# Galaxy

**make it run your jobs for you**

Martin Čech – Telč – 28.3.25

slides at: **ces.net/telc**

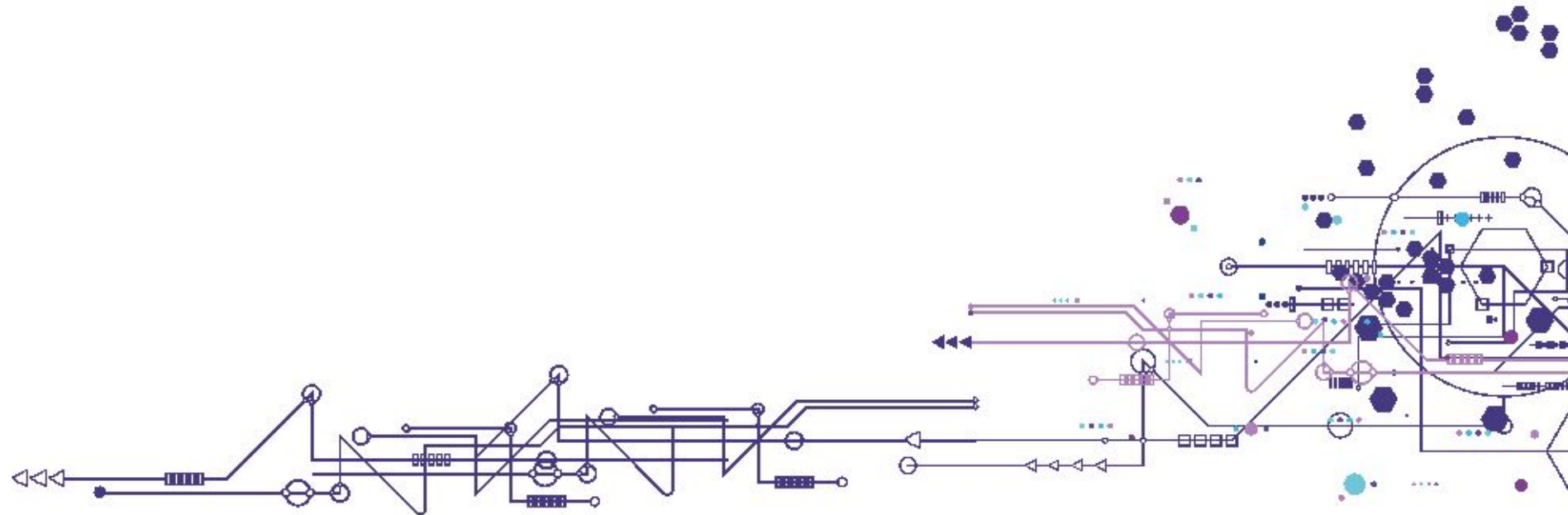cesnet    elixir **CZECH REPUBLIC**

# Outline

- Computing in Czechia with MetaCentrum
- Galaxy's purpose & capabilities
- Galaxy in Czechia

# Computing in Czechia

**with MetaCentrum**

# metacentrum

cesnet

- National Grid Infrastructure (NGI)
- Provider of computational resources and data storage
- Free (as in beer)
  - For employees and students in Czech Academia
  - But also for industry users (non-profit public research, upon individual request)
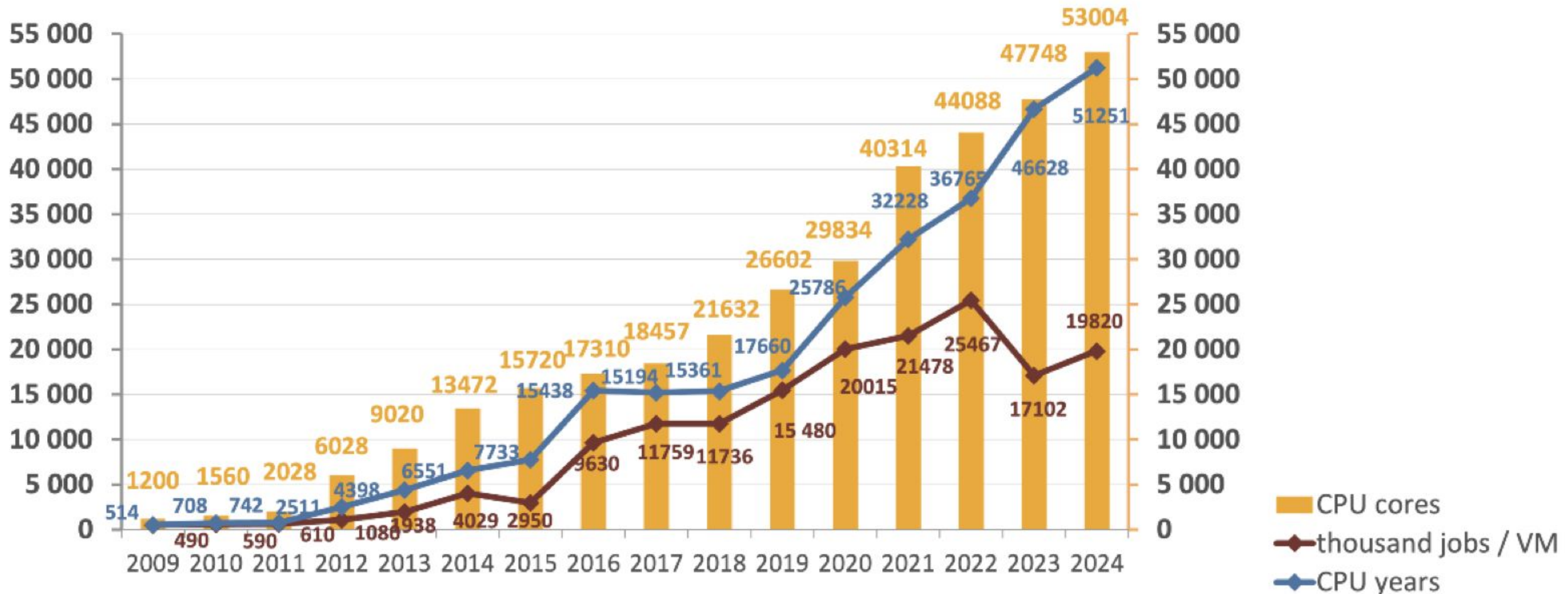
**cesnet**

**meta**centrum

Compute resources are pooled together by all partners (CESNET, universities, Czech Academy of Sciences) and…

- ... are centrally managed
- ... are shared among all users
- ... have privileged access for cluster owners
- ... are replaceable during an outage
- … include support for federated AAI
- ... are dedicated to grid HPC/HTC, containerised computing, cloud computing, data storage capacities

**cesnet stats**

Number of CPUs, executed jobs and corresponding CPU years
(PBS, cloud, K8s, EGI)

**cesnet**

**meta**centrum

MetaCentrum targets

- individual users (access to resources)
- projects (cooperation, sharing data in a group)
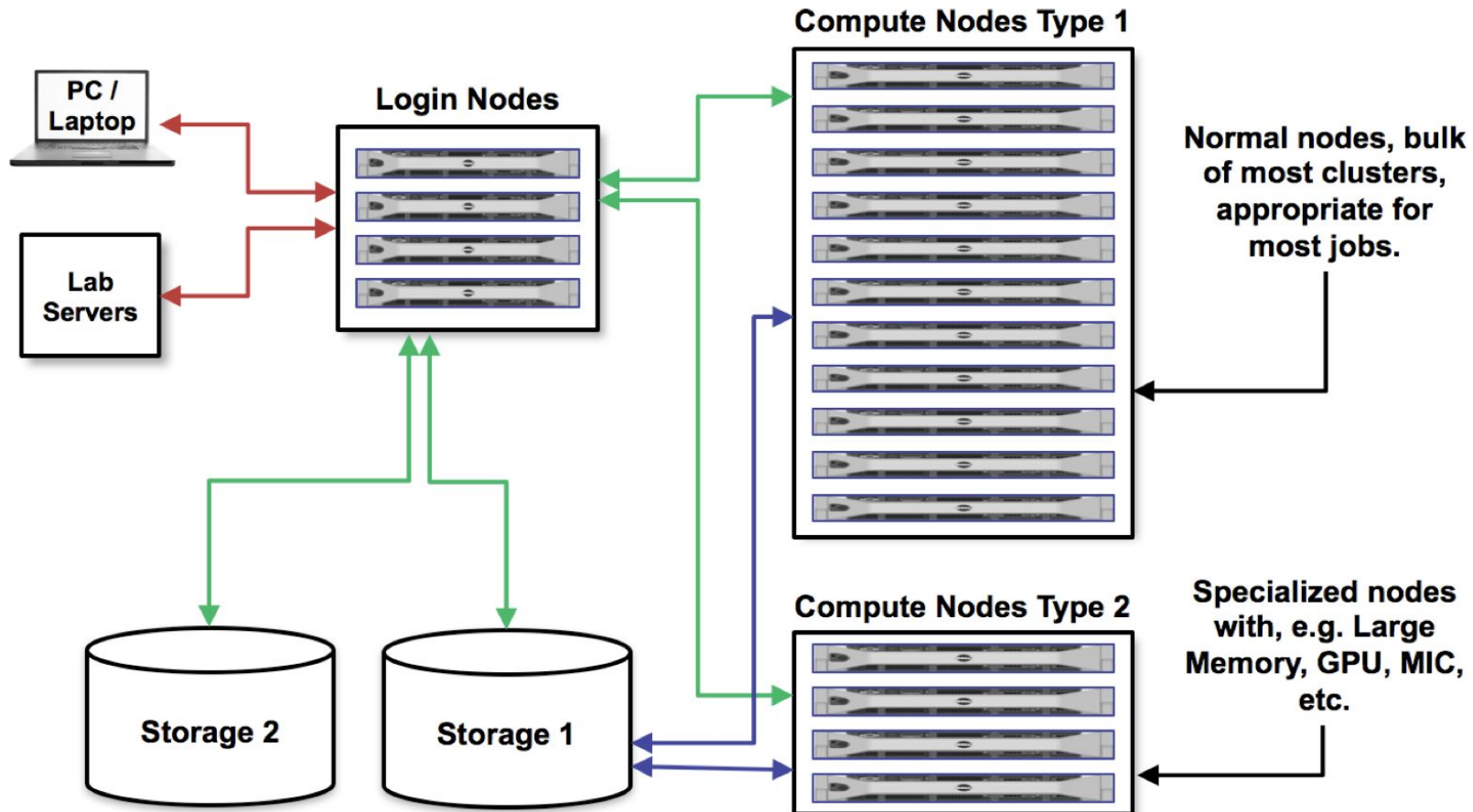- organisations (incorporate their resources under central management)

MetaCentrum offers

- immediate access to HW resources
- various application tools (commercial, free, open source)
- CPU/GPU resources, GUI applications and access, cloud services

Great starting point is: https://docs.metacentrum.cz/

# MetaCentrum batch job work setup



**PC / Laptop**

**Lab Servers**

**Login Nodes**

**Compute Nodes Type 1**

Normal nodes, bulk of most clusters, appropriate for most jobs.

**Compute Nodes Type 2**

Specialized nodes with, e.g. Large Memory, GPU, MIC, etc.

**Storage 2**

**Storage 1**

# MetaCentrum batch job work setup

**qsub Manual Page**
    NAME
        qsub - submit    job
    DESCRIPTION
        **To create a job is to**  **submit an executable script to a batch server.**  The batch server will be the default server unless the -q option is specified. Typically, the script is a shell script which will be executed by a command shell such as sh or    csh.

        Options on the qsub command allow the specification of attributes which affect the behavior of the job.

```bash
#!/bin/bash
#PBS -q default@meta-pbs.metacentrum.cz
#PBS -l walltime=24:0:0
#PBS -l select=1:ncpus=8:mem=100gb:scratch_ssd=50gb
#PBS -N my_awesome_job
#PBS -m e

# test if a scratch directory exists
# variable SCRATCHDIR is set automatically
test -n "$SCRATCHDIR" || { echo >&2 "Variable SCRATCHDIR is not set!"; exit 1;

# set a DATADIR variable
DATADIR=/storage/brno12-cerit/home/vorel/data/

# copy input file "data.fa" to the scratch directory
cp $DATADIR/data.fa $SCRATCHDIR

# move into the scratch directory
cd $SCRATCHDIR

# load a module for your application
module add blast-plus/blast-plus-2.12.0-gcc-8.3.0-ohlv7t4

# run the calculation
# do not forgeto to use reserved CPUs by '-num_threads' flag
# variable PBS_NCPUS is a number of CPUs requested for the entire job
blastp -query data.fa <other_parameters> -num_threads $PBS_NCPUS -out results.

#copy results
cp results.txt $DATADIR

# clean the scratch directory
clean_scratch
```
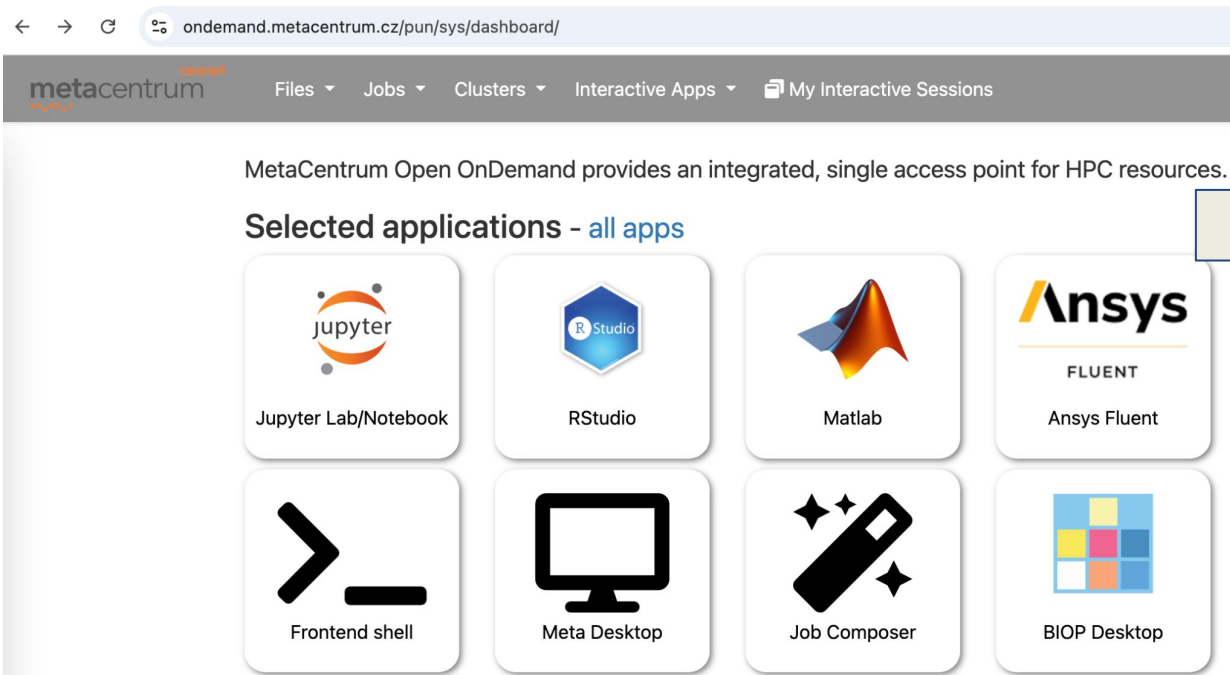
# OnDemand user experience

- I need 64 cores, 1 TB RAM, 4 hi-end GPU for my Matlab/Jupyter/RStudio/other calculation
- Wow, this is expensive…
- But I need it 6 hours per month only
- Go to https://ondemand.metacentrum.cz/
- Fill a formwith  what you need for how long
- Profit (I mean *research*)

# MetaCentrum supported pubs in 2024

**Research fields:**
- 40 Materials Science Multidisciplinary
- 32 Chemistry Multidisciplinary
- 19 Physics Applied
- 17 Astronomy Astrophysics
- 35 Chemistry Physical
- 22
- 27 Biochemistry Molecular Biology
- 17 Nanoscience Nanotechnology
- 15 Physics Atomic Molecular
- 21 Multidisciplinary Sciences

**Institutions:**
- 104 CZECH ACADEMY OF SCIENCES
- 49 MASARYK UNIVERSITY BRNO
- 22 BRNO UNIVERSITY OF TECHNOLOGY
- 19 INSTITUTE OF ORGANIC CHEMISTRY BIOCHEMISTRY OF THE CZECH ACADEMY OF SCIENCES
- 78 CHARLES UNIVERSITY PRAGUE
- 32 TECHNICAL UNIVERSITY OF OSTRAVA
- 19 INSTITUTE OF PHYSICS OF THE CZECH ACADEMY OF SCIENCES
- 18 CENTRE NATIONAL DE LA RECHERCHE SCIENTIFIQUE CNRS
- 28 CZECH TECHNICAL UNIVERSITY PRAGUE
- 16 PALACKY UNIVERSITY OLOMOUC

# IT4Innovations – supercomputers at your service

- Three times a year an open access grant competition
- But also Fast Track Access for smaller projects (4 months)

Computational resources distributed every 4 months (node hours):

- Barbora CPU: 500,000 n/h, GPU: 30,000 n/h, FAT: 3,400 n/h

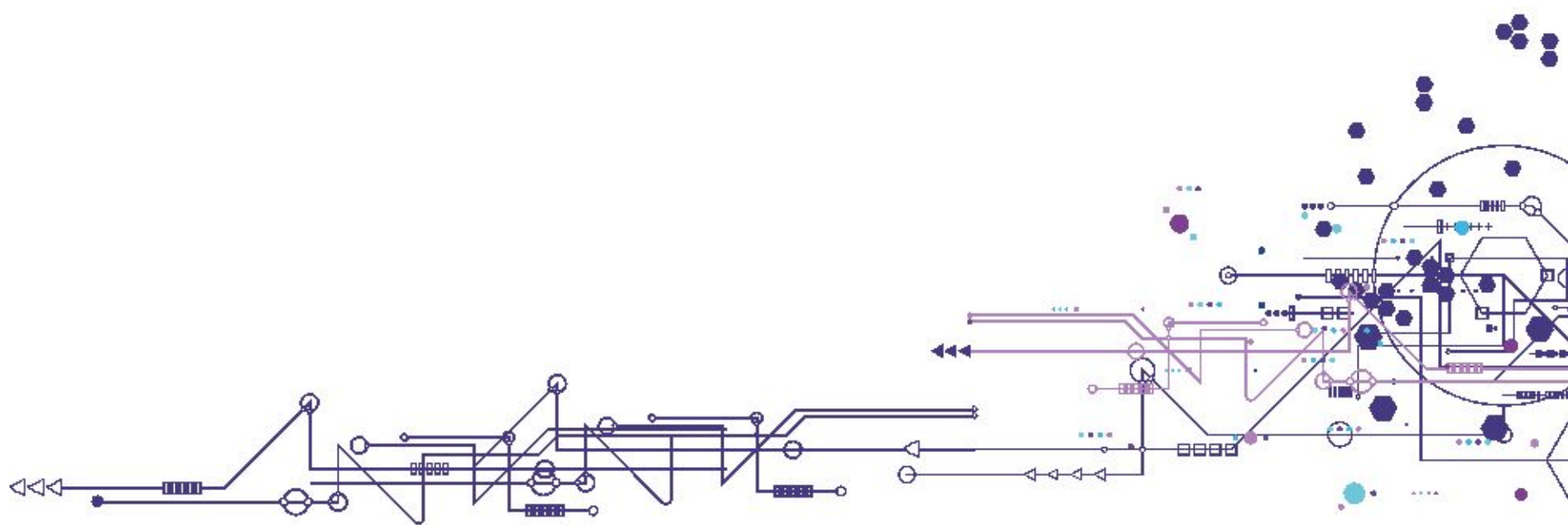- Karolina CPU: 1,012,000 n/h Karolina GPU: 83,000 n/h, FAT: 1,200 n/h
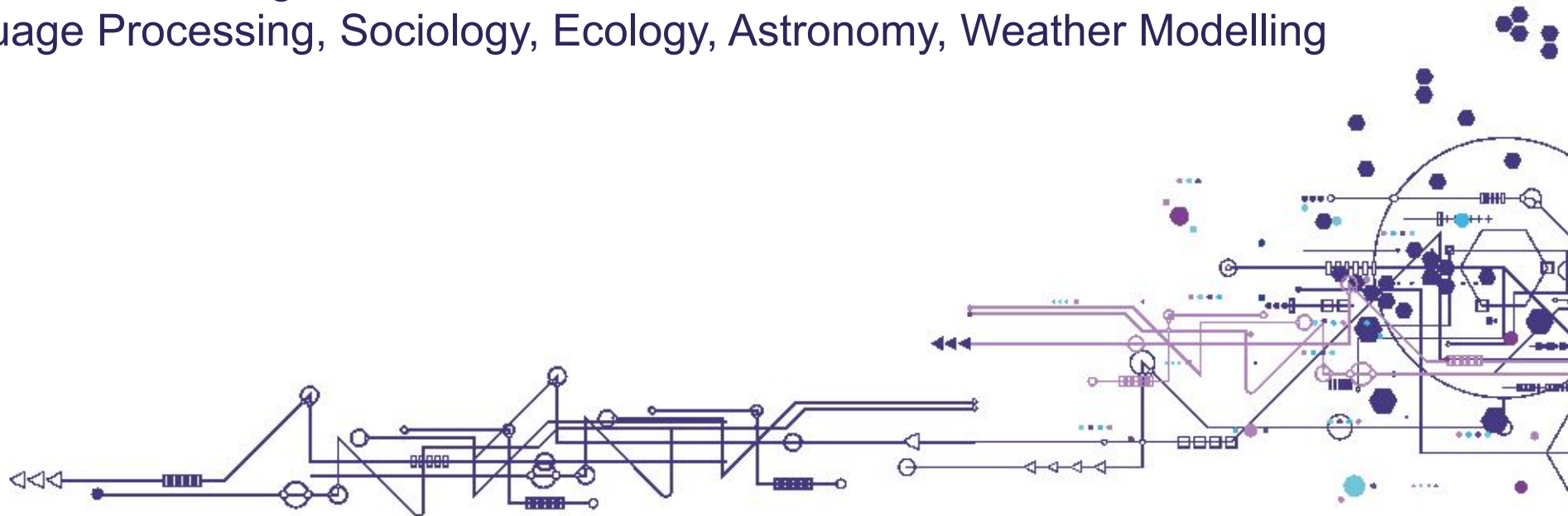
details: https://www.it4i.cz/en/for-users/open-access-competition

# Galaxy's purpose & capabilities

## simplified interface to software and compute infrastructure

# Galaxy - key purpose

- Analyze data using thousands of tools without installation and maintenance
- Learn or publish methods and techniques
- Create and share workflows, using an expressive graphical interface
- Distribute one's own tools, locally or globally
- Expose compute and storage infrastructure for easy access by users
- Born from biology but domain-agnostic
  - Natural Language Processing, Sociology, Ecology, Astronomy, Weather Modelling

# Capability - support scale

**Lex Nederbragt** @lexnederbragt · Apr 26, 2017

"Please learn command-line to be able to use docker to be able to make installing web-based tool to replace command-line (Galaxy) easier"

> **Björn Grüning** @bjoerngruening · Apr 26, 2017
>
> Replying to @lexnederbragt
>
> This helps many people with installation: github.com/bgruening/dock...
> - you still need competence to handle the data.

💬 5          ⟲ 7          ♡ 11          ⬆

**Devon Ryan** @dpryan79 · Apr 26, 2017

Those of us doing this aren't doing it for us, we're doing it for the bench scientists whose data we then needn't personally analyze.

💬 1          ⟲ 2          ♡ 4          ⬆

**Devon Ryan**
@dpryan79
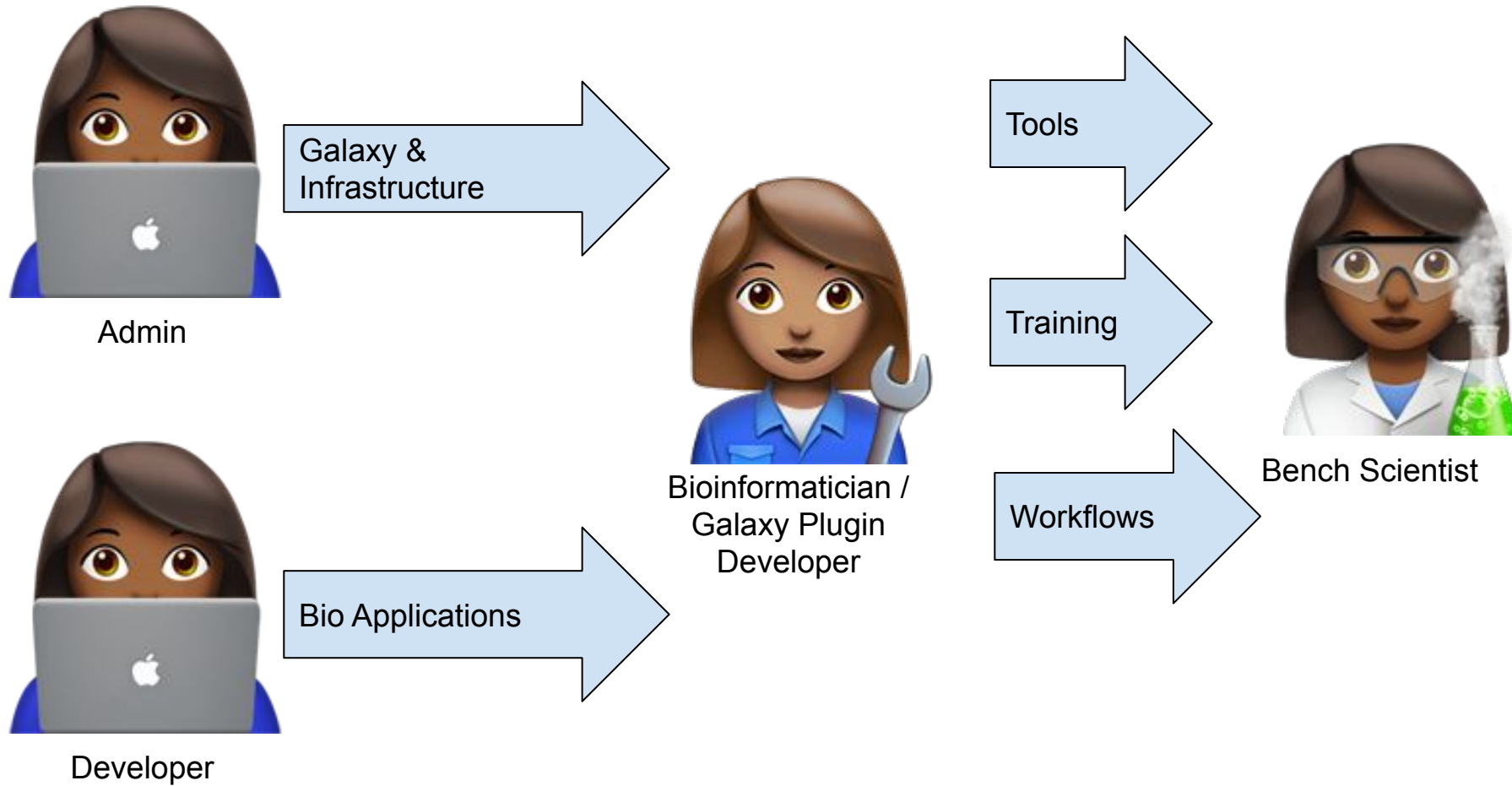
Replying to @dpryan79 and @lexnederbragt

Galaxy scales, my time does not.

3:45 PM · Apr 26, 2017 · TweetDeck

**4** Retweets    **8** Likes

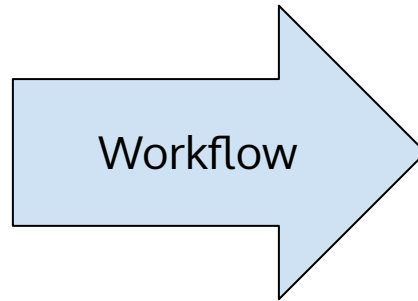💬          ⟲          ♥          ⬆

# Capability - support specializations

# Capability - hide underlying complexity

Methods Developer

Workflow

Pipeline Executor

Galaxy UI enables method developers without knowledge of scripting, etc...

No need to understand command-lines, etc.. Embedded visualizations, etc..

# How does it look like

# A little Galaxy demo

**if wifi holds**

# Galaxy Show case: VGP



VGP plan:

- Publish genome assemblies of 12 species every week and accelerate
- Until all the 75,923 extant vertebrate species genomes are known
- which is expected within a decade
- =~ petabytes of assembly data (ignoring raw data)
- *Reference point: Assembly of human genome took 10+ years and $3 billion dollars*

https://galaxyproject.org/projects/vgp

# Galaxy Show case: VGP



A PROJECT OF THE G10K CONSORTIUM

https://galaxyproject.org/projects/vgp

Galaxy provides:

- Integration of the Genome Ark on public Galaxy servers.
- A Galaxy platform with toolkits specifically tailored for Genome assembly
- Workflows available using the most up-to-date VGP pipelines.
- A list of publicly-available histories for each assembly completed on Galaxy as they are generated.
- Extensive training materials used to spread the load and accelerate

Using Galaxy helps make the assemblies reproducible and enables scaling

# Vertebrate genome assembly using HiFi, Bionano and Hi-C data - Step by Step

Authors: Delphine Lariviere · Alex Ostrovsky · Cristóbal Gallardo · Anna Syme · Linelle Abueg · Brandon Pickett · Giulio Formenti · Marcella Sozzoni

## Overview

**Questions:**
- What combination of tools can produce the highest quality assembly of vertebrate genomes?
- How can we evaluate the quality of the assembly in a reference-free way?

**Objectives:**
- Learn the tools necessary to perform a de novo assembly of a vertebrate genome
- Evaluate the quality of the assembly

**Requirements:**
- Introduction to Galaxy Analyses
- 🎞 Slides: Quality Control
- 🖥 Hands-on: Quality Control

⏳ **Time estimation:** 5 hours

🎓 **Level:** Intermediate 🎓🎓🎓

🎒 **Supporting Materials:**
  - 📄 Datasets   ⚙ Workflows   ❓ FAQs   🎬 Recordings ▾

📅 **Published:** Jun 4, 2021

📅 **Last modification:** Sep 27, 2024

⚖ **License:** Tutorial Content is licensed under Creative Commons Attribution 4.0 International License. The GTN

🔗 **PURL:** https://gxy.io/GTN:T00039

⭐ **Rating:** 5.0 (1 recent ratings, 3 all time)

⌀ **Revision:** 54

bit.ly/vgp-training

# Workflows Overview

Eight analysis trajectories are possible depending on the combination of input data. Decision on invocation of workflow 6 is based on the analysis of QC output of workflows 3, 4, or 5 (see below). Thicker lines connecting workflows 7, 8, and 9 represent the fact that these workflows are invoked separately for each phased assembly (once for maternal [or hap1] and once for paternal [or hap2]). **Solo** = data is only available for the sample whose genome is being assembled. In this case, you can make either a pseudohaplotype assembly, or a HiC-phased assembly if you have HiC data from the same individual.

**Trio** = parental information is available in the form of Illumina reads from each parent of the F1 being assembled.

# A peak into VGP workflows

**if wifi holds**

# Galaxy Show case: SARS-CoV-2

**galaxyproject/ SARS-CoV-2**

Four Goals:

- Continuously analysis of within-host sequence variants in high quality public read-level datasets
- Maintenance of curated workflows for the analysis of SARS-CoV-2 sequence data
- Development of continuously updated analysis page and dashboard summarizing latest insights from the variant.
- Providing access to all results in raw and aggregated form for immediate use.

# Galaxy Show case: SARS-CoV-2

https://observablehq.com/@spond/intrahost-dashboard

# Galaxy Show case: SARS-CoV-2

Everything implemented in unprecedented levels of open science

Repository with bot processing requests:
https://github.com/usegalaxy-eu/sars-cov-2-processing-requests

This project eventually analyzed nearly 500,000 public SARS-CoV-2 sequencing datasets until 2022.

# Galaxy in Czechia

**everybody can run Galaxy**

# Galaxies we run

- RepeatExplorer Galaxy
  - specialized instance for graph-based clustering and characterization of repetitive sequences in NGS
  - run Institute of Plant Molecular Biology in Ceske Budejovice
  - one of **ELIXIR CZ** services
- UMSA Galaxy
  - Untargeted Mass Spectrometry Analysis
  - operated on behalf of **MUNI RECETOX**
  - publish tools, workflows, and other research results
  - provide pipeline as a service for collaborators
- UseGalaxy.cz
  - all purpose instance (and we'll support your tools)
  - aims to make the resources of National Grid Infrastructure more accessible

# Galaxy Training Academy 2025

12-16 May 2025 (online)

- 5-day Global Online and Asynchronous learning event
- choose your own path from 400 trainings
- slides, hands-on, videos, screencasts
- daily sync times, continuous support on slack
- years of experience running these, they're really good
- everything lives on https://training.galaxyproject.org

| | | | |
|---|---|---|---|
| Proteomics | Assembly | Transcriptomics | Single Cell |
| Microbiome | Machine Learning | From Zero to Hero with Python | |

# Coming up

Events:
- 22 April 2025        **Galaxy Imaging Hackathon 2025** (Freiburg, GER + online)
- 12-16 May 2025       **Galaxy Training Academy 2025** (online)
- 27-29 May 2025       **12th Repeat Explorer Workshop** (Ceske Budejovice, CZ)
- 23-27 June 2025      **Galaxy and Bioconductor Community Conference** (NY, USA)
- 1-3 October 2025     **European Galaxy Days** (Freiburg, GER)

Projects:
- National Repository Platform for Research Data 2024-2028
  - imagine it as a national Zenodo instance per research area
  - Galaxy will provide integration with the repositories
- EuroScienceGateway
  - Pooling of research compute resources between EU countries with Galaxy Project.

# Take home message

- after PhD every project you work on will be bigger
- your time does not scale
- Excel, Matlab, and homegrown approaches will slow you down


- MetaCentrum gives you **the iron**
- Galaxy gives you **the utility**
- Galaxy Training Network gives you **the mastery**


- https://docs.metacentrum.cz
- https://usegalaxy.cz


p.s. If you take only one thing from here be it training.galaxyproject.org

# Credits and thank yous

- Delphine Lariviere – Galaxy, VGP
- Wolfgang Maier – Freiburg, SARS-CoV-2
- Nate Coraor – Galaxy
- John Chilton – Galaxy
- Jirka Vorel – CESNET
- Aleš Křenek – CESNET/MUNI

**& The Galaxy Community**

# Thank You

**Questions, please?**

Martin Čech – Telč – 28.3.25

@martenson

slides at: **ces.net/telc**

# Galaxy Project CVMFS network

- Stratum 0 servers
- Stratum 1 servers

cvmfs-stratum0.galaxyproject.eu
cvmfs1-ufr0.galaxyproject.eu

**ELIXIR, de.NBI, RZ Freiburg**

**EU JRC, Ispra**
galaxy.jrc.ec.europa.eu

**Digital Research Alliance of Canada**
cvmfs-s1-galaxy.computecanada.ca

cvmfs0-psu0.galaxyproject.org
cvmfs0-psu1.galaxyproject.org
cvmfs1-psu0.galaxyproject.org

**Penn State**

**ACCESS/Jetstream2, Indiana University**
cvmfs1-iu0.galaxyproject.org

cvmfs1-melb0.gvl.org.au

**Melbourne Bioinformatics**

**CyVerse, TACC**
cvmfs0-tacc0.galaxyproject.org
cvmfs1-tacc0.galaxyproject.org

**e-INFRA CZ**

Legend:
- **Supercomputer**
- **Grid HPC + HTC + Kubernetes**
- **HPC Cloud**
- **Object DataStorage**
- **Resources available in EGI**
- **Resources dedicated to ELIXIR CZ**
- **Virtalization platform**
- **NRP - plan 2025-2027**

Liberec — HPC Clusters PBS Pro

Ústí nad Labem — Repository NRP

Praha:
- Repository NRP
- Virtualization platform
- ELIXIR HTP Cluster PBS
- HPC Clusters HPC CUDA Clusters PBS
- HTC cluster Condor, DPM

Vestec — ELIXIR HPC Cluster PBS, GPFS

Plzeň:
- HPC Clusters PBS, GPFS
- Object storage CEPH

Jihlava — Object storage CEPH

Olomouc:
- ELIXIR HPC Cluster PBS, GPFS
- Repository NRP
- HPC Clusters PBS

Ostrava:
- Object storage CEPH OwnCloud
- Supercomputer IT4Innovation
- IT4I Cloud OpenStack platform
- Repository NRP

Brno:
- MetaCentrum Cloud EGI FedCloud + ELIXIR OpenStack platform Object storage CEPH
- HPC Clusters (CESNET/MU/MENDELU) HPC CUDA Clusters Kubernetes PBS, GPFS
- Virtualization platform
- Repository NRP
- Object storage CEPH
- ELIXIR HPC Cluster PBS, GPFS Galaxy SensitiveCloud

České Budějovice — HPC Clusters PBS

PIONIER

AMS-IX

NIX

GÉANT

ACONET

SANET

# GTN

Galaxy Training Academy this week (~3k ppl)

Per public stats (training.galaxyproject.org/training-material/st____)

419 trainings

400 contributors

55k unique visitors per month

Couples with Zenodo for input data management

## 28 Scientific Topics

**Tutorials per Topic**

| Topic | Count |
|---|---|
| Assembly | 19 |
| Climate | 12 |
| Computational chemistry | 9 |
| Ecology | 22 |
| Epigenetics | 10 |
| Evolution | 9 |
| FAIR Data, Workflows, and Research | 21 |
| Foundations of Data Science | |
| GMOD | 15 |
| Genome Annotation | 20 |
| Image analysis using Deep Learning | 3 |
| Imaging | 7 |
| Introduction to Galaxy Analyses | 13 |
| Materials Science | 1 |
| Metabolomics | 9 |
| Microbiome | 18 |
| One Health | 9 |
| Plants | 8 |
| Proteomics | 32 |
| SARS-CoV-2 | 9 |
| Sequence analysis | 8 |
| Single Cell | 35 |
| Statistics and machine learning | 18 |
| Synthetic Biology | 3 |
| Transcriptomics | 24 |
| Using Galaxy and Managing your Data | 23 |
| Variant Analysis | 17 |
| Visualisation | 5 |