# GE32/MM12 - Data Reuse

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

# Availability of slides

- All materials are freely available (CC BY) - after the lectures:
  - StudIP: **GE32/MM12**
  - GitHub: https://github.com/bpucker/teaching

- Questions: Feel free to ask at any time

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

Technische
Universität
Braunschweig

# What are the advantages of re-use?

# What are the advantages of re-use?

- Cost-effective

- Immediately available

- Extremely large datasets

# Challenges with data re-use?

- Lacking metadata

- Unknown details/issues

- Mislabeling possible

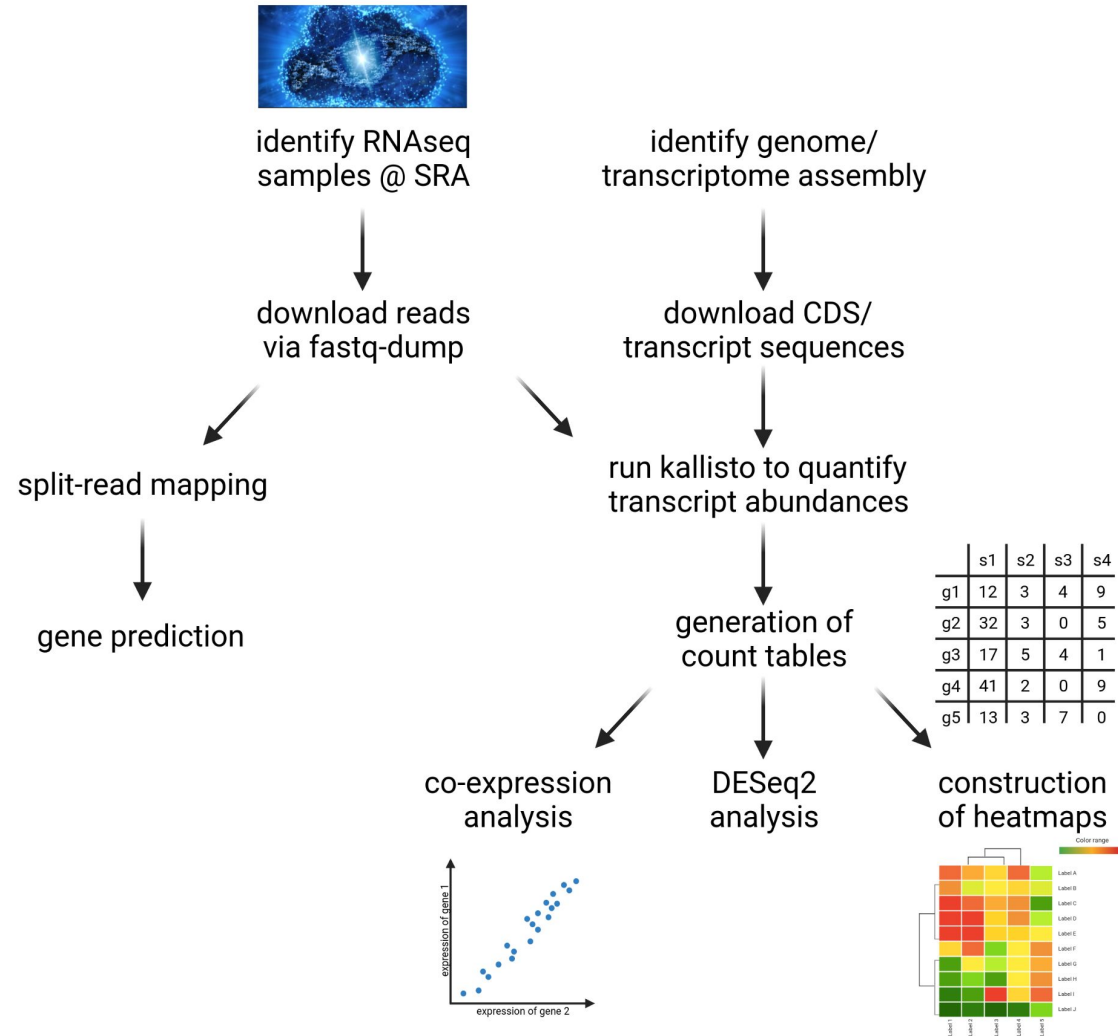- Not perfectly matching needs

- Technology outdated

# Find gene expression data sets

- SRA read selector

- Gene Expression Omnibus (GEO) search

- Publications: data availability statements & supplementary tables

# How to retrieve data?

- Preprocessed data sets (count tables @ GEO)
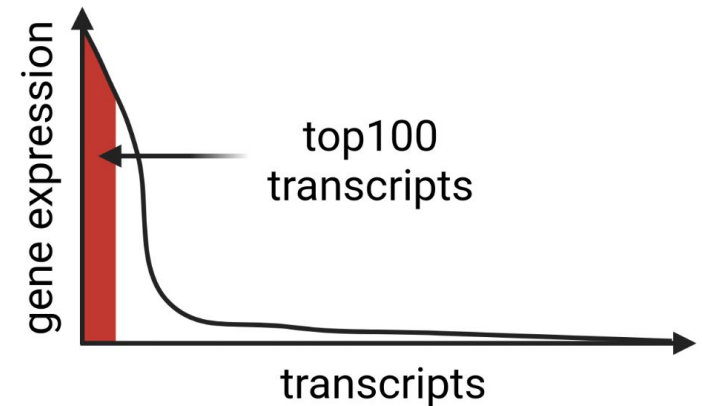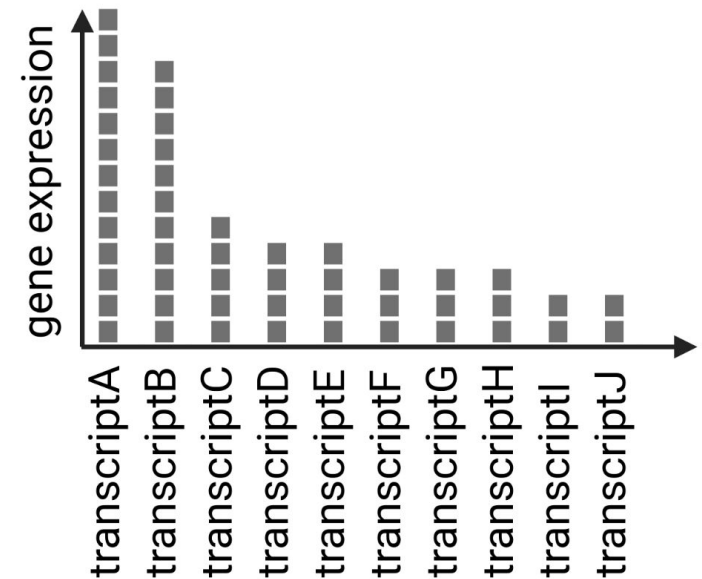
- Cloud solutions (Galaxy)

- Fastq-dump

# Workflow

# How to check RNA-seq data sets?

# RNA-seq quality control

- Percentage of reads mapped to individual mRNA sequences

- Distribution of abundance across transcripts
  - Substantial coverage of top100 transcripts

- Metadata assessment / marker gene check
  - Highly expressed marker genes like RuBisCO



top100 transcripts

**Technische Universität Braunschweig**

# How to analyze the distribution of species?

# GBIF

- Sources have heterogeneous quality

- Extensive filtering required

- What types of issues can be expected?

# GBIF

- Coordinates located in the sea

- Coordinates located in zoos/botanical gardens

- Coordinates located at the center of grids/in capitals

- Entries of fossils / too old entries

- Entries containing likely typos

# How to cite data sets?

# 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/

# DataCite

- Central service for generation and
  management of DOIs for research data sets

- Enable connection and reuse of data sets

Technische
Universität
Braunschweig

# Managing references

- Entries about publications / data sets are stored in a local database

- Convenient citation when writing manuscripts

- Support for comments and key words assigned to entries

- Examples:
  - Zotero: free
  - Mendeley: free, but belongs to Elsevier
  - Citavi: commercial, but campus license
  - EndNote: commercial

Technische
Universität
Braunschweig

# Data reuse examples

# Benchmarking of NOVOPlasty

- Public sequence read data sets are used for plastome assemblies

- Data sets can be selected based on specific criteria

- Pure bioinformatics groups can work with real data sets

- No costs for generation of data sets

Technische
Universität
Braunschweig

# Benchmarking SANPolyA

- SANPolyA detects Poly(A) signals through deep learning

- Comparison of SANPolyA results against results of other tools

- Freely available output of tools is required for systematic comparison

Technische
Universität
Braunschweig

# Pangenome of hexaploid bread wheat

- Pangenome = combination of all genomes of a species

- Integration of public data sets to identify presence/absence of genes in wheat cultivars

- Power of pangenomes increases with number of analyzed

Technische
Universität
Braunschweig

# Single plant GWAS + bulk segregant analysis

- spGWAS = single plant genome-wide association studies

- Bulk segregant analysis = comparisons of plant groups with contrasting phenotypes

- Comparison of different approaches to identify genetic basis of a trait (plant height)

- Evaluation of findings against previous reports

Technische
Universität
Braunschweig

# Co-expression gene network analysis

- Genes belonging to the same pathway show similar expression patterns

- Bamboo genome sequence was integrated with RNA-seq data sets

- Identification of co-expression modules associated with development

**Technische
Universität
Braunschweig**

# Integration of GWAS and co-expression analysis

- GWAS allows the identification of genomic regions associated with a trait

- Co-expression analysis can help to identify individual genes in these regions

- RNA-seq data sets of most species are available for free

Technische
Universität
Braunschweig

# Identification of RT-qPCR reference genes

- RNA-seq data sets can be analyzed to identify constantly expressed genes

- Reference genes with constant expression across samples is crucial for RT-qPCRs

- Reference genes can be specific for certain conditions/tissues



Hruz et al., 2011: 10.1186/1471-2164-12-156

# Electronic Fluorescent Pictograph browser

- Integration of existing microarray datasets

- Option to integrate other large scale datasets

- Basis for hypothesis generation

- Visualization of gene expression in many
  different plant tissues/conditions

Technische
Universität
Braunschweig

# Identification of conserved amino acid residues

- Identification of orthologous sequences in hundreds of species

- Comparison of sequences to identify highly conserved amino acid residues

- 3D modeling based on known structures of homologs



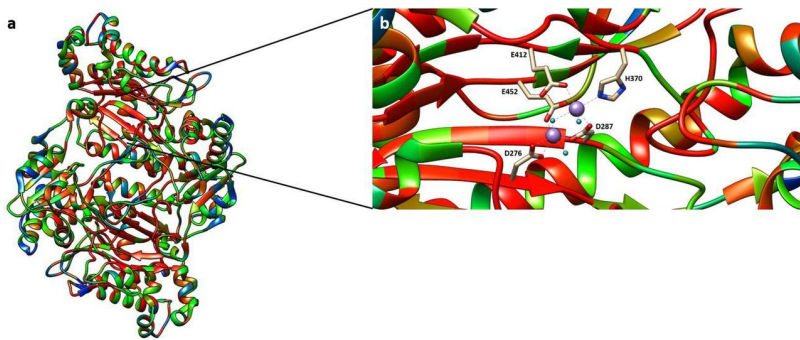| | Animals | Plants | Fungi | Archaea | Bacteria |
|---|---|---|---|---|---|
| D276 | 94 | 99 | 100 | 100 | 100 |
| D287 | 94 | 98 | 100 | 100 | 100 |
| H370 | 94 | 98 | 100 | 100 | 100 |
| E412 | 94 | 96 | 100 | 100 | 100 |
| E452 | 91 | 97 | 100 | 100 | 100 |
| T289 | 94 | 97 | 100 | 100 | 97 |
| T410 | 93 | 96 | 100 | 79 | 100 |
| H377 | 94 | 98 | 100 | 100 | 97 |
| R398 | 93 | 98 | 89 | 10 | 57 |
| W107 | 88 | 98 | 96 | 0 | 96 |
| Y241 | 94 | 96 | 100 | 2 | 90 |
| I244 | 93 | 98 | 97 | 88 | 100 |
| H255 | 94 | 98 | 100 | 100 | 100 |
| V376 | 89 | 1 | 38 | 81 | 94 |
| C58 | 58 | 64 | 0 | 0 | 0 |
| C158 | 40 | 1 | 0 | 0 | 0 |

# VANESA: enzyme properties for model

- Pathway structure taken from KEGG

- Integration of enzymatic properties ($K_M$, $V_{max}$) from BRENDA

- Freely accessible metabolic modeling solution

# Inference of metabolic pathways from metabolite data
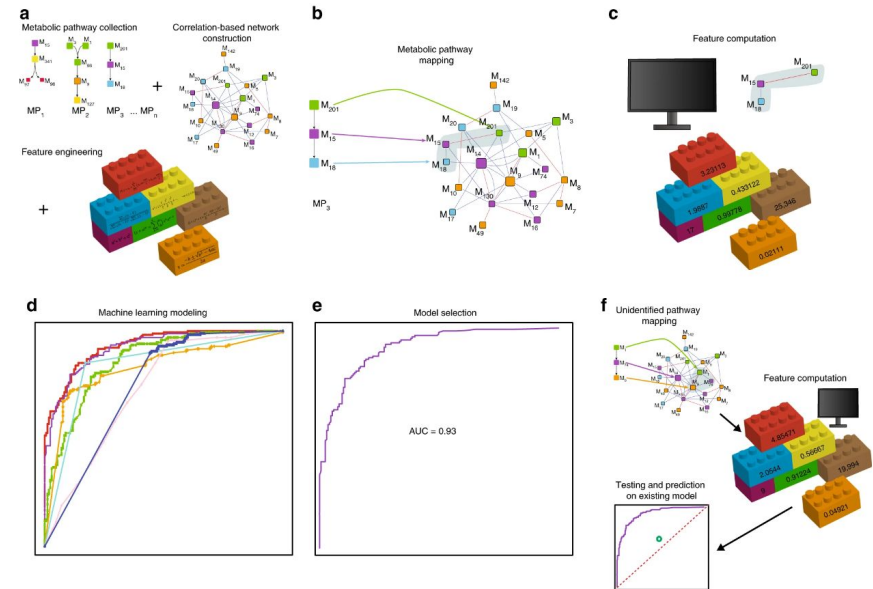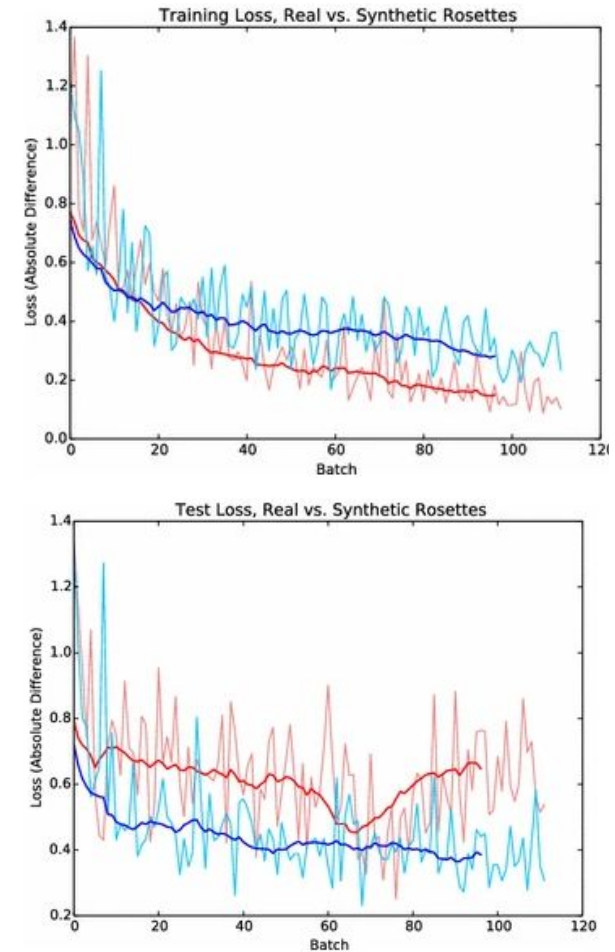
- Correlation-based inference of metabolic pathways from metabolite data sets

- Machine learning was deployed

- Identification of the most important features



Combined correlation-based network analysis and machine learning workflow. The workflow of the current study: **a** Metabolic pathways were gathered from existing repositories. In parallel, correlation-based networks of metabolites were constructed for the tissue of the organism of interest (here, the tomato pericarp). In addition, a vector of features was engineered based on network properties. **b** Metabolic pathways with partial to full coverage in the correlation networks were mapped to the networks. Each pathway was considered as a single instance. Training and test sets were proposed based on the existence of the pathways in the tomato. **c** A set of features was computed for each instance in the training set (for the current study 148 * 3 networks = 444 features in total). **d** The training set was used to generate different ML models. **e** The model that generated the best performance measures (the AUC) was selected. The ML model was validated in silico using cross-validation. **f** Test set instances were mapped onto the networks with subsequent feature computation. The proposed ML model was used to predict the potential existence of unidentified pathways in the tomato pericarp
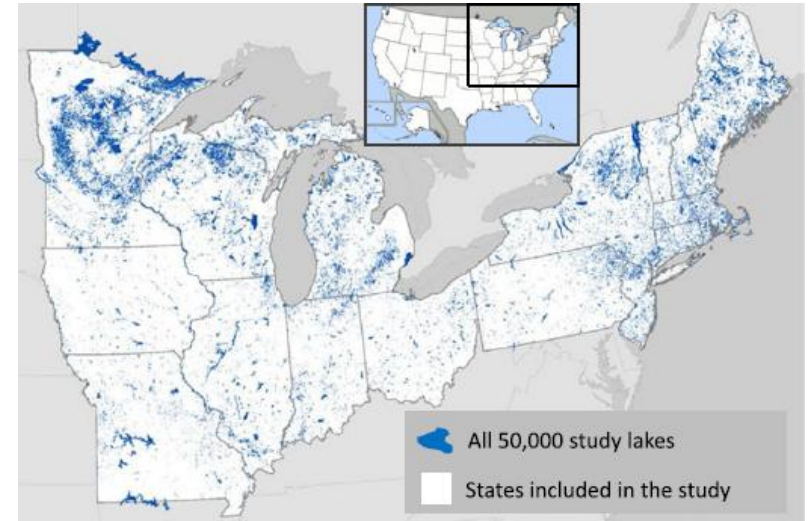
# Leaf counting

- Automatic leaf counting in rosette plants

- Deep convolutional neural networks require large and diverse data sets for optimization

- Synthetic plants are proposed as an alternative



Ubbens, et al., 2018: 10.1186/s13007-018-0273-z

Technische
Universität
Braunschweig

# LAGOS: Ecology database

- LAGOS = LAke multi-scaled GeOSpatial and temporal database

- Combination of site-based ecosystem datasets with national geospatial datasets

- Best practice example for lake related data sets covering about 50,000 lakes

- 100 individual datasets about water quality



All 50,000 study lakes

States included in the study

Technische
Universität
Braunschweig

# Time for questions!

# Questions

1. What are the advantages/disadvantages of reuse?
2. How can you assess the quality of RNA-seq data sets?
3. How to filter GBIF data sets?
4. How can you cite data sets?
5. Which studies used existing data sets to gain novel insights?