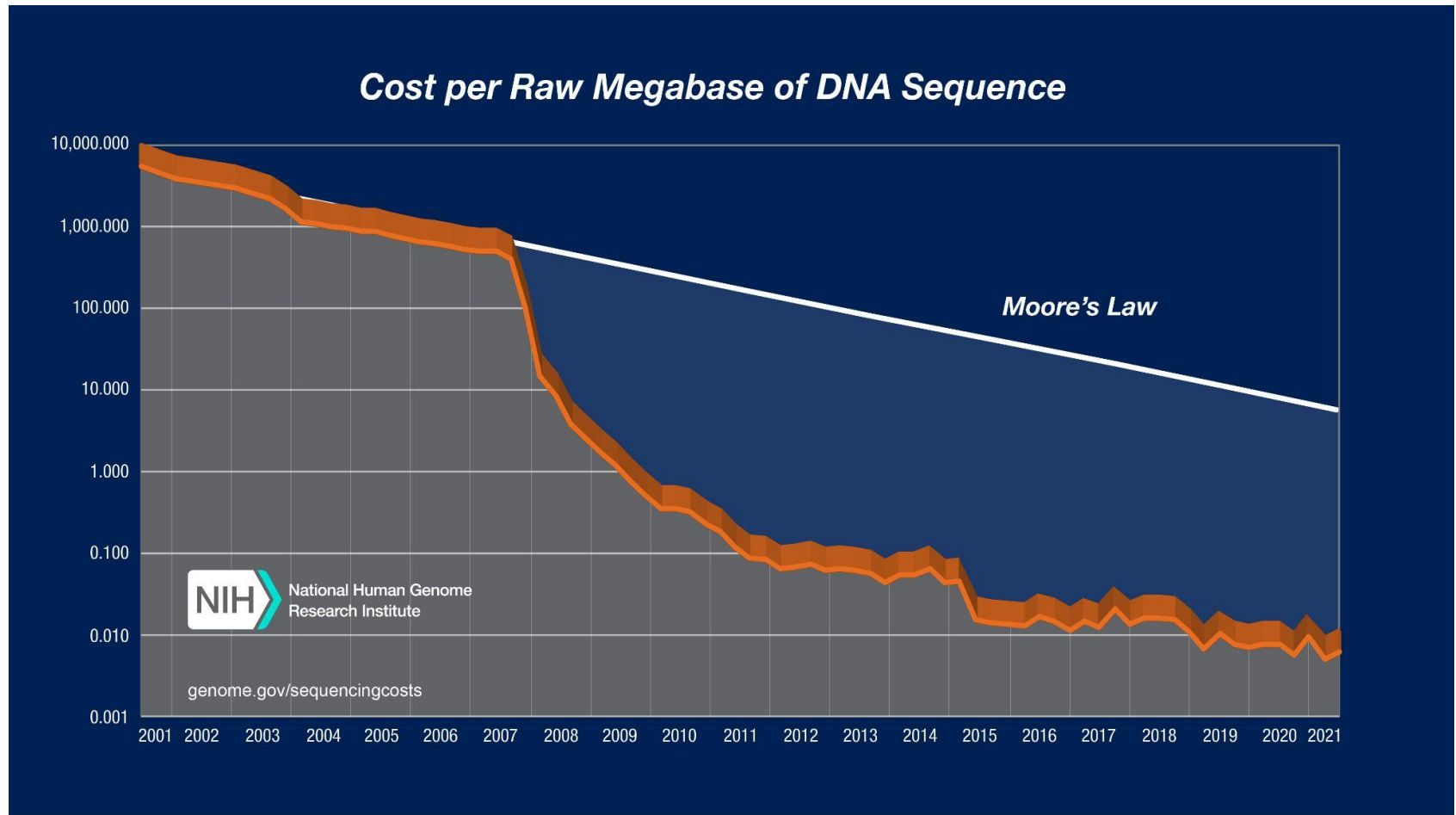# Introduction to data literacy in genomics

Prof. Dr. Boas Pucker and Katharina Wolff
(Plant Biotechnology and Bioinformatics)

# Availability of slides

- All materials are freely available (CC BY) - after the lectures:
  - StudIP: **Data Literacy in Genomics**
  - GitHub: https://github.com/bpucker/teaching

- Questions: Feel free to ask at any time

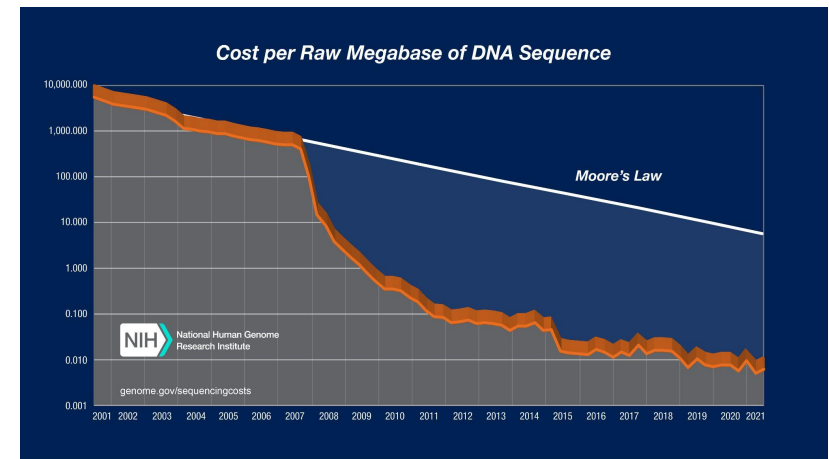- Feedback, comments, or questions: b.pucker[a]tu-bs.de

My figures and content can be re-used in accordance with CC BY 4.0, but this might not apply to all images/logos. Some figure were constructed using bioRender.com.

Technische
Universität
Braunschweig

# 'Big Data'



Cost per Raw Megabase of DNA Sequence

# Reasons for growth of databases

- Data generation costs are dropping
    - Rapid sequencing technology development
    - Resolution of pictures is increasing
    - Robotics supports data acquisition

- Data storage capacities are increasing
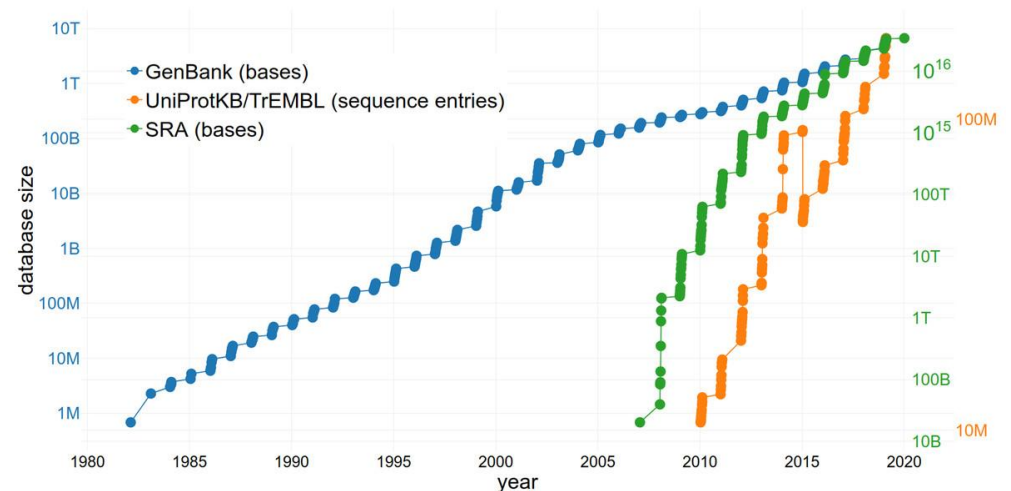
- Potential of data reuse is recognized



*Cost per Raw Megabase of DNA Sequence*

Technische
Universität
Braunschweig

# What is Data Literacy?

# Data Literacy

- 'Competency to handle data properly'

- 'Read, understand, create, communicate data as information' (wikipedia)

Technische
Universität
Braunschweig

# Why is Data Literacy important?

- We live in a world of data: important for science, business, and society

- Objective: data informed decision making

- Sizes and complexity of data sets are increasing ('BigData')

- Data is accumulating over time

# Why is Data Literacy important for you?

- Data (interpretation) is valuable

- Data scientists are needed

- Combination of different fields is particularly powerful

**Technische Universität Braunschweig**

# More definitions

- Data: collection of observations
    - 2020: 312 red flowers; 95 white flowers
    - 2021: 298 red flowers; 98 white flowers

- Knowledge: accumulation of facts and data
    - Ratio of red flowers to white flowers is 3:1

- Insights: grasping the underlying nature of knowledge; understanding general concepts
    - Flower color has a genetic basis
    - Red flower allele dominant over white flower allele

# Digital Object Identifier (DOI)

- DOI = Digital Object Identifier

- Unique and short way to point to a publication or a data set

- How to resolve a DOI? https://dx.doi.org/

# Data management

- Larger data sets require more efficient data management

- Data management plans required in project proposals

- Many services and organizations emerging (NFDI)



https://nfdi4plants.de/

# Data protection and data security

- Data protection is becoming a huge issue in the EU

- Avoid any personal data in your research data sets

- Data security is gaining relevance

- German universities are frequently attacked:
  - Gießen (2019): offline for weeks + 1.7 million direct costs
  - HHU Düsseldorf (2020): ransom attack on clinic
  - Bochum, Dresden, Freiburg, Berlin

Technische
Universität
Braunschweig

# Electronic Laboratory Notebooks (ELN)

- Links between data sets and documentation

- Avoids repetitive documentation (copy&paste); version control

- Automatic search

- No issues with handwriting

- Quick and regular backups possible

- Accessible from everywhere

- Collaborative

Technische
Universität
Braunschweig

# Technical solutions for electronic lab notebooks

- Simple wiki page

- Dedicated systems developed by research institutions
  - GABI-Kat LIMS
  - Chemotion
  - e!DAL
  - gitlab

- Commercial offers
  - Benchling
  - Dotmatic's
  - Signals Notebook (Perkin Elmer)
  - LabArchives

# Example: wiki

- Electronic lab notebook for collaborative learning

- iGEM wiki to document team projects

- Github wiki to document tools

# Laboratory Information Management System (LIMS)

- Workflows are represented and samples are tracked

- Information are linked between analysis

- Different levels of permissions/access

- Examples:
  - Handle a collection of T-DNA insertion lines (GABI-Kat)

  - Manage all samples submitted for sequencing

  - Manage all oligonucleotide orders



https://commons.wikimedia.org/wiki/File:Gnulims_lab_requests.png

Technische
Universität
Braunschweig

# Example: Chemotion

- Research data management tool for chemists

- Electronic Laboratory Notebook (ELN)

- Repository for research data (easy transfer from ELN)

- DOIs are assigned to datasets and protocols



https://www.chemotion.net/chemotionsaurus/index.html; © Chemotion - KIT

# Example: e!DAL

- Maintained at IPK; supported by de.NBI

- DOIs are assigned to submitted data sets

- Reporting about access statistics

- Version tracking

- FAIR compatible and data linkage



https://edal.ipk-gatersleben.de/

# Example: gitlab

- Free version control solution

- Can be used for software development, but also suitable as ELN

- GITZ offers a central service at TUBS: https://doku.rz.tu-bs.de/doku.php?id=server:gitlab

- Non-commercial alternative to github



https://gitlab.com/gitlab-org/gitlab



https://doku.rz.tu-bs.de/doku.php?id=server:gitlab

# Example: Benchling

- Cloud-based platform for biotechnology

- Templates for documentations

- Convenient to use

- Compatible with various platforms

- Suitable for certified processes

Technische
Universität
Braunschweig

# Lab 4.0: digital lab for higher efficiency

- Connection of lab and analysis processes

- Data from different instruments are synchronized through a cloud
  - Values and pictures are directly inserted into the lab book

- Automation of processes (robotics)
  - Pipetting robot for large scale sample processing
  - Automatic phenotyping facilities

- Samples are labeled with barcodes
  - Barcodes and scanning avoid human errors
  - Can also prevent fraud

# Summary: data is everywhere

- Size of data sets is growing

- More sophisticated data collection methods

- Databases enable dissemination/reuse

- Electronic documentation

- Digitalization makes processes more efficient

# Course challenge

# What causes the color difference?

# Previous knowledge?

- Are there previous reports of similar observations? (internet search)

- Are there publications about possible explanations? (literature search)

- Are there data sets that could be helpful for a study? (re-use)

# Hypotheses

- Research should be driven by hypothesis (hypothesis testing)

- Hypothesis must be specific and falsifiable

- Proofing hypothesis right is not possible

- Rejecting a hypothesis if possible based on contradicting evidence

- Different hypothesis must not be mutually exclusive

# Example: hypotheses

- pH value of the soil is responsible for the color

- Presence of metal ions in the soil is responsible

- Genetic factors determines the color

- Pollination status determines the color

- Plants belong to different species

**Technische Universität Braunschweig**

# What do we need to test a hypothesis?

- Example: genetic factor determines the color difference

- Genome sequence (long read sequencing & assembly)

- Structural annotation of the genome sequence (gene prediction)

- Functional annotation of the genes (function prediction)

- Data sets to compare white and red flowering plants (read mapping & variant calling)

- Analysis of gene activity (RNA-seq, data re-use)

**Technische
Universität
Braunschweig**

# Disseminate findings and data

- Submit data sets to appropriate repository

- Share documentation and scripts developed for analyses

- Share findings through talks, posters, or publications

# Data management tipps

# How to structure your data?

- Document every step (e.g. in a README)
  - Origin of data sets; versions of tools, parameters of analyses

- Keep raw data sets separated from scripts and results

- Sort data by project

- Structure analysis related data sets/results by date

**Technische
Universität
Braunschweig**

# Example: Linux file system



/prj/literature

/prj/results ⟶ /prj/results/20220627_assembly

/prj/scripts

/prj ⟶ /prj/data ⟶ /prj/data/20220627_seq1_3dX9
/prj/data/20220627_seq2_eQ7p

/ ⟶ /home ⟶ /home/ubuntu

/mnt ⟶ /mnt/vol1
/mnt/vol2

# Naming files

- File names should be informative

- Never use 'new' or 'final'; use version numbers instead

- Use date as file name prefix (e.g. 2022-06-27 or 20220627)

- Never use spaces in file or folder names (underscore or minus as replacement)

- Include your initials as suffix in collaborative projects

# Documentation

- Document as much as possible

- Others must be able to repeat your experiments/analyses

- Document dates of data acquisition and processing

- Document all steps of data processing

- Document versions and parameters of applied tools

# de.NBI cloud

- Virtual machine (VM) for data analysis

- Accounts are required for access
  (ORCID or TUBS SSO for login)

- Addition to project required for access

- User create pair of private and public
  keys for authentication

# Transferring files

- Filezilla: graphical user interface for file transfer protocols
    - https://filezilla-project.org/

- Scp (secure copy): command line file transfer method

- Wget: command line file transfer method
    - https://ftp.gnu.org/old-gnu/Manuals/wget-1.8.1/html_mono/wget.html

- Rsync: sophisticated file transfer method that avoids redundant transfers
    - https://wiki.ubuntuusers.de/rsync/

# Linux system structure

- Linux (Ubuntu) is operation system of choice for bioinformatics

- Hierarchical structure of directory ('/' is basis)

- Separation of tool installation and data sets

- File naming conventions:
  - Never use spaces in file or directory names
  - Include dates in file names (year-month-day)

- Indicating commands with prefix '$'

/vol/data/

/vol/

/

/home/ → /home/bin/

/home/scripts/

/home/downloads/

Technische
Universität
Braunschweig

# Permissions

- Files can have different permissions:
  - Read (r)
  - Write (w)
  - Execute (x)

- Users have full permissions to edit their own files

- Downloaded files are usually not executable without adjustment

- chmod XXX <FILE_NAME> … can be run to change file permissions

Technische
Universität
Braunschweig

# Linux introduction (1)

Connection to virtual machine (VM):

```
$ ssh -i /path/to/private_key ubuntu@123.133.7.49 -p 1234
(base) ubuntu@agilezuse-10552:~$
```

'$' is used to indicate that following text needs to go into terminal

'#' indicates comment (should not be transferred into terminal)

Frequent issues:
1) Path to private key file not correct
2) Private key file too public

# Linux introduction (2)

- Moving through the folder structure:
  - $ cd /full/path/to/folder    #change into specific folder
  - $ cd subfolder  #change into subfolder
  - $ cd ..      #change into parent directory

- Checking content of a folder:
  - $ ls  #shows content of current folder
  - $ ls -lh     #shows more details
  - -l triggers display of additional details
  - -h human readable
  - -a show also hidden files

# Running Python scripts

- Run script or open argument infos:
    - $ python <SCRIPT_NAME>

- Run script with arguments:
    - `$ python <SCRIPT_NAME> --argument1_name <ARGUMENT1> --argument2_name— <ARGUMENT2>`

- Scripts show help message if started with insufficient arguments

```
                           python3 ./KIPEs3.py --baits ./flavonoid_baits/ --positions
./flavonoid_residues/ --out ./ --subject ./croton_red.fasta --seqtype pep --scoreratio 0.3 --simcut
40.0 --minsim 0.4 --minres 0.0 --minreg 0.0 --possibilities 3 --cpus 1
```

Technische
Universität
Braunschweig

# Running other tools

- Show help message:
  - $ <NAME_OF_TOOL>
  - $ <NAME_OF_TOOL> -h
  - $ <NAME_OF_TOOL> -- help

- Providing arguments is different for each tool:

- Most tools show help message if provided with insufficient/wrong arguments

Technische
Universität
Braunschweig

# Time for questions!

# Questions

1. What do you know about the growth of (sequence) databases?
2. What is Data Literacy?
3. Why is Data Literacy important?
4. Which properties has a hypothesis?
5. What are important considerations when naming files?
6. What needs to be documented?
7. What are advantages of an ELN?
8. What is LIMS?
9. What characterizes lab 4.0?

Technische
Universität
Braunschweig