

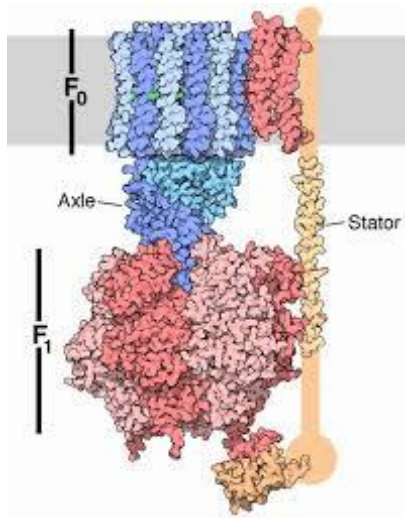
# Protein Design and Engineering with Multi-Modal Generative AI Models

**Yunan Luo**

School of Computational Science and Engineering  
College of Computing  
Georgia Institute of Technology

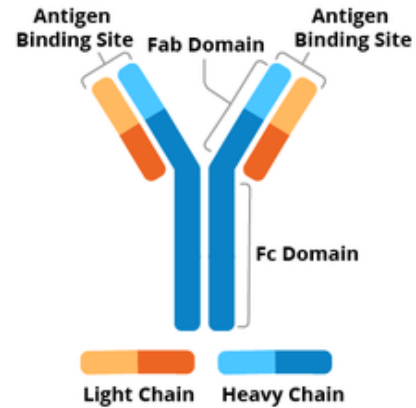


# Proteins perform diverse biological functions



**ATP Synthase**

Producing energy  
currency (ATP) of life



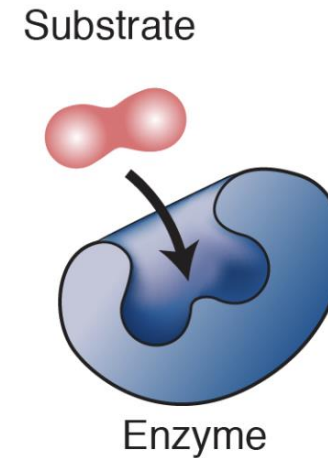
**Antibody**

Detecting and  
neutralizing viruses



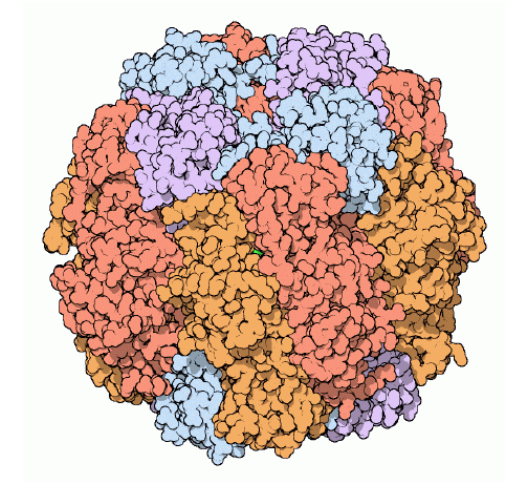
**Insulin**

Regulating blood  
sugar



**Enzyme**

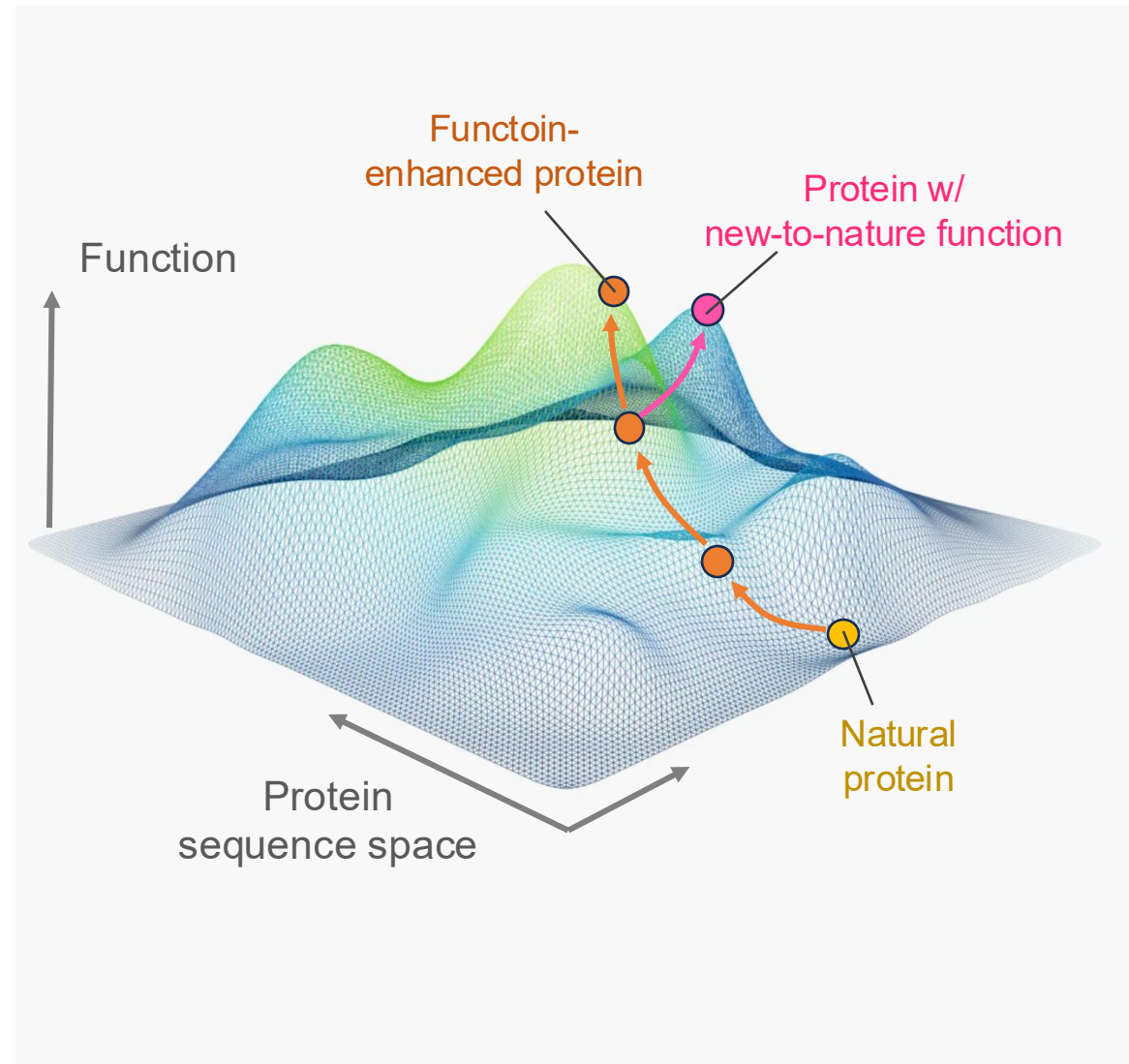
Speeding up  
chemical reactions



**Rubisco**

Capturing CO<sub>2</sub> in  
photosynthesis

# Protein design and engineering



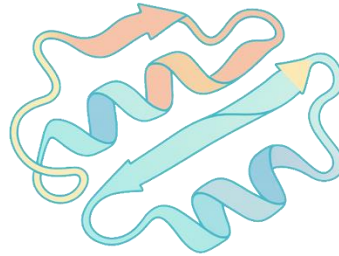
# How is a protein's function determined?

Protein's sequence-structure-function relationship

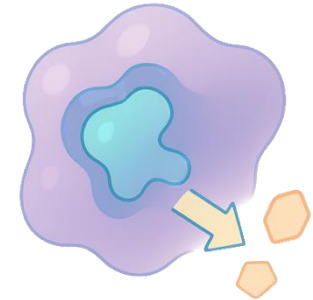
Protein sequence



Protein structure



Protein function



**Nobel Prize in Chemistry 2024**  
*for protein structure prediction (AlphaFold)*



Demis Hassabis



John Jumper

**This project:** *functional protein design*

# How to design functional proteins?

- Learning from natural protein evolution using GenAI

$$p(\text{functional protein}) \approx p(\text{natural protein})$$

MSKGEELFTGVVPIL  
MSK**C**EELFTGVVPIL  
MSKG**W**EELFTGVVPIL  
MSKGEELFT**S**VVPIL  
MSKGE**F**G**F**AGVVPIL  
M**D**YGEELFT**V**VVP**S**L  
...

Natural proteins



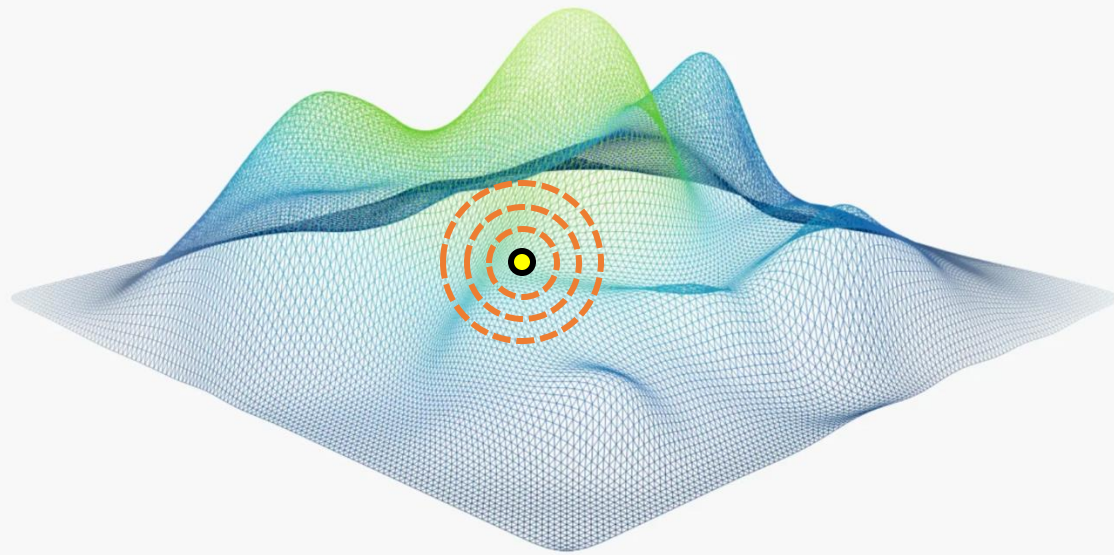
$p(\text{natural protein})$



**M****Y****K****R****E****L****W****G****V****S****T****I****S****Y**

Novel functional proteins

## Protein fitness landscape

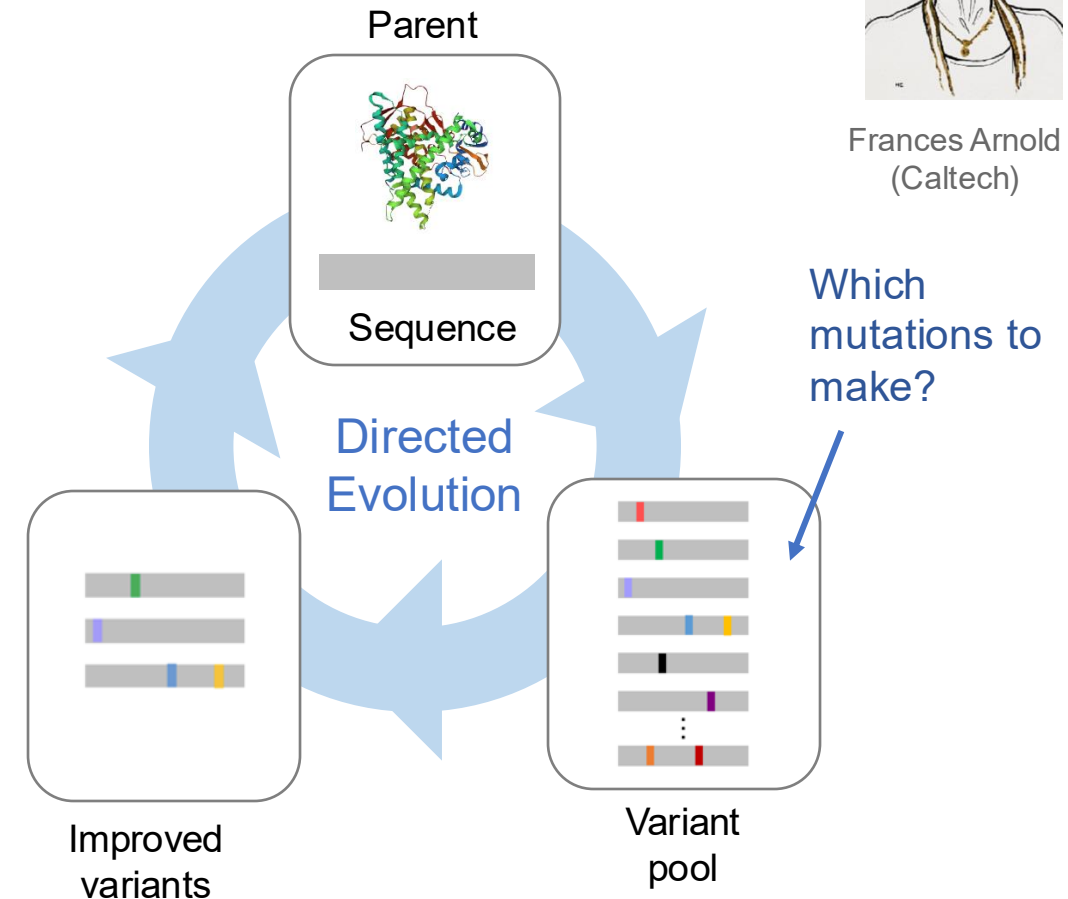


## How to find functional proteins?

Directed evolution for protein engineering  
(2018 Nobel Prize in Chemistry)



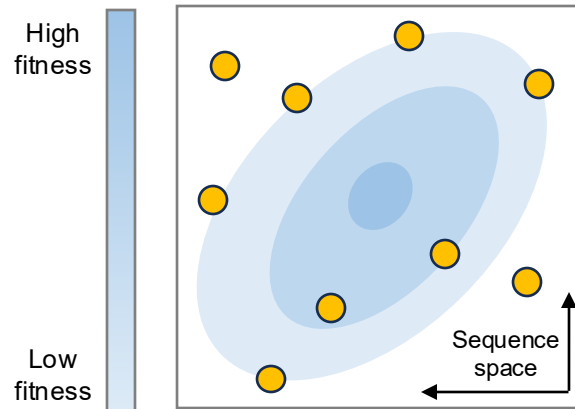
Frances Arnold  
(Caltech)



# Protein engineering for new-to-nature function

- **Task:** Design functional and diverse proteins
- **Challenge:** No existing fitness data (e.g., because the function is new-to-nature)

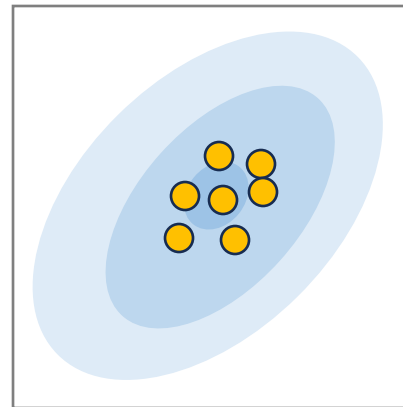
## Conventional starting library design (e.g., NNK library)



(random search, many are non-functional)

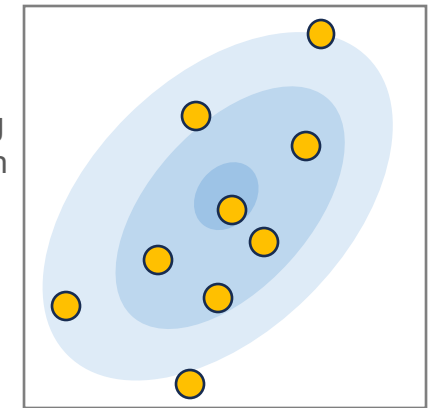
## Our approach: Pareto-optimal library design (MODIFY)

Goal 1: High *fitness*

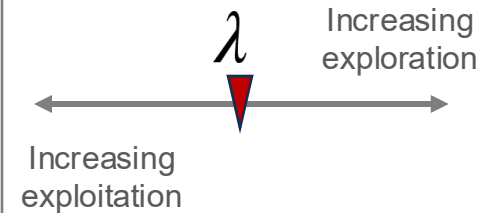


(most variants in the library are functional)

Goal 2: High *diversity*

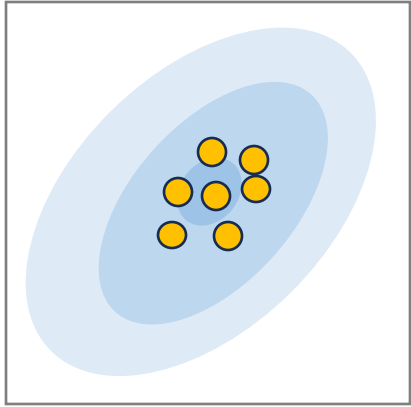


(the library contains novel protein variants)



# Zero-Shot Fitness Prediction using Foundation Models

Goal 1: High *fitness*

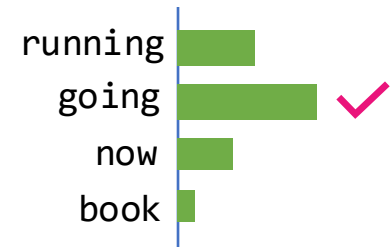
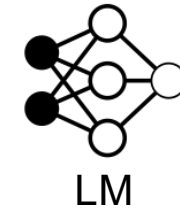


**Q:** Cold-start problem. No data to train a supervised fitness predictor.

**A:** Use pre-trained *protein language models* to make zero-shot fitness prediction

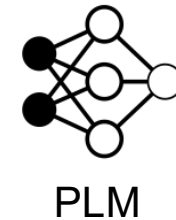
- *Natural language model (LM):*

Where are we going



- *Protein language model (PLM):*

MSTYSTFVLLGVEY

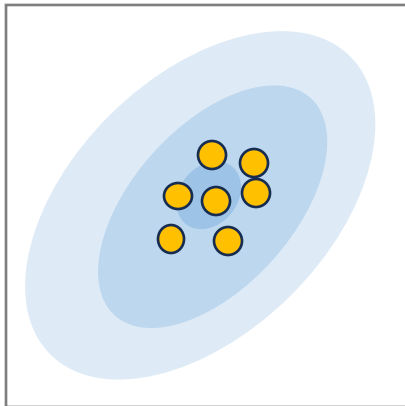


- PLMs are trained on protein sequences we observed in nature
- Evolutionary plausibility correlates with fitness

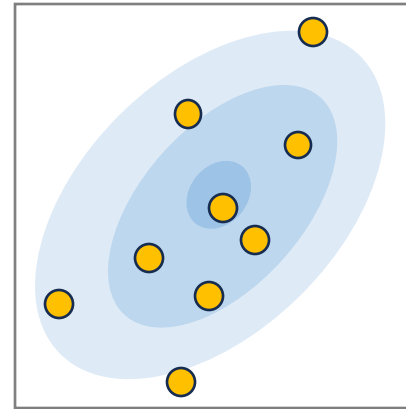


# Characterizing the library diversity

Goal 1: High *fitness*



Goal 2: High *diversity*



← Increasing exploitation  
 $\lambda$   
Increasing exploration →

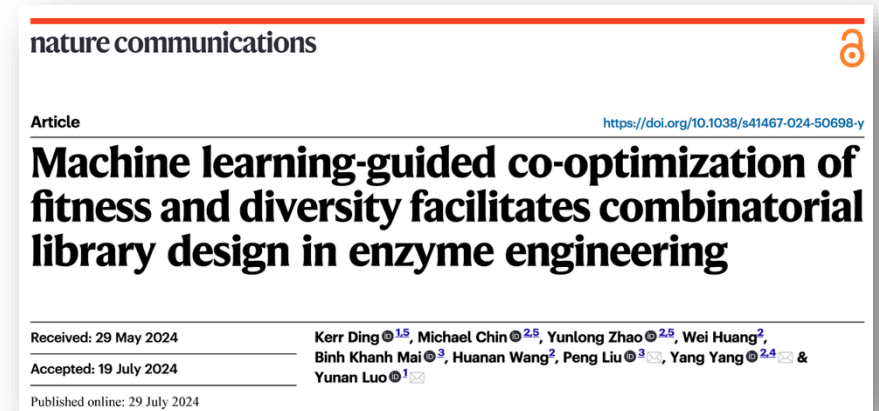
- **Goal:**

$$\max_{p \in \mathcal{P}} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \text{fitness}(\mathbf{x}) + \lambda \cdot \text{diversity}(p)$$

Predicted by protein  
language model

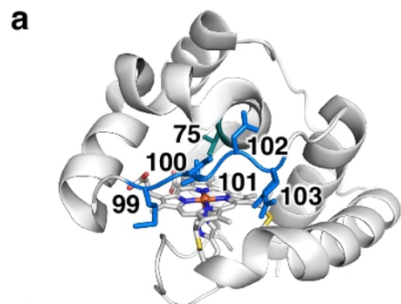
Quantified by sequence  
entropy

- A Pareto optimization problem
- Learned a probability distribution over protein sequence

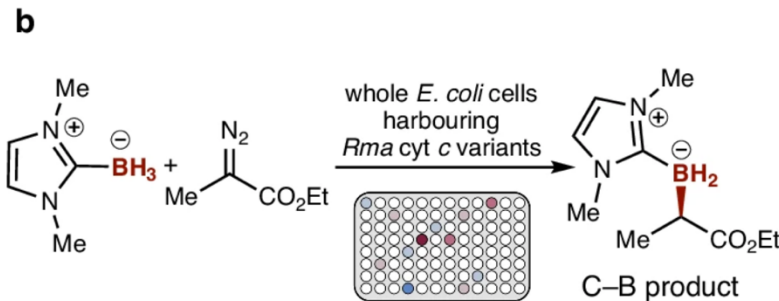


Ding, K., Chin, M., Zhao, Y. et al. "Machine learning-guided co-optimization of fitness and diversity facilitates combinatorial library design in enzyme engineering," *Nature Communications*, 2024

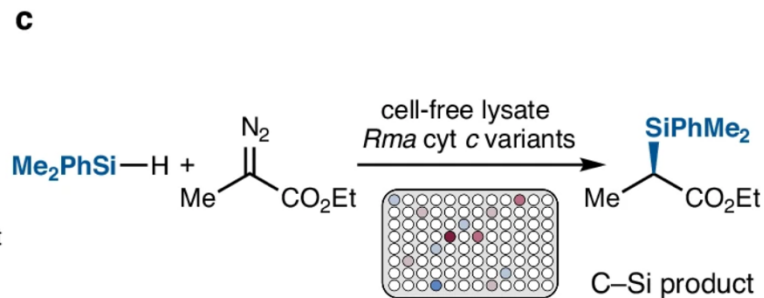
# Using AI to Engineer generalist new-to-nature biocatalysts



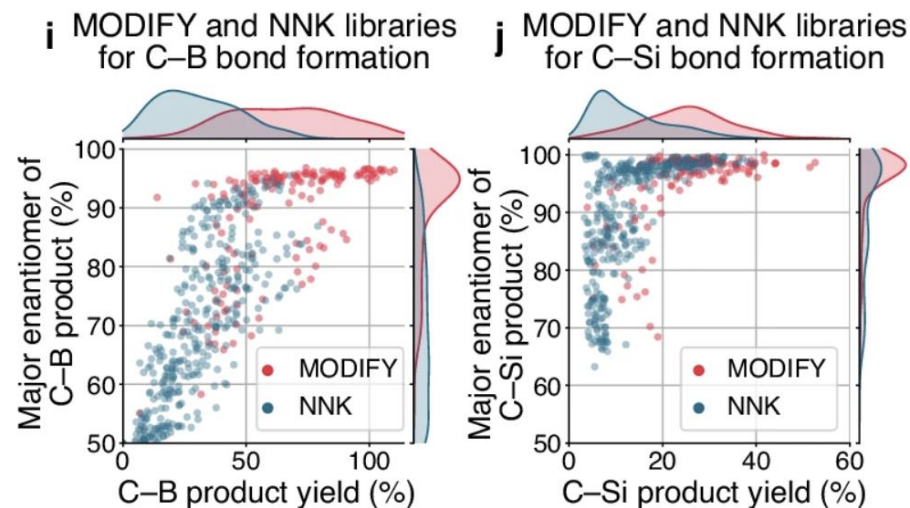
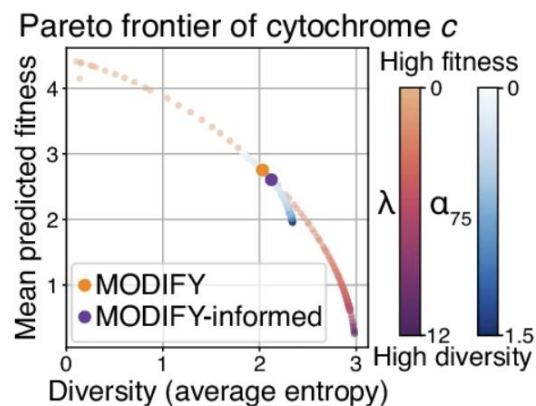
*Rma* cytochrome *c*



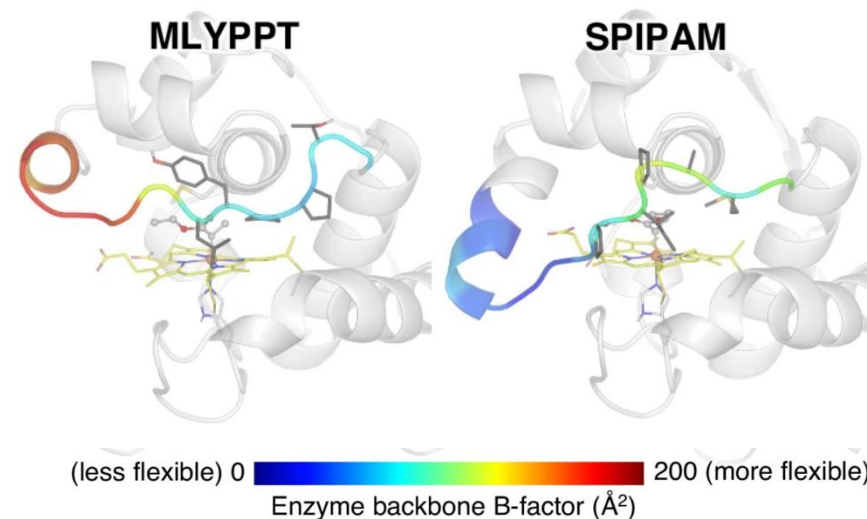
Carbon-boron (C-B) bond formation



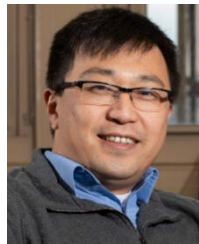
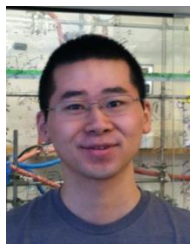
Carbon-silicon (C-Si) bond formation



Improved two objectives: yield and selectivity



GenAI designed novel proteins 6-mutation away from human-designed proteins



Yang Yang (UCSB Chemistry)

Peng Liu (U of Pitt Chemistry)

# How to design functional proteins?

- Learning from natural protein evolution using GenAI

$$p(\text{functional protein}) \approx p(\text{natural protein})$$

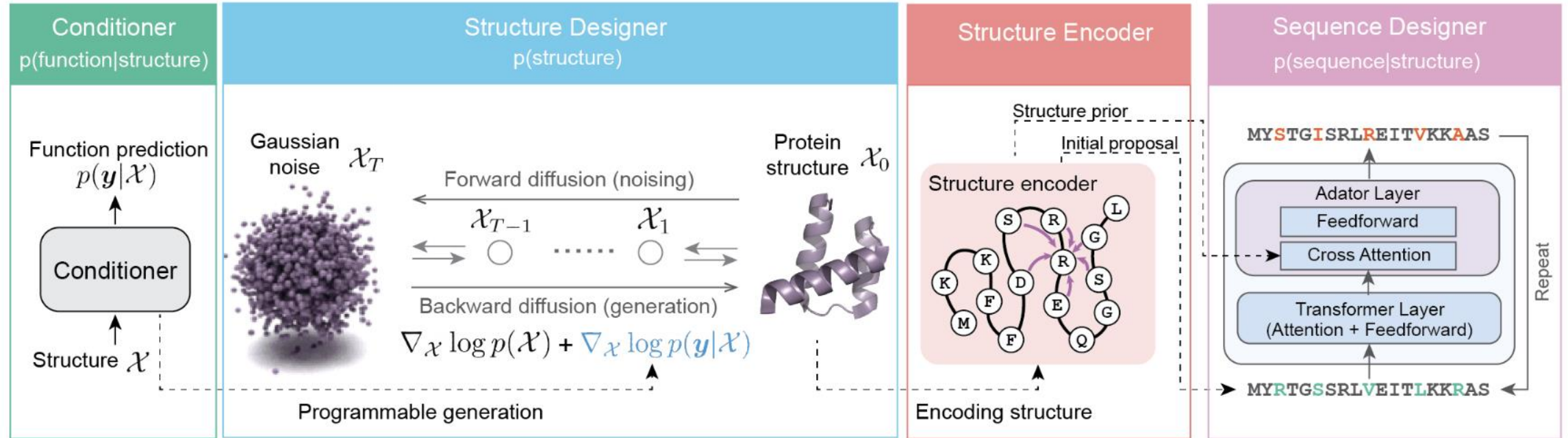
- Function-guided protein design with GenAI

$$p(\text{protein}|\text{function}) \propto p(\text{protein}) \times p(\text{function}|\text{protein})$$

↓  
Protein sequence/structure  
GenAI model

↓  
Protein function  
prediction model

# Programmable multi-modal protein design



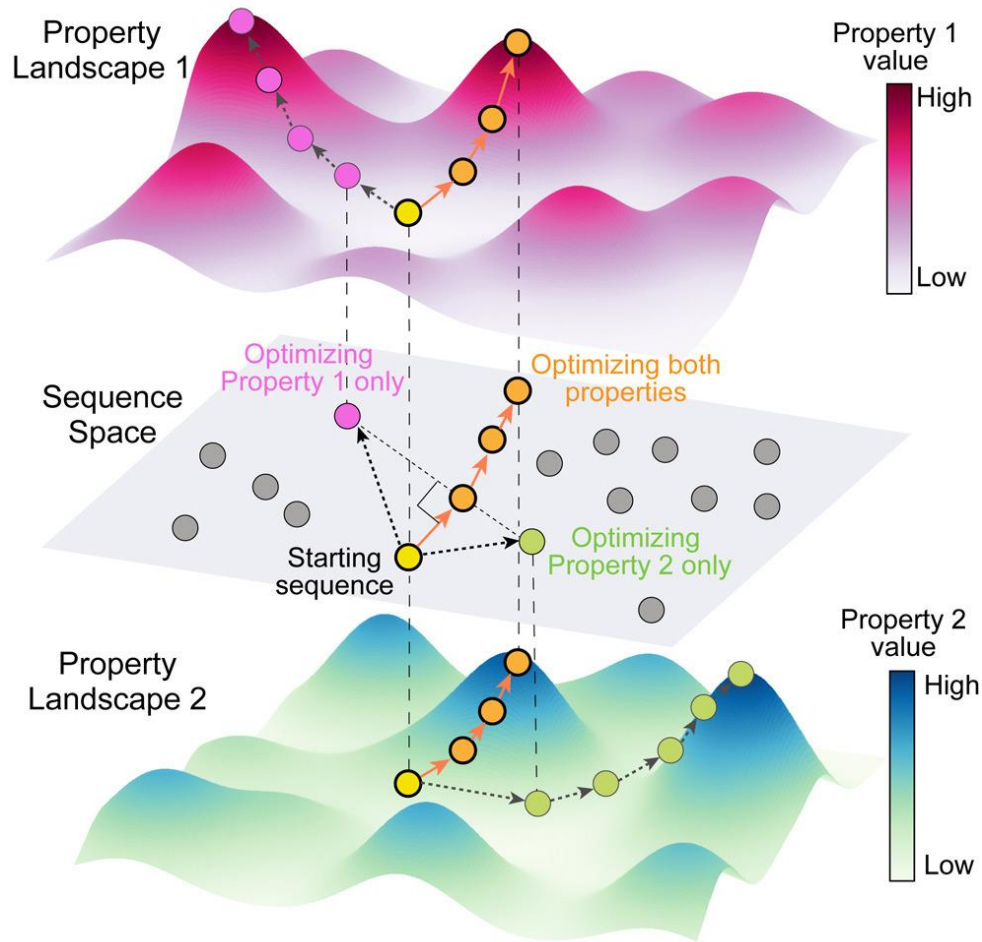
Protein function  
prediction model

Protein structure generative model

Protein sequence  
generative model

# Multi-objective protein design

## Recent results

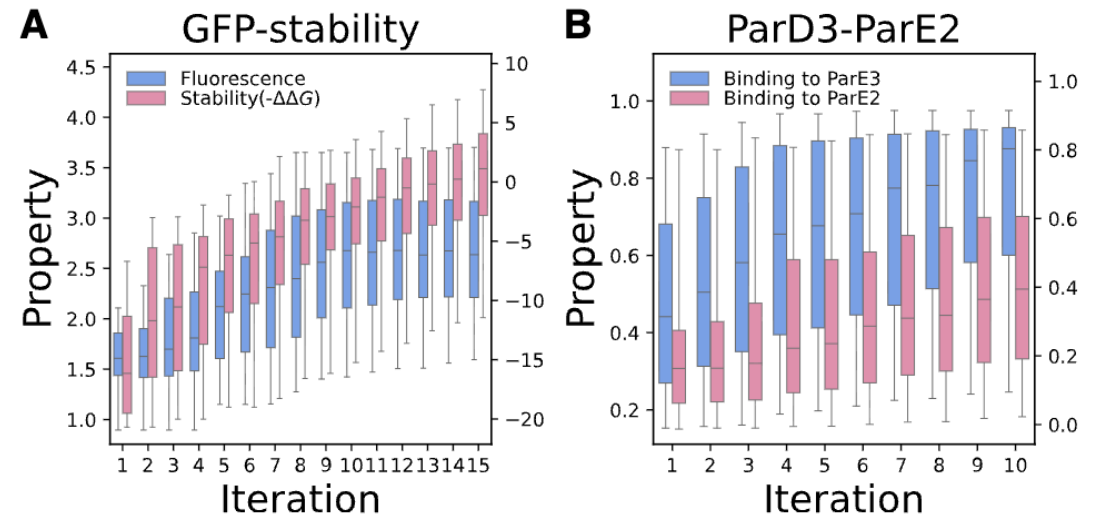


**iScience** **CellPress**  
OPEN ACCESS

**Article**  
**Pareto-optimal sampling for multi-objective protein sequence design**

Jiaqi Luo,<sup>1</sup> Kerr Ding,<sup>1</sup> and Yunan Luo<sup>1,2,\*</sup>

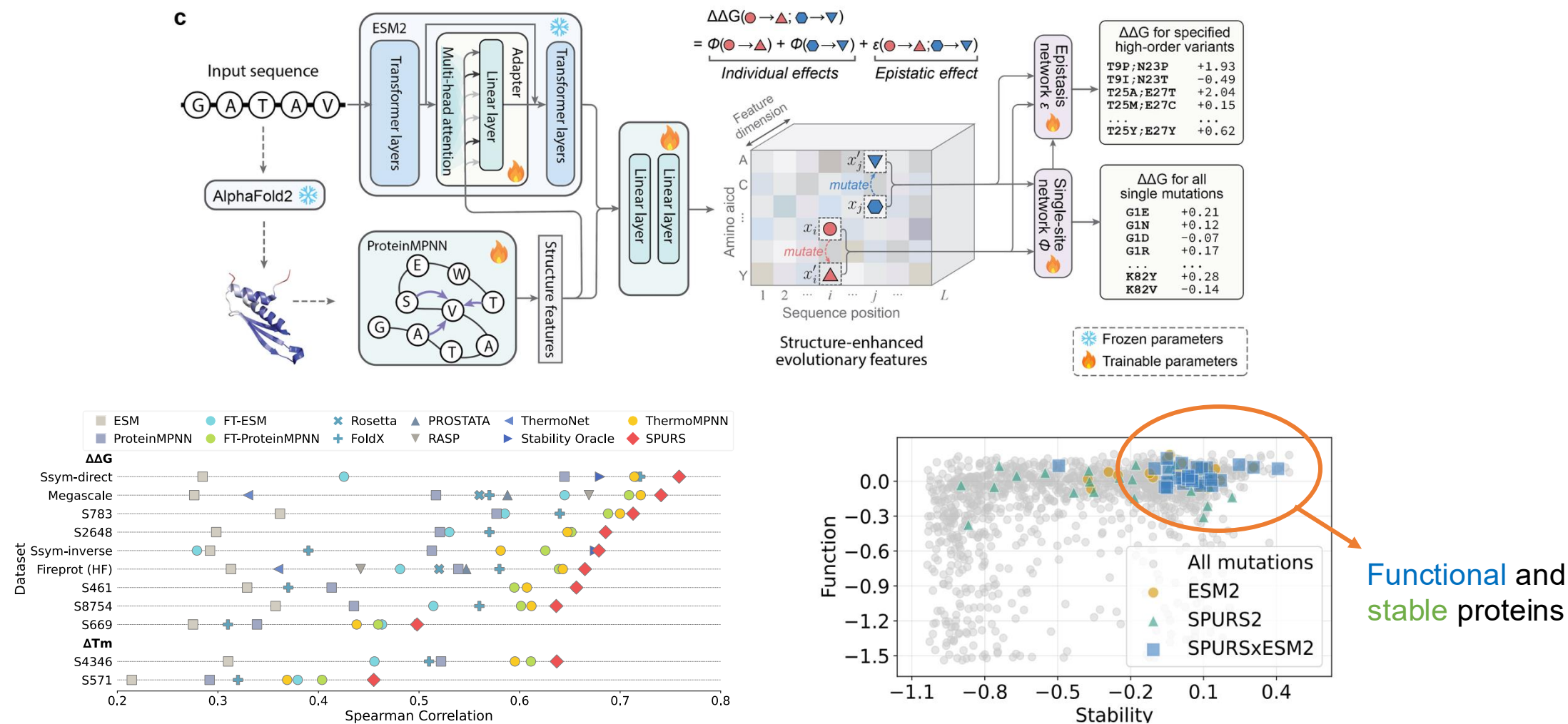
<sup>1</sup>School of Computational Science and Engineering, Georgia Institute of Technology, Atlanta, GA 30308, USA  
<sup>2</sup>Lead contact  
\*Correspondence: [yunan@gatech.edu](mailto:yunan@gatech.edu)  
<https://doi.org/10.1016/j.isci.2025.112119>





# Multi-modal, multi-objective protein design

## Preliminary results



# Conclusion

- Protein generative AI models capture evolutionary patterns of functional proteins
- Function-guided protein design improves hits rate
- Enables the design of novel, diverse functional proteins

## Acknowledgements

- **Group:** Kerr Ding, Ziang Li, Jiaqi Luo, Tony Tu, Shitong Dai
- **Collaborators:** Huimin Zhao, Yang Yang, Peng Liu, Tianhao Yu, Junming Zhao, Liupeng Zhao, Michael Chin, Yunlong Zhao, Wei Huang, Binh Khanh Mai, Huanan Wang, et al.
- **Fundings:** GaTech IDEaS x Microsoft CloudHub Seed Grant, NIH MIRA (R35GM150890), NSF CAREER (2442063)



Georgia Tech  
Institute for Data  
Engineering and Science



National Institutes  
of Health

