# Implementing Principles of Reproducible Research at Scale in a Bionformatics Core Facility: Challenges and Solutions

Hubert Rehrauer

# FGCZ Genome Informatics

**Vision**

- Enable design to result workflows using Next Generation Sequencing
- Generate exciting results from the analysis of omics data for and with ETHZ/UZH researchers
- fill the gap (if there is any) between user skills/knowledge and existing bioinformatic tools

**Notes**

- Reproducibility is never asked for but considered as granted
- Nobody ever asks us about reproducibility

# Mission of FGCZ Bioinformatics

## Data Processing

- support **data generation** by the wet lab units of the FGCZ
- operate **data processing** infrastructure and implement **data quality control**

## Data Analysis

- **operate data analysis infra-structure**
- **perform data analysis services** for ETH/UZH researchers
- collaborative data analysis within research projects
- training and education in omics areas

# Reproducibility in the context of research projects

- Funding bodies set the scene for research projects

- SNF encourages reproducibility of analysis results and long-term data re-use
- Responsibility is with the PI
  - have proper strategies for data analysis and data management (DMP)
  - make sure project members apply appropriate practices
  - make sure service providers (core facilities) are compliant
- FGCZ Genome Informatics mainly acts as provider for part of the analysis trail

# Research Cycle

- Bioinformatics service usually covers part of the research cycle
- FGCZ Bioinformatics
  - supports and consults on all steps
  - has full responsibility the *Data Processing* and *Data Study & Analytics*
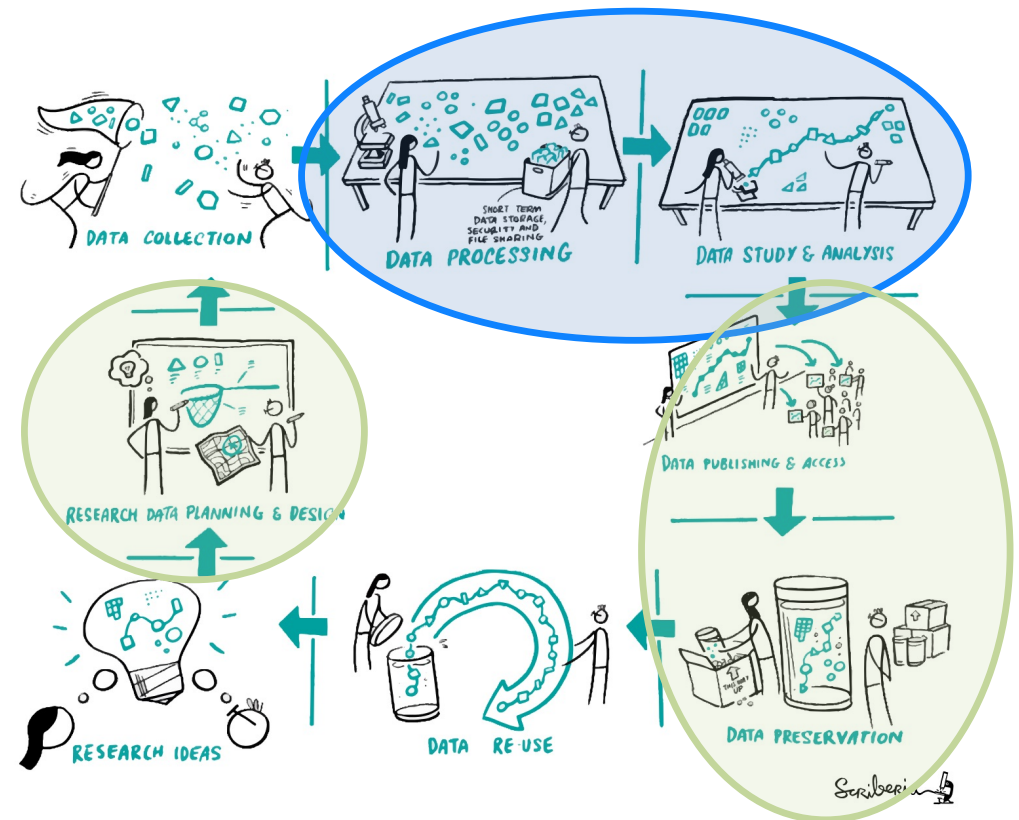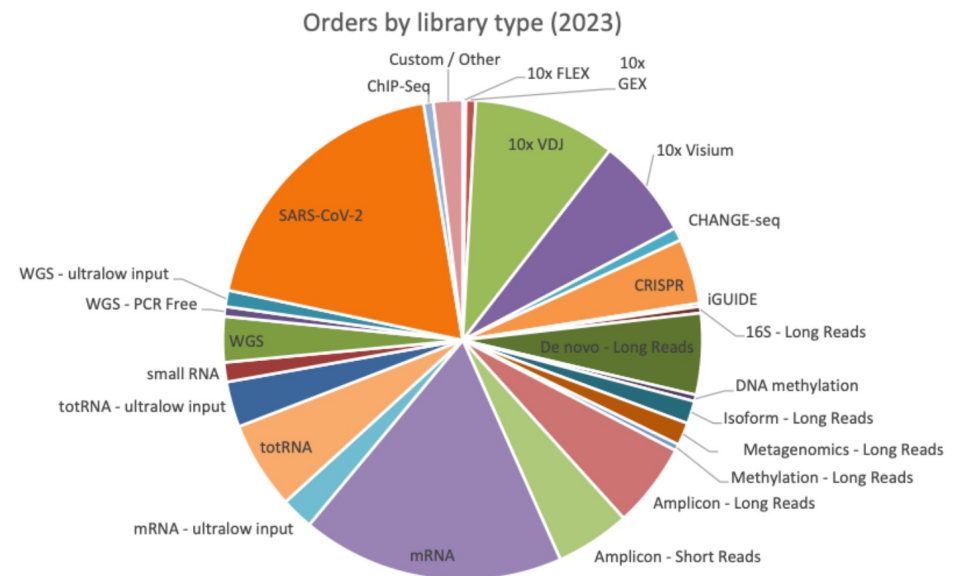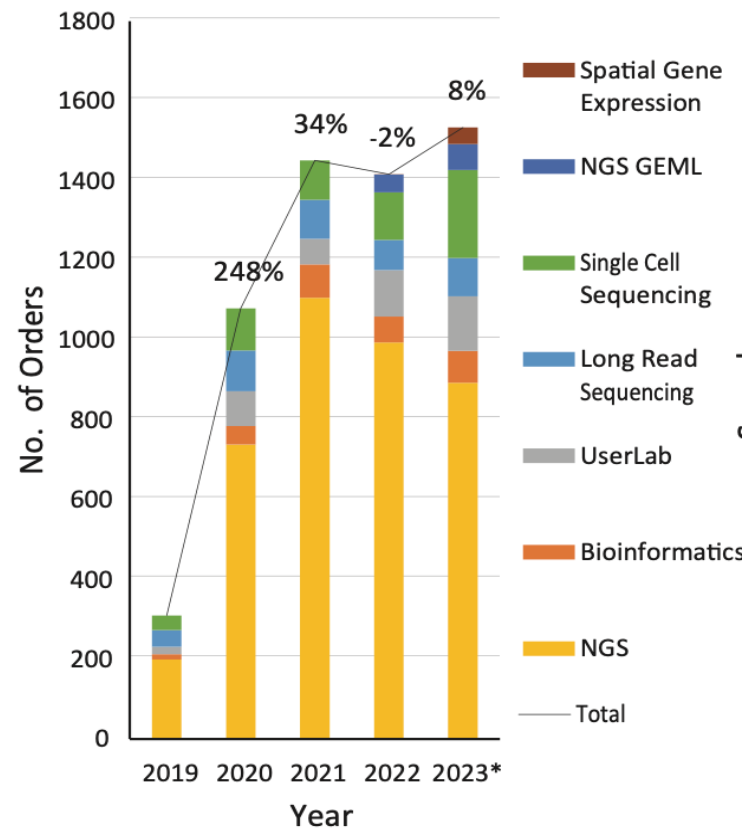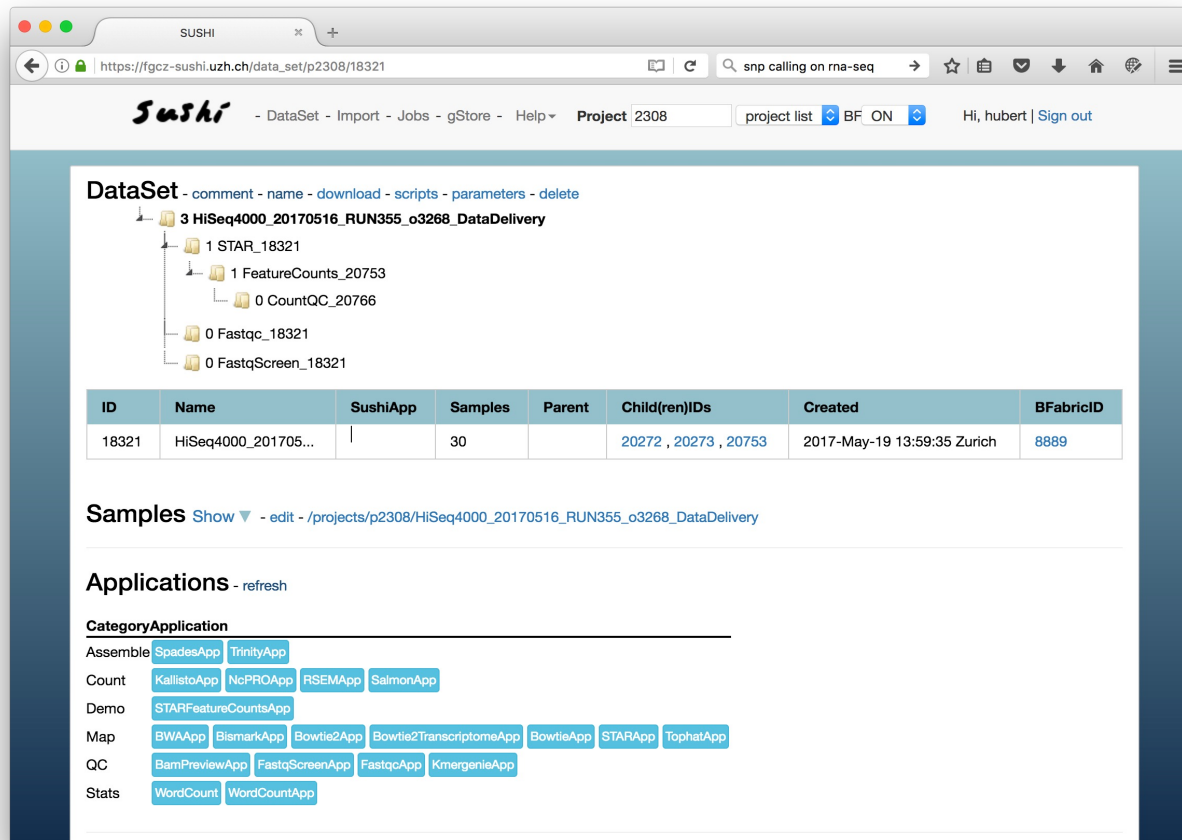  - supports *Planning* and *Data Preservation* according to DMPs



Fig. 8 *The Turing Way* project illustration by Scriberia. Used under a CC-BY 4.0 licence. DOI: 10.5281/zenodo.3332807. #

https://the-turing-way.netlify.app/reproducible-research/overview

# Core Facility: Throughput and Diversity





Orders by library type (2023)

- Large Number of experiments
- Many different types of analyses
- *Aspiration to be reproducible*

# SUSHI



- offers SUSHI apps to perform individual analysis steps
- ultimately generates an entire analysis trail
- relies on a web server and an associated database to run analysis steps

**http://fgcz-sushi.uzh.ch**

**DataSet**

| ID | Name | Samples | Parent | Child(ren)IDs | Created |
|----|------|---------|--------|---------------|---------|
| 736 | ventricles | 4 | | 8641 | 2014-02-13 19:34:56 |

- 📁 **1 ventricles** RNA-seq - genotype effect - full data set
  - 📁 1 Map_STAR_736
    - 📁 2 Count_FeatureCounts_8641
      - 📁 0 QC_CountQC_11632
      - 📁 0 Differential_Expression_DESeq2_11632

- web-based
- fully reproducible
- self-contained data sets
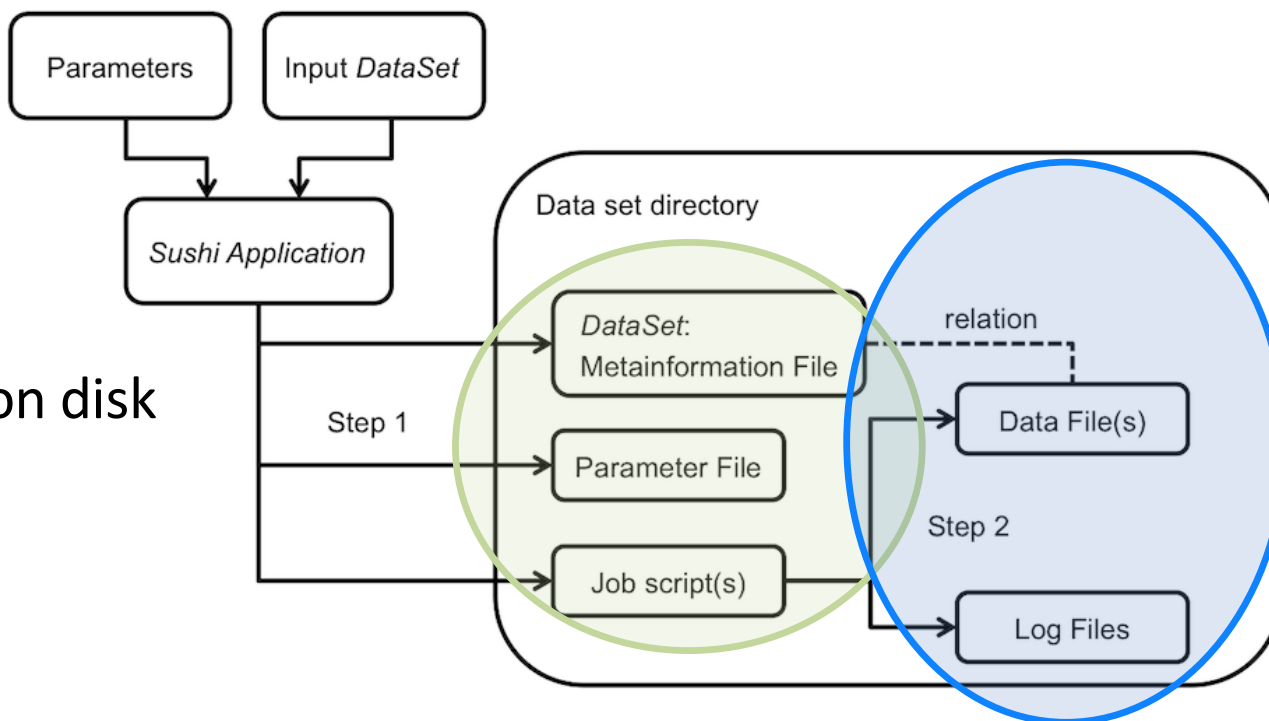
Hatakeyama et al. 2016

## SUSHI Informatics

Step 1:
- Analysis Generation
- Done by the SUSHI App
- Writes the intended analysis on disk

Step 2:
- Analysis Execution
- Performs the analysis
- Writes results on disk



**SUSHI Applications are not stored persistently**

**Analysis Generation is not long-term repeatable**

**Actual Data Analysis is repeatable**

## SUSHI Apps and Workflows

- SUSHI models only individual steps
- Workflows are built as a series of steps; workflows are not explicitly modelled and managed
- Nevertheless: An individual step may actually perform an entire workflow in one step (e.g if the step executes an nf-core workflow)

- With a focus on steps there are fewer things to maintain

# Available SUSHI Apps



- Currently ~80 Apps are available
- Apps are plugins that can be removed or updated based on needs

# Going stepwise

- Advantages:
    - modification of analysis trail can happen ad hoc
    - encourages revisiting of results of individual steps and adapt subsequent steps
    - number of steps much smaller than the number of sequences of steps, a.k.a. workflows
    - analysis steps re-used in different workflows

## FGCZ Genome Informatics – data analysis portfolio

- What do those job scripts do?
- Software environment:
  - R with > 1000 packages installed
  - module environment with ~60 tools where we support different versions
  - ~100 tools installed with one version
  - ~ 30 conda environments for group-wise use
  - numerous personal conda environments, for individual use

# Software Versions: Updates vs Consistency

- Conflicting interests:
  - run latest version vs consistent analysis within a project
- Update strategy:
  - Keep different version of tools within the modules; SUSHI can call specific versions of tools → support for different versions in different projects
  - Update R and all R-packages twice a year
  - All own code is git version-controlled and continuously updated; we keep track of the git tag
    → let's us continuously update our code

```
## ezRun tag: 34d4529f71ed91c9c6cc6fc3fb4e04c9ffee04aa
## ezRun github link: https://github.com/uzh/ezRun/tree/34d4529f71ed91c9c6cc6fc3fb4e04c9ffee04aa
##
## R version 4.4.2 (2024-10-31)
## Platform: x86_64-pc-linux-gnu
## Running under: Debian GNU/Linux 12 (bookworm)
##
```

# Software life cycle

- SUSHI offers push-button repeatability but only for a short term
  - applications in SUSHI are retired if unused or if superseded
    - retired applications are effectively removed
    - retired applications can no longer be run on new data
    - previous analysis can only be repeated on command line
- SUSHI doesn't keep any legacy!

# External Data: Genomes and Annotations

- Data analysis often uses additional data/knowledge from external variable sources, e.g. genome assemblies and annotation
- We keep
  - local versions: https://fgcz-gstore.uzh.ch/reference/Homo_sapiens/GENCODE/GRCh38.p13/Annotation/Release_34-2024-10-17/Genes/
  - script that generated the local version

- Generation of local genome copy is not repeatable if source files disappear from provider (NCBI, GENCODE, ….)

16

# External Data: Enrichr

- API calls to external systems, e.g. Enrichr that host pathway knowledge bases → reproducibility not guaranteed

# Where are we on the Reproducibility Spectrum?

- We commit to provide analysis code and, to limited extend, data

# The issue with Data Rights

- Data rights are managed by the Principal Investigator
- PI is responsible for consent
- Mandate of core facility is usually to process and analyse the data not to long-term store the data

- Long-term storage implies long-term costs and requires funding which is usually not provided by PI
- PIs usually provide only the minimal context information on the samples that is necessary to achieve the analysis goals for the NGS data
  BUT: without full context information, the data has only limited value

# Analysis Throughput & Requests for Reproduction



SUSHI Datasets

SUSHI App:
- JunctionSeqApp
- SCRNAVelocityApp
- VirDetectApp
- MageckTestApp
- MpileupApp
- DnaBamStatsApp
- SCReportMergingApp
- SCMultipleSamplesApp
- SCCountsApp
- VPipeApp
- SCCountQCApp
- SCReportApp
- Bowtie2App
- Cov19QcApp
- SCOneSampleApp
- KallistoApp
- ScSeuratCombineApp
- ScSeuratApp
- RnaBamStatsApp
- FastqScreen10xApp
- CellRangerApp
- Fastqc10xApp
- FeatureCountsApp
- STARApp
- CountQCApp
- DESeq2App
- EdgeRApp
- FastqScreenApp
- FastqcApp

- Requests for repetition:
  - very few
- Requests for reanalysis
  - dozens per year (updated tools; different thresholds for sensitivity/specificity, …)
- Failed reanalysis/repetition
  - none

# Other approaches: Galaxy

## SUSHI

- repeat: rerun the static job script from the SHELL
- workflows: not modelled

- requires compute environment available

## GALAXY

- repeatability: rerun the workflow in the GALAXY interface
- workflows: explicitly modelled
- keeps track of history of workflows

- requires GALAXY instance

# Repro Challenge: False Positives



**RESEARCH ARTICLE**

## Estimating the reproducibility of psychological science

**Open Science Collaboration** *,†

+ Author Affiliations

↵†Corresponding author. E-mail: nosek@vi

*Science* 28 Aug 2015:
Vol. 349, Issue 6251,
DOI: 10.1126/science.aac4716

# Summary

- Commitment to generate biologically interpretable and reproducible results
- Constraints:
  - dozens of workflows for hundreds of projects in a changing environment with data rights managed by users
  - analysis often needs to integrate in a larger workflow
- Minimalistic approach to Reproducibility
  - documented analysis workflows with scripts that enable long-term **repeatability** by a skilled bioinformatician (report the information necessary to understand and repeat the analysis)
  - push-button **repeatability** only available for ~1 year
- Results of analysis steps as self-contained and documented folders that enable reproducibility by other groups
- Long-term average of ~3 requests per year for repeated analysis