International Statistical Genetics Workshop

IBG

Boulder, Colorado

2024 Edition

# Intro to Unix / R / computing

## 8AM - Monday 4th March

José J. Morosoli

Elizabeth Prom-Wormley

with thanks to Sarah Medland

Please answer the survey before we start!

https://forms.office.com/e/ibLW7WNNec

# Welcome!

Getting the most out of the workshop:

► Ask questions!

► Try to sit next to someone you don't already know.

► Work with someone with different skillset and experience level.

► You will have access to your files after you leave.

► Come to the social functions.

► Do ask questions!!

✓ In person, email, or on the forum.

https://forms.office.com/e/ibLW7WNNec

https://isgw-forum.colorado.edu/

https://isgw-forum.colorado.edu

# A diverse community!

https://forms.office.com/e/ibLW7WNNec

**UCL**

Murcia Twin Registry

Diverse skillsets, backgrounds, research focus, experience levels… and timezones!

José J. Morosoli

QIMR Berghofer
Medical Research Institute

Be open, be kind, be respectful. We are all learners.

# Survey results!

# Getting started

► https://workshop.colorado.edu/

# Getting started



command-line interface (CLI)

graphical user interface (GUI)

# Why Unix? Why?



✓ It allows you to automate tasks.

✓ Replicability: one script, multiple re-runs.

✓ It is more efficient (i.e., fast), scalable and stable than other systems, not to mention open source (e.g., Windows, MacOS).

✓ Big data (e.g., genetic data) is usually stored and analyzed in high-performance computing (HPCs) environments, which for multiple reasons (see above) are based on this language.
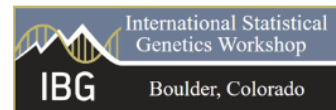
# Getting started

► Open the first exercise:

https://qualtrics.ucl.ac.uk/jfe/form/SV_0pHnuiW6juZ9ezs
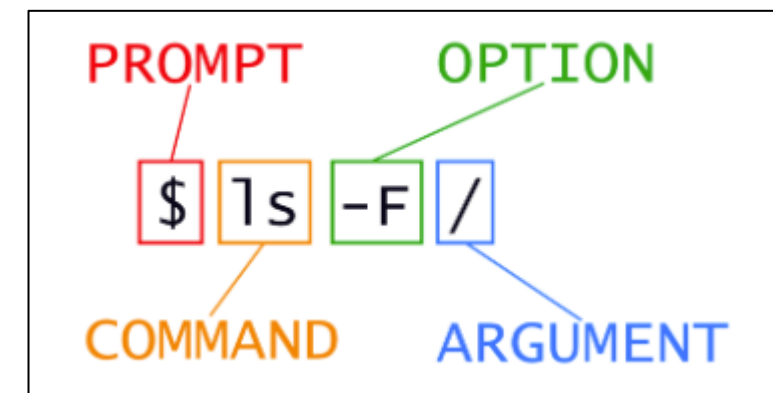
► Open the SSH tab and log in with your username.

# Intro to Unix: Glossary



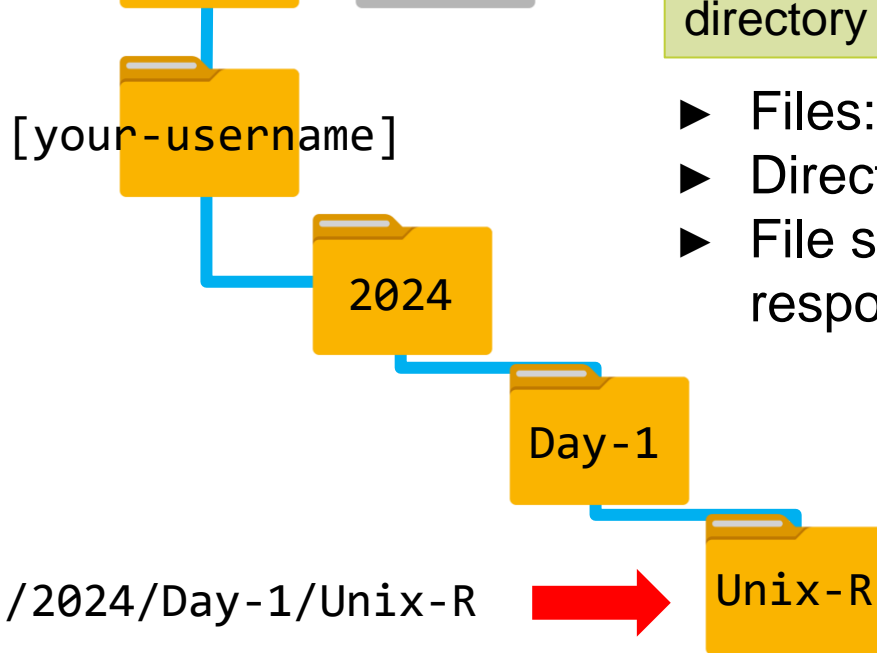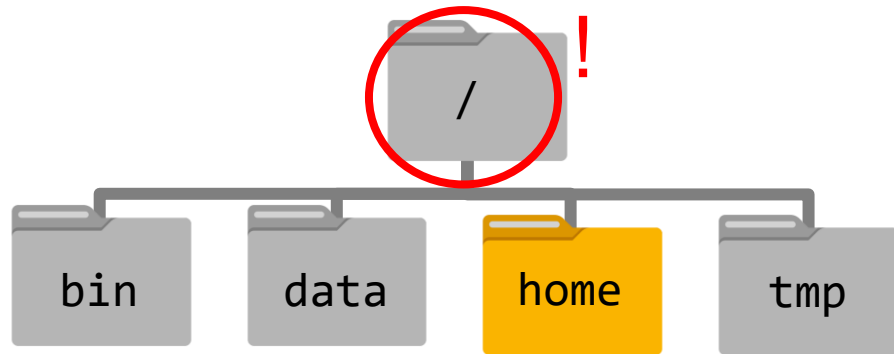► **Graphical User Interface (GUI):** platform to interact with a computer that involves point and clicking and using menus.

► **Command Line Interface (CLI):** text-based platform where you can input text commands to interact with a computer.

► **Unix**: an operative system (just like Windows or MacOS).

► **Shell**: program where users can type commands.

► **Bash**: most popular shell in Unix.

► **Prompt**: symbol that indicates that the shell is waiting for input.

► **Command**: Pre-defined "words" that tell the system what to do, they can be modified using **options** and sometimes require **arguments** to indicate what files to operate on, the paths to find them, etc.
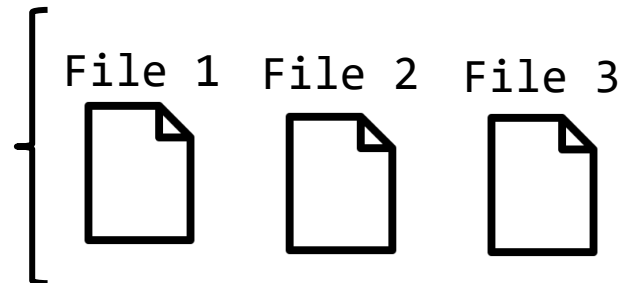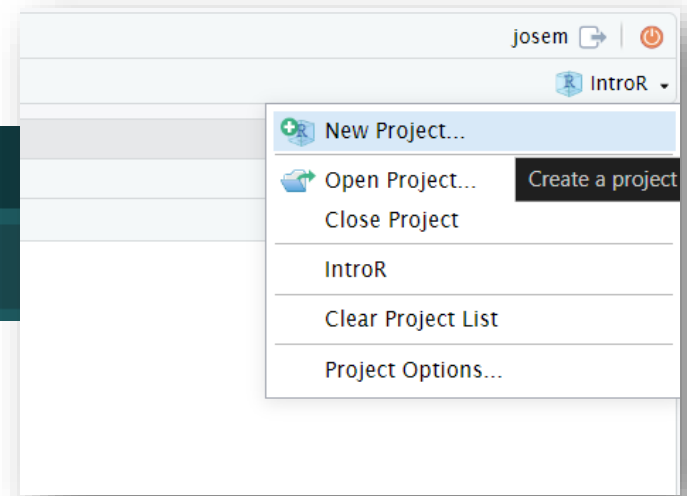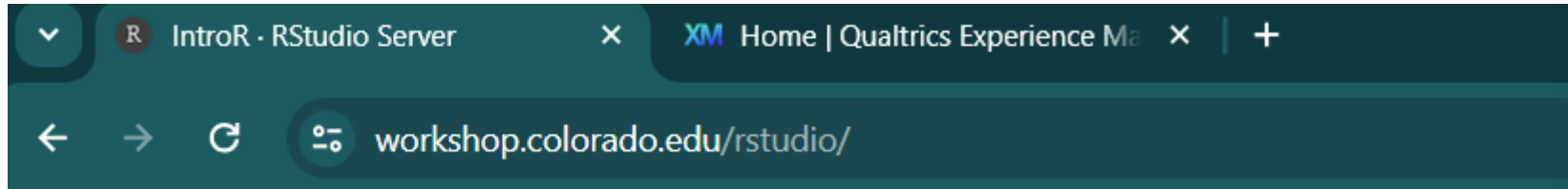
# Setting up our home directory



**Note:**
We need to include " / " in the path whenever we want to give Unix the full path to the file or directory we are interacting with or using.

► Files: store information.
► Directories: hold files or other directories.
► File system: part of the operating system responsible for managing files and directories
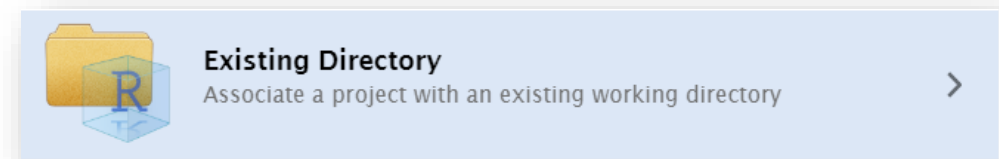
# Next stop: R

▶ Go to the next Qualtrics link:

  https://qualtrics.ucl.ac.uk/jfe/form/SV_6FGfSGwaABJQo1E

▶ Create a R project:
  ▪ Creates/finds a working directory for you.
  ▪ Remembers its location.
  ▪ Makes it easier to resume work after a break.

▶ Click on IntroToR2024.R to open it.

Directory: ~/2024/Day-1/Unix-R

# Similarities between Unix and R

✓Unix and R offer command-line interfaces.

✓Both support scripting for automation.

✓Active communities and extensive package.

✓Open-source.

✓Customizable.

✓Similar commands but not the same!

- E.g., `pwd` vs `getwd()`



**DeepAI** (b. 2016)
***Unix r programming languages are friends,*** 2024
Digital work

# Resources

▶ Software carpentry: https://software-carpentry.org/lessons/index.html

▶ https://stackoverflow.com/ and https://unix.stackexchange.com/

▶ UNIX cheatsheet in `/home/josem/2024/Day-1/Unix-R`.

▶ Specific resources for R:

1. R for SAS and SPSS Users: https://science.nature.nps.gov/im/datamgmt/statistics/R/documents/R_for_SAS_SPSS_users.pdf
2. An R Style Guide: http://google-styleguide.googlecode.com/svn/trunk/Rguide.xml
3. Rseek: A search engine for all things R related (http://www.rseek.org)
4. R-Bloggers.
5. Quick-R's section on graphics: http://www.statmethods.net/advgraphs/parameters.html
6. More information on data frames: http://www.r-bloggers.com/exploratory-data-analysis-useful-r-functions-for-exploring-a-data-frame.
7. Details on how to develop your own package: http://r-pkgs.had.co.nz

# R Code – Best practices



Keep names short (≤25 characters).

Choose names using 3-4 key unchanging pieces of information.

Use YYYY-MM-DD format for better sorting even over the span of many years.

All numeric fields should be zero-padded for equivalent width.

For better visibility, give preference to dashes over underscores.

Create a README.txt describing the file naming convention.

Add tags to files properties to enhance their findability in your workspace.

Be consistent!

Do not add spaces! They are often interpreted as delimiters and may cause problems.

Do not include special characters such as: " / \ [ ] : ; | = , < ? > & $ # ! ' { } () *.

Do not rely on case to distinguish filenames. Not all systems are case-sensitive.

Avoid unnecessary repetition in file names and file paths.

Avoid using words such as 'draft' or 'letter' at the start of file and folder names.

If files will be shared and edited by multiple people, avoid naming multiple versions. Consider using a version control system such as Git instead.

https://www.r-bloggers.com/2018/09/r-code-best-practices/

https://www.library.ucsb.edu/sites/default/files/dls-n01-2021-filenaming.pdf