# Class core values

1. Be **respect**ful to yourself and others
2. Be **confident** and believe in yourself
3. Always do your **best**
4. Be **cooperative**
5. Be **creative**
6. Have **fun**
7. Be **patient** with yourself while you learn
8. Don't be shy to **ask "stupid" questions**

Week 8, Lecture 1

# Machine learning in protein engineering
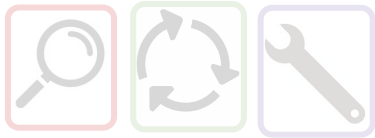
# Learning Objectives

1. Describe applications of machine learning in protein engineering
2. Identify challenges in the application of neural nets
3. Evaluate learning literature based on performance and application
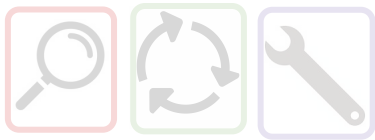
# Machine learning

# Machine learning

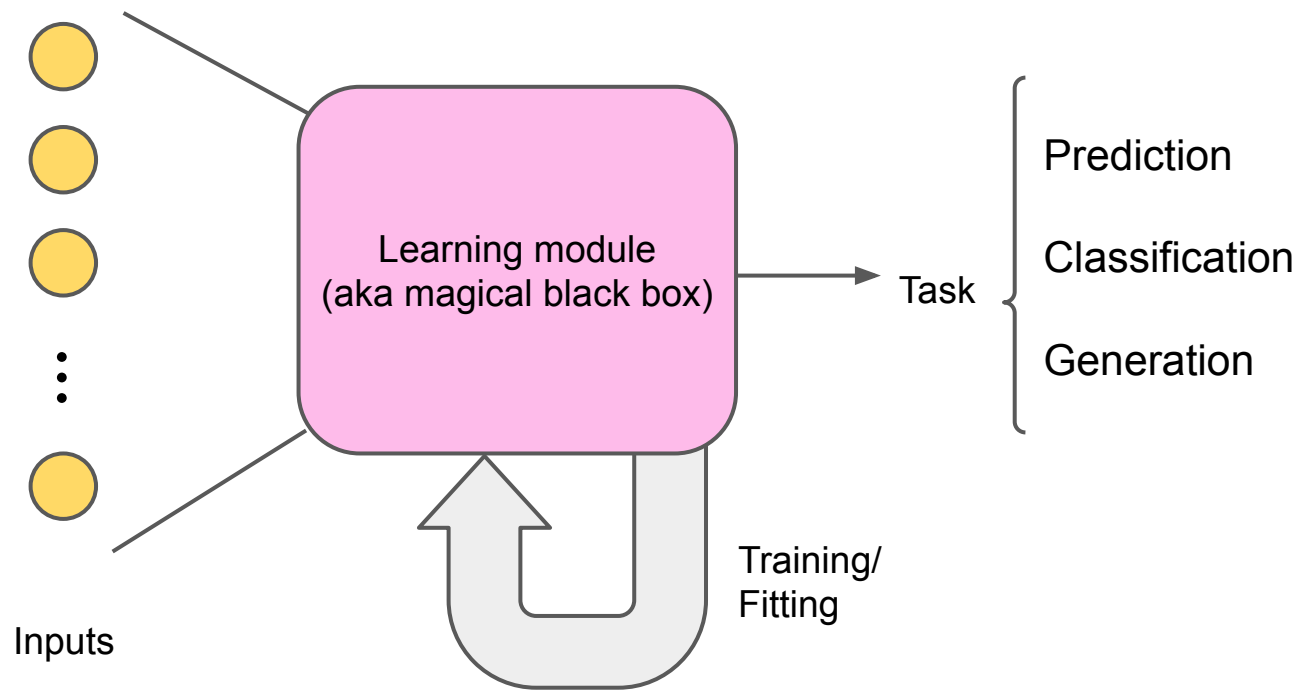learns from large databases for prediction, classification, generation

# Machine learning

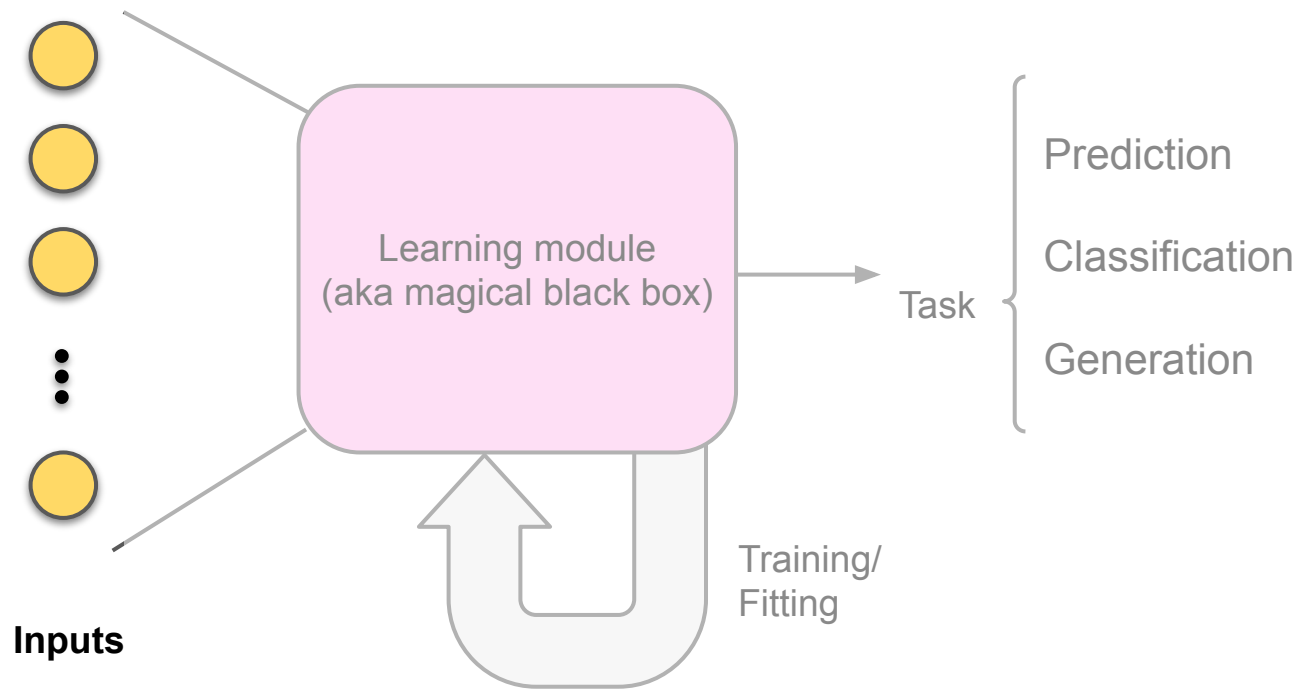learns from large databases for prediction, classification, generation



| Advantage | Disadvantage |
|---|---|
| The computational power is getting better | Needs large amount of data |
| It can uncover new patterns and rules that we can't | Can't be specific/target-based (great for general cases) |
| Fully reproducible | Can sample new sequences but in the neighborhood of what exist |
| Can take in multiple considerations into account | |

# Basic components of a learning module

# The inputs define the model to use



Learning module
(aka magical black box)

Inputs

Task

Prediction

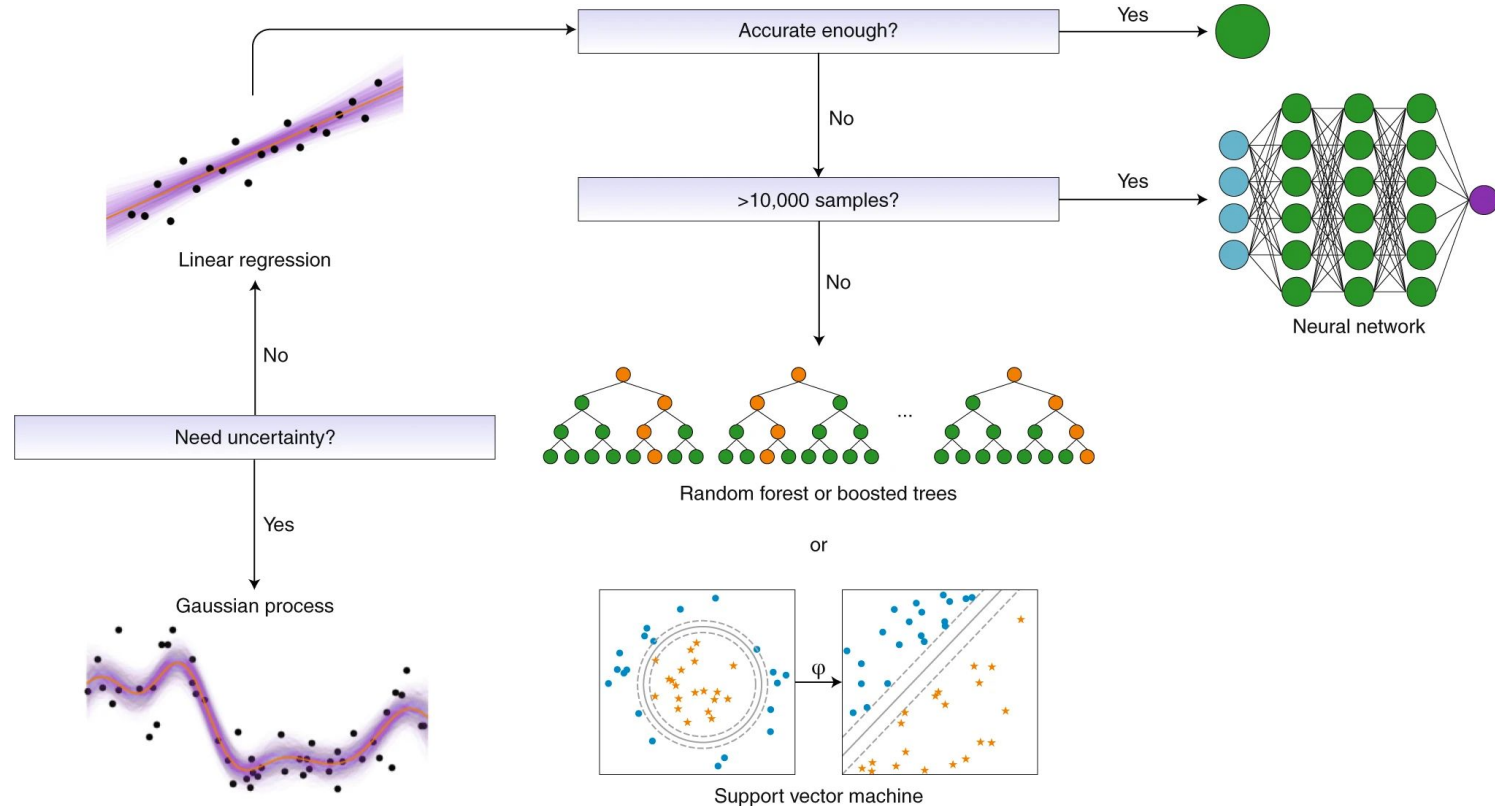Classification

Generation

Training/
Fitting

# What to check for in inputs

1. How much data do I have?

# The number of inputs defines types of ML model

# What to check for in inputs
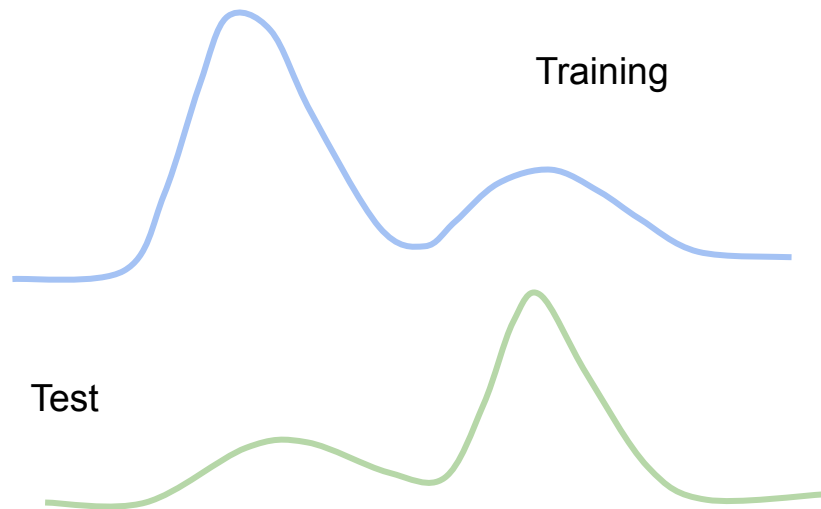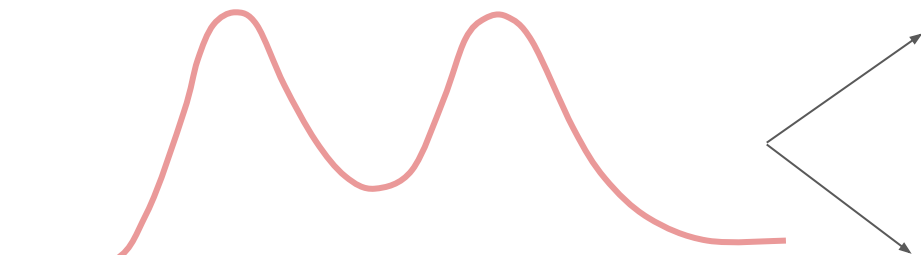
1. How much data do I have?
2. What is my data type?
   a. Sequence
   b. Structure
   c. Image
   d. …

# What to check for in inputs

1. How much data do I have?
2. What is my data type?
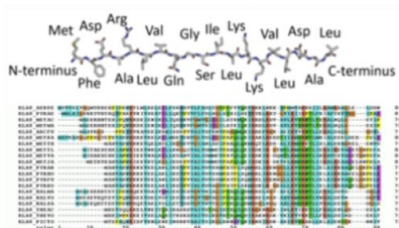3. How much noise do I have in my data?
   a. How can I clean it?

# What to check for in inputs

1. How much data do I have?
2. What is my data type?
3. How much noise do I have in my data?
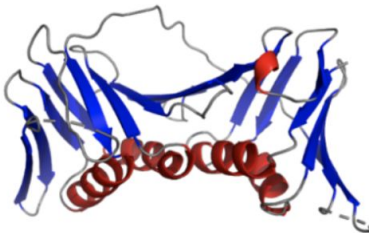4. What is my data distribution?



Training

Test

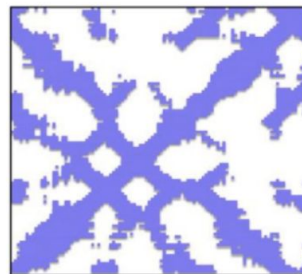# Proteins can be represented in many different ways



Sequence / MSA profile
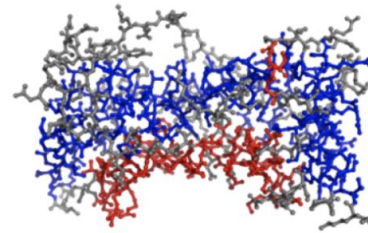many classical ML methods

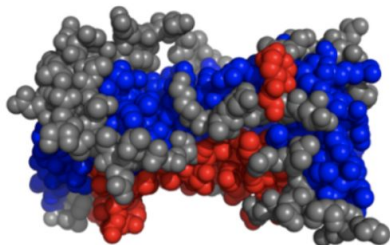Secondary structure elements

Distance / HB / Contact matrix
state-of-the-art structure prediction methods
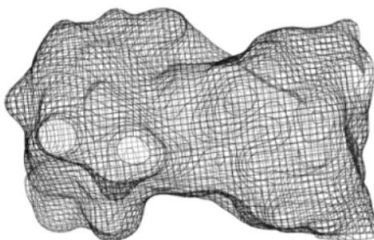
Molecular graph
Fout et al. NIPS 2017; Ingraham et al. NIPS 2019; Baldassarre et al. Bioinformatics 2020; Igashov et al. 2020
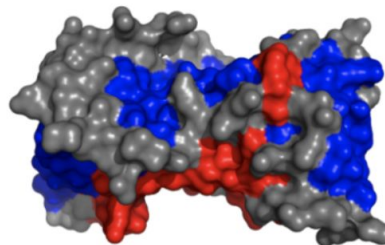
Set of balls / Point cloud
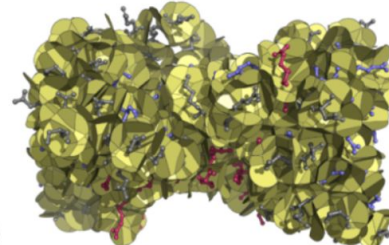classical statistical potentials; Eismann et al. 2020

Gaussian clouds
Derevyanko et al. Bioinformatics 2018; Pages et al. Bioinformatics 2019;
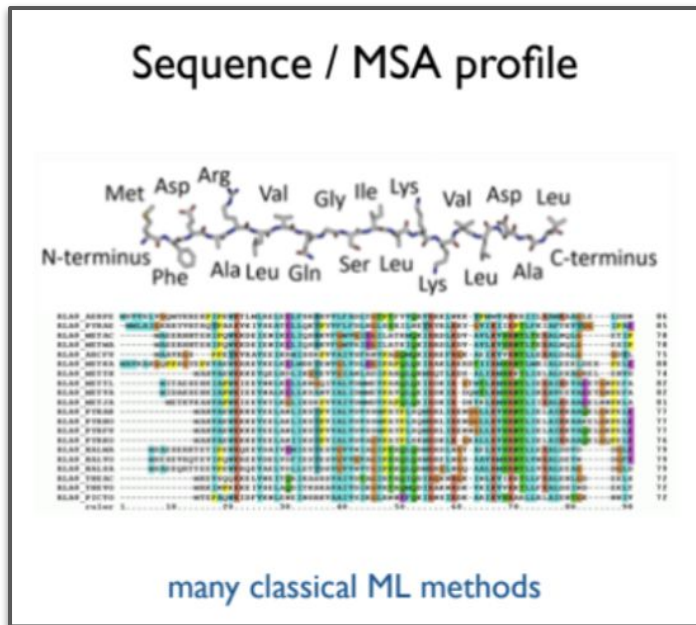
Molecular surface
Olechnovic & Venclovas, Proteins 2017; Correia, Bronstein et al. Nat Met 2020

3D tessellation
Igashov et al. Bioinformatics 2021; Olechnovic et al. Proteins 2021

# Each representation fits a different model better



Sequence / MSA profile

...ture elements    Distance / HB / Contact matrix    Molecular graph

Met Asp Arg Val Gly Ile Lys Val Asp Leu

N-terminus Phe Ala Leu Gln Ser Leu Lys Leu Ala C-terminus

many classical ML methods

...clouds

...et al. NIPS
...atics 2020;

Classical learning methods
    Support vector machine
    Random forests

Deep learning methods
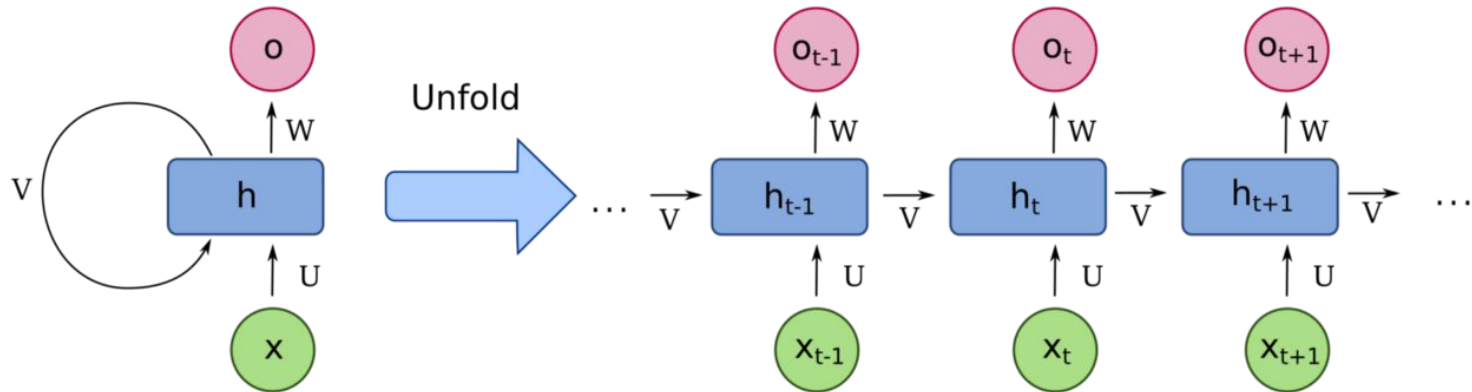    RNNs
    LSTMs
    Transformers

classical statistical potentials; Eismann et al. 2020

Derevyanko et al. Bioinformatics 2018; Pages et al. Bioinformatics 2019;

Olechnovic & Venclovas, Proteins 2017; Correia, Bronstein et al. Nat Met 2020

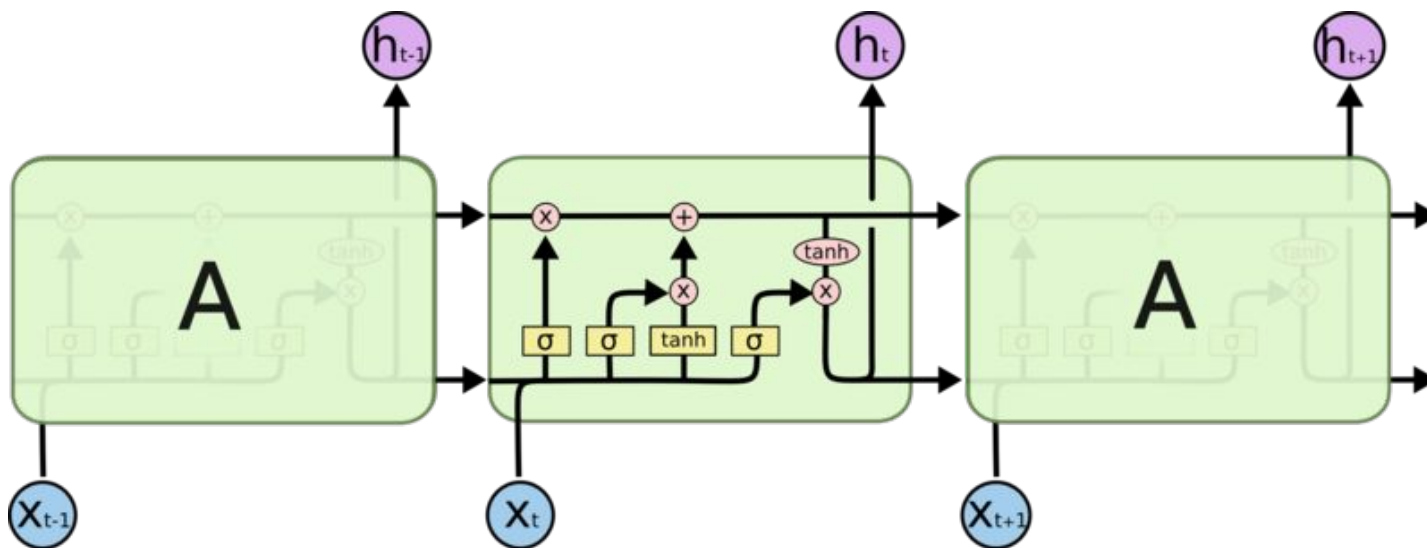Igashov et al. Bioinformatics 2021; Olechnovic et al. Proteins 2021

# Sequences are natural inputs for NLP (natural language processing) models
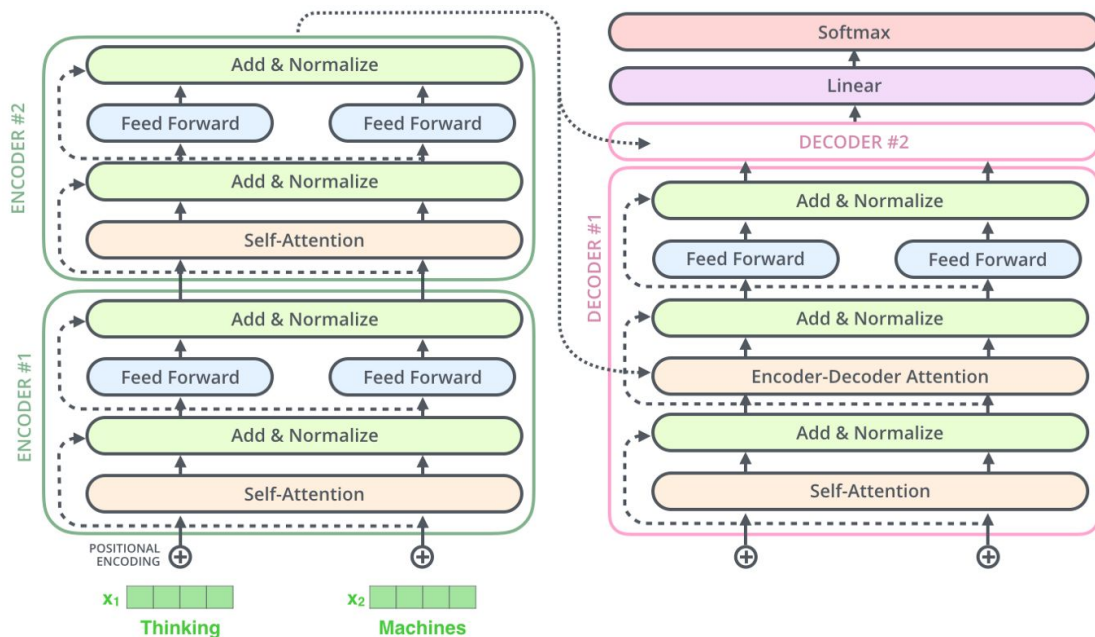


**R**ecurrent **N**eural **N**ets (RNNs)

# Sequences are natural inputs for NLP (natural language processing) models



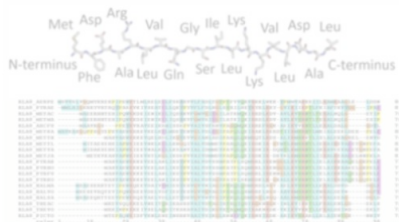**L**ong **S**hort **T**erm **M**emory (LSTM)

# Sequences are natural inputs for NLP (natural language processing) models
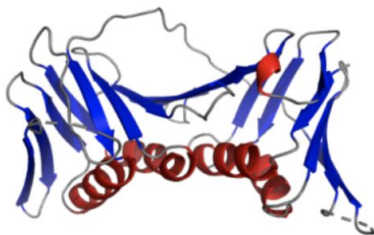


Transformer

# Each representation fits a different model better



Sequence / MSA profile

many classical ML methods

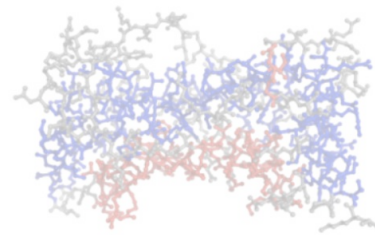**Secondary structure elements**

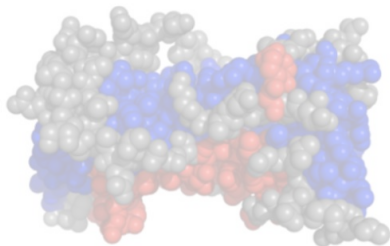Distance / HB / Contact matrix

state-of-the-art structure prediction methods
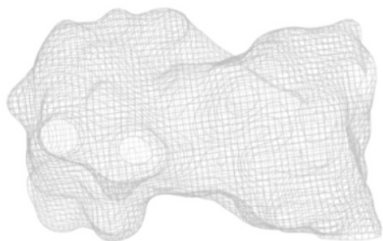
Molecular graph

Fout et al. NIPS 2017; Ingraham et al. NIPS 2019; Baldassarre et al. Bioinformatics 2020; Igashov et al. 2020
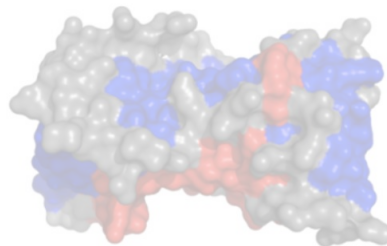
Set of balls / Point cloud

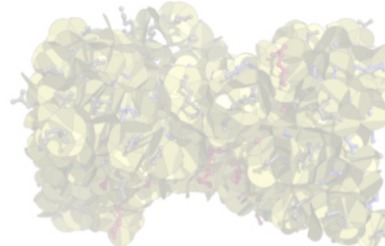classical statistical potentials; Eismann et al. 2020

Gaussian clouds

Derevyanko et al. Bioinformatics 2018; Pages et al. Bioinformatics 2019;

Molecular surface

Olechnovic & Venclovas, Proteins 2017; Correia, Bronstein et al. Nat Met 2020

3D tessellation

Igashov et al. Bioinformatics 2021; Olechnovic et al. Proteins 2021

# Each representation fits a different model better



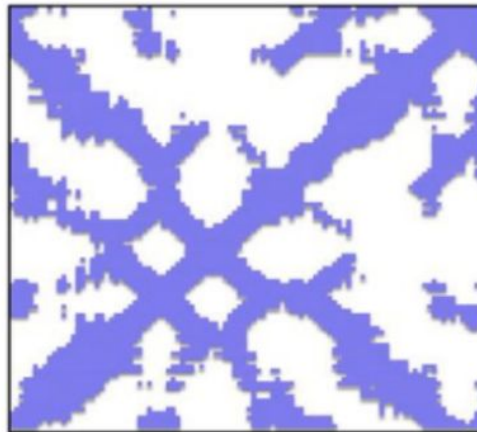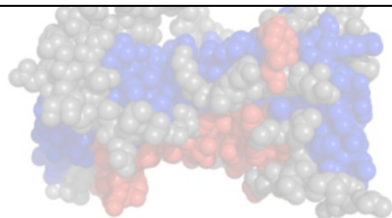Sequence / MSA profile    Secondary structure elements

Distance / HB / Contact matrix

state-of-the-art structure prediction methods

Classical learning methods
    Support vector machine
    Random forests

Deep learning methods
    CNNs

classical statistical potentials; Eismann et al. 2020
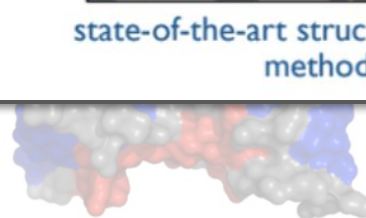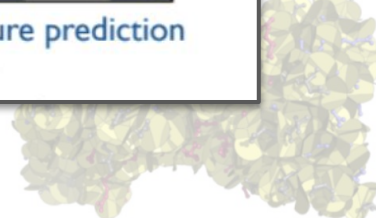
Derevyanko et al. Bioinformatics 2018; Pages et al. Bioinformatics 2019;
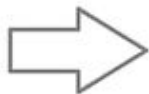
Olechnovic & Venclovas, Proteins 2017; Correia, Bronstein et al. Nat Met 2020

Igashov et al. Bioinformatics 2021; Olechnovic et al. Proteins 2021

# Convolutional neural nets (CNNs) are the preferred methods for images

# Convolutional neural nets (CNNs) are the preferred methods for images

| 1 | 1 | 0 |
|---|---|---|
| 4 | 2 | 1 |
| 0 | 2 | 1 |

⇨

| |
|---|
| 1 |
| 1 |
| 0 |
| 4 |
| 2 |
| 1 |
| 0 |
| 2 |
| 1 |

| $1_{\times1}$ | $1_{\times0}$ | $1_{\times1}$ | 0 | 0 |
|---|---|---|---|---|
| $0_{\times0}$ | $1_{\times1}$ | $1_{\times0}$ | 1 | 0 |
| $0_{\times1}$ | $0_{\times0}$ | $1_{\times1}$ | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| 4 | | |
|---|---|---|
| | | |
| | | |

Convolved Feature

# Each representation fits a different model better

Sequence / MSA profile    Secondary structure elements    Distance / HB /

Classical learning methods
 Support vector machine
 Random forests

Deep learning methods
 GCNs
 Message passing
 Link/node prediction

state-of-the-art str
meth

Molecular

## Molecular graph



Fout et al. NIPS 2017; Ingraham et al. NIPS 2019; Baldassarre et al. Bioinformatics 2020; Igashov et al. 2020

classical statistical potentials; Eismann et al. 2020

Derevyanko et al. Bioinformatics 2018; Pages et al. Bioinformatics 2019;

Olechnovic & Venclovas, Proteins 2017; Correia, Bronstein et al. Nat Met 2020

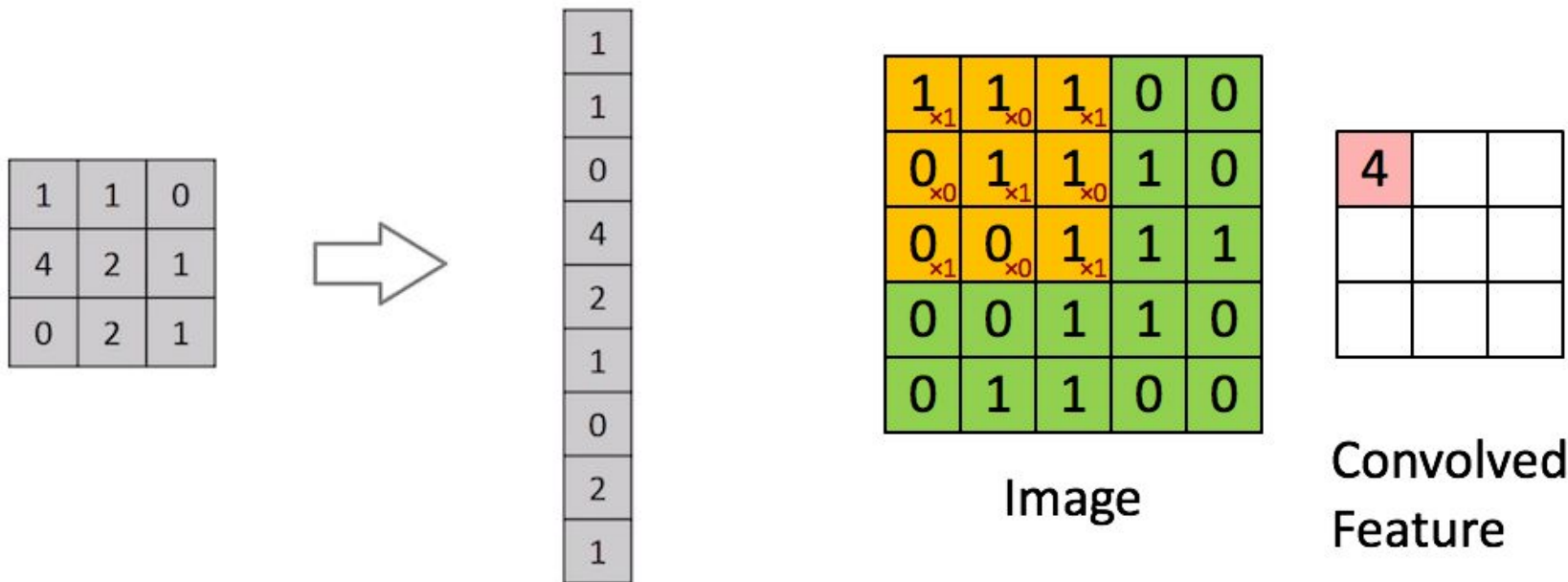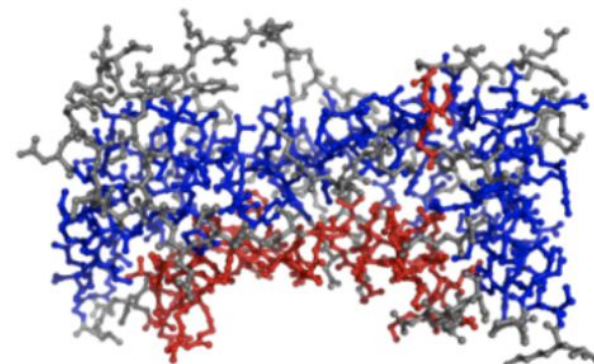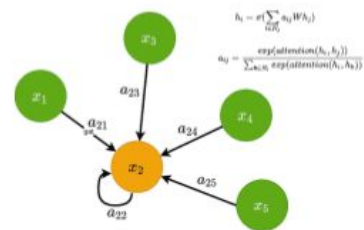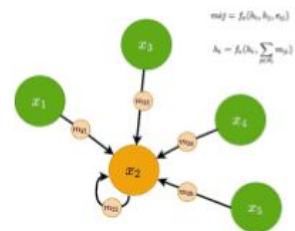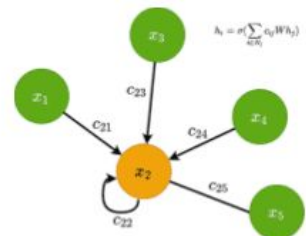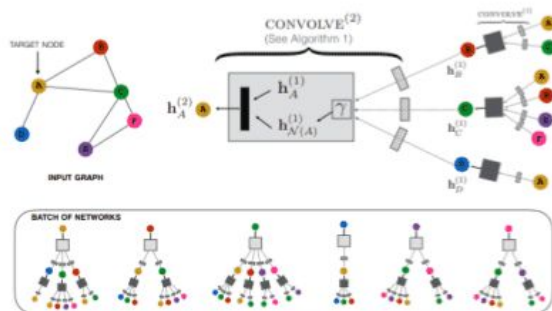Igashov et al. Bioinformatics 2021; Olechnovic et al. Proteins 2021

# The underlying idea behind convolution on graphs is message passing

# Each representation fits a different model better



Sequence / MSA profile   Secondary structure elements   Distance / HB / Contact matrix   Molecular graph

**Classical learning methods**
Support vector machine
Random forests

**Deep learning methods**
Geodesic learning
Point clouds

many classical ML methods

state-of-the-art structure prediction methods

Fout et al. NIPS 2017; Ingraham et al. NIPS 2019; Baldassarre et al. Bioinformatics 2020; Igashov et al. 2020

Set of balls / Point cloud   Gaussian clouds   Molecular surface   3D tessellation

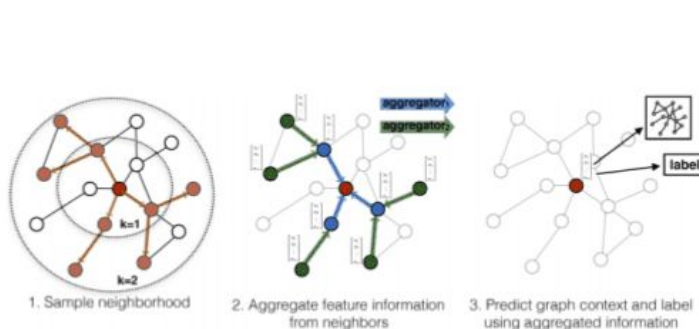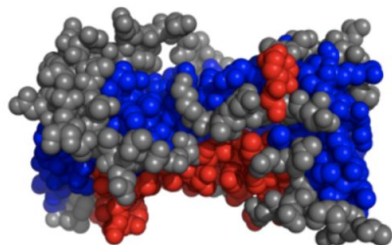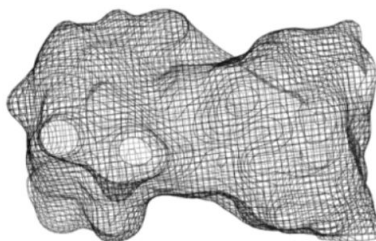classical statistical potentials; Eismann et al. 2020

Derevyanko et al. Bioinformatics 2018; Pages et al. Bioinformatics 2019;

Olechnovic & Venclovas, Proteins 2017; Correia, Bronstein et al. Nat Met 2020

Igashov et al. Bioinformatics 2021; Olechnovic et al. Proteins 2021

# Basic components of a learning module

# ML applications in protein engineering

# ML applications in protein engineering

Machine learning guided directed evolution

Representation learning

Unsupervised variant prediction

Generative models

**Predicting stability**

Predicting structure from sequence

Predicting interactions with other molecules

Classification and annotation

Predicting sequence from structure

# 3D CNN for prediction of stability

# ML-based models for the specific case of GPCRs



Energy Calculations from LitiCon Ensemble
Mutate Structural model
Generate multiple conformations from helix rigid-body sampling and Scrwl sidechain reassignment
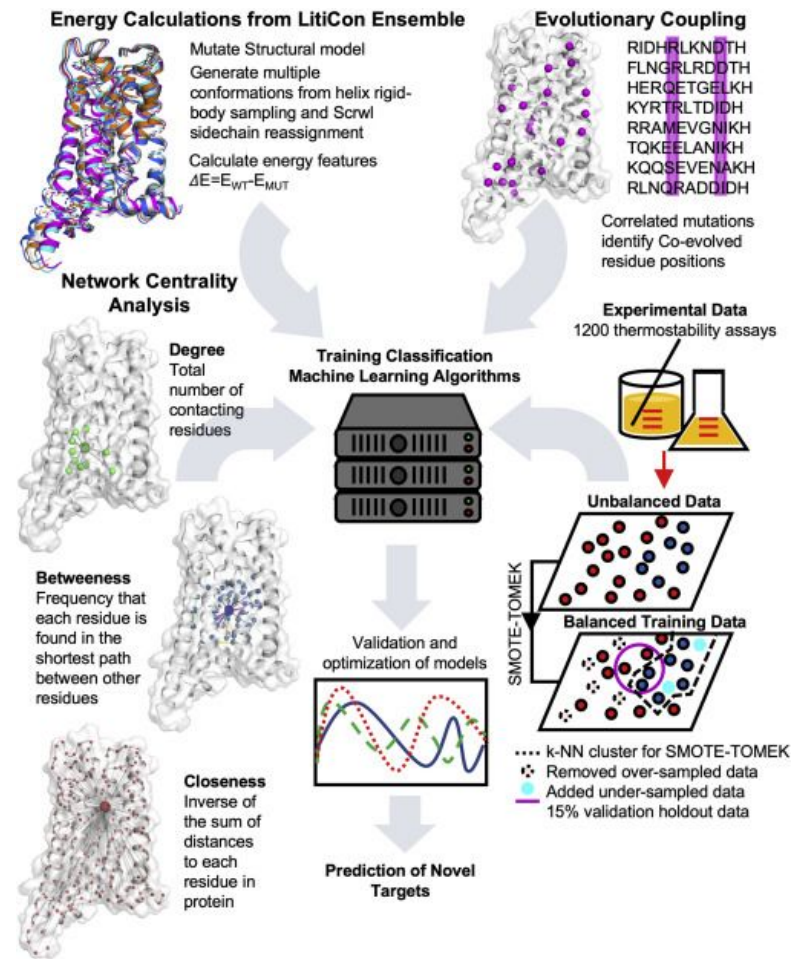Calculate energy features $\Delta E = E_{WT} - E_{MUT}$

Evolutionary Coupling
RIDHRLKNDTH
FLNGRLRDDTH
HERQETGELKH
KYRTRLTDIDH
RRAMEVGNIKH
TQKEELANIKH
KQQSEVENAKH
RLNQRADDIDH
Correlated mutations identify Co-evolved residue positions

Experimental Data
1200 thermostability assays

Network Centrality Analysis

Degree
Total number of contacting residues

Training Classification Machine Learning Algorithms

Betweeness
Frequency that each residue is found in the shortest path between other residues

Unbalanced Data

Balanced Training Data

Validation and optimization of models

SMOTE-TOMEK

Closeness
Inverse of the sum of distances to each residue in protein

Prediction of Novel Targets

···· k-NN cluster for SMOTE-TOMEK
Removed over-sampled data
Added under-sampled data
15% validation holdout data

# ML applications in protein engineering

Generative models

Machine learning guided directed evolution

Predicting interactions with other molecules

Representation learning

Predicting stability

**Classification and annotation**

Unsupervised variant prediction

Predicting structure from sequence

Predicting sequence from structure

# Using GCNs for function prediction

# Deep sequence model for protein classification

# ML applications in protein engineering

**Machine learning guided directed evolution**

Generative models

Predicting interactions with other molecules

Representation learning

Predicting stability

Classification and annotation

Unsupervised variant prediction
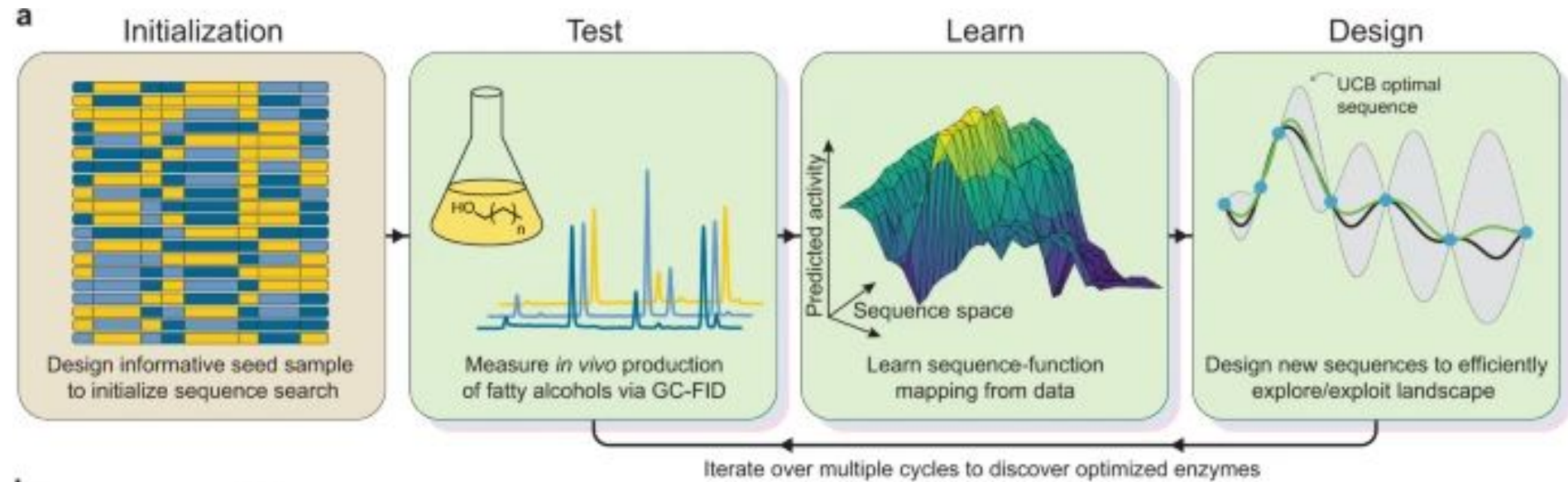
Predicting structure from sequence

Predicting sequence from structure

# Machine learning assisted directed evolution

# low-N protein engineering

# ML-accelerated protein sequence optimization
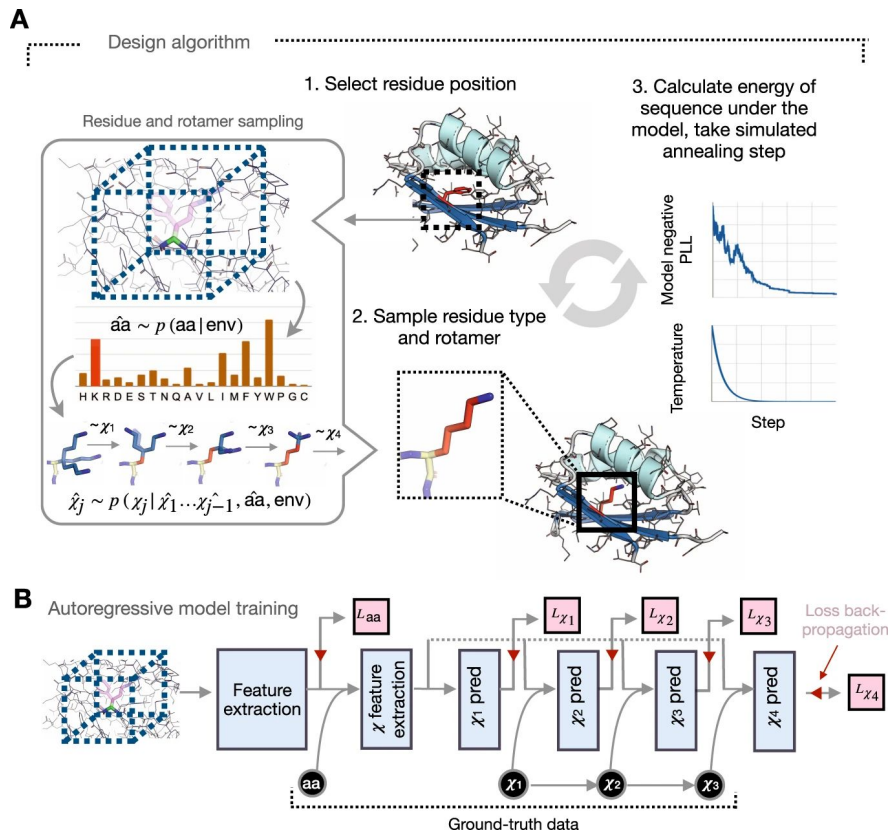


a

**Initialization** — Design informative seed sample to initialize sequence search

**Test** — Measure *in vivo* production of fatty alcohols via GC-FID

**Learn** — Predicted activity / Sequence space — Learn sequence-function mapping from data

**Design** — UCB optimal sequence — Design new sequences to efficiently explore/exploit landscape

Iterate over multiple cycles to discover optimized enzymes

b

# ML applications in protein engineering

Machine learning guided directed evolution

Generative models

Predicting interactions with other molecules

Representation learning

Predicting stability

Classification and annotation

Unsupervised variant prediction

Predicting structure from sequence

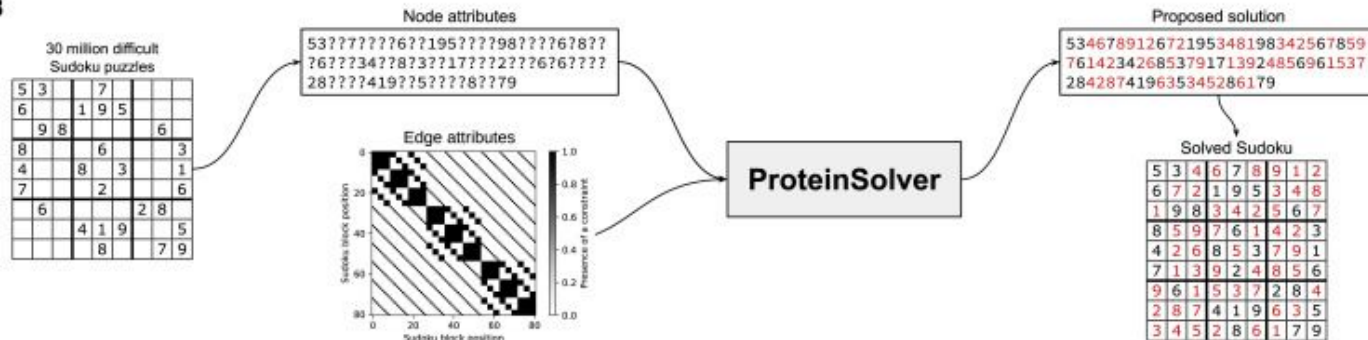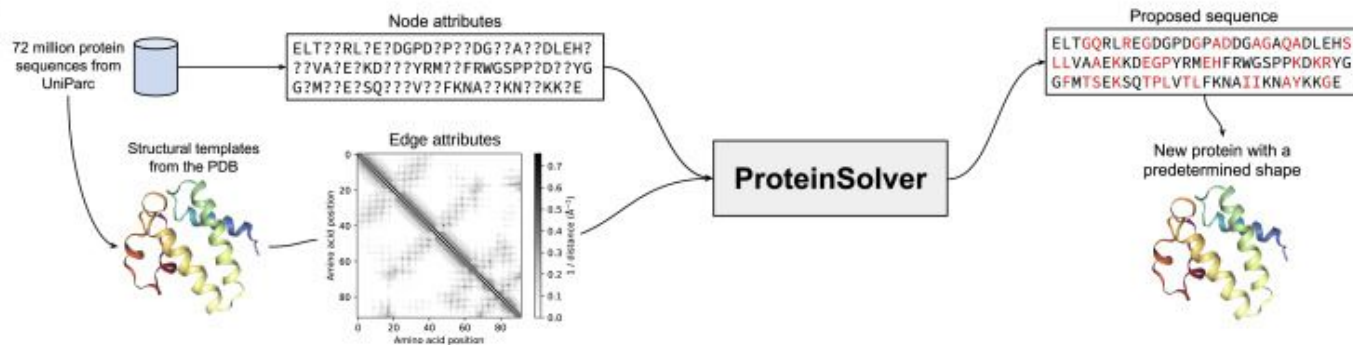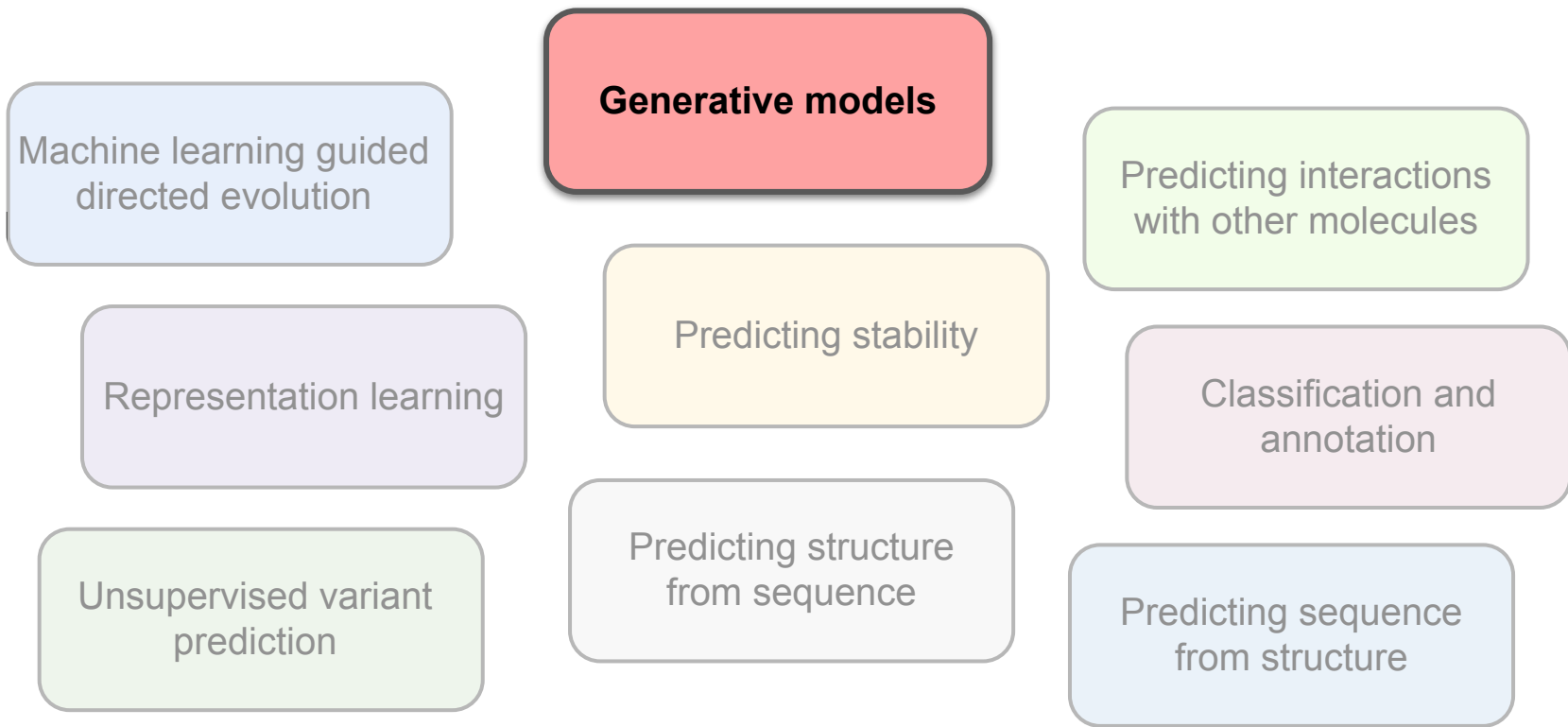**Predicting sequence from structure**

# Sequence design with a learned potential

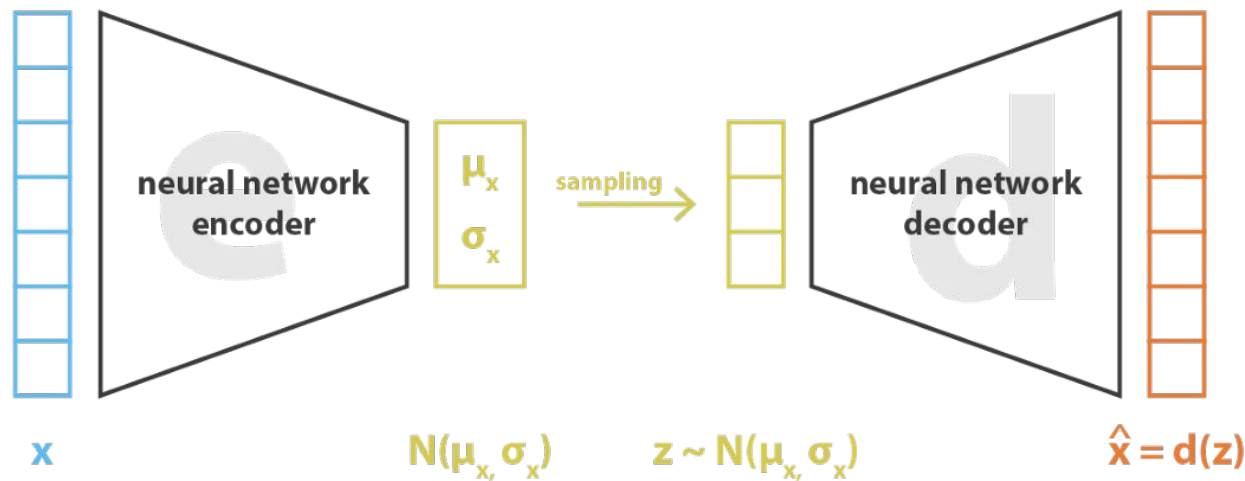# Autoregressive models

# Deep GCNs to design protein sequences

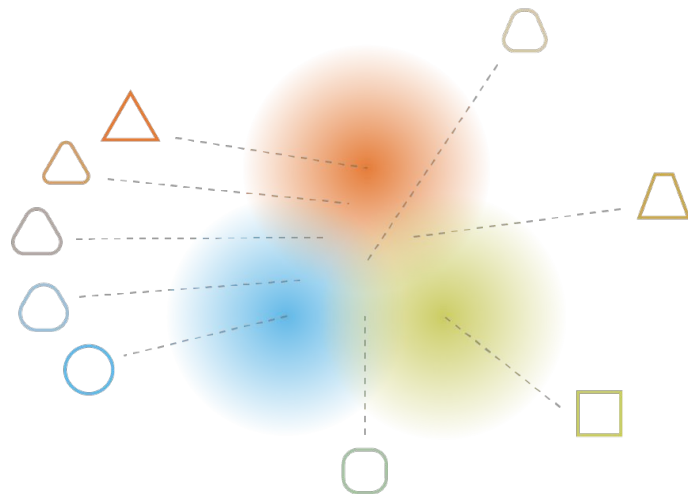# ML applications in protein engineering

Machine learning guided directed evolution

**Generative models**

Predicting interactions with other molecules

Representation learning

Predicting stability

Classification and annotation

Unsupervised variant prediction

Predicting structure from sequence

Predicting sequence from structure

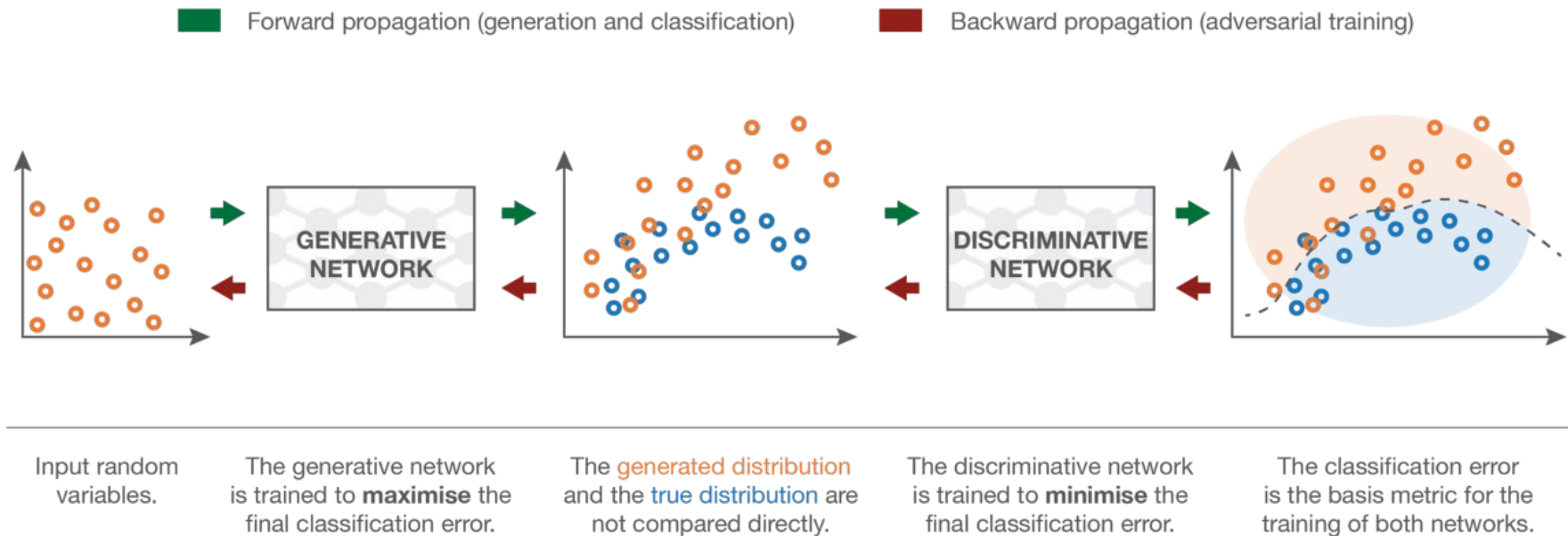# Variational Autoencoders (VAEs)

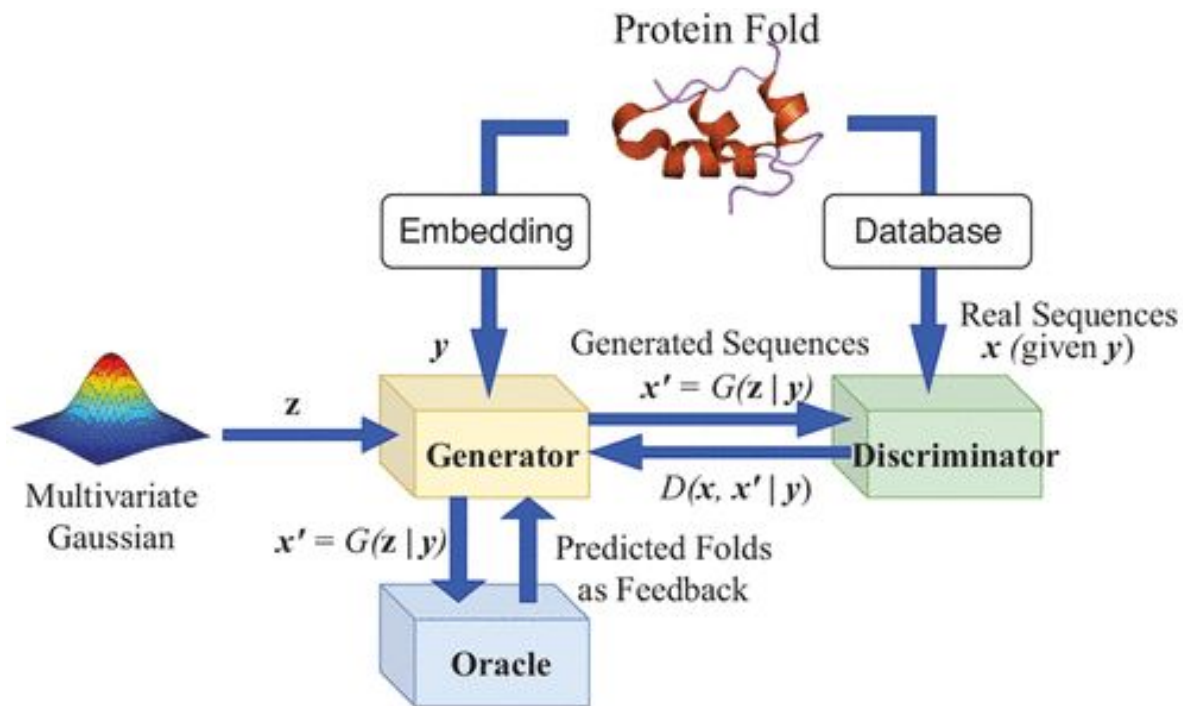# Variational Autoencoders (VAEs)

# VAEs for protein design

# Transformers for translating from protein to sequence space

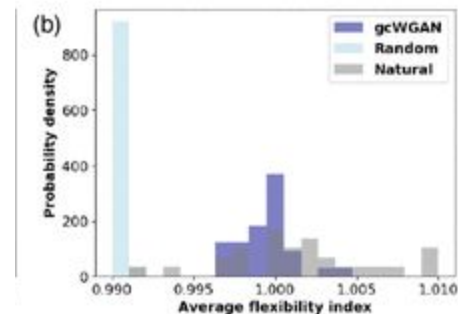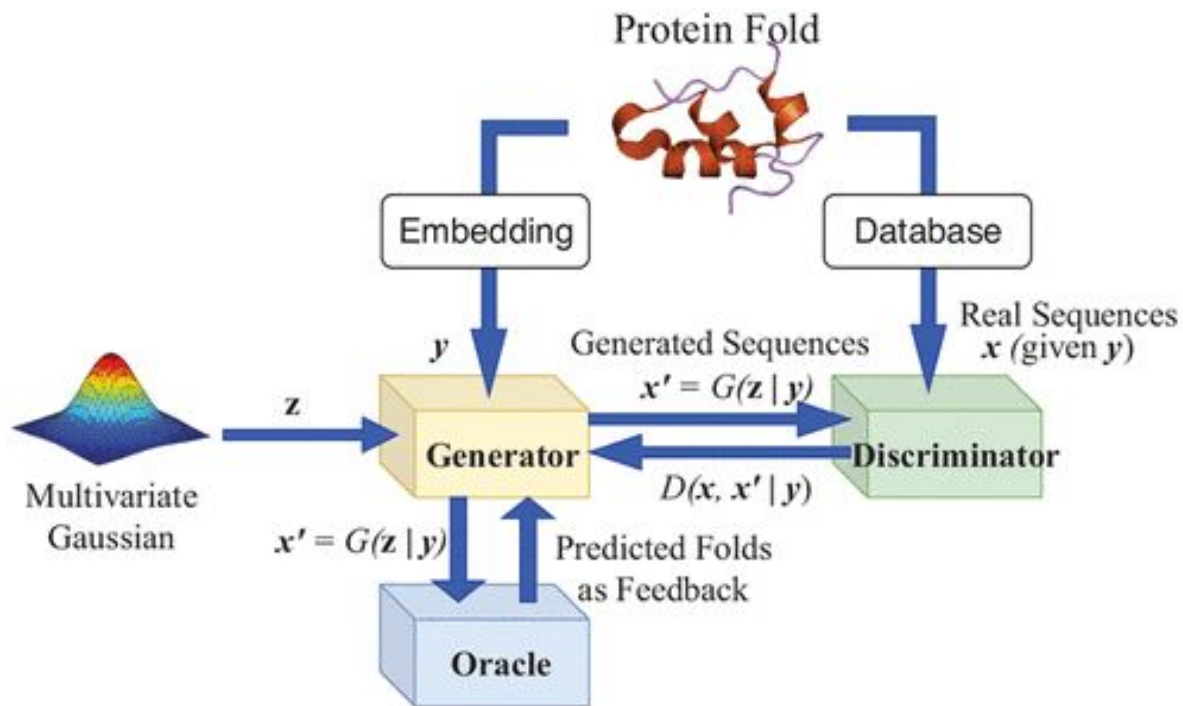# Generative Adversarial Networks (GANs)
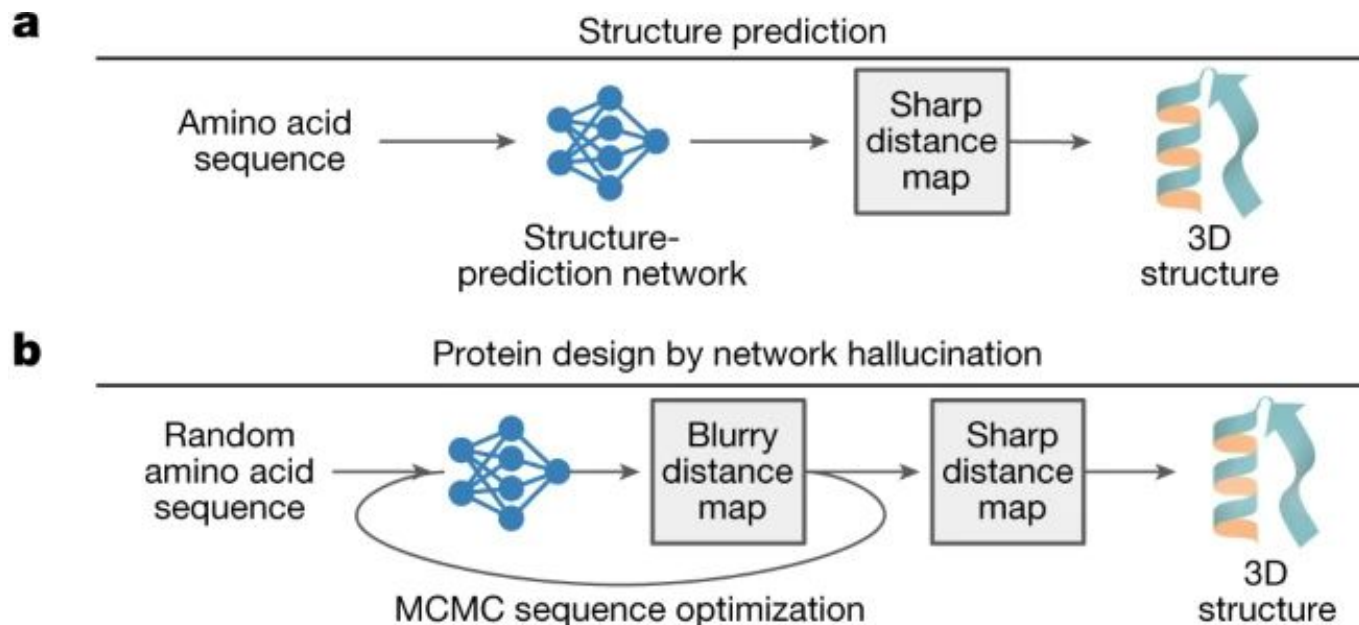
# GANs for protein design
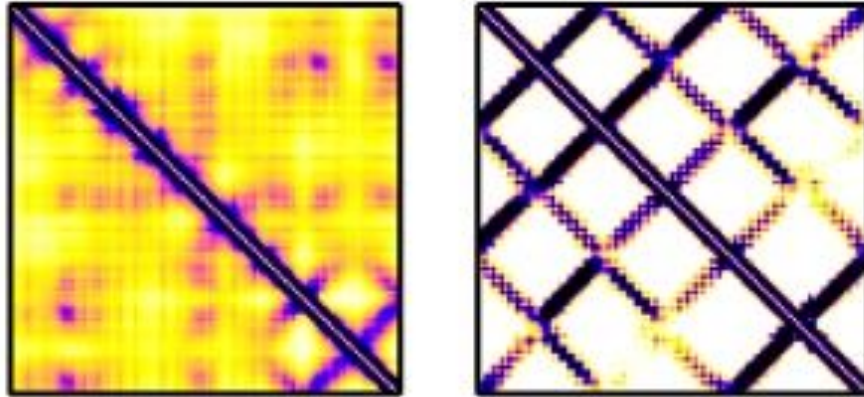
# GANs for protein design

# Protein design by reversing protein structure prediction

# In class activity:

Design new proteins using deep hallucination

# For the next lecture:

1. Read journal for the next lecture
   a. Moderated by **group IV**


2. W7L2 assignment due next lecture

# Next lecture:

## *Design of proteins with functional sites*