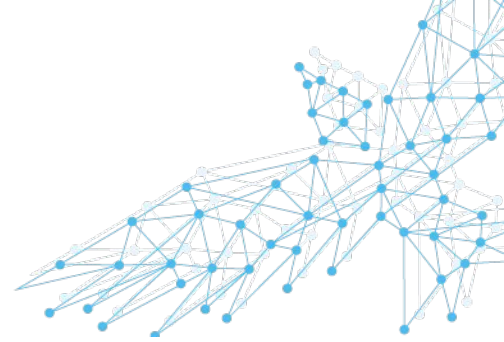


- Security
 - Tighter storage permissions
 - Private folders will stay private
 - explicit groups and folders required for sharing data
 - Improved system management
 - Fully scripted
 - DevOps
 - Multi-factor authentication
- CIT working towards ISO certification

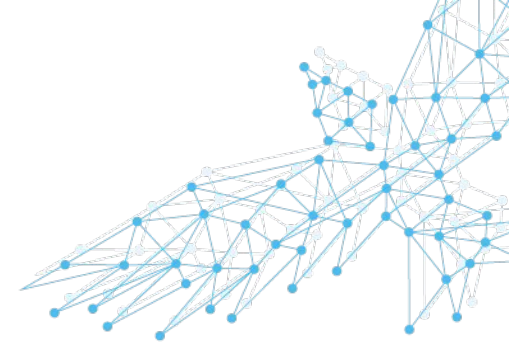


Better performing storage

- Full flash shared storage not feasible
- File system metadata on flash storage
- More capacity
 - 2.7 PB
 - Spinning disk
- Fast local disks in all nodes
 - SSD or NVMe
- Long term storage options
 - /projects & RDMS



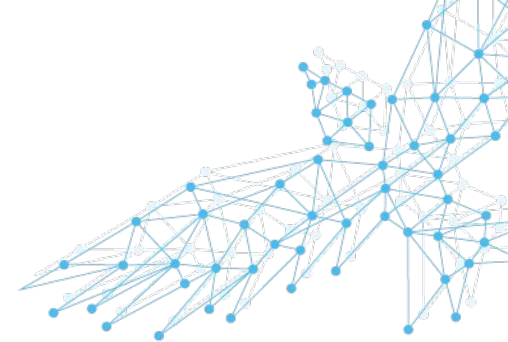
More Flexibility: Cloud backend v3.0 Bateleur



- Hábrók constructed within an on-premise cloud
- Software defined networking on switches
- Combination of bare metal and virtual machines
- Multitenancy all from the same environment:
 - Hábrók
 - Extra secure clusters
 - virtual workloads/HPC-PaaS platforms
- Infrastructure as code
- Config as code



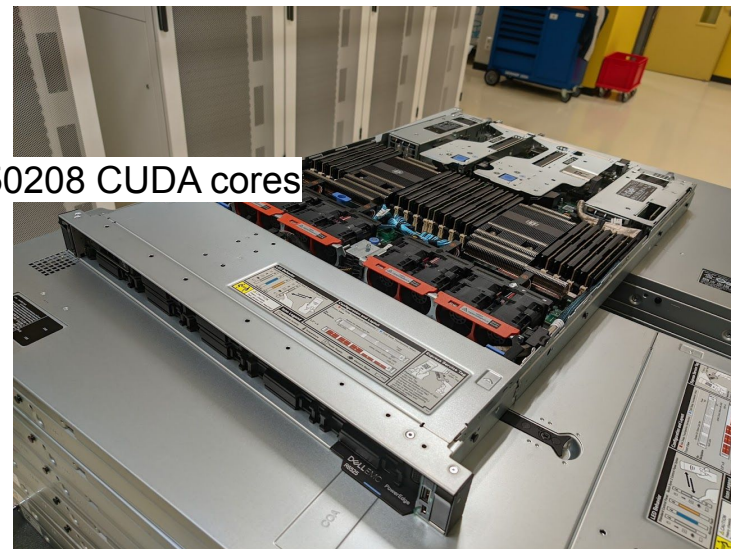
Cluster overview (raw numbers)



Type	#	Cores	Memory (GiB)	GPU	Partition	Local storage (TiB)
Regular	117	128	512	-	regular	3.5
Omnipath network	24	128	512	-	parallel	3.5
Memory	4	80	4096	-	himem	14
GPU1	6	64	512	4 x A100 (40 GiB)	gpu	12
GPU2	36	12	128	1 x V100 (32 GiB)	gpu	1
Gelifes	15	64	512	-	gelifes	15 (HDD)

19184 CPU cores

350208 CUDA cores

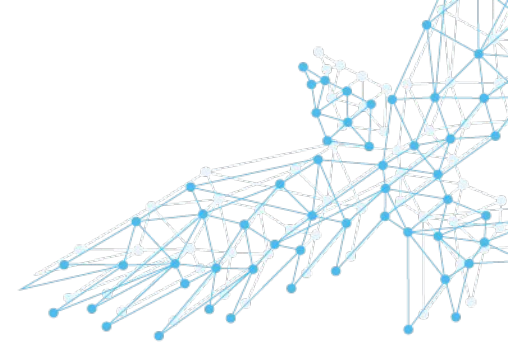


university of
 groningen

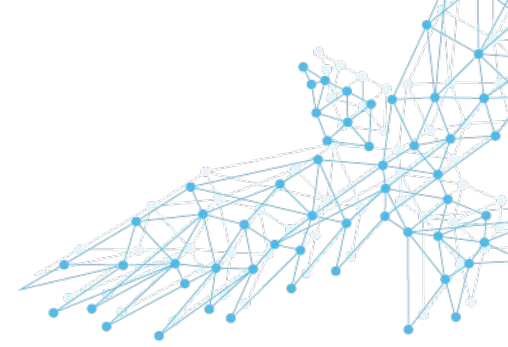
center for
 information technology

Migration to Hábrók

1. Request an account for Hábrók
 - a. Prevents migrating inactive accounts
 - b. <https://iris.service.rug.nl>
2. Migrate your data
 - a. Peregrine /data at /mnt/pg-data
 - b. Peregrine /home at /mnt/pg-home
 - c. Three months time
3. Report issues and missing software
 - a. hpc@rug.nl (leads to HPC team in IRIS)



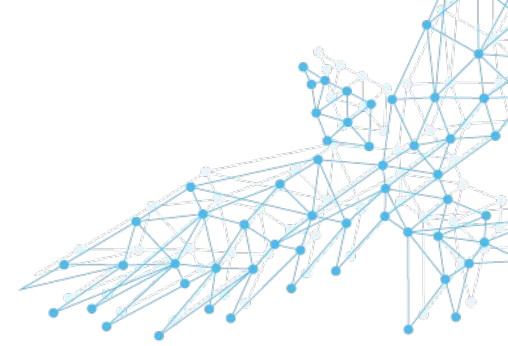
Current status



- Many accounts activated
 - Still many to go
- Groups need to be created
 - Groups will get access to their existing data on Peregrine
 - Gelifes access
 - Should be done by end of week
- Group member migration
 - Automatically if your account was migrated before creating the groups
 - Possible manual intervention afterwards



Changes: Accounts (increased security)



- RUG accounts only
 - UMCG staff needs to use p-account
 - Permissions on existing data will be adjusted
 - f-accounts are harder to obtain (RUG policy)
 - s-accounts for courses
- Multi factor authentication
 - Standard RUG policy
 - 8 hour window
 - May not work for f-accounts without email address



Changes: Logging in

Multiple login and interactive nodes:

- `login1.hb.hpc.rug.nl`
- `login2.hb.hpc.rug.nl`

- `interactive1.hb.hpc.rug.nl`
- `interactive2.hb.hpc.rug.nl`

CPU and memory limits on all nodes

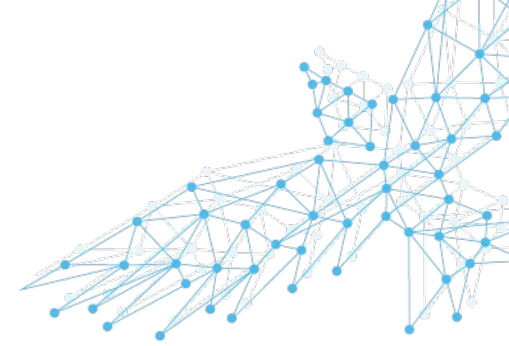
Web portal not available yet.

Login nodes with limited capacity
light usage

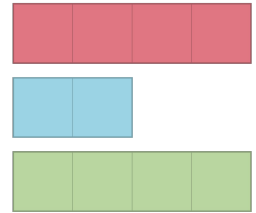
Nodes for interactive work
testing, small runs



Changes: Software stack, job environment

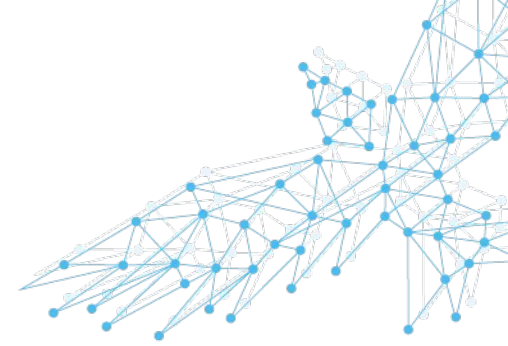


- Needed to be rebuilt
- Optimized for all architectures
 - AMD Zen3 in regular nodes
 - Intel Ice Lake in GPU and big memory nodes
- Improved distribution mechanism
 - http with proxies and local cache
- Most commonly used software is available, rest will be added on request
- Modules loaded on login nodes not transferred to job, explicit module load `mysoftware/x.y-foss-2022a` required



Changes: Storage areas

- Home directories
 - multiple NFS servers /home1 - /home4
 - Use ~ or \$HOME
- Large /scratch space
 - No data cleanup, just quota, fair use policy
- Fast local disk
 - \$TMPDIR
- /projects for medium term storage
 - Fair use policy with payment for large quota
 - Discuss backup needs
- RDMS connection



Quota and 'scratch' space

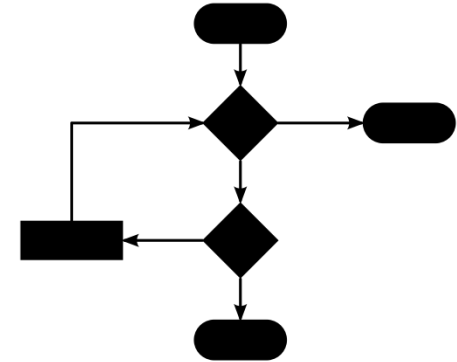
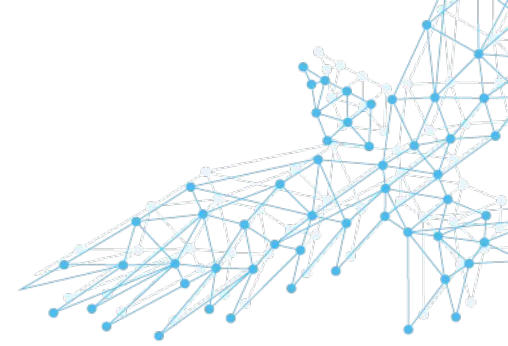


- Quota described on the wiki, a 'hbquota' tool will be created
- Quota for projects and scratch can be increased upon request
 - You will start with the default
- **scratch means what it sounds like**
- **You are responsible for making sure important data is safely stored elsewhere**
- Store important data in /projects or RDMS
 - Fair use, payment required for large amounts of data

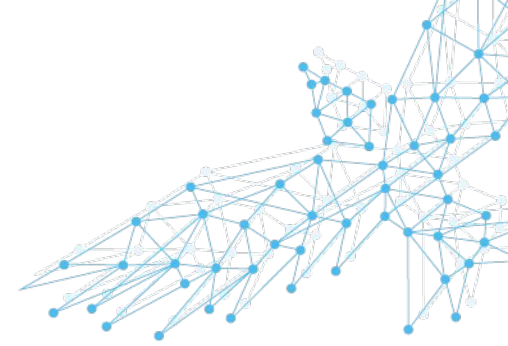


Changes: Suggested workflow - Large files

1. Prepare:
 - a. Upload input data to /projects storage
 - b. Copy working set to /scratch storage
2. Jobs:
 - a. Submit jobs, that use the data /scratch
3. Finalizing:
 - a. Copy results back to /projects
 - b. Copy results to RDMS



Changes: Suggested workflow - Many small files



1. Prepare:
 - a. Upload data to /projects storage as archive file
 - i. e.g. tar, zip, etc.
 - b. Copy working set to /scratch storage
2. Jobs:
 - a. Submit jobs
 - i. Within job extract data set from archive to \$TMPDIR on /local
 - ii. Copy results back at the end of the jobscript
 1. In archive files if many files
3. Finalizing:
 - a. Copy results back to /projects
 - b. Copy results to RDMS

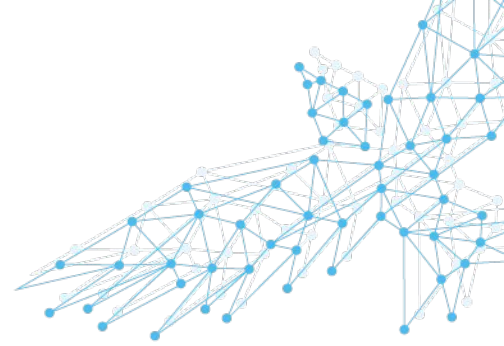


Changes: Partitions

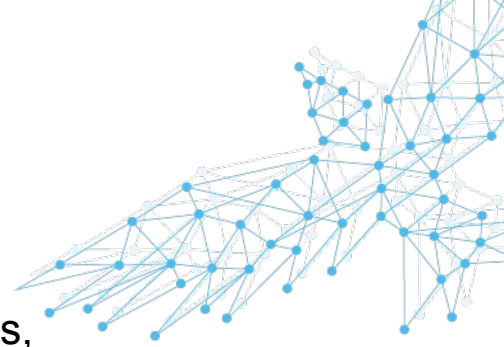
Standard partitions based on hardware classes

- regular: standard nodes
- parallel: OmniPath nodes
- gpu: GPU nodes
- himem: Big memory nodes
- gelifes: GELIFES nodes

Subpartitions -short, -medium -long to prevent whole cluster being occupied with 10 day jobs.



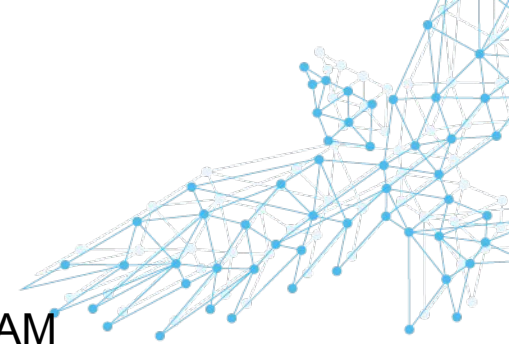
Changes: Selecting partitions



- No partition selected:
 - regular or himem based on core and memory requirements,
 - gpu when a GPU was requested
- short, medium and long based on time
- parallel needs to be specified explicitly
 - Only useful if you need the fast interconnect



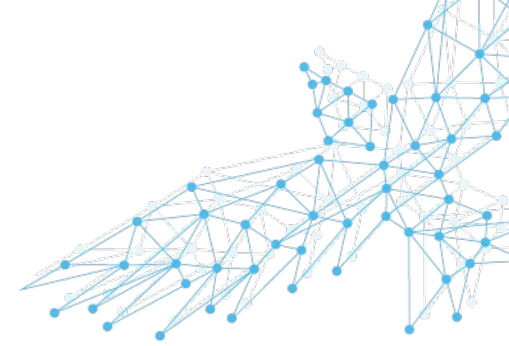
Changes: GPUs



- Multiple NVIDIA GPU types
 - 32 x V100 32 GB, 1 per node, 12 cores, 128 GB server RAM
 - Under construction
 - 24 x ½ A100 20 GB, 8 per node 20 GB, 512 GB server RAM
 - Single GPU per application
 - 12 x A100 40 GB, 4 per node, 512 GB server RAM
 - Available for multi-GPU runs
- gpu partition
 - Default 20 GB GPU
 - Specify type (a100, a100.20gb, v100) or constraint GPU_xxGB
 - Slurm flags:
#SBATCH --gpus-per-node=[type:]<number>



Known issues



- Missing software
 - Report things you need
- Tools must be ported
 - hbquota, jobinfo
- Web portal not yet available
- V100 nodes
 - Still in Peregrine
- GELIFES nodes
 - Still in Peregrine, interim capacity in Hábrók
- See: https://wiki.hpc.rug.nl/habrok/additional_information/known_issues



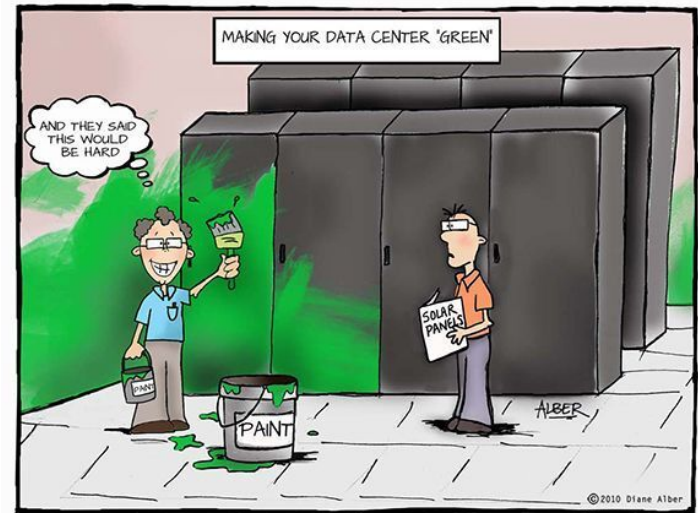
Future changes

- Monitoring energy consumption per job
- MultiXscale project for external optimized scientific software stack
- Account and group lifetime
- Security
- Handling external accounts
- Better user representation



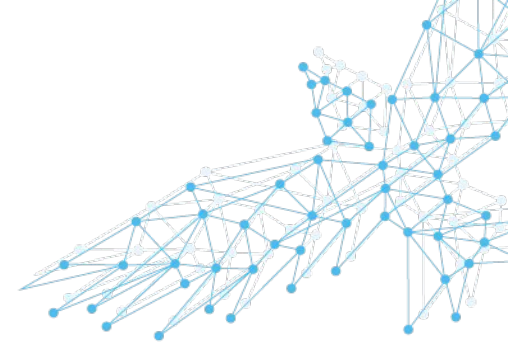
university of
 groningen

center for
 information technology



Vacancies

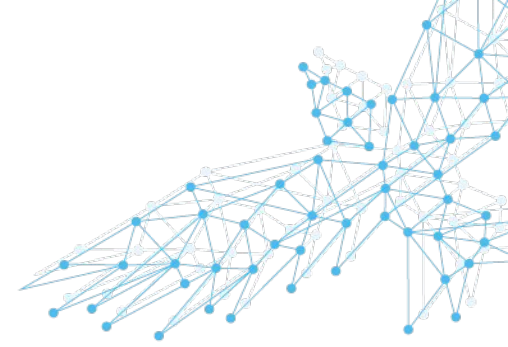
- Looking for new colleagues
 - DevOps engineer
 - Cloud & storage
 - MultiXscale project
 - Software distribution
- Student positions
 - User support
 - Depends on expertise



university of
 groningen

center for
 information technology

Questions?



university of
groningen

center for
information technology