

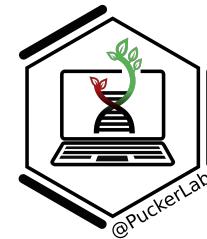
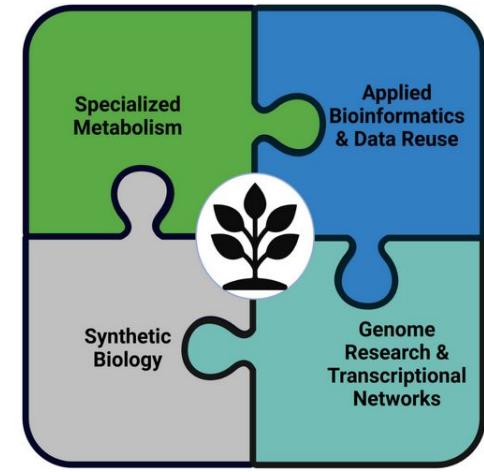
Prof. Dr. Boas Pucker

# **Introduction to Plant Bioinformatics**

# Round of introduction

- Name?
- Study program?
- Semester?
- Previous experiences with biochemistry?
- Previous experiences with bioinformatics?
- Expectations?

- Biochemistry at HHU Düsseldorf
- (Systems) Biology at Bielefeld University
- Doctoral student (CeBiTec, Bielefeld University)
  - Genomics & Bioinformatics; synthetic biology (iGEM)
- Post doc (Ruhr-University Bochum)
- Post doc (Department of Plant Sciences, Cambridge, UK)
- Plant Biotechnology & Bioinformatics, TU Braunschweig (since 2021)
  - Specialized plant metabolites, applied bioinformatics



# Availability of slides

- All materials are freely available (CC BY) - after the lectures:
  - eCampus: WBIO-A-08
  - GitHub: <https://github.com/bpucker/teaching>
- Questions: Feel free to ask at any time
- Feedback, comments, or questions: [pucker\[a\]uni-bonn.de](mailto:pucker[a]uni-bonn.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.

# Overview

UNIVERSITÄT BONN

- 1) Introduction to bioinformatics (today)
- 2) Identification of biosynthesis pathways
- 3) Pathway databases
- 4) Gene expression & co-expression analyses
- 5) Phylogenetic analyses
- 6) Synteny & biosynthetic gene clusters
- 7) Metabolic flux & modeling
- 8) Repetition, questions & quiz

# DOI - finding the literature

- Suggestions for detailed literature will be included on slides
  - DOI = Digital Object Identifier
  - Unique and short way to point to a publication
  - How to resolve a DOI? <https://dx.doi.org/>



bioRxiv posts many COVID19-related papers. A reminder: they have not been formally peer-reviewed and should not guide health-related behavior or be reported in the press as conclusive.

New Results

**Apicomplexan FMS1 originated from F9H through tandem gene duplication**

• **Boris Pukac, Massimo Iezzoni**  
doi: <https://doi.org/10.1101/2022.06.48.50750>  
bioRxiv preprint doi: <https://doi.org/10.1101/2022.06.27.520000>; this version posted June 27, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

                                                                                                                                                                                   <img alt="bioRxiv icon" data-bbox="111 66

## Resolve a DOI Name

doi: 10.1101/2022.02.16.480750

Go

Type or paste a DOI name into the text box. Click Go. Your browser will take you to a Web page (URL) associated with that DOI name.

Send questions or comments to doi-help@doi.org

[Further documentation is available here](#)

DOI Resolution Documentation

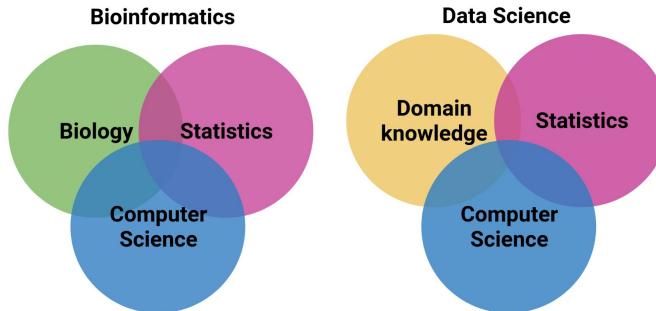
 doi:10.5281/pelagic-ecology-and-the-DOE® environmental assessment of the International DOE-Ecosystems

# What is bioinformatics?



# What is bioinformatics?

- Subdiscipline of biology, statistics, and computer science
- Acquisition, storage, analysis, and dissemination of biological data
- Bioinformatics is a biology-specific data science version

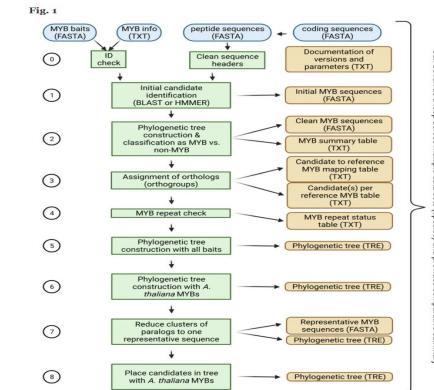
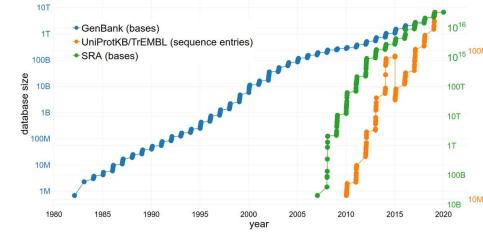
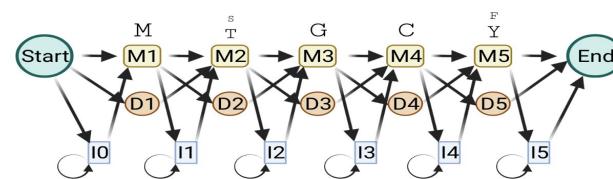


# Why do we need bioinformatics?

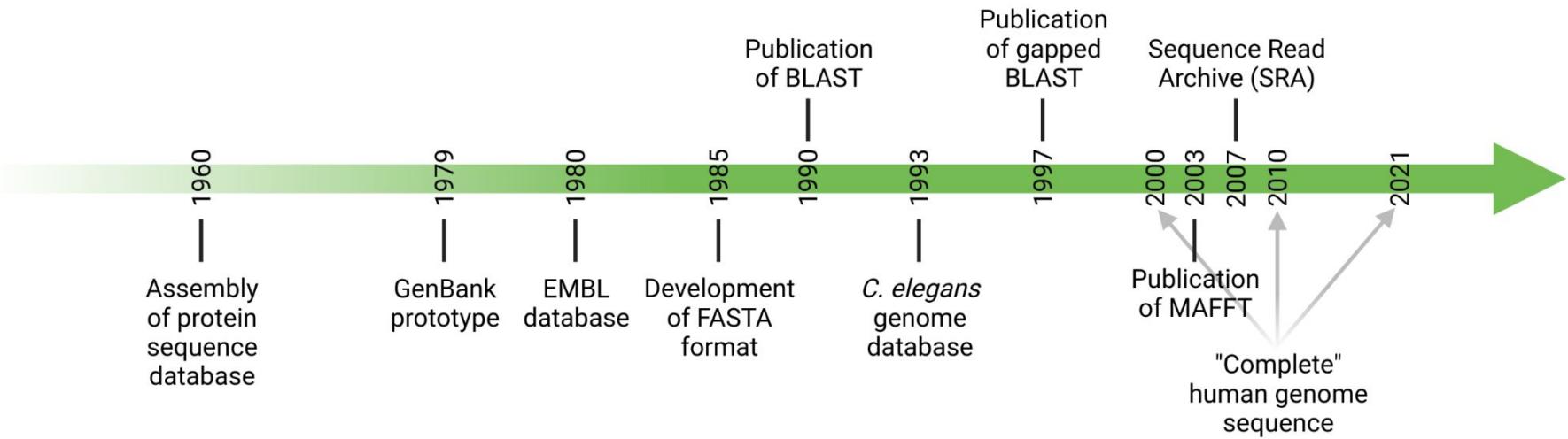
?

# Why do we need bioinformatics?

- Large data sets
- Complex models
- Automatic analyses



# History of bioinformatics

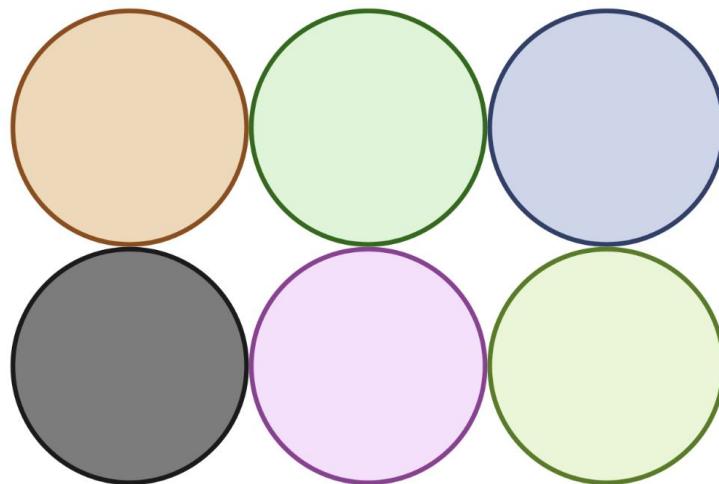


- Experience with ChatGPT (OpenAI)?
- Do you see potential and risk?
- AI does not only work for text, but also for code, images, movies, games, gene prediction, pattern recognition, .... but there are limitations!!!

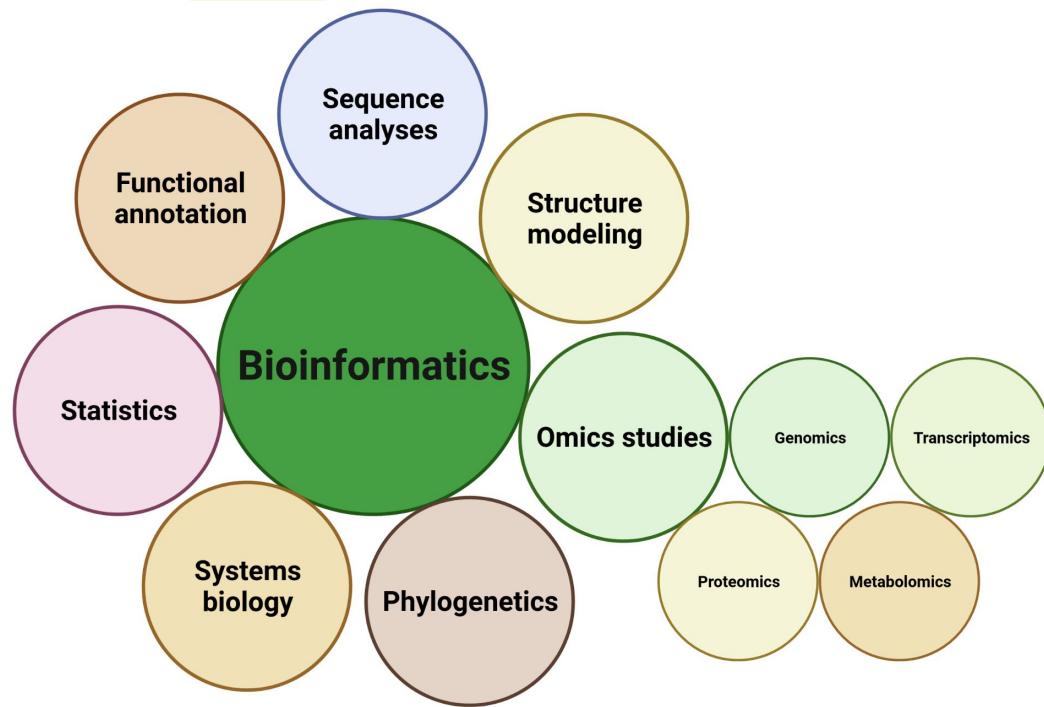
Future Reading & Related Work (selected publications from the Pucker Lab):

1. **Pucker, B. et al. (2024).** "A comprehensive survey of anthocyanin pathway disruptions in diverse plant species." *BMC Plant Biology*, 24, 5316. <https://doi.org/10.1186/s12870-024-05316-w>
2. **Pucker, B. et al. (2021).** "Automatic identification of genes with loss-of-function mutations in plant pigment biosynthesis pathways." *Genes*, 12(3), 452.
3. **Pucker, B., & Brockington, S. F. (2022).** "Insights into the evolution of pigment pathways from comparative genomics of Caryophyllales." *Frontiers in Plant Science*, 13, 870263.

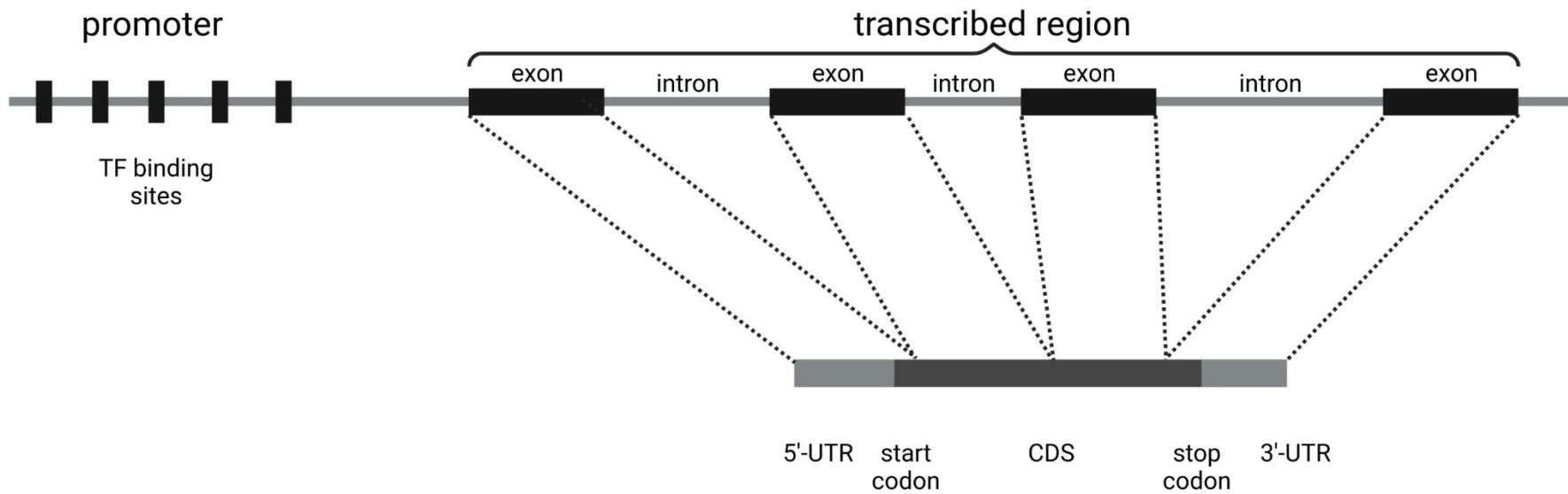
# What are the different (sub)fields of bioinformatics?



# What are the different (sub)fields of bioinformatics?



# Sequence analysis



# Sequence analysis II

- Comparison of sequences to find similarities
- Analysis of sequence composition
- Important methods:
  - BLAST = Basic Local Alignment Search Tool
  - MAFFT = Multiple Alignment using Fast Fourier Transform

## chalcone synthase [Chenopodium quinoa]

Sequence ID: [XP\\_021762431.1](#) Length: 392 Number of Matches: 1

Range 1: 2 to 389	<a href="#">GenPept</a>	<a href="#">Graphics</a>	<a href="#">▼ Next Match</a>	<a href="#">▲ Previous Match</a>
Score	Expect	Method	Identites	Positives
.666 bits(1719)	0.0	Compositional matrix adjust.	328(83%)	362/388(93%)
				/0/388(0%)
Query 3	TPSVQEIRDQASNPGPATIAGTATPANEWYOAEPYDFYFRVTKEHMKELKKFKRMC			
Sbjct 2	T+S+TEIR AQR++GPAFITAGTA +Y0++PD+YFRVTKEHMKELKKFKRMC			
Query 63	DNSMIKKRYMHVTEELLEENPHLCDFNASSLDTRODILATEVPKLGEAAKAIKEWGP			
Sbjct 62	D SHI KRYMH+TEE L+ENP+ +SLLDTROD+ +EV+LGEAAKAIKEWGP			
Query 123	RSKITHWIFCTTSQGVMPGRPSVKRFMLYQOGCYAGGTVLRLAKDIAEN			
Sbjct 122	+SKITHWI CTTSQGVMPGRADYOLTLKLLGLRPSV+RFLMYQOGC+AGGTVLRLAKD+AEN			
Query 183	NRGARLVWCAETTTICPRGPTQJHLDGMVGOAEdanqanav/vgt/PPESTERPTFOLV			
Sbjct 182	NRGARLVWCAETT TICPRGPT+ HLDGMVGOAEd+GAGA+IVGADPDEIERPLFKMV			
Query 243	WAAQTILPGSEGADGHHLREVGTLTFHLLKDVPGLISKNIKEALVEAFPOIGIDWNNSIFW			
Sbjct 242	WAAQTILP SEGADGHHLREVGTLTFHLLKDVPGLISKNIKEALVEAFPLGIDWNNSIFW			
Query 303	VAHPGGRAILDDVESKLGKEDKLKTRHVLSEYGNMSSACVLFILDEMRRKAMKEGAT			
Sbjct 302	+AHGG ALD VE+KLGKE-KL TR+VLSF-GMNSSACVLFILDEMRRK++MKEG AT			
Query 363	TGEGLENGVLFGFPGLTVETVMLHSVP	390		
Sbjct 362	TGDDLDNGVLFGFPGLTVETVVLHSVP	389		

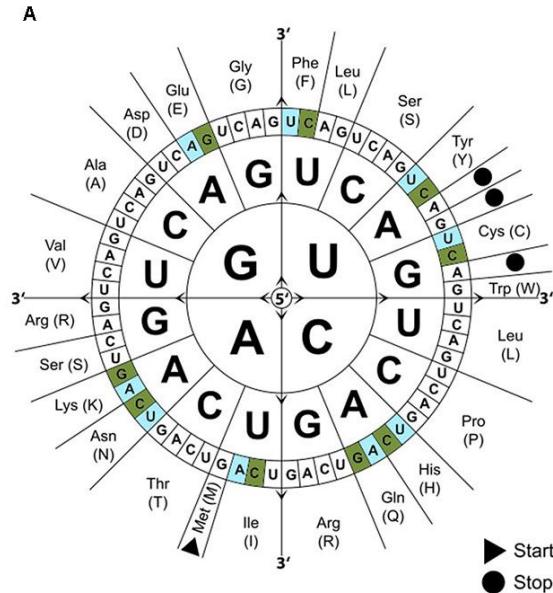
## CLUSTAL format alignment by MAFFT FFT-NS-i (v7.487)

```

BvCHS          M---ATPSVQEIRDQASNPGPATIAGTATPANEWYOAEPYDFYFRVTKEHMKELKK
AtCHS          MMAGASSLDEIRQAORADGPAGILAIGTANPENHVLOAEYPDYFRITNSEHMTDLKEK
* .:;:****:***:*** ***** *.:*****:***:*****:***:*
FKRMDKNSMIKKRYMHVTEELLEENPHLCDFNASSLDTRODILATEVPKLGEAAVKAIK
FKRMDKSTIRKRRMHMELTEEFKENPHMCAYMAPSLDTRODIZVVVEPKLGEAAVKAIK
*****:*****:*****:*****:*****:*****:*****:*****:*****:*****:*****
EWGOPRSKITHWIFCTTSQGVMPGRADYOLTLKLLGLRPSVKRFLMYQOGCYAGGTVLRLAK
EWGOPRSKITHWIFCTTSQGVMPGRADYOLTLKLLGLRPSVKRFLMMYQQGCFAAGTVLRIAK
*****:*****:*****:*****:*****:*****:*****:*****:*****:*****:*****
DIAENNRGARLVWCAETTICPRGPTQIHLDSMVQALFGGGAGAGAVIVGADPDESI-ER
DIAENNRGARLVWCAETTAVTFRGPSPDHLDLSVGQALFSGAIDHLIVGSDPDTSGEK
.*****:*****:*****:*****:*****:*****:*****:*****:*****:*****:*****
PIFOLWVAAAQTILPGSEGADGHHLREVGTLTFHLLKDVPGLISKNIKEALVEAFQPIGIDD
PIFEMVWAAAQTILPDSDGAIDGHHLREVGTLTFHLLKDVPGLISKNIVKSLDEAFKPLGJSD
*****:*****:*****:*****:*****:*****:*****:*****:*****:*****:*****
WNSIFWVVAHPGGRAILDDVESKLGKEDKLKTRHVLSEYGNMSSACVLFILDEMRRKAM
WNSLFWIAHPPGPAILDVEIKLGLKEEMRATRHVLSEYGNMSSACVLFILDEMRRKSA
*****:*****:*****:*****:*****:*****:*****:*****:*****:*****:*****
KEGMATTGEGLEWGVLFGLFGFPGLTVETVMLHSVPAN
KGDVATTGEGLEWGVLFGLFGFPGLTVETVVLHSVP--
```

# Sequence analysis III

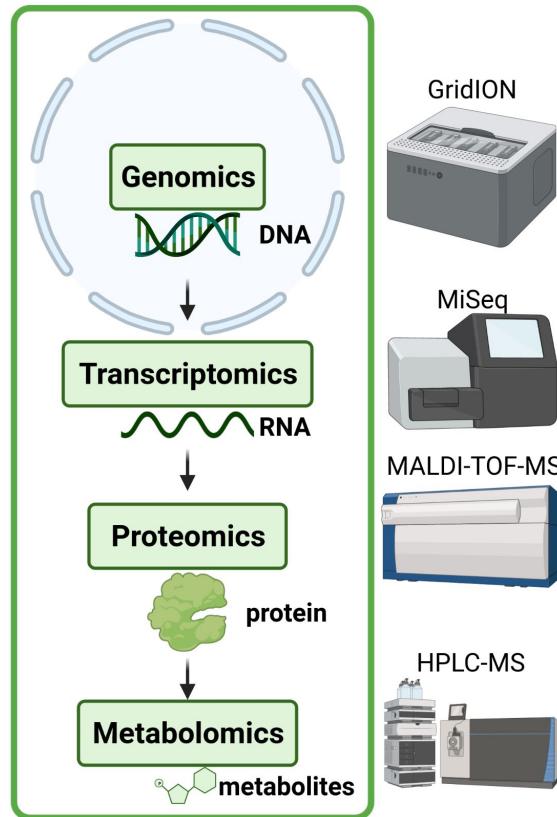
- Codons define amino acids to integrate
- Genetic code is redundant i.e. multiple codons for each amino acid
- Organisms can have a preference for certain codons for a given amino acid
- Optimization of sequences for heterologous expression



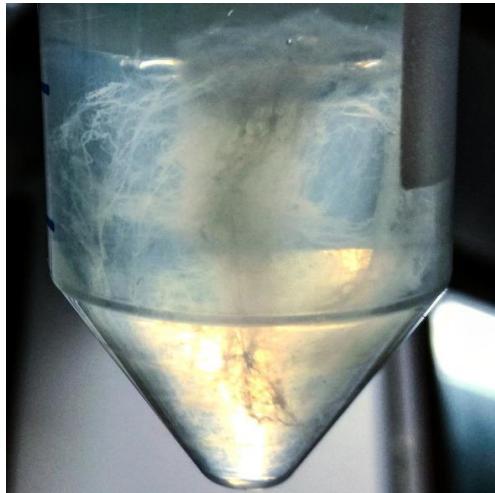
**B**

Amino acid	Codon	<i>D. patens</i>	<i>A. thaliana</i>	<i>S. cerevisiae</i>	<i>S. pombe</i>	<i>H. Sapiens</i>
Cys	UGC	++	-	+	+	++
	UGU	-	-	+	-	-
Glu	GAA	-	-	+	-	-
	GAG	++	-	+	+	++
Phe	UUC	++	++	++	++	++
	UUU	-	-	-	-	-
His	CAC	++	++	++	++	++
	CAU	-	-	-	-	-
Ile	AUA	-	-	/	-	-
	AUC	++	/	-	-	+
	AUU	-	-	-	-	-
Lys	AAA	-	-	-	-	-
	AAG	++	++	++	++	++
Asn	AAC	++	++	++	++	++
	AAU	-	-	-	-	-
Gln	CAA	-	-	-	+	-
	CAG	++	++	-	+	-
Tyr	UAC	++	++	++	++	++
	UAU	-	-	-	-	-

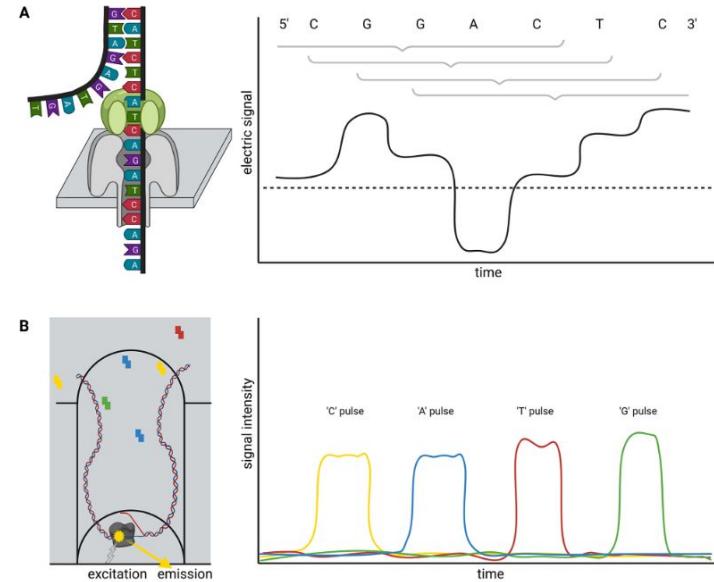
# Omics levels



# What is a genome?

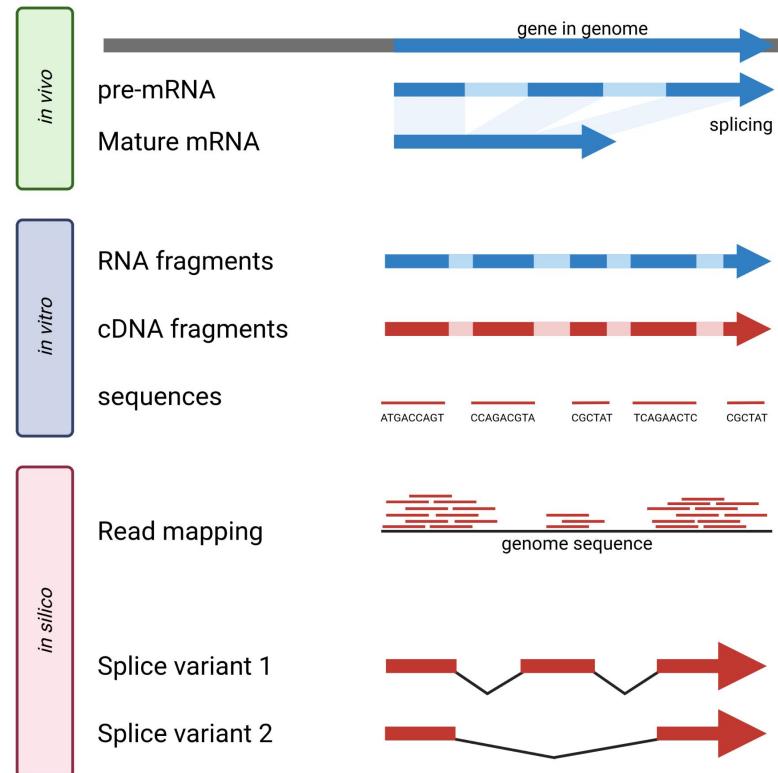


- Genome = sum all of genetic information in a cell
- Plants have DNA in nucleus (nucleome), chloroplasts (plastome), and mitochondria (chondromes)
- Sequencing plant genomes: Oxford Nanopore Technologies (ONT) and Pacific Biosciences (PacBio)
- Genome sequence assembly: HiCanu, FALCON, Shasta, NextDenovo2 ...
- Comparison of genome sequences
- Handling of large data sets (>1TB) necessary for most projects

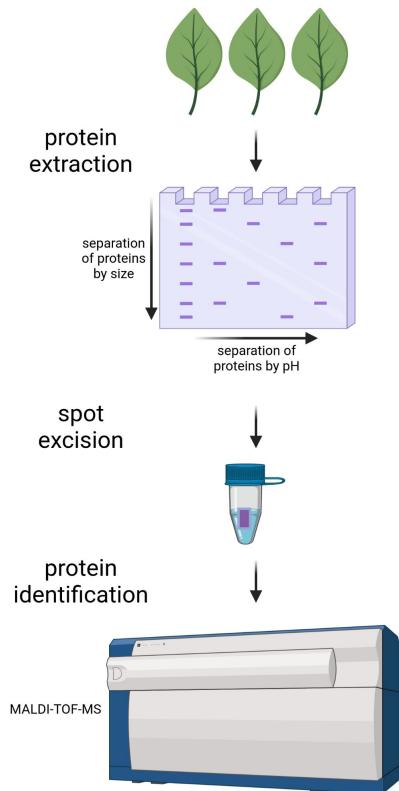


# Transcriptomics

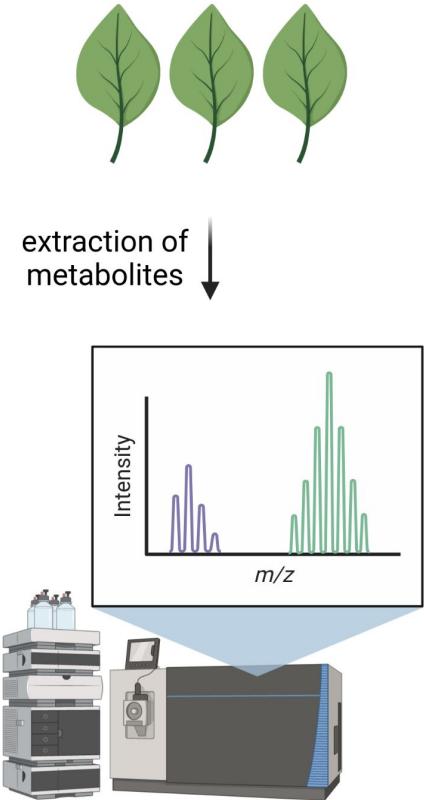
- Transcriptome = sum of transcripts in a cell/tissue under defined conditions and at a defined time point
- Systematic investigation of gene expression:
  - Microarrays: method of choice for many years
  - RNA-Seq: current method of choice to study gene expression systematically
- Comparison of different samples (tissues, conditions, genotypes, ....)



- Proteome = sum of all proteins in a cell/tissue under defined conditions at a defined time point
- Methods for systematic investigation:
  - High Performance Liquid Chromatography (HPLC)-Mass Spectrometry (MS)
  - Separation of proteins via 2D gels and investigation of spots by MS
- Heterogenous properties of proteins make systematic analysis challenging

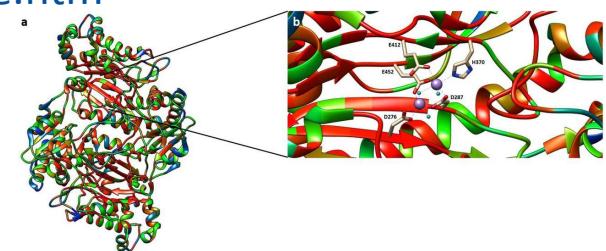


- Metabolom = sum all of metabolites in a cell/tissue under defined conditions at a defined time point
- Methods for systematic metabolite analysis:
  - HPLC-MS/MS
  - GC-MS/MS
- Chemical diversity of metabolites poses a challenge for the analysis
- Identification of metabolites remains a challenge



# Structure modeling

- Modeling of RNA structures
  - Collection of tools: [https://molbiol-tools.ca/RNA\\_analysis.htm](https://molbiol-tools.ca/RNA_analysis.htm)
  - BiBiServ: <https://bibiserv.cebitec.uni-bielefeld.de/rna>
  - RNAfold:  
<http://rna.tbi.univie.ac.at/cgi-bin/RNAWebSuite/RNAlign.cgi>
- Modeling of 3D protein structures
  - Collection of tools:  
[https://molbiol-tools.ca/Protein\\_tertiary\\_structure.htm](https://molbiol-tools.ca/Protein_tertiary_structure.htm)
  - PredictProtein: <https://predictprotein.org/>
  - AlphaFold2 (10.1021/acs.jcim.1c01114)



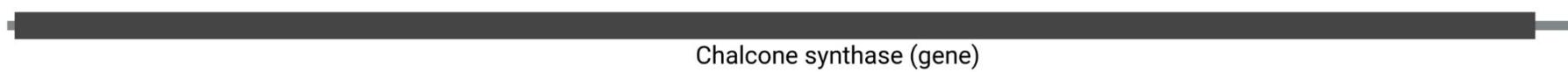
# Structural vs. functional annotation

- Structural annotation = position of features (e.g. exons) in a genome sequence
- Functional annotation = biochemical function of a feature (e.g. gene)

## Structural annotation

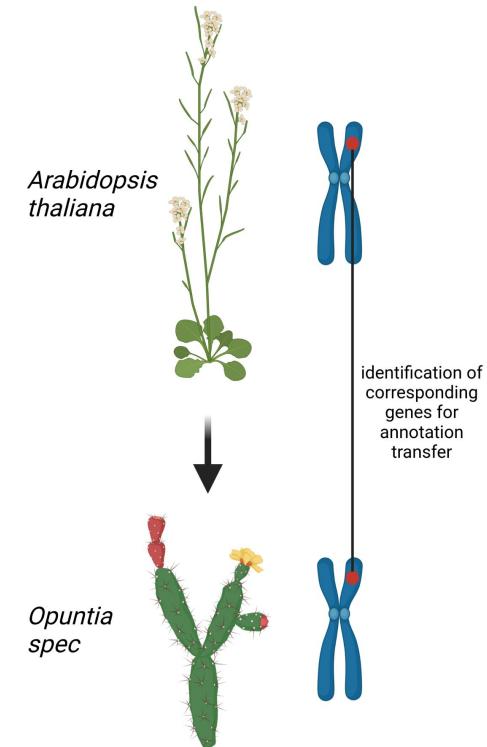


## Functional annotation

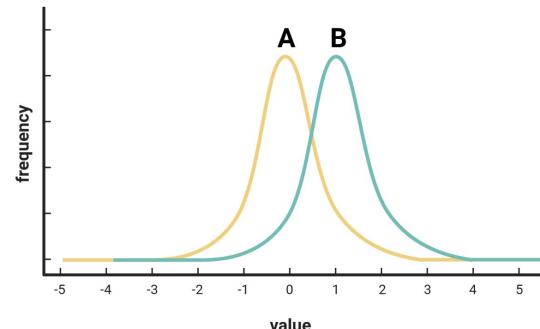
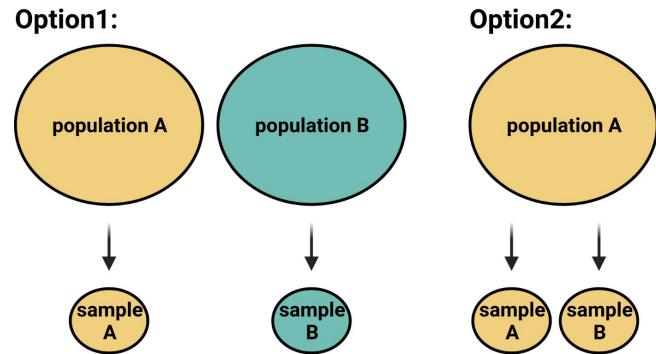


# Functional annotation

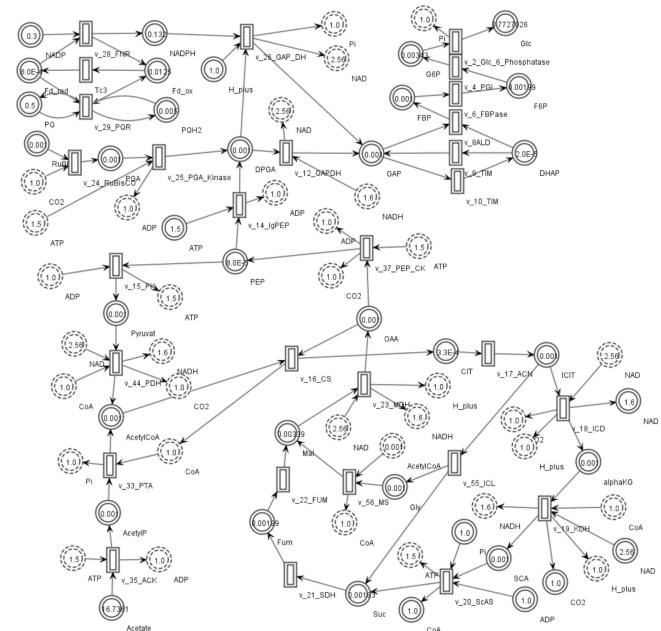
- Assumption: similar sequences (genes) have similar functions
- Sequence similarity indicates origin from a shared common ancestor
- Ancestor is likely to have inherited the same function to offspring
- Transfer of functional annotation based on sequence similarity
- *Arabidopsis thaliana* gene functions are well studied and usually serve as reference



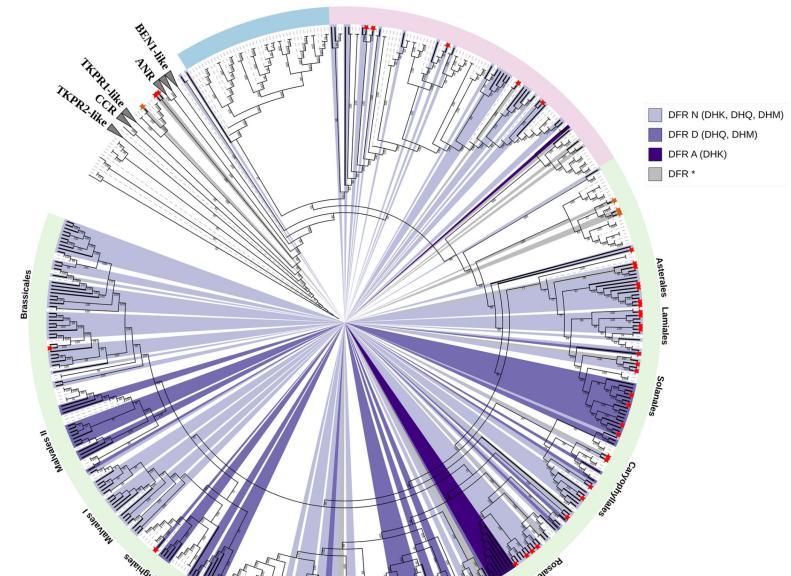
- Statistical tests to assess data sets
- Important for analysis of omics data sets
- Concept of a statistical test:
  - $H_0$  = difference of two samples can be explained by random effects
  - $H_1$  = certain factor(s) determine the difference of groups



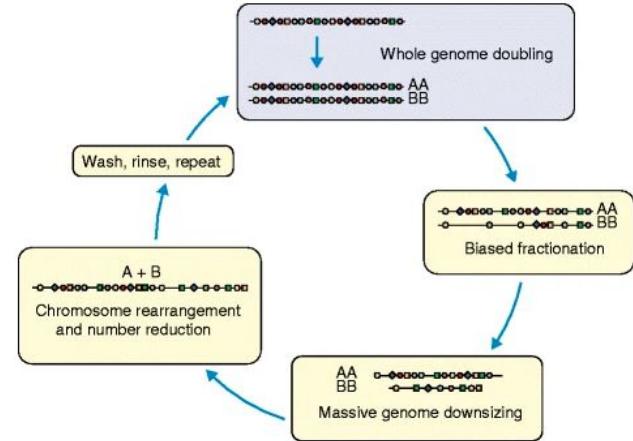
- Description of biological systems with mathematical methods
- Michaelis-Menten kinetics is basis of many metabolic models
- Integration of multiple omics data sets
- Modeling of metabolic flux through pathways
- Example: Modeling of aerobic carbon metabolism in *Chlamydomonas reinhardtii*



- Analysis of the evolutionary relationships between sequences/species
- Reconstruction of possible models to explain these relationships
- Not limited to sequence similarity, but considered (most likely) ancestral stages



- Genomes evolve over time
- Whole genome duplication is crucial source for new sequences in plants
- Polyploidy = multiple copies of all chromosomes
- Transposable Elements (TEs) can increase genome size
- Huge differences in chromosome numbers and genome sizes



**INSERT NANCY's FIGURE**

# Do you know any bioinformatics tools?

Required  
subject

 Choose File No file chosen

seqtype

 pep

Optional  
scoreratio

 0.3

simcut

 0.4

minsim

 0.4

minres

 0.0

minreg

 0.0

possibilities

 3

Privacy

I would like to be notified when my job status changes. Please contact me via the following email address:

\_\_\_\_\_

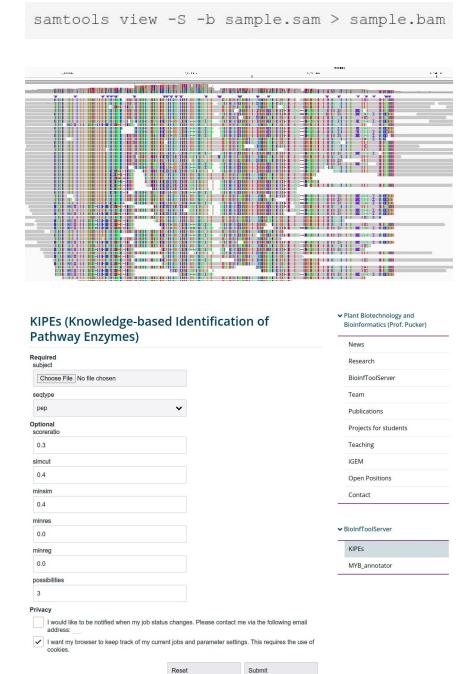
I want my browser to keep track of my current jobs and parameter settings. This requires the use of cookies.

Reset

Submit

# Types of bioinformatics tools

- Command line: require a basic understanding of Linux, but allow the processing of large data sets
  - Example: long read-based t-DNA analysis (loreta);  
<https://doi.org/10.1186/s12864-021-07877-8>
- Graphical user interface (GUI): easier to use, but usually restricted to smaller data sets
  - Example: Integrated Genome Viewer (IGV),  
<https://igv.org/>
- Web services: easy to use, but usually restricted to small data sets
  - Example: Knowledge-Based Identification of Pathway Enzymes (KIPES); <http://ppb-tools.de/KIPES>



Command line input: samtools view -S -b sample.sam > sample.bam

Large genome visualization showing multiple chromosomes with colored tracks.

Navigation sidebar:

- Plant Biotechnology and Bioinformatics (Prof. Pucker)
- News
- Research
- BioinfToolServer
- Team
- Publications
- Projects for students
- Teaching
- IGM
- Open Positions
- Contact
- KIPES
- MtP\_Annotator

Form fields (from left to right):

- Required file: Choose File [No file chosen]
- Format: jpg
- Optimal scenario: 0.3
- Smooth: 0.4
- Minmax: 0.4
- Minres: 0.0
- Minring: 0.0
- Possibleities: 3
- Privacy checkboxes:
  - I would like to be notified when my job status changes. Please contact me via the following email: [input field]
  - I want my browser to keep track of my current jobs and parameter settings. This requires the use of cookies. [checked checkbox]
- Reset button
- Submit button

# How to run command line tools

- Trimmomatic:

<http://www.usadellab.org/cms/?page=trimmomatic>

```
java -jar <path to trimmomatic jar> SE [-threads <threads>] [-phred33 | -phred64] [-trimlog <logFile>] <input> <output> <step 1> ...
```

- Samtools:

<http://www.htslib.org/doc/samtools-view.html>

```
samtools view -S -b sample.sam > sample.bam
```

# Do you know any programming languages?

**string integer float  
float list dict**  
**sorted key value  
def function print  
else elif**

# Script and programming languages

- Perl: one of the first script languages used in bioinformatics (sequences)
  - Examples: AUGUSTUS and helper scripts
  - <https://www.perlfoundation.org/>
- Python: established script language for bioinformatics (sequences)
  - Examples: BUSCO, KIPEs, MYB\_annotator
  - <https://www.python.org/psf/>
- Julia: script language that is gaining attention (structures)
  - Examples: biojulia
  - <https://julialang.org/>
- R: established script language for bioinformatics (numbers/statistics)
  - Examples: DESeq2, ggplot2,
  - <https://www.r-project.org/about.html>
- Java: powerful for the analysis of large data sets
  - Integrated Genomics Viewer (IGV), SnpEff, samtools, Trimmomatic, ...
  - <https://www.java.com/de/> (commercial)

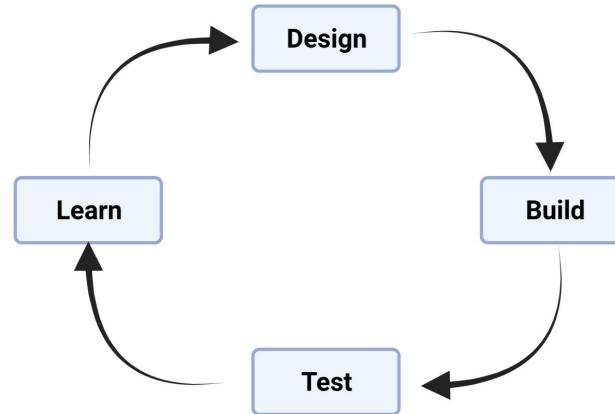


# How would you optimize a bioinformatics tool?



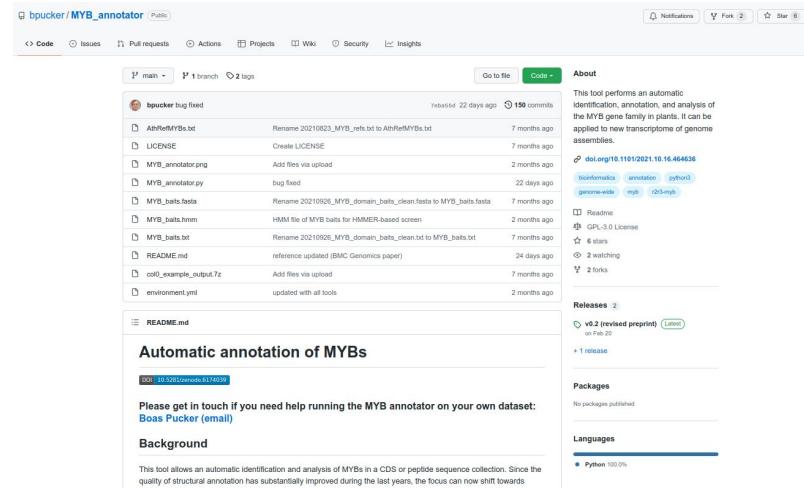
# Design-Build-Test-Learn Cycle

- Software (any product) can be improved through the DBTL cycle:
  - Design how the software could work
  - Build the software according to the design
  - Test how well the software performs
  - Learn from the testing and start again with an improved design



# Version control and repositories

- Version control is important when improving successively
- Different version control systems:
  - Git, CVS, SVN, ...
- Online repositories that allow sharing of code:
  - GitHub, Bitbucket, GitLab, SourceForge, ...



The screenshot shows a GitHub repository page for 'bpucker / MYB\_annotator'. The repository has 150 commits, 2 branches, and 2 tags. The 'Code' tab is selected. The repository description states: 'This tool performs an automatic identification, annotation, and analysis of the MYB gene family in plants. It can be applied to new transcriptome of genome assemblies.' It includes links to doi.org/10.1191/2021.10.16.464636, bioinformatics, annotation, python3, genome-wide, myb, QTL-myb, Readme, GPL-3.0 License, 6 stars, 2 watching, and 2 forks. The 'Releases' section shows v0.2 (revised preprint) and v0.1 (released 2021-06-20). The 'Packages' section shows no packages published. The 'Languages' section shows Python at 100.0%. The README.md file contains instructions for running the tool on one's own dataset.

- MIT: leanest licence (everything is possible)
- Apache: similar to MIT, but lengthy
- GPL (General Public License): ensures that derived work remains open
- BSD (Berkeley Software Distribution): similar to MIT, but more cases specified

#### MIT License

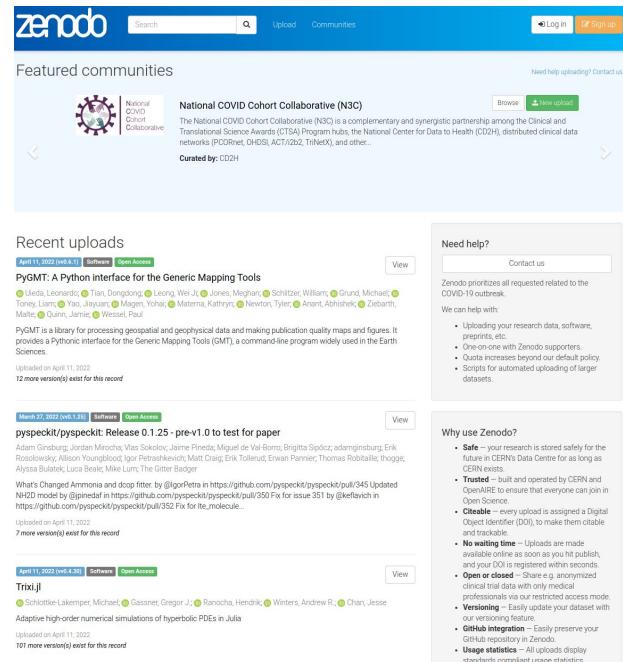
Copyright (c) [year] [fullname]

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

- Data permanently stored in CERN's data centre
- OpenAIRE allows everyone to engage in OpenScience
  - socio-technical infrastructure for scholarly communication
- DOI makes entries citeable
- Uploads are published immediately
- Integration with GitHub possible



The screenshot shows the Zenodo homepage. At the top, there is a search bar, upload and community links, and a log in/sign up button. Below the header, "Featured communities" are listed, with one entry for the "National COVID Cohort Collaborative (N3C)". This entry includes a thumbnail, the name, a brief description, and a "View" button. In the center, "Recent uploads" are displayed, with a prominent entry for "PyGMT: A Python Interface for the Generic Mapping Tools". This entry shows the date (April 11, 2022), version (v0.1.20), type (Software), and access level (Open Access). It includes a detailed description, a list of contributors, and a "View" button. To the right of the uploads, a "Need help?" section provides links to contact us and a forum, and a "Why use Zenodo?" section lists benefits such as safety, trustworthiness, citability, and versioning.

- International open-access repository
- Suitable for large data sets, but also an option for software
- CCO licence makes everything completely re-useable by other
- DOIs make entries citeable



- Bioinformatic tools often have dependencies
- Dependency = module/tool that needs to be installed already
- Conda enables installation without administrator privileges



*Package, dependency and environment management for any language—Python, R, Ruby, Lua, Scala, Java, JavaScript, C/C++, FORTRAN, and more.*

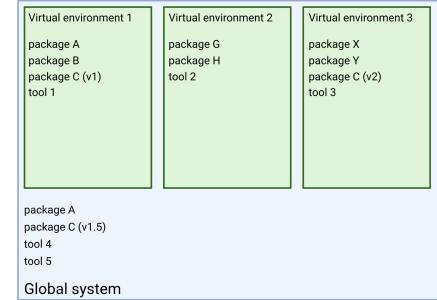
Conda is an open source package management system and environment management system that runs on Windows, macOS and Linux. Conda quickly installs, runs and updates packages and their dependencies. Conda easily creates, saves, loads and switches between environments on your local computer. It was created for Python programs, but it can package and distribute software for any language.

Conda as a package manager helps you find and install packages. If you need a package that requires a different version of Python, you do not need to switch to a different environment manager, because conda is also an environment manager. With just a few commands, you can set up a totally separate environment to run that different version of Python, while continuing to run your usual version of Python in your normal environment.

In its default configuration, conda can install and manage the thousand packages at [repo.anaconda.com](http://repo.anaconda.com) that are built, reviewed and maintained by Anaconda®.

# Virtual environment

- Virtual environment allows controlled installation of tools without interfering with system
- Tool and all dependencies are installed in a box
- Different versions of packages/modules are required for different tools
- Different virtual environments can be used for different tools



## 12.2. Creating Virtual Environments

The module used to create and manage virtual environments is called `venv`. `venv` will usually install the most recent version of Python that you have available. If you have multiple versions of Python on your system, you can select a specific Python version by running `python3` or whichever version you want.

To create a virtual environment, decide upon a directory where you want to place it, and run the `venv` module as a script with the directory path:

```
python -m venv tutorial-env
```

This will create the `tutorial-env` directory if it doesn't exist, and also create directories inside it containing a copy of the Python interpreter and various supporting files.

A common directory location for a virtual environment is `~/venv`. This name keeps the directory typically hidden in your shell and thus out of the way while giving it a name that explains why the directory exists. It also prevents clashing with `.env` environment variable definition files that some tooling supports.

Once you've created a virtual environment, you may activate it.

On Windows, run:

```
Tutorial-env\Scripts\activate.bat
```

On Unix or MacOS, run:

```
source tutorial-env/bin/activate
```

(This script is written for the bash shell. If you use the `csh` or `fish` shells, there are alternate `activate.csh` and `activate.fish` scripts you should use instead.)

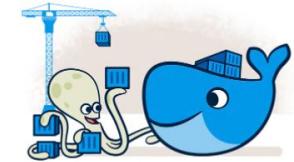
Activating the virtual environment will change your shell's prompt to show what virtual environment you're using, and modify the environment so that running `python` will get you that particular version and installation of Python. For example:

```
$ source /envs/tutorial-env/bin/activate
(tutorial-env) $ python
Python 3.5.1 (default, May  6 2016, 10:59:36)
>>> import sys
>>> sys.path
['/usr/local/lib/python35.zip', ...,
 '/envs/tutorial-env/lib/python3.5/site-packages']
>>>
```

- Solution to prepare an environment for running tools
- Build: developer of tool generates an image of the required environment
- Share: image is shared with users of the tool
- Run: user can mount the image to run a tool in a defined environment
- Alternatives: podman, OpenVZ, VirtualBox, Kubernetes, LXC, ZeroVM, ...

## Build

- Get a head start on your coding by leveraging Docker Images to efficiently develop your own unique applications on Windows and Mac. Create your multi-container application using Docker Compose.
- Integrate with your favorite tools throughout your development pipeline
  - Docker works with all development tools you use including VS Code, CircleCI and GitHub.
- Package applications as portable container images to run in any environment consistently from on-premises Kubernetes to AWS ECS, Azure ACI, Google GKE and more.



## Share

- Leverage Docker Trusted Content, including Docker Official Images and Images from Docker Verified Publishers from the Docker Hub repository.
- Innovate by collaborating with team members and other developers and by easily publishing Images to Docker Hub.
- Personalize developer access to Images with roles based access control and get insights into activity history with Docker Hub Audit Logs.



## Run

- Deliver multiple applications hassle free and have them run the same way on all your environments including design, testing, staging and production – desktop or cloud-native.
- Deploy your applications in separate containers independently and in different languages. Reduce the risk of conflict between languages, libraries or frameworks.
- Speed development with the simplicity of Docker Compose CLI and with one command, launch your applications locally and on the cloud with AWS ECS and Azure ACI.



- Graphical user interface (web-based platform) for bioinformatic workflows
- Open source enables local installation (e.g. on compute cluster)
- Supported by de.NBI, elixir, and many others

Galaxy is an open source, web-based platform for data intensive biomedical research. If you are new to Galaxy start here or consult our help resources. You can install your own Galaxy by following the tutorial and choose from thousands of tools from the Tool Shed.



The logo for the Galaxy Community Conference (GCC) 2022 in Minneapolis. It features a stylized city skyline silhouette with a cherry on the right side. The text "GCC2022" is prominently displayed above "MINNEAPOLIS". Below the city graphic, there are three key dates: "May 12 Early registration ends", "June 3 Poster/demo abstracts due", and "June 14 Full registration ends". A "Key Dates" button is also present. At the bottom, there is a link to "Donate to the James P. Taylor Foundation for Open Science" and a "Learn More" button.



The development and maintenance of this site is supported by NIH NHGRI award U24 HG006620 and NSF award 1929694. Additional support is provided by NIH awards AI134384 and HG010263, as well as NSF award 1931533.

This is a free, public, internet accessible resource. Data transfer and data storage are not encrypted. If there are restrictions on how your research data can be stored and used, please consult your local institutional review board or the project PI before uploading it to any public site, including this Galaxy server. If you have protected data, large data storage requirements, or short deadlines you are encouraged to set up your own local Galaxy instance or run Galaxy on the cloud.

- Workspace for data-driven science
- Cloud storage to exchange files
- Cloud computing to run analyses on data sets
- Training courses in data science



#### Analyze & Share

Do all your research in one place with our easy to use bioinformatics tools, image analyses, cloud services, and resources for reproducibility, automation, storing, and sharing.

[Learn More](#)



#### Build Your Skills

Access our educational webinars, hands-on workshops, self-paced tutorials, and more to teach yourself and your students how to do open, collaborative science.

[Learn More](#)



#### Bring Your Own

Connect your own storage or compute power, build off our APIs, or even install your own version of CyVerse. We can help you find out what is right for your project and goals.

[Learn More](#)

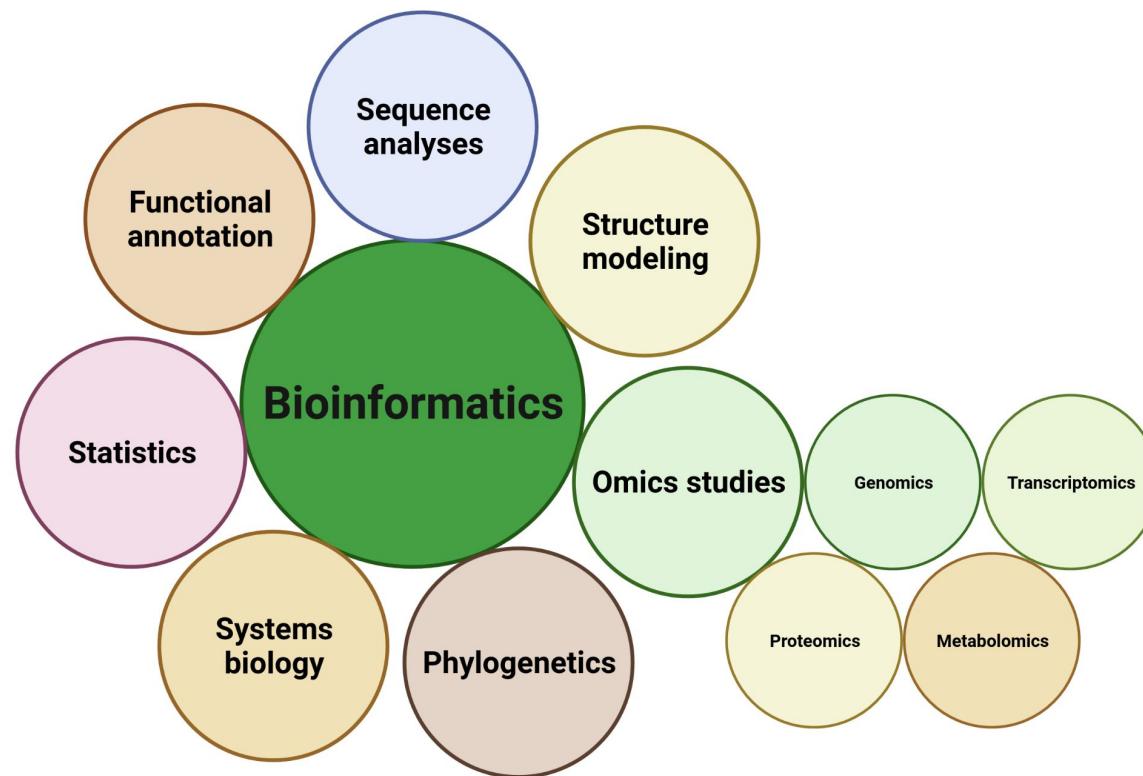
- German Network for Bioinformatics Infrastructure Service, Training, Cooperations & Cloud Computing
- Centrally funded support for bioinformatic analyses in Germany
- Almost 100 tools provided and supported by de.NBI
- Training events for next generation scientists
- de.NBI cloud enables researchers in Germany to conduct bioinformatics research



- Intergovernmental organisation of life scientists and computer scientists in Europe
- Provides resources for European life scientists
- 23 ELIXIR nodes in different countries
- EMBL-EBI heads ELIXIR
- de.NBI is German ELIXIR node



# Summary - introduction to bioinformatics



# Time for questions!

# Questions

1. Which fields are part of bioinformatics?
2. What is the definition of a genome, transcriptome, proteome, and metabolome?
3. What is the full name of 'BLAST'?
4. Which (sequencing) method can be used to analyze a genome?
5. Which method can be used to analyze a transcriptome?
6. Why do biologists apply statistical methods?
7. Which languages are used for the development of bioinformatics tools?
8. Which steps are involved in software development/improvement?
9. Which platforms enable developers to share their software?
10. Which licence can be used to make software freely available to other researchers?
11. Where can you find (online) computational resources for bioinformatics?