# Welcome to the Seminar!



www.biocomicals.com

# Overview of Bioinformatics

Seminar "Informatics in Biochemistry"

By Lukas Jarosch and Leonhard Kohleick

# Question Round

**In what semester are you?**

**What is your name?**

**Let's proceed with [www.menti.com](www.menti.com)**

# What even is "Bioinformatics"?

# Definition of Bioinformatics

" 

*Bioinformatics is an interdisciplinary field that develops methods and software tools for understanding biological data*
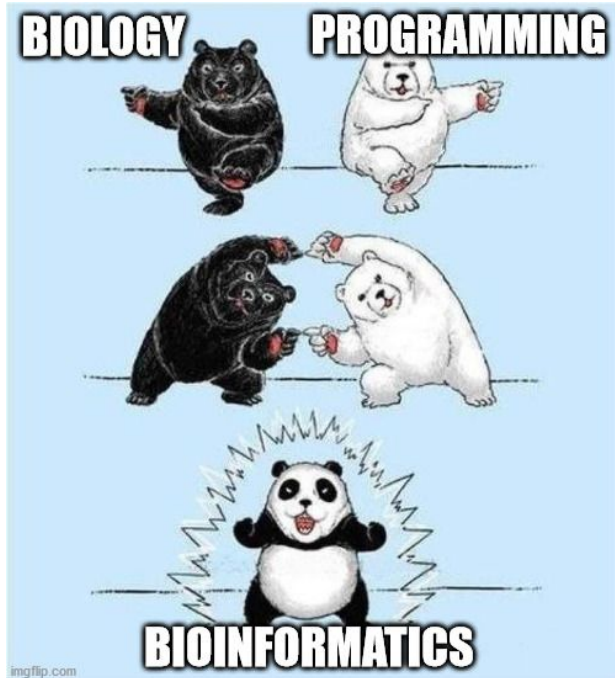
*- Wikipedia*

"

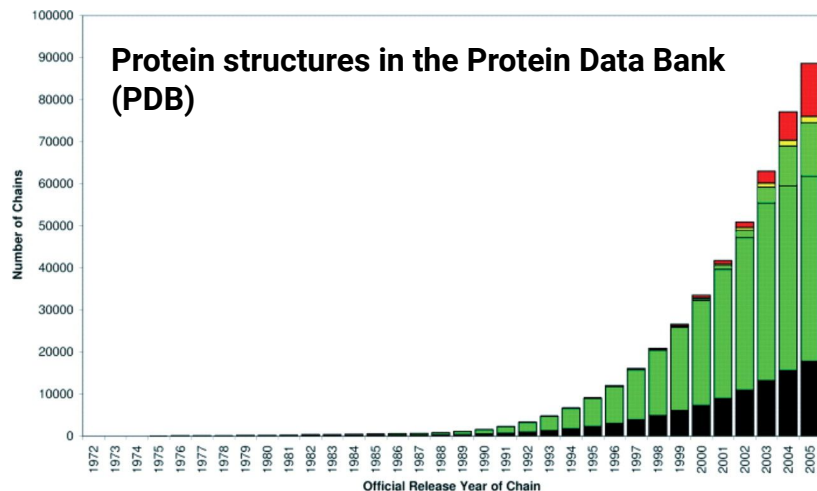# Research fields in Bioinformatics

- **Genomics (DNA)**
  - genotype-phenotype patterns
  - protein homology
  - ancestral relationships

- **Proteomics (proteins)**
  - protein-protein interactions
  - post-translational modifications
  - biomarkers

- **Transcriptomics (RNA)**
  - identifying expressed genes
  - understanding disease mechanisms
  - gene regulatory relationships

- **Structural modeling**
  - protein structure prediction
  - simulating protein function
  - drug identification

- **Image Analysis, Systems Biology, ...**
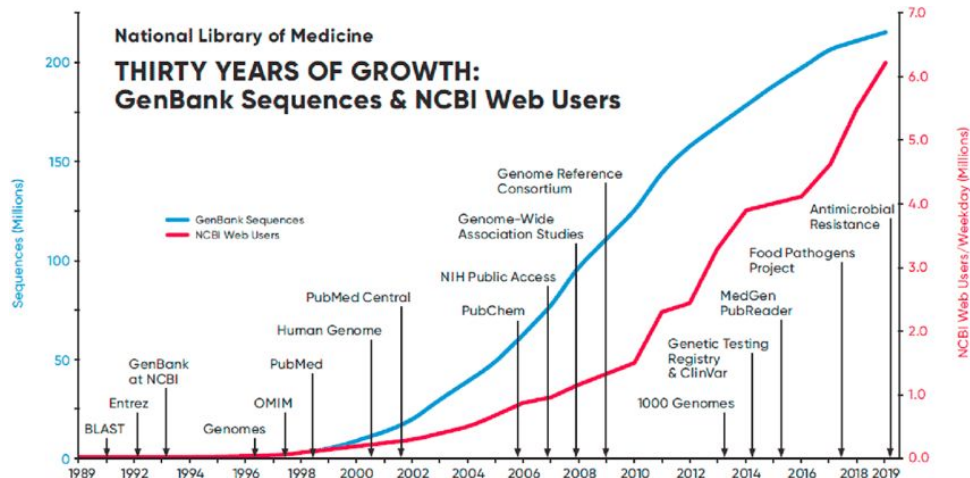
# What skills does a Bioinformatician need?



- Programming (mostly Python & R)

- Statistics

- Biology

- Physics

# Why should you consider learning Bioinformatics?



**Protein structures in the Protein Data Bank (PDB)**

Greene et al., 2017



https://www.nlm.nih.gov/about/2021CJ_NLM.pdf

12

# Why should you consider learning Bioinformatics?

- Current trend towards **high-throughput technologies** that generate massive amounts of data

- Lots of exciting **new research fields**

- **PhD applications** increasingly require coding skills

- Coding is (mostly) **reproducible**

- Programming can **automate** a lot of tedious tasks

13

# Sequence Alignments & Genomics

# Sequence alignment

Protein 1

`KINLKVIKNTLLFRAL`

Protein 2

`GKALLRVRNTLIELAI`

# Sequence alignment

Aligned sequences

K-INLKVIKNTLLFRAL
GKALLRV-RNTLIELAI

# Multiple sequence alignment

https://upload.wikimedia.org/wikipedia/commons/thumb/7/79/RPLP0_90_ClustalW_aln.gif/1024px-RPLP0_90_ClustalW_aln.gif

# Phylogenetics



Theys et al., 2019

18

# Large–scale genome sequencing



https://www.genome.gov/about-genomics/fact-sheets/Sequencing-Human-Genome-cost

http://www.genome.gov/sites/default/files/genome-old/images/content/1000genomes.jpg

# Genome-wide association studies (GWAS)



GWAS of breast cancer in Japanese population

Low et al., 2013

# Modeling

# Protein engineering with Rosetta

# Rosetta can predict protein folding and also has various design features



A

TS

U  search

N

Protein with a novel globular protein fold:

TOP7
PDB-ID: 1QYS

(Englander and Mayne, 2014) & (Kuhlman et al., 2003)  23

# Designing a protein with Rosetta



Mutate residue 11 to phenylalanine
Mutate residue 34 to cysteine

Refine the model

(Leman et al., 2020)

# Vaccines designed with Rosetta (for RSV)



(Marcandalli et al., 2019)

# A new quadrivalent vaccine against influenza



(Boyoglu-Barnum et al., 2021)

# Rosetta has a lot of capabilities



Structure prediction

Protein–protein docking

Ligand docking

Loop modeling

Biomineral surface docking

Protein design

Modeling with experimental data

Tasks

Systems

Symmetric assemblies

Non-canonical chemistries

Antibodies

Membrane proteins

RNA and DNA

Peptides

Carbohydrates

(Leman et al., 2020)   27

# Molecular Dynamics

# Studying physical movements of atoms and molecules



Hemoglobin

# Molecular dynamics searches for the lowest energy conformation

# Use of Molecular Dynamics in Biochemistry

- Using Newton's laws of motion, we predict the position of each atom as a function over time



bonded interactions    non-bonded interactions

# Studying conformational flexibility and stability



https://www.
youtube.co
m/watch?v=
7AhQ19m2o
k4

# Modeling the lipid bilayer to better understand cell behaviour

https://www.youtube.com/watch?v=SbWh_XgCHyw

# Drug screening

# Normal drug screening workflow



construction of the reporter gene

recombined mild CoVs

compound libraries

high-throughput drug screening against CoV

structural optimization

Hits

evaluation agai...

Leads

(Liu et al., 2020) 35

# In silico screening reduces the number of molecules of interest



(Ou-Yang et al. 2012)    36

# Deep learning yields new antibiotics candidates



10^8 molecules screened

(Stokes et al., 2020)

# Introduction to programming languages

# From machine code to Python

```
b8      21 0a 00 00
a3      0c 10 00 06
b8      6f 72 6c 64
a3      08 10 00 06
b8      6f 2c 20 57
a3      04 10 00 06
b8      48 65 6c 6c
a3      00 10 00 06
b9      00 10 00 06
ba      10 00 00 00
bb      01 00 00 00
b8      04 00 00 00
cd      80
b8      01 00 00 00
cd      80
```

print("Hello World")

Languages suited for general applications

Languages suited for research

Java

C/C++

C#

Python

R

Matlab

# Why Python?

- Easy to Learn and Use
- Big Community
- Established in the Corporate World ( primarily Google and YouTube)
- Versatility, Efficiency, Reliability, and Speed
- Hundreds of Libraries and Frameworks

# Python vs. R

|  | **Python** | **R** |
| --- | --- | --- |
| **Purpose** | multi purpose language | mainly statistical and data science applications |
| **Learning curve:** | easy to learn, linear learning curve | easy to learn, advanced functionalities can be difficult to use |
| **Used by:** | industry, academia, engineering | academics, scientists without programming skills, (few industries) |
| **Visualization of Data** | can be difficult | easy, straightforward |

# How to run Python?

The most relevant ways are:
- Simply executing .py files in your command line
- Using an IDE (Integrated Development Environment)
- The standard Python shell
- What we use in the course: **Jupyter notebook**

# What is the "Jupyter Notebook"?

- Integrated framework to write:
    - Python code
    - Text/Markdown
    - Figures, Tables, animations

- Runs inside your browser
- Makes research reproducible and allows others to understand your code better

# Jupyter Notebook

Editor that allows you to combine Python code with formatted text, plots, tables, and even animations.

Runs in your web-browser and is compatible with all operating systems

➡️ **reproducible research**
➡️ **great for sharing**



https://jupyter.org/assets/jupyterpreview.png

# Our roadmap for today:

| Day 1   9:00-18:00  (Lunch Break 13:00-13:45) | | |
|---|---|---|
| **Topic** | **Time** | **Break** |
| Talk "Overview of Bioinformatics" | 1h | |
| Basic Variables and DataTypes | 2h | 15min |
| Conditional Clauses | | |
| Lists and Tuple | 2h | 15min |
| Sets and Dictionaries | | |
| Loops | 2h | 15min |
| Functions | | |

# Schedule for tomorrow

| Day 2 10:00 -15:00   (Lunch Break 12:45 - 13:30 ) | | |
|---|---|---|
| **Topic** | **Time** | **Break** |
| Working with .txt and .csv files | 2h | 15 min |
| Data plotting | | |
| Big Exercise 1 | 1h | 10 min |
| Big Exercise 2 | 1h | |
| Guest talk by Jannik Buhr: "Data Science with R" | 30min | |

# References and further reading material

# References: Modeling

Kuhlman, B., Dantas, G., Ireton, G.C., Varani, G., Stoddard, B.L., and Baker, D. (2003). Design of a Novel Globular Protein Fold with Atomic-Level Accuracy. Science *302*, 1364–1368.

Englander, S.W., and Mayne, L. (2014). The nature of protein folding pathways. PNAS *111*, 15873–15880.

Boyoglu-Barnum, S., Ellis, D., Gillespie, R.A., Hutchinson, G.B., Park, Y.-J., Moin, S.M., Acton, O.J., Ravichandran, R., Murphy, M., Pettie, D., et al. (2021). Quadrivalent influenza nanoparticle vaccines induce broad protection. Nature *592*, 623–628.

Marcandalli, J., Fiala, B., Ols, S., Perotti, M., de van der Schueren, W., Snijder, J., Hodge, E., Benham, M., Ravichandran, R., Carter, L., et al. (2019). Induction of Potent Neutralizing Antibody Responses by a Designed Protein Nanoparticle Vaccine for Respiratory Syncytial Virus. Cell *176*, 1420-1431.e17

Leman JK et al. . Macromolecular modeling and design in Rosetta: recent methods and frameworks. Nat Methods. 2020 Jul;17(7):665-680. doi: 10.1038/s41592-020-0848-2. Epub 2020 Jun 1. PMID: 32483333; PMCID: PMC7603796.

Hub JS, Kubitzki MB, de Groot BL. Spontaneous quaternary and tertiary T-R transitions of human hemoglobin in molecular dynamics simulation. PLoS Comput Biol. 2010 May 6;6(5):e1000774. doi: 10.1371/journal.pcbi.1000774. PMID: 20463873; PMCID: PMC2865513.

Ou-Yang SS, Lu JY, Kong XQ, Liang ZJ, Luo C, Jiang H. Computational drug discovery. Acta Pharmacol Sin. 2012 Sep;33(9):1131-40. doi: 10.1038/aps.2012.109. Epub 2012 Aug 27. PMID: 22922346; PMCID: PMC4003107.

Liu J, Li K, Cheng L, Shao J, Yang S, Zhang W, Zhou G, de Vries AAF, Yu Z. A high-throughput drug screening strategy against coronaviruses. Int J Infect Dis. 2021 Feb;103:300-304. doi: 10.1016/j.ijid.2020.12.033. Epub 2020 Dec 14. PMID: 33333250; PMCID: PMC7832824.

Stokes JM et al. A Deep Learning Approach to Antibiotic Discovery. Cell. 2020 Feb 20;180(4):688-702.e13. doi: 10.1016/j.cell.2020.01.021. Erratum in: Cell. 2020 Apr 16;181(2):475-483. PMID: 32084340; PMCID: PMC8349178.