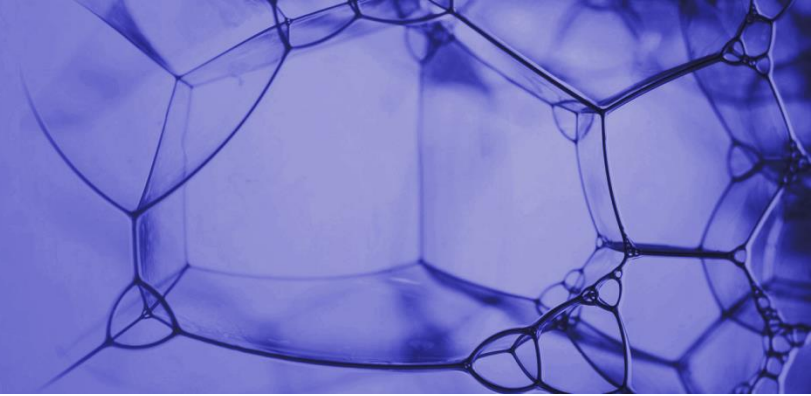# 12. Artificial Intelligence

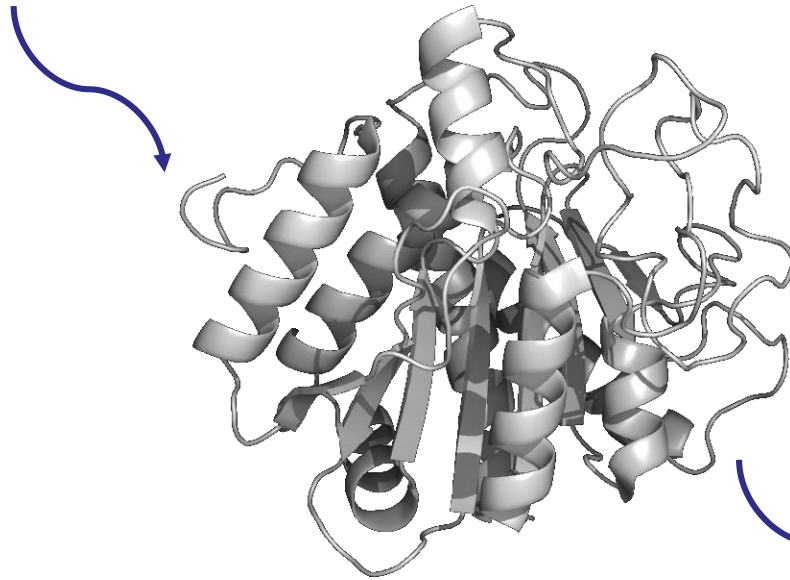# in Life Sciences

# Outline

- ❑ **Motivation**

- ❑ **Introduction to AI and ML**

- ❑ **Modern challenges in Bioengineering**

- ❑ **Basics of ML**
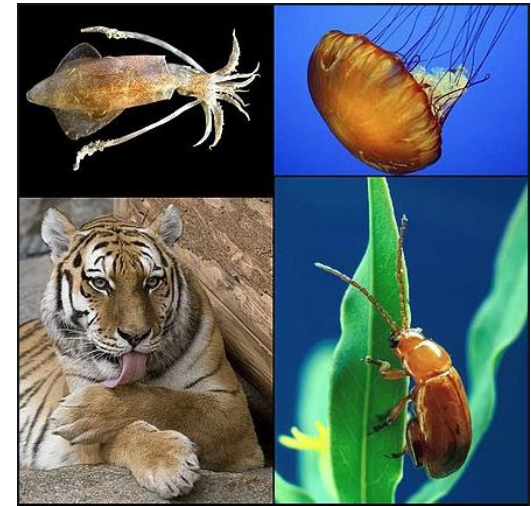
- ❑ **Recent applications**

# Motivation

MSLGAKPFGEKKFIEIKGRRMAYIDEGTGDPILFQHGNPTSSYLWRNIMPHC
AGLGRLIACDLIGMGDSDKLDPSGPERYAYAEHRDYLDALWEALDLGDRVV
LVVHDWGSALGFDWARRHRERVQGIAYMEAIAMPIEWADFPEQDRDLFQ
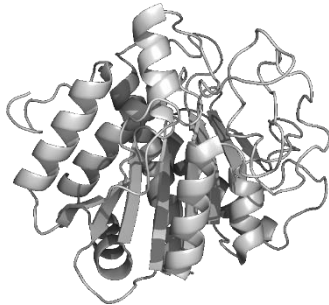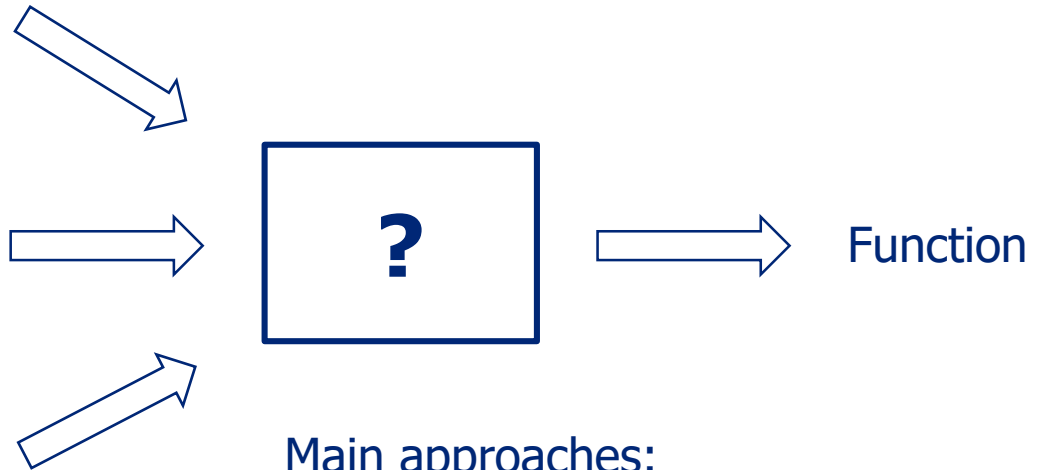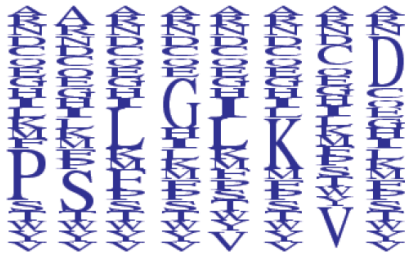AFRSQAGEELVLQD



**Function**

# Motivation

## Sequence

```
MSLGAKPFGEKKFIEIKGRRMAYIDEGTG
DPILFQHGNPTSSYLWRNIMPHCAGLGR
LIACDLIGMGDSDKLDPSGPERYAYAEHR
DYLDALWEALDLGDRVVLVVHDWGSAL
GFDWARRHRERVQGIAYMEAIAMPIEW
ADFPEQDRDLFQAFRSQAGEELVLQD
```

## Structure



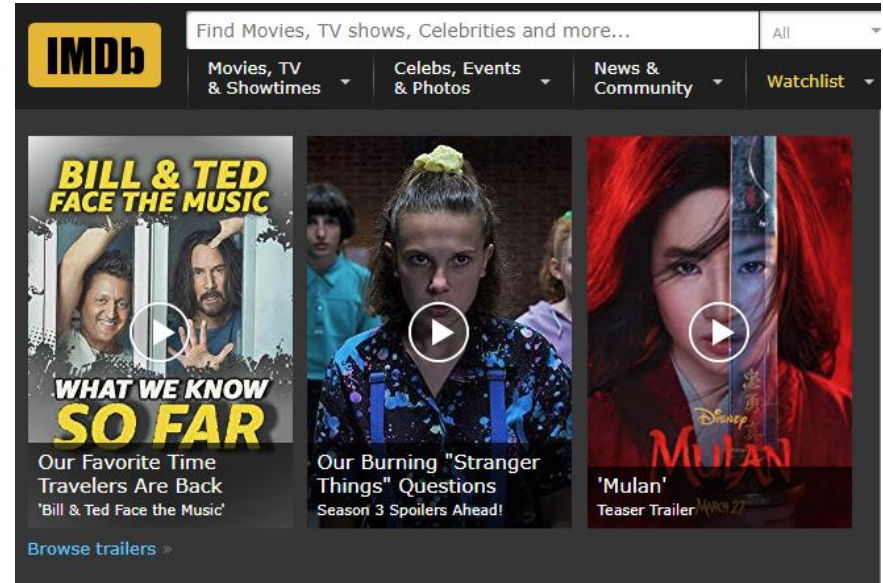## Evolution



**?** → Function

Main approaches:

- Experimental
- Rule-based
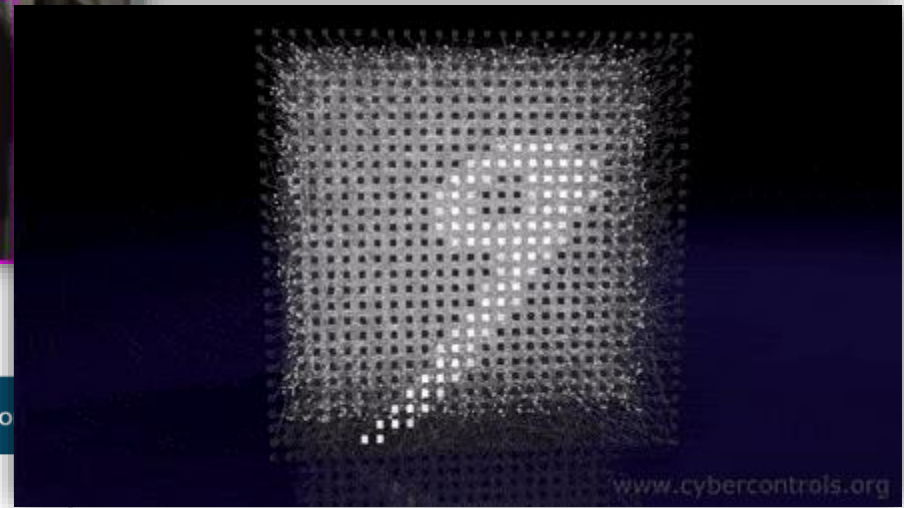- Machine learning

# Introduction to AI and ML

# Introduction to AI and ML

- Recommendation engines

- Gaming

- Image & speech recognition

- Anomaly detection
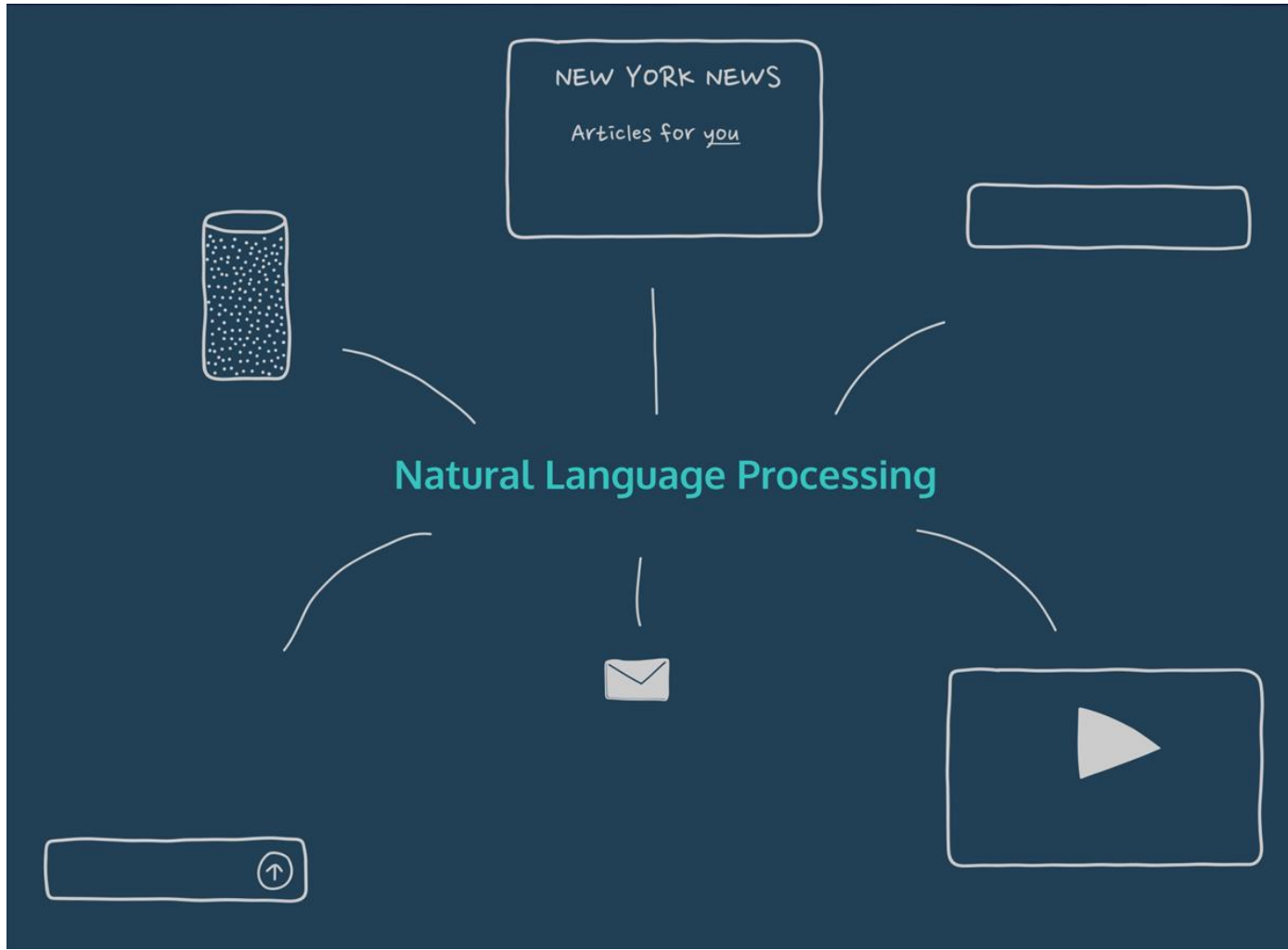
- Natural language processing

- Data mining

- …

| Login | Non- Transactional Activities | | Transactio |
|---|---|---|---|
| • Challenges | • View balance | • Add new user | • ACH |
| • Device | • View history | • Change limits | • Wire |
| • Cookie | • Updated address | • Set up batch | • Bill Pay |
| • IP Address | • Update email | • Set up template | • Loan Draw |
| • Time of day | • Update password | • Add payees | |
| • Network | | | |

www.cybercontrols.org

NEW YORK NEWS

Articles for you

**Natural Language Processing**

# Introduction to AI and ML



Passes: Center Defenders

# Modern challenges

# in Bioengineering

# Modern Challenges

free drug
ligand

free target
protein

bound drug-target
co-complex

## Chemical Space

**92M compounds**

**~9,600 Drugs**

## Protein Space

**20244 human proteins**

~6,200 3D structures

~2700 Drug targets

- **How to annotate:**

  - **protein structure**

  - **protein function**

  - **protein interactions**

  - **…**

- **How to increase:**

  - **enzymatic activity**

  - **stability & solubility**

  - **substrate specificity**

  - **enantioselectivity**



0.2 ns

# Modern Challenges



**March 2000** — Expert panel proposes cohort study of 500,000 adults.

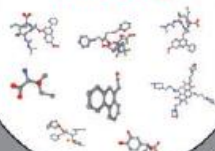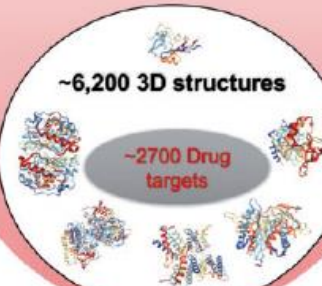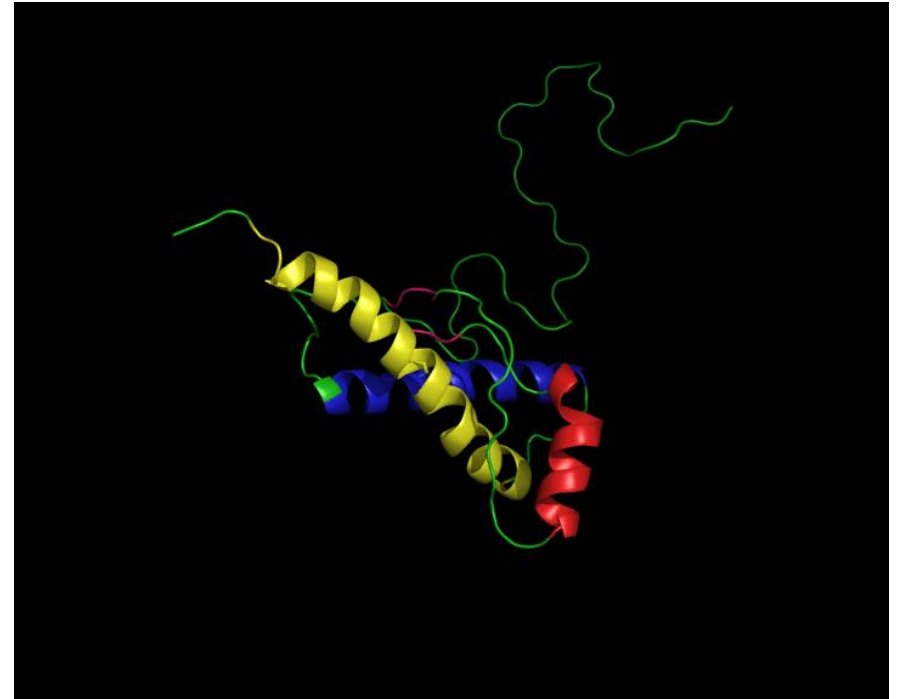**April 2002** — Wellcome Trust and U.K. government announce initial funding of £45 million.

**March 2012** — UK Biobank resource launches.

**May 2015** — Genotyping data on 150,000 released.

**October 2015** — Imaging data available for 5000.
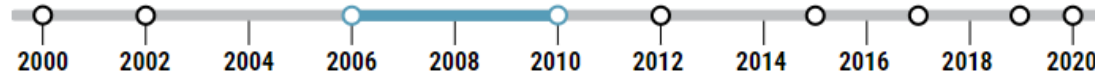
**July 2017** — Genotyping data on 500,000 released.

**March 2019** — Exome data on 50,000 to be released.
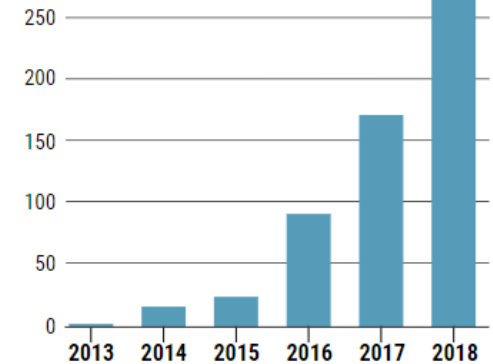
**2020** — All exome data released.

Recruitment of participants

N. DESAI/*SCIENCE*

## Engine of productivity

Published papers based on the UK Biobank's bounty of health and genetics data are piling up fast, in part because the data are freely available.
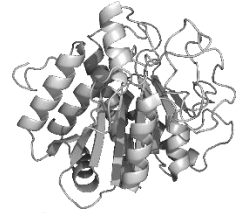
(GRAPHIC) N. DESAI/*SCIENCE*; (DATA) UK BIOBANK

UK Biobank Principal Investigator Rory Collins stands amid stored biospecimens from the project's half-million participants. NIGEL HILIER

# Huge trove of British biodata is unlocking secrets of depression, sexual orientation, and more

By Jocelyn Kaiser, Ann Gibbons | Jan. 3, 2019 , 1:20 PM

1.  **A large number of <u>relevant features</u>;**

    ▪ **From bio-physico-chemical to textual;**

    ▪ **Heterogeneous data (e.g. clinical, imaging, and genomic data);**

2.  **<u>Complex tasks</u> and large parameter space;**

    ▪ **A single 300-amino-acid-long protein will have 300·19=5700 single-point variants!**

3.  **<u>Large datasets</u> available and new data are collected.**

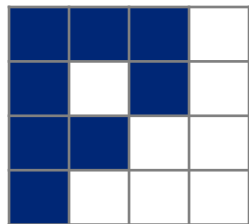**Ideal ML setup = complex tasks + relevant features + abundant data**

# Basics of ML

- **Feature vector:** $x = (x_1, x_2, \dots, x_n)$

- **All features must be converted to numbers:**

**yes =1, no = 0**

A C G T

1 0 0 0
0 0 0 1
0 0 1 0
0 0 0 1
1 0 0 0
0 1 0 0
0 0 0 1
0 0 1 0
1 0 0 0

A T G T A C T G A

One-hot encoding →

| 1 | 1 | 1 | 0 |
|---|---|---|---|
| 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 |

- **Unsupervised learning**

  - **only features are available;**

  - **goal: <u>cluster the data</u> or <u>reduce their dimensionality</u>;**

- **Supervised learning**

  - **features and <u>labels</u> are available;**

  - **goal: <u>learn to predict the label based on the features</u>;**

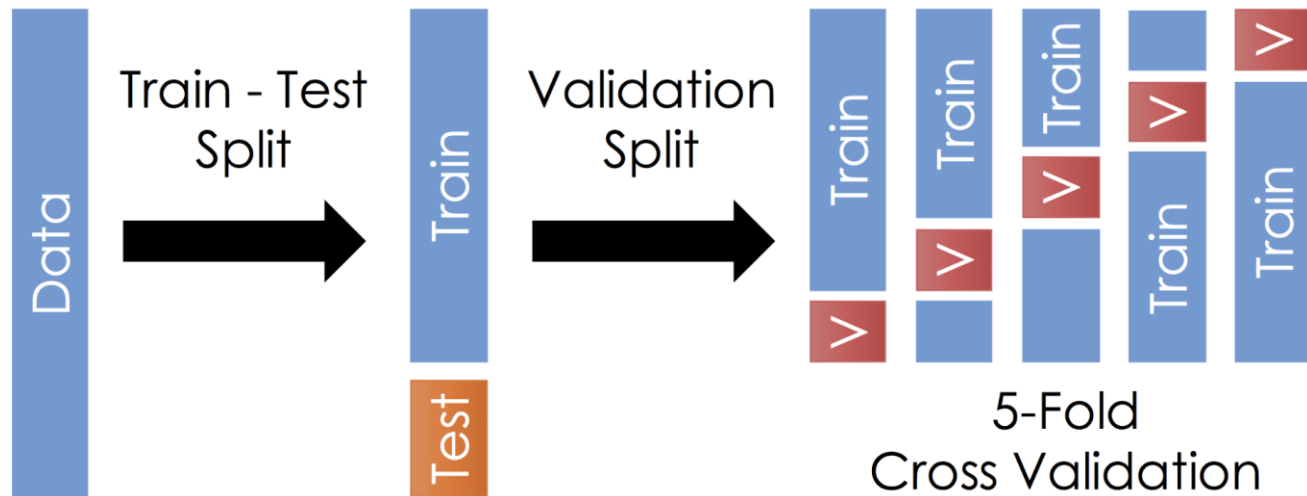- **Reinforcement learning, feature learning, anomaly detection, etc…**

# Basics of ML: unsupervised learning

- **We want ML models to be generalizable = good at predicting labels for <span style="color:red">previously unseen data</span>;**

- **It is essential to split the data into <u>training set</u> and <u>test sets</u> and use the latter for final evaluation only!**

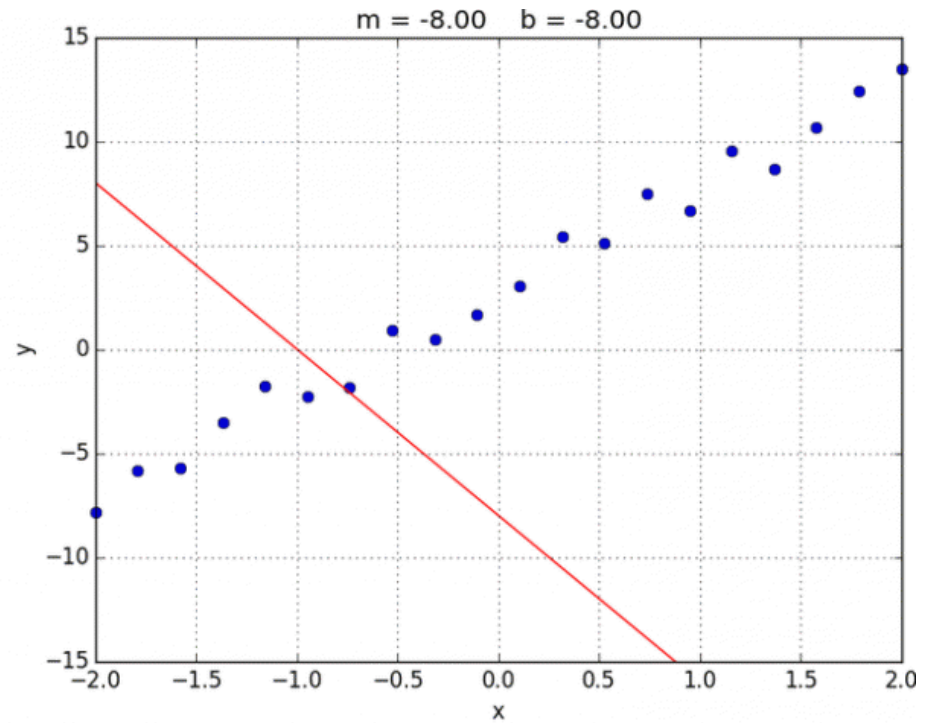- **To fine-tune an algorithm, <u>K-fold cross validation</u> is implemented:**

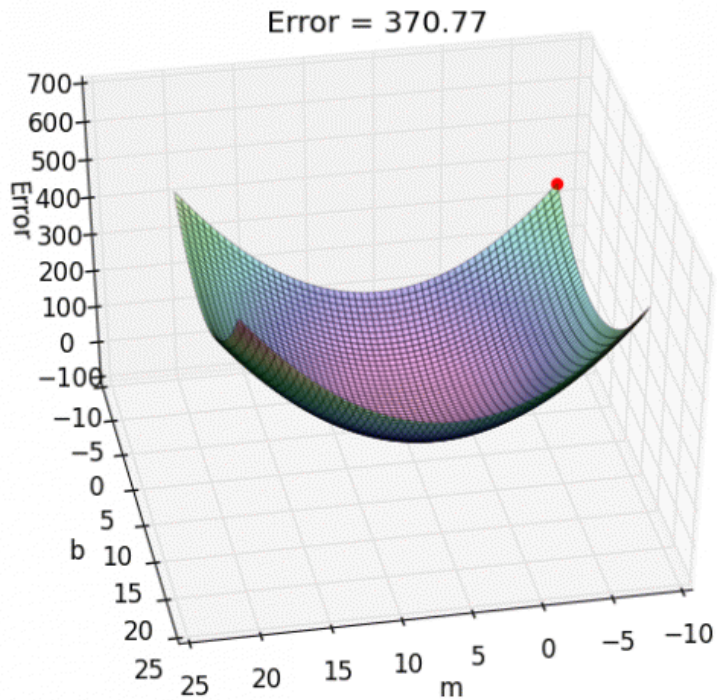- **Let $\left(x^{(i)}, y^{(i)}\right)$ be our data set, where $x$ are <u>feature values</u>, and $y$ are <u>labels</u>;**

- **Any ML predictor is fundamentally a function $f(x)$:**

$$f(\text{feature values}) = \text{label}$$

- **Usually, a generic group of functions $f(x, \beta)$ is chosen, where $\beta$ is a set of parameters;**

- **Then we "train" the ML predictor: pick such $\beta^*$ that $f\left(x^{(i)}, \beta^*\right)$ is as close to $y^{(i)}$ as possible.**

$$f(x) = m \cdot x + b$$

# Recent applications

Figure 1 A hypothetical example of how a decision tree might predict protein-protein interactions.

- **QSAR modeling for ligand binding activity;**

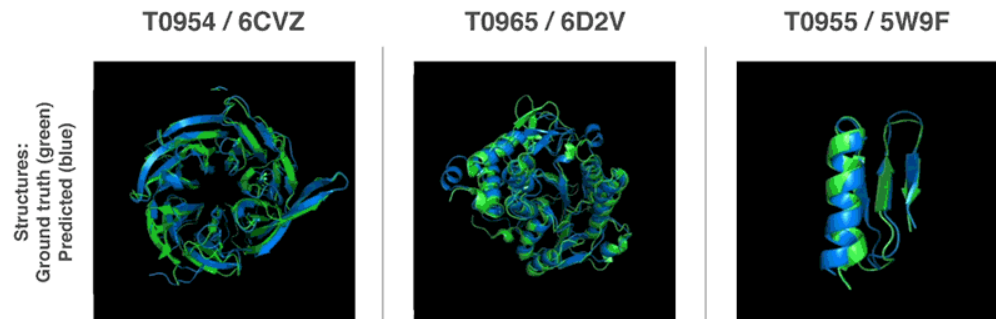- **Structure generation;**

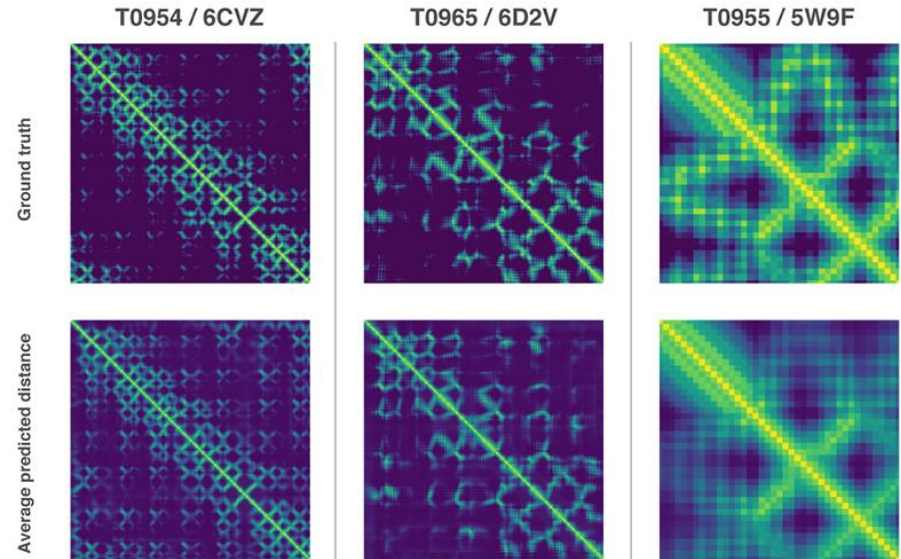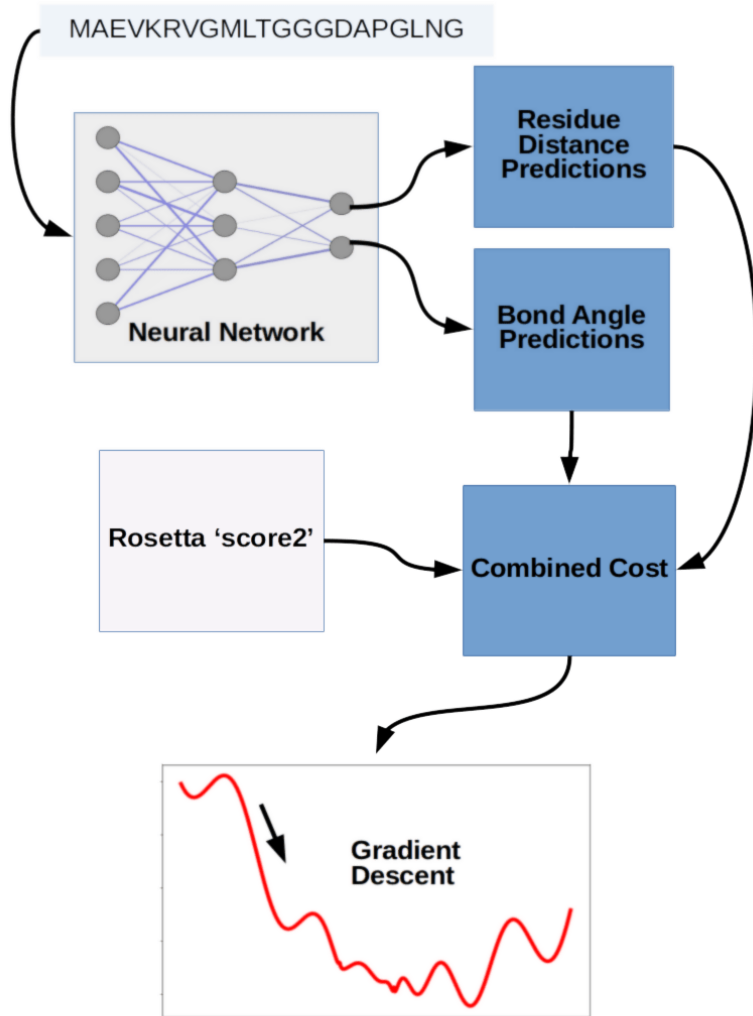- **Synthetic pathway generation;**



Figure 2: Illustration of the Chemception framework. After a SMILES to structure conversion, the 2D images are mapped onto an 80 x 80 image that serves as the input image data for training a deep neural network to predict toxicity, activity, and solvation properties.



**Figure 2.** Retrosynthetic reaction prediction task and an example of a possible retrosynthetic disconnection for a target molecule.



Goh et al. 2017 Chemception: a deep neural network with minimal chemistry knowledge matches the performance of expert-developed QSAR/QSPR models
Segler et al. 2018 Planning chemical syntheses with deep neural networks and symbolic AI. Nature

❑ **Eraslan G et al. 2019 "Deep learning: new computational modelling techniques for genomics." Nature Reviews Genetics. (especially pages 2-6)**



REVIEWS

## Deep learning: new computational modelling techniques for genomics

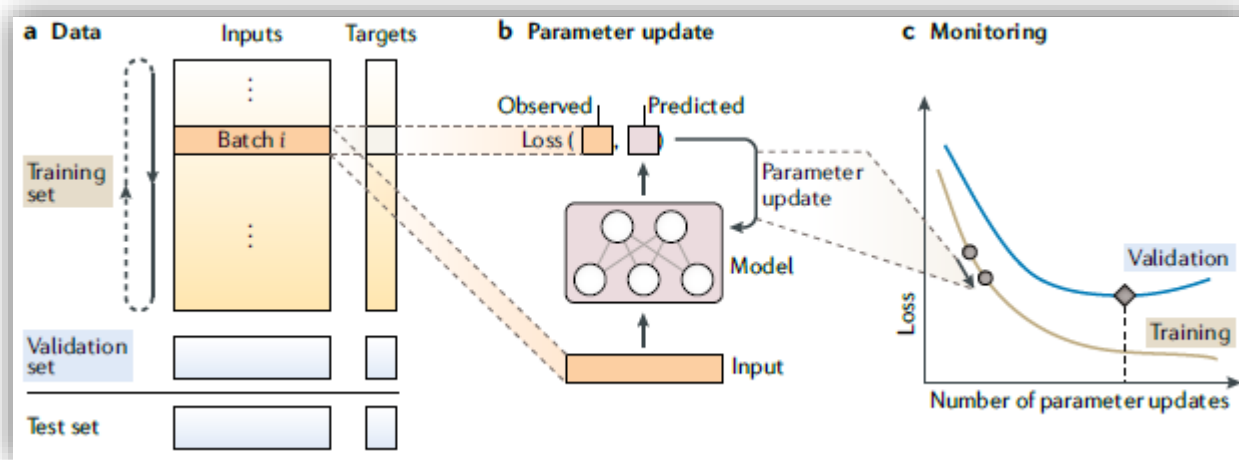Gökcen Eraslan [1,2,5], Žiga Avsec[3,5], Julien Gagneur[3]* and Fabian J. Theis [1,2,4]*

Abstract | As a data-driven science, genomics largely utilizes machine learning to capture dependencies in data and derive novel biological hypotheses. However, the ability to extract new insights from the exponentially increasing volume of genomics data requires more expressive machine learning models. By effectively leveraging large data sets, deep learning has transformed fields such as computer vision and natural language processing. Now, it is becoming the method of choice for many genomics modelling tasks, including predicting the impact of genetic variation on gene regulatory mechanisms such as DNA accessibility and splicing.