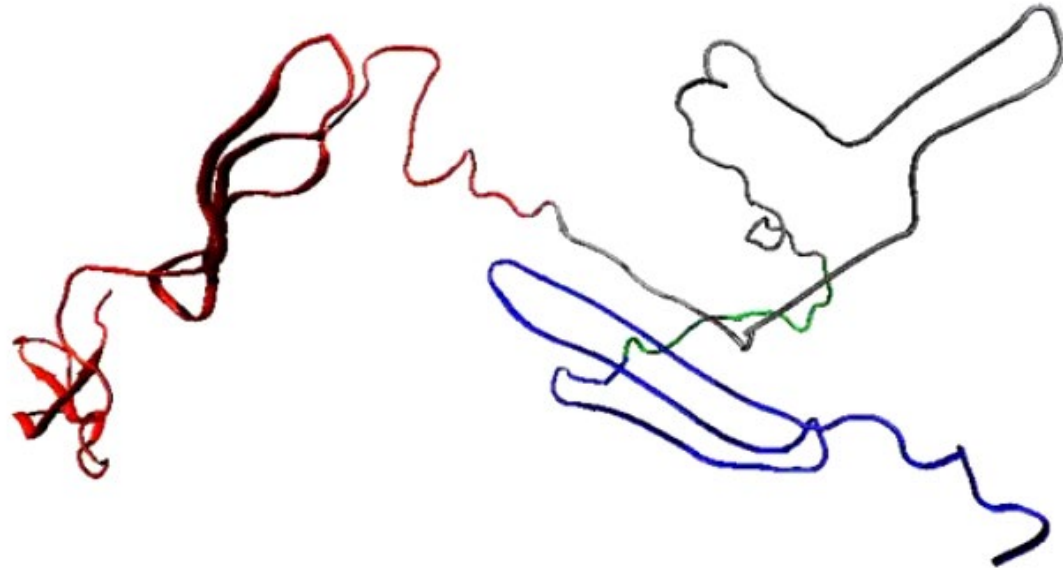# CS 189/289

Some applications of AI in biology:
1. protein structure prediction
2. protein design

# CS 189/289

Some applications of AI in biology:
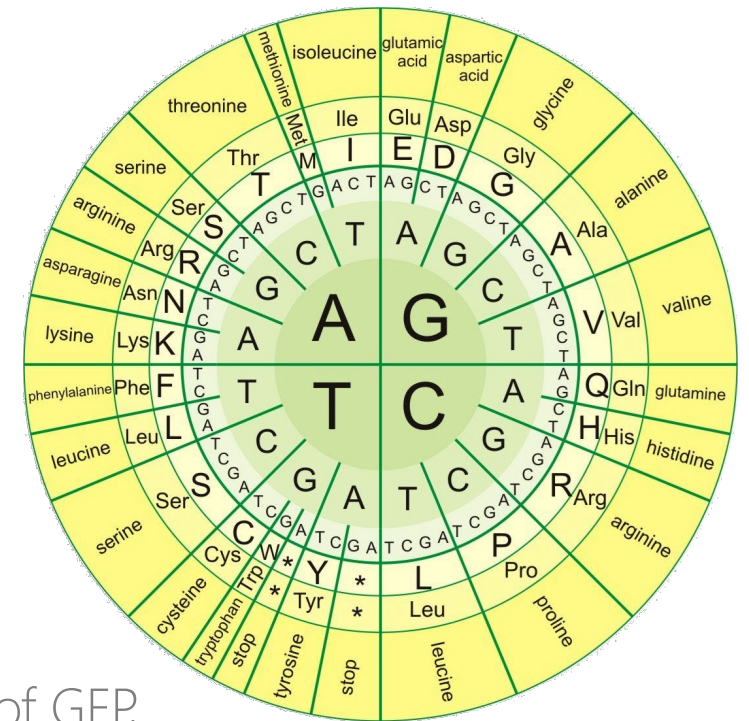
1. protein structure prediction
2. protein design

# Proteins are strings of nucleotides

238 length amino acid sequence:
MSKGEELFTGVVPILVELDGDVNGHKFSVSG
EDFFKS...NSHNVYIMADKQKNGIKVNFKIRH
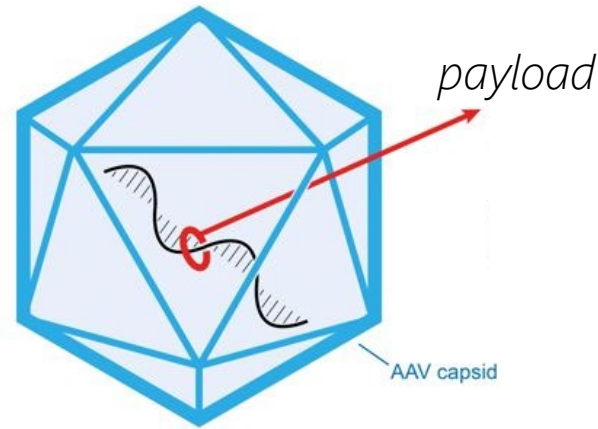
Green fluorescent protein
(GFP) folding itself

[2008 Nobel in chemistry for discovery and development of GFP,
Osamu Shimomura, Martin Chalfie and Roger Y. Tsien]

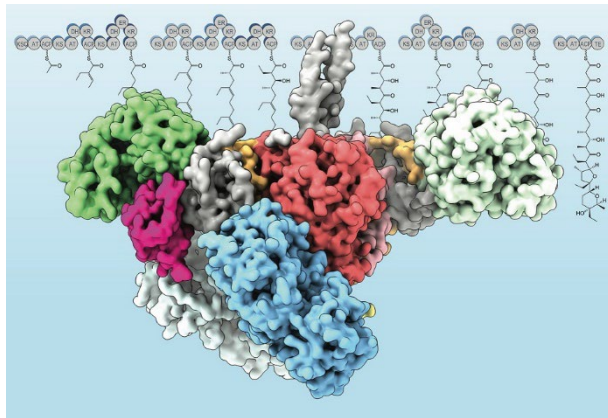# Protein engineering: therapeutics, environment, *etc.*
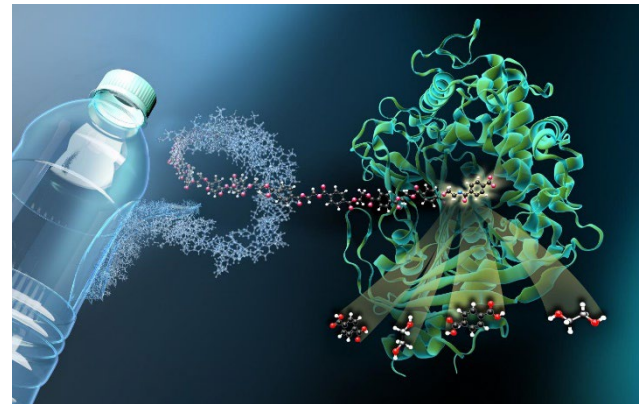

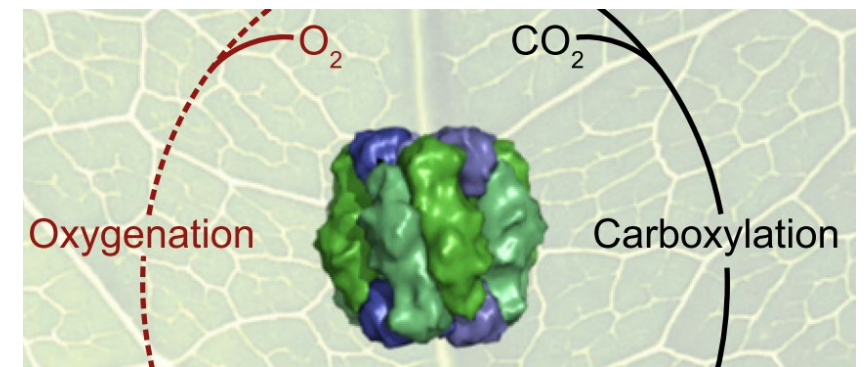antibody therapeutics


gene therapy virus delivery (AAV)


gene editing (CRISPR/Cas9)


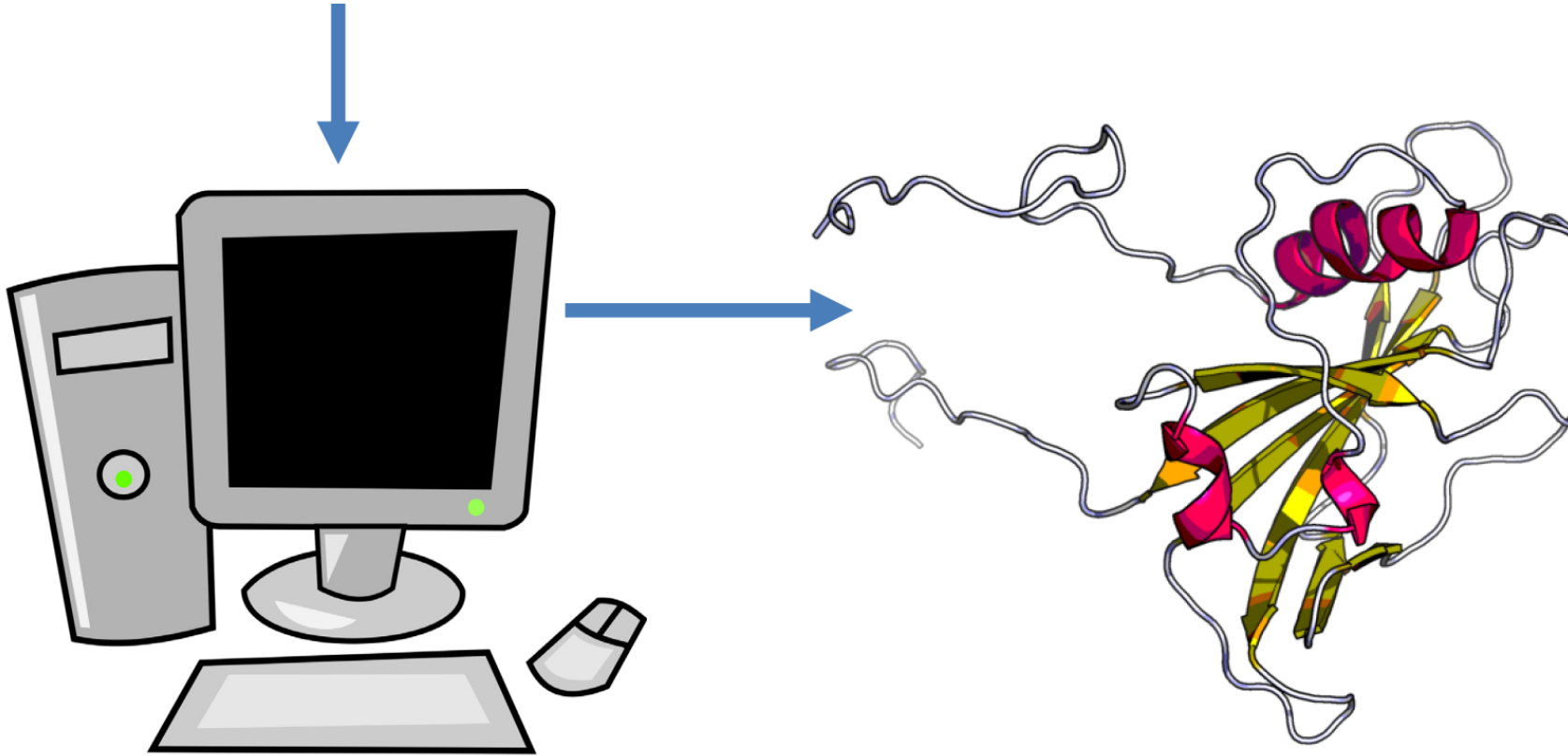antibiotics & biofuel production (PKS)


plastic recycling (PETase)
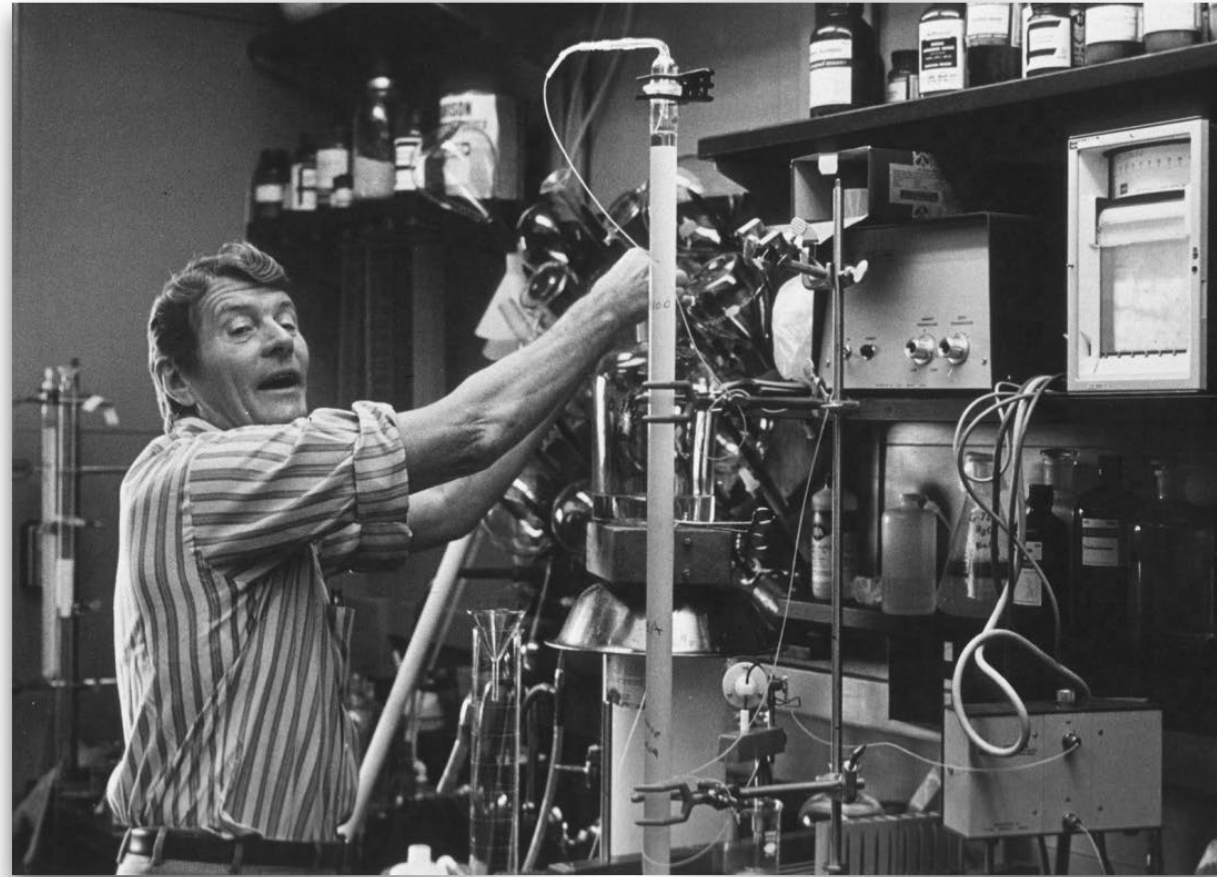

$CO_2$ biosequestration (RuBisCO)

# Protein Structure Prediction

MEKVNFLKNGVLRLPPGFRFRPTDEELVVQYLKRKVFSFPLPASIIPEVEVYKSDPWDLPGDMEQEKYFFSTK
EVKYPNGNRSNRATNSGYWKATGIDKQIILRGRQQQQQLIGLKKTLVFYRGKSPHGCRTNWIMHEYRLAN
LESNYHPIQGNWVICRIFLKKRGNTKNKEENMTTHDEVRNREIDKNSPVVSVKMSSRDSEALASANSELKK



## Has been studied several decades

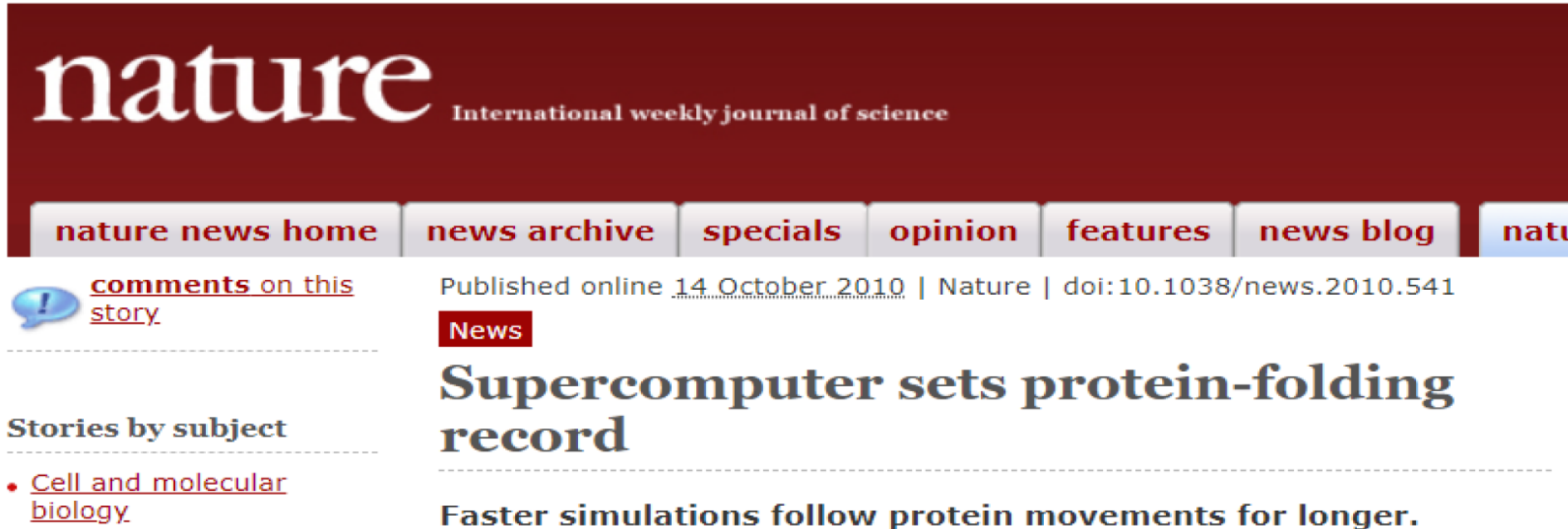# Amino acid sequence determines protein 3D structure



**Christian Anfinsen**
**Nobel Prize in Chemistry 1972**

# Protein Structure Prediction

## State of the Art Until 2015

- A lot of computing power needed

- Success rate is low even for small proteins



nature
International weekly journal of science

| nature news home | news archive | specials | opinion | features | news blog | natu |

comments on this story

Published online 14 October 2010 | Nature | doi:10.1038/news.2010.541

**News**

Stories by subject

- Cell and molecular biology

### Supercomputer sets protein-folding record

Faster simulations follow protein movements for longer.

[slide from Jinbo Xu, TTI]

2020

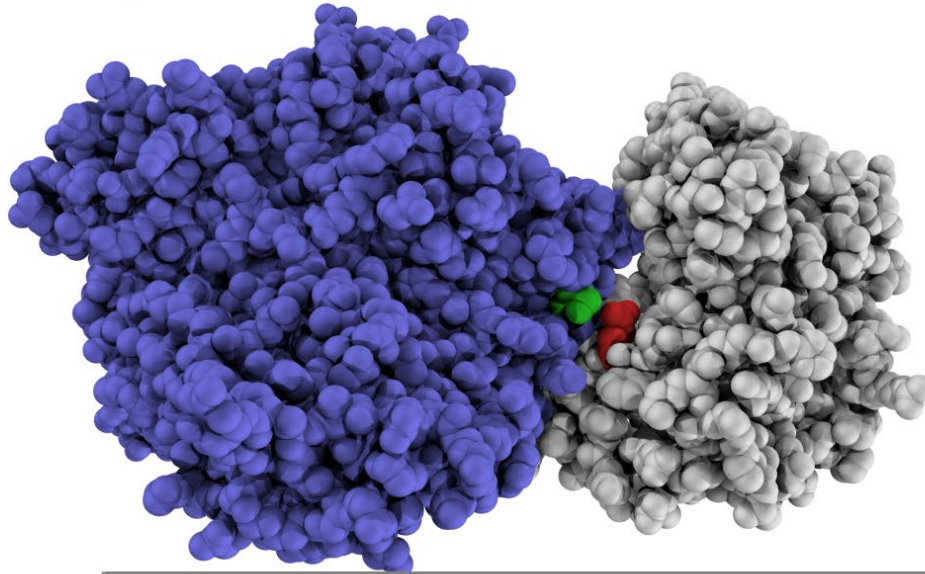State-of-the-art is deep learning based:

# AlphaFold2 relies on previous key insights



Amino acids in direct physical contact tend to covary or "coevolve" across related proteins

For example, a mutation that causes one amino acid to get bigger is more likely to preserve protein structure and function (and thus survive) if another amino acid gets smaller to make space

...GANPMHGRDQ**S**GAVASLTSVA...
...GANPMHGRDQ**E**GAVASLTSVA...
...GANPMHGRDE**K**GAVASLTSVG...
...GANPMHGRDS**H**GWLASCLSVA...
...GANPMNGRDV**K**GFVAAGASVA...
...GANPMHGRDR**D**GAVASLTSVA...
...GANPMHGRDQ**V**GAVASLTSVA...
...GANPMHGRDQ**E**GAVASLTSVA...

...VEDLMK**E**VVTYRHFMNASGG...
...VEALMA**R**VLSYRHFMNASGG...
...VATVMK**Q**VMTYRHYLRATGG...
...VARAMR**E**IGKYAQVLKISRG...
...VPELMQ**D**LTSYRHFMNASGG...
...ADHVLR**R**LSDFVPALLPLGG...
...FERART**A**LEAYAAPLRAMGG...
...VPEVMK**K**VMSYRHYLKATGG...

# AlphaFold2 "almost end-to-end" neural network

# AlphaFold2 "almost end-to-end" neural network

# AlphaFold2 "almost end-to-end" neural network



(uses an equivariant attention architecture)

# AlphaFold2 "almost end-to-end" neural network
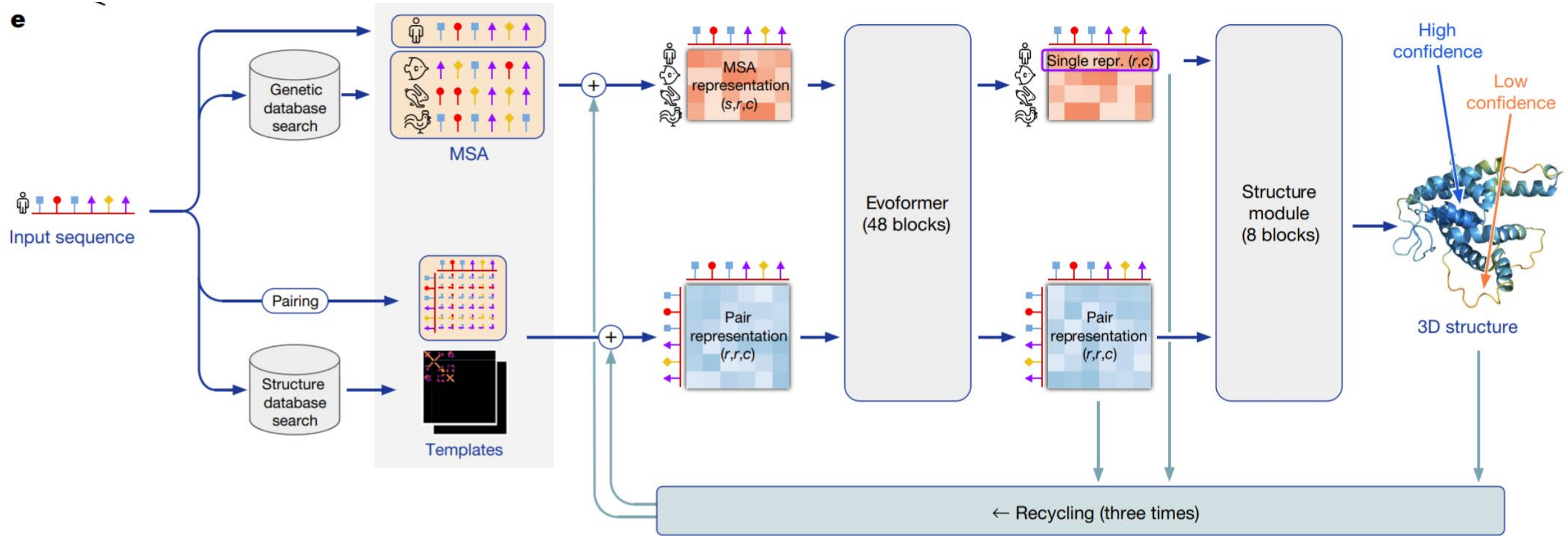
- Can end up with atom positions in violation of physics.
- Thus relies on old style energy-based approaches to refine the predicted 3D coordinates.



**Relaxation**

© 2020 DeepMind Technologies

→ The end result of iterative refinement is not guaranteed to obey all stereochemical constraints

→ Violations of these constraints are resolved with coordinate-restrained gradient descent

→ We use the Amber ff99SB force field[1] with OpenMM[2]

Steric violation

Orange: pre-relax
Blue: post-relax

# AlphaFold2 "almost end-to-end" neural network



From great blog by Mohamed Alquraishi:

https://moalquraishi.wordpress.com/2020/12/08/alphafold2-casp14-it-feels-like-ones-child-has-left-home/

# Some thoughts on AlphaFold2

- DeepMind took on a long-tackled, well-defined problem, with clear data, clear benchmarks, and a clear way to demonstrate improvement.

- Expense of protein structure data used for AlphaFold2, conservatively estimated at ~US$20 billion (Burley et al., 2023).

- They relied **heavily** on years of prior work in protein folding research: "template-based modelling", "evolutionary co-evolution modelling", "contact prediction", energy-functions.
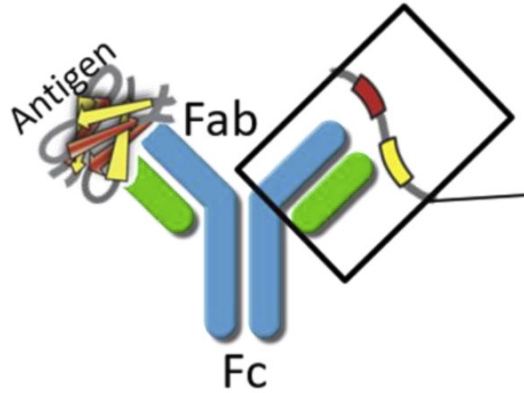
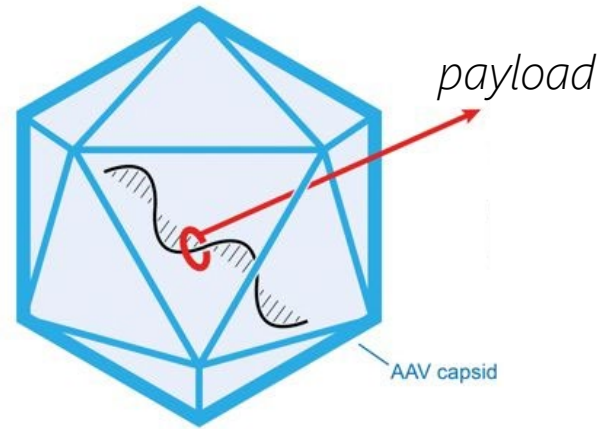# CS 189/289

Some applications of AI in biology:
1. protein structure prediction
2. protein design

# Protein engineering: therapeutics, environment, *etc.*
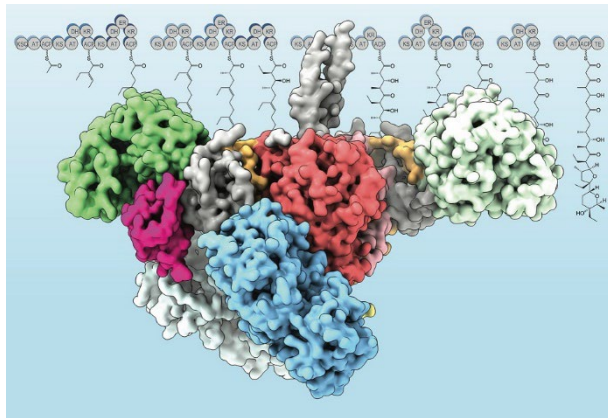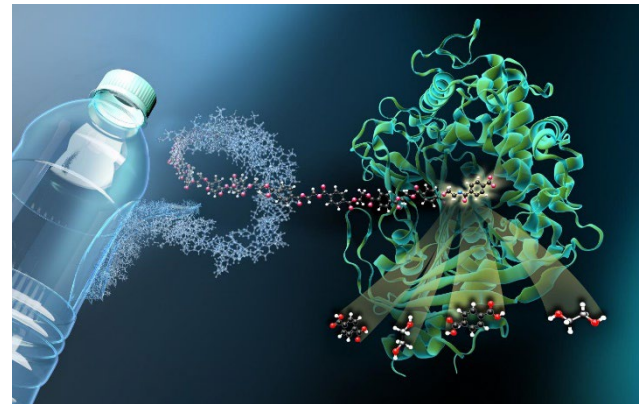

antibody therapeutics


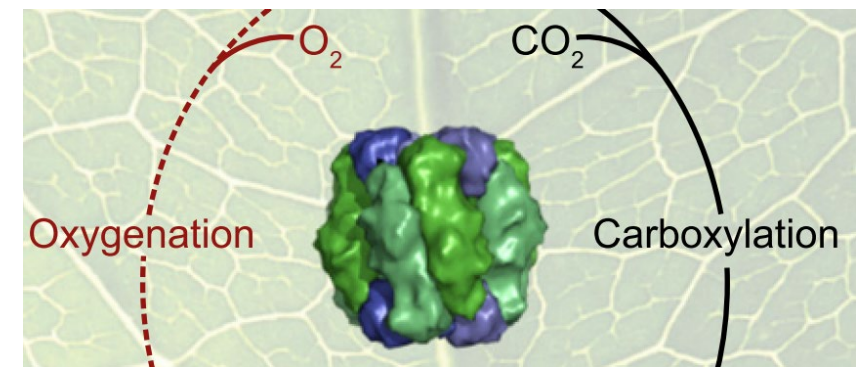gene therapy virus delivery (AAV)


gene editing (CRISPR/Cas9)


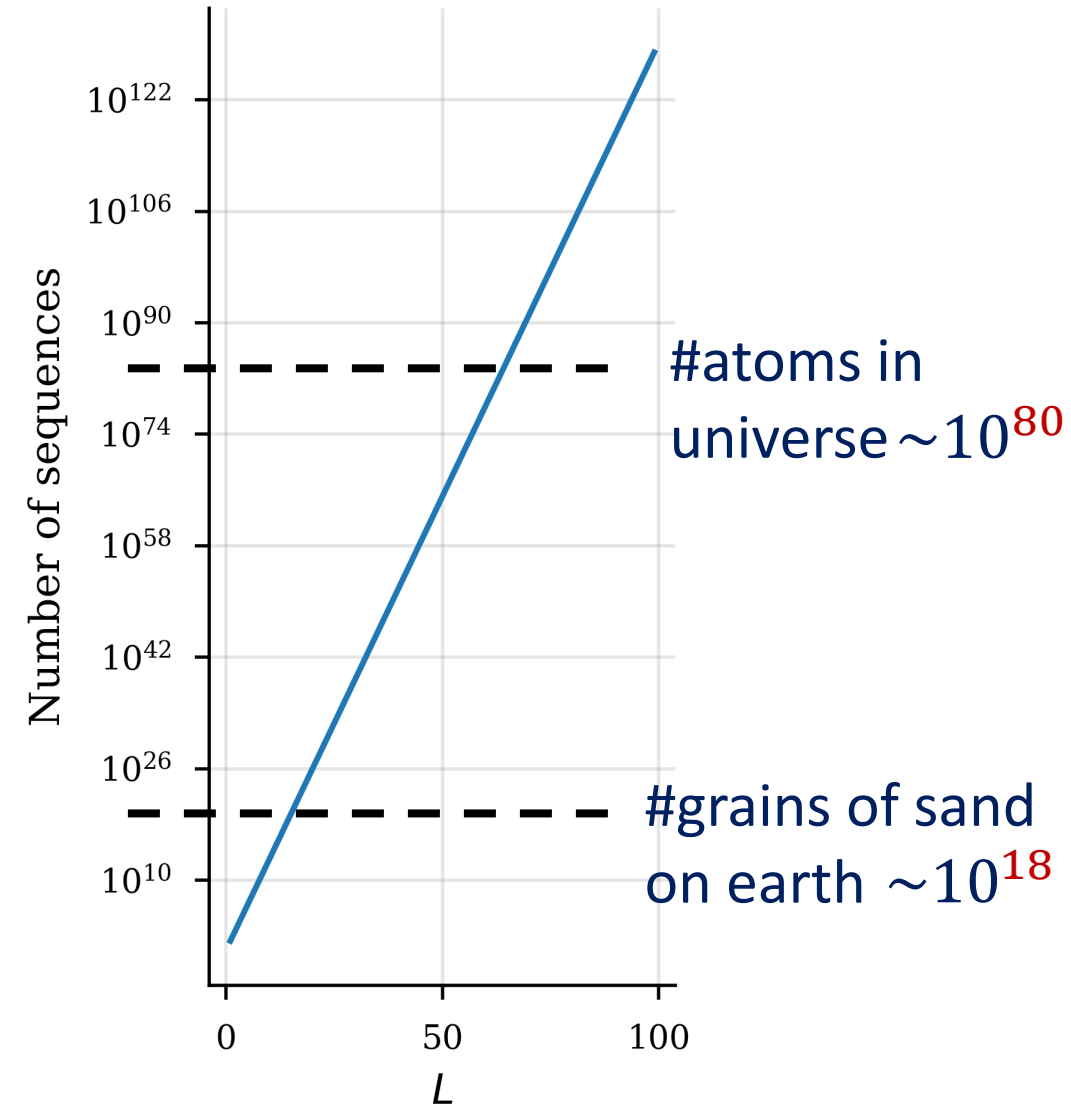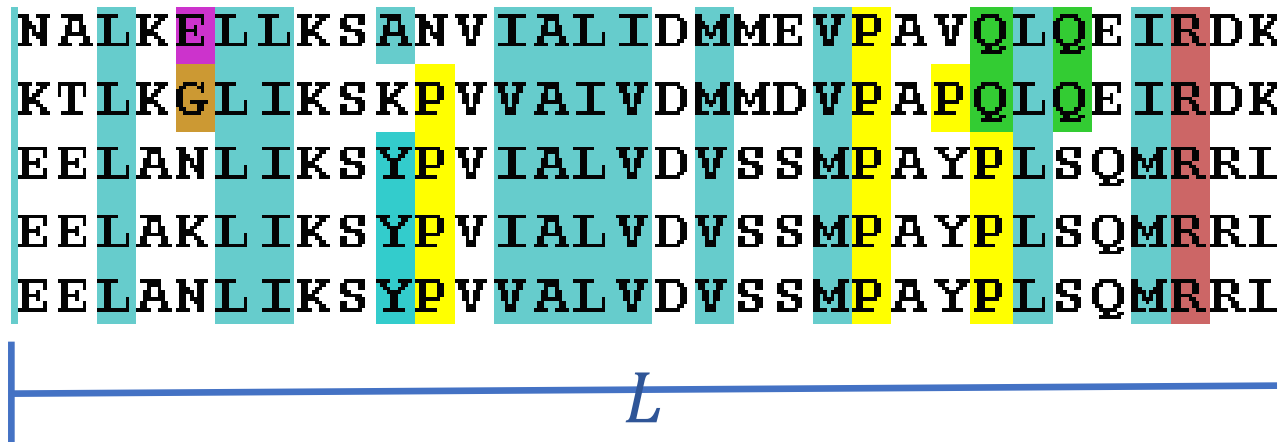antibiotics & biofuel production (PKS)


plastic recycling (PETase)


$CO_2$ biosequestration (RuBisCO)

# Fundamental difficulty: design space is nearly infinite

- Also highly rugged design space

$\Longrightarrow$ size scales as $\sim 20^L$

- Discrete search space (no gradients)

# Successes in navigating this complex space

## 1. Nature: via evolution *over millions of years.*



MSKGEELFTGVVPILV
ELDGDVNGHKFSVSG
EGEGDATYGKLTLKFIC
TTGKLPVPWPTLVTTF
SYGVQCFSRYPDHMK
QHDFFKSAMPEGYVQ
ERTIFFKDDGNYKTRA
EVKFEGDTLVRIELKGI
DFKEDGNILGHKLEYN
YNSHNVYIMADKQKN
GIKVNFKIRHNIEDGSV
QLADYQQNTPIGDGPV
LLPDNHYLSTQSALSK
DPNEKRDHMVLLEFVT
AAGITHGMDELYK

green fluorescent
protein folding itself

# Successes in navigating this complex space

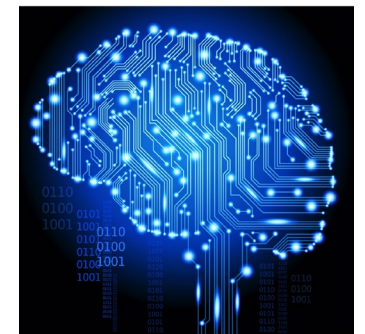1. Nature: via evolution *over millions of years.*

2. Various protein engineering strategies.

# Protein engineering strategies emerging

i. <u>Computation ("data free")</u>: **physics-based energy functions** (*e.g.*, Rosetta) to model **protein structure**, and protein binding. *~1997-2023'ish* (almost R.I.P.)

ii. <u>Wetlab</u>: **directed evolution** to iteratively directly design property of interest. *~1993-present* [2018 Nobel Prize]

iii. <u>Machine learning (augmented)</u>: generative models; function prediction; structure prediction, etc. *~2018(?)-present*

# One strategy: ML-based Directed Evolution



2018 Nobel Prize
in Chemistry

*Goal:* get same results with fewer measurements, and/or, get better result than pure DE.

1. Replace assay with predictive model.
2. Replace search with intelligent search.

# Did AlphaFold2 "solve" protein engineering?

## DeepMind's AI predicts structures for a vast trove of proteins

AlphaFold neural network produced a 'totally transformative' database of more than 350,000 structures from *Homo sapiens* and 20 model organisms.

Ewen Callaway

*sequence→ structure*



©nature

Global distance test (GDT_TS; average) vs Contest year. A score above 90 is considered roughly equivalent to the experimentally determined structure. AlphaFold (2018), AlphaFold 2 (2020).
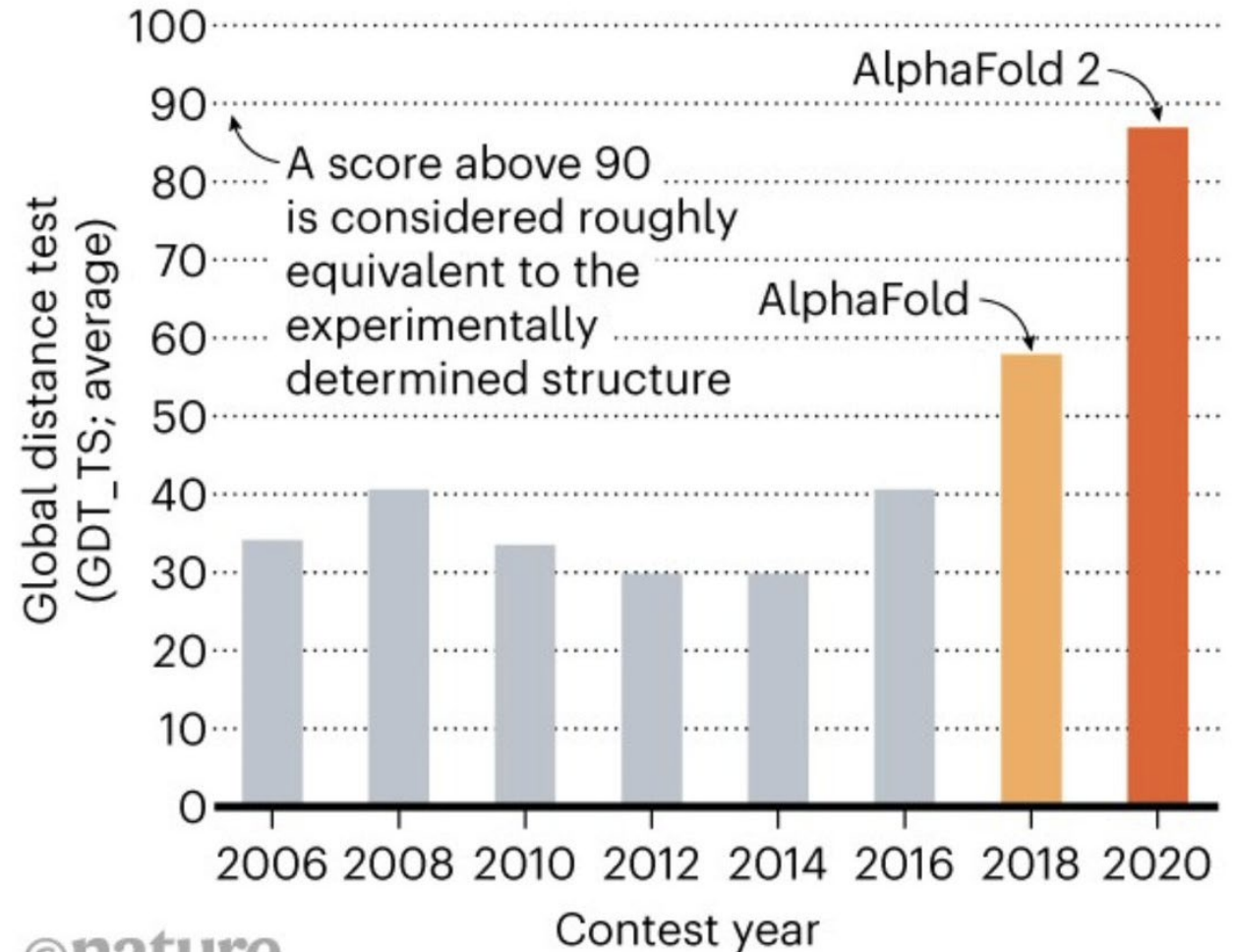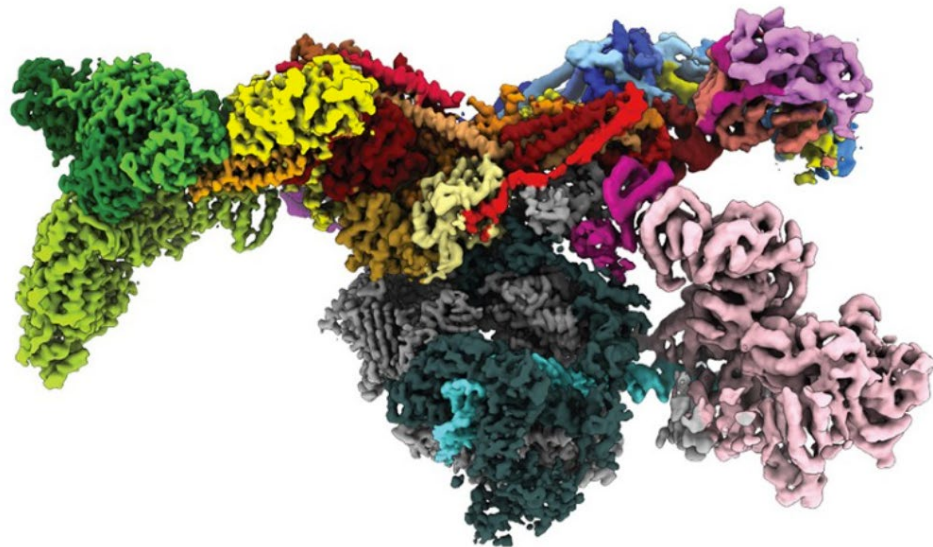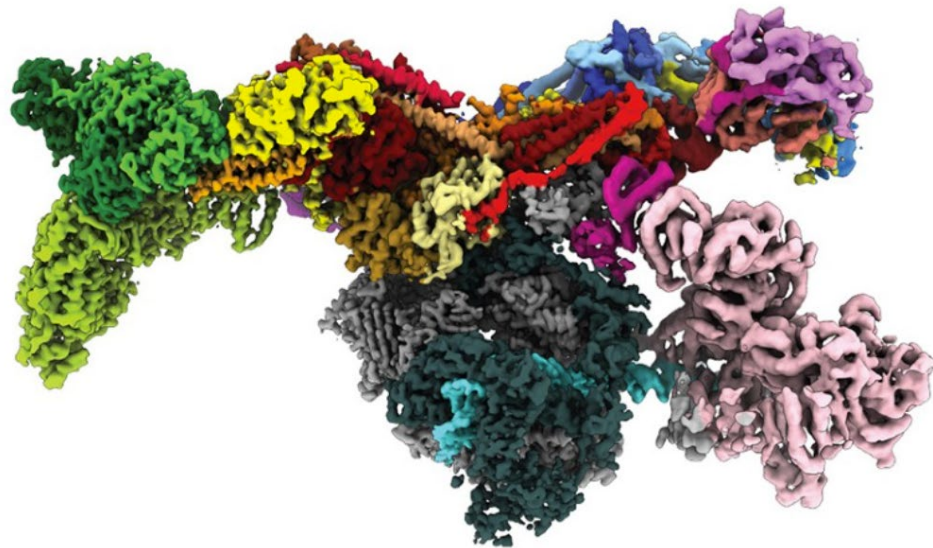
# Did AlphaFold2 "solve" protein engineering?

**DeepMind's AI predicts structures for a vast trove of proteins**

AlphaFold neural network produced a 'totally transformative' database of more than 350,000 structures from *Homo sapiens* and 20 model organisms.
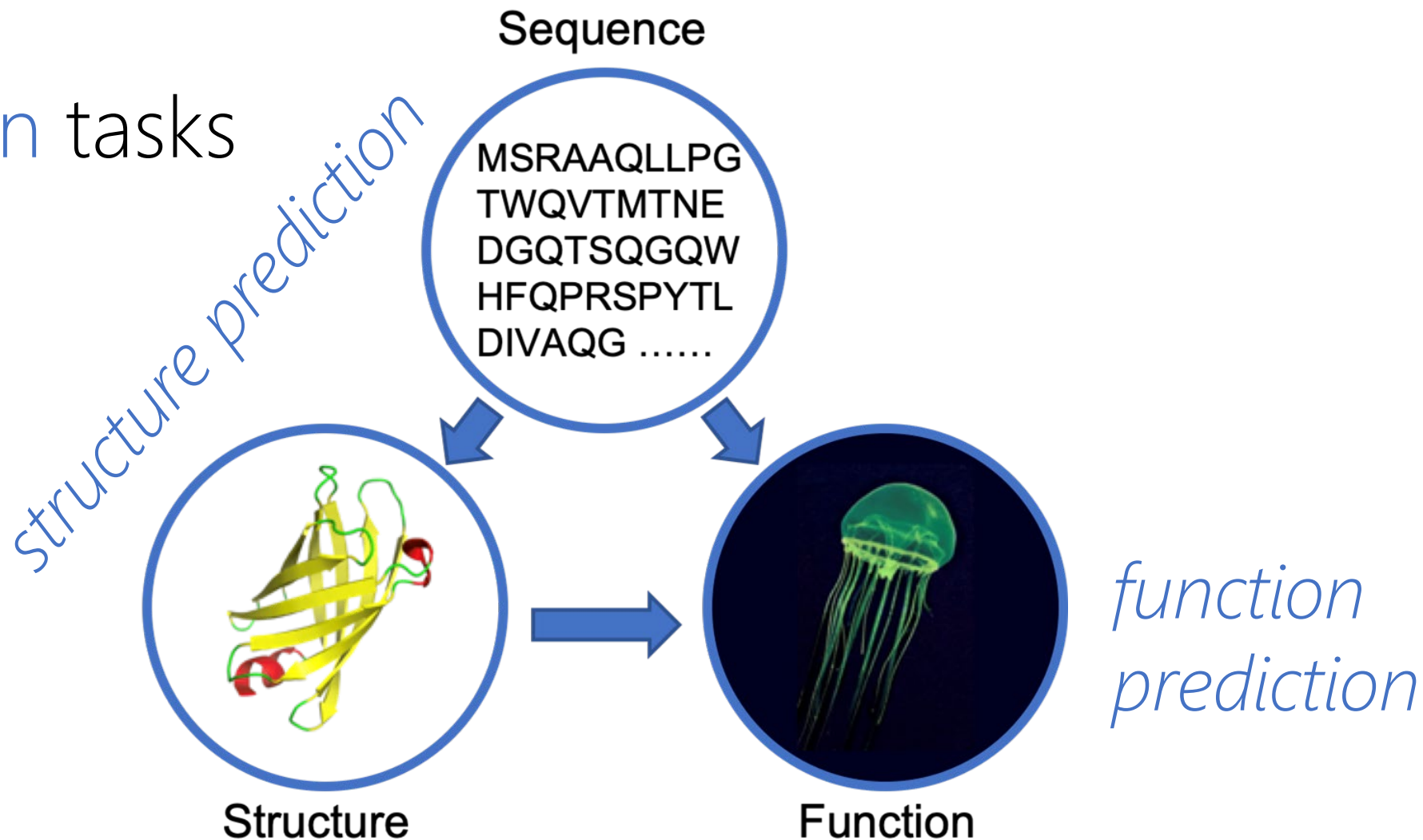
Ewen Callaway

*sequence→ structure*

- No: don't typically know which protein structures we need.

- If did, would need: *structure→sequence.* (decent ML solutions exist).

- <u>Bottleneck challenge</u>: predict which protein have the function we desire.

- AlphaFold2 *was* a breakthrough, and will surely be useful.
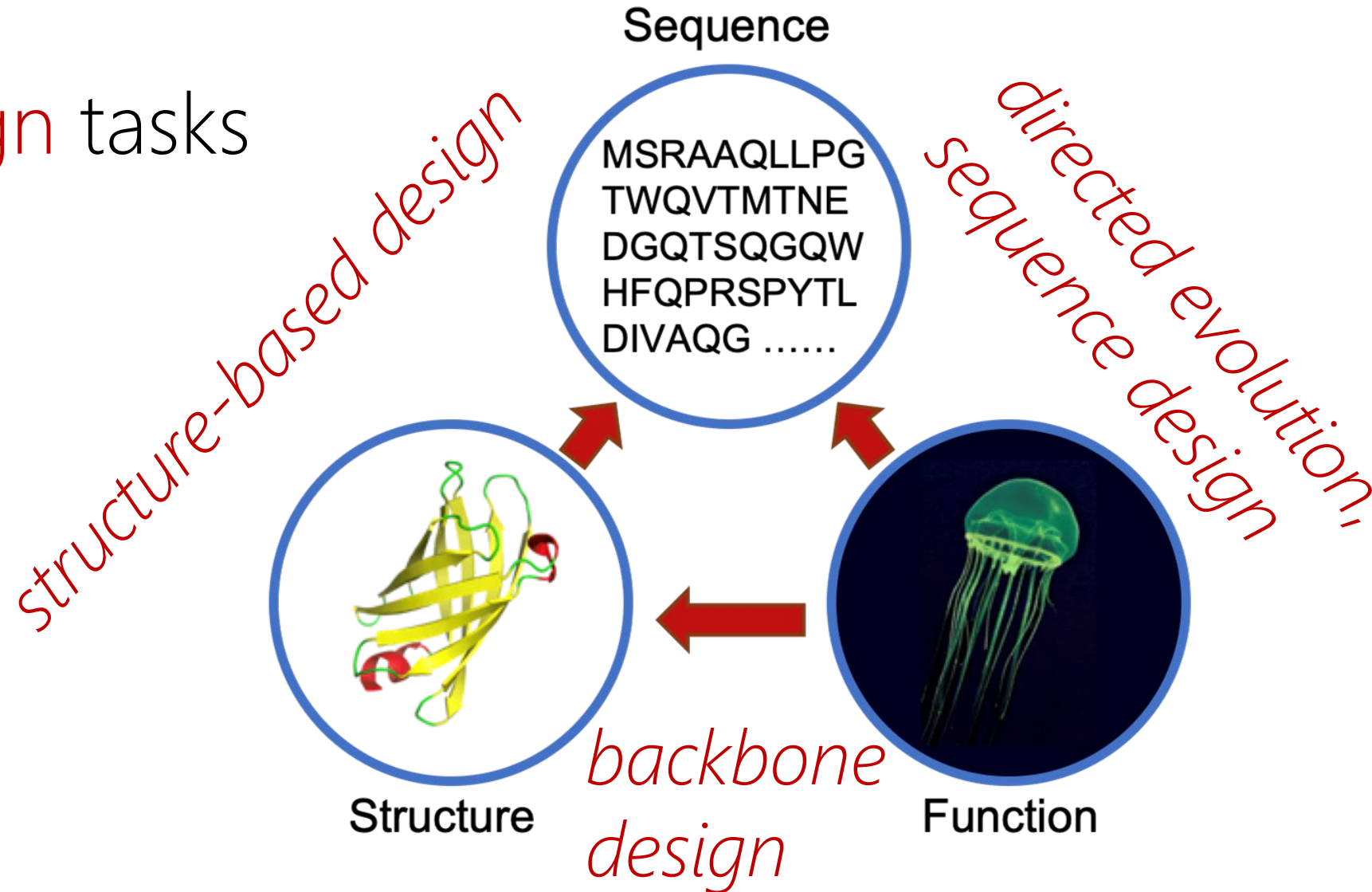
# A suite of ML protein engineering problems

Prediction tasks

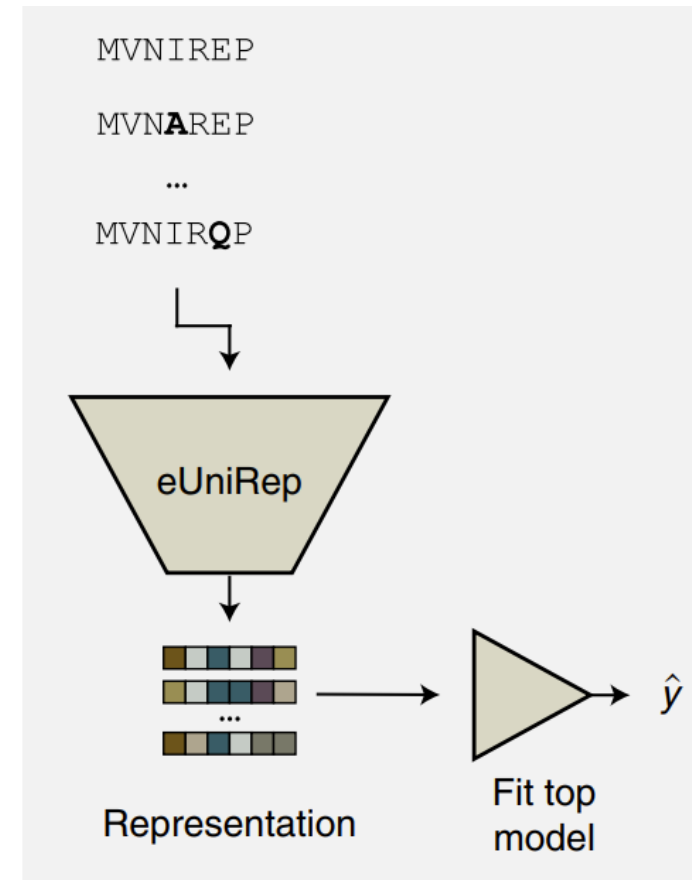# A suite of ML protein engineering problems

Design tasks



Sequence

MSRAAQLLPG
TWQVTMTNE
DGQTSQGQW
HFQPRSPYTL
DIVAQG ......

*structure-based design*

*directed evolution, sequence design*

*backbone design*

Structure

Function

# Some trends in ML + protein engineering

1. Representation learning: *un(self)supervised learning on* large-scale databases (millions of natural proteins, with *e.g.*, Transformers), or families.

   - This is really (approx.) *density estimation,* $p_\theta(\textbf{\textit{sequence}})$ through a bottleneck.


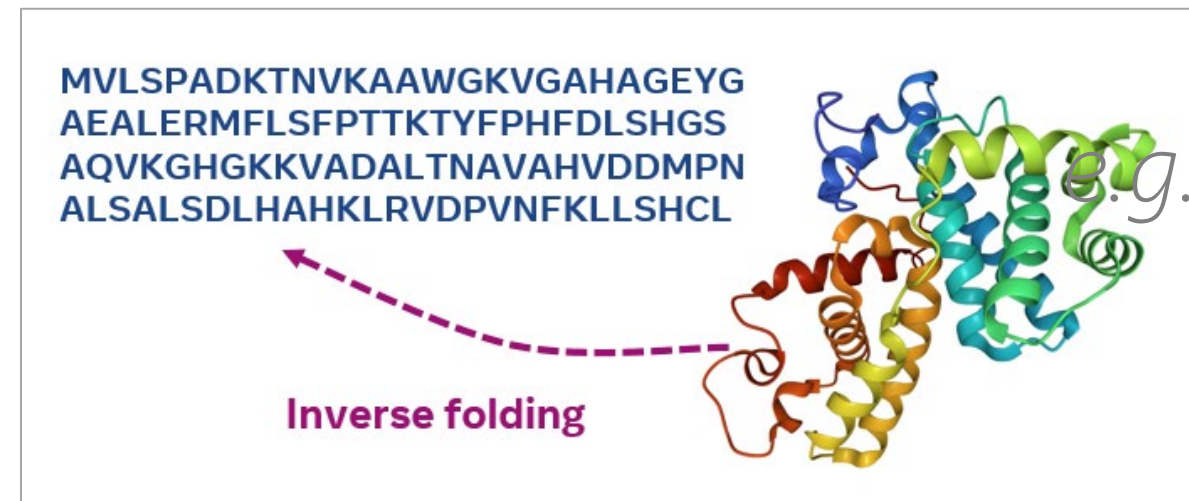
[Biswas *et al.*, *Nat. Meth.* 2021]

# Some trends in ML + protein engineering

2. (Conditional) generative models for <u>sequences</u>.

This is really (conditional) density estimation, $p_\theta(\textbf{sequence}|\textbf{C})$, (*e.g.* auto-regressive Transformer, Potts/VAE).
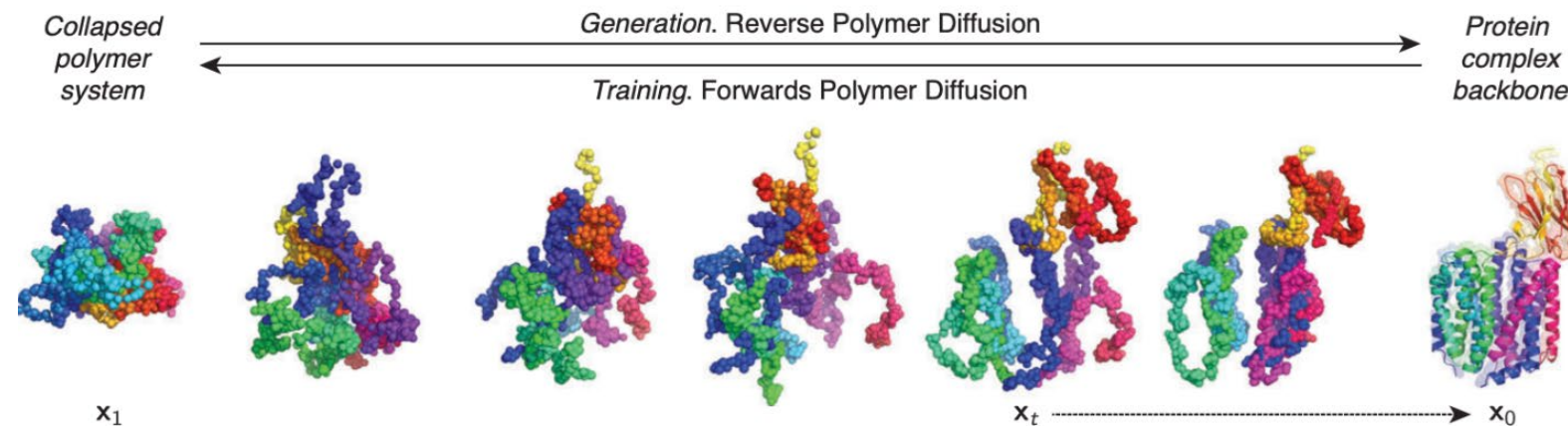
a) structure-conditioned,
   aka "inverse folding"

b) "control tag" conditioned,
   protein family



MVLSPADKTNVKAAWGKVGAHAGEYG
AEALERMFLSFPTTKTYFPHFDLSHGS
AQVKGHGKKVADALTNAVAHVDDMPN
ALSALSDLHAHKLRVDPVNFKLLSHCL

*e.g.,*

**Inverse folding**

# Some trends in ML + protein engineering

3. (Conditional) generative models for <u>structure</u>.

- This is really (conditional) density estimation, $p_\theta(\mathbf{backbone|F})$, (*e.g.* "Diffusion" models latest trend).
- Only as good as function prediction, $p(F|backbone)$.
- Paired with inverse-folding to get sequence.



Collapsed polymer system

*Generation. Reverse Polymer Diffusion*

*Training. Forwards Polymer Diffusion*

Protein complex backbone
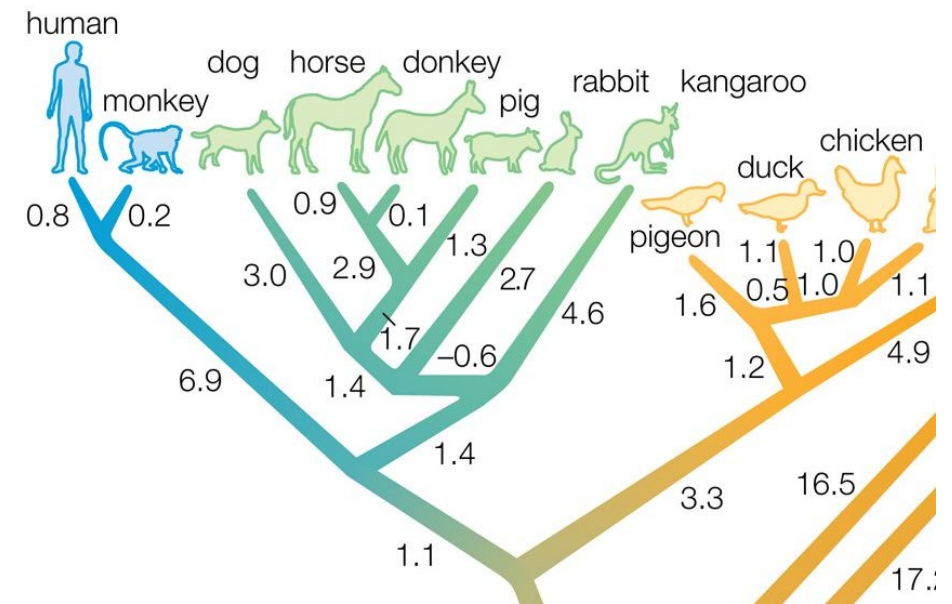
$\mathbf{x}_1$      $\mathbf{x}_t$      $\mathbf{x}_0$

[Ingraham *et al. bioRxiv* 2022]

# Some trends in ML + protein engineering

4. ML to estimate function from sequence and/or function:

- *e.g.,* $p_\theta(F|sequence)$.
- Few or no labelled data.
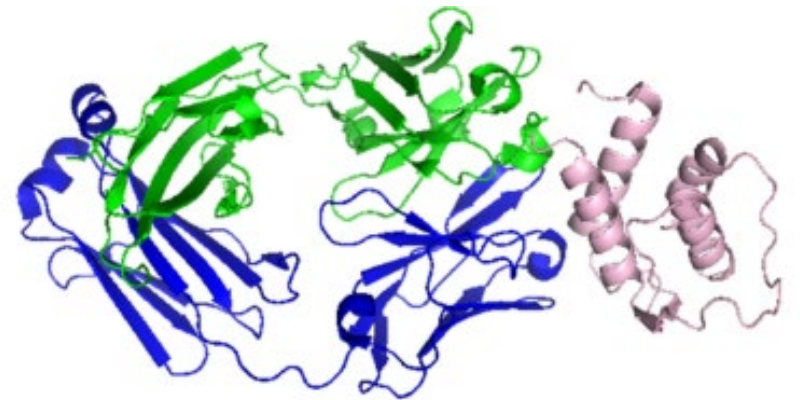- *Leverage evolutionary information\**, or large unsupervised models on pan-proteomic database.

*\*key part of AlphaFold2*

# Some trends in ML + protein engineering

5. Structure prediction: filling the gaps left by AlphaFold2

- Orphan proteins (with *no/few homologs*).
- Proteins *in bound form.*
- Protein dynamics and conformational distributions.
- Protein-protein binding.
- Protein-DNA/RNA binding

# ML focus of my group: "ML-based design":

A. Natural tension between extrapolation vs. trustworthiness. [1-4].

B. Related to caus... ...certainty (whereas we typically think ...

C. Suitable protei... ...NLP) [4-7].

D. Design of distr... ...al sequences [1,2,8,9].



1. Brookes *et al* ICLM 2019
2. Fannjiang *et al NeurIPS*
3. Fannjiang *et al PNAS* 20
4. Nisonoff *et al* arXiv 202
5. Aghazadeh *et al* Nat. Comm. 2021 (sparse
6. Brookes *et al PNAS* 2022 (funct. pred.)
7. Hsu *et al Nat. Biotech*. 2022 (function pred
8. Zhu, Brookes, *et al* ,*bioRxiv*. (opt. design)
9. Busia & Listgarten, *bioRxiv* (log enrichmen
10. Fannjiang & Listgarten, *arXiv* (overview)

# Analogy: can we trust "banana" design?
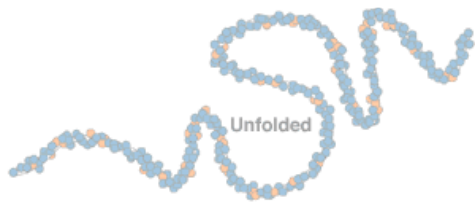
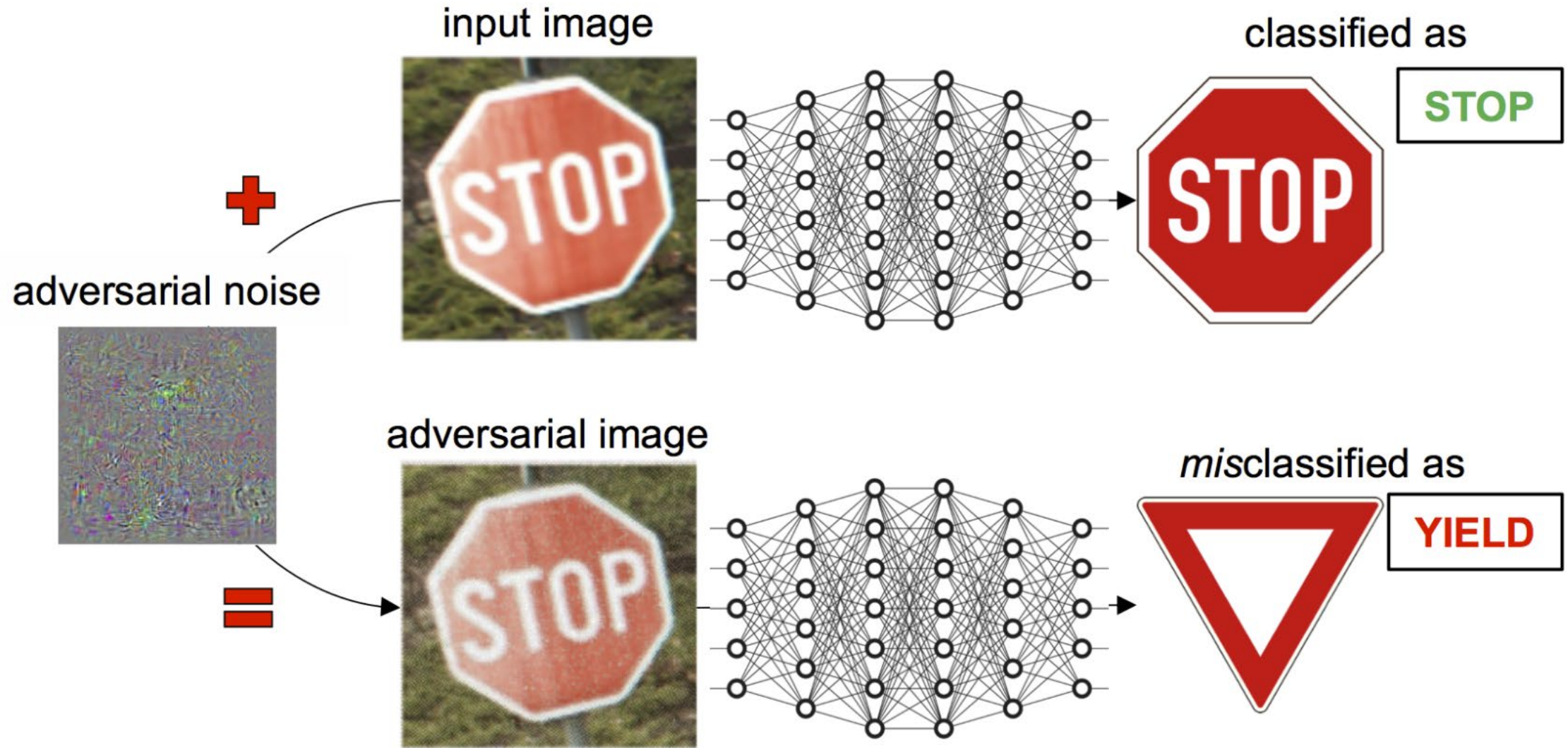# Naïve design yields abstract art.



*desired property*

*non-folding protein*

catalytic efficiency

1. Brookes *et al ICLM* 2019 (CbAS)
2. Fannjiang *et al NeurIPS* 2020 (autofocus)

# Pathologies of DNNs: in design, we're the adversary



input image

classified as

STOP

adversarial noise

+

adversarial image

=
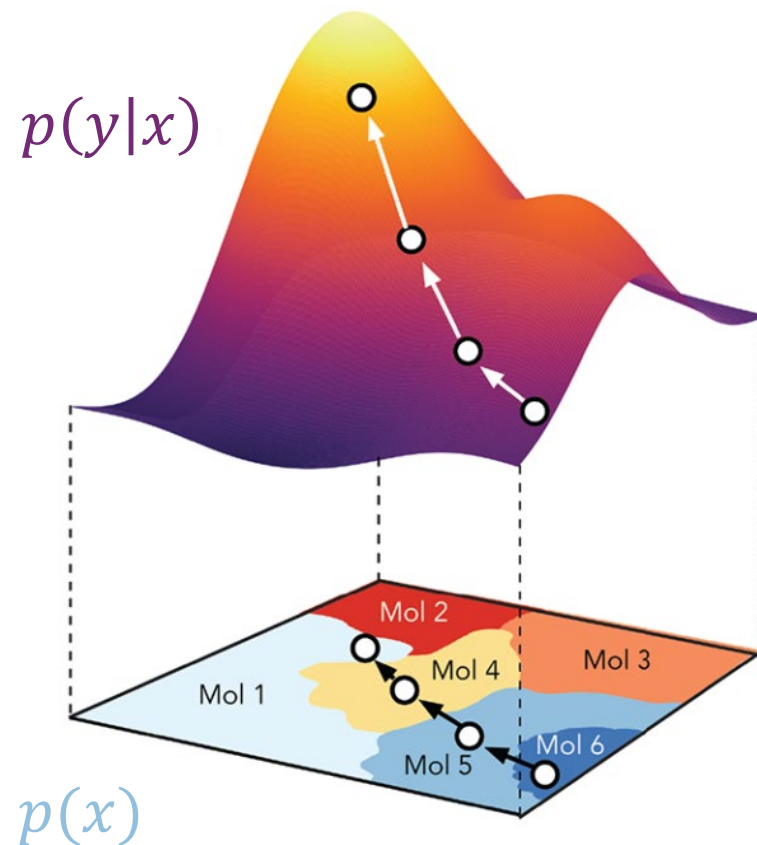
misclassified as

YIELD

# *Conditioning by Adaptive Sampling for Robust Design (CbAS)*

How to handle a pathology in design?

Leverage prior knowledge, $p(x)$, by modeling:

1. Where training data lie.

2. "Protein-likeness", e.g. stability via biophysics, or implicitly via large pan-proteome unsupervised models.
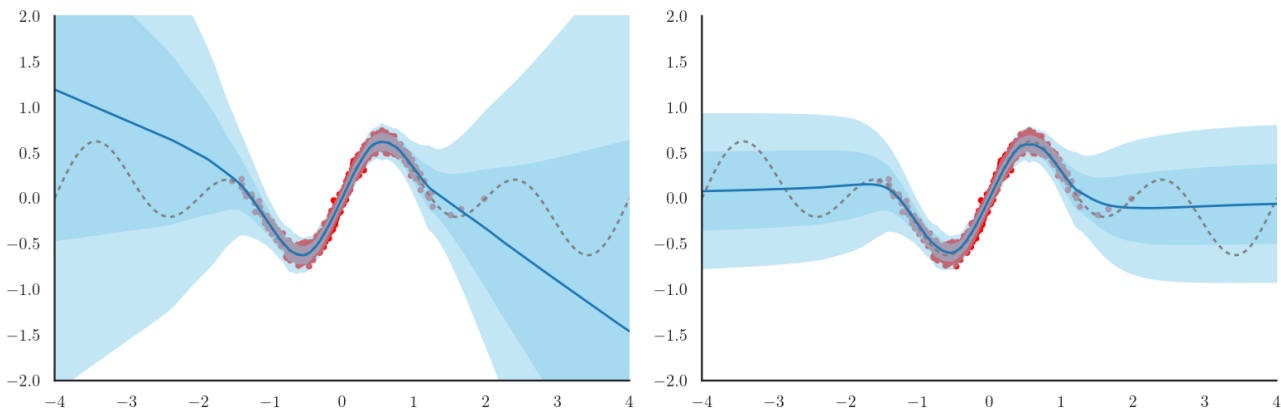


$p(y|x)$

$p(x)$

*[Gomez-Bombarelli, ACS Cent. Sci. 2018.]*

Brookes, Park & Listgarten *ICML* 2019

David Brookes

# Augmenting Neural Networks with Priors on Functional Values

Coherent blending of <u>function value prior information</u>, such as biophysical models, to Bayesian Neural Networks (BNN).
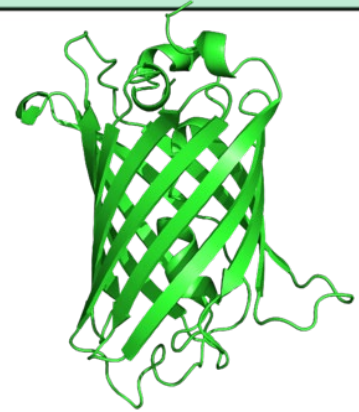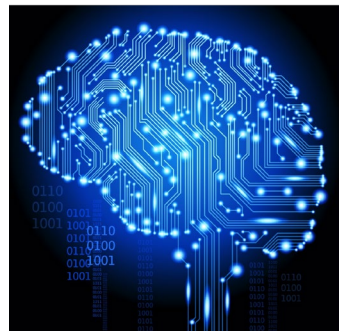
Easy to implement, zero added cost.



regular BNN

function-value augmented BNN

| METHOD | LOG-LIKELIHOOD |
|---|---|
| NN | $-8.33 \pm 0.66$ |
| BNN | $-5.73 \pm 0.18$ |
| STACKING: BNN+NON-FUNCTIONAL PRIOR | $-8.63 \pm 0.33$ |
| STACKING: BNN+STABILITY PRIOR | $-8.61 \pm 0.34$ |
| *fv-BNN* (NON-FUNCTIONAL PRIOR) | $-1.82 \pm 0.00$ |
| *fv-BNN* (STABILITY PRIOR) | $-1.53 \pm 0.00$ |

<u>Nisonoff</u>, Wang, Listgarten, *bioRxiv*

Hunter Nisonoff

# The real deal: testing+developing our ideas with wetlab collaborators

- David Schaffer (UC Berkeley; AAV for gene therapy)
- David Savage (UC Berkeley; CRISPR-Cas9 system)
- Chris Garcia (Stanford, protein-protein interactions)
- Phil Romero (U Wisconsin; enzymes for plastic degradation)
- Secure and Robust Biosystems Design Group (LL National Labs, Columbia University, University of Maryland, University of Minnesota)

# Engineering AAV for gene therapy delivery

The Adeno-associated virus (AAV) is a <u>non-pathogenic virus</u> that shows promise for <u>delivering gene therapies</u> (*e.g.* deliver blindness therapy to outer retina).

*UC Berkeley: Chem. & Bio. Engineering*



David Schaffer

Bonnie Zhu

David Brookes
(now at Dyno)

Akosua Busia
<u>Now on job market!</u>

<u>Zhu</u>, <u>Brookes</u>, <u>Busia</u>,…, Nowakowski, <u>Listgarten</u>, <u>Schaffer</u>, *bioRxiv*

# Promising AAV clinical trials

Recent clinical trial success:

Leber's congenital amaurosis (AAV)

Spinal muscular atrophy (AAV)

Hemophilia B (AAV)

Lipoprotein lipase deficiency (AAV)



Many diseases targets are still beyond the reach of current gene delivery technology

# Ongoing challenges for AAV-based therapeutics

- Inefficient delivery to target tissues/cells.
- Non-specific delivery.
- Pre-existing immunological neutralization.
- Inefficient uptake into target cells.

First AAV project goal, "library design":
- Obtain optimal starting "library" for all these engineering goals.
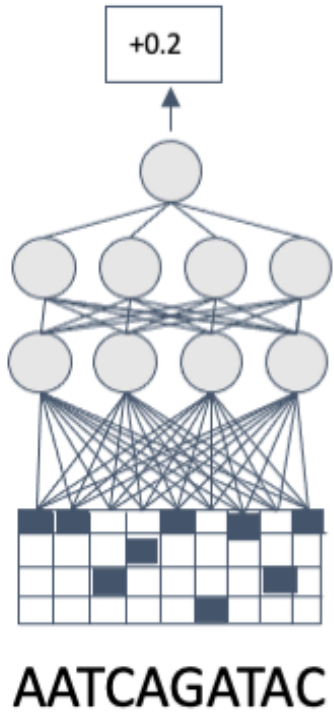- *i.e.*, fix the huge amount of library that gets wasted because doesn't "package".
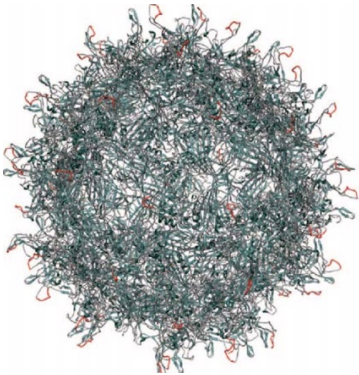
# AAV library design

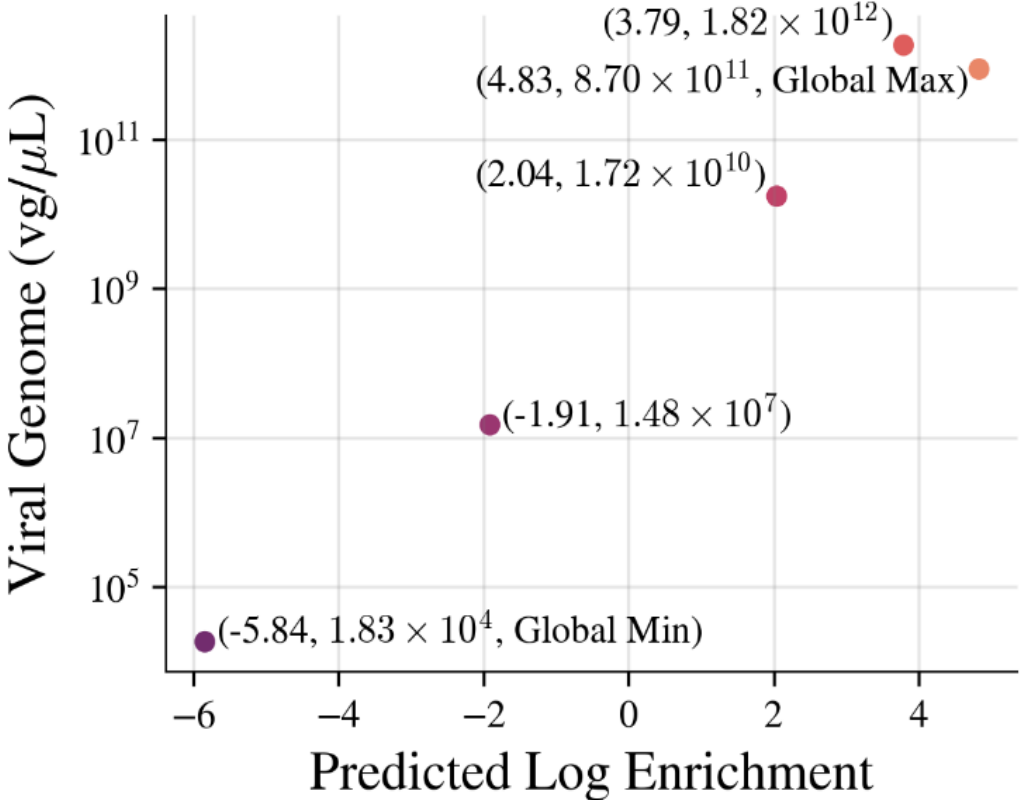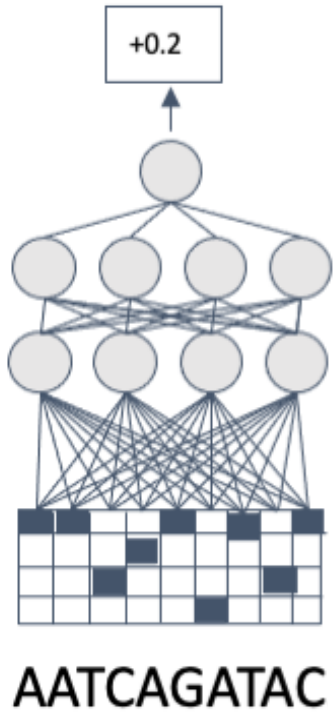1. Build predictive model and test (*sequence→packaging* fitness).

# AAV library design
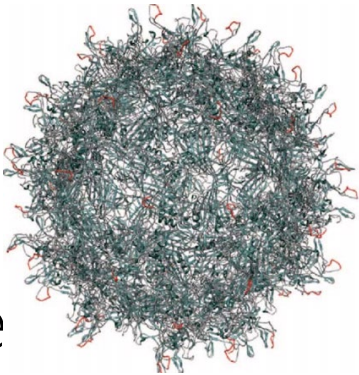
2. Wetlab validate model (measure titer directly)



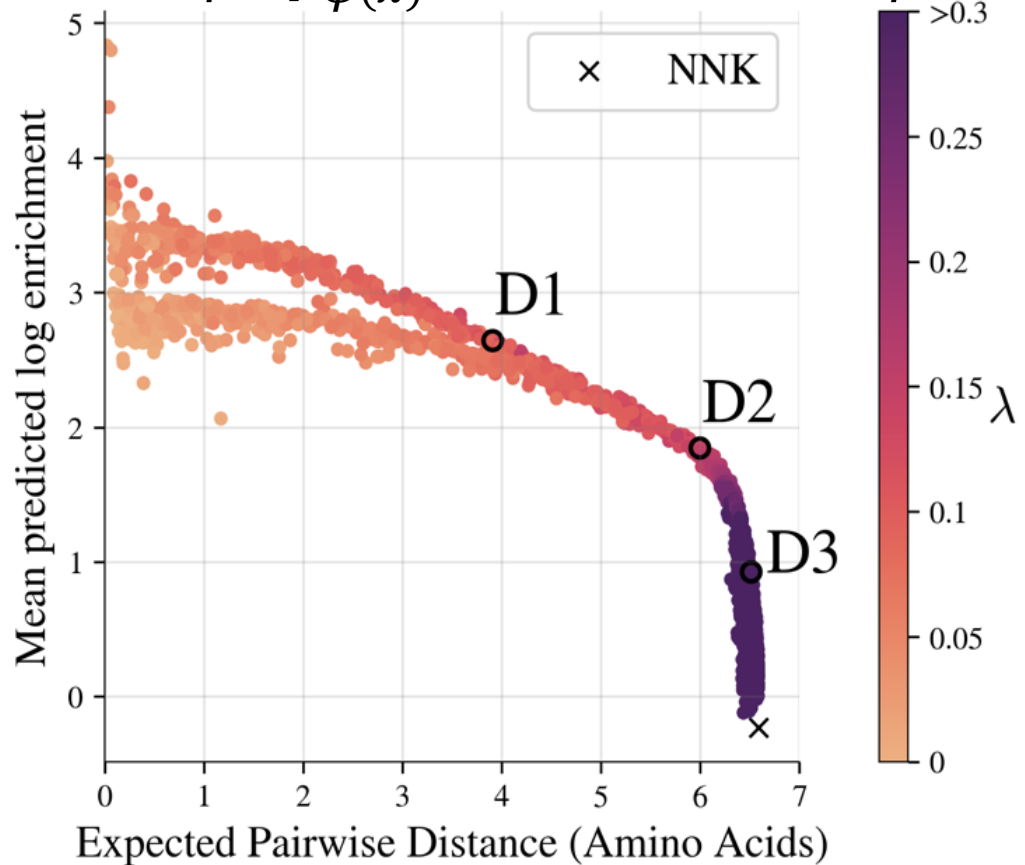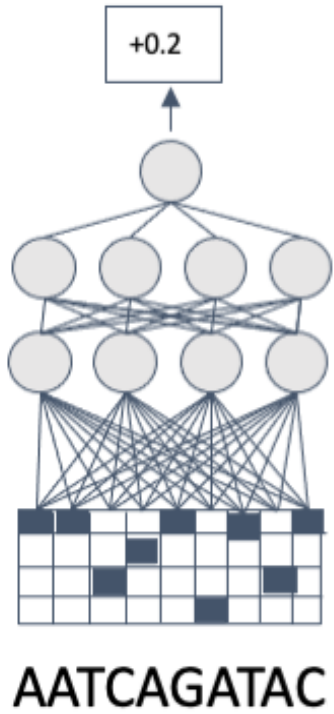| Sequences | Predicted Log Enrichment | Experimental Viral Titer (vg/$\mu$L) |
|---|---|---|
| LSSTTAA | 4.834 | $8.70 \times 10^{11}$ |
| DSRLSGT | 3.793 | $1.82 \times 10^{12}$ |
| LEPDAAL | 2.044 | $1.72 \times 10^{10}$ |
| IRWRATG | (-) 1.91 | $1.48 \times 10^{7}$ |
| RWPRRVL | (-) 5.84 | $1.83 \times 10^{4}$ |

# AAV library design
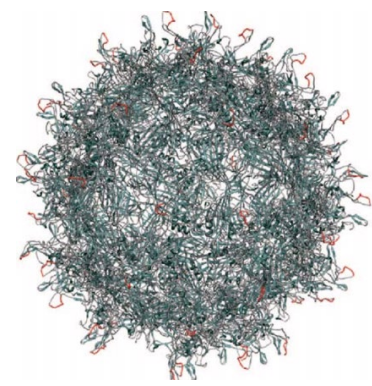
3. Invert ML predictive model to get diversity-fitness optimality curve
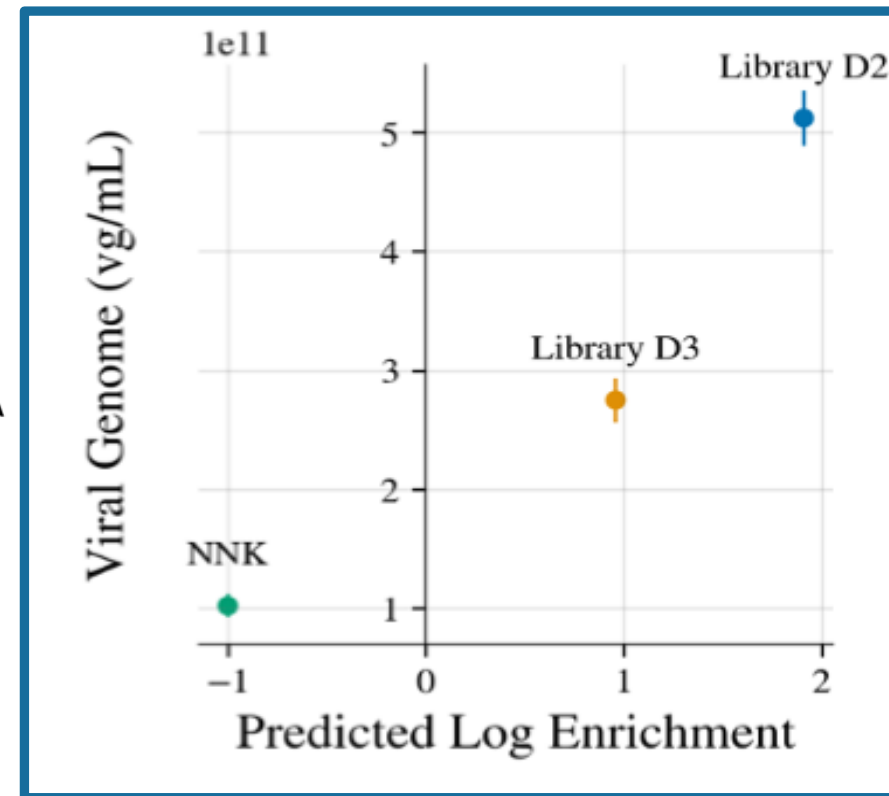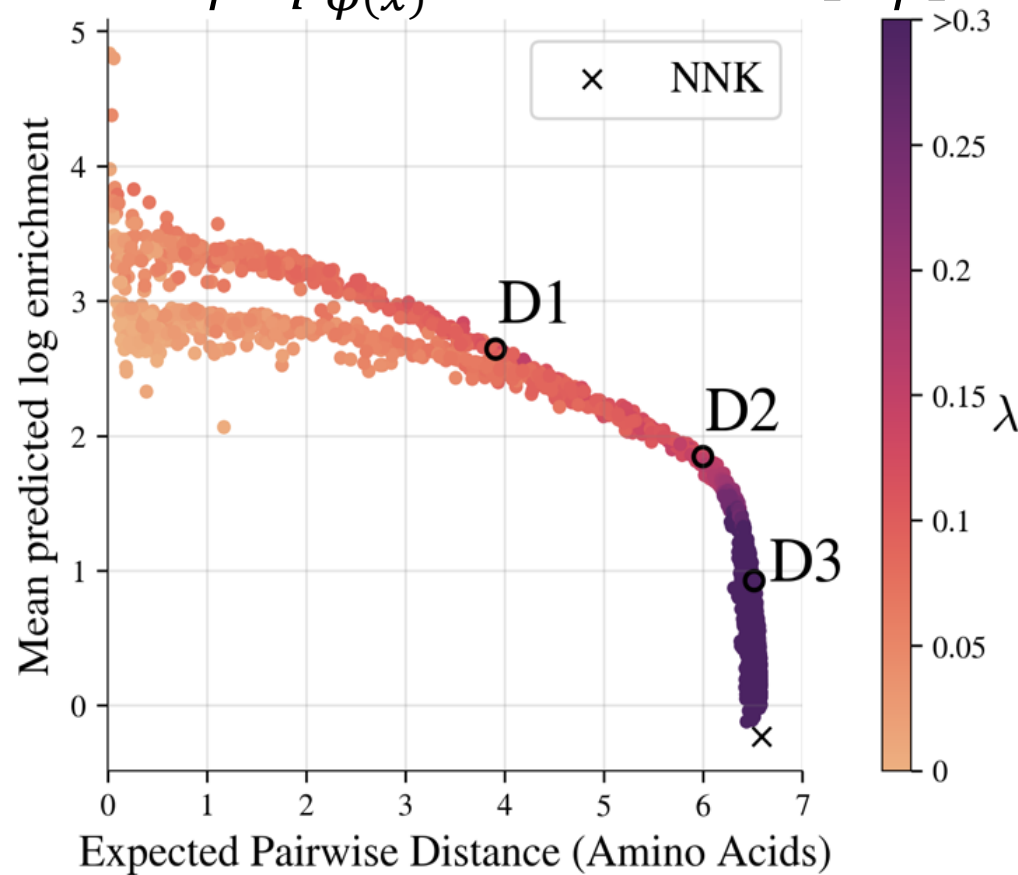
$$argmax_\phi \mathbb{E}_{p_{\phi(x)}}[f(x)] + \lambda H[p_\phi]$$

# AAV library design
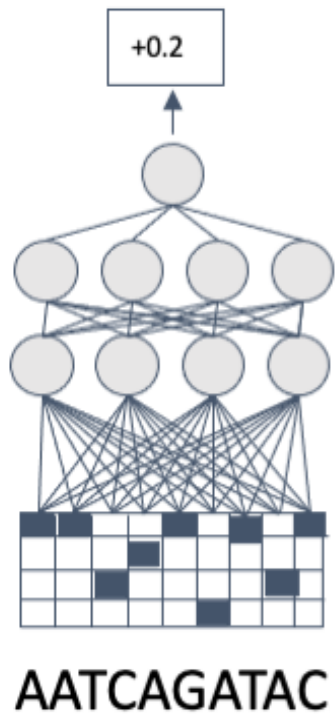
4. Validate in the lab.

$$argmax_\phi \mathbb{E}_{p_{\phi(x)}}[f(x)] + \lambda H[p_\phi]$$

# AAV library design

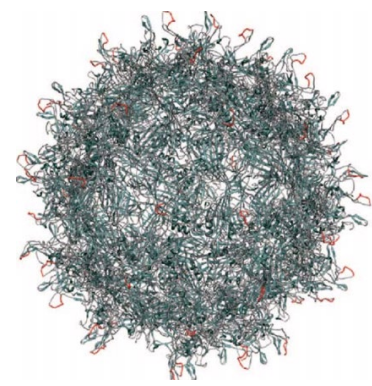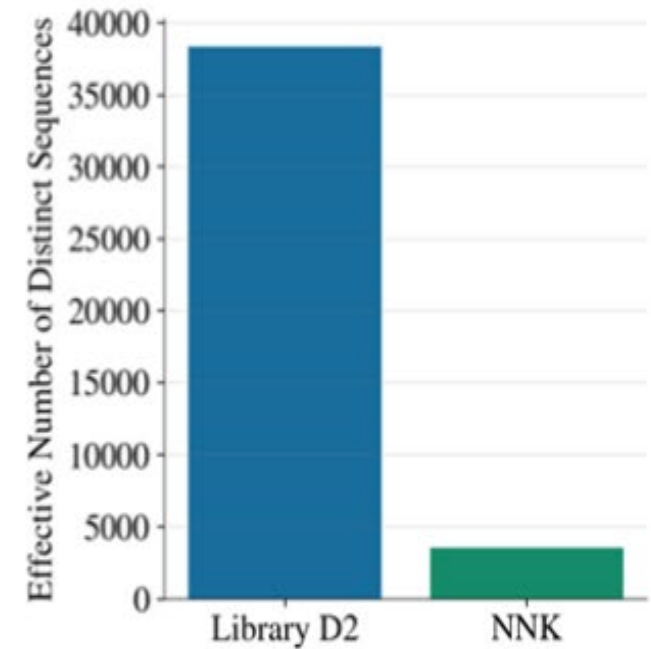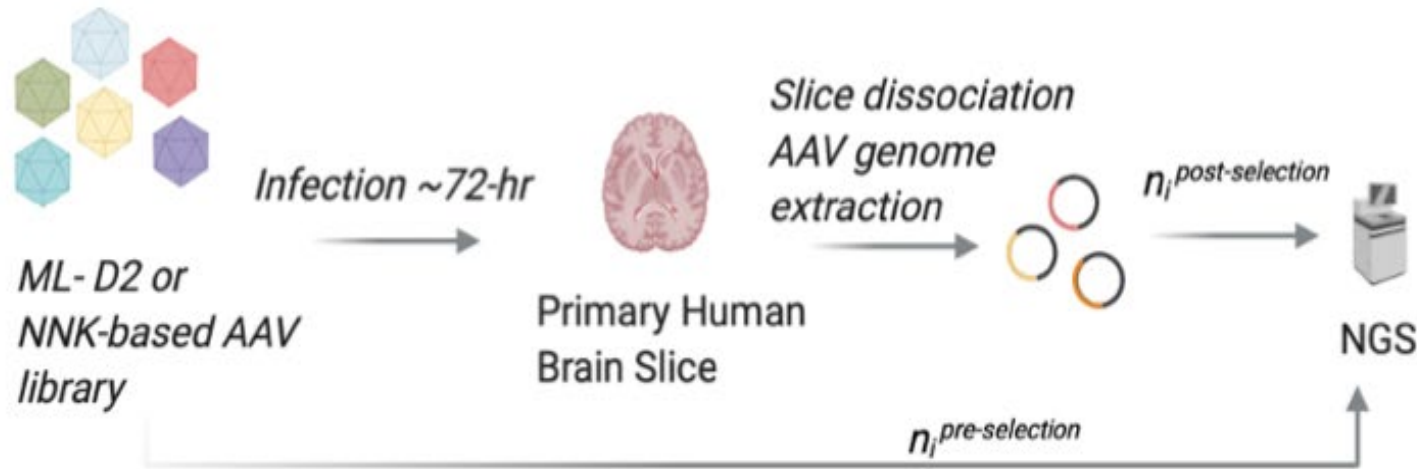5. Demonstrate better downstream selection (human brain cell infectivity), that it *was not specifically designed for.*

# Parting thoughts: ML + protein engineering

1. Exciting times!

2. Are we close to ChatGPT4 for protein engineering? No.

3. Far less data than in text, vision—will need to be much more clever for the answers to "emerge" (unless same functions).

4. AlphaFold2 and progeny will help advance protein engineering.

5. Predicting function (generally) will remain difficult problem for a long time.

6. Whiplash---this field is moving quickly, hard to tell what is real/ useful.

# The perpetual motion machine of AI-generated data and the distraction of "ChatGPT as scientist"

Jennifer Listgarten

EECS Department
University of California, Berkeley
Technical Report No. UCB/EECS-2023-239
November 30, 2023

Since ChatGPT works so well, are we on the cusp of solving science with AI? Isn't AlphaFold2 suggestive that the potential of LLMs in biology and the sciences more broadly is limitless? Can we use AI itself to bridge the lack of data in the sciences in order to then train an AI? Herein we present a discussion of these topics.