

UNIVERSITÄT
HEIDELBERG



Introduction

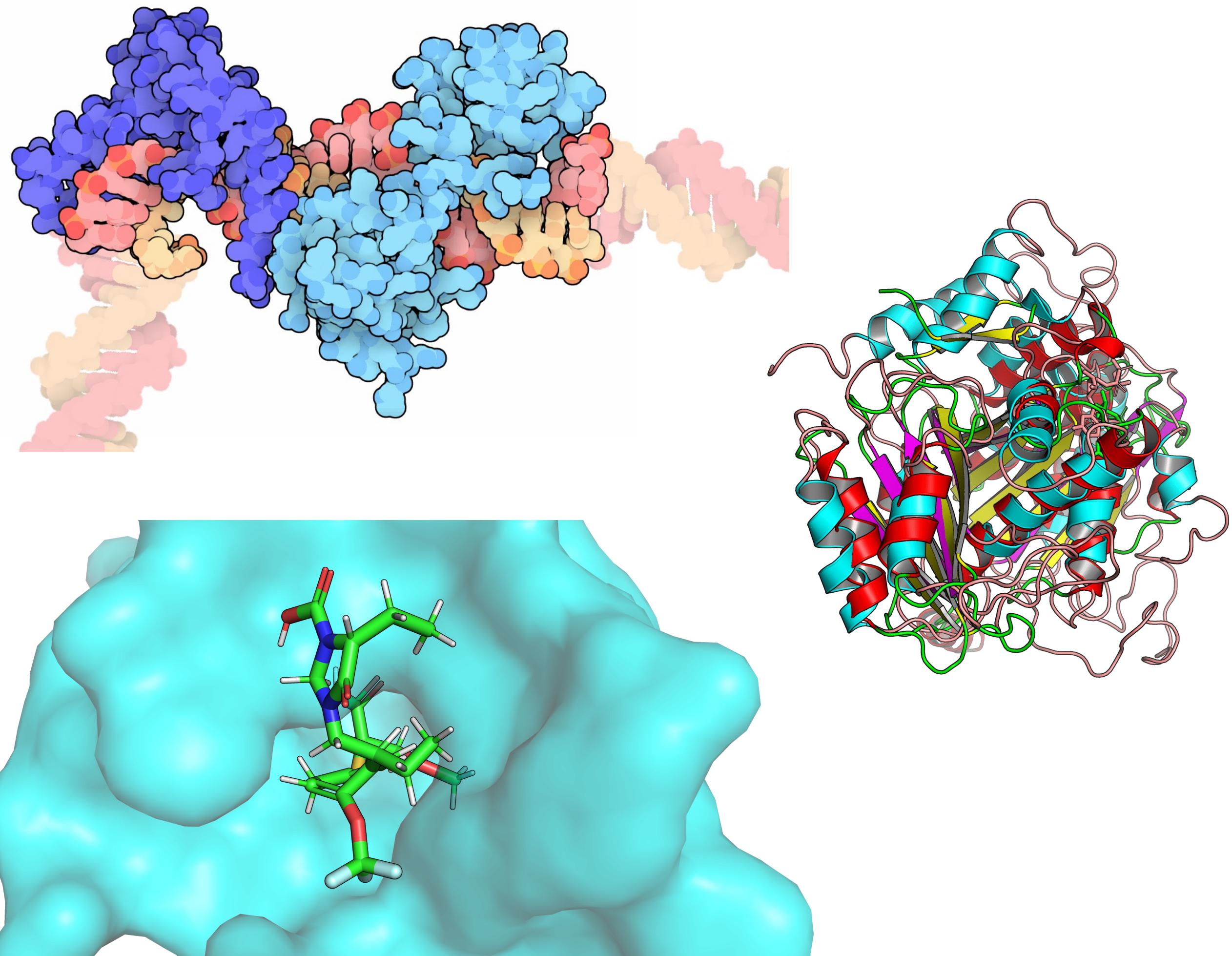
L1, Structural Bioinformatics

WiSe 2023/24, Heidelberg University

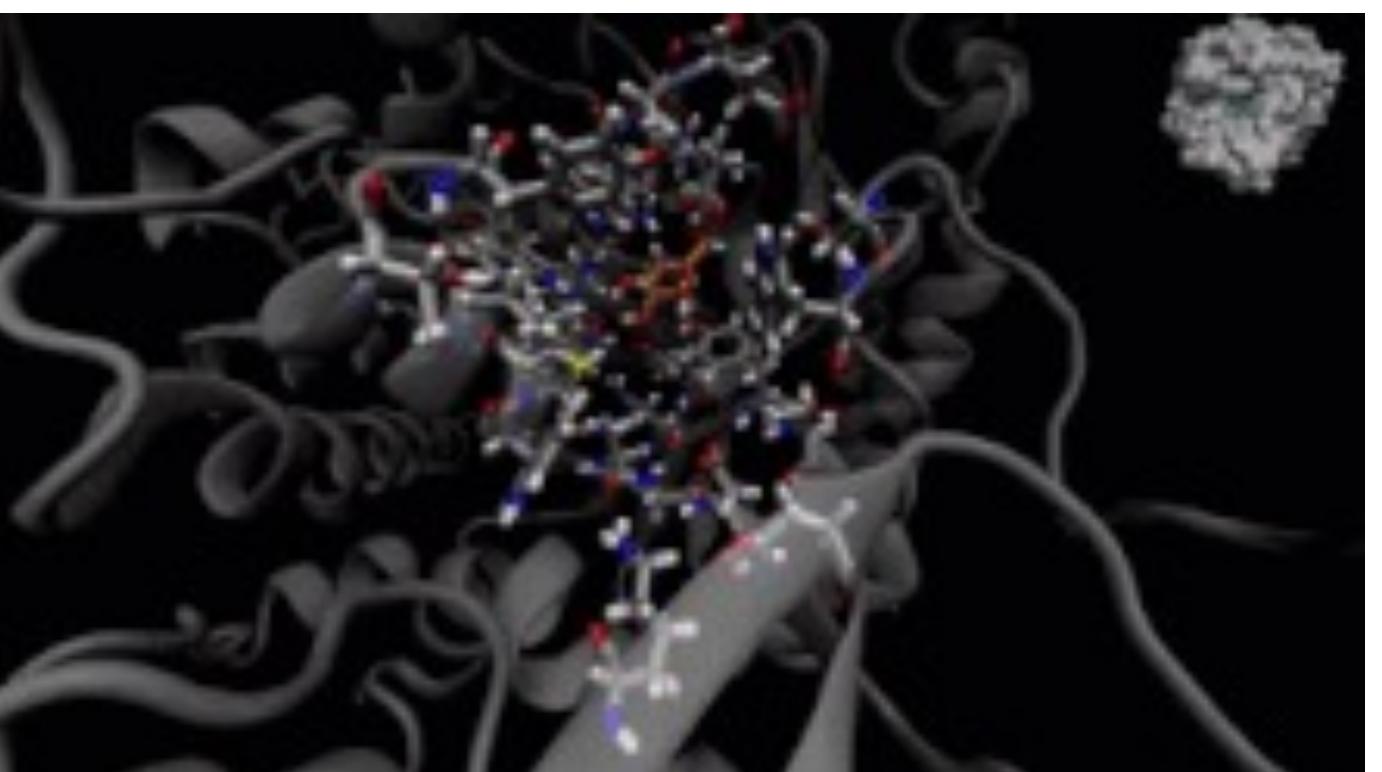
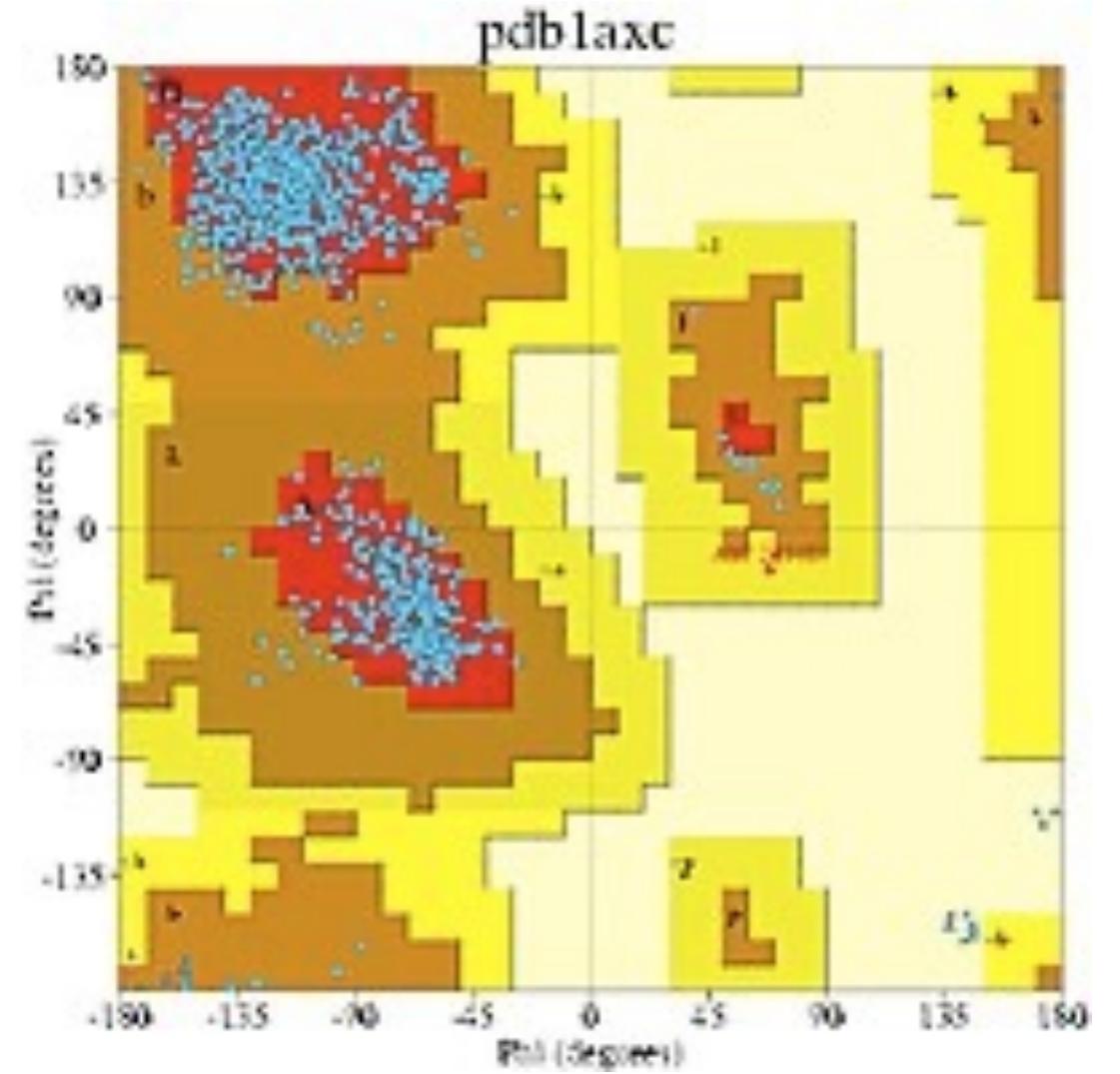
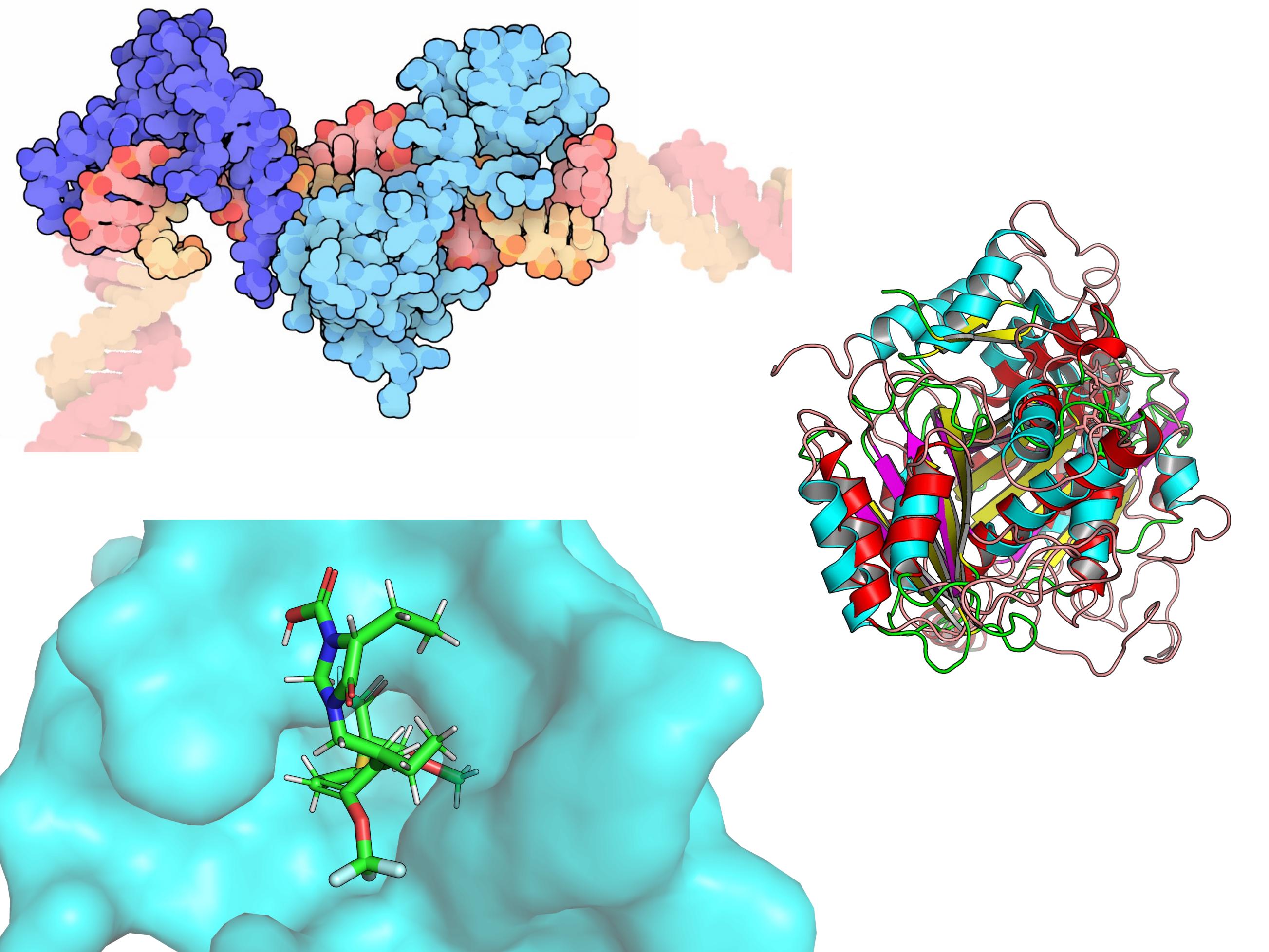
Kieran Didi

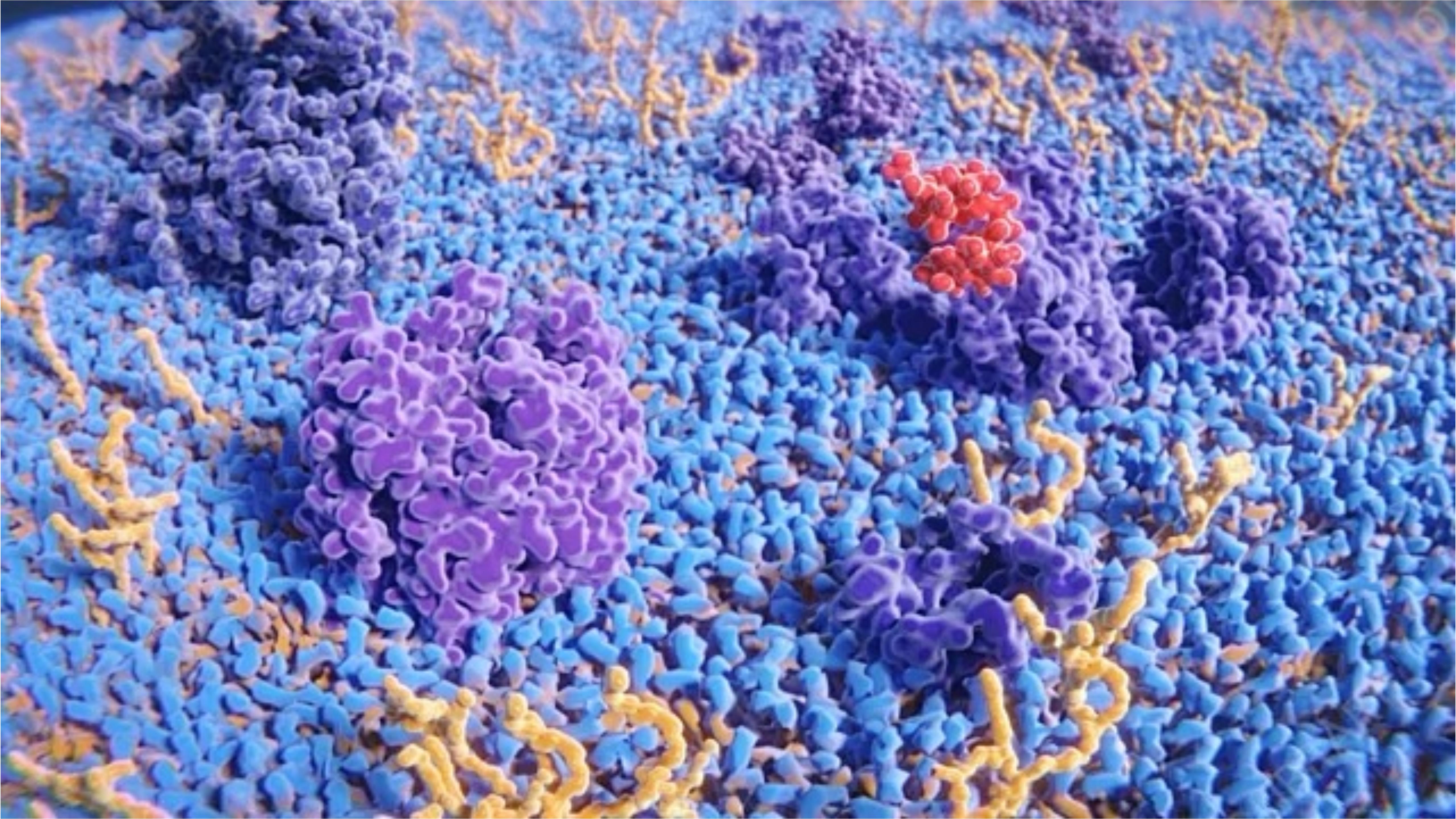
Structural **Bio**informatics

Structural Bioinformatics



Structural Bioinformatics



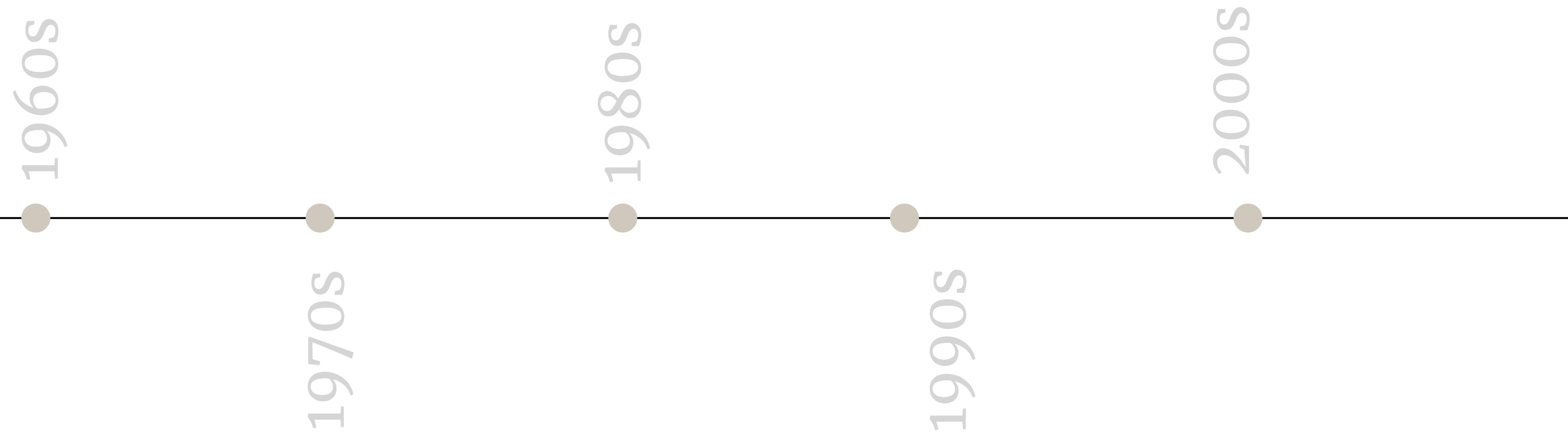


Overview

- 1. A Brief History of the Field**
- 2. Where we are and where we are headed**
- 3. This course**
- 4. To-Dos for you!**

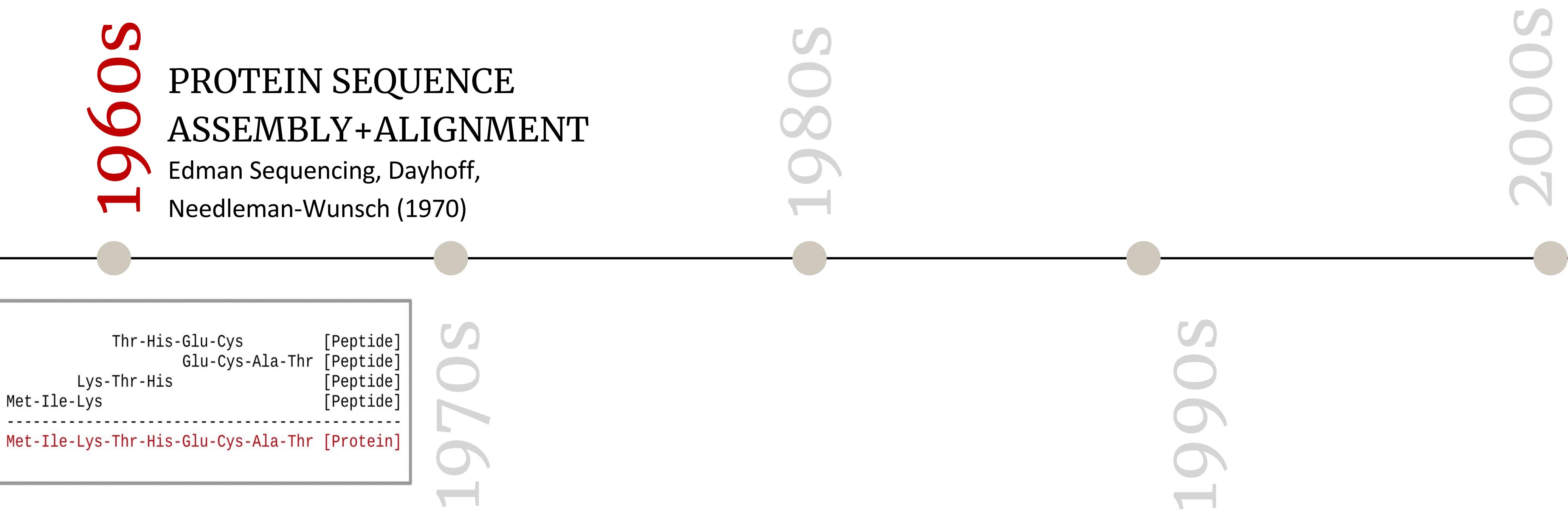
Where do we come from?

Bioinformatics did not start structural, and not with DNA



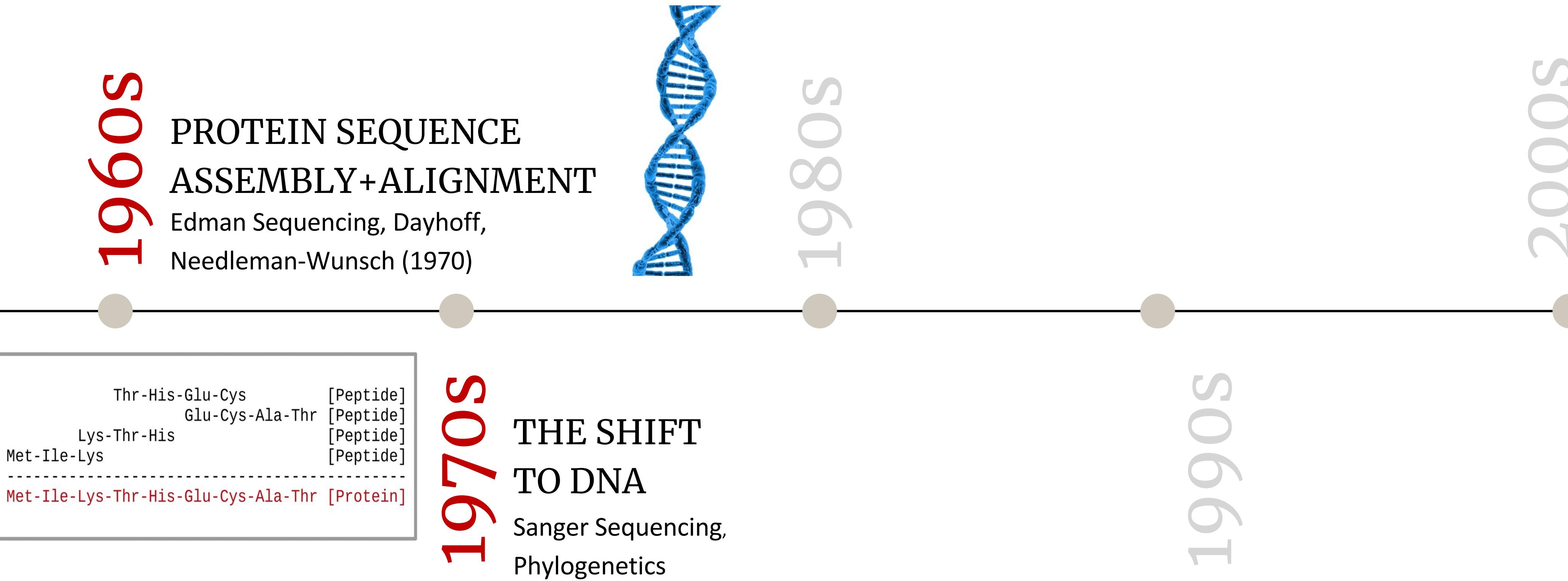
Where do we come from?

Bioinformatics did not start structural, and not with DNA



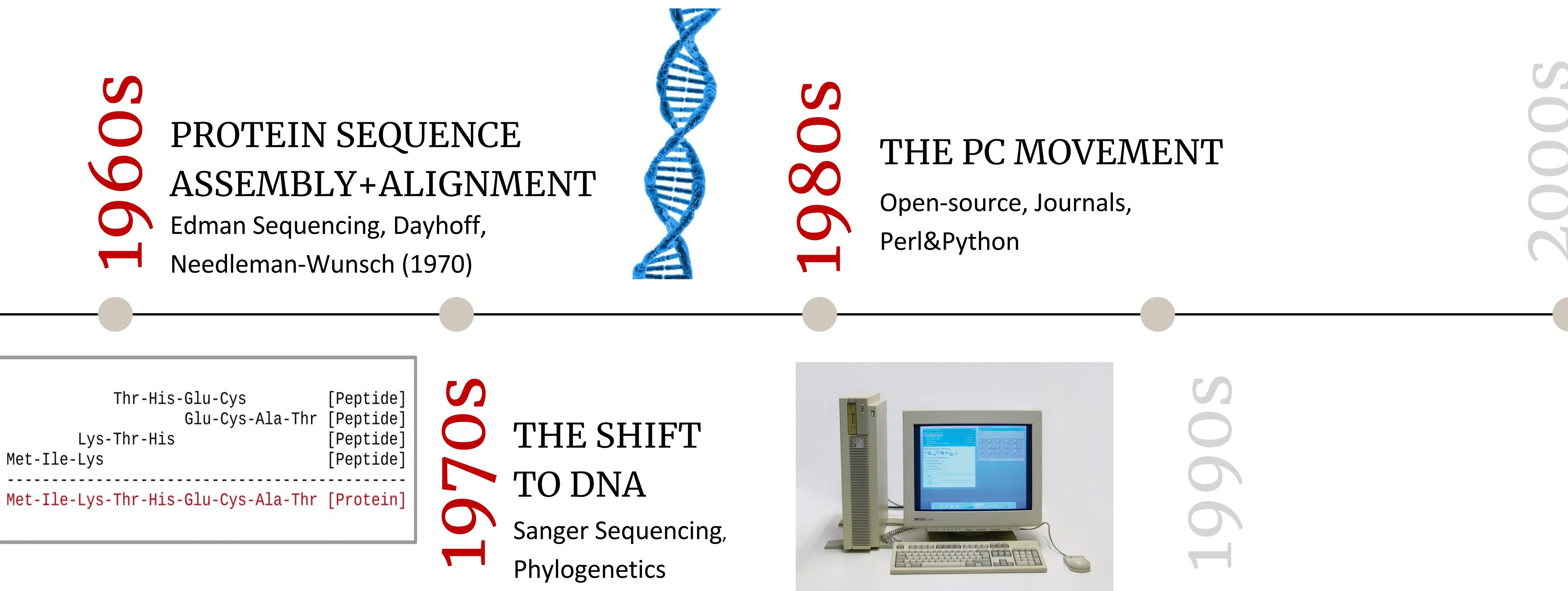
Where do we come from?

Bioinformatics did not start structural, and not with DNA



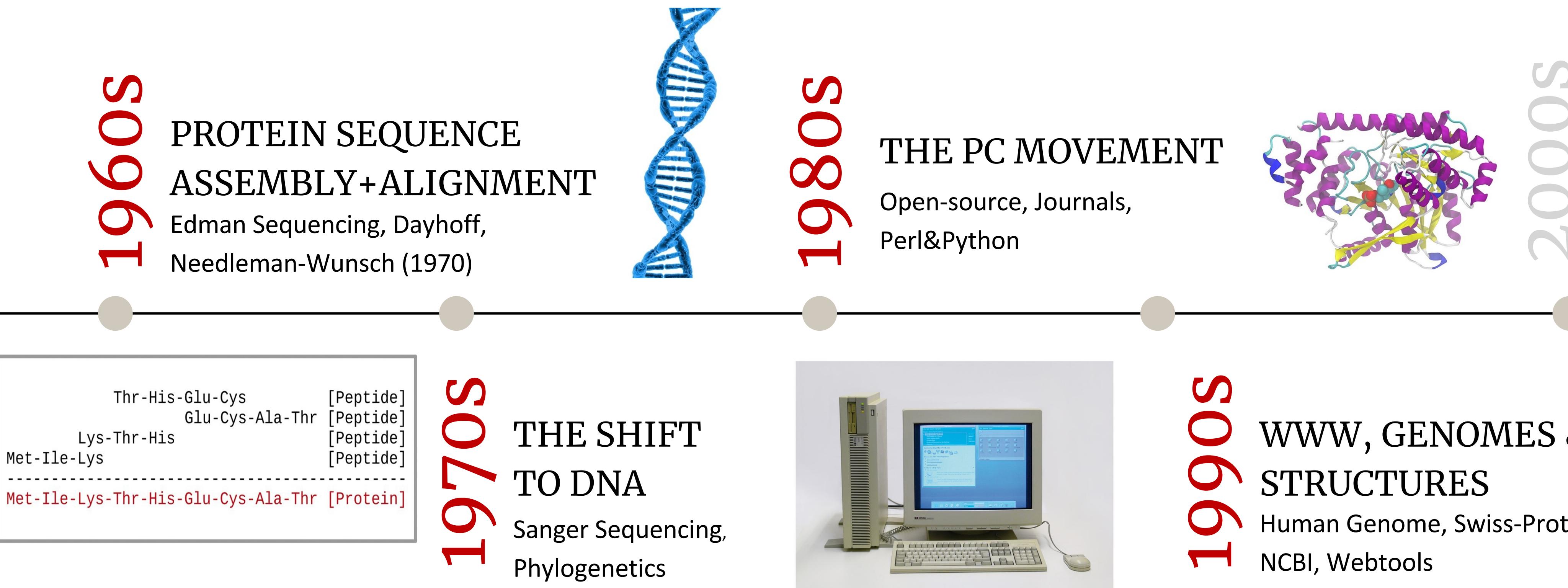
Where do we come from?

Bioinformatics did not start structural, and not with DNA



Where do we come from?

Bioinformatics did not start structural, and not with DNA



Where do we come from?

Bioinformatics did not start structural, and not with DNA

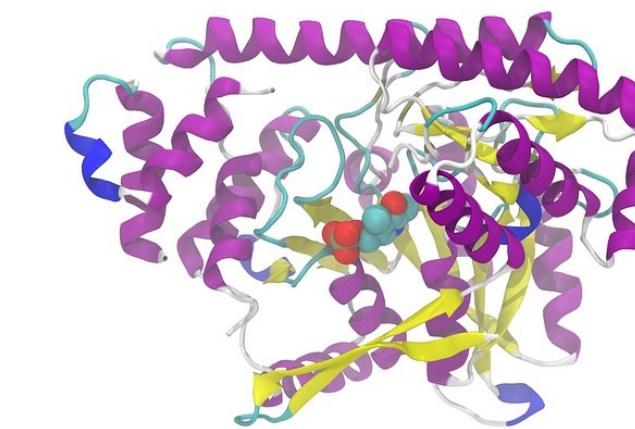
1960S

PROTEIN SEQUENCE
ASSEMBLY+ALIGNMENT
Edman Sequencing, Dayhoff,
Needleman-Wunsch (1970)



1980S

THE PC MOVEMENT
Open-source, Journals,
Perl&Python



2000S

HIGH THROUGHPUT
NGS, Compute Clusters,
PDB 3000 -> 8000 entries

Thr-His-Glu-Cys	[Peptide]
Glu-Cys-Ala-Thr	[Peptide]
Lys-Thr-His	[Peptide]
Met-Ile-Lys	[Peptide]
<hr/>	
Met-Ile-Lys-Thr-His-Glu-Cys-Ala-Thr [Protein]	

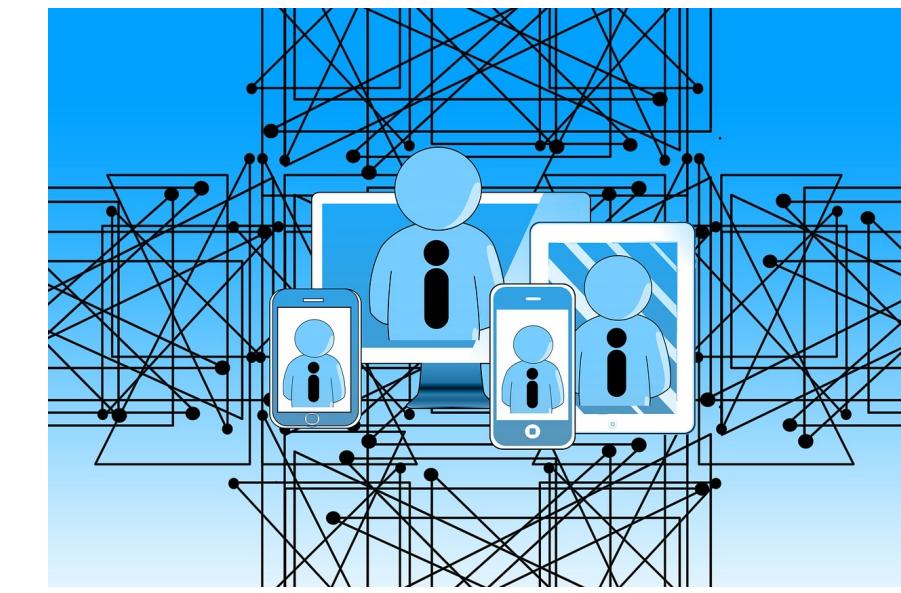
1970S

THE SHIFT
TO DNA
Sanger Sequencing,
Phylogenetics



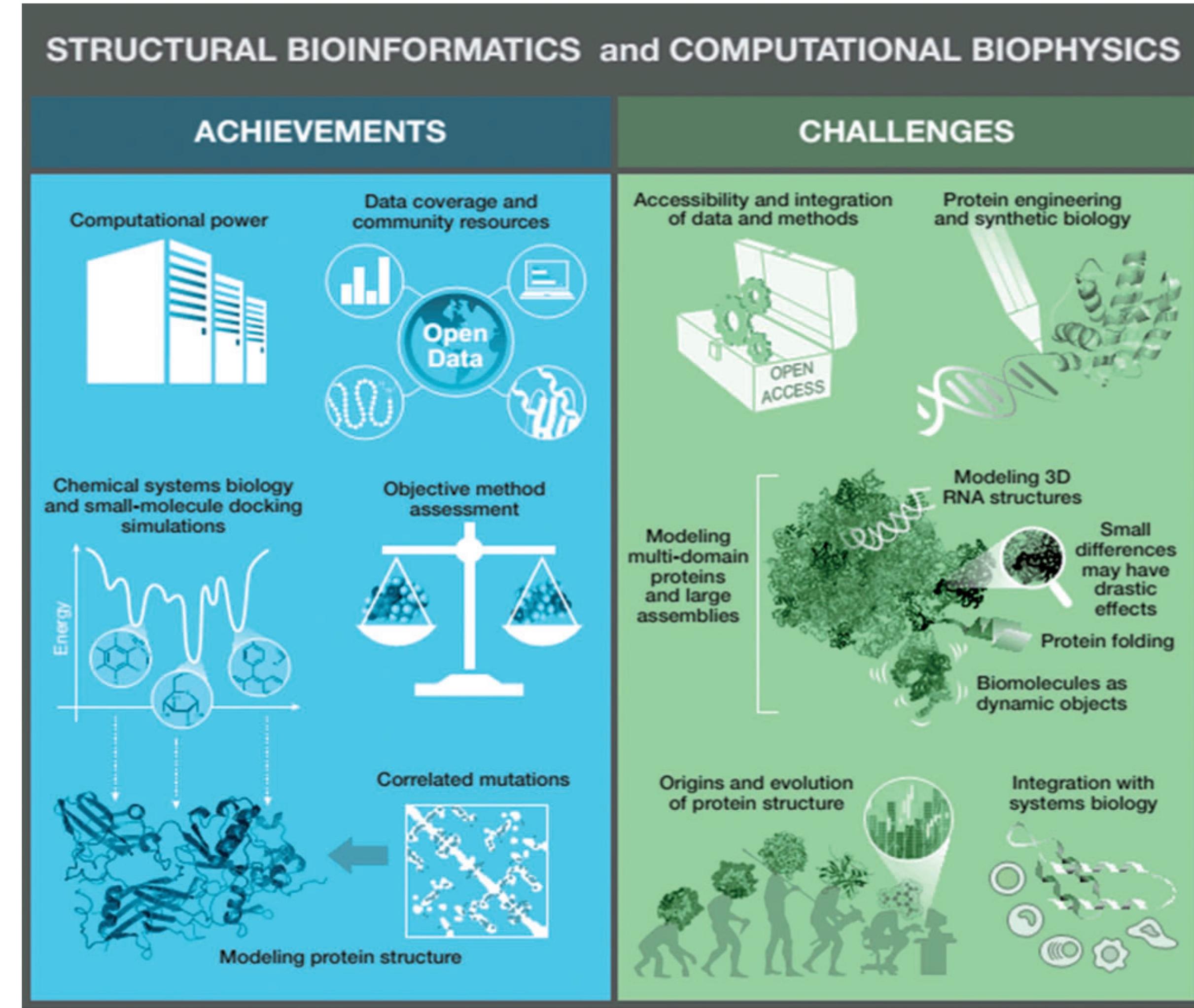
1990S

WWW, GENOMES &
STRUCTURES
Human Genome, Swiss-Prot,
NCBI, Webtools



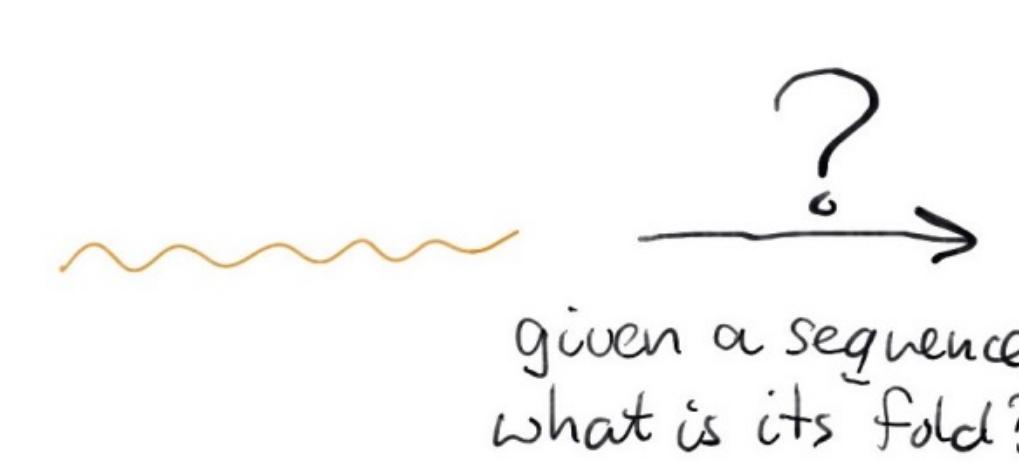
Where are we today?

Achievements in Structural Bioinformatics (2014)

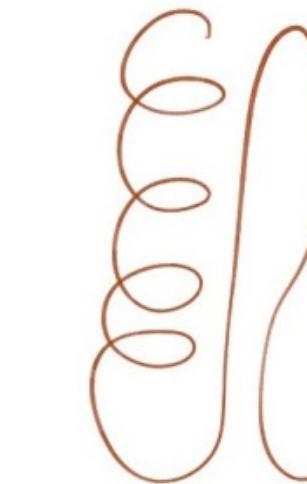


What are the questions we ask?

And how can we answer them?



given a sequence,
what is its fold?

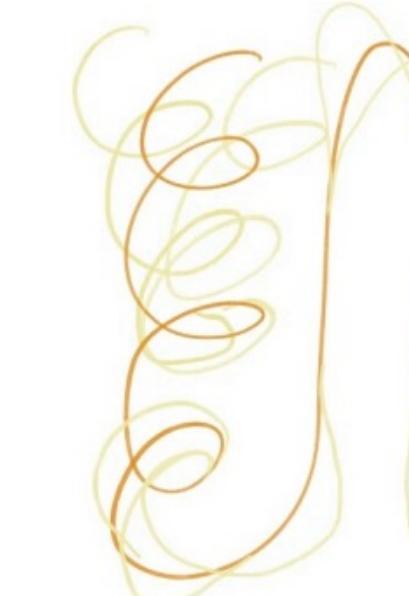
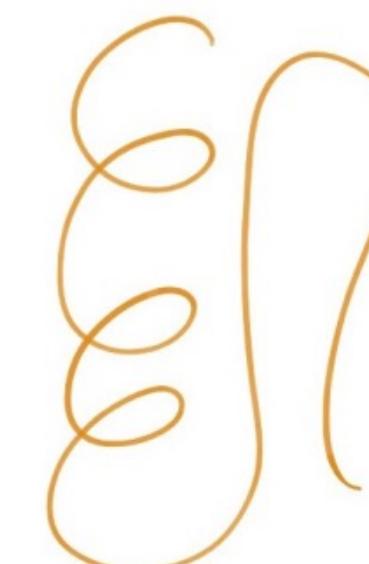


Protein Structure
Prediction
AlphaFold2



how similar
are these structures?

Sequence Alignments
Structure Alignments
Classification

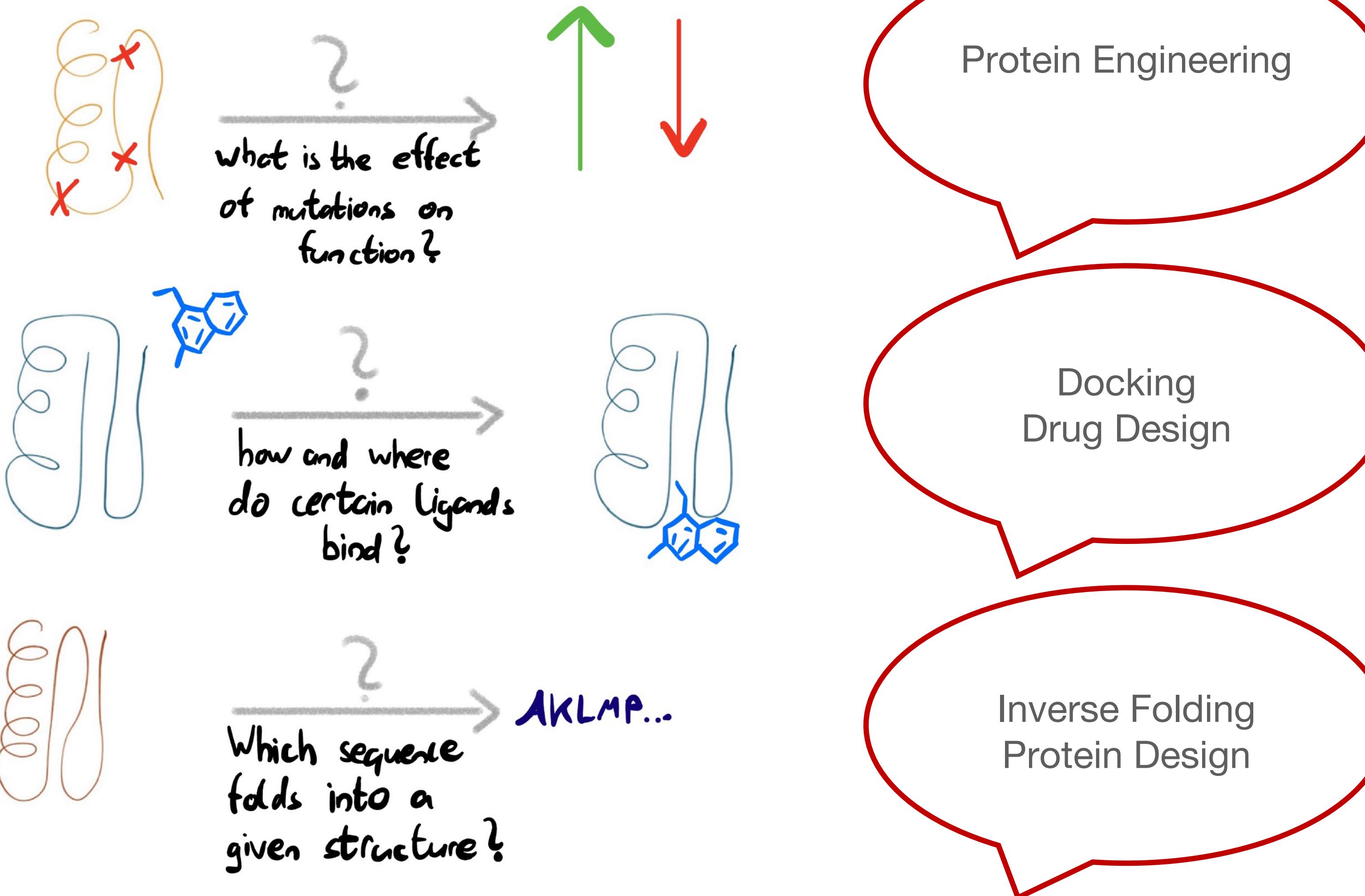


how do dynamics
affect function?

MD Simulations
Protein Design

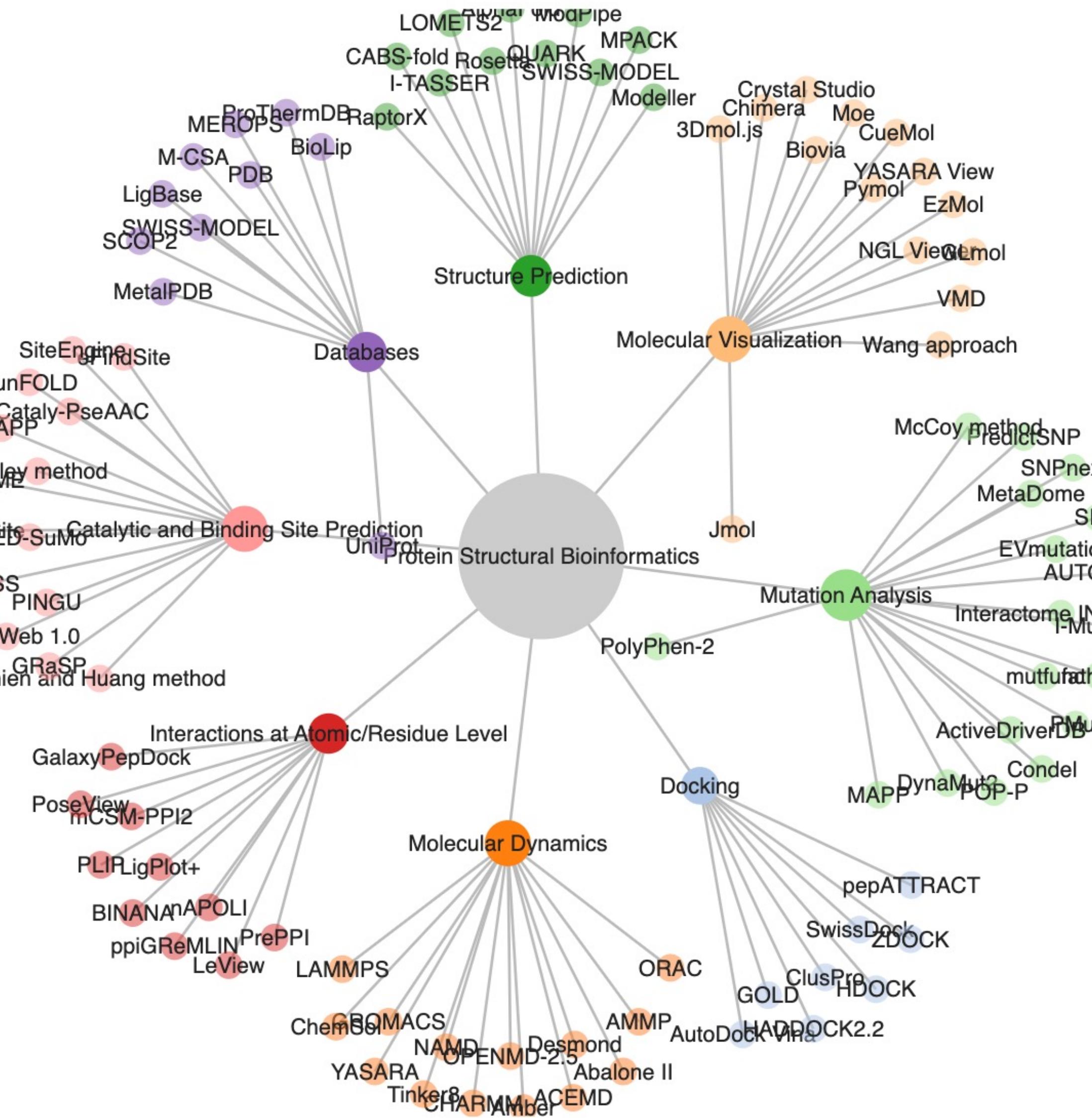
What are the questions we ask?

And how can we answer them?



What are the tools at our disposal?

PreStO



This course

What will we talk about?

1 Intro

2 ML Basics

3 ML
Architectures

4 Graph ML

5 Structure
Prediction

6 Generative
Models

7 Protein
Design

8 Simulations

9 Evolution &
Bioinformatics

10 Drug
Design

11 Graphics&
Animations

12 Future
Directions

Know your tools

Pymol – Python - Proteins



PyTorch

To-Dos for you!

1. Enter the Discord server
2. Install Pymol
3. Read the [Python Post](#) and do Google Colab Intro
4. Do the first exercises!

UNIVERSITÄT
HEIDELBERG



Math Primer 1: Linear Algebra

L1, Structural Bioinformatics

WiSe 2023/24, Heidelberg University

Kieran Didi

Overview

- 1. A Brief History of the Field**
- 2. Where we are and where we are headed**
- 3. This course**
- 4. To-Dos for you!**

Scalars

- A scalar is a single number
- Integers, real numbers, rational numbers, etc.
- We denote it with italic font:

a, n, x

Vectors

- A vector is a 1-D array of numbers:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}. \quad (2.1)$$

- Can be real, binary, integer, etc.
- Example notation for type and size:

\mathbb{R}^n

Matrices

- A matrix is a 2-D array of numbers:

The diagram shows a 2x2 matrix with elements $A_{1,1}$, $A_{1,2}$, $A_{2,1}$, and $A_{2,2}$. A green oval encloses the first row ($A_{1,1}$ and $A_{1,2}$), with a black arrow pointing to it labeled "Row". An orange oval encloses the first column ($A_{1,1}$ and $A_{2,1}$), with a black arrow pointing to it labeled "Column". The matrix is enclosed in square brackets with a period at the end.

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}. \quad (2.2)$$

- Example notation for type and shape:

$$A \in \mathbb{R}^{m \times n}$$

Tensors

- A tensor is an array of numbers, that may have
 - zero dimensions, and be a scalar
 - one dimension, and be a vector
 - two dimensions, and be a matrix
 - or more dimensions.

Matrix Transpose

$$(A^\top)_{i,j} = A_{j,i}. \quad (2.3)$$

The diagram shows a 3x2 matrix A with elements $A_{1,1}, A_{1,2}, A_{2,1}, A_{2,2}, A_{3,1}, A_{3,2}$. A curved arrow starts from the top-left element $A_{1,1}$ and points to the bottom-right element $A_{3,2}$, illustrating that the transpose operation reflects the matrix across its main diagonal.

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow A^\top = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

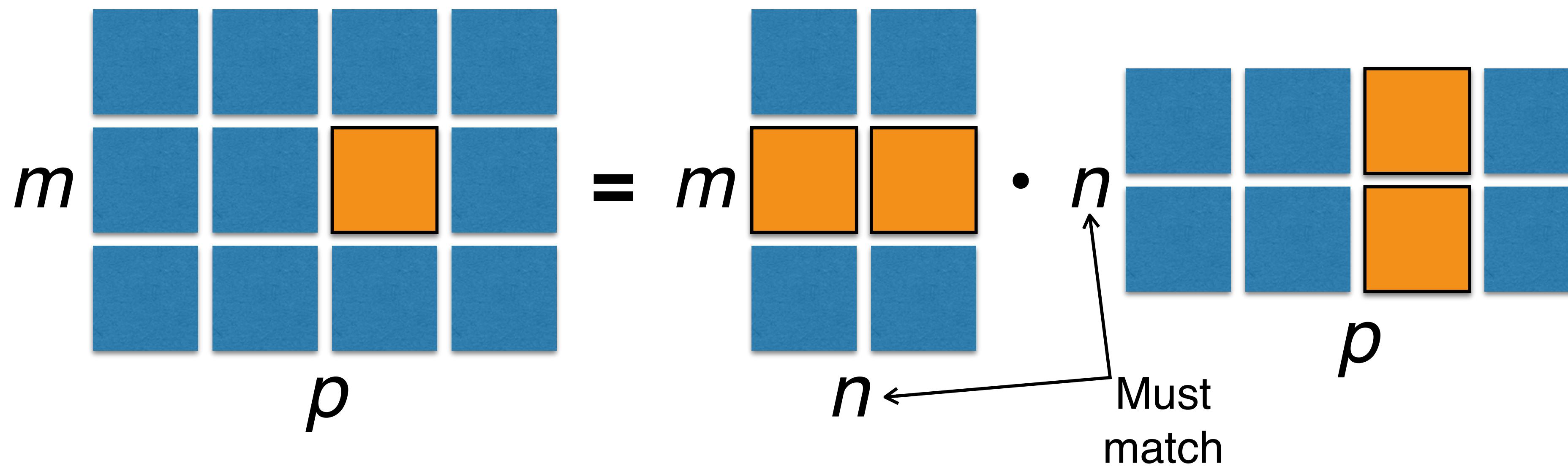
Figure 2.1: The transpose of the matrix can be thought of as a mirror image across the main diagonal.

$$(AB)^\top = B^\top A^\top. \quad (2.9)$$

Matrix (Dot) Product

$$C = AB. \quad (2.4)$$

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}. \quad (2.5)$$



Identity Matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Figure 2.2: *Example identity matrix:* This is I_3 .

$$\forall \mathbf{x} \in \mathbb{R}^n, I_n \mathbf{x} = \mathbf{x}. \quad (2.20)$$

Systems of Equations

$$Ax = b \tag{2.11}$$

expands to

$$A_{1,:}\mathbf{x} = b_1 \tag{2.12}$$

$$A_{2,:}\mathbf{x} = b_2 \tag{2.13}$$

$$\dots \tag{2.14}$$

$$A_{m,:}\mathbf{x} = b_m \tag{2.15}$$

Solving Systems of Equations

- A linear system of equations can have:
 - No solution
 - Many solutions
 - Exactly one solution: this means multiplication by the matrix is an invertible function

Matrix Inversion

- Matrix inverse:

$$A^{-1}A = I_n. \quad (2.21)$$

- Solving a system using an inverse:

$$Ax = b \quad (2.22)$$

$$A^{-1}Ax = A^{-1}b \quad (2.23)$$

$$I_nx = A^{-1}b \quad (2.24)$$

- Numerically unstable, but useful for abstract analysis

Invertibility

- Matrix can't be inverted if...
 - More rows than columns
 - More columns than rows
 - Redundant rows/columns (“linearly dependent”, “low rank”)

Norms

- L^p norm

$$\|x\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}$$

- Most popular norm: L2 norm, $p=2$

$$\|x\|_1 = \sum_i |x_i|. \tag{2.31}$$

$$\|x\|_\infty = \max_i |x_i|. \tag{2.32}$$

UNIVERSITÄT
HEIDELBERG



Math Primer 2: Probability

L1, Structural Bioinformatics

WiSe 2023/24, Heidelberg University
Kieran Didi

Probability Mass Function

Describing discrete event space

- The domain of P must be the set of all possible states of x .
- $\forall x \in X, 0 \leq P(x) \leq 1$. An impossible event has probability 0 and no state can be less probable than that. Likewise, an event that is guaranteed to happen has probability 1, and no state can have a greater chance of occurring.
- $\sum_{x \in X} P(x) = 1$. We refer to this property as being **normalized**. Without this property, we could obtain probabilities greater than one by computing the probability of one of many events occurring.

Example: uniform distribution: $P(X = x_i) = \frac{1}{k}$

Probability Density Function

Describing continuous event space

- The domain of p must be the set of all possible states of x .
- $\forall x \in \mathcal{X}, p(x) \geq 0$. Note that we do not require $p(x) \leq 1$.
- $\int p(x)dx = 1$.

Example: uniform distribution: $u(x; a, b) = \frac{1}{b-a}$.

The Sum Rule of Probability

How to calculate a marginal

$$\forall x \in X, P(X = x) = \sum_y P(X = x, Y = y). \quad (3.3)$$

$$p(x) = \int p(x, y) dy. \quad (3.4)$$

Conditional Probability

A slice through the distribution

$$P(y = y \mid x = x) = \frac{P(y = y, x = x)}{P(x = x)}.$$

The Chain Rule of Probability

How to factor a joint distribution

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)}). \quad (3.6)$$

(Conditional) Independence

When can we consider events separately?

$$\forall x \in X, y \in Y, p(X = x, Y = y) = p(X = x)p(Y = y). \quad (3.7)$$

$$\forall x \in X, y \in Y, z \in Z, p(X = x, Y = y \mid Z = z) = p(X = x \mid Z = z)p(Y = y \mid Z = z). \quad (3.8)$$

Expectation

A weighted average of all possible outcomes

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x), \quad (3.9)$$

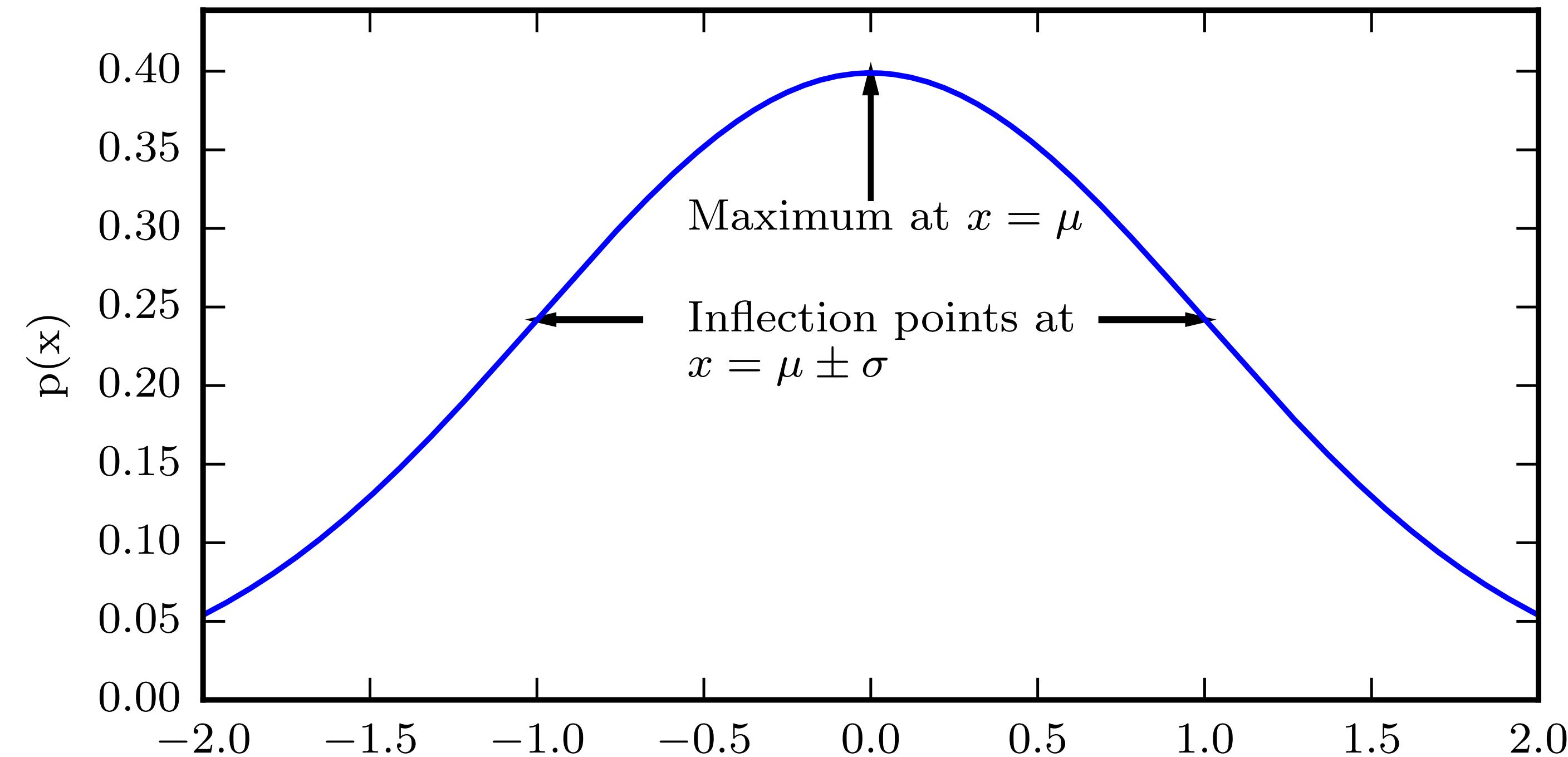
$$\mathbb{E}_{x \sim p}[f(x)] = \int p(x)f(x)dx. \quad (3.10)$$

linearity of expectations:

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)], \quad (3.11)$$

Gaussian Distribution

The bread-and-butter of ML



Gaussian Distribution

The bread-and-butter of ML

$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (3.21)$$

$$\mathcal{N}(x; \mu, \Sigma) = \sqrt{\frac{1}{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1} (x - \mu)\right). \quad (3.23)$$

Bayes' Rule

Incorporating prior knowledge

$$P(x | y) = \frac{P(x)P(y | x)}{P(y)}. \quad (3.42)$$