# From Data Chaos to Data Harmony: Managing NGS Data in a Wet Lab

*Jose Alejandro Romero Herrera*

DeIC conference 2023-11-08

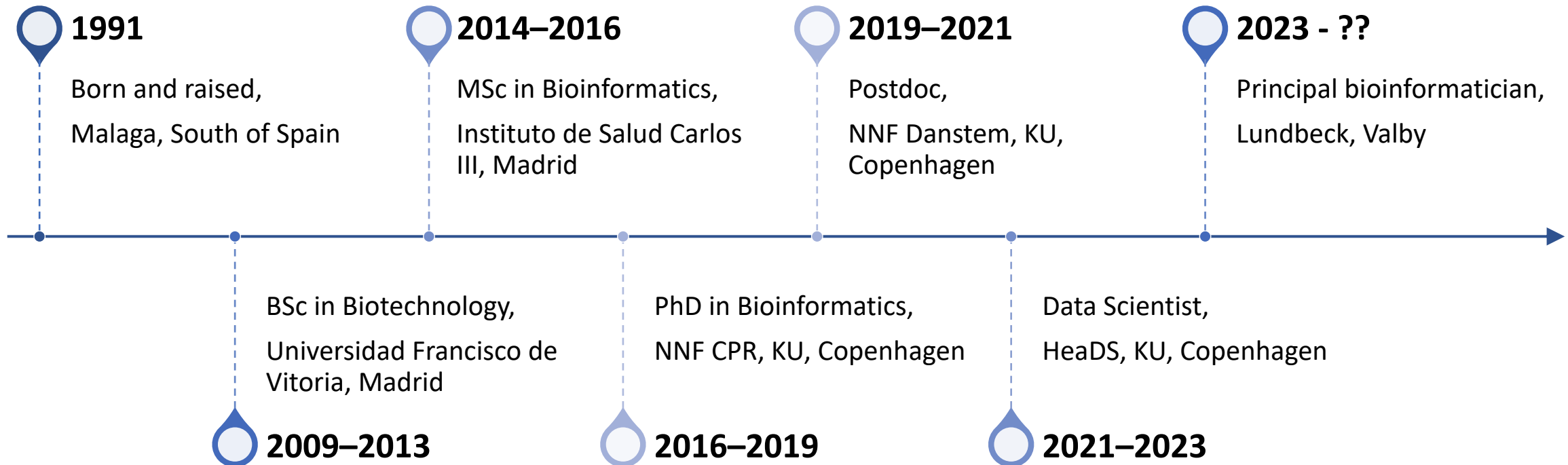UNIVERSITY OF COPENHAGEN

HeaDS

# Agenda

About

Intro

The problem

Simple solutions

Future work

Acknowledgments

# About HeaDS

- Center for Health Data Science, KU
  - Strengthen health data science at KU
  - Serve as a hub for researchers
  - Provides consulting services
  - Teaching in various topics

Anders Krogh

- Working for the Sandbox project
  - Multi-disciplinary
  - Multi-institutional
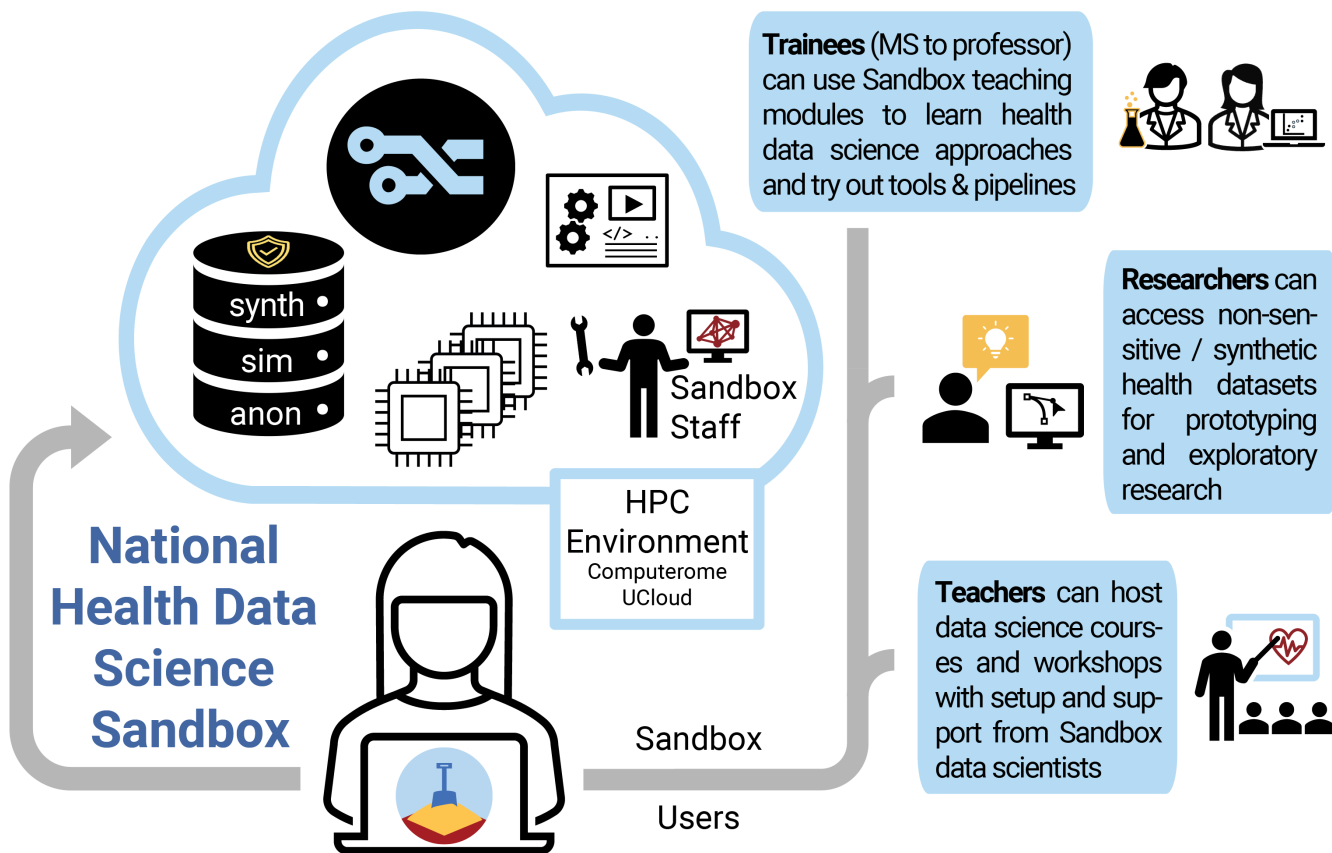  - Virtual Environment for Health Data Science

Jennifer Bartell

UCPH HeaDS

National Health Data Science Sandbox

HeaDS

UNIVERSITY OF COPENHAGEN

# About Sandbox

**Trainees** (MS to professor) can use Sandbox teaching modules to learn health data science approaches and try out tools & pipelines

**Researchers** can access non-sensitive / synthetic health datasets for prototyping and exploratory research

**Teachers** can host data science courses and workshops with setup and support from Sandbox data scientists

synth
sim
anon

Sandbox Staff

HPC Environment
Computerome
UCloud

**National Health Data Science Sandbox**

Sandbox Users

**Aarhus University**
Genomics

**Southern Denmark University**
Proteomics

**Copenhagen University**
Transcriptomics

**Denmark DTU**
Supercomputing

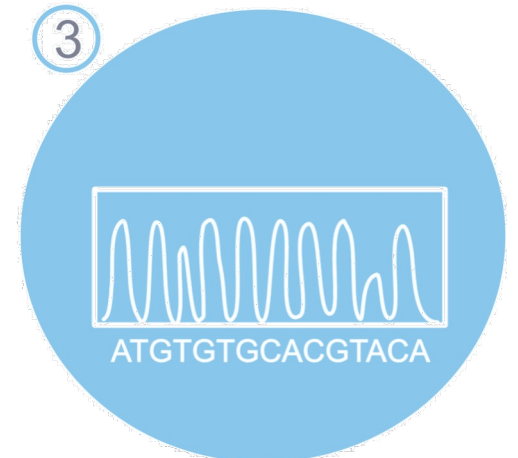**Aalborg University**
Health records & predictive modeling

UNIVERSITY OF COPENHAGEN

HeaDS

# NGS data

- ## Next Generation Sequencing Data
  - ### Determine genetic information

- ## Applications:
  - ### Genomic research
  - ### personalized medicine
  - ### forensic science, and more

- ## Millions of DNA sequences generated
  - ### Really big datasets (Gb per sample)
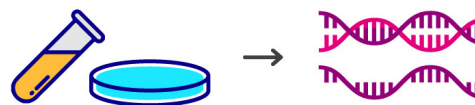
① Extraction

② Library

③ Sequencing

ATGTGTGCACGTACA

HeaDS

UNIVERSITY OF COPENHAGEN

# RDM for wet lab



- RDM supported options
  - ELNs and LIMS

- Tracks:
  - Experiments, protocols, samples
  - Order of chemicals
  - Instruments and physical storage

- **Bioinformatics data and its analysis is not tracked at all!**



STEP 1:
Extraction

STEP 2:
Library
Prep

Template
Fragmentation    PCR or RT-PCR
Fragmented DNA    Amplicons
Adapter Ligation PCR    Adapters
Sequencing Library

iRepertoire

STEP 3:
Sequencing

STEP 4:
Analysis

FAST Q    Align Reads    BAM    Identify Variants    VCF

>read1
aacgctcgtacttagctct
agctacggatcgctacgga
ctaggtcactcgctatctata
aaaactccgctctttctgcgc
gcgatcgactcgatctacgc
ggttggtaccgcatcactacg
ccgatctagc

UNIVERSITY OF
COPENHAGEN

HeaDS

# Bioinformatics in a wet lab

# Data chaos

Digital version of a chaotic desktop

- Data accumulated for several years
- No common structure
- No common file naming conventions
- No metadata or documentation
- No provenance
- Massive problems when staff leaves



*Wonderlane on Unsplash*

# Data chaos

## Me, myself and only I to blame

# Solution: FAIR principles

- Many guidelines on FAIR principles

- Not much applied to NGS data
  - Some few examples, but not applications

- Simple guidelines for any bioinformatician
  - Basic command line experience
  - Version control with git/GitHub



FAIR DATA PRINCIPLES

AH!

FINDABLE

DOI 10.1017/8.787

ACCESSIBLE

HOW DO YOU OPEN A .XZQ FILE?

INTEROPERABLE

HERE

REUSABLE

HeaDS

UNIVERSITY OF COPENHAGEN

# Simple rules for NGS RDM

1. Adhere to **folder structures** and **naming conventions** using **templates**
2. Fill a **DMP template** that it is **prefilled** with common information
3. Create and fill **metadata file** and **README file** in each folder
4. Make **a browsable database** from all metadata files
5. Use **community-curated workflows** for data preprocessing
6. **Version control** data analysis with git/Github
7. Display **data analysis reports** with GitHub Pages
8. **Archive** data and data analysis in **repositories** (Zenodo-GEO/Annotare)

HeaDS

- Create custom templates using **COOKIECUTTER**

- Command line utility
  - Very flexible
  - Simple, but can do complex things

- Two templates
  - *Assay* folders: NGS data
  - *Project* folders: Data analyses

UNIVERSITY OF COPENHAGEN

HeaDS

# 1. Folder templates and naming conventions

**Data folder ("assays")**

- Subfolders → NGS experiment data

- Read only, no duplicates

- Unique ID, human readable

- Metadata file
  - ID, keywords, tech, author, date, etc.
  - Controlled vocabularies

- README.md: additional details

- Data organization:
  - raw, processed, pipeline
  - pipeline: community curated workflows

```
(ngs)
sarahlu at sarahlu-ThinkPad-T14-Gen-2i in ~/Documents
$ tree RNA_20230628
RNA_20230628
├── checklist.md
├── description.yml
├── metadata.yml
├── pipeline.md
├── processed
│   ├── bam
│   ├── bed
│   └── counts
└── raw
    ├── fastq
    └── samplesheet.csv

6 directories, 5 files
```

# 1. Folder templates and naming conventions

**Project folder ("projects")**

- Subfolders → research project

- Simlink multiple *Assays*, no copies!

- Version controlled

- ID: AUTHOR_DESC_YYYYMMDD

- Metadata and README file

- Folder organization:

  - data, notebooks, reports, scripts, results, docs

- Naming conventions for results and figures

  - Heatmap_DEA_TreatVsControl_YYYYMMDD.tiff



```
sarahlu at sarahlu-ThinkPad-T14-Gen-2i in ~/Documents/Lundregan_RNAex
$ tree
.
├── data
│   ├── assays
│   │   └── RNA_20230628 -> /home/sarahlu/Documents/RNA_20230628
│   └── processed
├── documents
│   ├── Hamilton_etal_2019_Nature.pdf
│   └── project_overview.docx
├── notebooks
│   ├── 01_preprocessing.qmd
│   └── 02_differential_expression.qmd
├── _quarto.yml
├── README.md
├── reports
│   ├── 01_preprocessing_files
│   │   └── figure-html
│   │       ├── fig-01_model_fit.png
│   │       └── fig-01_PCA.png
│   ├── 01_preprocessing.html
│   ├── 02_differential_expression_files
│   │   └── figure-html
│   │       ├── fig-02_heatmap_0h_vs_2h.png
│   │       └── fig-02_volcano_0h_vs_2h.png
│   ├── 02_differential_expression.html
│   ├── index.html
│   ├── README.html
│   └── search.json
```

UNIVERSITY OF COPENHAGEN

HeaDS

# 2. Prefilled DMP template

**DMP**ONLINE

## Non-sensitive NGS research project template

| Project Details | Plan overview | Write Plan | Share | Download |

expand all | collapse all                                    32/41 answered

1. Data Summary (6 / 6)                                              +

2. FAIR data (24 / 24)                                               +

3. Other research outputs (2 / 2)                                    +

- Prefilled with repetitive info
  - GEO, Github, Labguru…
  - Metadata and standards used
  - How it adheres to FAIR

- Streamline the process of writing your DMP

- Shared publicly

- More templates can be created depending on data

  - Sensitive datasets
  - Other omics data
  - Imaging data

UNIVERSITY OF COPENHAGEN

HeaDS

# 3. Metadata and README

- Cookiecutter template will require you to fill metadata fields

- Collected metadata will be saved in a metadata.yml file

- Short descriptions in README.txt

| Metadata field | Convention | Example |
| --- | --- | --- |
| assay_type | - | ChIP-seq |
| owner | <Initials> | JARH |
| creation_date | <YYYYMMDD> | 20231108 |
| platform | - | Illumina |
| organism | <Genus species> | Homo sapiens |
| nsamples | <integer> | 9 |

```
# NGS Analysis Project: Exploring Gene Expression in
Human Tissues

## Aims

This project aims to investigate gene expression
patterns across various human tissues using Next
Generation Sequencing (NGS) data. By analyzing the
transcriptomes of different tissues, we seek to uncover
tissue-specific gene expression profiles and identify
potential markers associated with specific biological
functions or diseases.
```

HeaDS

UNIVERSITY OF
COPENHAGEN

# 4. Browsable database

- Collect metadata files

- Create a tsv or SQL database

- Browse it with Shiny R app

- Very useful for all lab members

# 5. Community-curated workflows

nf-core

- A community effort to collect a curated set of analysis pipelines
- Uses **nextflow** language, specialized for reproducible workflows
- Many gold-standard bioinformatics and NGS pipelines

# 6. Version control data analysis

- Use version control framework for your data analysis: Git
  - NGS data is too big though

- Online repositories: GitHub
  - Create your lab organization

- Benefits
  - Enhance collaboration
  - Public sharing of the analysis of your data
  - Easier tracking of changes and results

# 7. Data analysis reports

- Lab webpage using GitHub Pages
  - Intro to lab
  - Feature papers or data analyses

- Display your data analysis
  - Public after publication
  - Transparency of results and analysis

- Could be used as tool documentation

## 4.2.1 Principal component plot of the samples

PCA plot using the first two components

```
pcaData <- plotPCA(vsd, intgroup=c("ShortLabel"), returnData=TRUE)
percentVar <- round(100 * attr(pcaData, "percentVar"))
```
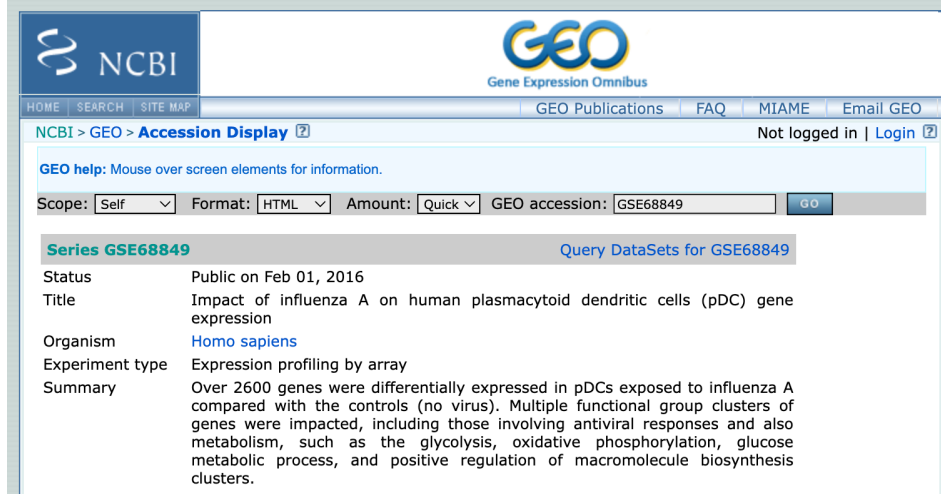
```
ggplot(pcaData, aes(PC1, PC2, color=ShortLabel, label=ShortLabel)) +
  geom_point(size=3) + geom_text_repel() +
  xlab(paste0("PC1: ",percentVar[1],"% variance")) +
  ylab(paste0("PC2: ",percentVar[2],"% variance")) +
  coord_fixed() + theme_bw()
```



HeaDS

UNIVERSITY OF COPENHAGEN

# Online course

https://hds-sandbox.github.io/RDM_NGS_course/

### RDM for NGS data workshop

Introduction   Course contents   DTU workshop   Keyword index   Contributors

📋 **Overview**

📖 **Syllabus:**

1. What is Research Data Management and why it is important

2. What is NGS data

3. Data Life Cycle

4. Open Science and FAIR principles

5. Data Management plans

6. Folder and file structures applied to NGS data

7. Metadata applied to NGS data

8. Create a database of your data and projects

9. Version control of your data analysis

10. Archiving and repositories

HeaDS

UNIVERSITY OF COPENHAGEN

# Future work

- Metadata ontologies and controlled vocabularies
  - Not easy to implement and enforce
  - Some examples are provided for different metadata fields

- Version control of the actual NGS data

- Need of command line experience and work for cookiecutter
  - Difficult for experimentalists

- Interaction between LIMS + ELN systems
  - Cross-linking data analysis with wet lab RDM
  - One place to go for everything?

HeaDS

# Acknowledgements

## Brickman lab



## HeaDS center



UNIVERSITY OF
COPENHAGEN

HeaDS