



HARVARD
MEDICAL SCHOOL



Kempner
INSTITUTE

For the Study of Natural
& Artificial Intelligence
at Harvard University



BROAD
INSTITUTE

AI in Drug Design: Part I **Molecular Drug Discovery**

Marinka Zitnik

Department of Biomedical Informatics, Harvard Medical School

Kempner Institute for the Study of Natural and Artificial Intelligence, Harvard University

Broad Institute of Harvard and MIT

Harvard Data Science

zitniklab.hms.harvard.edu

Biomedical research in the age of AI

Experiments and simulation

Data acquisition and measurements at scale

Accelerated discovery & hypothesis generation



The Economist

Dealing with Europe's hard right
Tim Scott, the non-victim
Zelensky on the long war
The hunt for green metals

SEPTEMBER 16TH - 22ND 2023

HOW AI CAN REVOLUTIONISE SCIENCE



Review

Scientific discovery in the age of artificial intelligence

<https://doi.org/10.1038/s41586-023-06221-2>

Received: 30 March 2023

Accepted: 16 May 2023

Published online: 2 August 2023

Check for updates

Hanchen Wang^{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000}

Artificial intelligence (AI) is being increasingly integrated into scientific discovery to augment and accelerate research, helping scientists to generate hypotheses, design experiments, collect and interpret large datasets, and gain insights that might not have been possible using traditional scientific methods alone. Here we examine breakthroughs over the past decade that include self-supervised learning, which allows models to be trained on vast amounts of unlabelled data, and geometric deep learning, which leverages knowledge about the structure of scientific data to enhance model accuracy and efficiency. Generative AI methods can create designs, such as small-molecule drugs and proteins, by analysing diverse data modalities, including images and sequences. We discuss how these methods can help scientists throughout the scientific process and the central issues that remain despite such advances. Both developers and users of AI tools need a better understanding of when such approaches need improvement, and challenges posed by poor data quality and stewardship remain. These issues cut across scientific disciplines and require developing foundational algorithmic approaches that can contribute to scientific understanding or acquire it autonomously, making them critical areas of focus for AI innovation.

The foundation for forming scientific insights and theories is laid by how data are collected, transformed and understood. The rise of deep learning in the early 2010s has significantly expanded the scope and ambition of these scientific discovery processes. Artificial intelligence (AI) is increasingly used across scientific disciplines to integrate massive datasets, refine measurements, guide experimentation, explore the space of theories compatible with the data, and provide actionable and reliable models integrated with scientific workflows for autonomous discovery.

Data collection and analysis are fundamental to scientific understanding and discovery, two of the central aims in science¹, and quantitative

methods and emerging technologies, from physical instruments such as microscopes to research techniques such as bootstrapping, have long been used to reach these aims². The introduction of digitization in the 1950s paved the way for the general use of computing in scientific research. The rise of data science since the 2010s has enabled AI to provide valuable guidance by identifying scientifically relevant patterns from large datasets.

Although scientific practices and procedures vary across stages of scientific research, the development of AI algorithms cuts across traditionally isolated disciplines (Fig. 1). Such algorithms can enhance the design and execution of scientific studies. They are becoming

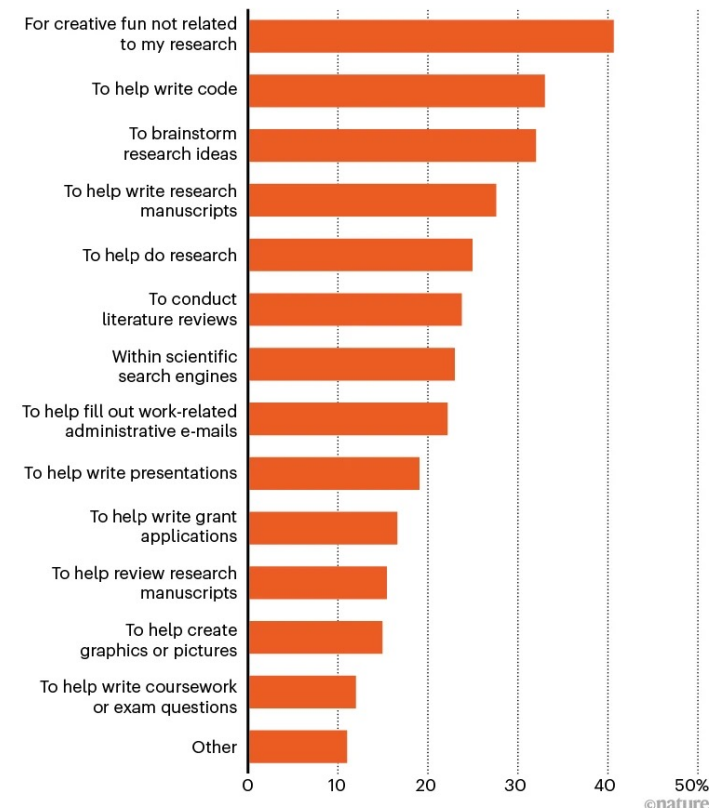
¹Department of Engineering, University of Cambridge, Cambridge, UK; ²Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA, USA; ³Department of Computational Science and Engineering, Georgia Institute of Technology, Atlanta, GA, USA; ⁴Department of Computer Science, Cornell University, Ithaca, NY, USA; ⁵Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA; ⁶Department of Computer Science, Stanford University, Stanford, CA, USA; ⁷Department of Physics, Massachusetts Institute of Technology, Cambridge, MA, USA; ⁸Department of Biology, University of California, Berkeley, Berkeley, CA, USA; ⁹Department of Mathematics, University of Toronto, Toronto, ON, Canada; ¹⁰Department of Earth, Environmental and Planetary Sciences, Brown University, Providence, RI, USA; ¹¹Data Science Institute, Brown University, Providence, RI, USA; ¹²MIT, Cambridge, MA, USA; ¹³Department of Computer Science, Princeton University, Princeton, NJ, USA; ¹⁴Department of Physics, Carnegie Mellon University, Pittsburgh, PA, USA; ¹⁵Department of Physics and Center for Data Science, New York University, New York, NY, USA; ¹⁶Google DeepMind, London, UK; ¹⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ¹⁸Department of Computer Science, University of Washington, Seattle, WA, USA; ¹⁹Department of Computer Science, University of Texas at Austin, Austin, TX, USA; ²⁰Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, USA; ²¹Department of Computer Science, University of Michigan, Ann Arbor, MI, USA; ²²Department of Computer Science, University of California, San Diego, San Diego, CA, USA; ²³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ²⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ²⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ²⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ²⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ²⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ²⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ³⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁴⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁵⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁶⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁷⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁸⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹¹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹²Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹³Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁴Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁵Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁶Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁷Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁸Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ⁹⁹Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA; ¹⁰⁰Department of Computer Science, University of Wisconsin-Madison, Madison, WI, USA.

Nature | Vol 620 | 3 August 2023 | 47

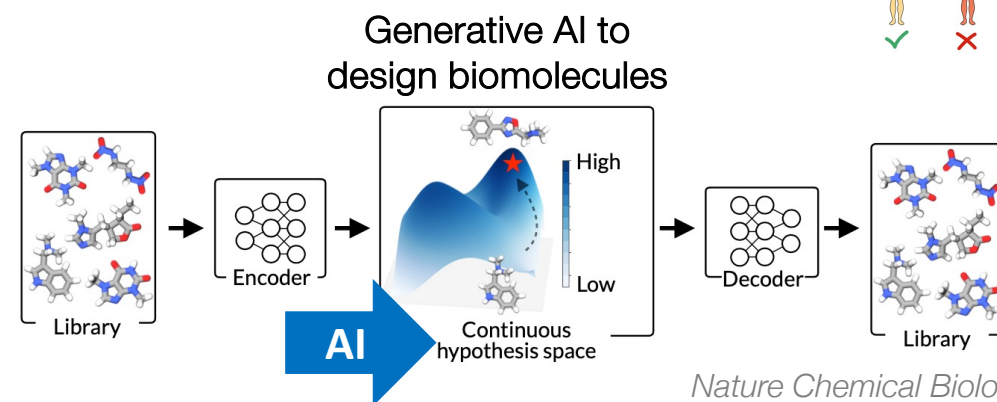
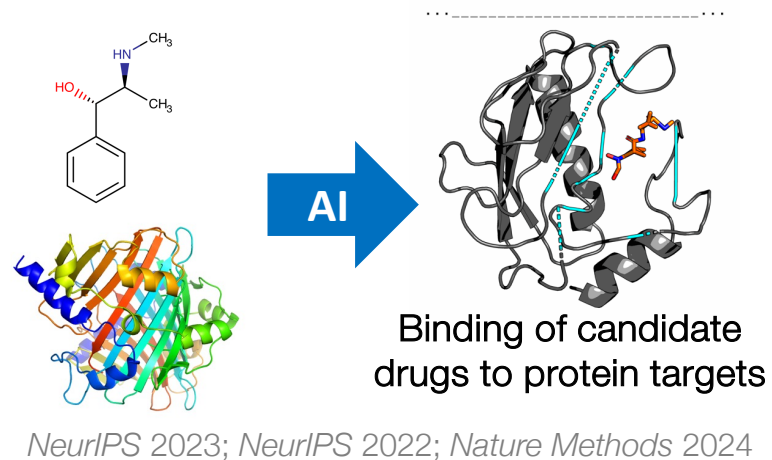
Biomedical research in the age of AI

HOW RESEARCHERS USE LARGE LANGUAGE MODELS

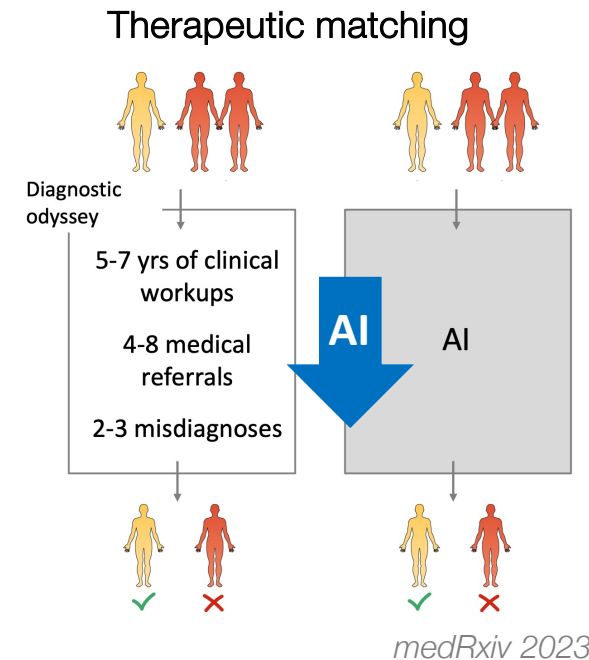
Q: What do you use generative AI tools (such as ChatGPT and other large language models) for? (Choose all that apply.)



Generative AI is changing the way science is done



AI is used to augment research, providing insights that might not have been possible using traditional methods alone

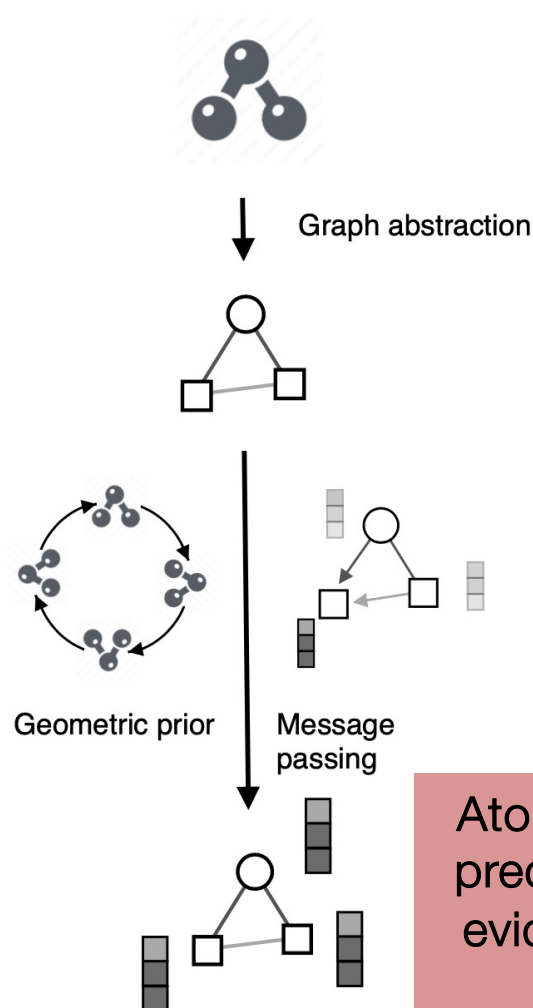


Potential for AI in drug design

Treatment	Organization	Description	Phase	Lead indication
REC-2282	Recursion	Small molecule pan-HDAC inhibitor	2/3	Neurofibromatosis type 2
REC-994	Recursion	Small molecule superoxide scavenger	2	Cerebral cavernous malformation
REC-4881	Recursion	Small molecule inhibitor of MEK1 and MEK2	2	Familial adenomatous polyposis
INS018_055	InSilico Medicine	Small molecule inhibitor	2	Idiopathic pulmonary fibrosis
BEN-2293	BenevolentAI	Topical pan-tyrosine kinase inhibitor	2a	Atopic dermatitis
EXS-21546	Exscientia and Evotec	A _{2A} receptor antagonist	1b/2	Solid tumors carrying high adenosine signatures.
RLY-4008	Relay Therapeutics	Inhibitor of FGFR2	1/2	FGFR2-altered cholangiocarcinoma
EXS-4318	Exscientia	PKC- θ inhibitor	1/2	Inflammatory and autoimmune conditions
BEN-8744	BenevolentAI	Small molecule PDE10 inhibitor	1	Ulcerative colitis
Undisclosed	Recursion	Small molecular inhibitor of RBM39, a CDK12-associated protein	Pre-clinical	HRD-negative ovarian cancer

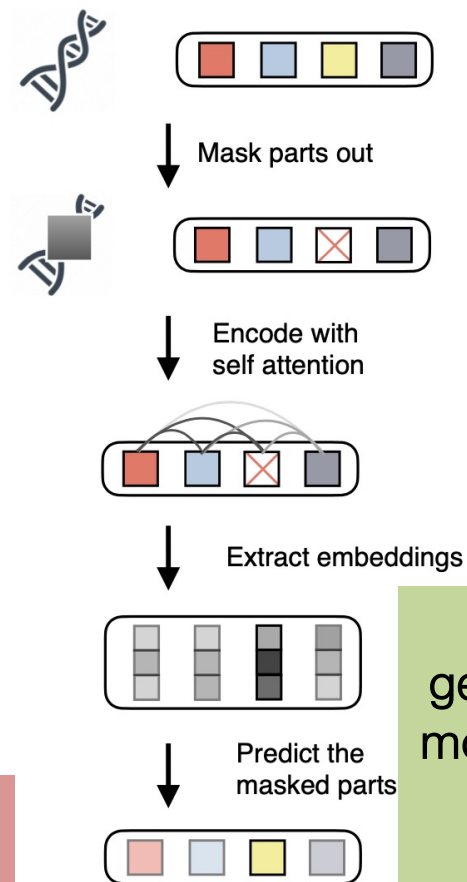
AI in drug design: Key innovations

Geometric learning



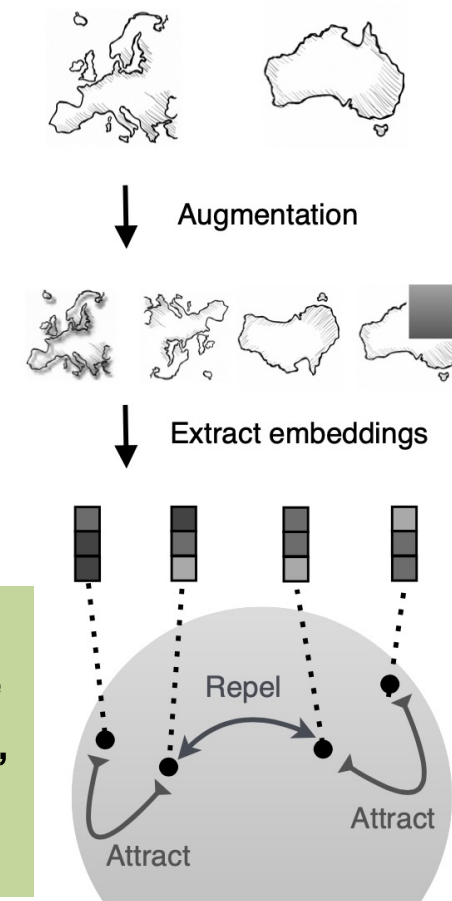
Atom-level molecular prediction, real-world evidence knowledge graphs

Self-supervised learning



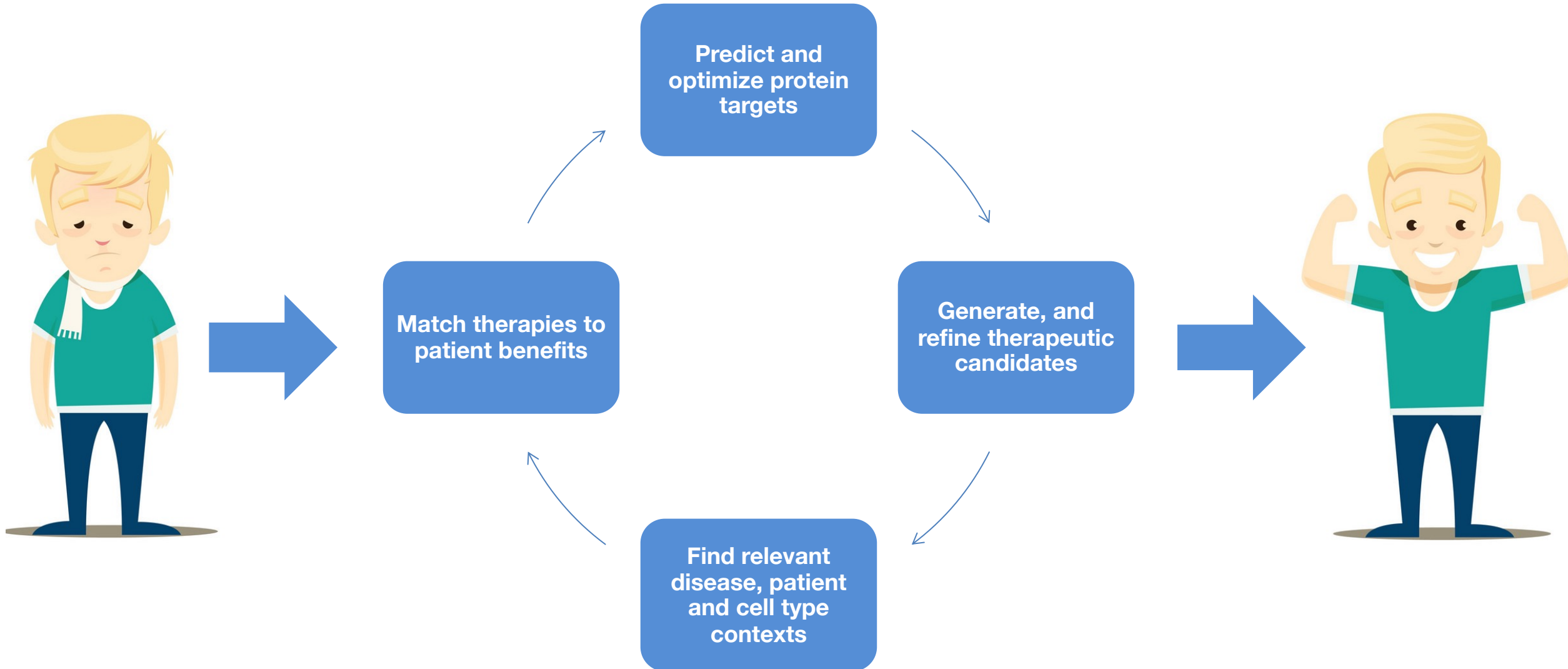
Target discovery, genotype-phenotype modeling, DNA, RNA, AA sequence modeling

Generative AI

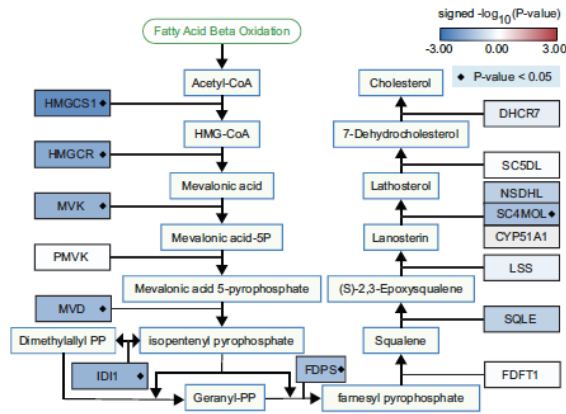


chatbots, copilots, molecular and drug design

Our vision: Lay the foundations for AI to enhance the understanding of medicine and drug design, eventually enabling AI to learn on its own

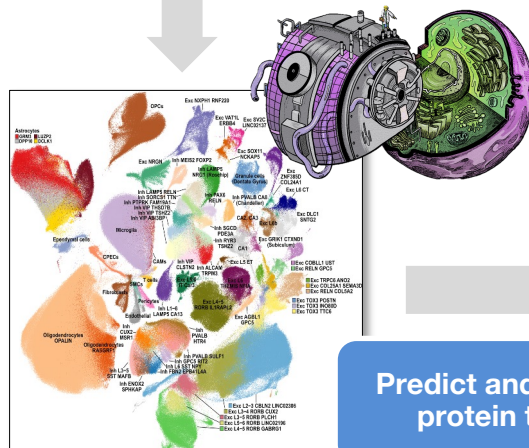


Disease circuitry across diseases and individuals



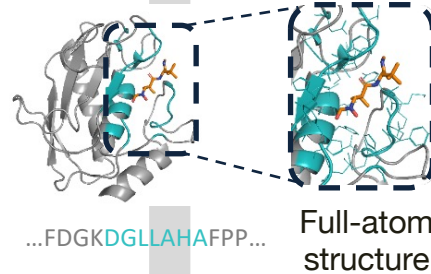
Key disease mechanisms, shared effects, interaction effects

Multimodal representation learning models
PINNACLE, PDGrapher

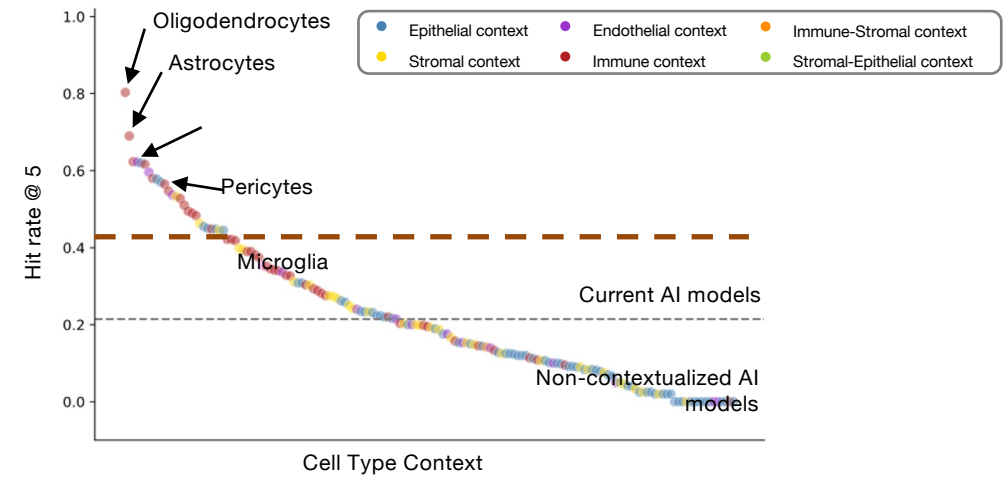


Predict and optimize protein targets

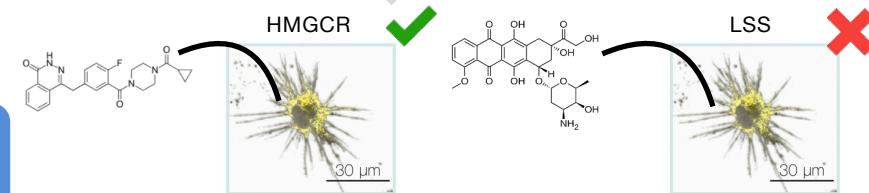
Generative geometric deep learning models
PocketGen, FAIR, and TxPLM



Generate, and refine therapeutic candidates



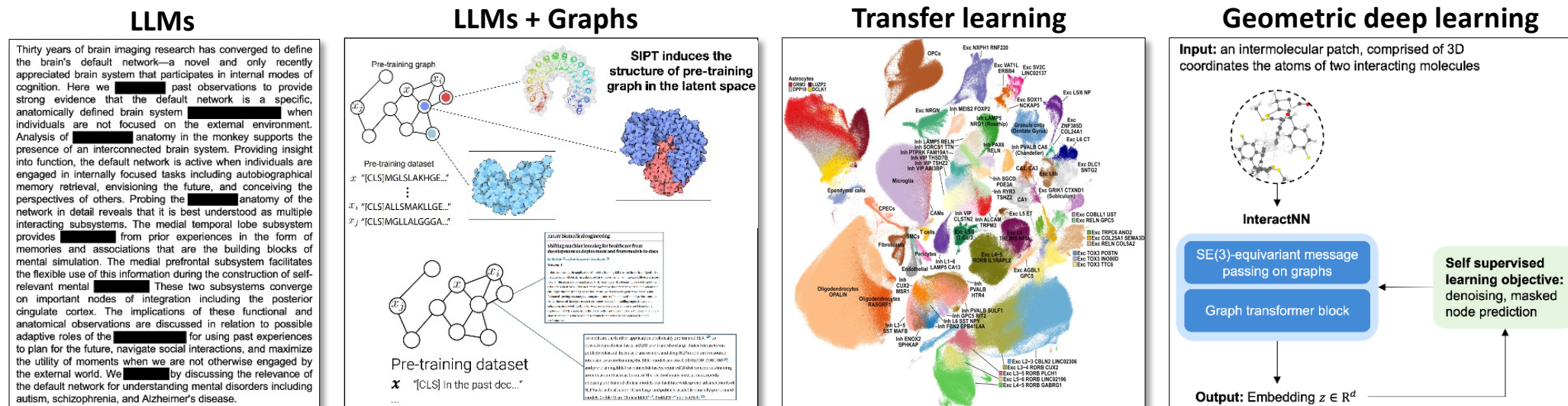
Rank-ordered lists of molecules to modulate disease circuitry in a cell-type specific manner



Chemistry optimization, synthesis, automation

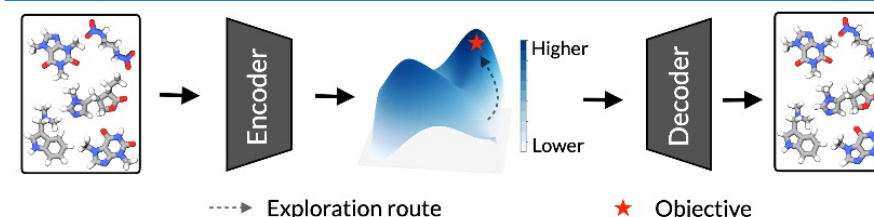


Methods: Geometric deep learning, LLMs + graphs, transfer learning



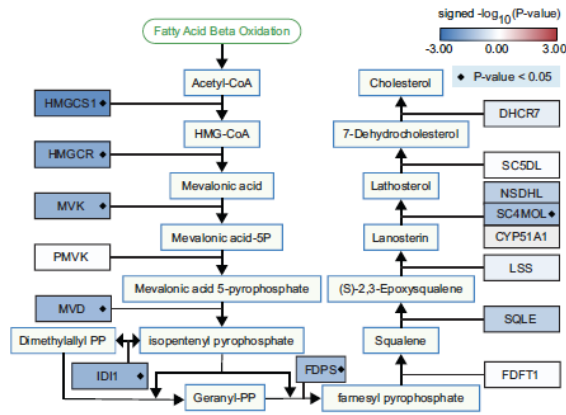
- **Multimodal LLMs** to leverage multimodal biological sequences and scientific knowledge
- **Knowledge-graph based models** that train self-supervised models on broad data at scale without pre-defined labels
- **Geometric and generative AI models** that create action plans for experiments and produce new designs such as small molecule drugs and proteins from experimental data

Generative AI: Better drug design, molecule optimization, guiding high-throughput perturbation & interaction screening



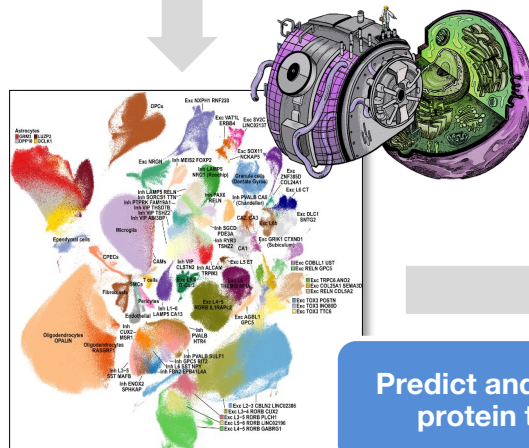
Li et al., *Nature Methods* 2024; Extetaie et al., *Nature Mach Intel* 2023; McDermott et al., *Nature Mach Intel* 2023; Li et al., *Nature Biomed Eng* 2022; Zhang et al., *ICLR* 2022; Zhang et al., *NeurIPS* 2022; Agarwal et al., *NeurIPS* 2022; Huang et al., *NeurIPS* 2021; Alsentzer et al., *NeurIPS* 2020; Zhang et al., *NeurIPS* 2020; Huang et al., *NeurIPS* 2020; Queen et al., *NeurIPS* 2023; He et al., *ICML* 2023; Jiali et al., *ICLR* 2023; Zhang et al., *NeurIPS* 2023; Scott et al., *Nature Mach Intell* 2023; Sanders et al., *Nature Mach Intell* 2023; Agarwal et al., *AISTATS* 2022; Zhong et al., *CVPR* 2024

Disease circuitry across diseases and individuals



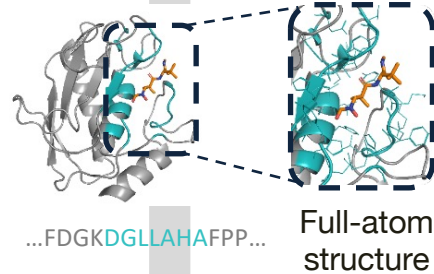
Key disease mechanisms, shared effects, interaction effects

Multimodal representation learning models
PINNACLE, PDGrapher

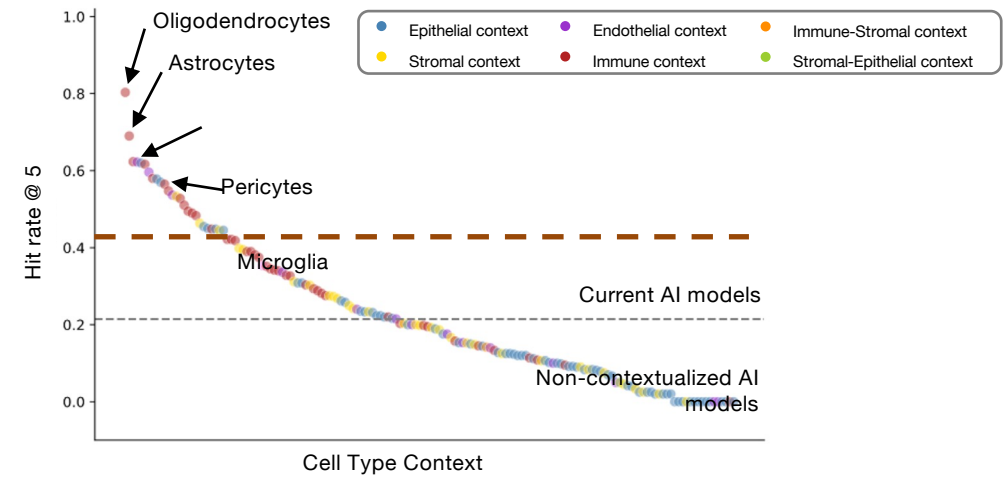


Predict and optimize protein targets

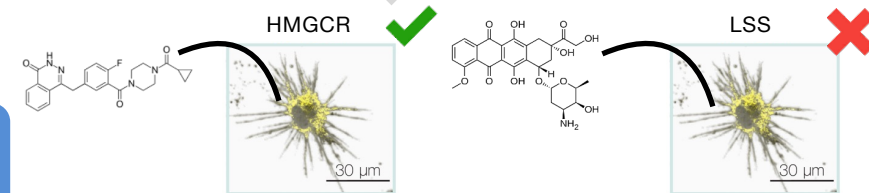
Generative geometric deep learning models
PocketGen, FAIR, and TxPLM



Generate, and refine therapeutic candidates



Rank-ordered lists of molecules to modulate disease circuitry in a cell-type specific manner



Chemistry optimization, synthesis, automation



“apple” is a **polysemic** word...



🔍 grow an apple

🔍 buy an apple|

... whose **particular meaning** is resolved via **sentence context**



🔍 grow an apple

🔍 grow an apple **tree**

🔍 grow an apple **tree from seed**

🔍 grow an apple **tree in a pot**

🔍 grow an apple **tree indoors**



🔍 buy an apple|

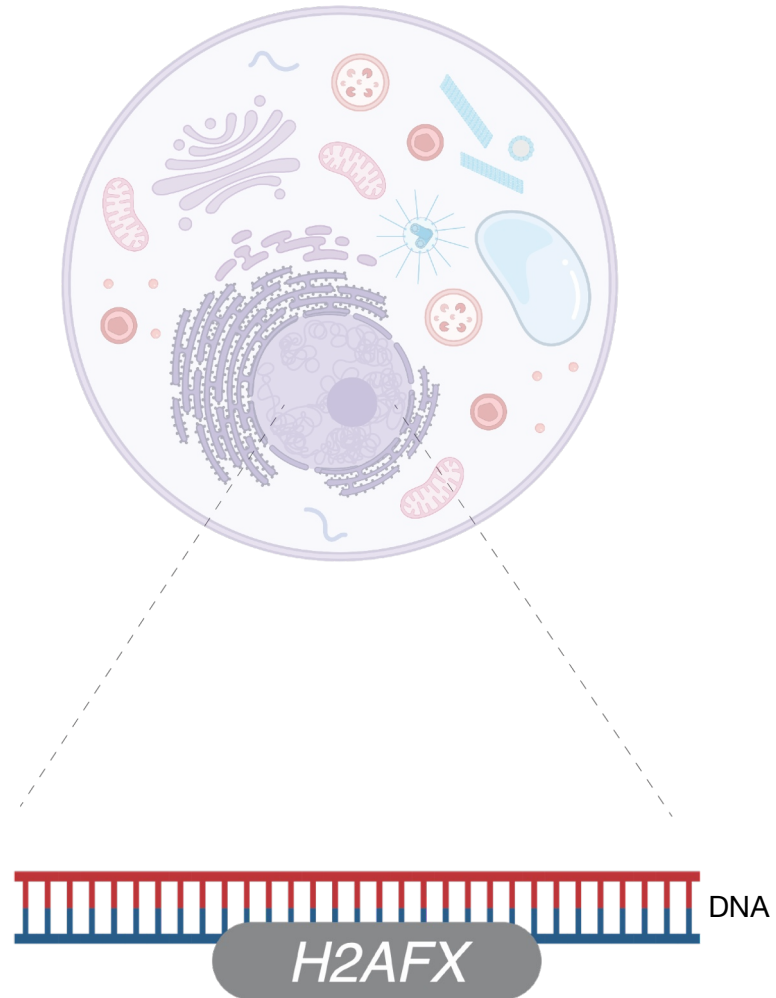
🔍 buy an apple **watch**

🔍 buy an apple **gift card**

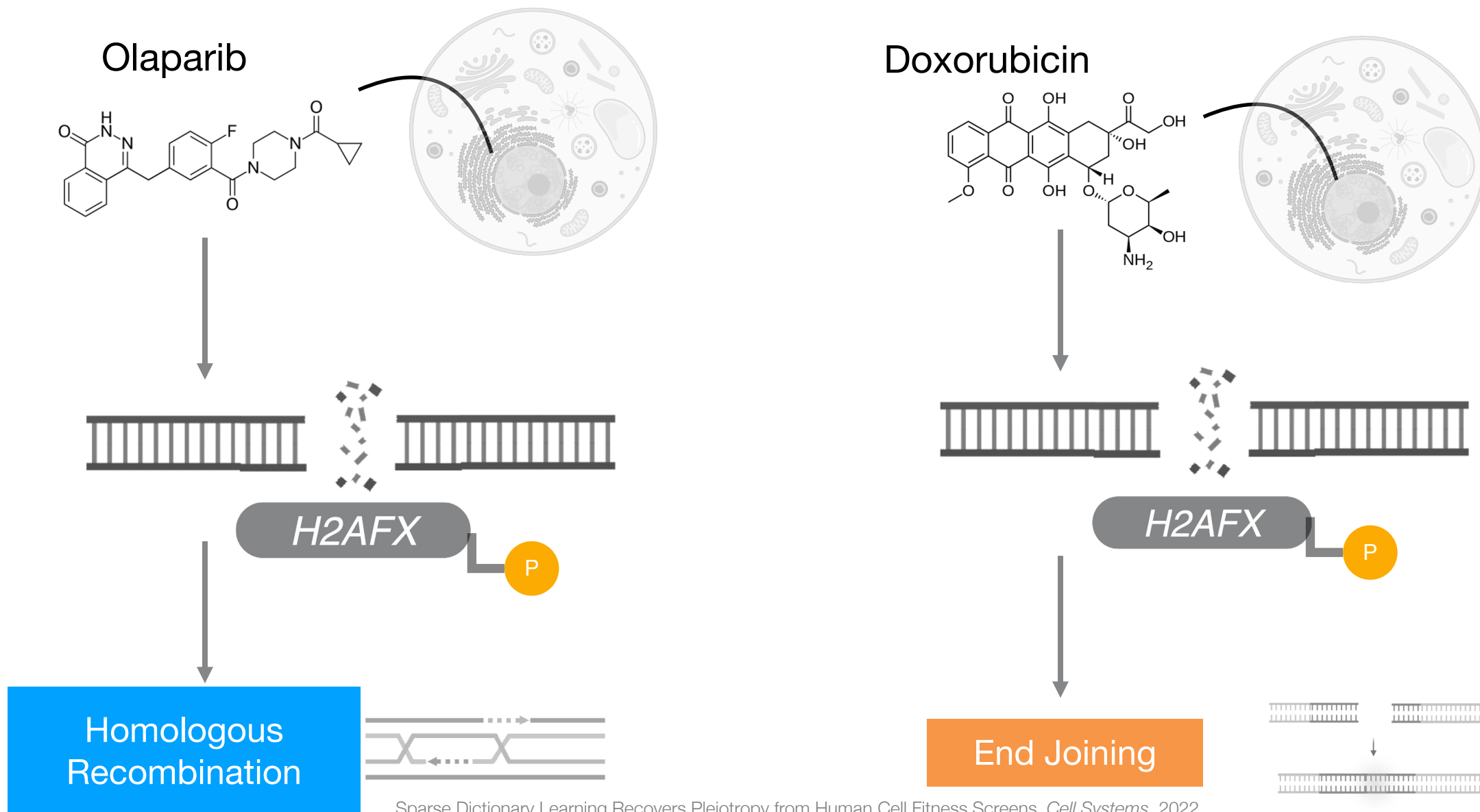
🔍 buy an apple **tv**



H2AFX is a **pleiotropic** gene...

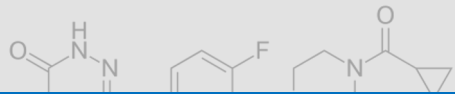


... whose **particular function** is resolved via **cell context**

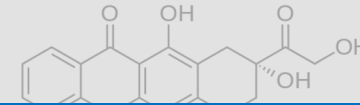


... whose **particular function** is resolved via **cell context**

Olaparib



Doxorubicin

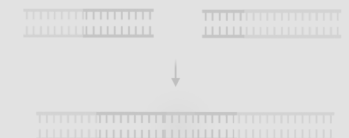


Can we develop models that dynamically
adjust their outputs to biological
contexts in which they operate?

Homologous
Recombination



End Joining

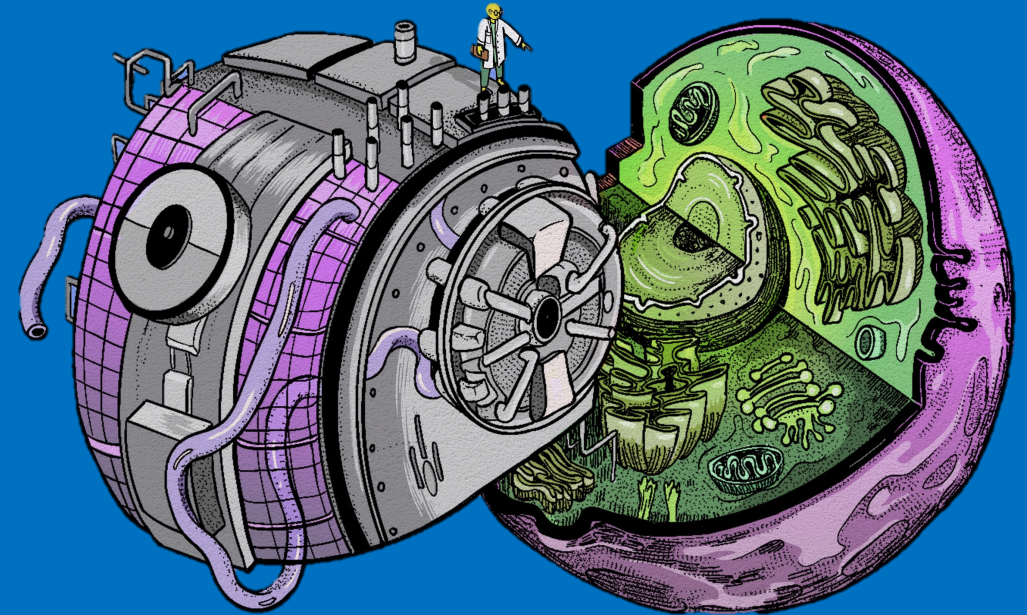


Sparse Dictionary Learning Recovers Pleiotropy from Human Cell Fitness Screens, *Cell Systems*, 2022

Contextualizing Protein Representations Using Deep Learning on Protein Networks and Single-Cell Data, *Nature Methods*, 2024 (in press)

PINNACLE AI

Precise and cell-type specific protein representation learning



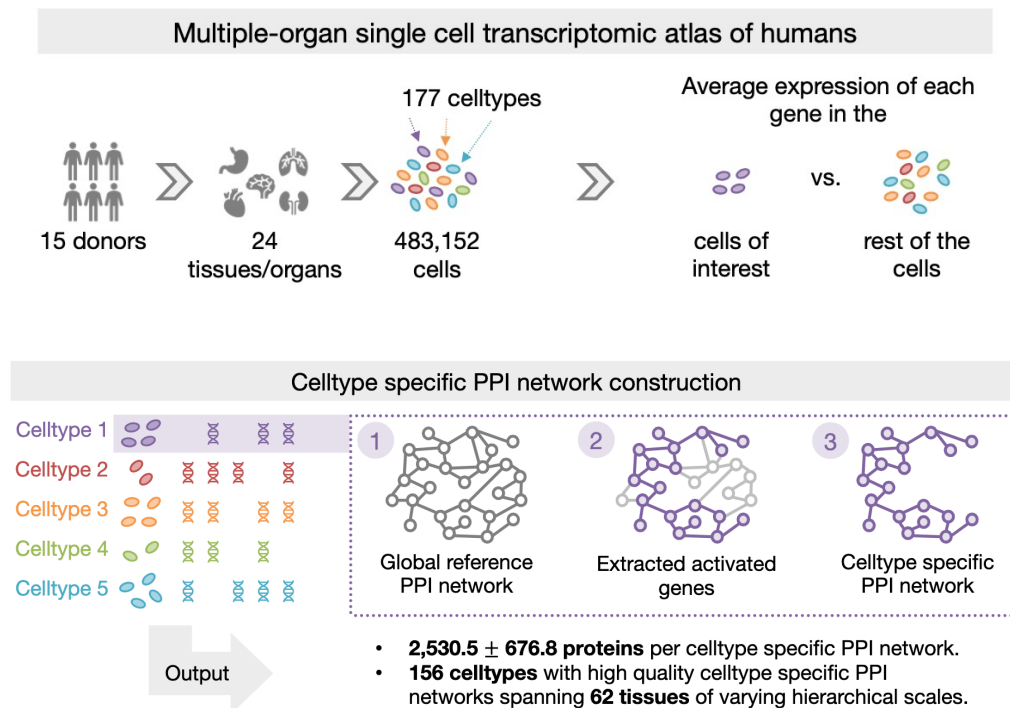
Providing outputs tailored to biological contexts is essential for broad use of foundation models in biology

PINNACLE models support a broad array of tasks:

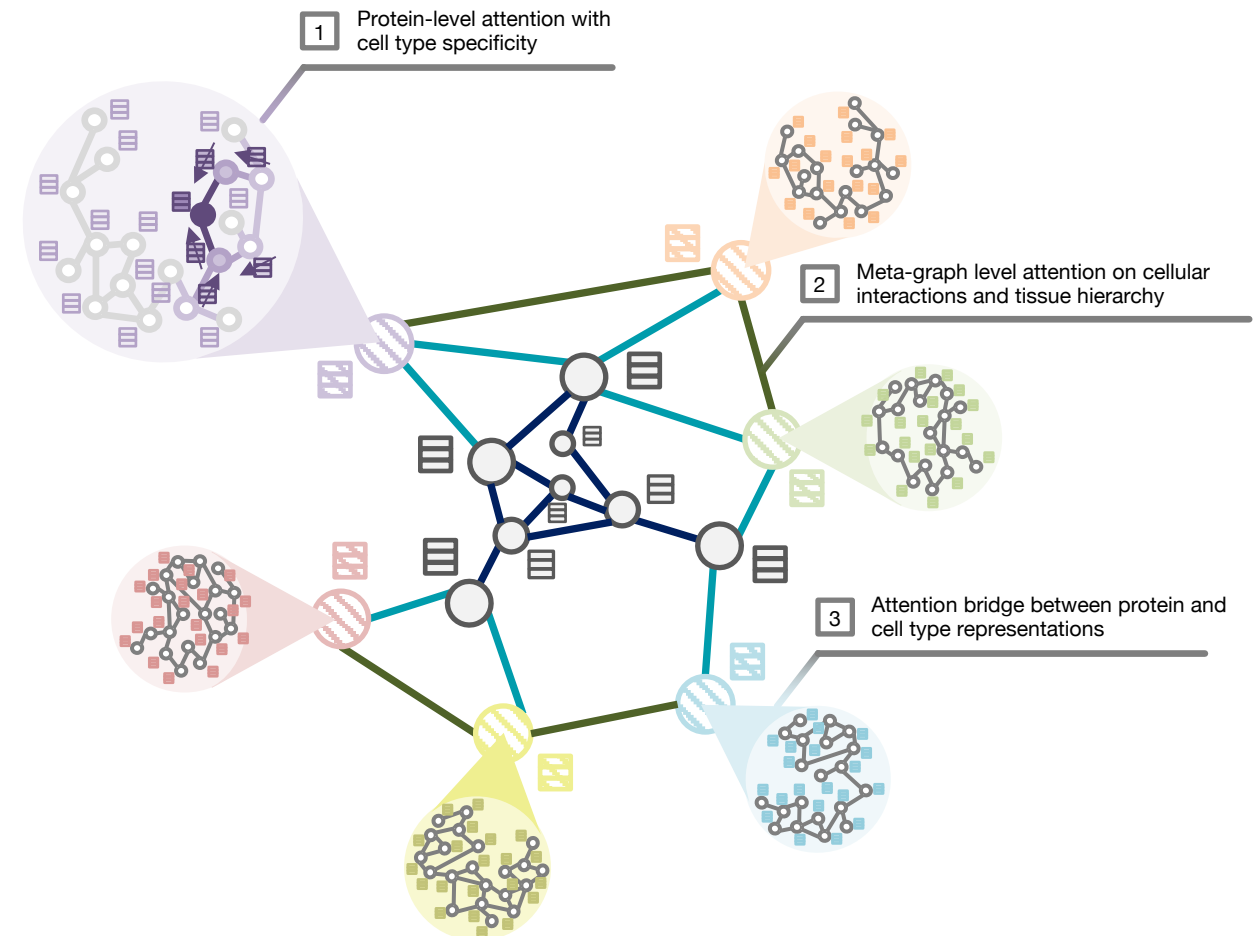
- 🧪 Enhance 3D structural protein representations
- 💊 Study effects of drugs across cell-type contexts
- 🎯 Nominate therapeutic targets in cell-type specific manner
- 🌲 Zero-shot retrieval of tissue hierarchy

PINNACLE: Geometric deep learning for precise protein representation learning and cell type- and state-specific prediction

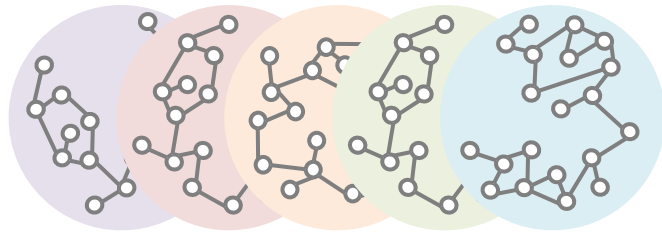
Data: Protein networks and single-cell transcriptomic data



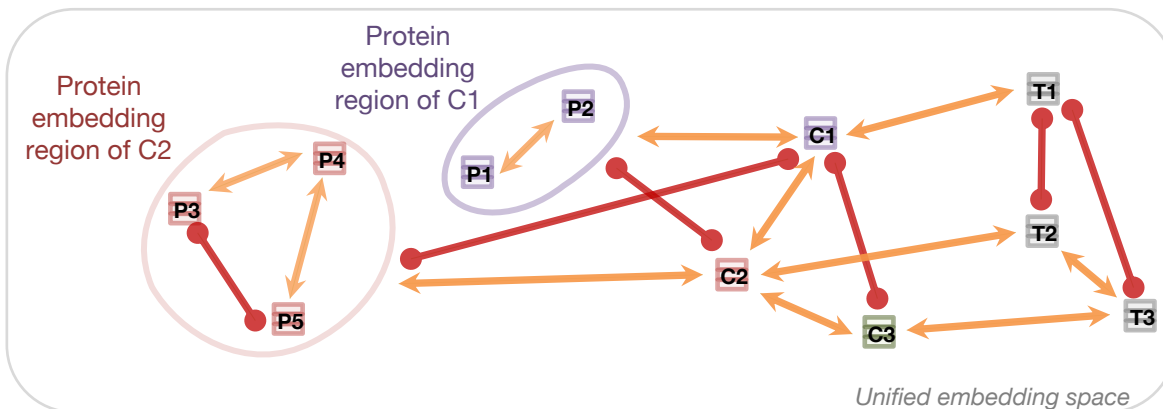
Model: Self-supervised contextual GNN



PINNACLE: Building our intuition



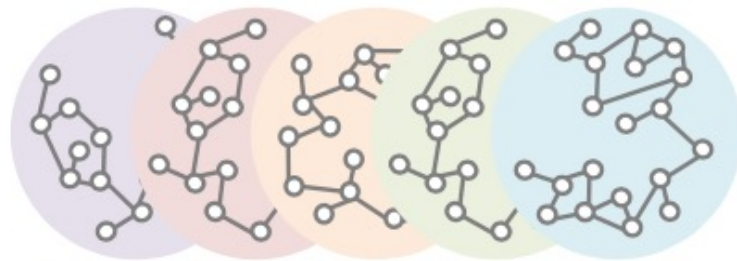
Protein networks across 156 cell type contexts spanning 62 tissues of varying hierarchical scales



Unified embedding space

Self-supervised learning to learn a general encoder of protein embeddings tailored to cellular contexts

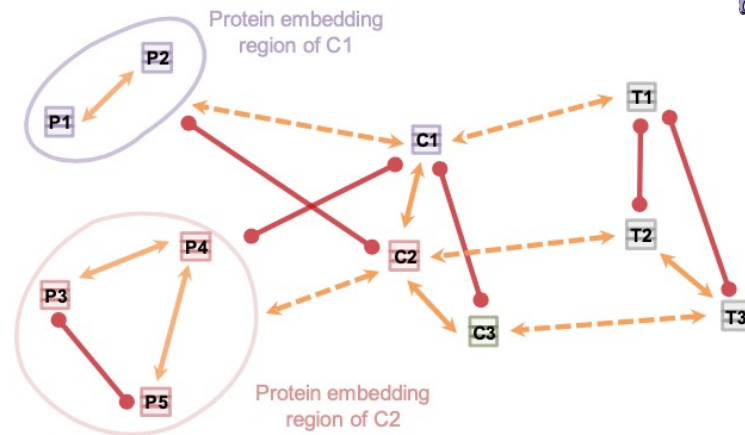
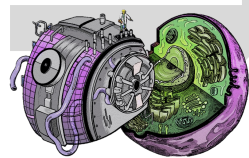
PINNACLE produces context-specific protein representations that are tailored to biological contexts, cell types and cell states



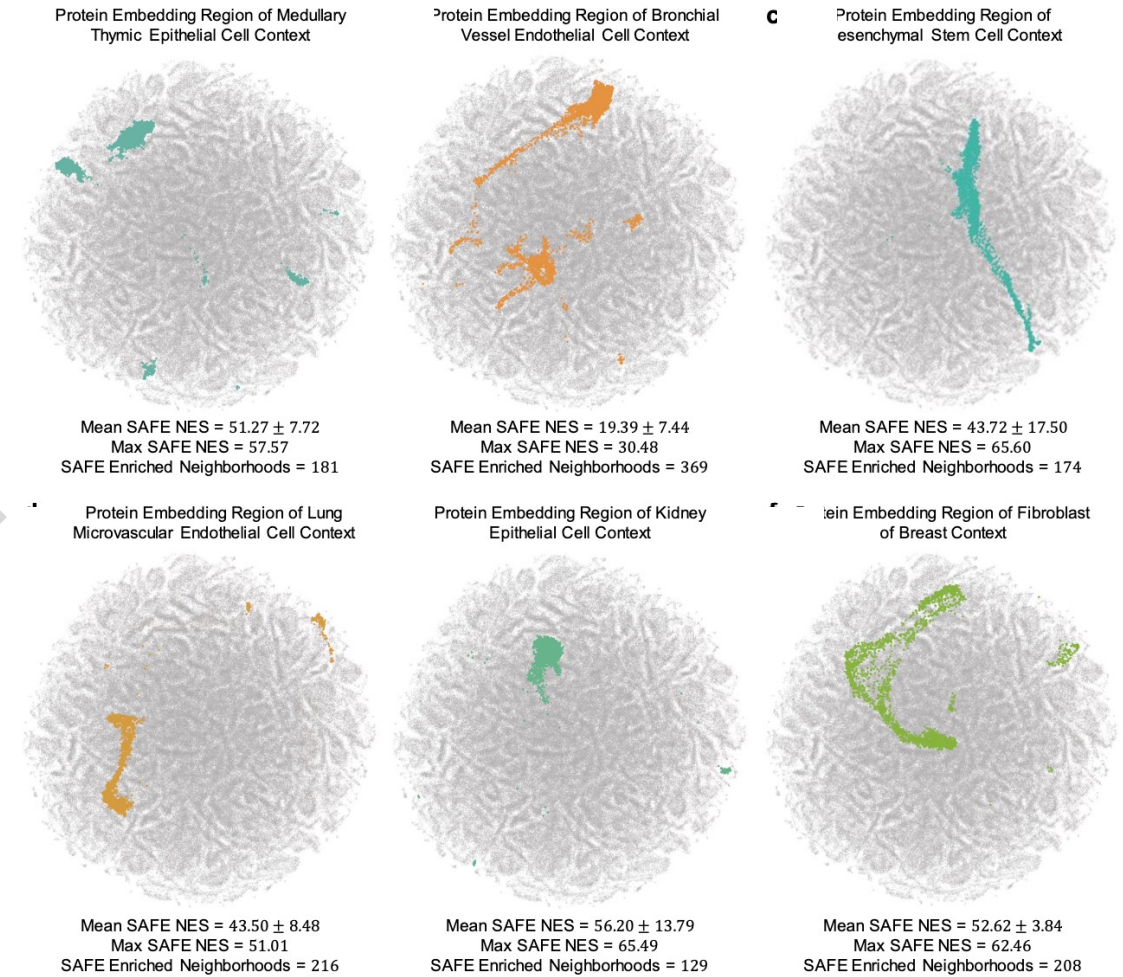
Protein networks across 156 cell type contexts spanning 62 tissues of varying hierarchical scales



PINNACLE

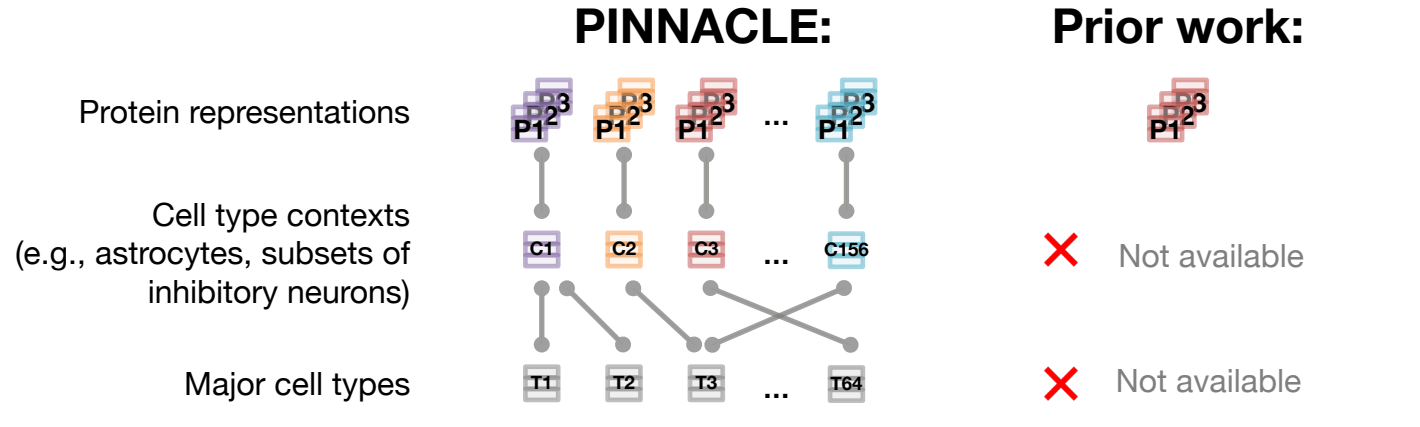


Self-supervised learning to learn a general encoder of protein embeddings tailored to cellular contexts

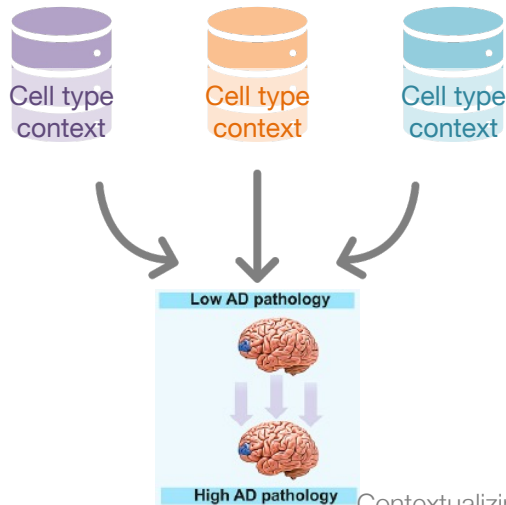


500,000 protein latent representations contextualized to 156 cross-organ cell types

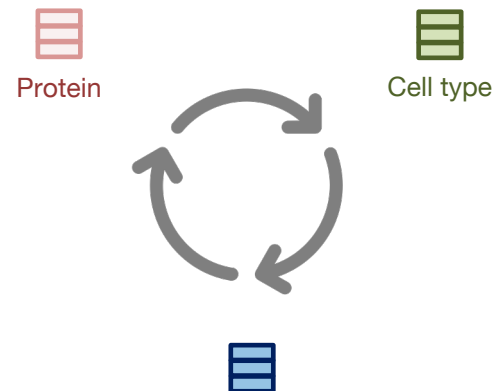
PINNACLE enables contextualized, precise predictions of drug effects across cell types and cell states



Multi-modal deep learning to identify changes at the single-cell level predictive of cognitive and behavioral phenotypes

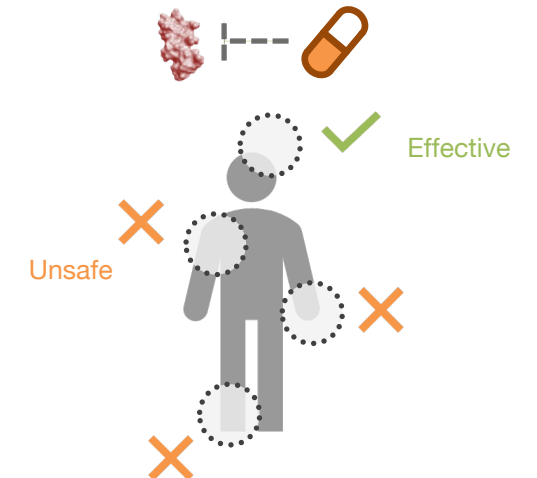


Transfer learning across cellular contexts to predict if candidate drugs affect disease-relevant cell types

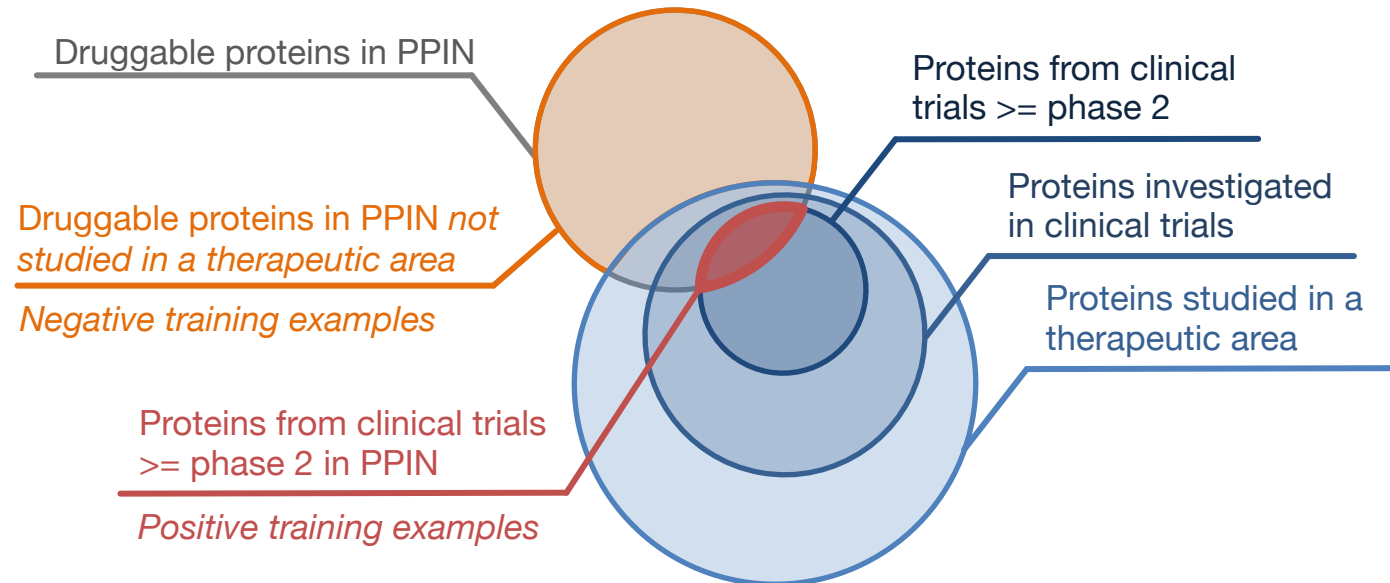
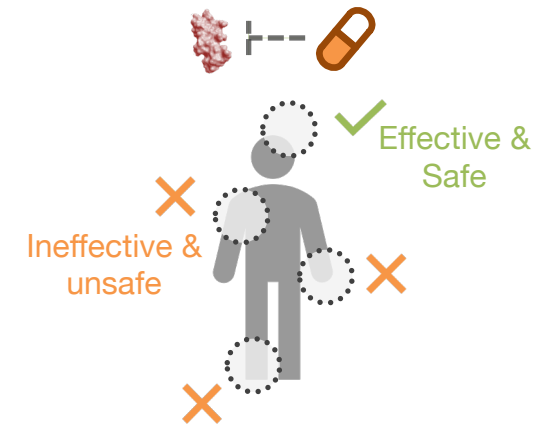
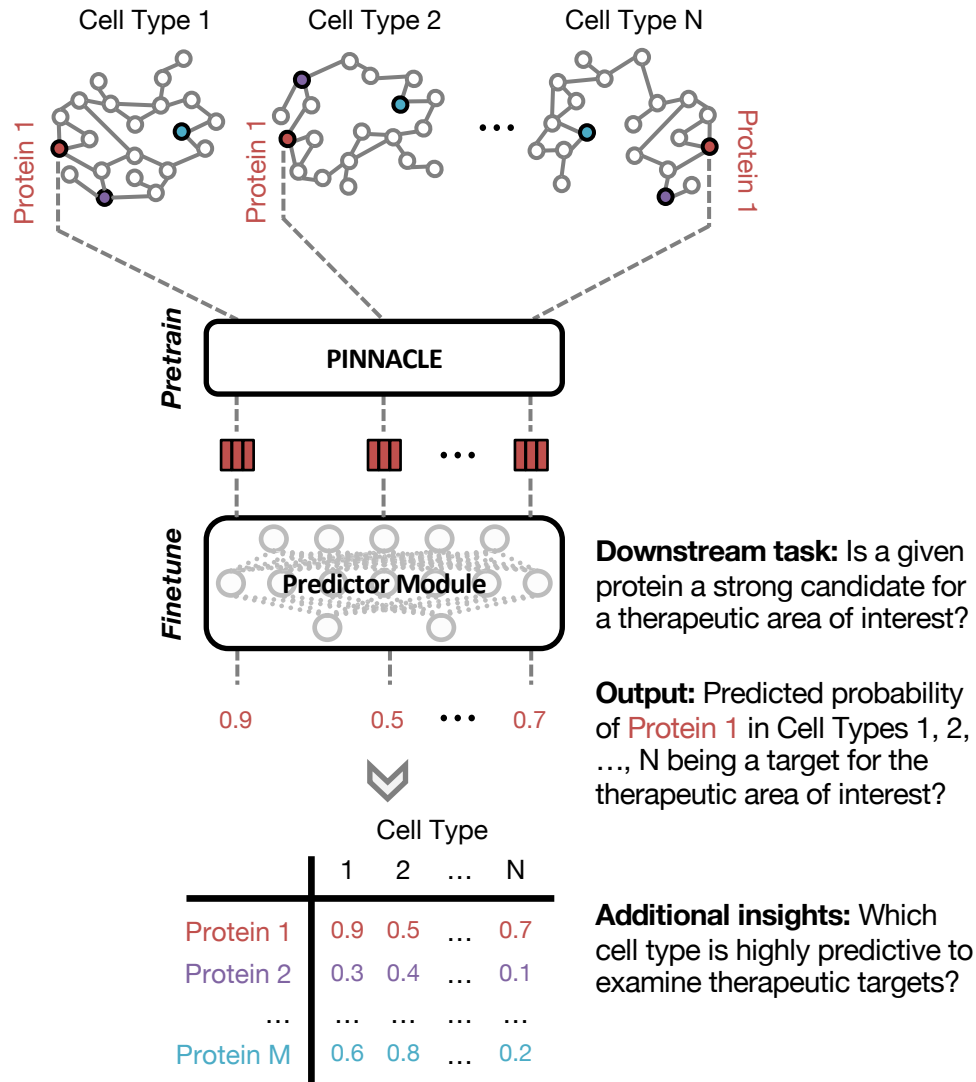


Functions: synaptic signaling, lipid metabolism, mitochondrial function, lipid and cholesterol biosynthesis

Identify molecules that are most effective at modulating astrocytes



Contextualized prediction: Setup

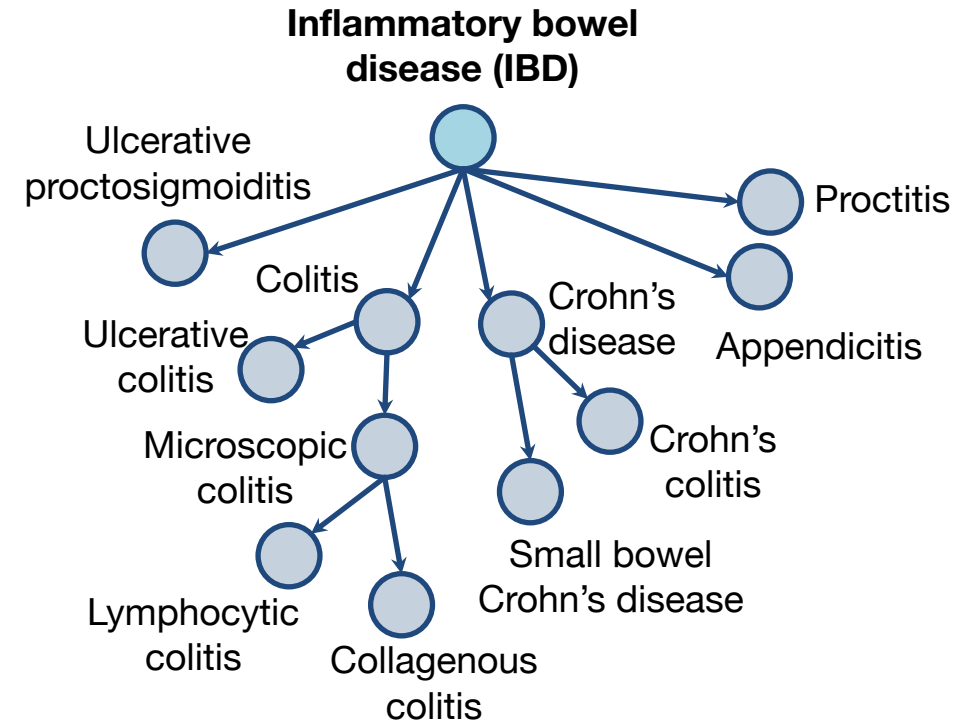
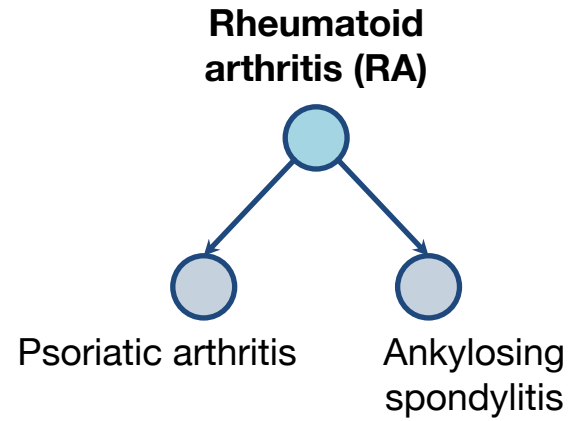


Dataset & experimental setup

1. Specify therapeutic area
(seed disease and its descendants in a disease ontology)

2. Curate clinical trials for diseases
(at least one completed clinical phase II or more)

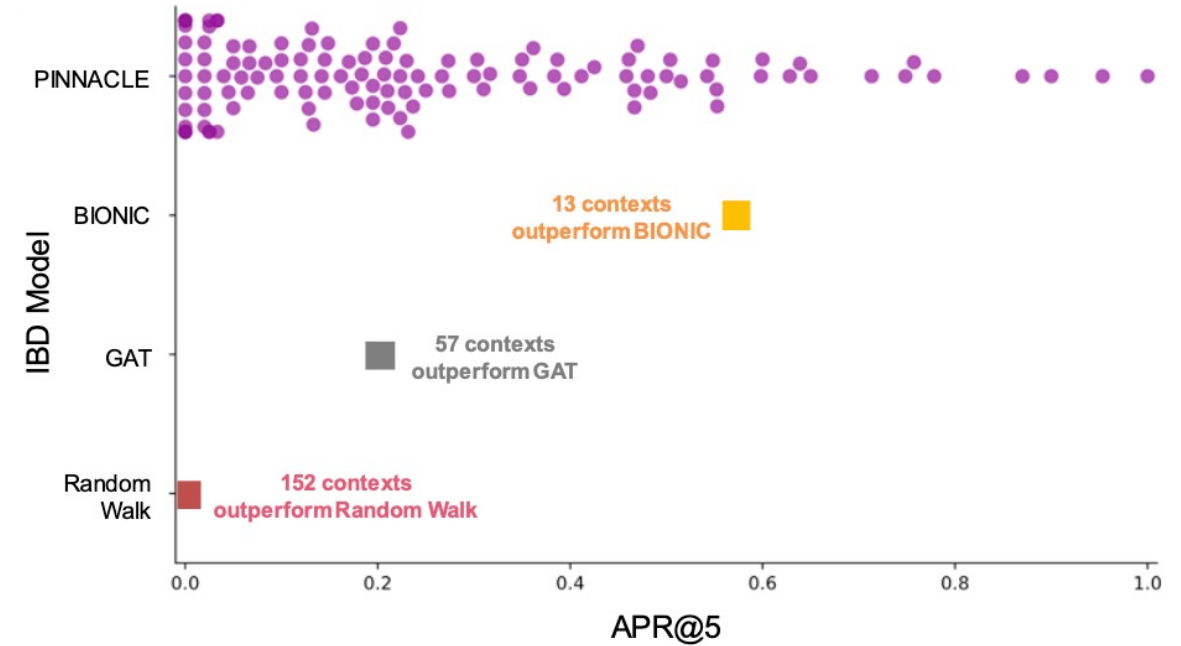
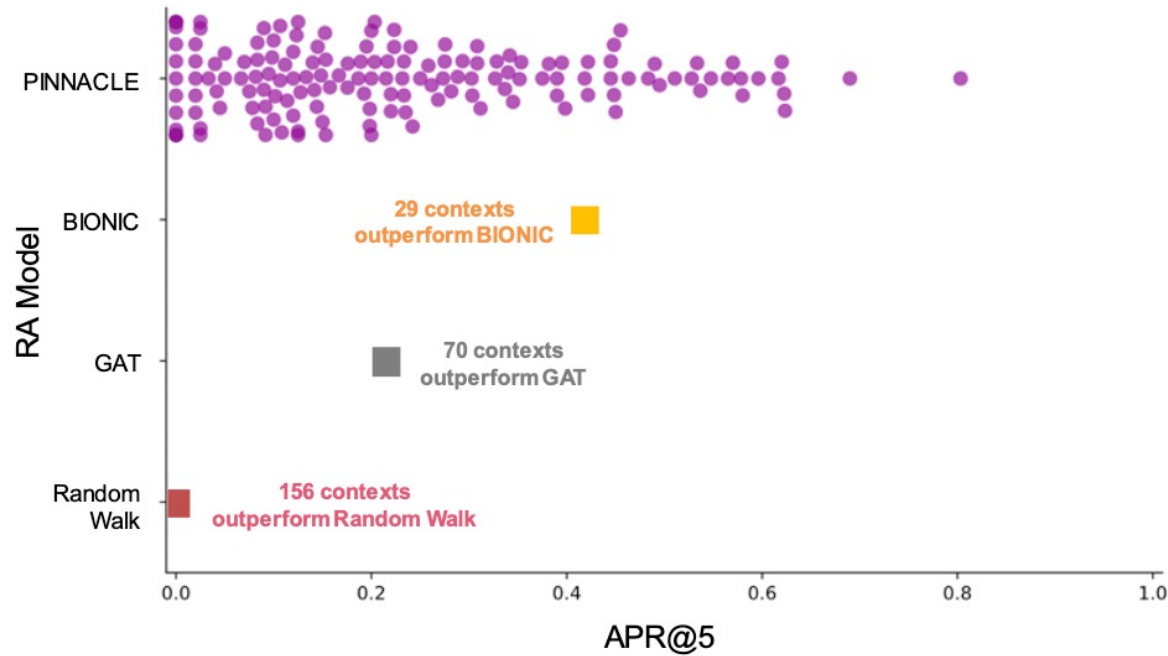
3. Curate list of candidate targets
(proteins targeted by a drug in a clinical trials)



Clinical phase	Unique drugs	Unique proteins
2	81	110
3	27	26
4	94	49

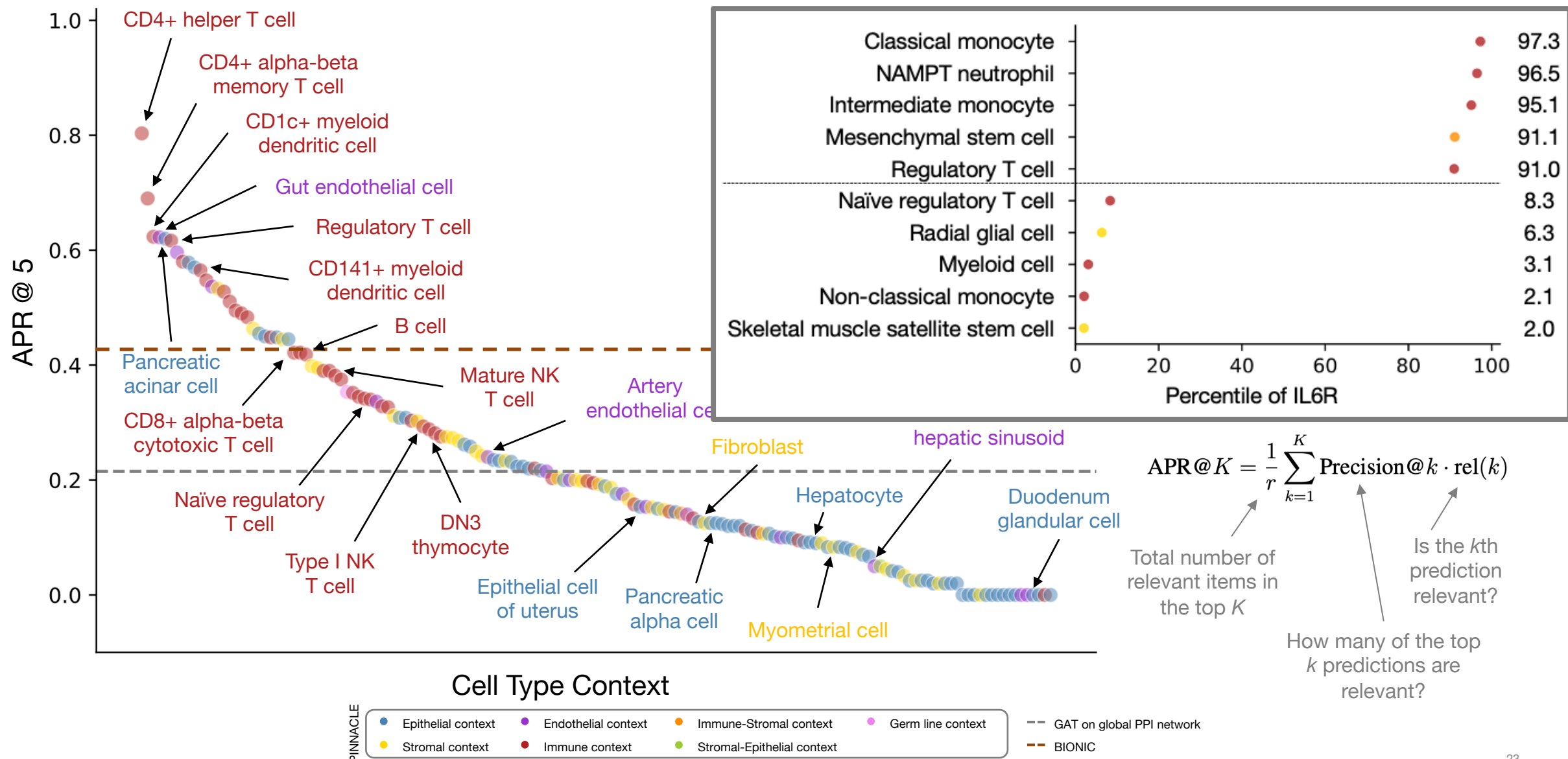
Clinical phase	Unique drugs	Unique proteins
2	41	67
3	21	26
4	59	46

PINNACLE accurately nominates therapeutic targets for RA and IBD in a cell-type specific manner, whereas existing data integration models conflate cell-type specific information, leading to poor performance

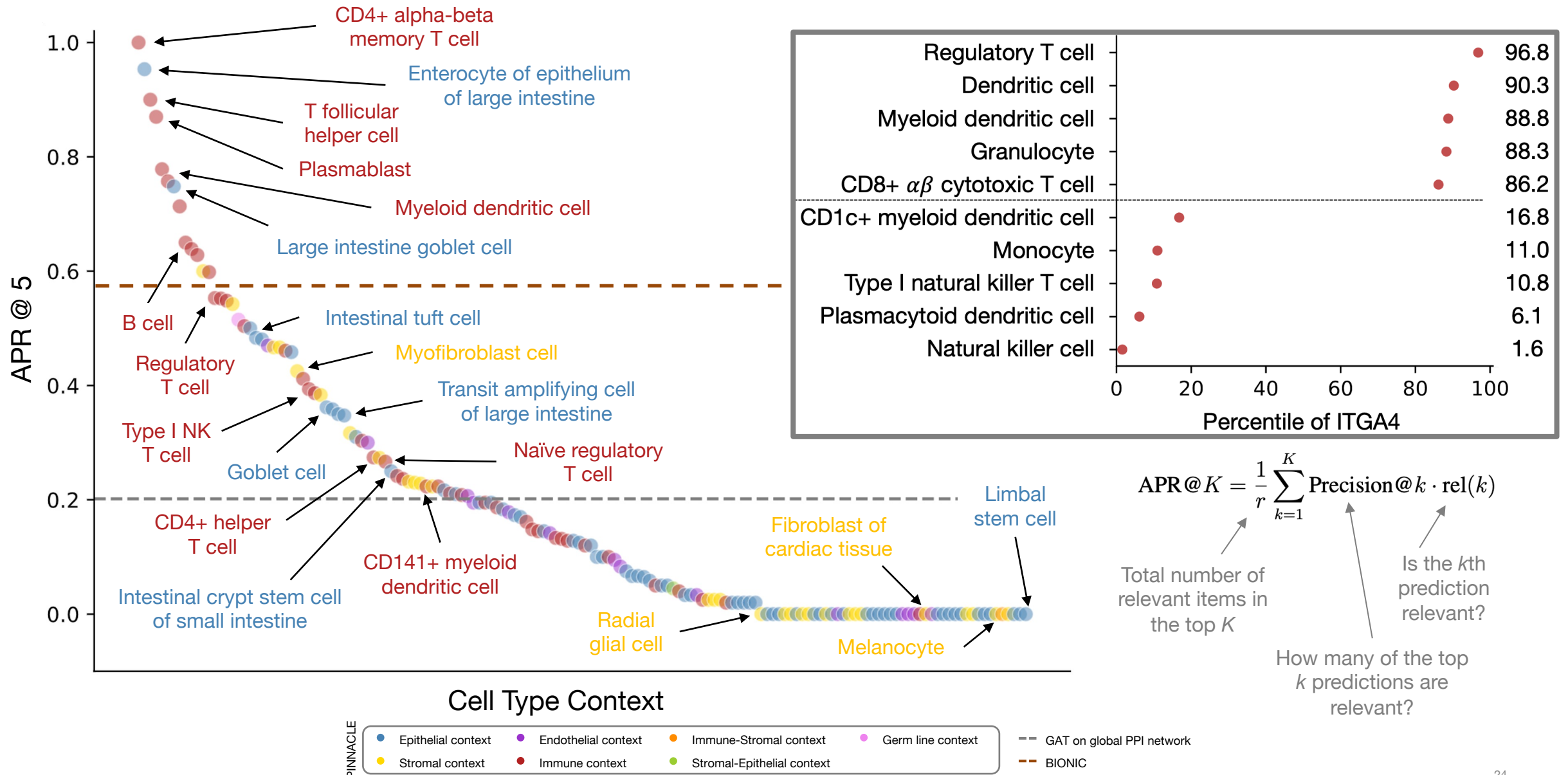


PINNACLE's representations **predict drug targets for RA and IBD significantly better than context-free methods**. PINNACLE also finds most predictive cell type contexts for investigating protein targets. PINNACLE supports the study of drug effects across cell type contexts.

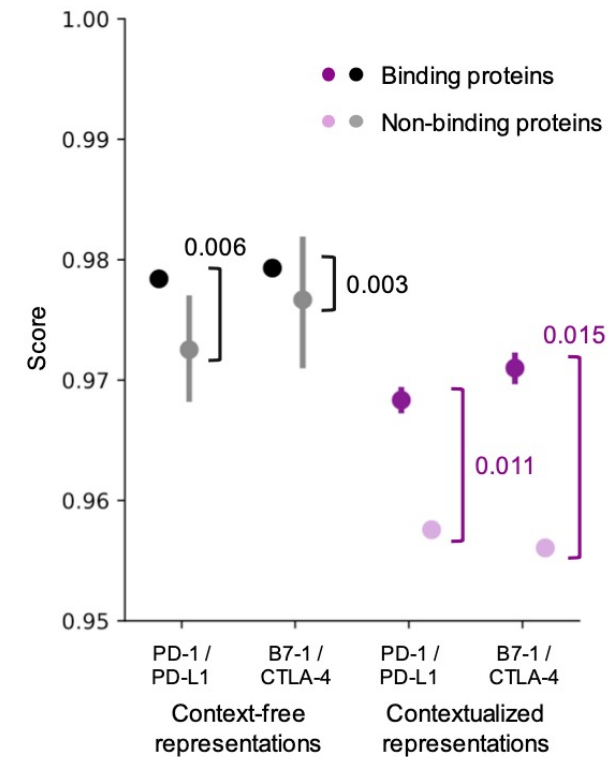
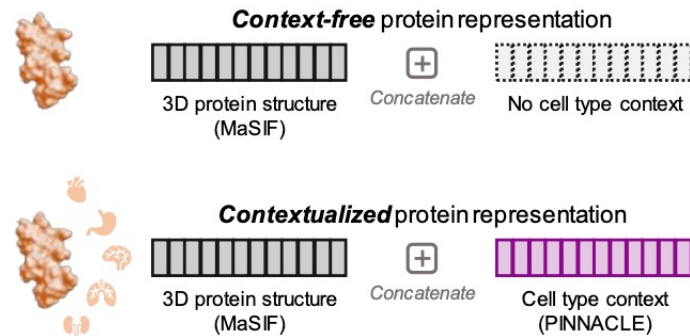
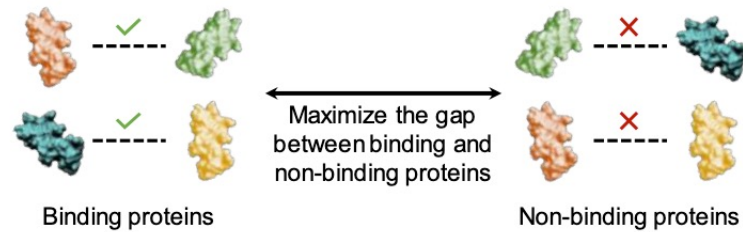
PINNACLE considers core disease processes in rheumatoid arthritis to identify candidate targets in cell type-specific manner



PINNACLE considers core disease processes in inflammatory bowel disease to nominate candidate targets in cell type-specific manner

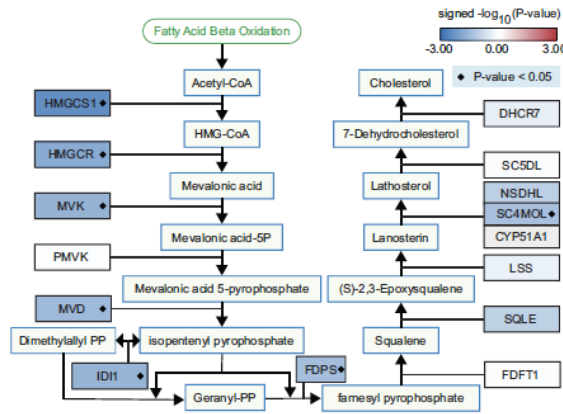


Contextualized predictions by integrating PINNACLE with 3D structures



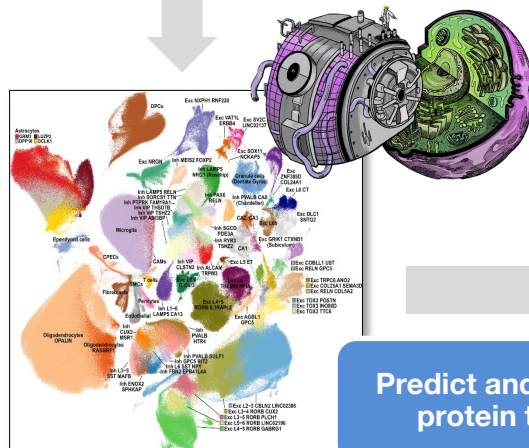
PINNACLE's representations **enhance 3D structure-based protein representations** for important protein interactions in immuno-oncology (PD-1/PD-L1 and B7-1/CTLA-4)

Disease circuitry across diseases and individuals



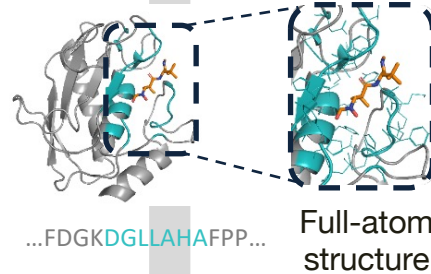
Key disease mechanisms, shared effects, interaction effects

Multimodal representation learning models
PINNACLE, PDGrapher

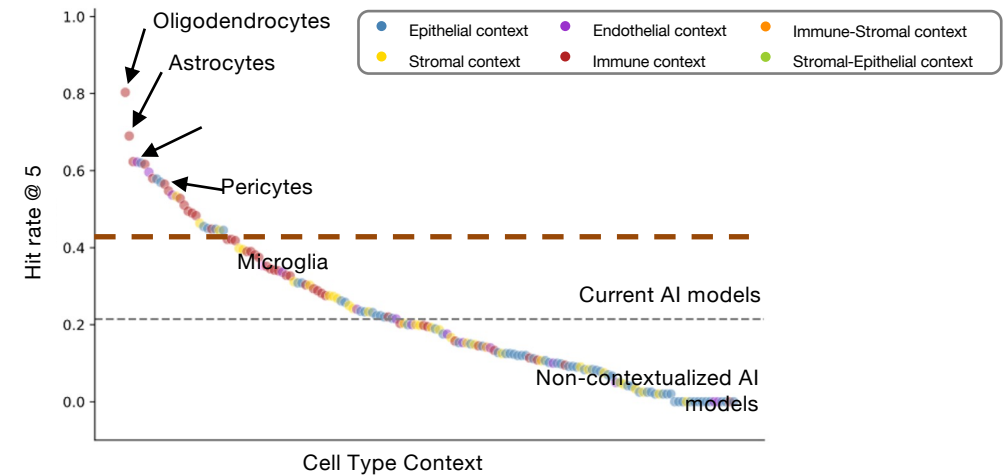


Predict and optimize protein targets

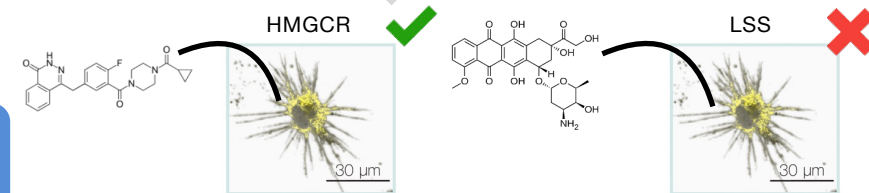
Generative geometric deep learning models
PocketGen, FAIR, and TxPLM



Generate, and refine therapeutic candidates



Rank-ordered lists of molecules to modulate disease circuitry in a cell-type specific manner



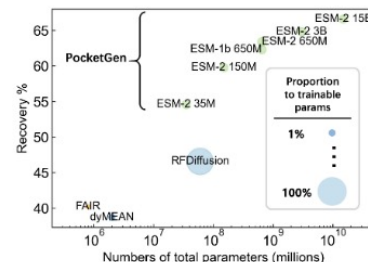
Chemistry optimization, synthesis, automation



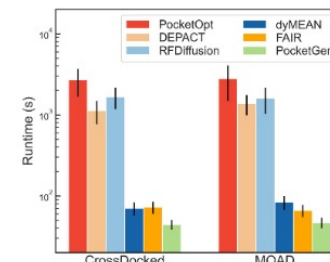
Generative sequence-structure models enable atom-level predictions of ligands binding to biological targets

Generative models:

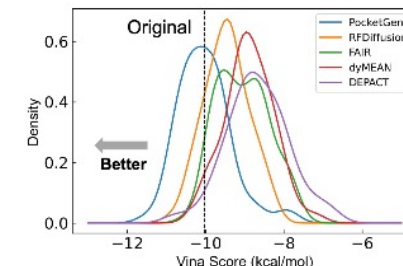
- **co-design of protein pocket sequence and 3D structure**
- **selective small molecule ligands**
- **optimized PPI interfaces**



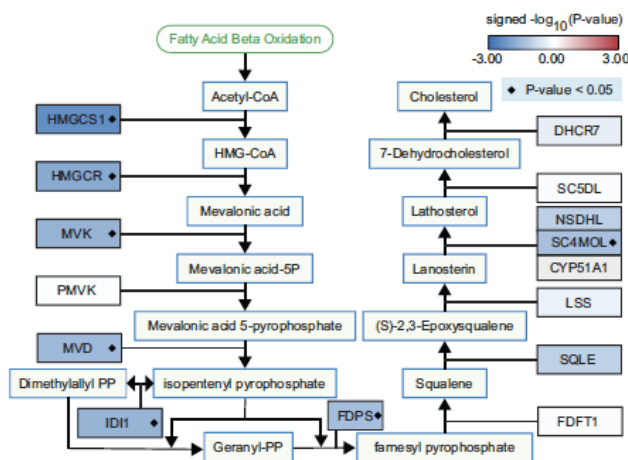
Iterative refinement based on side-chain effects, ligand flexibility, and sequence-structure consistency



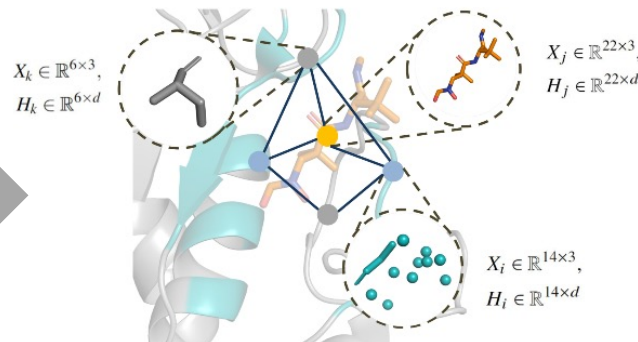
10x faster than current AI, 15% better accuracy (AAR, RMSE, docking score)



45% better hit rate than current AI, need to generate fewer molecules to find a hit



Key disease mechanisms, shared effects, interaction effects



Unified Representation

Generative sequence-structure models

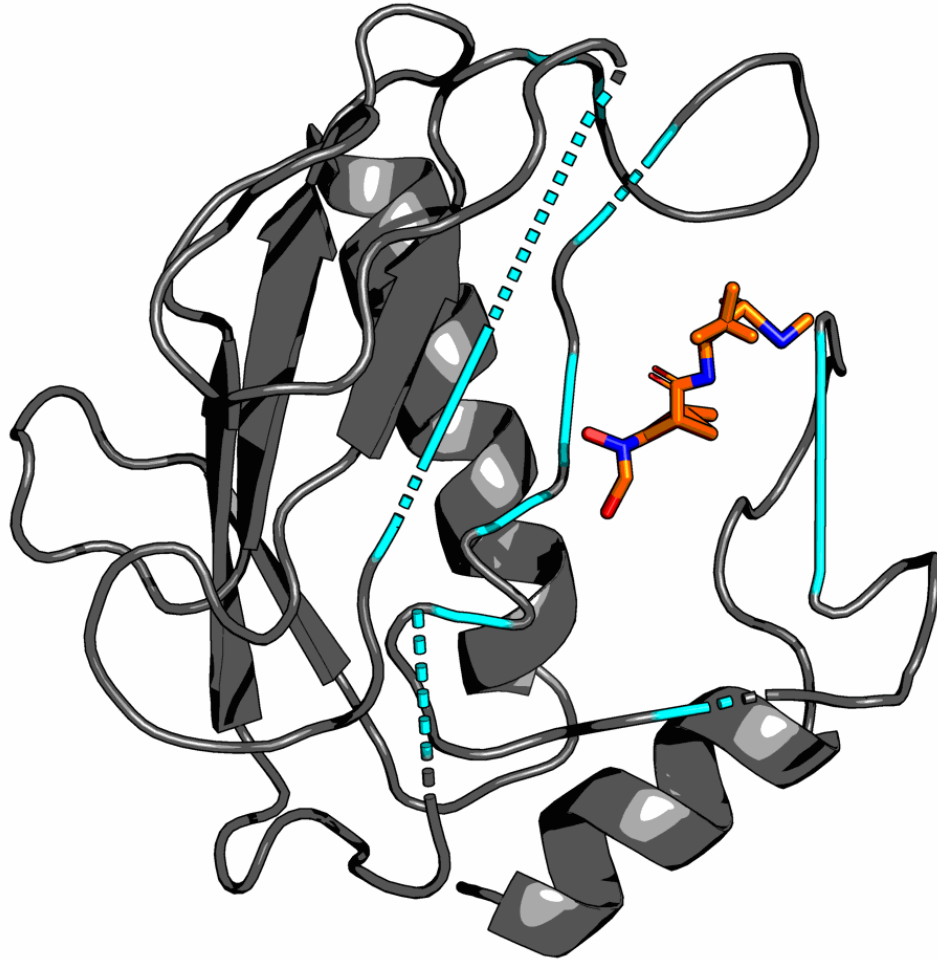
Molecular property predictors

- LSS-MSMO1 binding
- HMGCR-targeted
- PROTAC comprising a VHL ligand
- ...

Predictions

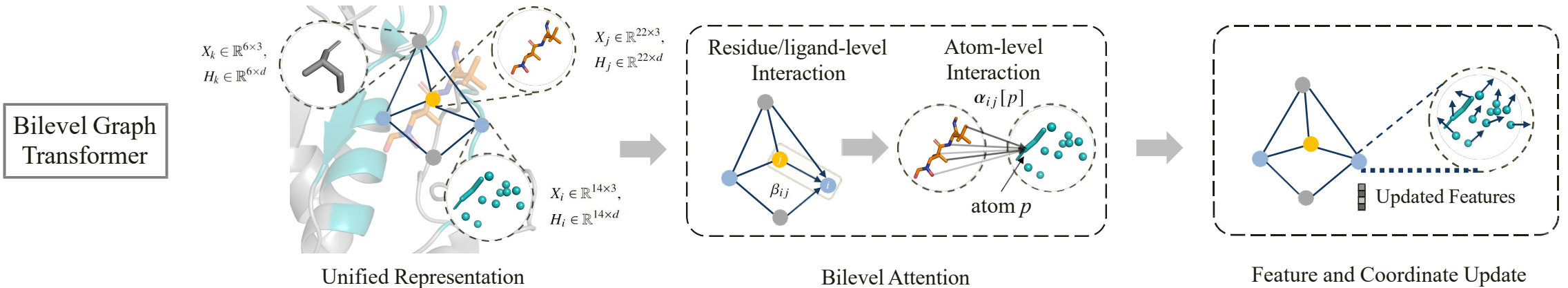
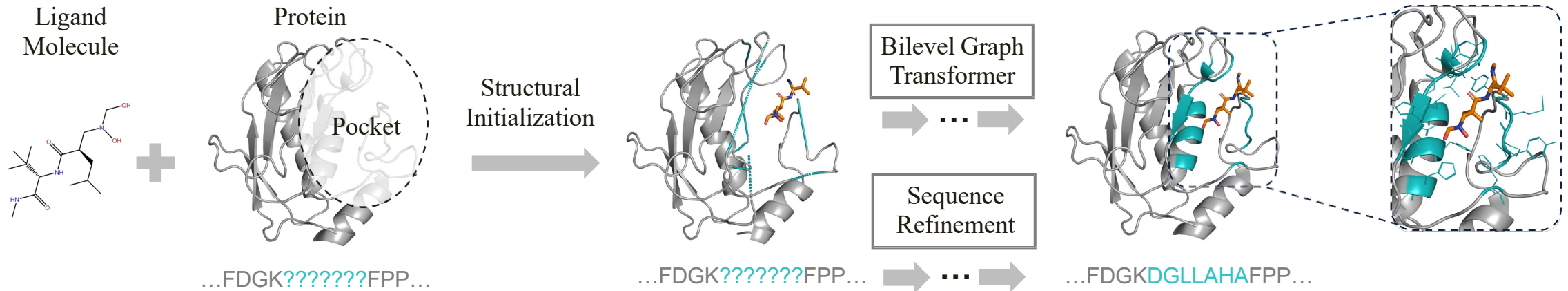
Priority lists of generated molecular structures

Sequence-structure co-generation of protein pockets

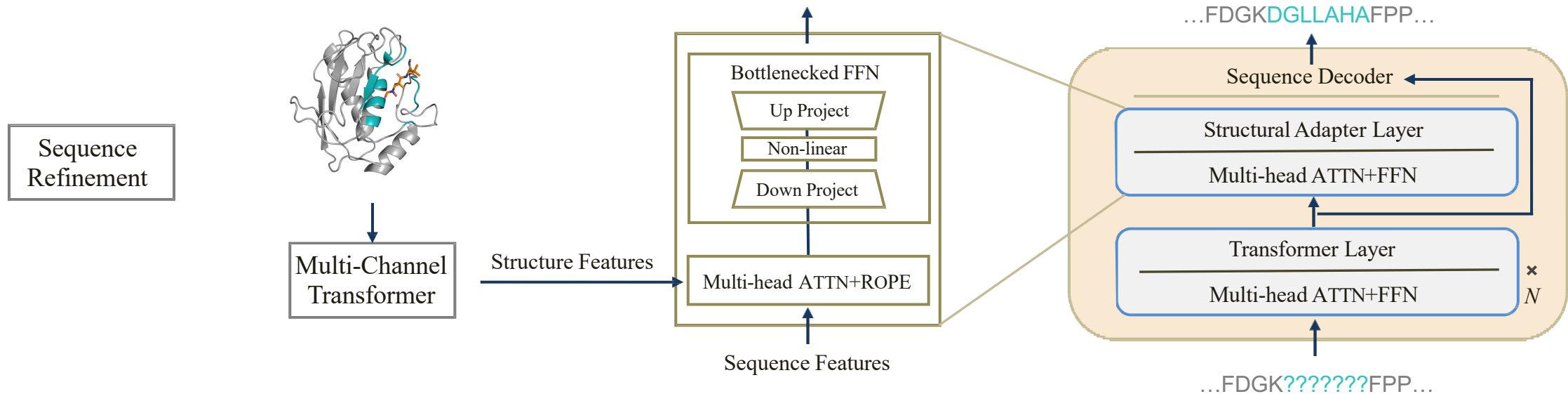
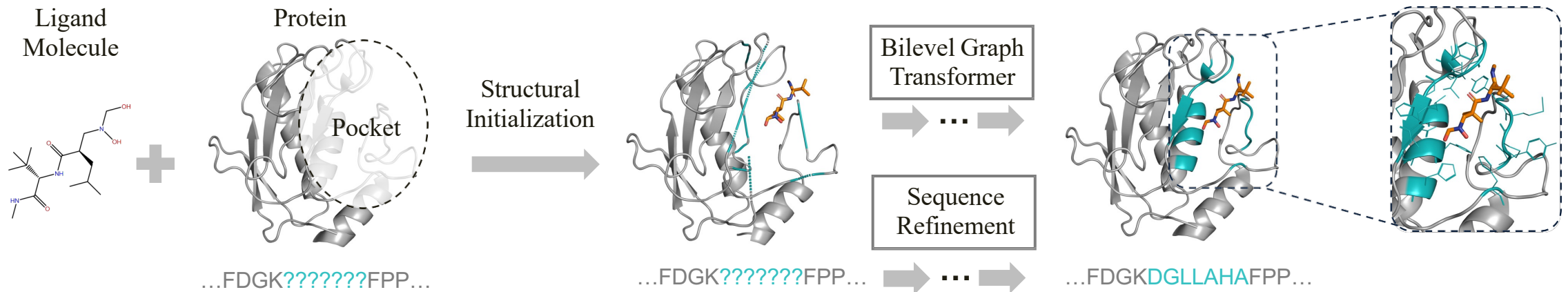


- Generating **high-fidelity protein pockets**—an area where a protein interacts with a ligand molecule
 - Complex interactions between ligand molecules and proteins
 - Flexibility of ligands and AA side chains
 - Complex sequence-structure dependencies
- PocketGen generates residue sequence and full-atom structure within protein pocket region

Iterative refinement of both sequence and structure in the protein pocket to maximize binding affinity with small molecule ligand



Iterative refinement of both sequence and structure in the protein pocket to maximize binding affinity with small molecule ligand

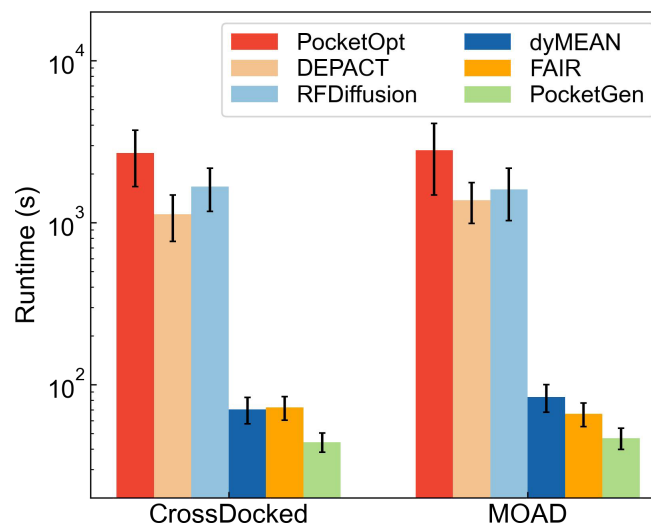


PocketGen generates protein pockets with higher binding affinity and structural validity than existing models

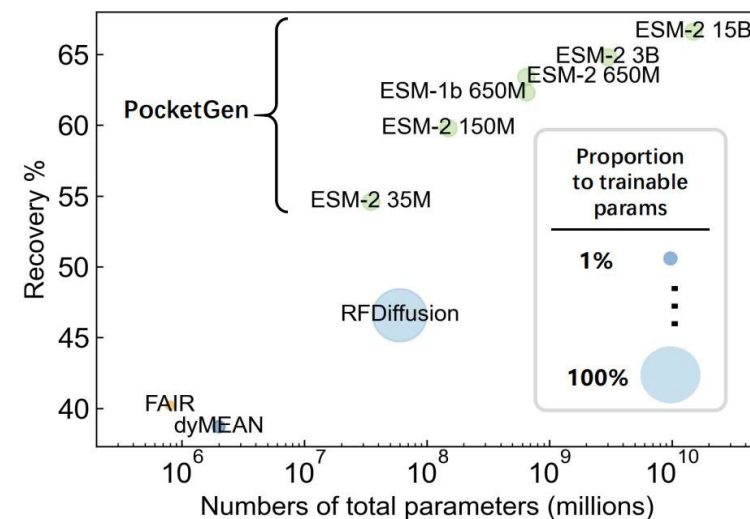
	PocketOpt	DEPACT	dyMEAN	FAIR	RFDiffusion	PocketGen
Top-1 generated protein pocket						
Vina score (↓)	-9.216	-8.527	-8.540	-8.792	-9.037	-9.655
Success Rate (↑)	0.92	0.75	0.76	0.80	0.89	0.97
RMSD (↓)	-	1.47	1.44	1.39	1.13	1.21
pLDDT (↑)	-	82.1	83.3	83.2	84.5	86.7
scTM (↑)	-	0.901	0.906	0.899	0.924	0.937
Top-3 generated protein pockets						
Vina score (↓)	-8.878	-8.131	-8.196	-8.321	-8.876	-9.353
RMSD (↓)	-	1.45	1.43	1.40	1.18	1.24
pLDDT (↑)	-	81.9	82.8	83.1	84.6	86.2
scTM (↑)	-	0.896	0.892	0.897	0.929	0.934
Top-5 generated protein pockets						
Vina score (↓)	-8.702	-7.786	-7.974	-7.943	-8.510	-9.239
RMSD (↓)	-	1.46	1.45	1.42	1.25	1.22
pLDDT (↑)	-	82.2	82.9	83.3	84.3	86.1
scTM (↑)	-	0.892	0.903	0.886	0.926	0.935
Top-10 generated protein pockets						
Vina score (↓)	-8.556	-7.681	-7.690	-7.785	-8.352	-9.065
RMSD (↓)	-	1.53	1.44	1.41	1.26	1.28
pLDDT (↑)	-	81.5	82.7	83.0	84.2	85.9
scTM (↑)	-	0.895	0.896	0.884	0.924	0.931

Improved structural validity, amino acid sequence recovery, and affinity with target ligands

Model	CrossDocked			Binding MOAD		
	AAR (↑)	RMSD (↓)	Vina (↓)	AAR (↑)	RMSD (↓)	Vina (↓)
Test set	-	-	-7.016	-	-	-8.076
DEPACT	31.52±3.26%	1.59±0.13	-6.632±0.18	35.30±2.19%	1.52±0.12	-7.571±0.15
dyMEAN	38.71±2.16%	1.57±0.09	-6.855±0.06	41.22±1.40%	1.53±0.08	-7.675±0.09
FAIR	40.16±1.17%	1.46±0.04	-7.015±0.12	43.68±0.92%	1.37±0.07	-7.930±0.15
RFDiffusion	46.57±2.07%	1.44±0.07	-6.936±0.07	45.31±2.73%	1.45±0.10	-7.942±0.14
PocketGen	63.40±1.64%	1.36±0.05	-7.135±0.08	64.43±2.35%	1.32±0.05	-8.112±0.14



Better generation efficiency



Performance wrt protein LM size

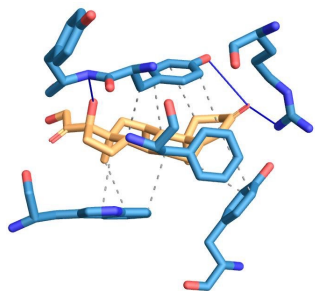
PocketGen can redesign pockets of antibodies, enzymes, and biosensors for target ligand molecules

- Protein
- Ligand
- Aromatic Ring Center
- Hydrophobic Interaction
- Hydrogen Bond
- ... π -Stacking (parallel)
- ... π -Stacking (perpendicular)
- ... π -Cation Interaction

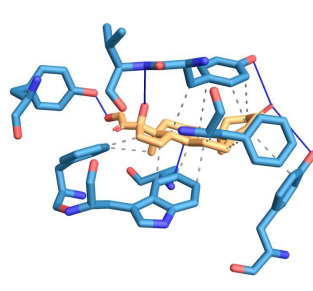
Cortisol (HCY)

...WFRYYDTMY...

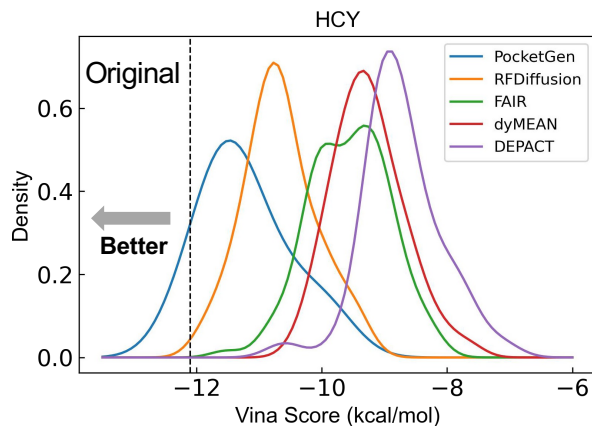
...WFRYVYSFY...



Original
(HP 12, HB 3)



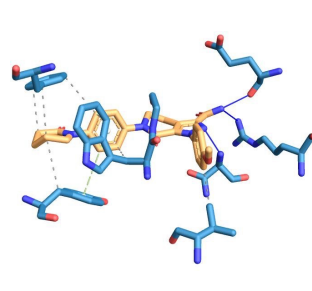
PocketGen
(HP 12, HB 5)



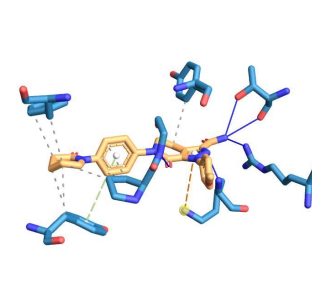
Apixaban (APX)

...YREFQVWGG...

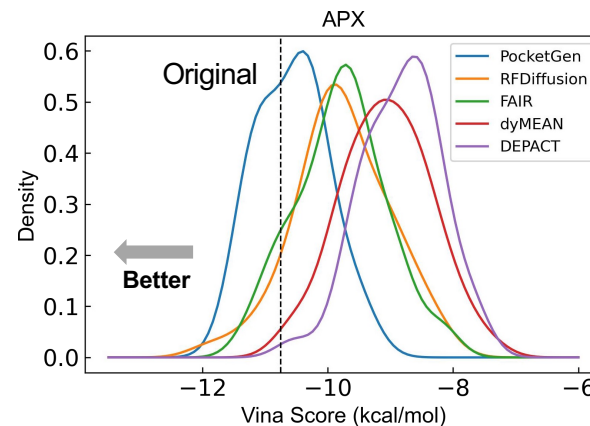
...YRTFKVPGY...



Original
(HP 7, HB 4, π 1)



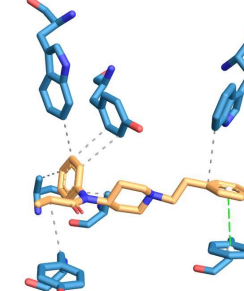
PocketGen
(HP 9, HB 5, π 2)



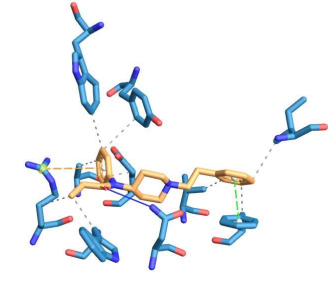
Fentanyl 7V7)

...FVWAWYWAY...

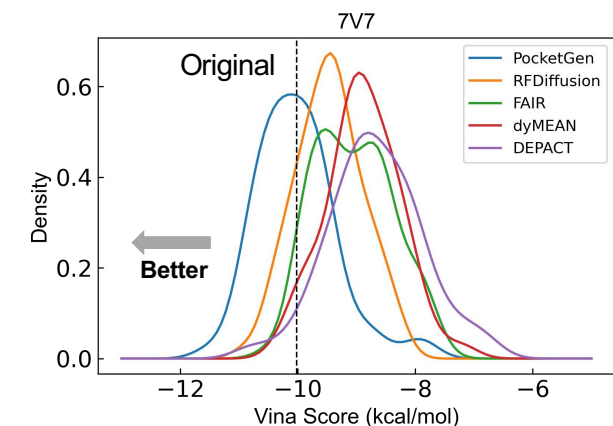
...FVRAWYWEW...

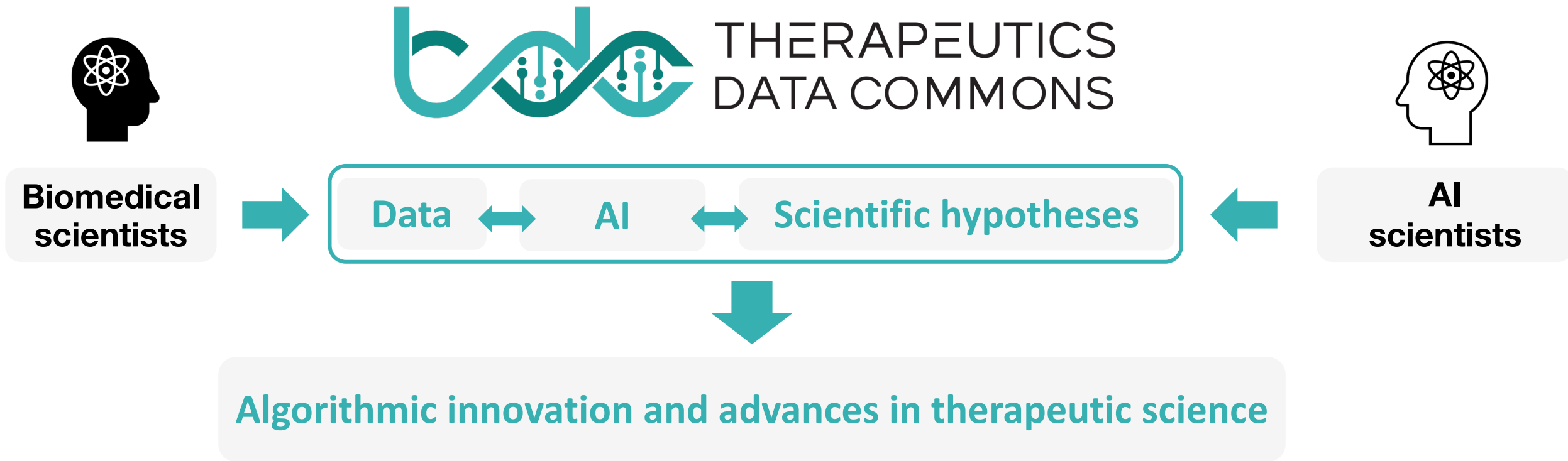


Original
(HP 7, HB 0, π 1)



PocketGen
(HP 9, HB 1, π 2)





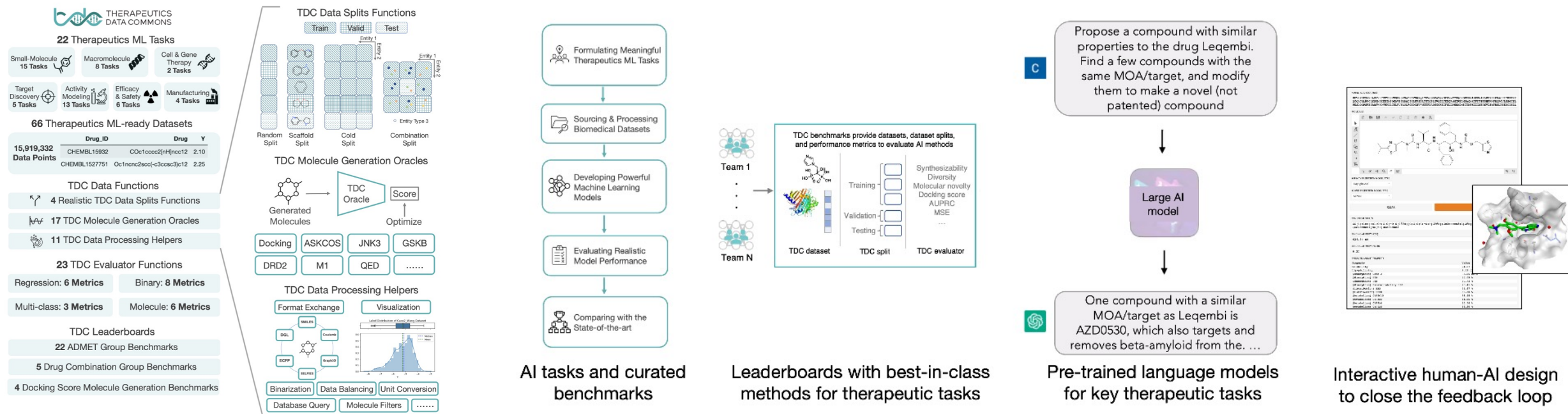
Global initiative to access and evaluate AI across therapeutic modalities and stages of drug discovery

270,000 active use cases of AI for drug design and therapeutic use prediction / 90,000 users worldwide

Ongoing collaborations with drug designers in **immuno-oncology**, **rare genetic diseases** and **neurodegenerative diseases**



We develop Therapeutics Commons, a data and AI model hub across therapeutics modalities and stages of discovery



- Molecular property prediction and optimization
- Binding of novel drugs to candidate therapeutic targets
- Molecule generation and AI-driven drug design
- Therapeutic use prediction, manufacturing, efficacy & safety
- AI models cover small molecules, proteins, peptides, miRNAs, and gene editing

270,000 active use cases of AI for drug design and therapeutic use prediction / 90,000 users worldwide

Partners validating AI models for **immuno-oncology**, **rare genetic diseases** and **neurodegenerative diseases**

Capabilities of the Commons

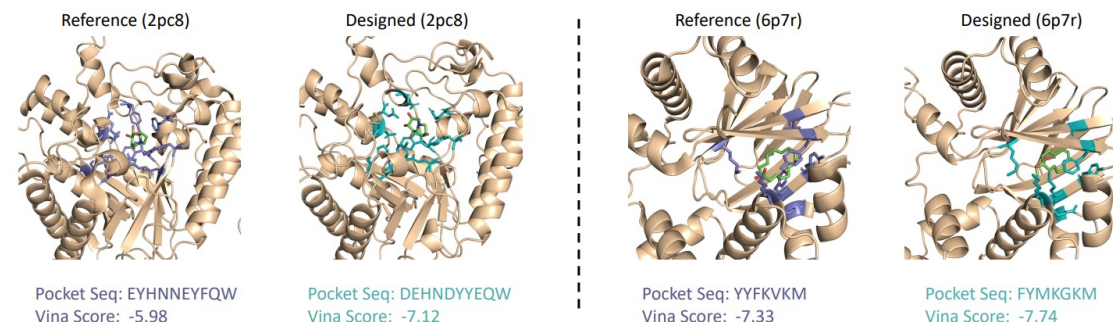
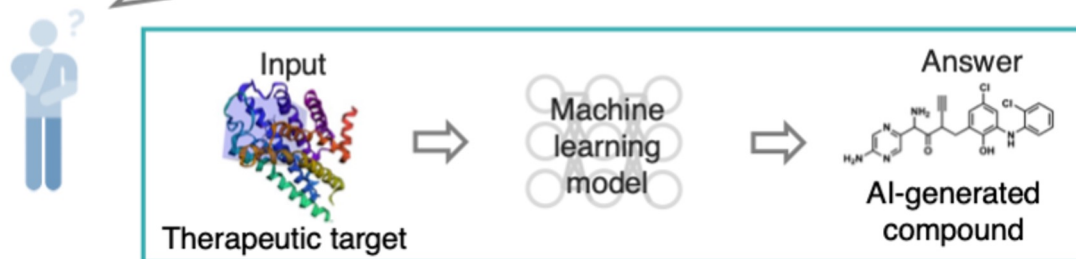
Predictive AI:

- Biological target nomination
- Target validation
- Molecular interaction screening

Generative AI:

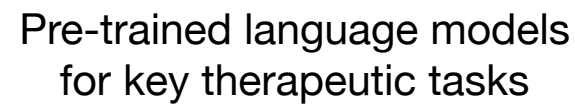
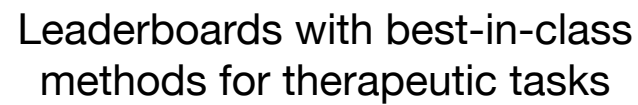
- Mutation effect prediction
- Binding of ligands to targets
- Molecular design and optimization

I want to generate a highly potent compound that effectively binds a therapeutic target.

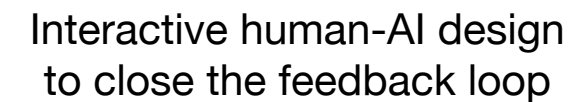


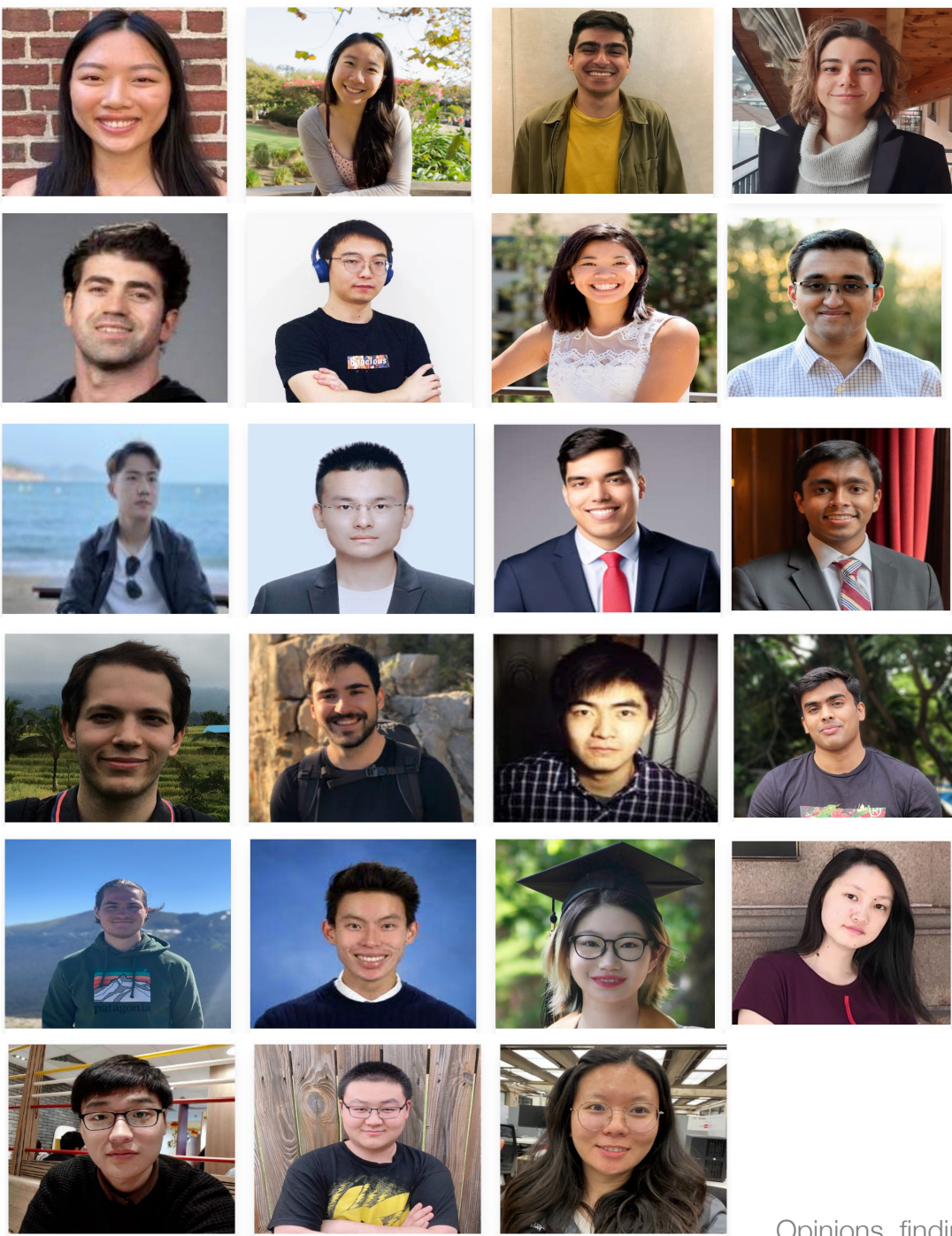
Model	CrossDocked			Binding MOAD		
	AAR (↑)	RMSD (↓)	Vina (↓)	AAR (↑)	RMSD (↓)	Vina (↓)
PocketOptimizer	27.89±14.9%	1.75±0.08	-6.905±2.39	28.78±11.3%	1.68±0.12	-7.829±2.41
DEPACT	22.58±8.48%	1.97±0.14	-6.670±2.13	26.12±8.97%	1.76±0.15	-7.526±2.05
HSRN	31.62±10.4%	2.15±0.17	-6.565±1.95	33.70±10.1%	1.83±0.18	-7.349±1.93
Diffusion	34.62±13.7%	1.68±0.12	-6.725±1.83	36.94±12.9%	1.47±0.09	-7.724±2.36
MEAN	35.46±8.15%	1.76±0.09	-6.891±1.86	37.16±14.7%	1.52±0.09	-7.651±1.97
FAIR	40.17±12.6%	1.42±0.07	-7.022±1.75	43.75±15.2%	1.35±0.10	-7.978±1.91

AI tasks and curated benchmarks



One compound with a similar MOA/target as Dasatinib is AZD0530, which also inhibits Fyn kinase ...





marinka@hms.harvard.edu



AI models, datasets and papers
zitniklab.hms.harvard.edu

Therapeutics Commons
tdcommons.ai

Opinions, findings, conclusions or recommendations expressed here do not necessarily reflect the views of the funders.