# Data Analysis

Speakers: **Siri Kallhovd** (UiB),
**Espen Robertsen** (UiT), and
**Erik Hjerde** (UiT)
**Moderator: Alexandra Gade** (NCMM, UiO)

**15 June 2022**

**Siri Kallhovd** (UiB)

# Computing Resources

✓ **Learning Objectives**

In this session, we will learn:

- what High-Performance Computing (HPC) is
- how to get access to HPC systems provided by Sigma2 for Norwegian research
- how to find other local computing resources available at your university

**⊓ NRIS**

# HPC Systems

Sigma2 owns the HPC Systems Saga, Fram and Betzy and NRIS operates the systems.

NRIS is a collaboration of the four BOTT universities and Sigma2 to pool competencies, resources and services.

Apply for resources:  https://www.sigma2.no/high-performance-computing

# Saga

System: HPE Apollo 2000

Number of cores: 16064

Number of nodes: 364 + 8 bigmem + 8 GPU nodes

Max floating point performance, double: 0.65 Petaflop/s

Total memory: 97.5 TiB

Total disk capacity: 5.3 PB (++)

Network: FDR IB (56 Gbit)

File System: BeeGFS, HPE 4150

# Software modules

Need for an isolated environment on a shared HPC resource

(Scientific) software should not be installed in the global PATH

Need some way of bringing software in and out of your environment

Many different solutions to this:

Lmod (modules) + EasyBuild (installation)

Virtualenv + pip

Conda

Singularity containers

We use Lmod + EasyBuild for the scientific software stack on all our HPC machines (Fram, Saga and Betzy)

Users can use any of the above to manage their own software stack on our machines

# EasyBuild

EasyBuild can be used as wrapper around different ways of installing software

     Targeting HPC systems

     Automation: less prone for human errors

     Reproducibility: rebuild software stack, (usually) portable to other systems

     Performance: (usually) build from source for particular hardware

     Very pedantic with versioning and dependencies

EasyBuild is tightly connected to the module system and automatically generates module files after installation

# Slurm jobs optimize the usage of HPC systems

Job scripts specify:

memory: Asking for too much can mean that you block idle CPUs and get charged for them

cores: Asking for too few can lead to underused nodes or longer run time

time: Asking for too little cuts the job

# General local computing resources

NREC (Norwegian Research and Education Cloud, nrec.no)

    -operated by UiB and UiO, the platform is built on OpenStack, users can create virtual machines

Fox: local HPC machine at UiO

Idun: local HPC machine at NTNU

# **Galaxy** **Introduction to Galaxy**

## ✓ **Learning Objectives**

At the end of this session, you will be able to:

- start exploring the Norwegian Galaxy instance: <u>usegalaxy.no</u>
- test 2000 analysis tools
- disseminate information to other users about this service

# Galaxy platform - browser accessible workbench for scientific computing



https://galaxyproject.org/blog/2017-10-public-galaxy-dashboard/

# usegalaxy.* – Federation of free public Galaxy servers

- Common core set of tools and reference genomes
- Open to anyone to use
- Backed by significant computational resources
- Support from multiple national infrastructure providers
- Galaxy Training Network



https://academic.oup.com/nar/advance-article/doi/10.1093/nar/gkac247/6572001

# usegalaxy.no – The Norwegian national Galaxy server

- Provide ~2000 bioinformatic tools and workflows
- Open to all Norwegian user and collaborators
- Directly connected to the NeLS storage
- Backed by significant computational resources
- End-user training and support through ELIXIR Norway support desk

# ELIXIR Norway - analysis connected with storage

# usegalaxy.no – The basics

- Important features:
- Tool menu with ~2000 tools sorted in sections
- Current disk usage (default is 200 GB total personal disk space)
- Server alerts
- Quick start guide
- Contact support
- Q&A forum

Main menu



Tool menu

Main window

History

# Your account and saved data

- Important features:
- Tool menu with ~2000 tools sorted in sections
- Current disk usage (default is 200 GB total personal disk space)
- Server alerts
- Quick start guide
- Contact support
- Q&A forum

# Shared data

- Data shared by other users or ELIXIR-NO with all users of usegalaxy.no
- E.g. workflows and complete histories
- You can import shared data to you user
- Instructions how to use ELIXIR-NO supported workflows are also here

# Shared data: workflows

- List of all workflows that are shared with all usegalaxy.no users
- You can import shared data to your user
- By selecting any workflow you can run data analysis, import into your user or save it on another computer

# Workflows

- Your workflows. These are the imported or the workflows you have made
- You can create new workflows here

# Import data from your computer or from the web

- Drag and drop, or browse and select file
- Alternatively, paste url for data available on the web
- Specify data type if you know (e.g. Fasta)
- The two imported files will appear as two datasets in your history

# Import data from NeLS

- Import data from Personal or Project folders in NeLS
- Redirect to the NeLS portal (require login)
- Files are selectable
- Imported data from NeLS will appear in your history
- Note: the yellow colour of files as this jobs are being processed (green = job complete)

# Use galaxy histories to organise data

- The current history is "your current work space"
- The history panel displays datasets in the order in which they were created
- You can make as many histories as you want and switch between them
- Typically, you can have one history for each project or analysis

# Galaxy tools

- Tools are available from the Tool menu
- Organised under sub-menus
- Possible to browse and search by name
- You can make your own list of favourite tools

# Galaxy tools = command line tools

- Command line tools are wrapped into Galaxy so they become accessible with a GUI
- Some Galaxy tools may have reduced the number of optional parameter settings for the tool
- Example here is the assembly tool called SPAdes

# Galaxy tools = command line tools

# Tool output

- It is possible to preview the output (result), view it in the main window or download the dataset
- You can also copy the dataset over to another history

# Galaxy workflows

- A workflow in Galaxy is basically a string of tools, where the output from one tool becomes the input for the next

| Tool 1 |
|--------|
| Output |

| Tool 2 |
|--------|
| Output |

| Tool 3 |
|--------|
| Output |

# Galaxy workflows

- The "nodules" indicate which output file from one acts as input for the next tool
- Each workflow has a name and version
- Additional text that describe the workflow and tags can be added

# Galaxy workflows

- New tools can be added by clicking on the tool in the Tool menu
- The tool will appear in the workflow editor
- Tool parameter settings can be pre-set or made up to the user to set when running the workflow

# Galaxy workflows

- The output from another tools can be connected as an input for the new tool

# Galaxy workflows

- Select input files
- Tools in the workflow will run successively
- You can view and download the result files

# Reproducibility & Transparency

## Dataset Information

| | |
|---|---|
| Number | 37 |
| Name | DESeq on data 33 and data 32: PCA plot |
| Created | Friday Dec 4th 8:16:43 2020 UTC |
| Filesize | **5.8** KB |
| Dbkey | ? |
| Format | png |
| File contents | contents |
| History Content API ID | 7ba0a550406c9f38 (9272) |
| History API ID | 39520843fe032123 (438) |
| UUID | a22dce33-aaf2-4b0b-be3f-8b5241195694 |
| Full Path | /data/part0/008/dataset_8860.dat |

## Dataset Information

| | |
|---|---|
| Number | 37 |
| Name | DESeq on data 33 and data 32: PCA plot |
| Created | Friday Dec 4th 8:16:43 2020 UTC |
| Filesize | **5.8** KB |
| Dbkey | ? |
| Format | png |
| File contents | contents |
| History Content API ID | 7ba0a550406c9f38 (9272) |
| History API ID | 39520843fe032123 (438) |
| UUID | a22dce33-aaf2-4b0b-be3f-8b5241195694 |
| Full Path | /data/part0/008/dataset_8860.dat |

## Job Information

| | |
|---|---|
| Galaxy Tool ID: | galaxy-ntnu.bioinfo.no/toolshed_/uit_deseq2_wrapper/0.0.1 |
| Command Line | Rscript /srv/galaxy/var/shed_tools/galaxy-ntnu.bioinfo.no/too... |
| Tool Standard Output | Loading required package: S4Vectors Loading required package:... |
| Tool Standard Error | *empty* |
| Tool Exit Code: | 0 |
| Job API ID: | 2ade379cb271bcc4 (5499) |

# Why use usegalaxy.no?

- Easy to use
- Connected directly with storage (NeLS)
- Backed by significant computational resources
- Training and support

- Transparent research
- Reproducible research
- Enable easy sharing of data and tools

### Tool categories executed by usegalaxy users

Text manipulation (23.6%)
Mapping (5.0%)
Computational chemistry (0.7%)
Data transfer (15.6%)
Aggregation (1.6%)
QC (3.2%)
Phylogenetics (0.1%)
Sequence analysis (9.6%)
Variant analysis (18.9%)
Reformatting (1.1%)
Transcriptomics (7.8%)
Galaxy operations (2.2%)
Assembly (0.3%)
Not categorized (10.0%)

*Nucleic Acids Res*, gkac247, https://doi.org/10.1093/nar/gkac247

# Quiz questions

Introduction to Galaxy:

1.  usegalaxy.no is a web platform for:
    a.  Storing life science data
    b.  Analysing life science data
    c.  Developing bioinformatic tools
    d.  Generating data management plans

# Useful Resources

- [usegalaxy.no](http://usegalaxy.no)
- [usegalaxy.eu](http://usegalaxy.eu)

**Espen Robertsen** (UiT)

# Workflow Management Systems

## ✓ Learning Objectives

At the end of this session, you will be able to:

- Conveniently use tonnes of ready-made Nextflow pipelines to analyse your data in a FAIR manner
- Utilize auxiliary resources such as Biocontainers, Workflowhub and nf-core to find containers and workflows enabling FAIR analyses

# What are workflow management systems (WfMS)?

- "Standards for describing computational data-analysis workflows"

- Streamline routine processes for optimal efficiency

- Handles input / output / execution of tools in a DAG-like manner

- Comes in various colors and flavours

# What are workflow management systems (WfMS)?

- Directed acyclic graphs (DAG) can be used to describe workflows visually



("Autoget" nextflow backend processing pipeline @ UiT)

# WfMS' enhances the principles of "FAIR"

- WfMS's provide an analysis description in their own right

- Enables portability and reproducibility in different computational environments

- Versioning and container integration ensures uniform data analysis

# Most WfMS integrates with containers, which makes Biocontainers.pro very useful!

- 10.6K tools, 45.5K versions, 225.3K containers and packages
- Dockerfile recipes / Conda recipes to automatically build containers in BioContainers.
- Biocontainers Registry is a hosted registry of all BioContainers images that are ready to be used



**BioContainers Flow**

SOFTWARE → CONTAINER → WORKFLOW

**BIOCONDA**®

# Workflowhub covers you workflow needs and then some!

# Nextflow.io

- Supports Bash, Python, R ...

- Dataflow handled by Groovy (Java super-set, "python for Java")

- Reproducibility

- Lots of "premium" features for free

- https://www.nextflow.io/docs/latest/

```
nextflow.enable.dsl=2


process sayHello {
    input:
      val cheers
    output:
      stdout

    """
    echo $cheers
    """
}

workflow {
    channel.of('Ciao','Hello','Hola') | sayHello | view
}
```

# Some examples of Nextflow usage in ELIXIR



- Mar-databases backend processing
  - Results integrated with and available at portal pages
  - https://gitlab.com/uit-sfb/autoget-nextflow
- FHI-pipeline utilizing TSD
  - Desensitization pipeline for Norwegian covid19-samples affected by GDPR
  - https://gitlab.com/uit-sfb/fhi-desensitize
- Metapipe
  - Metagenomic analysis pipeline developed at UiT (ELIXIR)
  - https://gitlab.com/uit-sfb/metapipe
  - https://mmp2.sfb.uit.no/metapipe/

**nf-core is a great resource!**

- Out-of-the-box high quality, production-ready, curated analysis

  pipelines built using Nextflow. (66 as of Jun, 2022)

- Maintained and validated releases ensure reproducibility.

- nf-core-tools for convenience (list, launch, create templates etc.)

- Excellent documentation!

# nf-core is a great resource!

**Nextflow tower**

- "Centralized command post"; monitoring, logging & observability of distributed workflows conveniently available in your browser

- Simplifies the deployment of pipelines on any cloud, cluster or laptop.

- Quite expensive, but some essential features are free

# Enough talk! Let's play around with nextflow!

If we have time :)

- Quick tour / usage of nf-tower

- Launch nextflow pipeline -with-tower

https://gitlab.com/uit-sfb/autoget-nextflow

```
# Install nextflow
curl -s https://get.nextflow.io | bash
mv nextflow ~/bin/

# Launch the RNAseq pipeline
nextflow run nf-core/rnaseq \
    --input samplesheet.csv \
    --genome GRCh37 \
    -profile docker

# Install nf-core tools
pip install nf-core

# List all nf-core pipelines and show available updates
nf-core list
```

## Useful Resources

- Nextflow, https://nextflow.io
- nf-core, https://nf-co.re
- nf-tower, https://tower.nf
- Biocontainers, https://biocontainers.pro
- Marine Metagenomics Portal, https://mmp2.sfb.uit.no
- Metapipe, https://mmp2.sfb.uit.no/metapipe/
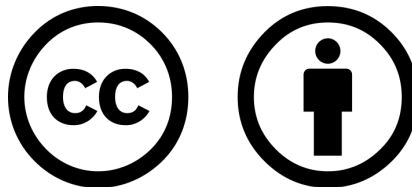
# Thank you!

🌐 **elixir-norway.org**

🐦 **@elixirnorway**

✉ **contact@bioinfo.no**

**Except where otherwise noted, this work is licensed under a**

Creative Commons Attribution 4.0 International License

**https://creativecommons.org/licenses/by/4.0/**