

The background is a dark, textured field filled with numerous thin, glowing blue lines that resemble molecular bonds or data streams. On the right side, there is a faint, circular, golden-brown pattern that looks like a complex molecular structure or a stylized logo.

VANTAI

ILUMINATE THE
MOLECULAR
WORLD

Neo-1 | 2025

Computational approaches for small molecule drug discovery

Introduction: computers and drugs

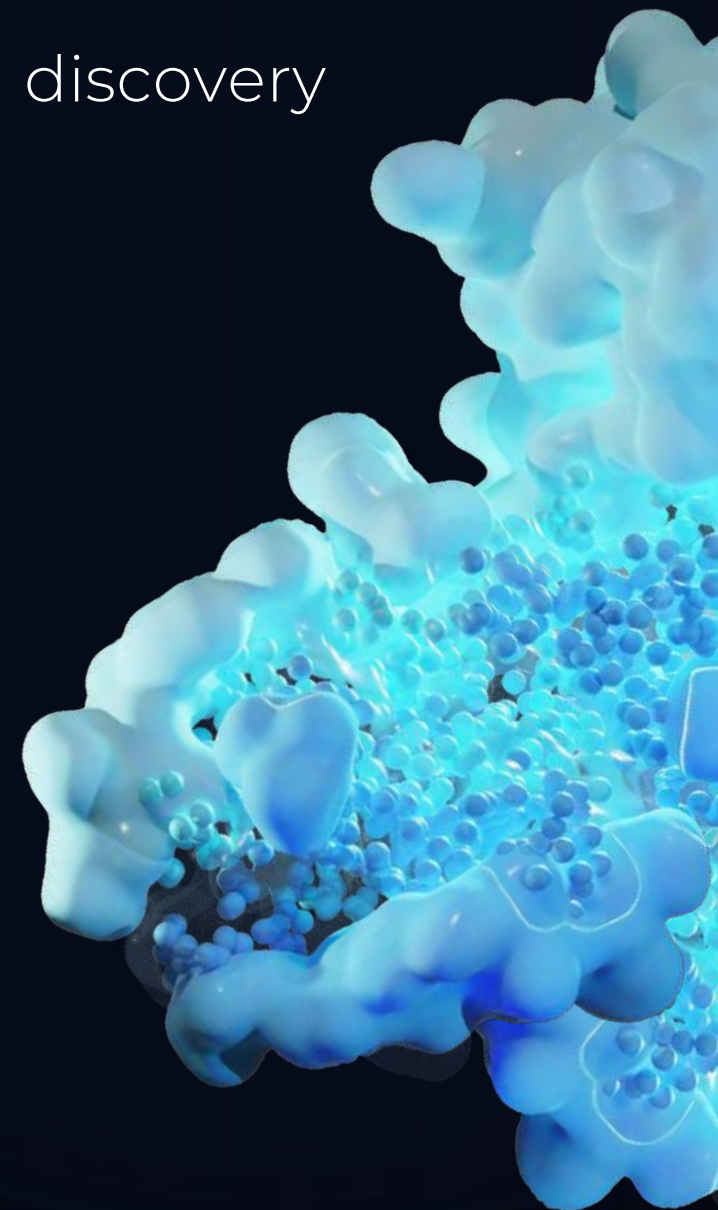
- Introduction
- The problem: the art and science of making drugs
- A quick history of computational drug discovery

The “classic” era: task-specific tools and models

- Single task focused methods
- Breakthrough in protein structure prediction: the road to AF2 and beyond
- De-novo molecular generators

The new era: foundation models trained on black box data – Neo and beyond

- Neo-1: unifying all-atom structure prediction and de-novo generation for the first time
- The promise of black-box data & NeoLink



Who are we?



Luca Naef
Co-Founder &
CTO

- BSc/MSc ETH Zurich – Molecular Bio & Deep Learning
- Research in Stanford, UNSW, TokyoTech
- Software Engineer & first tech startup during BSc
- Diverse roles in Biotech – Regeneus (AUS), CJ Partners (JPN)
- QuantumBlack & Mck – AI in Drug Discovery across Fortune 100/500
- Co-founded VantAI in 2019



Vladas Oleinikovas
Director of Comp Chem

- BA/MSci Cambridge – Nat. Sci. / Chem.
- PhD UCL – Chemistry
- Senior Scientist at UCB Pharma
- Acting Head of CADD at Monte Rosa Tx
- Co-inventor of clinical VAVI degrader
- Joined VantAI in 2024

Finding new medicines is hard

~1-3B

Cost to find a new drug¹

90+%

Chance of failure **after** entering human trials²

13+ yrs

Typical development **after** cause of disease identified

1. Estimates vary – e.g. from 0.88B (Eastern Research Group, “Drug Development Final Report”, Sept. 2024, for U.S. Department of Health and Human Services) to 2.6B (DiMasi et al., J Health Econ 2016)

2. Smietana et al., Nature Reviews Drug Discovery, 2016

3. Paul et al., Nature Reviews Drug Discovery, 2010



And has been getting harder

Moore's Law: Transistors

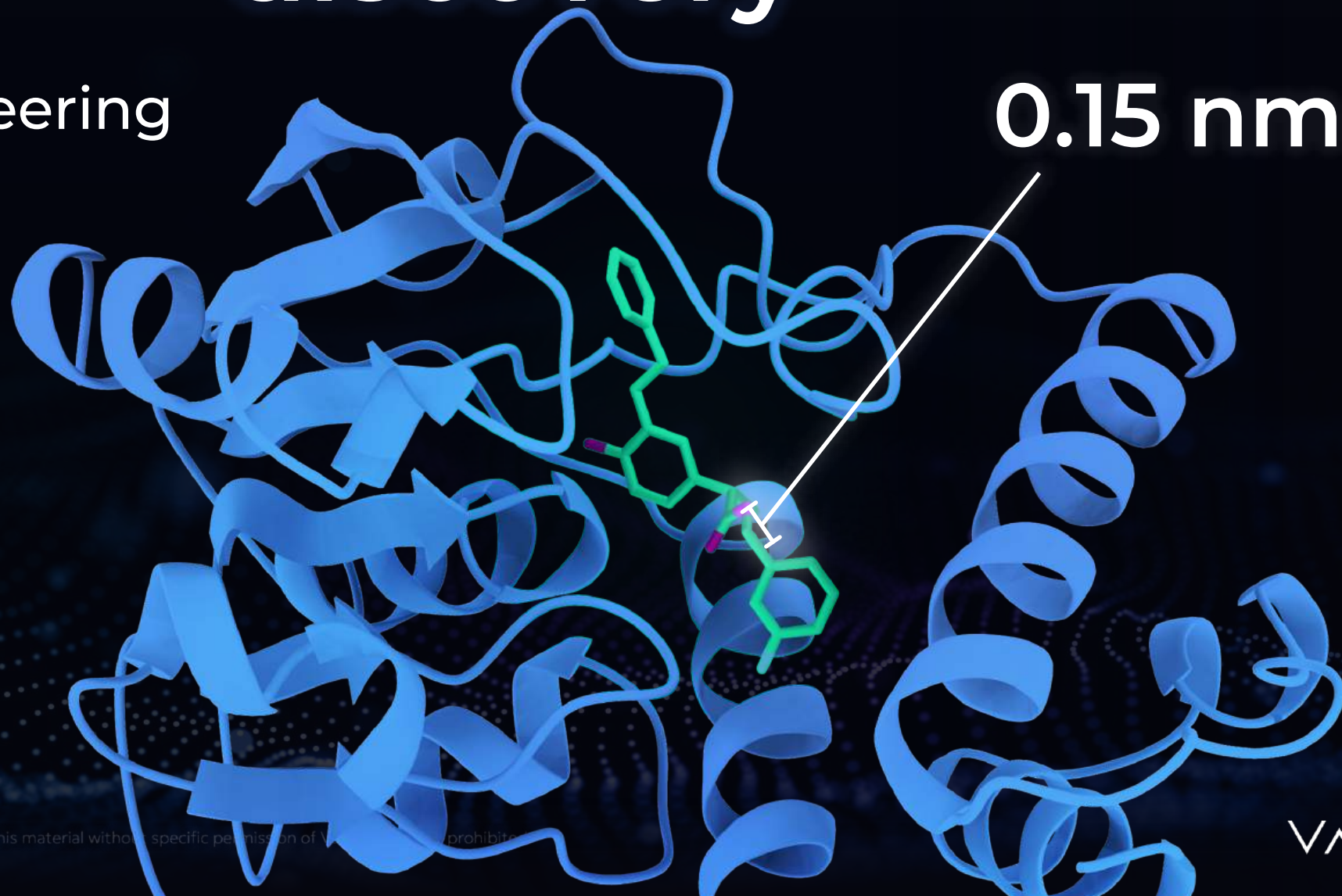


Eroom's Law: Medicines



Goal: Rational Drug discovery

Precision
nano-engineering
of therapies



This is, by no means, a new idea!

¹Medicinal Chemistry Department.
²Pharmacology Department.

Hitoshi Oinuma,^{1,2} Kazutoshi Miyake,¹
Motosuke Yamanaka,¹ Ken-Ichi Nomoto,²
Hiroshi Kato,² Kohji Sawada,²
Mitsumasa Shino,² Sachiyo Hamano
Tsukuba Research Laboratories
Eisai, Co., Ltd.
5-1-3, Tokodai, Tsukuba, Ibaraki 300-26, Japan
Received August 24, 1989

Neural Networks Applied to Structure-Activity Relationships

The neural network model was applied to the prediction of the activity of 100 compounds. The model was trained with 70 compounds and tested with 30 compounds. The model was able to predict the activity of the test compounds with a correlation coefficient of 0.85. The model was also able to predict the activity of new compounds. The model was able to predict the activity of new compounds with a correlation coefficient of 0.85. The model was able to predict the activity of new compounds with a correlation coefficient of 0.85.

Figure 1. Schematic diagram of the neural network model.

Table 1. Data Set and Classification of Compounds

Compound	Activity	Classification
1	0.1	Low
2	0.2	Low
3	0.3	Low
4	0.4	Low
5	0.5	Low
6	0.6	Low
7	0.7	Low
8	0.8	Low
9	0.9	Low
10	1.0	Low
11	1.1	Low
12	1.2	Low
13	1.3	Low
14	1.4	Low
15	1.5	Low
16	1.6	Low
17	1.7	Low
18	1.8	Low
19	1.9	Low
20	2.0	Low
21	2.1	Low
22	2.2	Low
23	2.3	Low
24	2.4	Low
25	2.5	Low
26	2.6	Low
27	2.7	Low
28	2.8	Low
29	2.9	Low
30	3.0	Low
31	3.1	Low
32	3.2	Low
33	3.3	Low
34	3.4	Low
35	3.5	Low
36	3.6	Low
37	3.7	Low
38	3.8	Low
39	3.9	Low
40	4.0	Low
41	4.1	Low
42	4.2	Low
43	4.3	Low
44	4.4	Low
45	4.5	Low
46	4.6	Low
47	4.7	Low
48	4.8	Low
49	4.9	Low
50	5.0	Low
51	5.1	Low
52	5.2	Low
53	5.3	Low
54	5.4	Low
55	5.5	Low
56	5.6	Low
57	5.7	Low
58	5.8	Low
59	5.9	Low
60	6.0	Low
61	6.1	Low
62	6.2	Low
63	6.3	Low
64	6.4	Low
65	6.5	Low
66	6.6	Low
67	6.7	Low
68	6.8	Low
69	6.9	Low
70	7.0	Low
71	7.1	Low
72	7.2	Low
73	7.3	Low
74	7.4	Low
75	7.5	Low
76	7.6	Low
77	7.7	Low
78	7.8	Low
79	7.9	Low
80	8.0	Low
81	8.1	Low
82	8.2	Low
83	8.3	Low
84	8.4	Low
85	8.5	Low
86	8.6	Low
87	8.7	Low
88	8.8	Low
89	8.9	Low
90	9.0	Low
91	9.1	Low
92	9.2	Low
93	9.3	Low
94	9.4	Low
95	9.5	Low
96	9.6	Low
97	9.7	Low
98	9.8	Low
99	9.9	Low
100	10.0	Low

The computer program LUDI: A new method for the de novo design of enzyme inhibitors

Hans-Joachim Böhm

BASF AG, Central Research, D-6500 Ludwigshafen, Germany

Received 27 May 1991

Accepted 16 August 1991

Key words: Enzymes; Enzyme inhibitors; Molecular modeling; Drug design; De novo design

64'

81'

82'

90'

92'

95'

2015

2020

2024

Generations

QSAR Models via regression:

Corwin Hansch and Toshio Fujita

Captopril: first "rational" structure-based Drug approved in 1981

Early Docking tools:

DOCK, Irwin D. Kuntz

NN for QSAR. Toshihisa Aoyama et al. (J. Med. Chem., 1990)

"De-novo" generative method via rules-based fragment assembly. Boehm, 1992, Journal of Computer-Aided Molecular Design

Computational modelling drives new class of anti-HIV drugs: saquinavir (approved 1995) indinavir (1996), zidovudine (1996), and nelfinavir (1997)

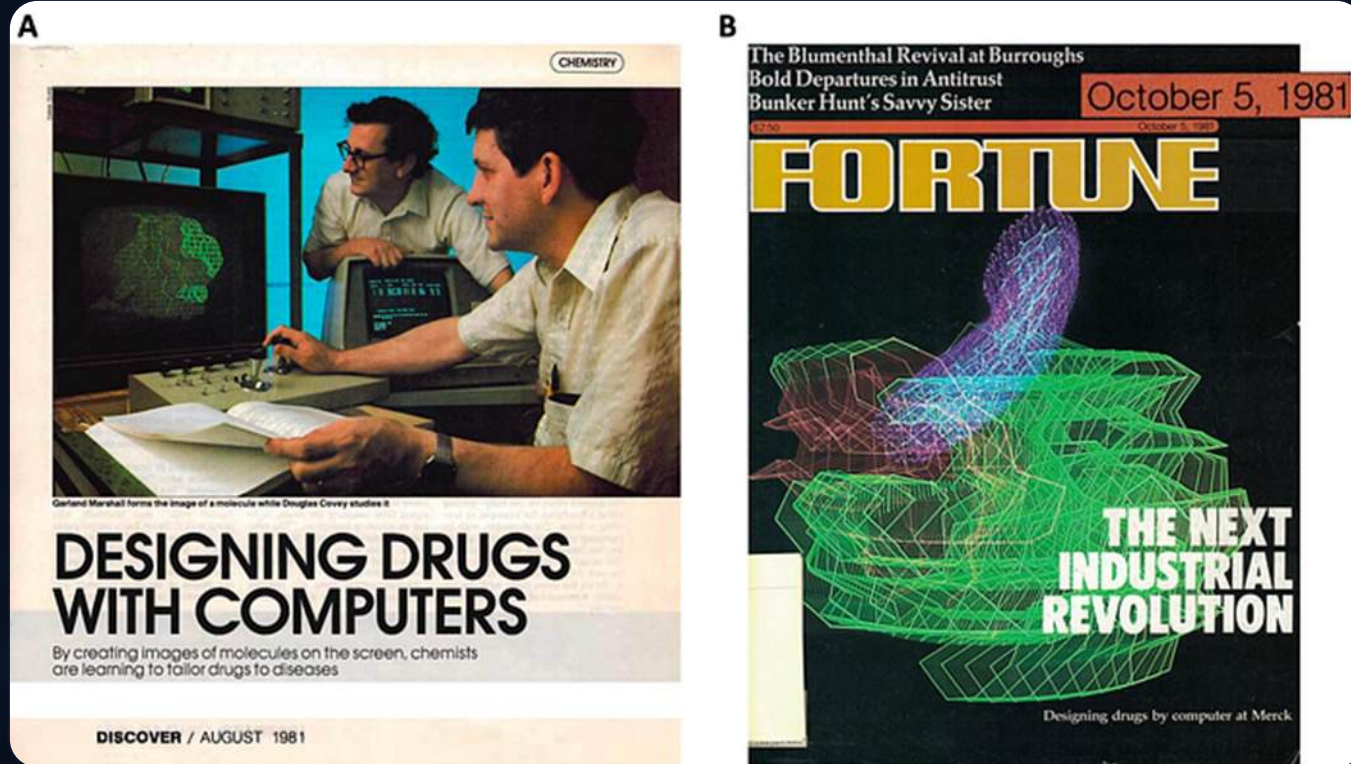
AtomNet.

Abraham Heifets & Izhar Wallach. First CNN for docking rescoring

Exscientia & DSP launch first Ph I of AI de-novo designed drug

Insilico Medicine announces first-in-class fibrosis drug (TNIK inhibitor INS018_055) to reach PhII

1981



Xu, The path to the next computational transformation of drug discovery, Medium, 2022

However, we should not forget that **irrational** drug discovery is incredibly successful



haloperidol
1958



"It was a matter of life or death — a matter of survival. I don't believe we really had a clear-cut strategy. We were simply doing whatever we could, and there weren't many things that we could do then. We didn't have much money and there were not many researchers. We had to make a lot of simple compounds as quickly and possible and screen them using very simple methods"

Paul Janssen

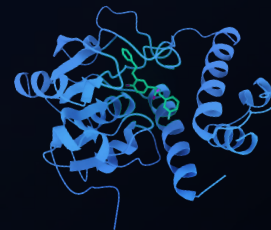
The classical modelling tools focus to solve one task

**Structure
+ Molecule/**



Small molecule docking

Molecule



- Optimized to solve one task efficiently, eg. docking: sampling + scoring
- Relies on specific (often rigid) input, that may be prohibitively expensive (eg. X-ray structure) and limited by applicability

Monomers



Protein-Protein
Docking

Complex



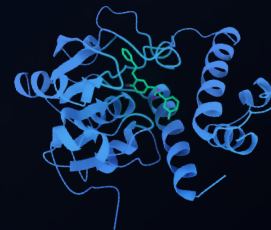
The classical modelling tools focus to solve one task

Structure + Molecule/



Small molecule docking

Molecule



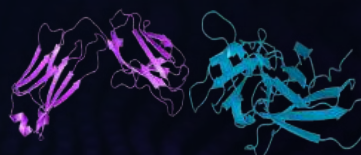
- Optimized to solve one task efficiently, eg. docking: sampling + scoring
- Relies on specific (often rigid) input, that may be prohibitively expensive (eg. X-ray structure) and limited by applicability

Monomers

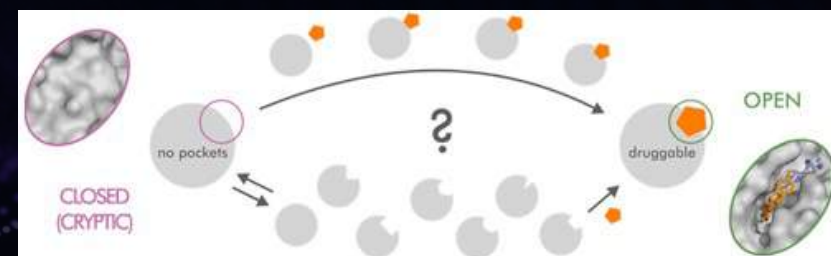


Protein-Protein
Docking

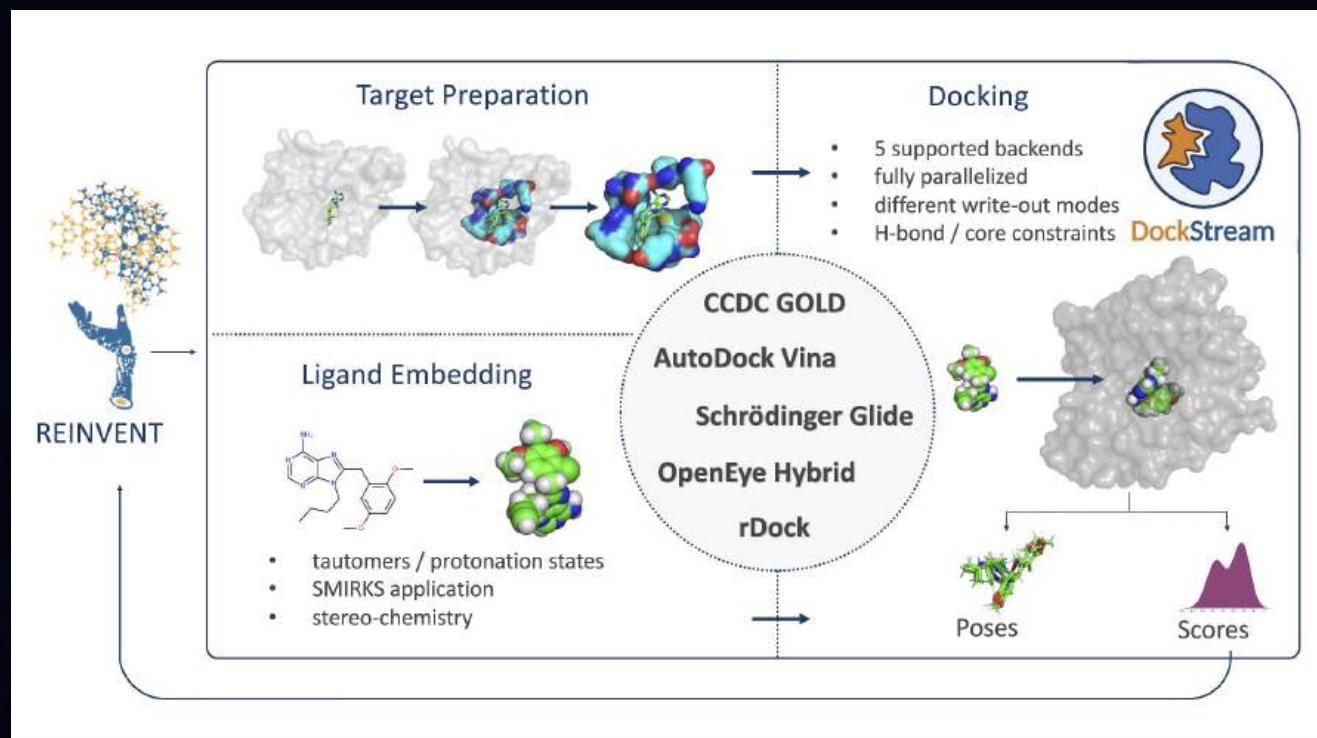
Complex



- Scaling with Moore's law:
 - ✓ run more iterations, larger libraries / systems
 - ✗ sampling limited by inputs (esp. receptor flex.)



Complex methods limited to sequential combination of tasks



<https://github.com/MolecularAI/DockStream>
<https://pubs.acs.org/doi/10.1021/acs.jcim.0c01451>

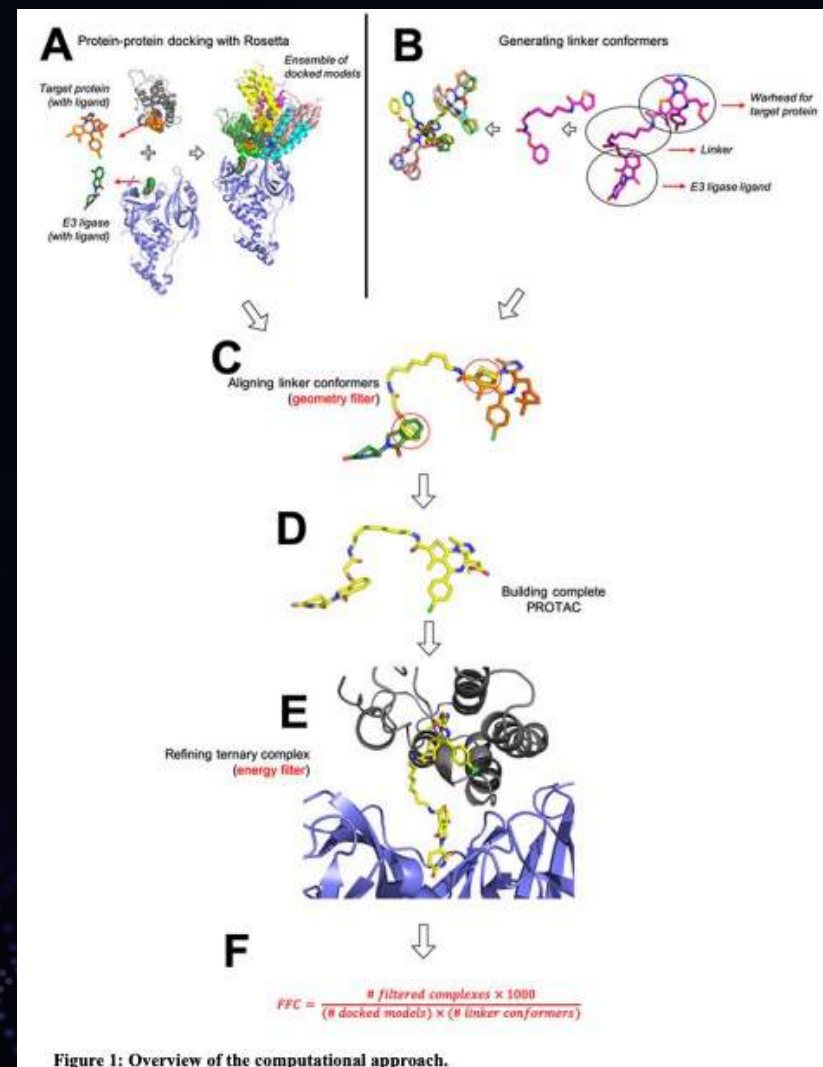
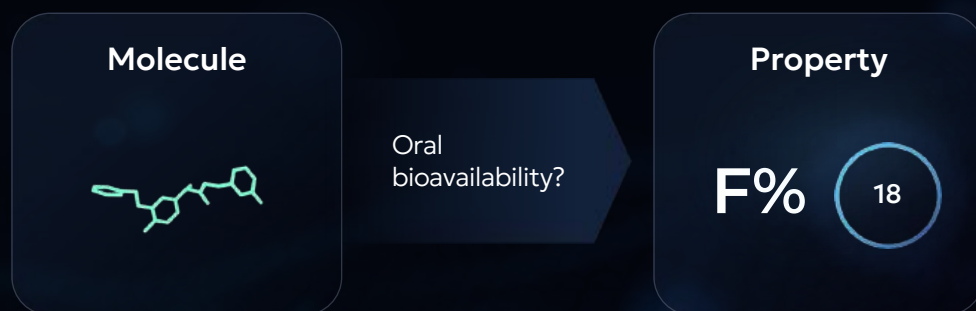


Figure 1: Overview of the computational approach.

There's broadly two flavors of machine learning that are used

Discriminative



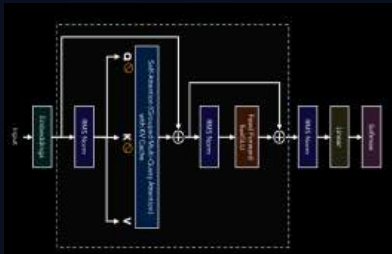
Generative



Two methods have been particularly critical
in the emergence of generative methods

Language models

CC(O)...



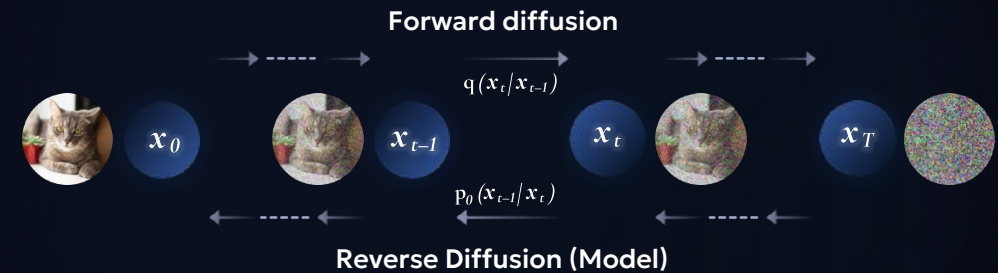
C: 0.2

N: 0.4

...: 0.X

- Trained to predict token: either next token (e.g. ChatGPT, “autoregressive”) or some masked tokens in middle of sequence (e.g ESM3)
- Breaks down problem into simpler steps: one token at a time
- Tokens require discrete data
- Became SOTA for language ~= with GPT-1 in 2018

Diffusion models



- Trained to remove simple noise added during training process
- Breaks down problem into simpler steps: only predict “a bit” of noise – makes problem easier since e.g. high and low noise steps often very different problems
- Usually used for continuous data: images, coordinates, expression levels, ...
- Usually a lot more fancy math and many extensions (Flows, Schrodinger Bridges, ...)
- Became SOTA for images ~= with ADM in 2021

Diffusion: A breakthrough approach for generative modelling

- Recent breakthroughs leverage Diffusion models
- Inspired by physical diffusion processes (brownian motion)
- Unsupervised training - stepwise noise addition to ground truth
- Model learns going from noised sample to ground truth
- Noise depends on problem at hand – Gaussian noise on pixels, latent noise or more complex (SO(3)-subspace diffusion)

Generative Adversarial Networks
(2015)

Images of flowers



GANs & Roses

Diffusion Models
(2020)



“Show a molecular glue holding together two planets”

OpenAI

DALL·E 3



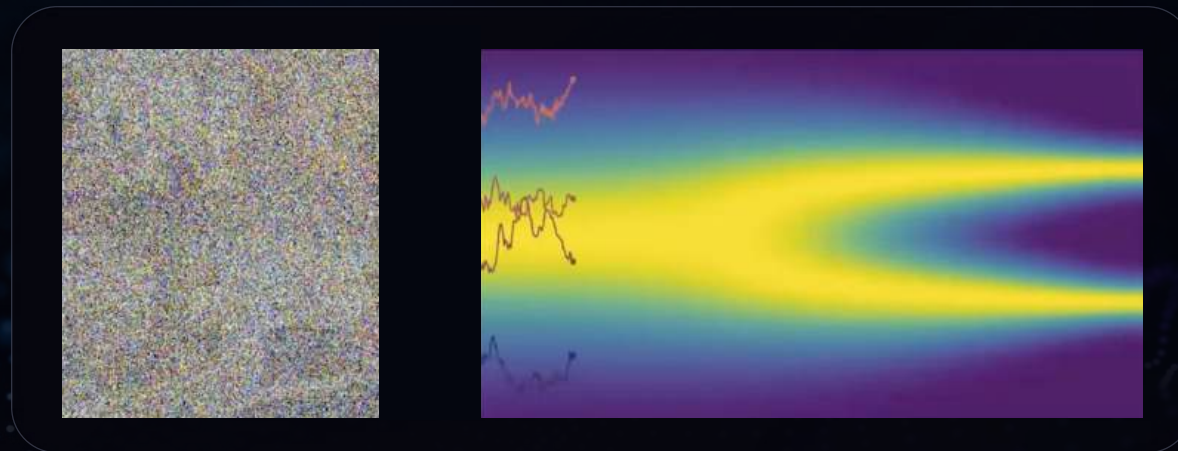
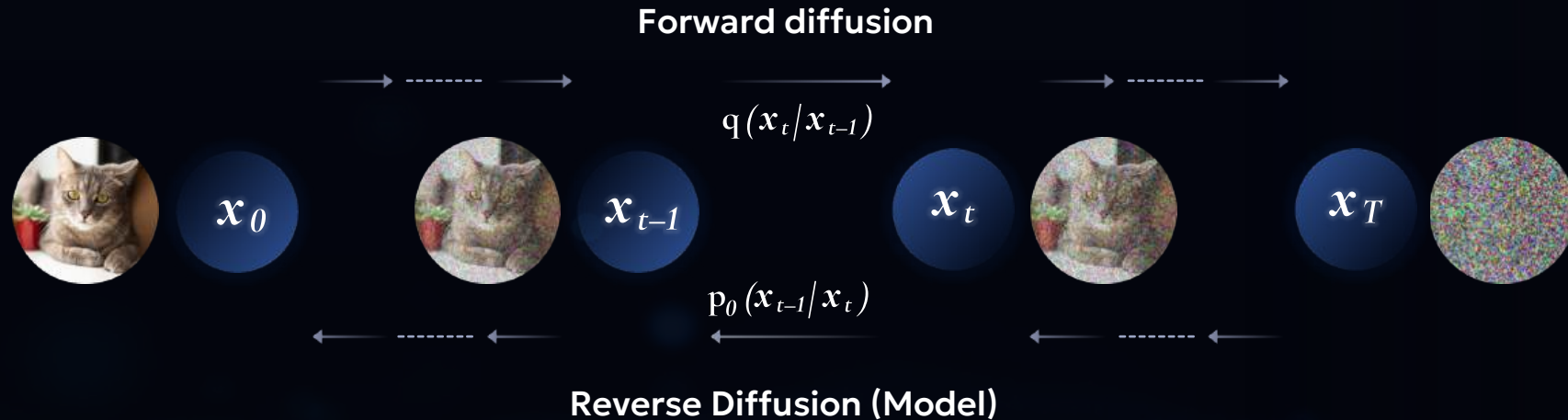
Stable Diffusion



Midjourney

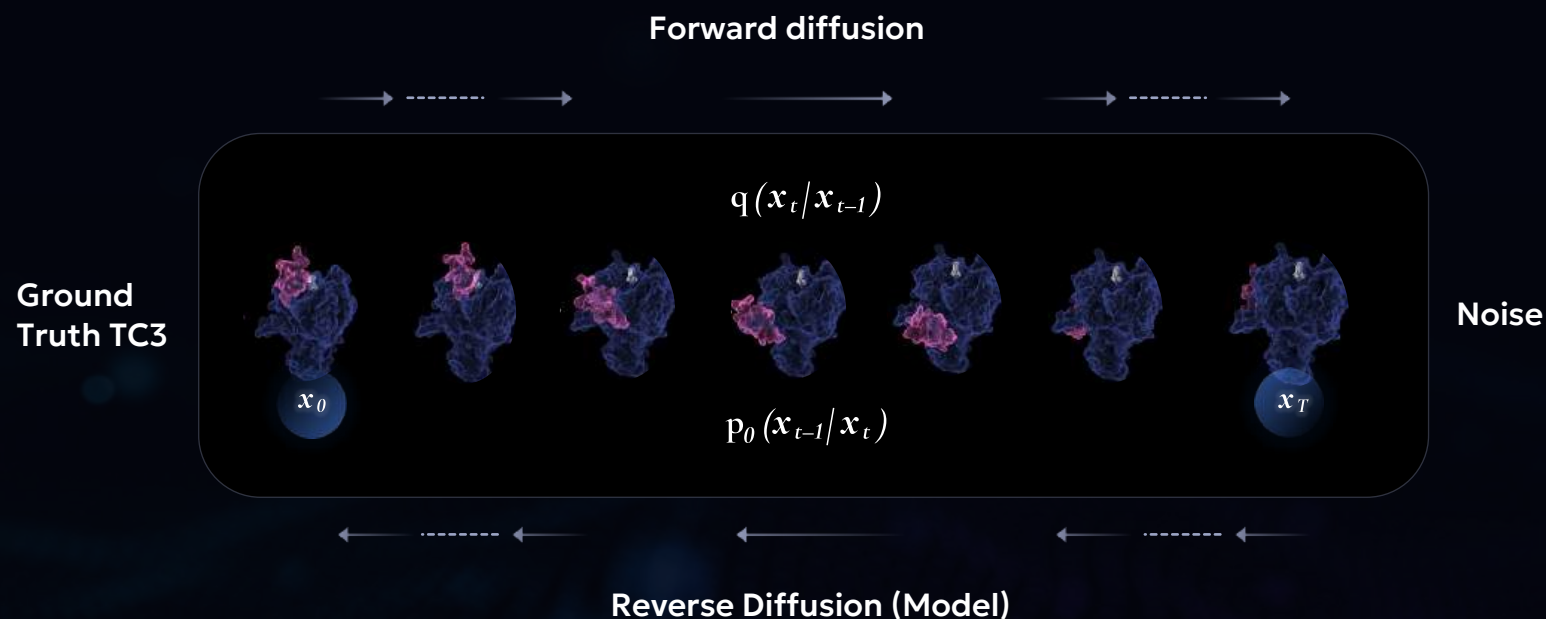
Imagen

Diffusion Models are trained unsupervised by adding noise to training data and predicting ground truth



Diffusion probabilistic models - Jascha Sohl-Dickstein, Google Brain Talks
<https://www.youtube.com/watch?v=XCUIhHP1TNM>

In molecular modelling, “noise” can be added in highly flexible ways

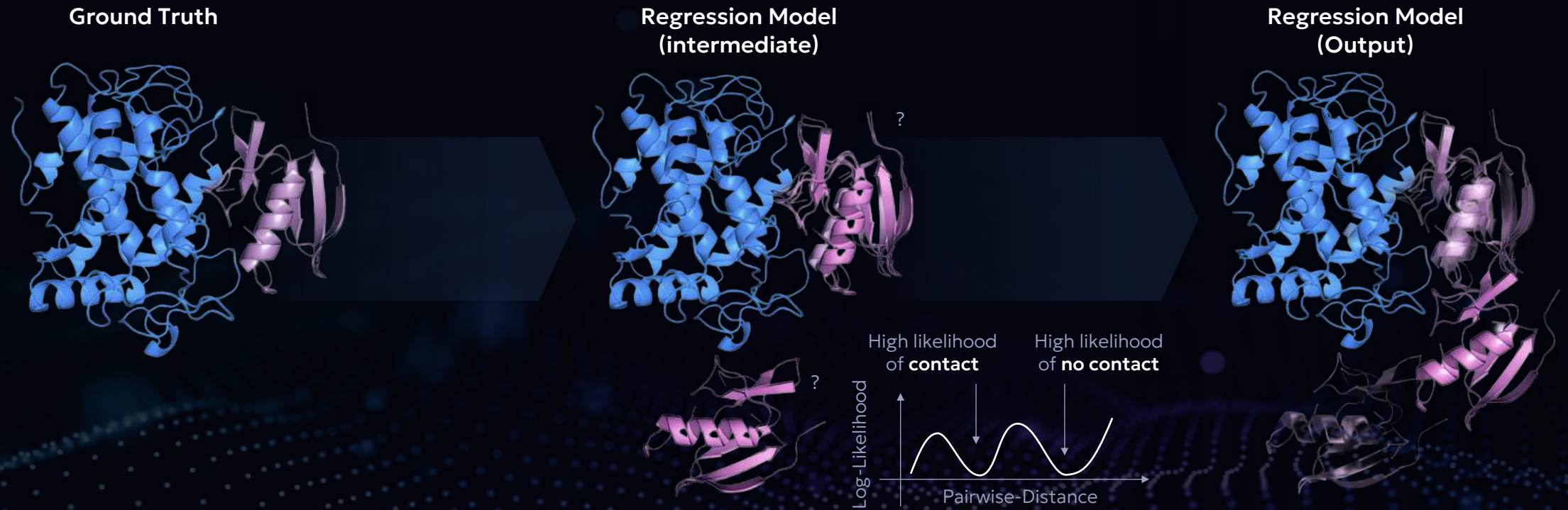


Noise types

- Example: SE(3) diffusion (rotation, translation): **DiffMaSIF**
- Euclidean diffusion: **ApolloDiff**
- Latent Diffusion: **LatentDiff**

Why use Diffusion for structure prediction problems?

Apart from the fact the protein structures are not static,
Regression models (e.g. AF2) have a major issue

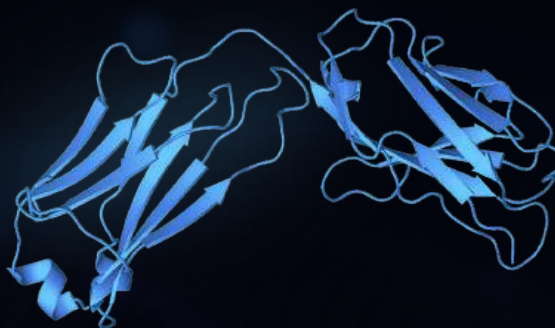


AI is generally used to predict one or multiple of the following given one or multiple of them as input

Sequence

MGRLQLVVLGREDAHFIYENKDVSQ...

Structure

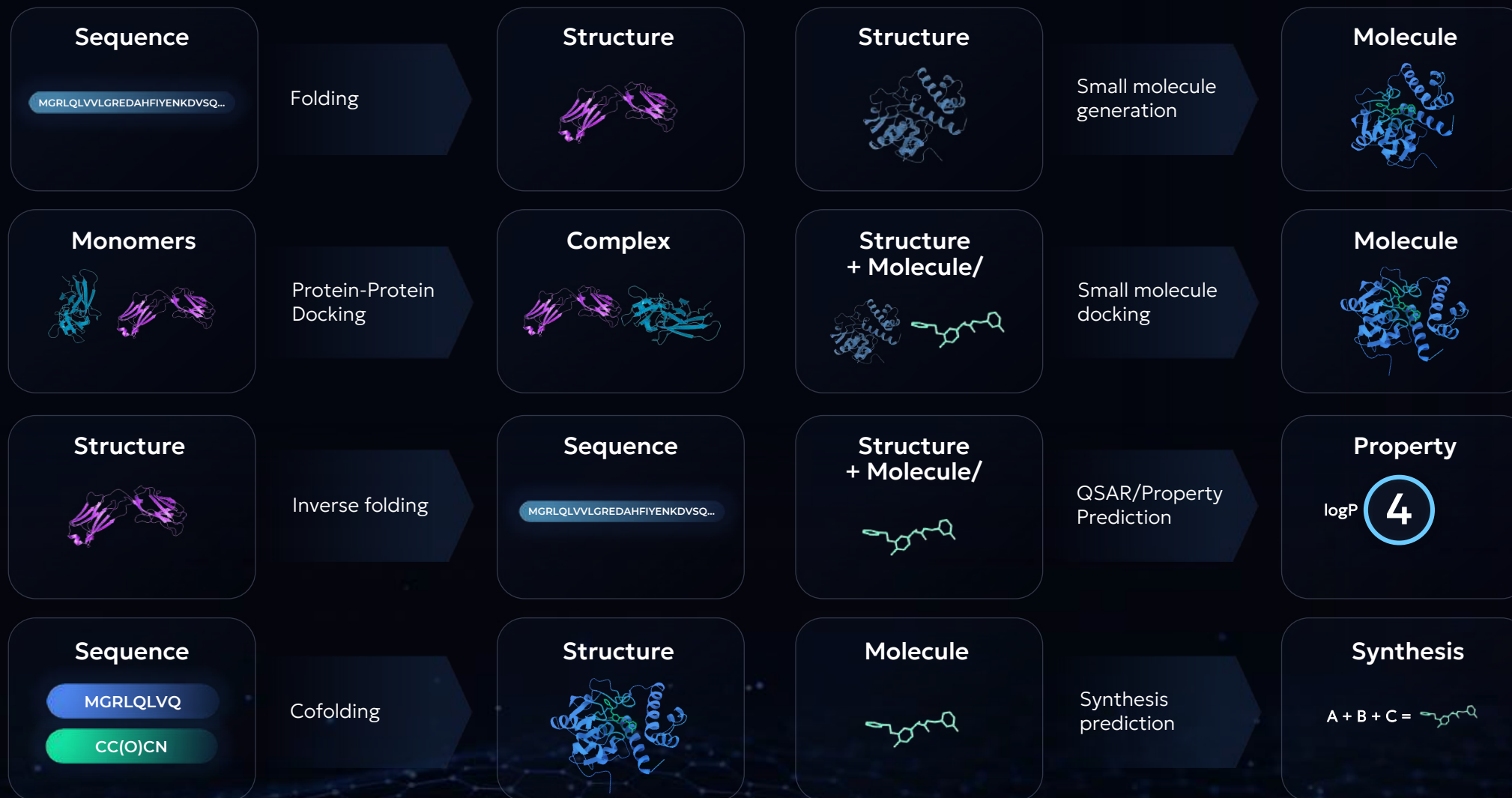


Property

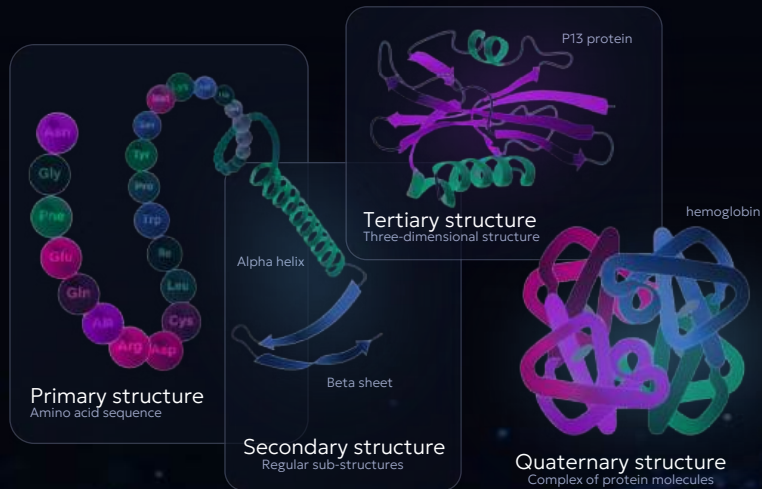
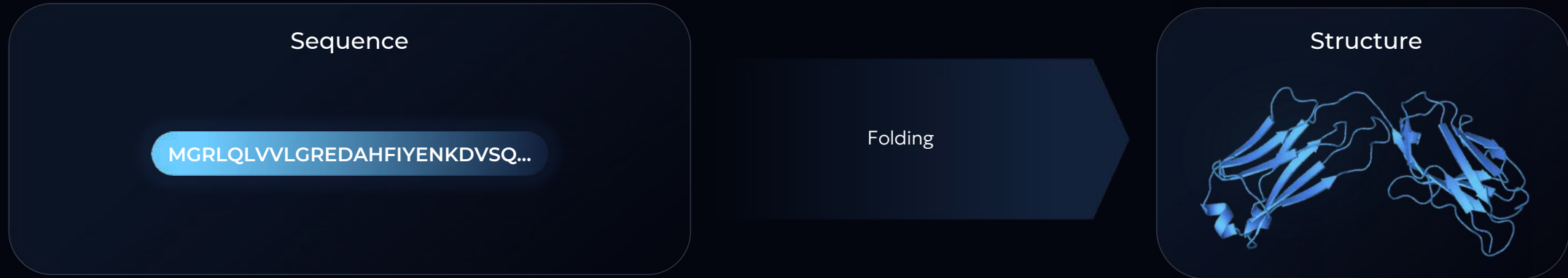
Kinase

Early ML methods mostly followed the classic methods focusing on the same tasks following the same limitations...

This has led to a large "Zoo" of models in use



Breakthrough: Protein folding



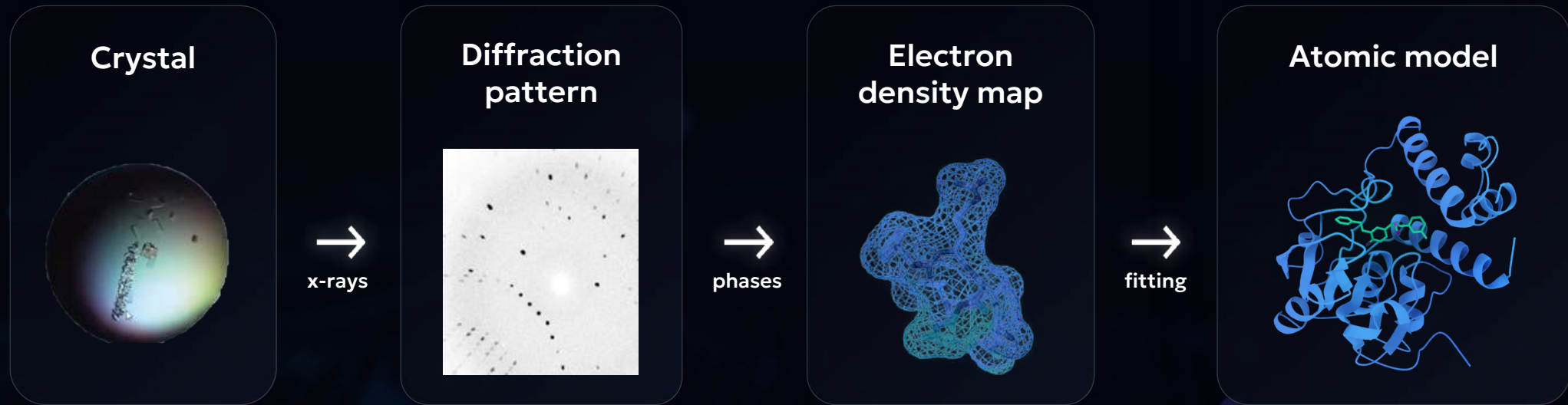
**A protein's structure
is uniquely determined
by its sequence**

**“Anfinsen’s Dogma”
– 1972 Nobel Prize**



Traditional route:

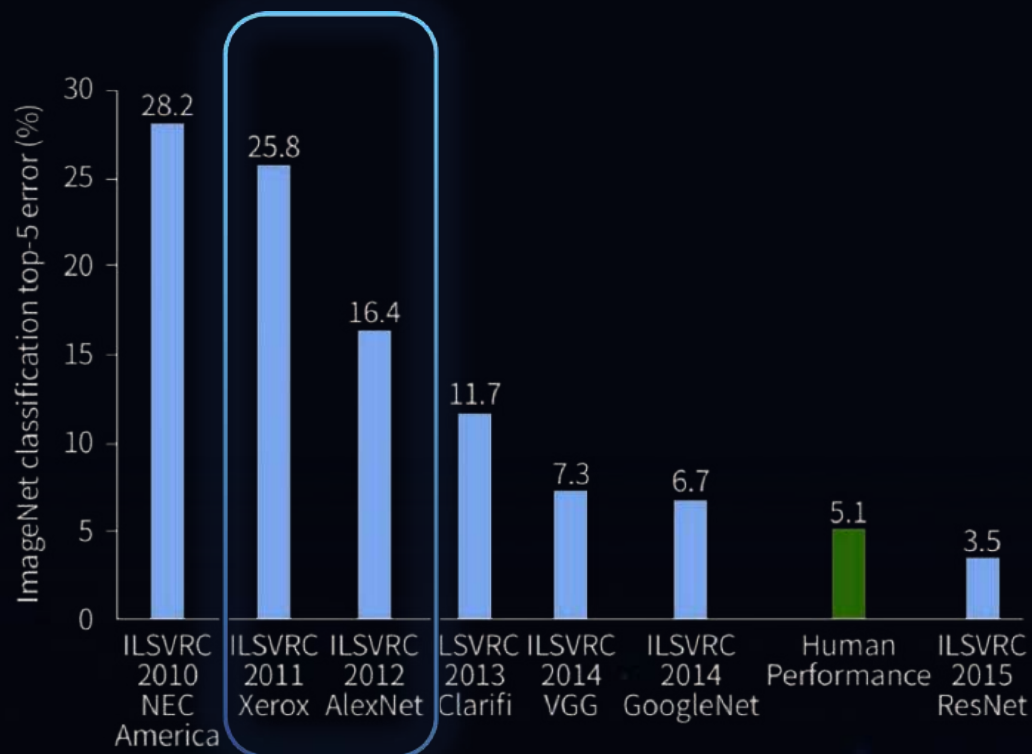
X-ray crystallography



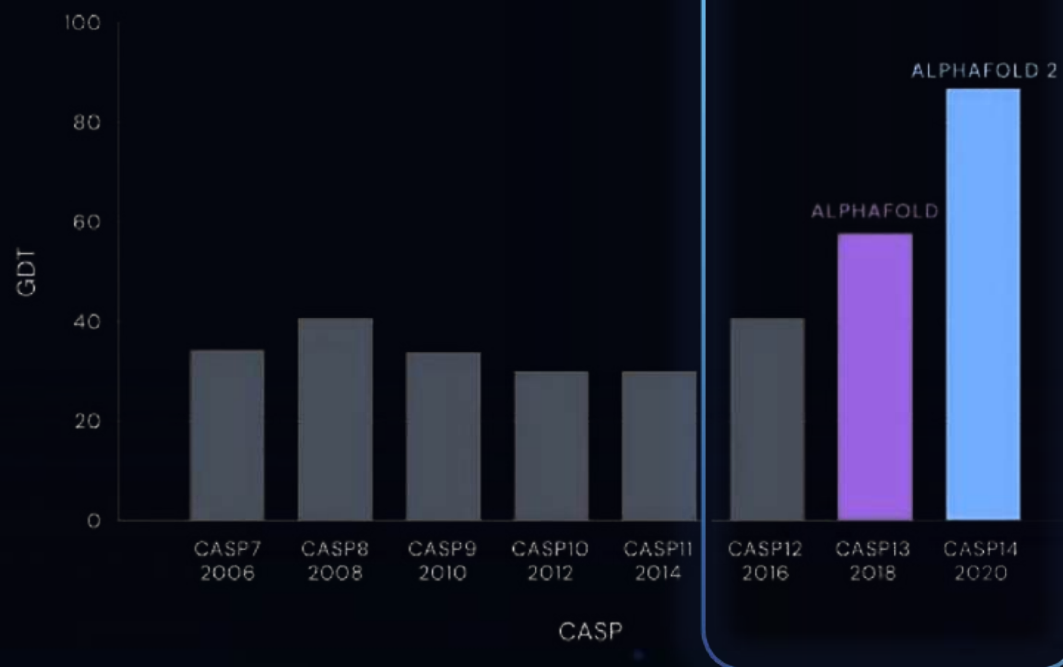
Up to **1M USD**

Can take **years**

AF2 was an “AlexNet” like moment



Median Free-Modelling Accuracy



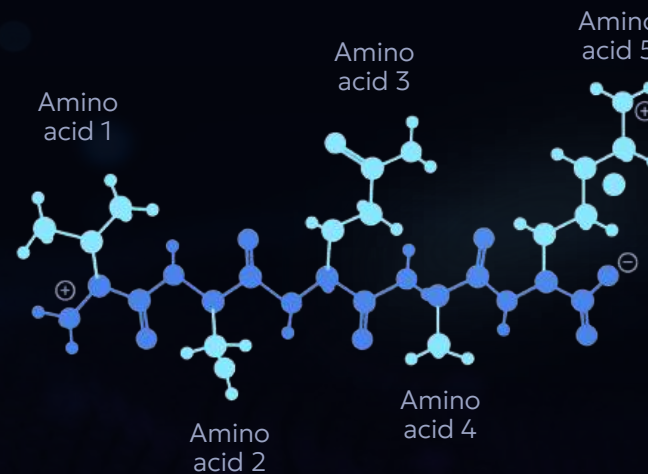
Proteins are polymers comprised of amino acids



UniRef90_P02057 Hemoglobin subunit beta-1/2 n=6
Tax=Euarchontoglires TaxID=314146 RepID=HBELRABIT
MVHLSSEEKSAVTALWGKVNVEEVGGEALGRLLVVYPWT
QRFFESFGDLSSANAVMNWKVKAHGKKVLAASFEGLSHL
DNLKGTFAKLSELHCDKLHVPENERLLGNVLVIVLSHHFG
KEFTPQVQAAYQKVAVAGVANALAHKYH

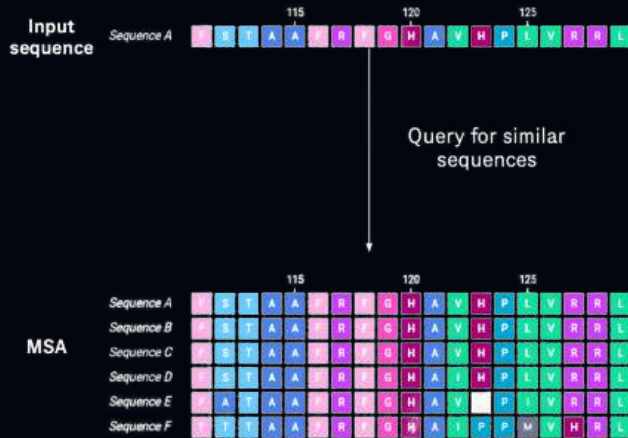
Backbone

Side chain

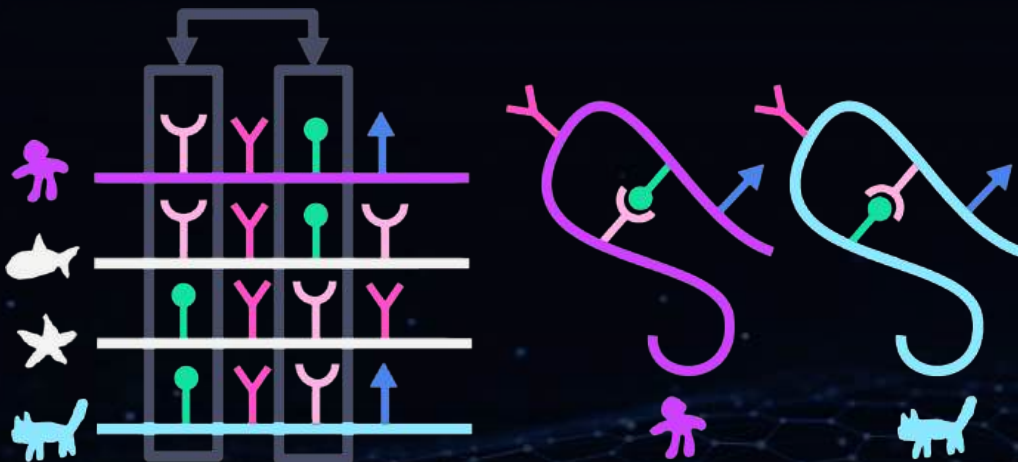


- We usually split structures in backbone (mainchain) & sidechain
- Backbone structure often assumed +/- independent of specific sidechains
- E.g. for protein design: often design backbones and then predict different residues that could take this shape

However, technically AlphaFold2 doesn't fold protein sequences, it predicts structures based on Multiple Sequence alignment



- MSA gives us **co-evolution**
- Residue pairs that co-evolve (e.g. charge switch on one always observed with opposing charge switch on other) are likely in contact
- Often ignored: it's also a form of **RAG** (retrieval augmented generation) – giving large number of sequences that likely fold into same structure makes sequence-based fold retrieval more robust

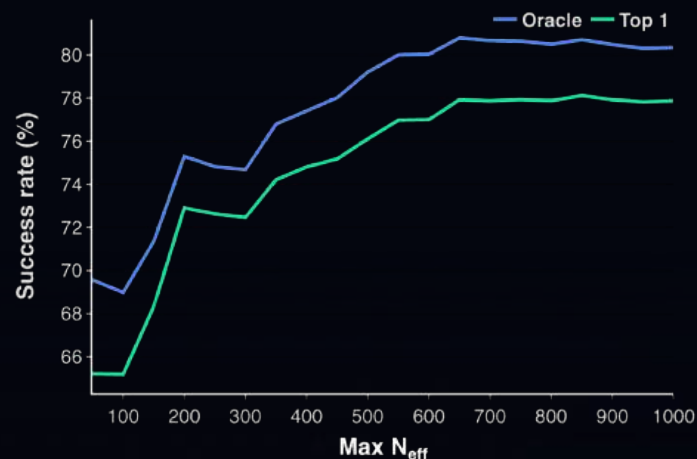


In fact, co-evolutionary information is absolutely critical



Source: Zeming Lin et al., 2021

AF2 PINDER performance by MSA depth

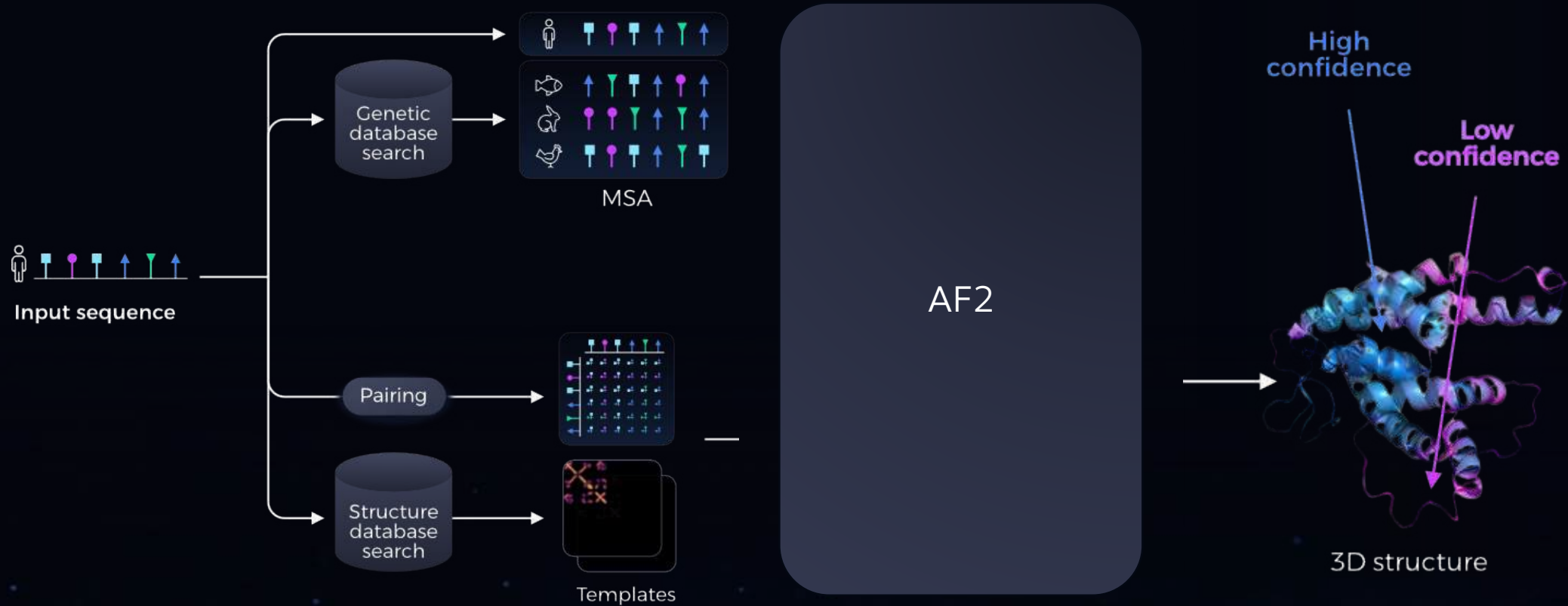


Source: Kovtun et al. PINDER, 2024 (VantAI)

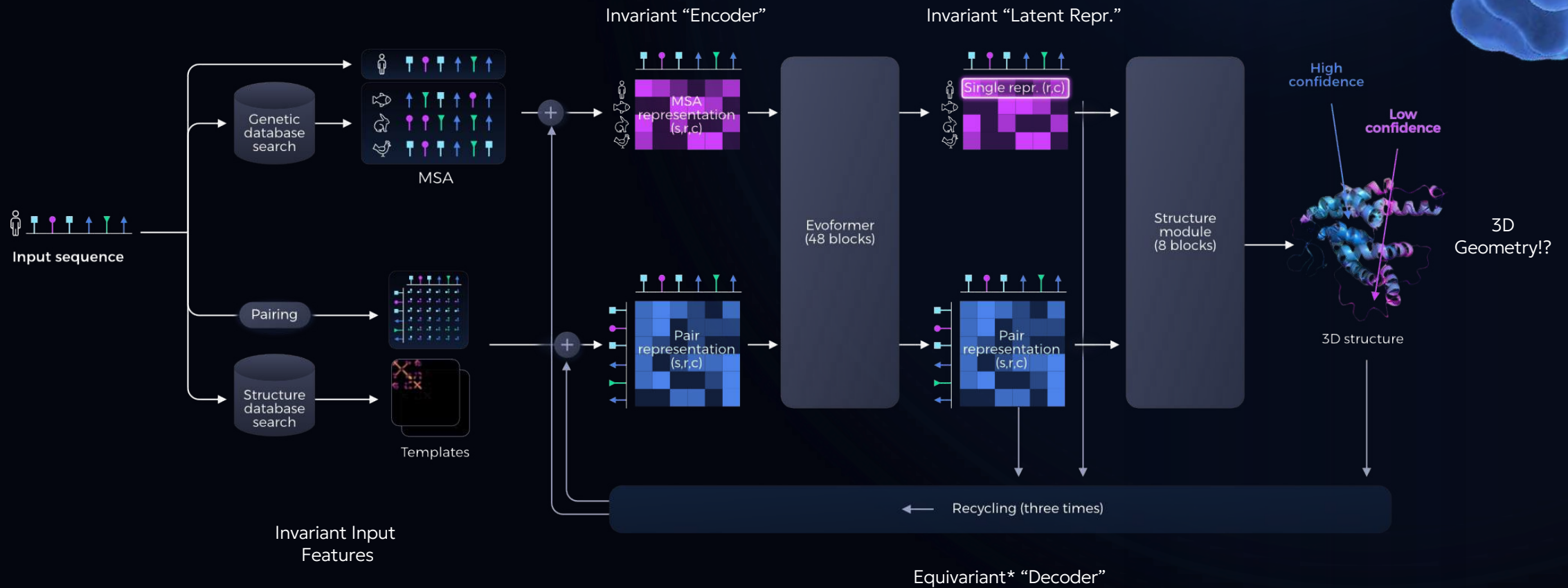
Notes:

- For both monomers and protein-protein interfaces, MSA are absolutely critical information
- While current language models such as ESM3 have been shown to learn the pairwise residue covariance implicitly, despite 100B+ they still underperform a sequence alignment

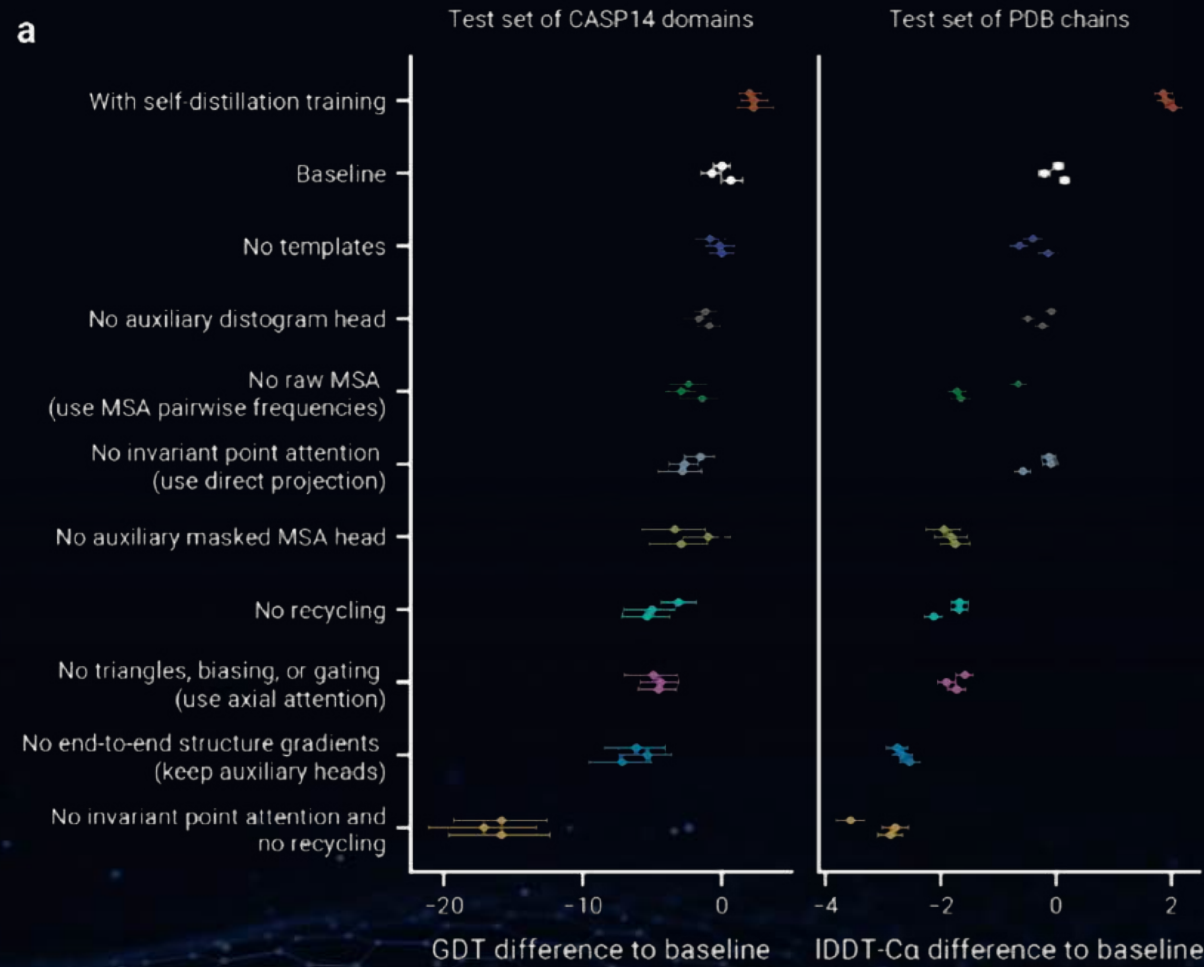
AF2 is a model that uses sequence and co-evolution to predict a protein's structure



It uses a transformer inspired architecture to first process sequence & MSA features and then reason over coordinates

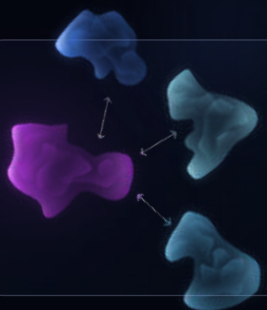


Critically – AF2 is an “engineering marvel” with many, many performance critical innovations



De-novo molecular generators for variety of use cases

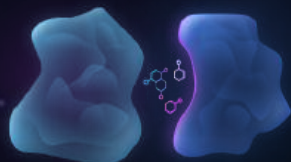
Example use-cases
for Proximity Modulators:



De-novo
generative
binder or
glue design



Generative
Linker Design



Generative
Lead opt

Generative Algorithms



Diffusion/Flow Matching
Probabilistic Models



LLMs +
RL/RLHF



Genetic algorithms
or search-based
generative methods

Scoring functions



Shape-based



Free Energy
& interaction-based



ADME/PK
Property classifiers

MaSIF & dMaSIF have been highly successful for protein-design applications

De novo design of protein interactions with learned surface fingerprints

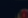
<https://doi.org/10.1038/s41586-023-05993-x>

Received: 16 June 2022

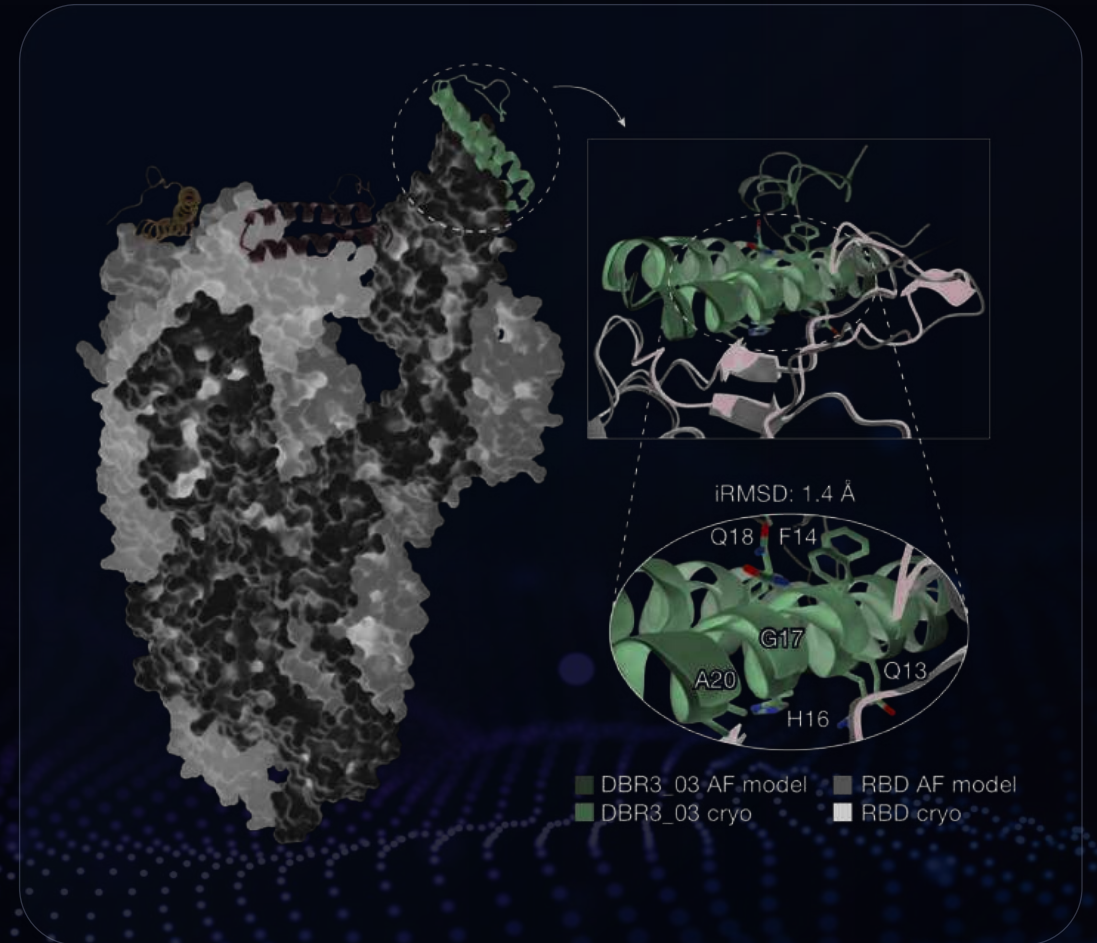
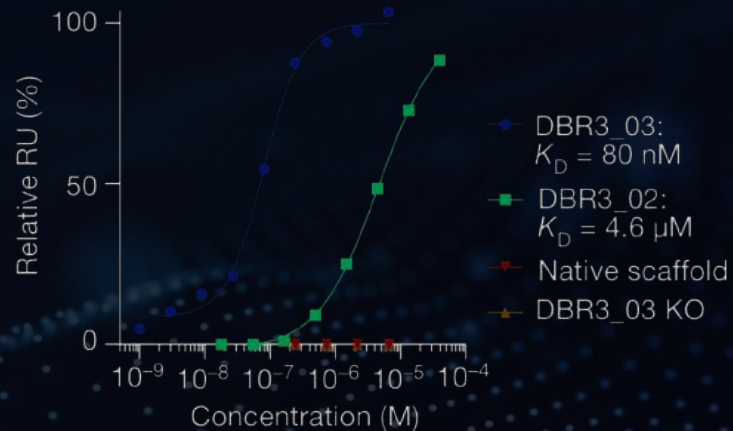
Accepted: 21 March 2023

Published online: 26 April 2023

Open access

 Check for updates

Pablo Gainza^{1,2,10}, Sarah Wehrle^{1,2,11}, Alexandra Van Hall-Beauvais^{1,2,11}, Anthony Marchand^{1,2,11}, Andreas Schreck^{1,2,11}, Zander Harteveld^{1,2}, Stephen Buckley^{1,2}, Dongchun Ni^{1,4}, Shuguang Tan⁵, Freyr Sverrisson^{1,2}, Casper Goverde^{1,2}, Priscilla Turelli⁶, Charlene Raclot⁶, Alexandra Teslenko⁷, Martin Pacesa^{1,2}, Stephane Rosset^{1,2}, Sandrine Georgeon^{1,2}, Jane Marsden^{1,2}, Aaron Petruzzella⁸, Kefang Liu², Zepeng Xu⁹, Yan Chai³, Pu Han⁹, George F. Gao⁹, Elisa Oricchio⁶, Beat Fierz¹, Didier Trono⁶, Henning Stahlberg^{1,4}, Michael Bronstein^{9,12} & Bruno E. Correia^{1,2,13}



A variety of generative algorithms have been used,
usually split into 2D or 3D representations

Generative Algorithms



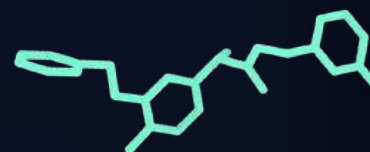
Diffusion/Flow Matching
Probabilistic Models



LLMs +
RL/RLHF



Genetic algorithms
or search-based
generative methods

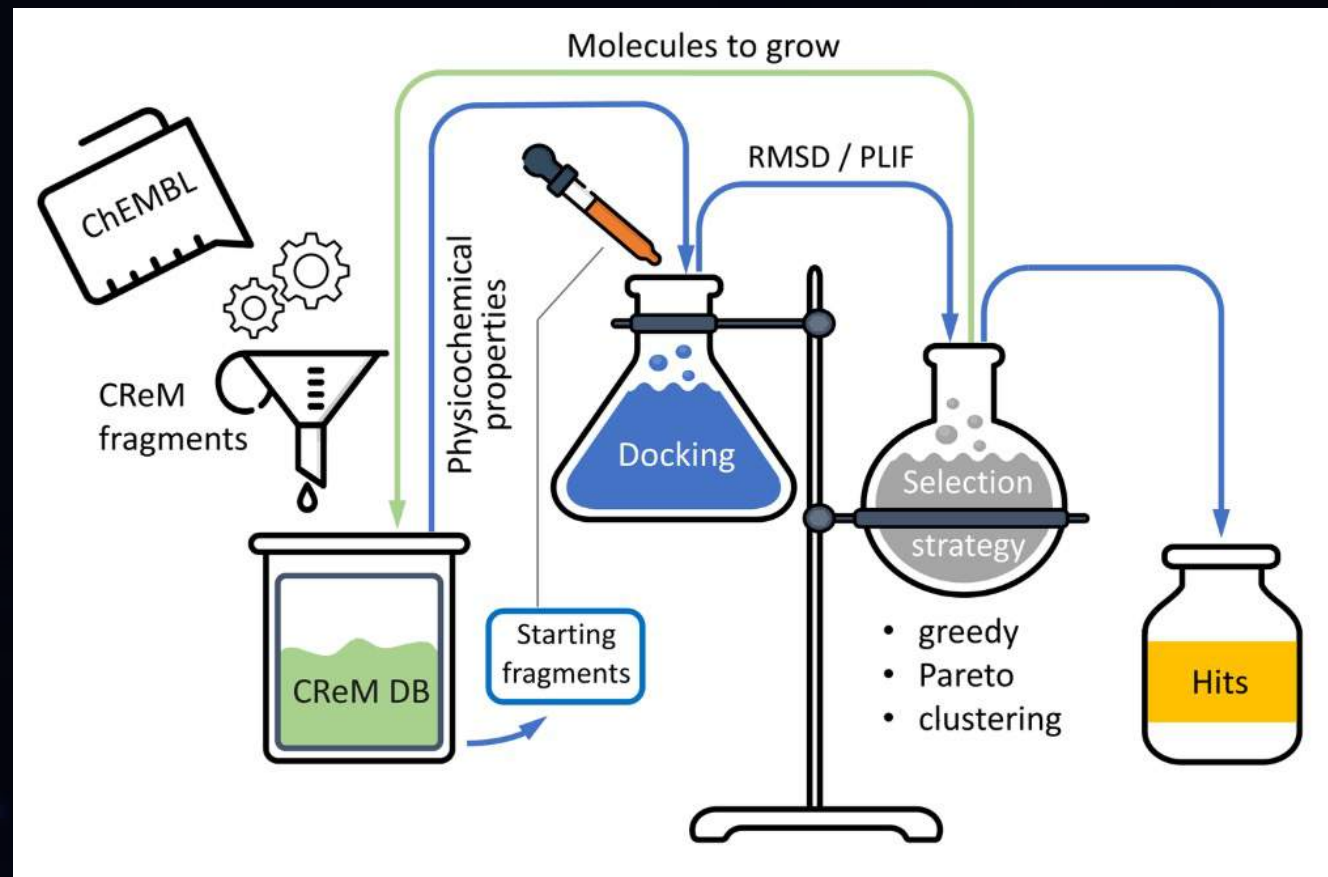
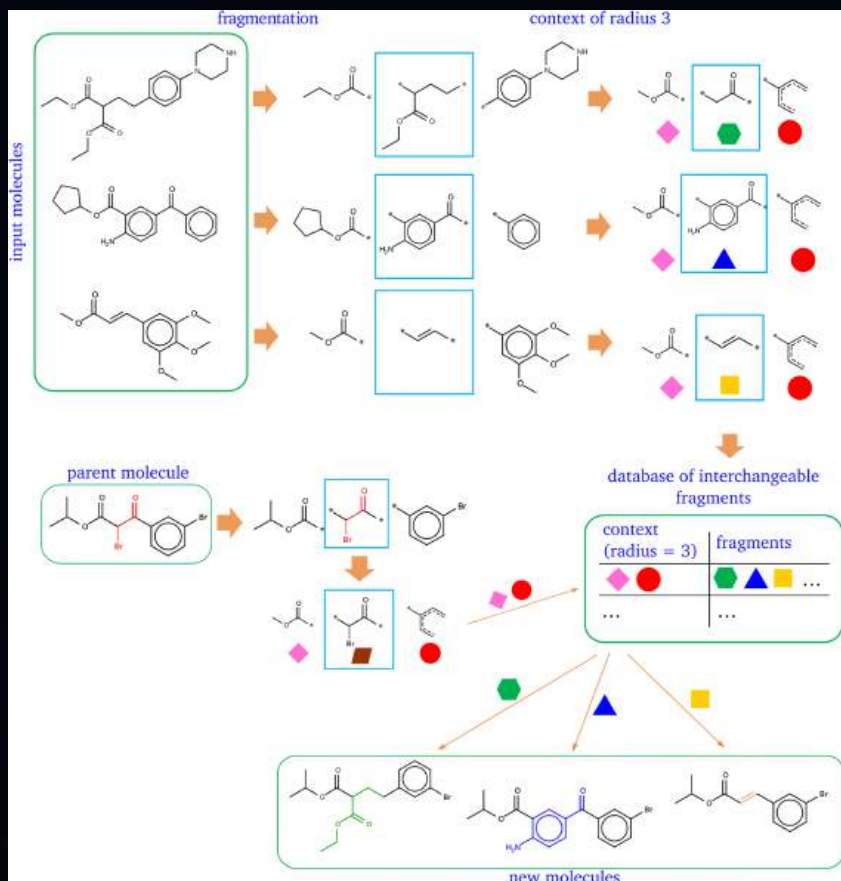


3D

CC(O)CN

2D

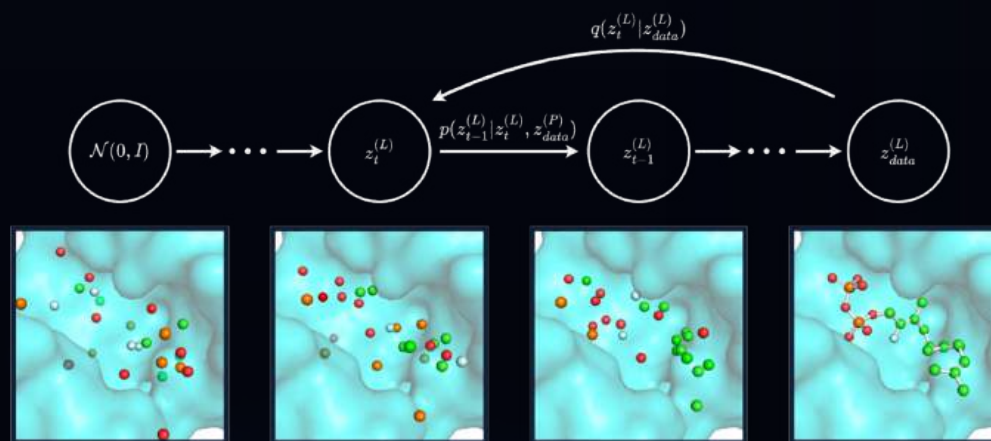
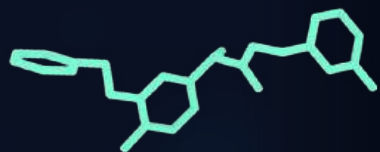
Classical approaches: rule-based (iterative) enumerators



<https://github.com/ci-lab-cz/crem-dock>

A variety of generative algorithms have been used, usually split into 2D or 3D representations

3D



Structure-based Drug Design with Equivariant Diffusion Models. Arne Schneuing, Charles Harris, Yuanqi Du, Kieran Didi, Arian Jamasb, Ilija Iqashov, Weitao Du, Carla Gomes, Tom Blundell, Pietro Lio, Max Welling, Michael Bronstein, Bruno Correia. ArXiv, 2022

2D

CC(O)CN

Graph:



SMILES:

ClCc1c[nH]cn1

One-hot encoding:

	Cl	C	c	1	c	nH	c	n	1
C	0	1	0	0	0	0	0	0	0
c	0	0	1	0	1	0	1	0	0
n	0	0	0	0	0	1	0	1	0
1	0	0	0	1	0	0	0	0	1
nH	0	0	0	0	0	1	0	0	0
Cl	1	0	0	0	0	0	0	0	0

Fig. 3 Three representations of 4-(chloromethyl)-1H-imidazole. Depiction of a one-hot representation derived from the SMILES of a molecule. Here a reduced vocabulary is shown, while in practice a much larger vocabulary that covers all tokens present in the training data is used

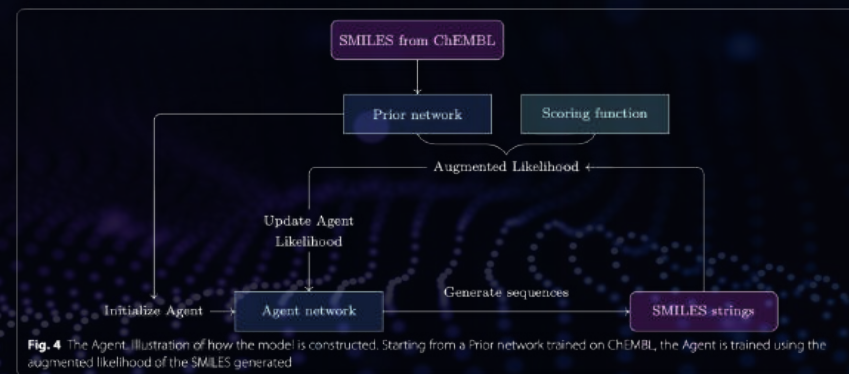


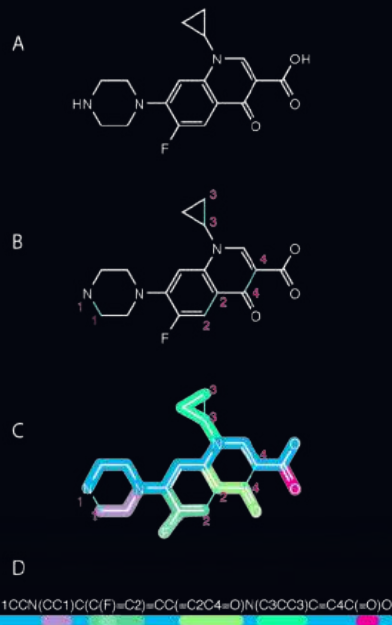
Fig. 4 The Agent. Illustration of how the model is constructed. Starting from a Prior network trained on ChEMBL, the Agent is trained using the augmented likelihood of the SMILES generated

Molecular de-novo design through deep reinforcement learning
Marcus Olivecrona, Thomas Blaschke, Ola Engkvist and Hongming Chen. J Cheminform (2017)

SMILES is a string representation of a molecular graph which allows standard LLM tokenization approaches

2D

CC(O)CN



Graph:

SMILES:

ClC1c[nH]cn1

One-hot encoding:

	Cl	C	c	1	e	nH	c	n	1
C	0	1	0	0	0	0	0	0	0
c	0	0	1	0	0	0	0	0	0
1	0	0	0	1	0	0	0	0	0
e	0	0	0	0	1	0	0	0	0
nH	0	0	0	0	0	1	0	0	0
c	0	0	0	0	0	0	1	0	0
n	0	0	0	0	0	0	0	1	0
1	0	0	0	0	0	0	0	0	1

Fig. 3 Three representations of 4-(chloromethyl)-1H-imidazole. Depiction of a one-hot representation derived from the SMILES of a molecule. Here a reduced vocabulary is shown, while in practice a much larger vocabulary that covers all tokens present in the training data is used.

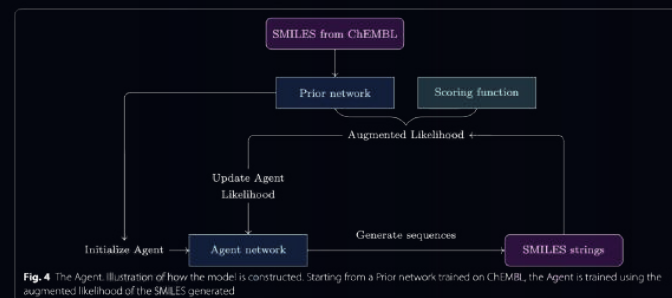
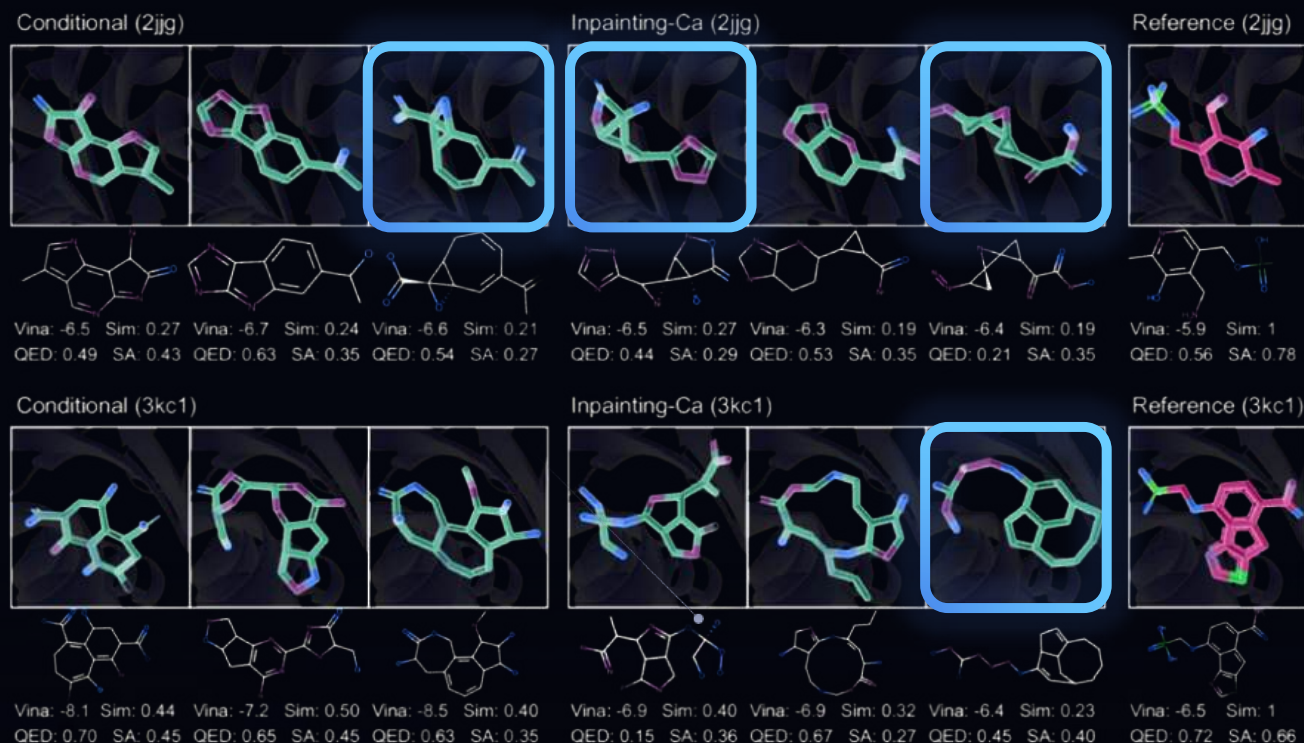


Fig. 4 The Agent. Illustration of how the model is constructed. Starting from a Prior network trained on ChEMBL, the Agent is trained using the augmented likelihood of the SMILES generated.

Molecular de-novo design through deep reinforcement learning
Marcus Olivecrona, Thomas Blaschke, Ola Engkvist and Hongming Chen. J Cheminform (2017)

- SMILES (Simplified Molecular Input Line Entry System) is a depth-first linearized representation of a molecular graph (i.e. converts [cyclic] graphs into a linear sequence of characters expressed in ASCII
- ASCII symbols can then be encoded as in classical LLMs
- Many early RNN-based systems such as REINVENT (AstraZeneca, 2017) which tie SMILES-based LLMs to a learnable reward via RL (REINFORCE) are still competitive today

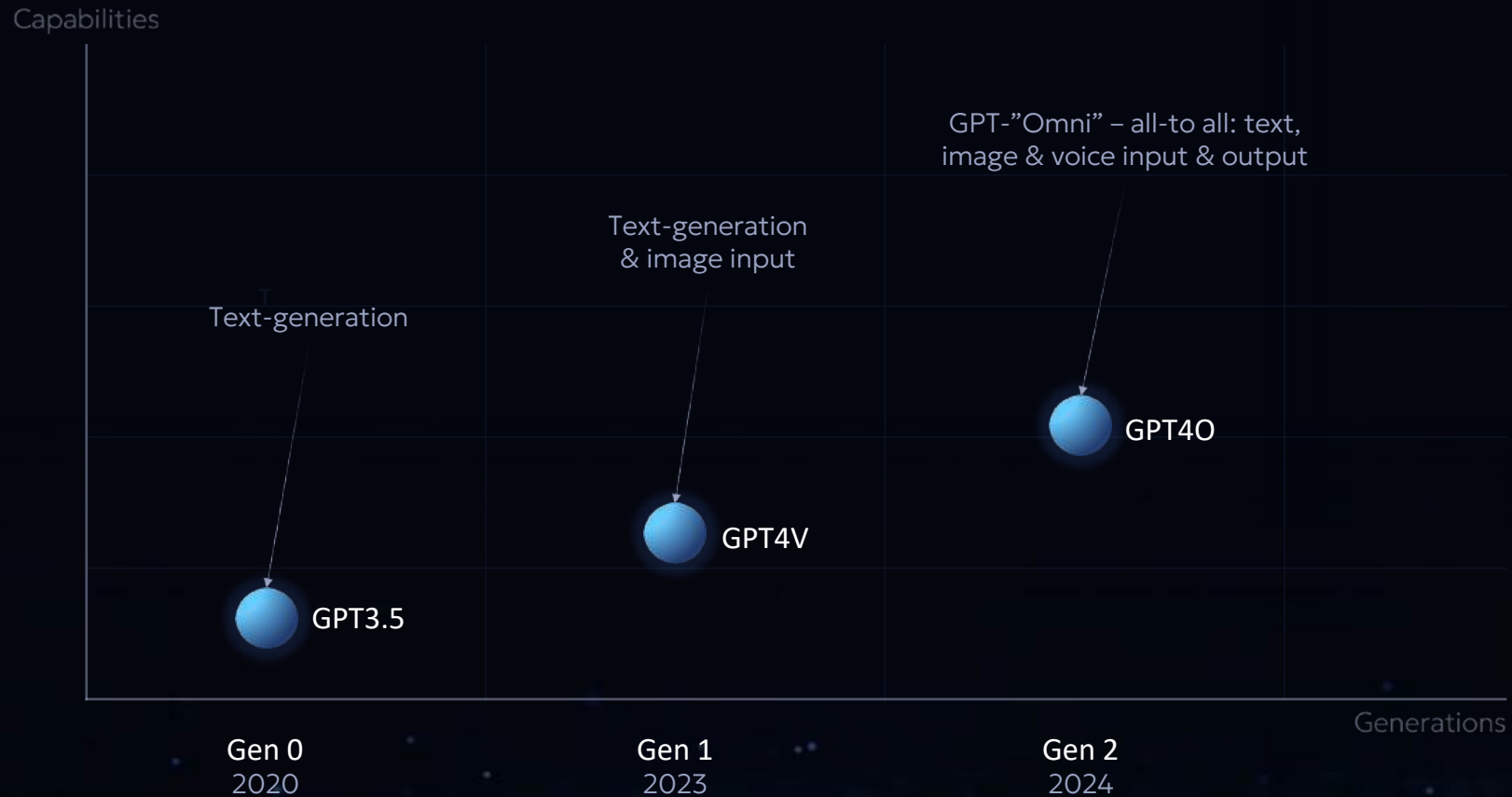
3D-based models still struggle to produce valid molecules



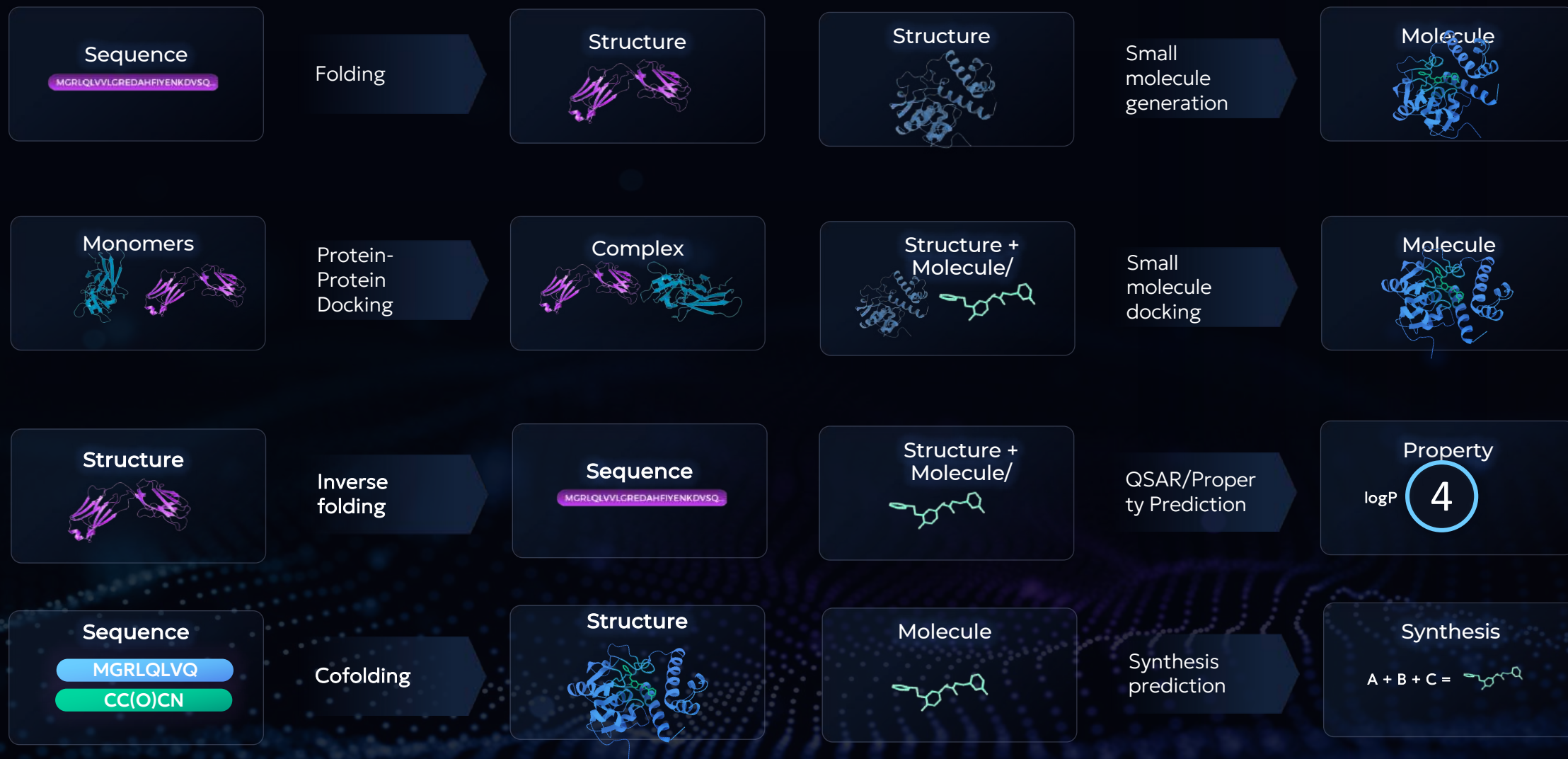
- Strained bond angles
- Highly strained or overextended rings
- Disconnected aromatic bonds
- Large bent overall molecular shapes
- Many methods are still in the “6-Finger phase” of generative models
- Whether same recipe from image world (larger models, more data: scale) will suffice or if explicit inductive biases will be required remains to be seen



The big breakthroughs in AI have been driven through ever more general models



Current Bio AI models, in contrast, still consist of a large number of specialist models



However, this is wasteful and limiting

- Models re-learn same things anew (“grammar” of natural sequence, physics of interactions, generative process, ...)
- Different and small datasets used to re-learn across tasks
- In real life, we often want to “feed” models with whatever information we have
- Currently, this requires models to be chained or “hacked” to include information they are not trained to use (e.g. MSA sampling in AF2), creating and compounding errors

With AF3, all atoms of life (DNA, Protein, Molecules)
could be predicted for the first time

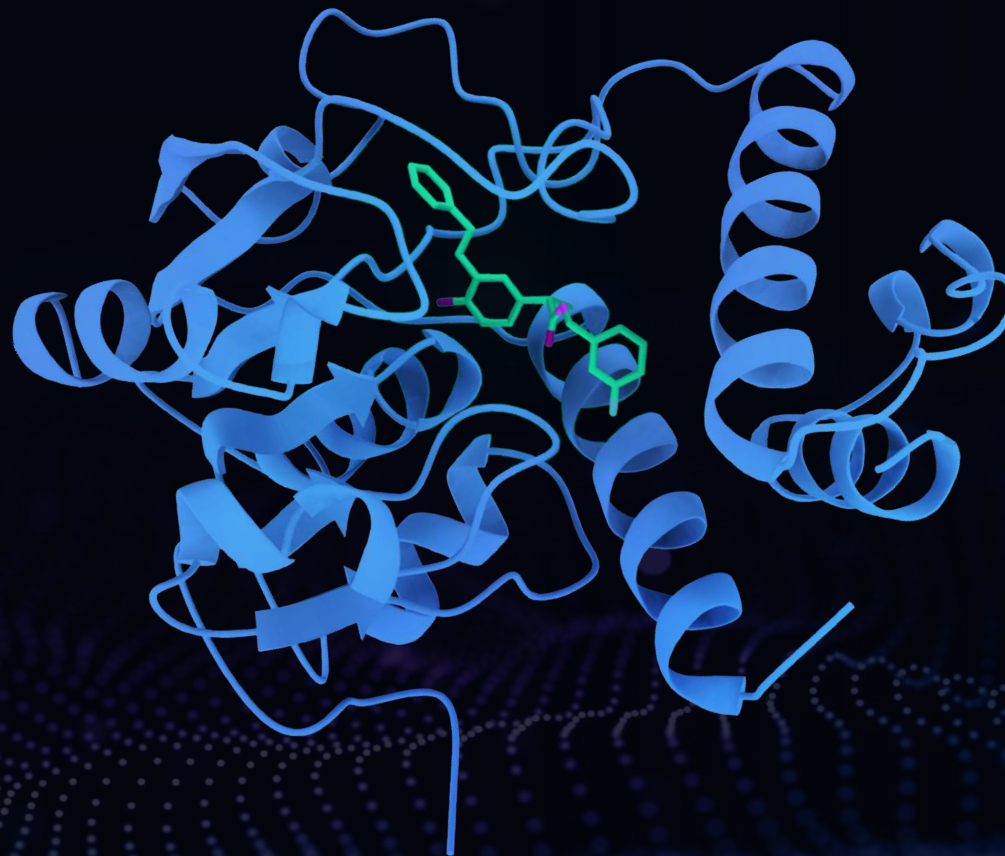


Co-Folding:

Predict atomistic
structure given
sequence of
molecular tokens

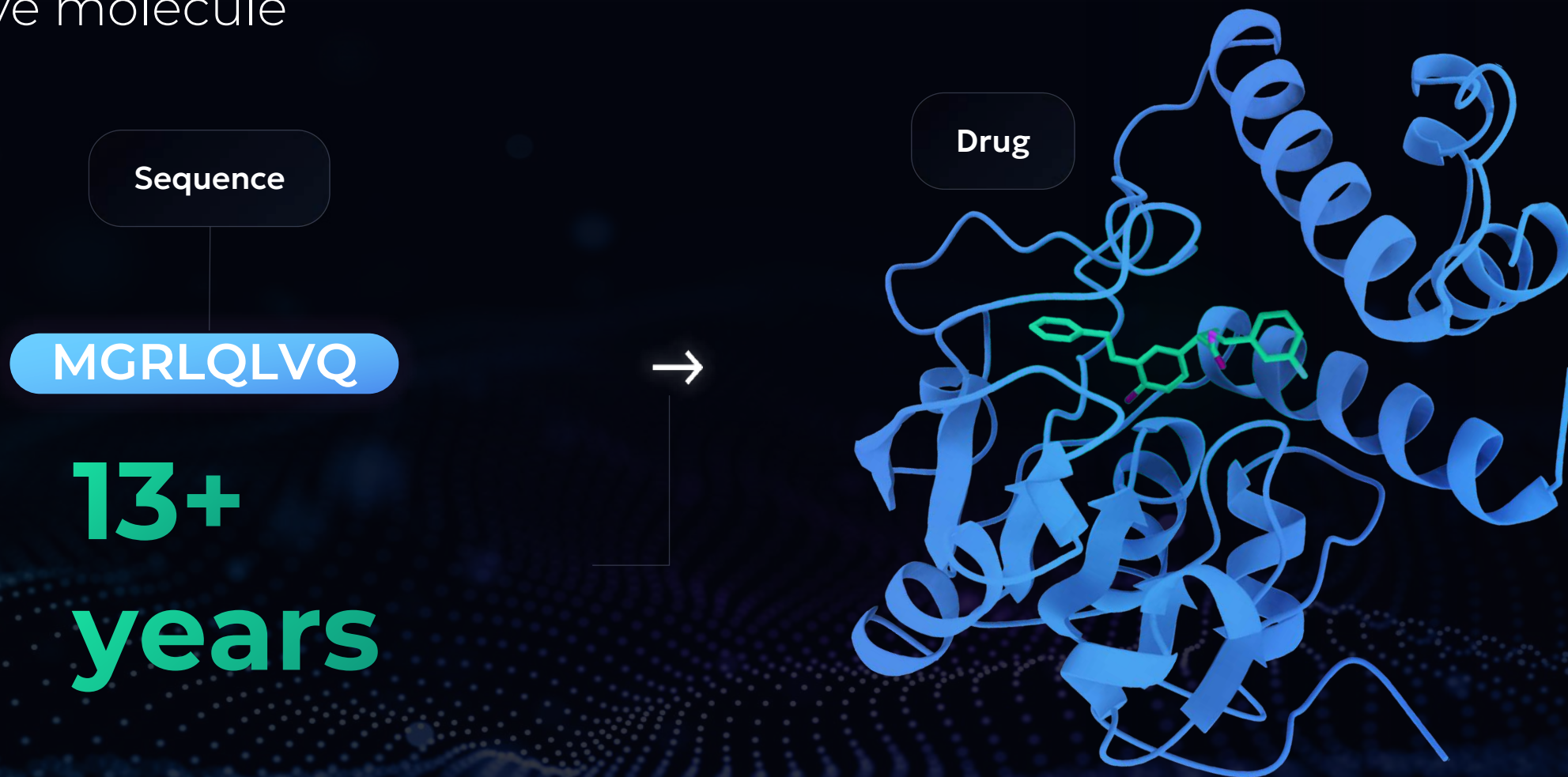
MGRLQLVQ

CC(O)CN



The challenge Neo set out to solve:

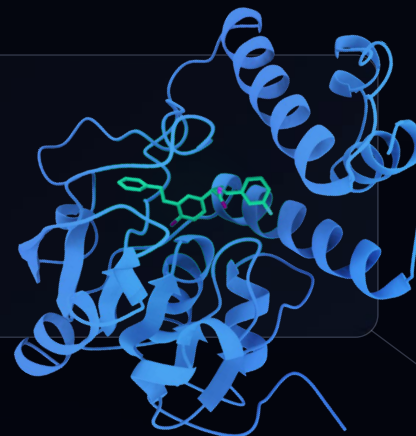
Decade+ process of finding an effective molecule



Challenge: Design & decode medicines

atom-by-atom

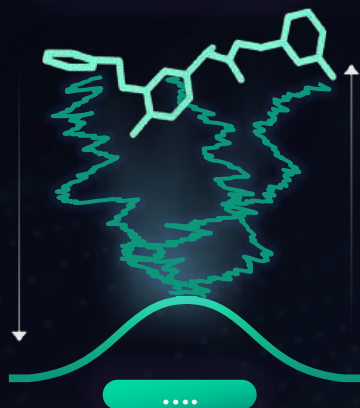
MGRLQLVQ



1. Design

(Molecular generation)

CC(O)CN

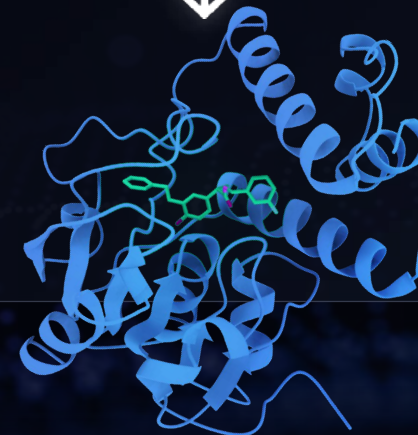


2. Decode

(Folding)

MGRLQLVQ

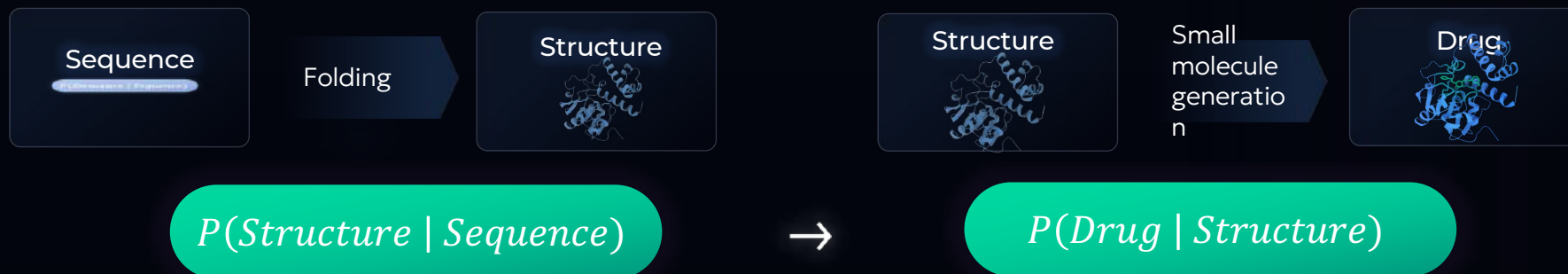
CC(O)CN



And the next frontier was clear



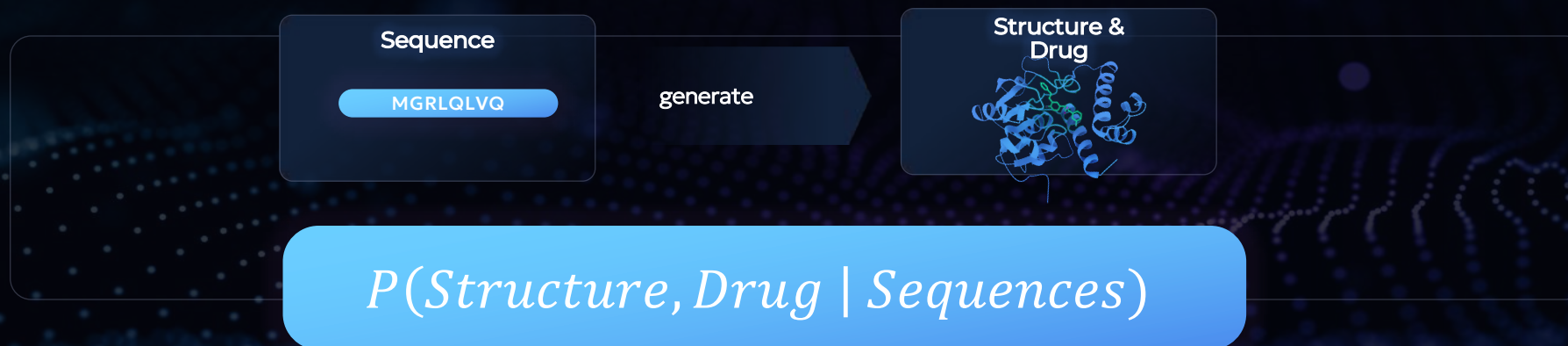
Existing methods



However, in ProMods

$$\text{Structure} = f(\text{Drug}, \text{Sequences})$$

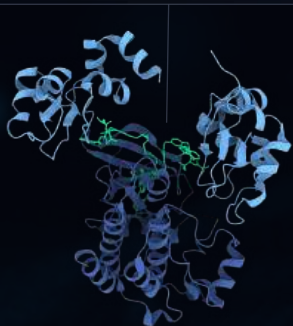
Needed



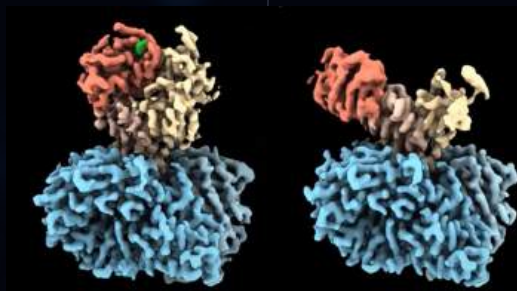
ProMods: molecule defines structure

often **doesn't form** stably **in absence of drug**

Molecule-induced interface

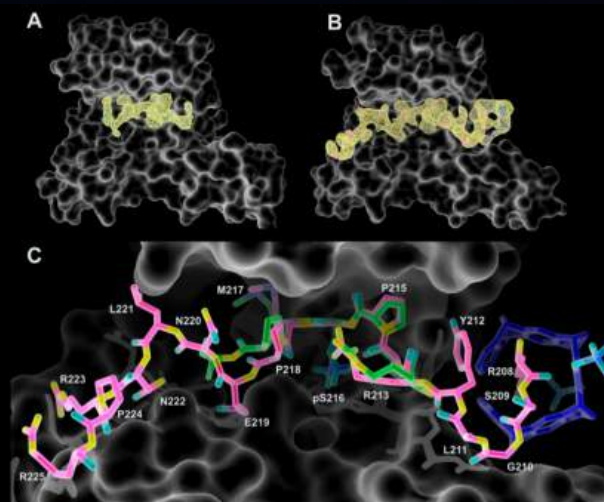


Molecule-induced conformational shift



Pomalidomide-induced conformational shift of CRBN. Gabe Lander, DFCI TPD Series, 2023

Molecule-induced disorder to order conversion



The Molecular Tweezer CLR01 Stabilizes a Disordered Protein-

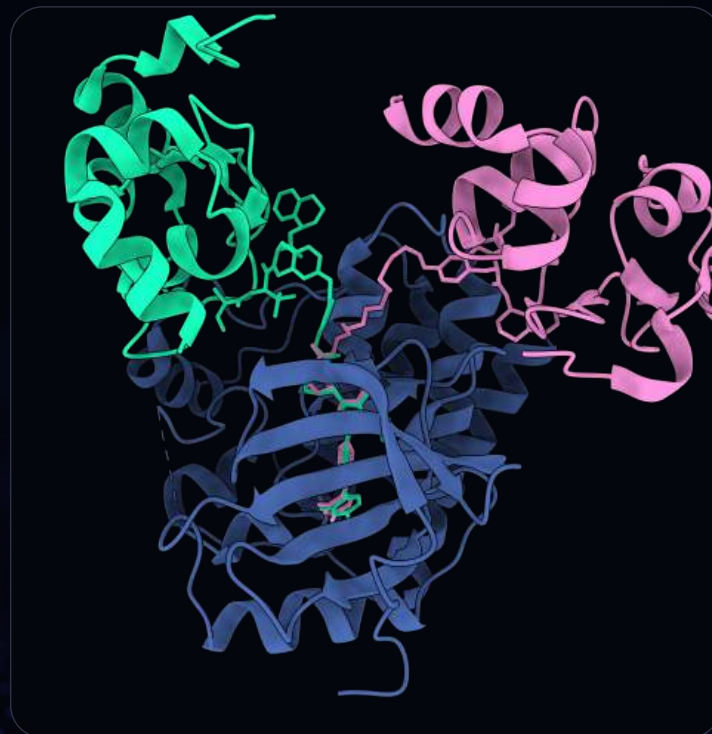
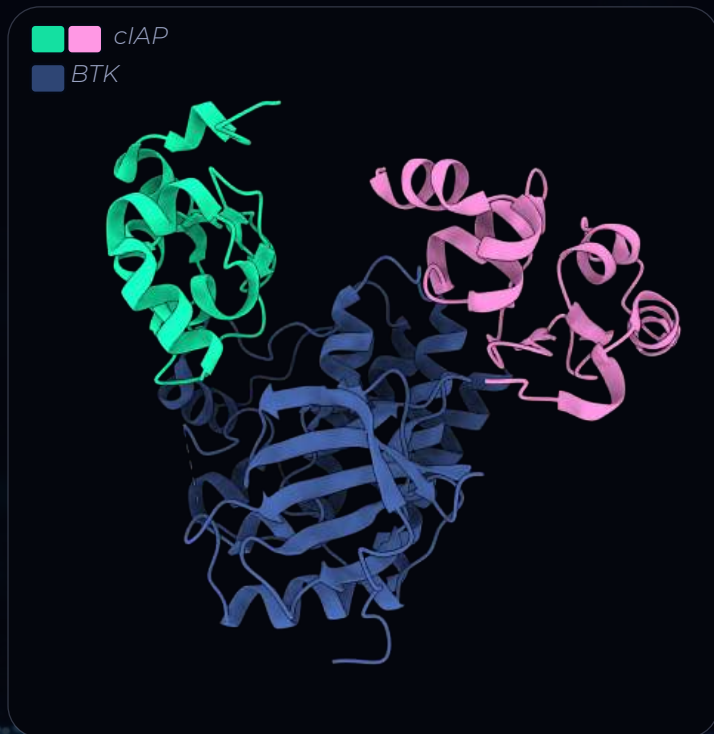
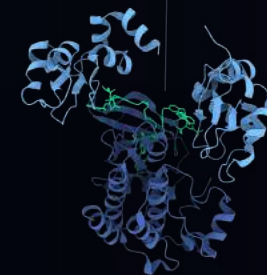
Protein Interface. Bier et al. JACS 2017

Protein-Protein & Protein-Ligand docking invalidly assume interface and monomers exist stably without molecule – simultaneous co-folding & de-novo design required

$$\text{Structure} = f(\text{Drug}, \text{Sequences})$$

Which of these two cIAP-BTK Protein-interfaces is the right one?

Molecule-induced interface

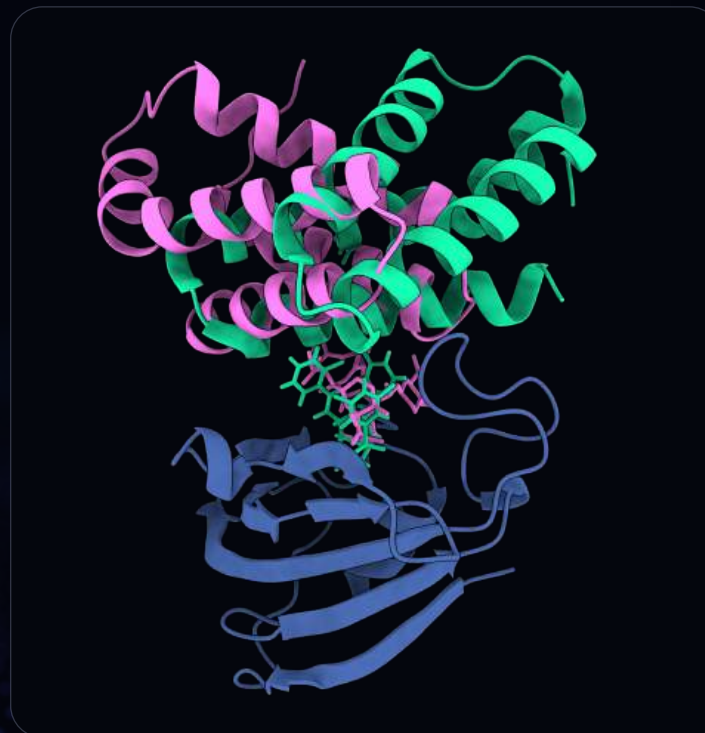
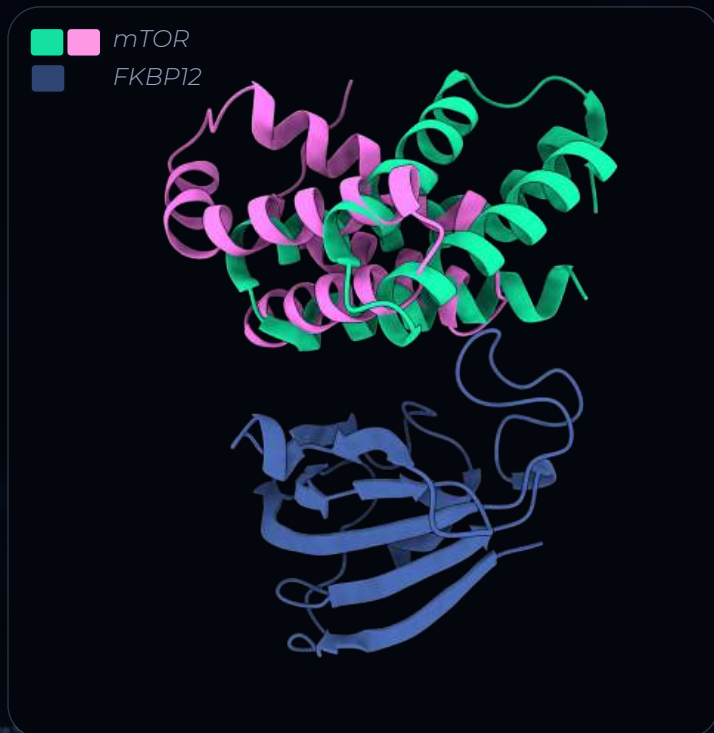
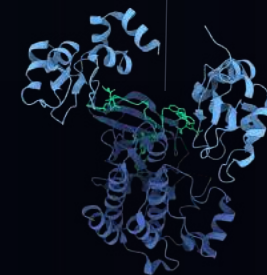


Both! Depending on molecule

6w8i & 6w70

Which of these two FKBP12-mTOR Protein-interfaces is the right one?

Molecule-induced interface

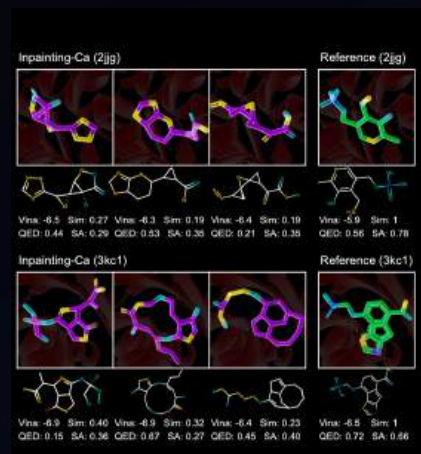


Both! Depending on molecule

1FAP & 8PPZ

De-novo design:

De-novo small molecule design still lag behind:
a) methods assume knowledge of bound-state protein & b) struggle with valid designs



Structure-based Drug Design with Equivariant Diffusion Models. Schneuing, ..., Bronstein, Correia. 2022

Introducing



VANTAI
NEO-1

vant.ai/neo-1

Introducing Neo-1: the worlds most advanced atomistic foundation model





To decode and design
all atoms of life