

NYGC Events

cover slide -- Saved to my Mail

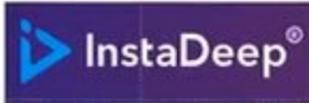
Home Insert Draw Design Transitions Animations Slide Show Record Review View

Record Comments Share

Picture Text Box Arrange Quick Stories Security Add-ins Designer




Thank you to our sponsors


And helpers!

Sarah Curtiss , Kristen Weatherley
 Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra
 Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott
 Adamson, Trevor Christensen, Yijie Kang

mlcb.org for schedule

Click to add notes

Slide 1 of 1 English (United States) Accessibility: Investigate

Notes Comments

Zoom Meeting

0:23 AM ...

applying

0:37 AM ...

zOOM

cover slide -- Saved to my Mac --

Search (Cmd + Ctrl + U)

Record Comments Share

Home Insert Draw Design Transitions Animations Slide Show Record Review View

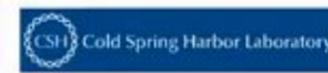
Play from Start Play from Current Slide Presenter View Custom Show Rehearse with Coach Set Up Slide Show Hide Slide Rehearse Timings Record Play Narrations Show Media Controls Subtitle Settings

1




Thank you to our sponsors





And helpers!

Sarah Curtiss , Kristen Weatherley
 Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra
 Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott
 Adamson, Trevor Christensen, Yijie Kang

mlcb.org for schedule

Click to add notes

Slide 1 of 1 English (United States) Accessibility: Investigate

Notes Comments 120%

Zoom Meeting

Screen Act

0:23 PM png 2024-09-11 8:53 AM png

All Bookmarks

MLCB

1:23 AM ...

applying

1:37 AM ...



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra
Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott
Adamson, Trevor Christensen, Yijie Kang

NYGC Events



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley

Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra
Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott
Adamson, Trevor Christensen, Yijie Kang

NYGC Events

Understanding and designing regulatory DNA using machine learning

Jacob Schreiber

Genomics and Computational Biology
UMass Chan Medical School



jmschreiber91



@jmschrei

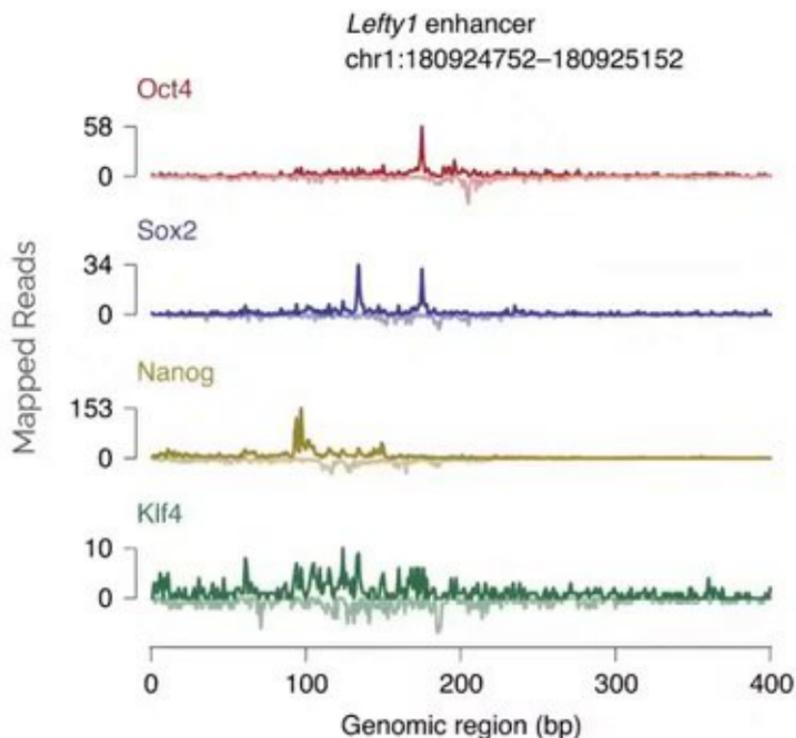


@jmschreiber91

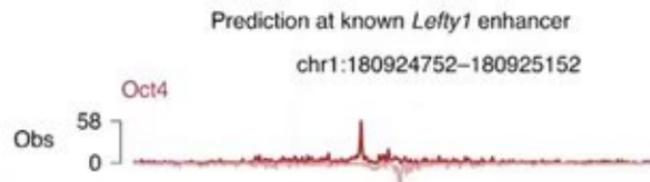
The first assembly of a human reference genome left us with ...

GAAAACAAGTACCATTTTCAACATTAACCTGATGCCTTGGCTTCATGCTATAATGCCATGTTGTGTTTCACTATAACCTCAGAGTGA
ATGAAAGAGGAAAATGGAGCTAGTTGAAATTTCTGCCTAAACTAGCCAGATTTTGAGACACTAAGTTATCTCAAATCAAGAAATCA
CCCTAATGAGAAATTTCAATAACCTCAGGAATTTAAGGTGCATGCATCCCCACCCCCCTTTTTTTTTTGGAGACGTAGTCCCGCT
CTGTTGCCCAGGCTGGAGTACAGTGGCGGATATCGGCTCACCACAACCTCTGCCTCCAGGTTCAAGGGATTCTCCCGCCTCAGC
TTCCAGAGTAGCTGGGACTACAGACACCCACCACCATGCGTGGCTAATTTTTGTATTTTTAGTAGAGAGGGGGTTTCGCCATGTTG
GCCAGGCTGGTTCAAACCTCTGACTTCAGGTGATCCGCTGCCACGGCTCCCAATTTACTGGGATTACAGGGGTGGGCCACCGC
GCCCCGCTTTTTCTTAATTTTTAAAAATATTAAGTTTTATCCATTCTGTTGAACCATATTCCTGATTTAAAAGTTGGAAACG
TGGTGAACCTAGAAGTATTTGTTGCTGGGTTTGTCTTCAGTTCTGTTGCTCGGTTTTCTAGTTCCACCTAGTCTGGGTTACTC
TGCAGCTACTTTTGCATTACAATGGCCTTGGTGAGACTGGTAGACGGGATTAACGAGAATTCACAAGGGTGGGTGAGTAGGGGGT
GTGCCCCCAGGAGGGGTGGGTCTAAGGTGATAGAGCCTTCATTATAAATCTAGAGACTCCAGGATTTTAACTGTTCTGCTGGACTG
AGCTGGTTGCCTCATGTTATTATGCAGGCAACTCACTTTATCCCAATTTCTTGATACTTTTCCTTCTGGAGGTCTATTTCTCTAA
CATCTTCCAGAAAAGTCTTAAAGCTGCCTAACCTTTTTCCAGTCCACCTCTTAAATTTTTCTCCTCTTCTCTATACTAACA
TGAGTGTGGATCCAGCTTGTCCCCAAAGCTTGCCTTGCTTTGAAGCATCCGACTGTAAAGAATCTTACCTATGCCTGTGATTTGT
GGGCTGAAGAAAATATCCATCCTTGCAAATGTCTTCTGCTGAGATGCCTCACACGGAGACTGGTAAGAAAAGAAATTTATCCTTG
AAAGGCCAAGTTCCTTAAAGGAAAAGAGAGAAGGAGAGAGGGTTAAGGGATCATTTCCCTCTTGAGCAATGATGGACCATTACTAT
AAAGAAGTGTATTATCAACTAATCCTCTGGAAACCCCTTTTTCCATTATAACTGGTGGCACCTGCCCTTTGAACTATGTCCCAG
GTCTCAGGAGTGTGCATTGAGTTGAAGGACACAGAATTCGGCAGTTGAACAGTGTGCAGTAAGTTTGAACCTATGGGCTTAGGC
ATGGTGGAAACAAAAATGTATCGTTATAGTTAAATGAAGGTGATGTGTACATCTTACATAGTGTGGACACATGTGAATAAATAG
CAGATTTATTGCTAATTAGCCAGAAGACCTAACGTCATAGCTCAGGGATGAGCATGATTTTGTGTTTGGCAAAAATGGCATGGCAA
TGACGATGAGATTTCTGTAATACATAATTTGGGTAATTTCTTCTATGTGAGTAACGGCTGTCTCTTCTCCATTCTCTGGGTTTGTG

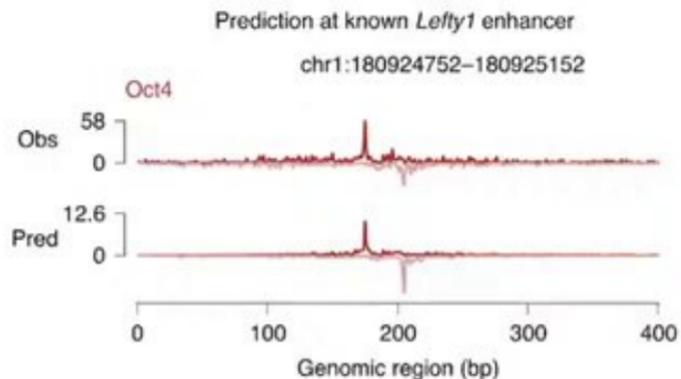
These assemblies were augmented with high-throughput assays measuring cell type-specific signals



Sequence-based machine learning models can be trained to make **predictions**...



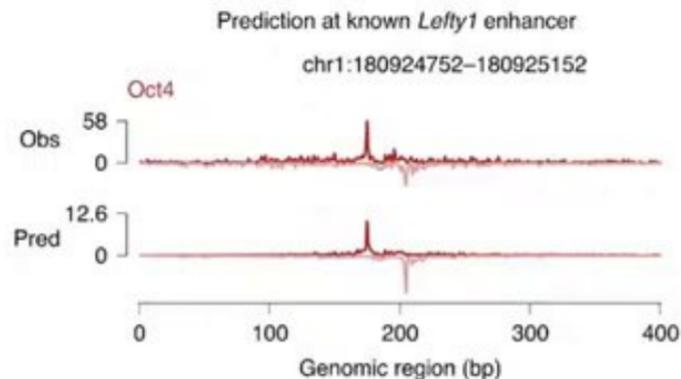
Sequence-based machine learning models can be trained to make **predictions**...



Sequence-based machine learning models can be trained to make **predictions...** and highlight nucleotides driving activity with **attributions**



JASPAR (Pou5f1-Sox2, MA0142.1)



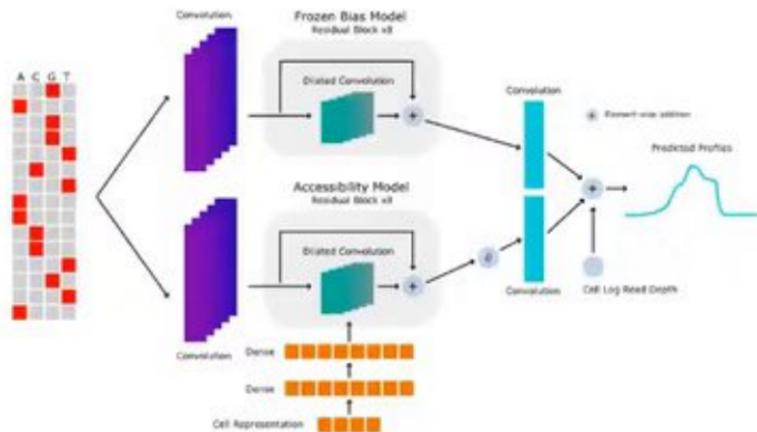
The same idea can then be used to **design** sequences with desired properties



Desired Output



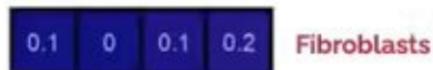
Understanding *cis*-regulatory elements at single-cell resolution with DragoNNFruit



Improvements in single-cell tech allow us to investigate continuous biological processes such as cellular reprogramming



Log Expression (tp10k)



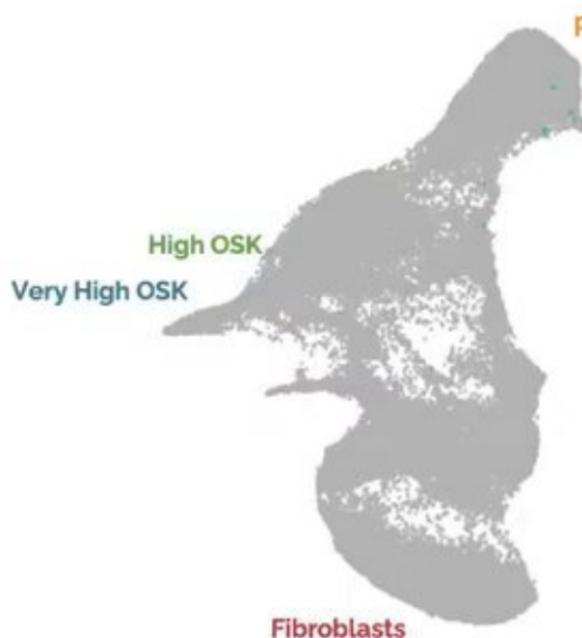
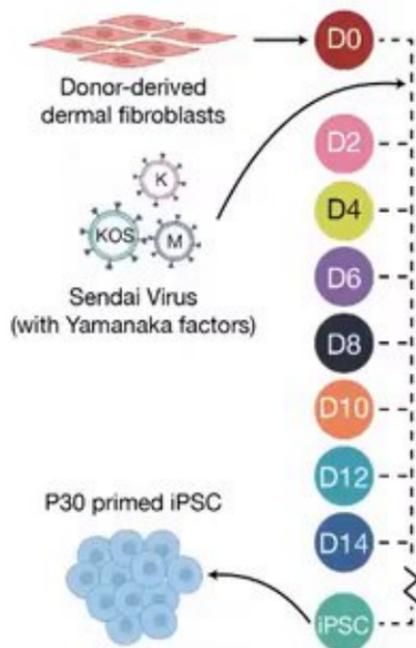
Ö Š Ķ Ą

Gene expression from harmonized parallel scRNA-seq, not an exact mapping

Surag Nair, Mo Ameen



Improvements in single-cell tech allow us to investigate continuous biological processes such as cellular reprogramming



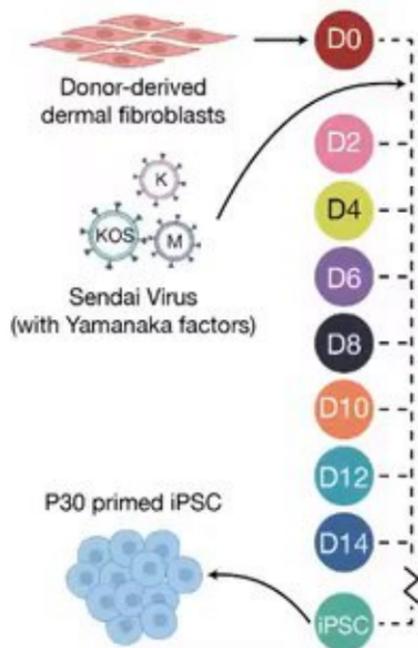
0.1	0	0.1	0.2	Fibroblasts
2.4	0.6	0.5	1	pre-/iPSC
3.2	2.1	2.1	2.6	High OSK
4	3	3.1	2.8	Very High OSK
Ô	Š	K	M	



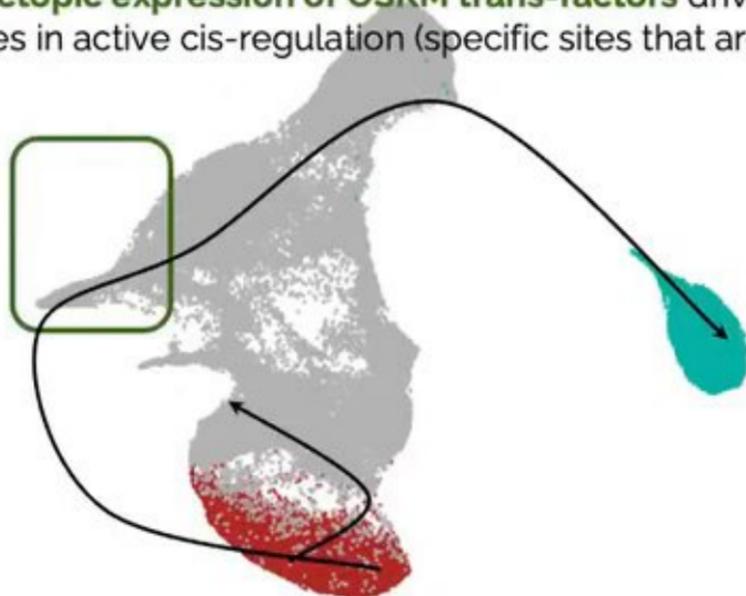
Gene expression from harmonized parallel scRNA-seq, not an exact mapping



Most processes are driven by a complex interplay between cis- (sequence-based) and trans- (protein- and other molecule-based) regulation



High ectopic expression of OSKM trans-factors drive rapid changes in active cis-regulation (specific sites that are bound)

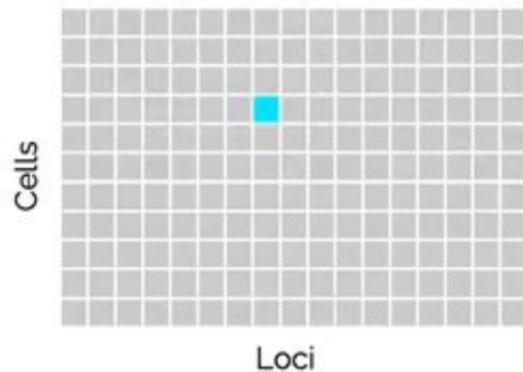


However, most analysis pipelines involve clustering+pseudobulking the data and treating it as a set of discrete states

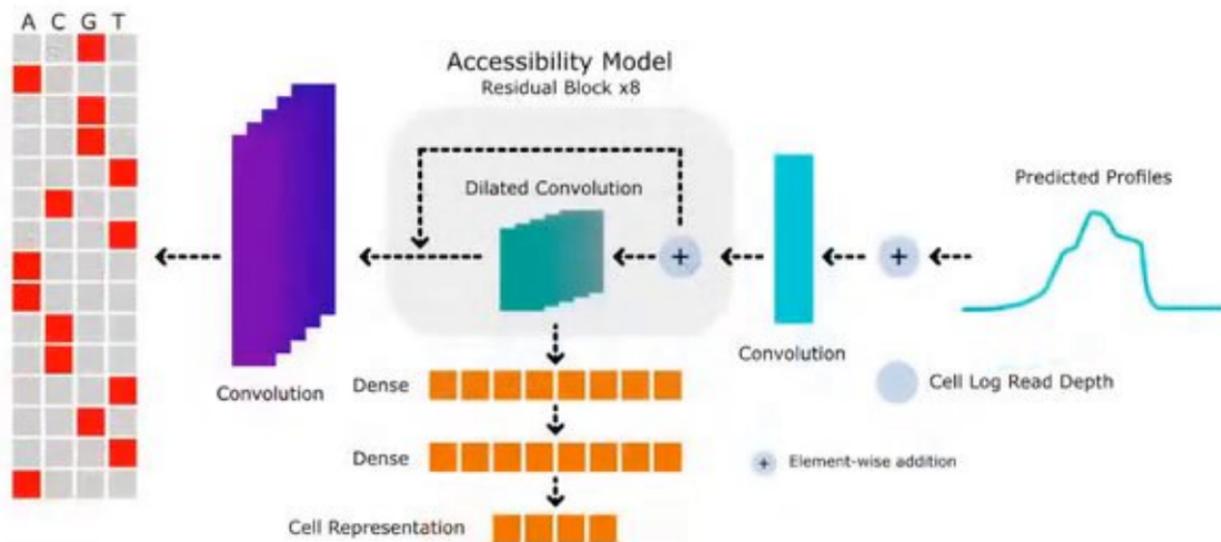


- Fibroblast
- Fibroblast-like
- Fibroblast-like
- Fibroblast-like
- Fibroblast-like
- Keratinocyte-like
- hOSK
- xOSK
- Intermediate
- Partially-reprogrammed
- Intermediate
- Intermediate
- Pre-iPSC
- Pre-iPSC
- iPSC

DragonNnFruit extends sequence-based modeling frameworks to single-cell data



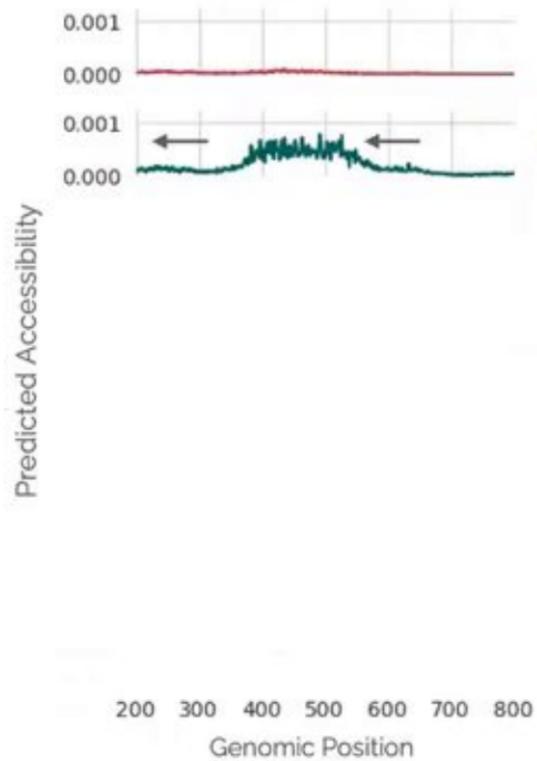
DragonNMFruit extends sequence-based modeling frameworks to single-cell data



DragoNNFruit accurately predicts accessibility as it changes across reprogramming



Some peaks appear to shift coordinates as proteins alter their binding



Some sites exhibit persistent accessibility but change their reasons why

Observed Reads

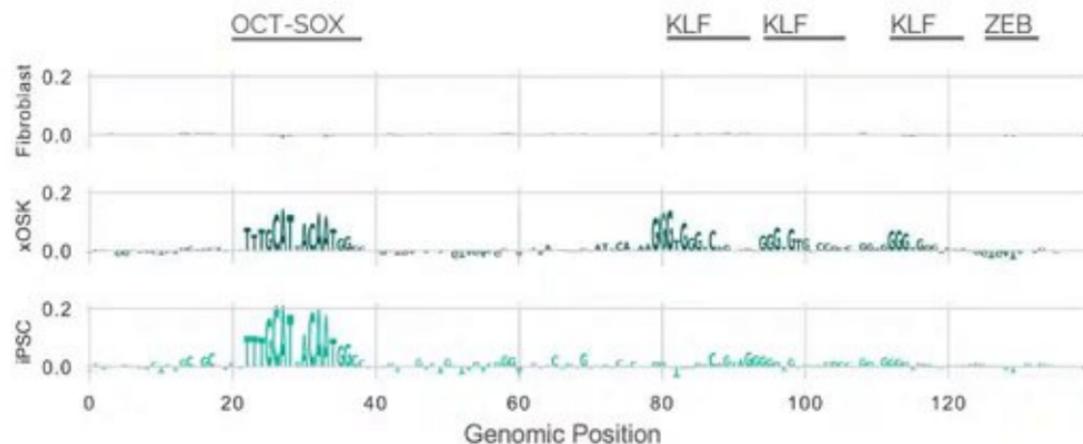
Predicted Accessibility

chr3:197,480,781-197,483,781

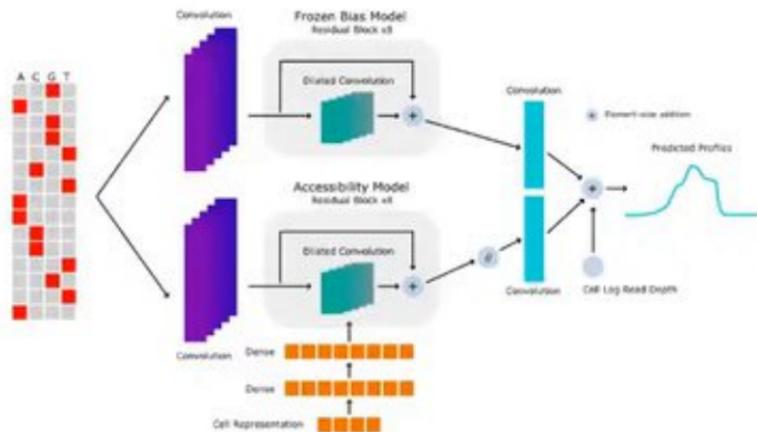


Accessibility at the NANOG promoter is driven by the usual suspects...

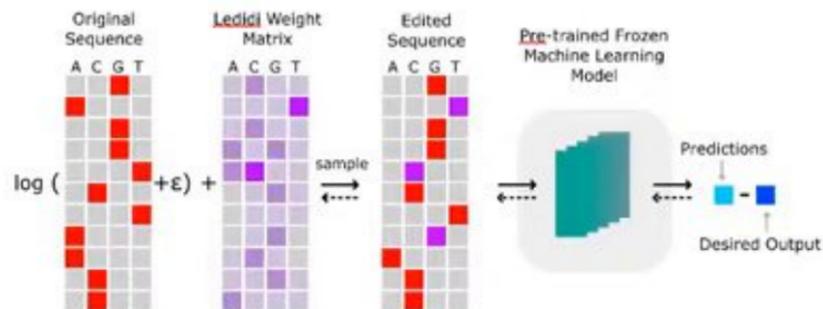
chr12:7,789,855-7,790,035 Nanog Promoter



Understanding *cis*-regulatory elements at single-cell resolution with DragoNNFruit



Designing *cis*-regulatory elements using Ledidi



Inverting the normal machine learning paradigm to **design** sequences with desired properties

Initial Sequence

A	C	G	T
		█	
			█
		█	
		█	
	█		
	█		
			█
		█	
█			
	█		
	█		

Ledidi turns the edit design task into a simple optimization problem

Initial Sequence



Ledidi Weight
Matrix

A C G T

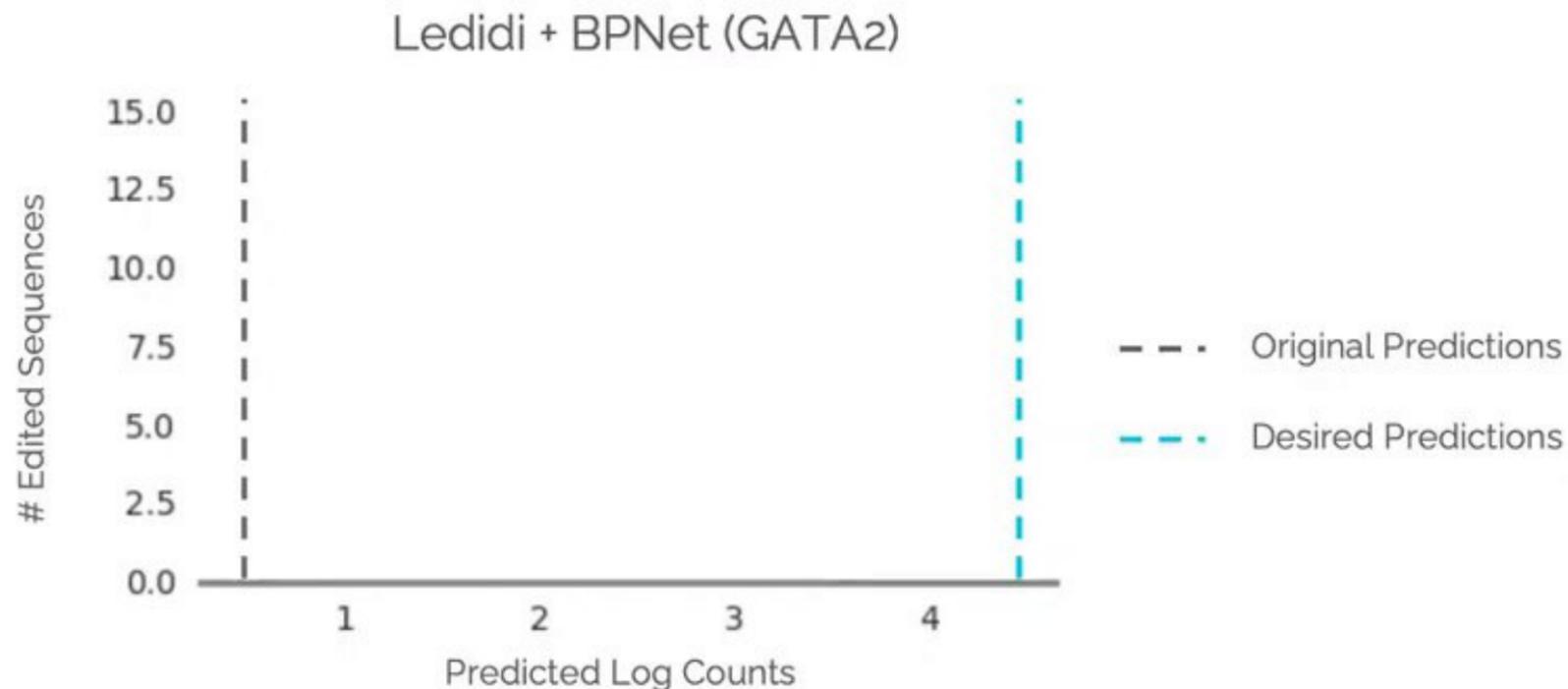
Edited
Sequence

A C G T

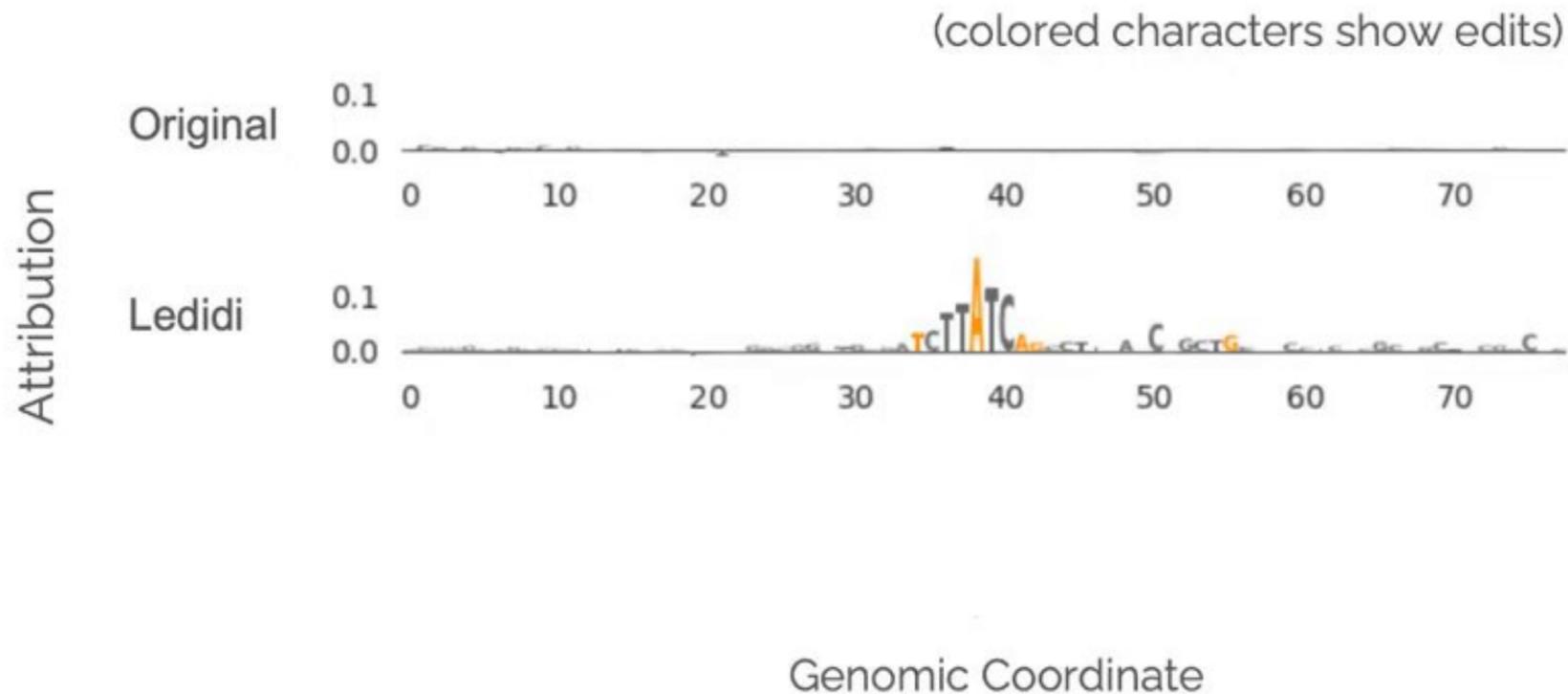
Pre-trained Frozen
Machine Learning
Model



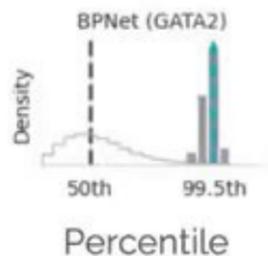
An initial test shows Ledidi can design edits predicted to greatly increase GATA binding



- ▶ The number of edits is significantly lower than expected by precisely targeting regions that are "poised" to become binding sites



Ledidi can design edits yielding desired characteristics in a variety of settings



--- Original Predictions

--- Desired Predictions

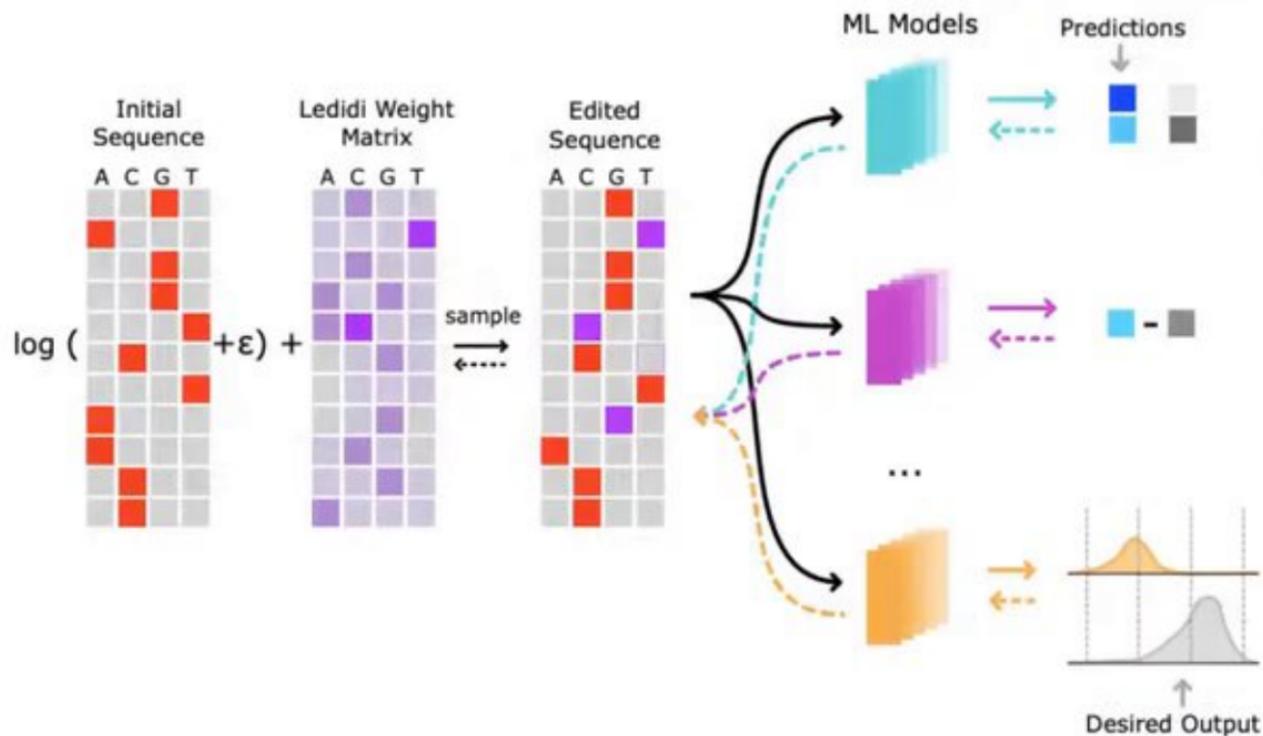


Predicted Genomic Distribution

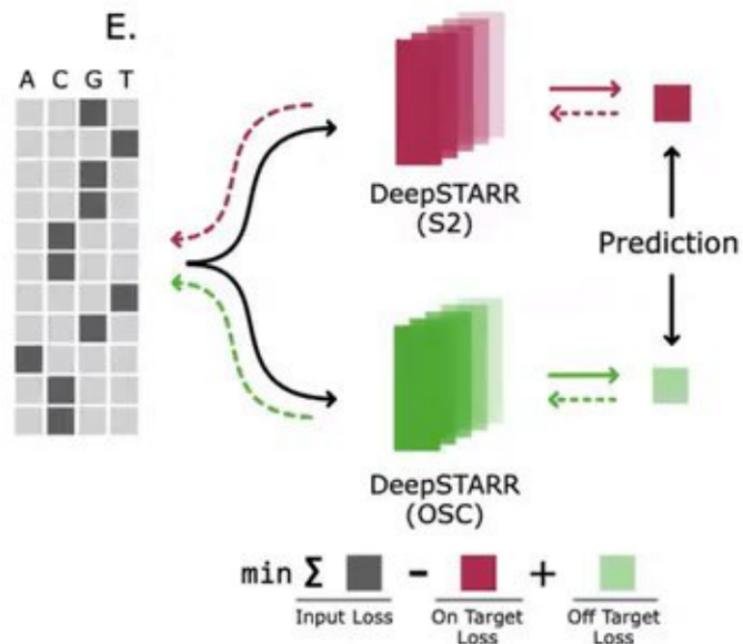


Predictions on Ledidi-edited Sequences

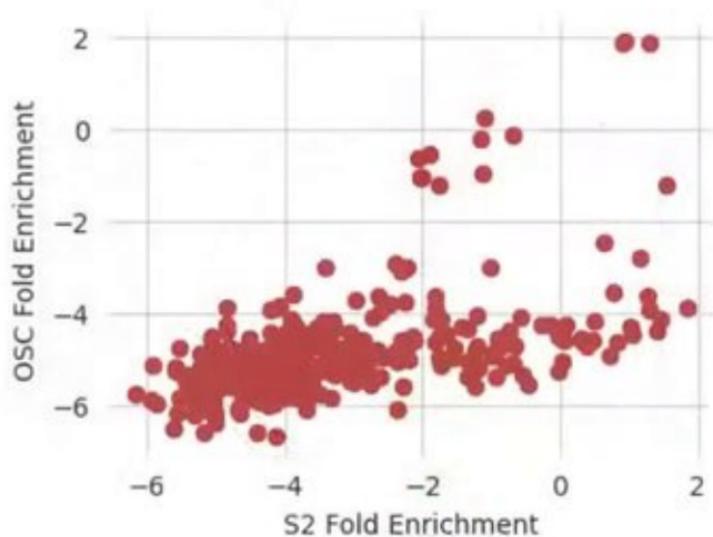
Ledidi allows multiple models to be used simultaneously for design



Ledidi can use multiple DeepSTARR models to design cell type-specific enhancers



Ledidi can use multiple DeepSTARR models to design cell type-specific enhancers



Initial Activity

Target Activity

Non-Peak

S2 Peak ■

OSC Peak ■

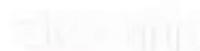
Dual Peak ■ ■

Franziska Lorbeer



Real STARR-seq data!

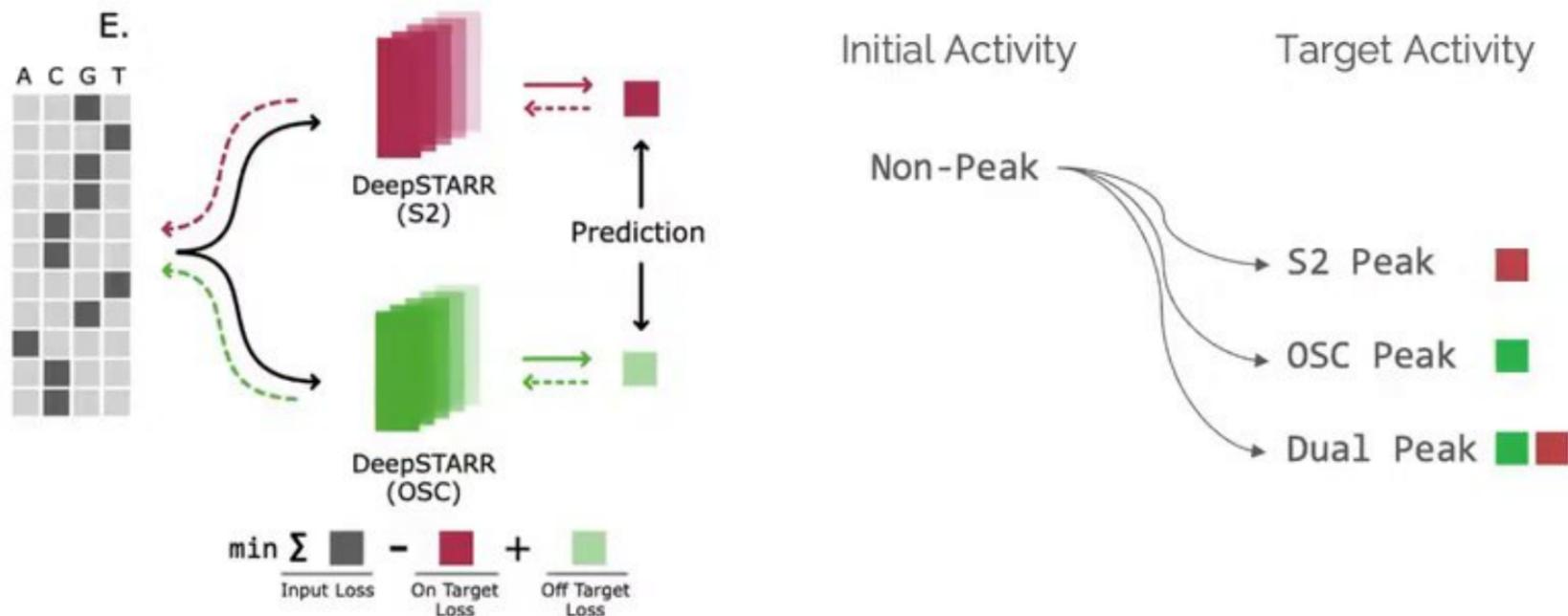
Each dot is a designed regulatory element



These designed enhancers intentionally exhibit a range of activities

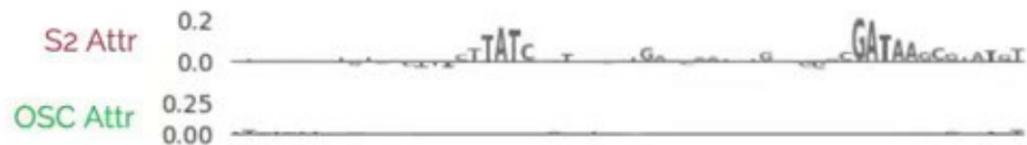


Ledidi can use multiple DeepSTARR models to design cell type-specific enhancers

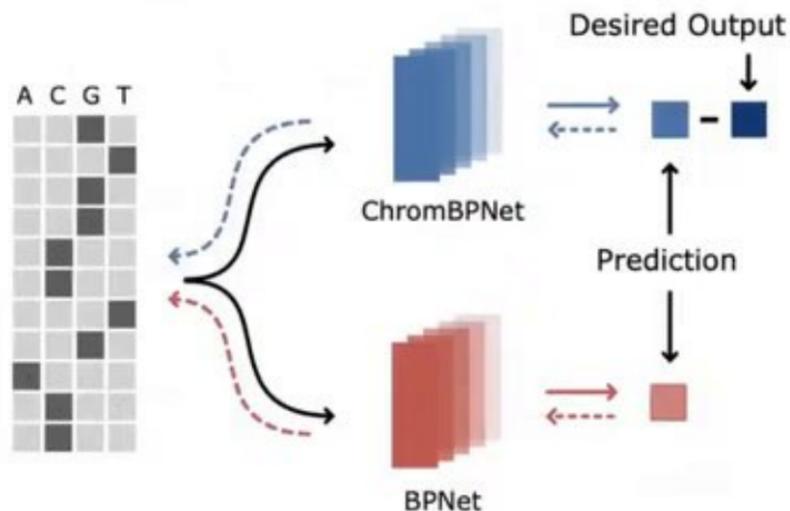


Leadid can turn a weak S2-specific enhancer into an strong OSC-enhancer

Original S2 Enhancer

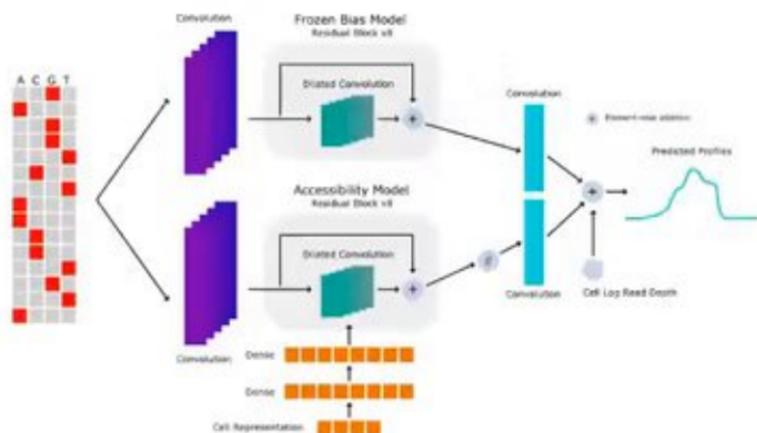


We can use multiple models to programmatically control which TFs drive accessibility

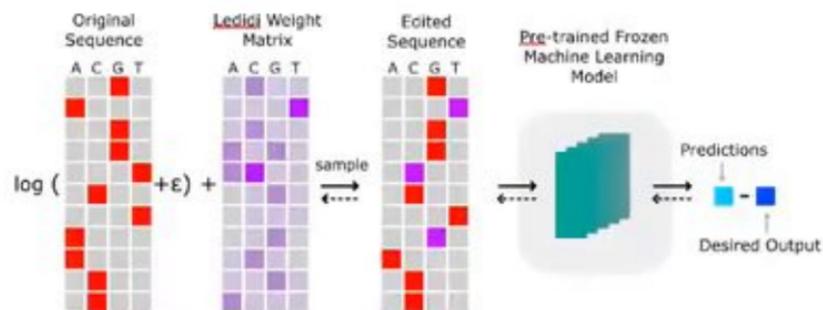


$$\min \sum \text{Input Loss} + \text{ChromBPNet Loss} - \text{BPNet Loss}$$

Understanding *cis*-regulatory elements at single-cell resolution with DragoNNFruit



Designing *cis*-regulatory elements using Ledidi





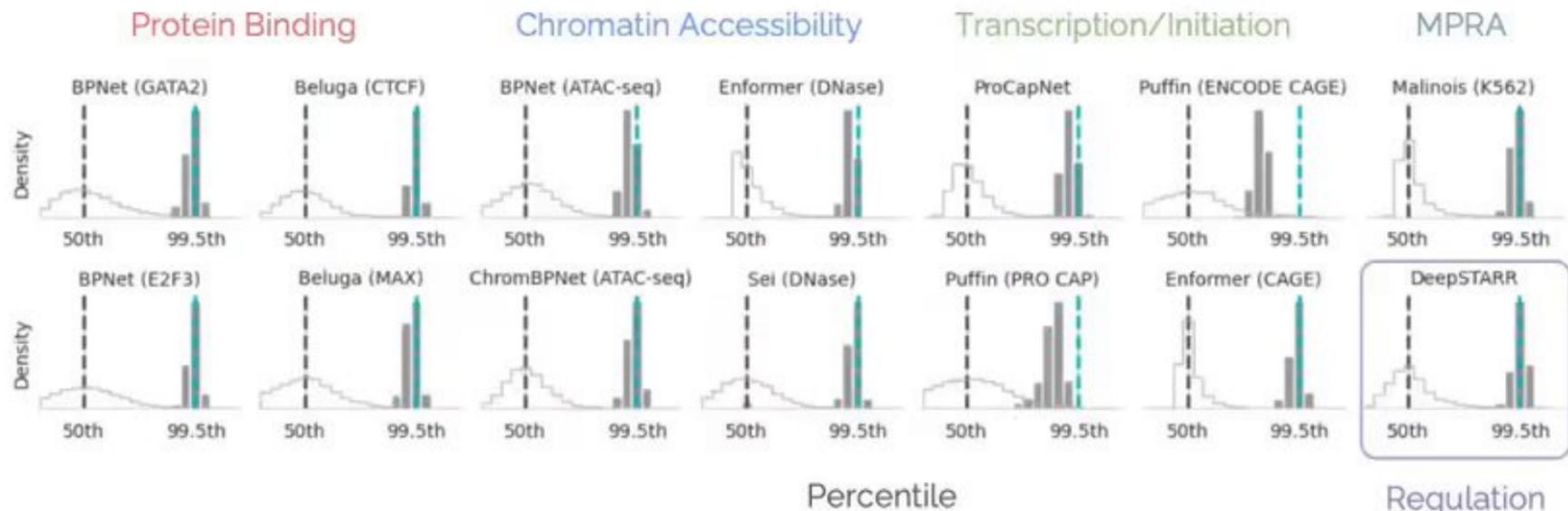
Now recruiting post-docs and software engineers to join me at UMass Chan Medical School!!

If interested, please reach out with your CV and reference contact info!

UMASS CHAN MEDICAL SCHOOL

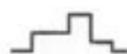


Ledidi can design edits yielding desired characteristics in a variety of settings



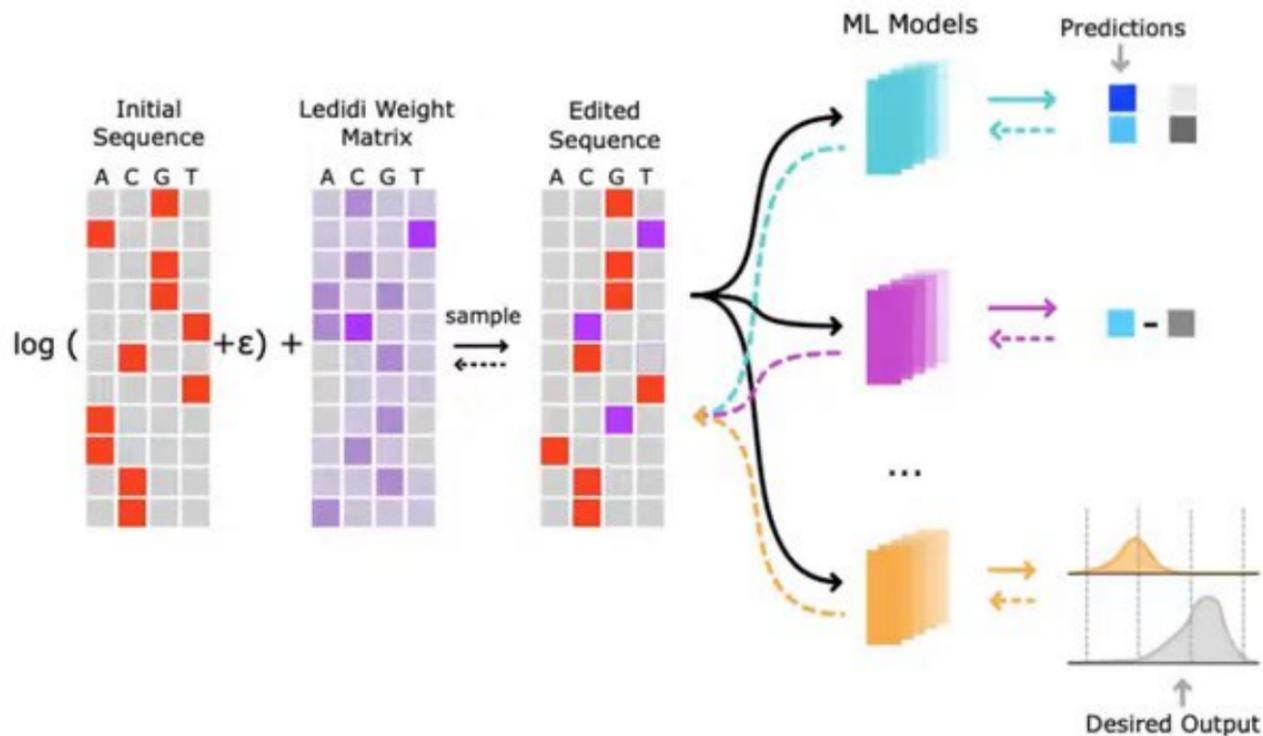
--- Original Predictions

... Desired Predictions

 Predicted Genomic Distribution

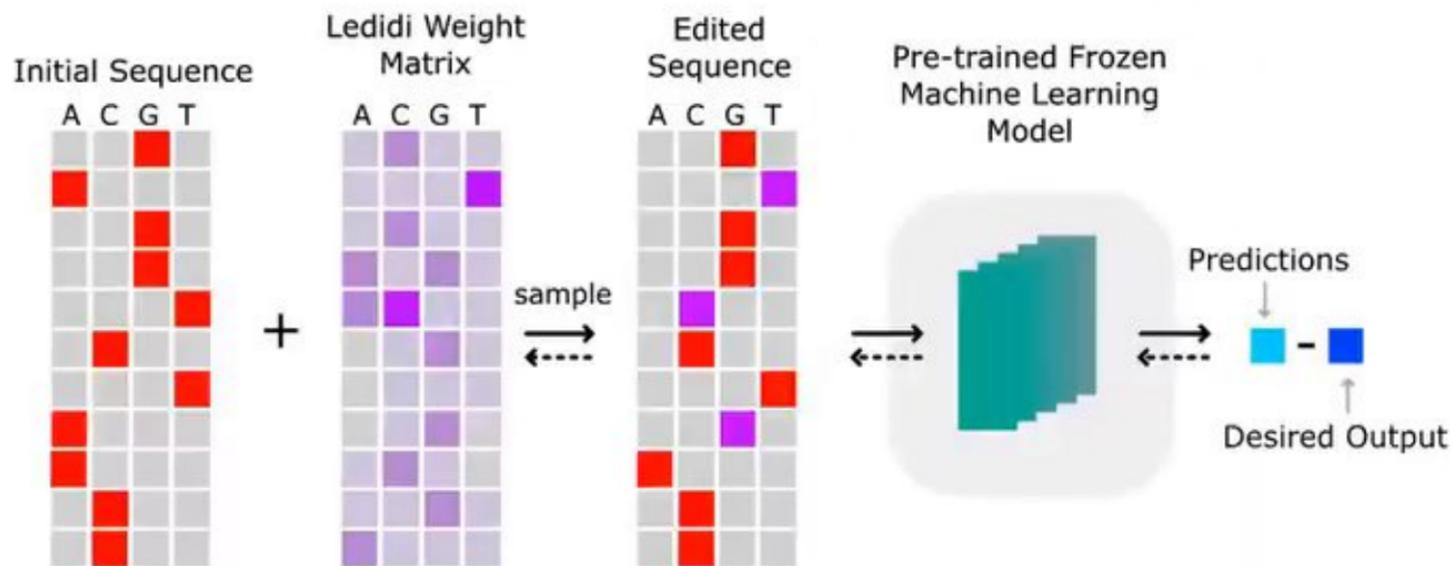
 Predictions on Ledidi-edited Sequences

Ledidi allows multiple models to be used simultaneously for design

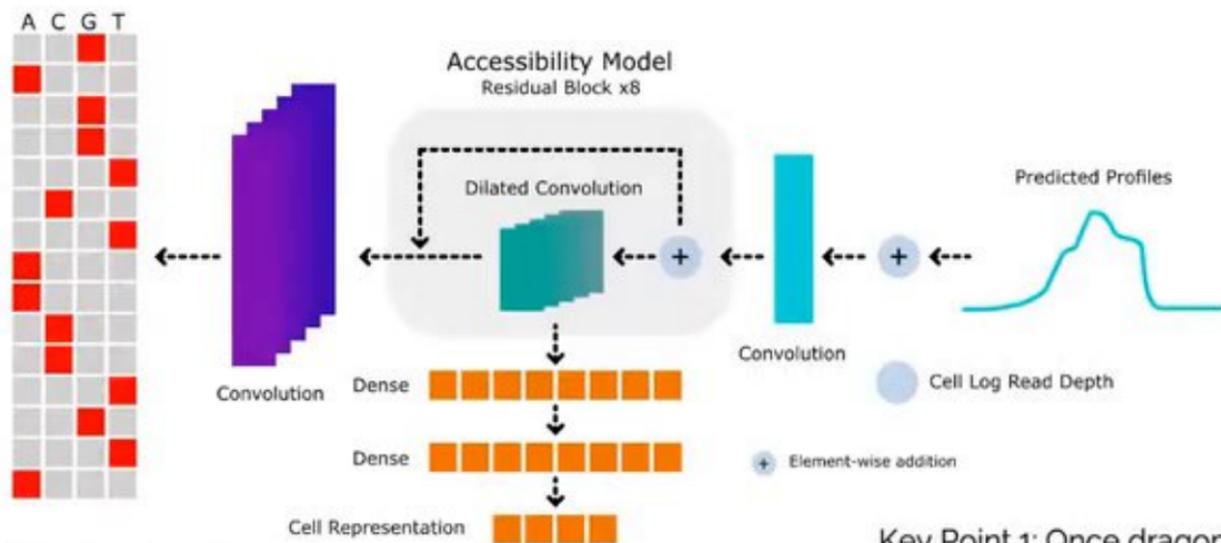


Ledidi turns the edit design task into a simple optimization problem

$$\text{ledidi loss} = \text{input_loss} * \lambda + \text{output_loss}$$



DragonNNFruit extends sequence-based modeling frameworks to single-cell data



Key Point 2: The cell representation can be anything, e.g. scATAC LSI, gene expression, spatial coordinates, protein abundances, batch-effect corrected embeddings..

Key Point 1: Once dragonnfruit is trained it literally produces a model for each cell and do anything a (pseudo-)bulk model can do.

NYGC Events

- 1 Simulation-free Structure Learning for Stochastic Dynamics
- 2 Motivation
- 3 Importance of Dynamics in Structure Learning
- 4
- 5 Problem Setup
- 6 Conclusion

Simulation-free Structure Learning for Stochastic Dynamics

Noah El Rimawi-Fine*, Adam Stecklov*, Lucas Nelson*, Alexander Tong, Mathieu Blanchette, Stephen Y. Zhang, Lazar Atanackovic



Chrome File Edit View History Bookmarks Profiles Tab Window Help

about:blank

00:00:01 Pause Reset

AUDIENCE TOOLS SPEAKER NOTES

Slide 1

Simulation-free Structure Learning for Stochastic Dynamics

Start 12:56. Introduce yourself+Adam

Mila McGill BROAD

ERIC AND WENDY SCHMIDT CENTER AT BROAD INSTITUTE

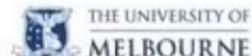
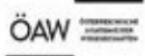
Next

Slide 1 of 41

Start 12:56. Introduce yourself+Adam

Structure Learning for Stochastic Dynamics

, Lucas Nelson*, Alexander Tong,
Lazar Atanackovic



Simulation-free Structure Learning for Stochastic Dynamics

Noah El Rimawi-Fine*, Adam Stecklov*, Lucas Nelson*, Alexander Tong, Mathieu Blanchette, Stephen Y. Zhang, Lazar Atanackovic

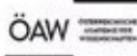


McGill

BROAD
INSTITUTE



ERIC AND WENDY
SCHMIDT CENTER
AT BROAD INSTITUTE



AITHYRA



THE UNIVERSITY OF
MELBOURNE

Dragon
data

Key Point 2:
can be anyth
gene expres
protein abur
corrected er

Simulation-free Structure Learning for Stochastic Dynamical Systems

Noah El Rimawi-Fine*, Adam Stecklov*, Lucas Nelson*, Alexander Tononi, Mathieu Blanchette, Stephen Y. Zhang, Lazar Atanackovic



Simulation-free Structure Learning for Stochastic Dyn

Noah El Rimawi-Fine*, Adam Stecklov*, Lucas Nelson*, Alexar
Mathieu Blanchette, Stephen Y. Zhang, Lazar Atanackovic



protein abund
corrected er



Simulation-free Structure Learning for Stochastic Dyn

Noah El Rimawi-Fine*, Adam Stecklov*, Lucas Nelson*, Alexar
Mathieu Blanchette, Stephen Y. Zhang, Lazar Atanackovic



McGill



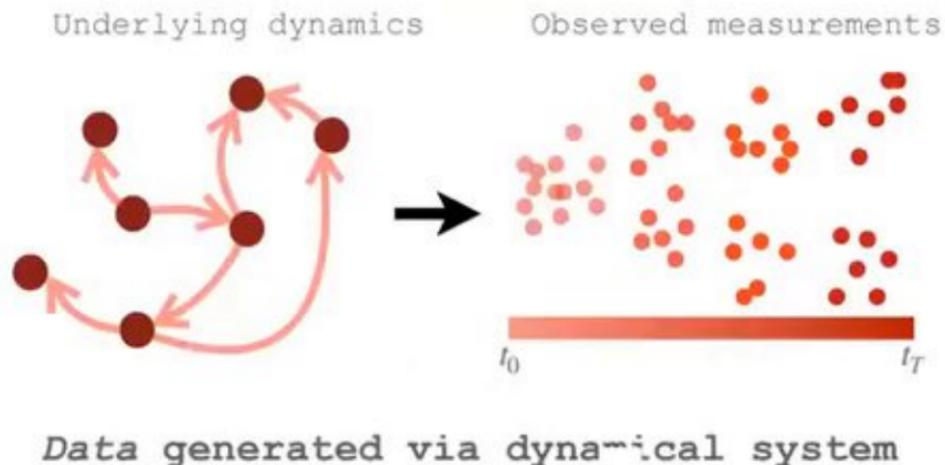
BROAD
INSTITUTE



ES
A

Motivation

In many scientific problems we want to model the dynamic evolution (e.g. cells) over time



Inference of Dynamics & Structure

How can we address both problems jointly?

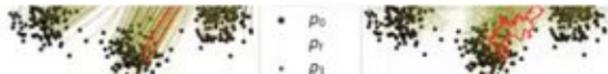
with the ones that do having limitations

Challenges

There are still significant challenges for addressing the problems of model dynamics and inferring the network structure of these systems

Simulation Free Training (flow + s

SF2M (Tong et al, 2024) infers dynamics from noisy snapshots by learning population flow



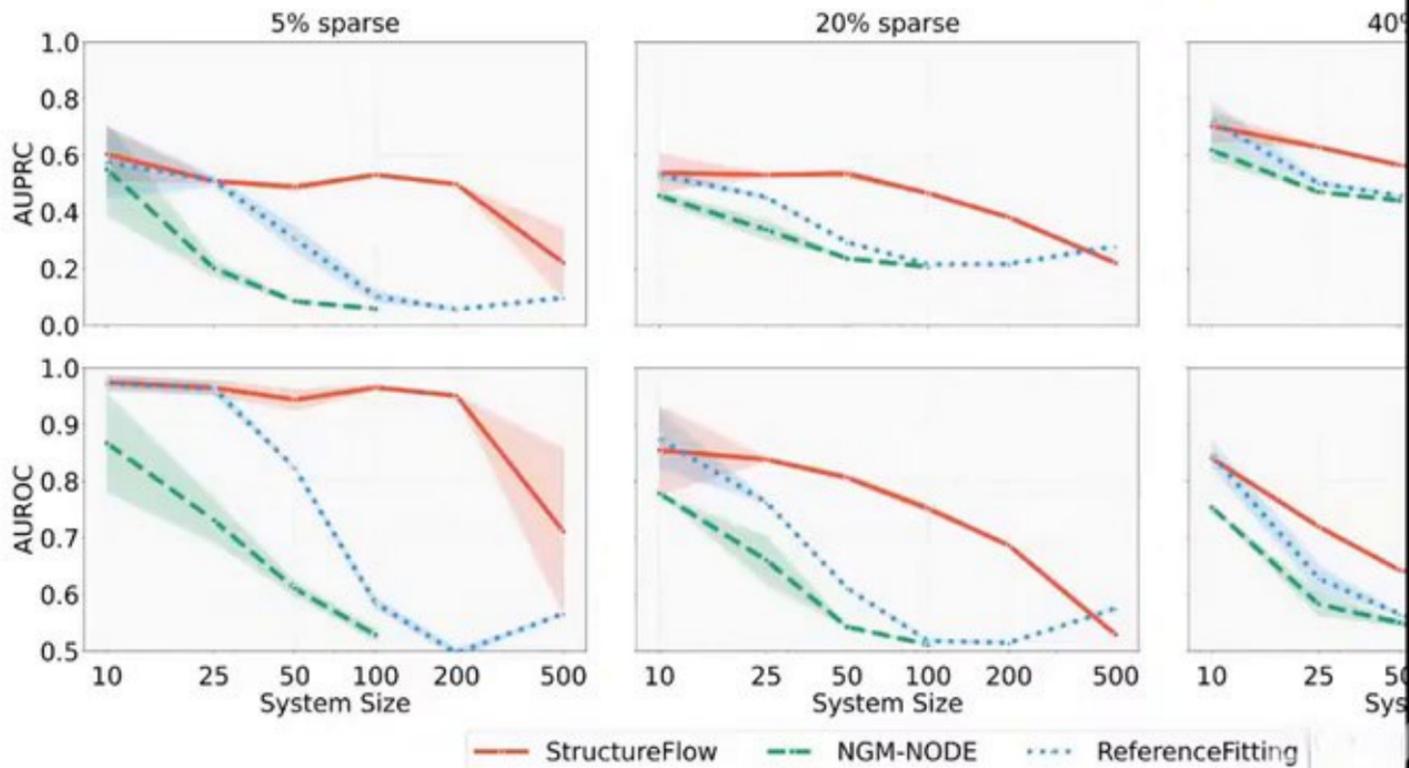
StructureFlow (our method)

However, Tong et al. doesn't exactly address this joint learning task, or mode

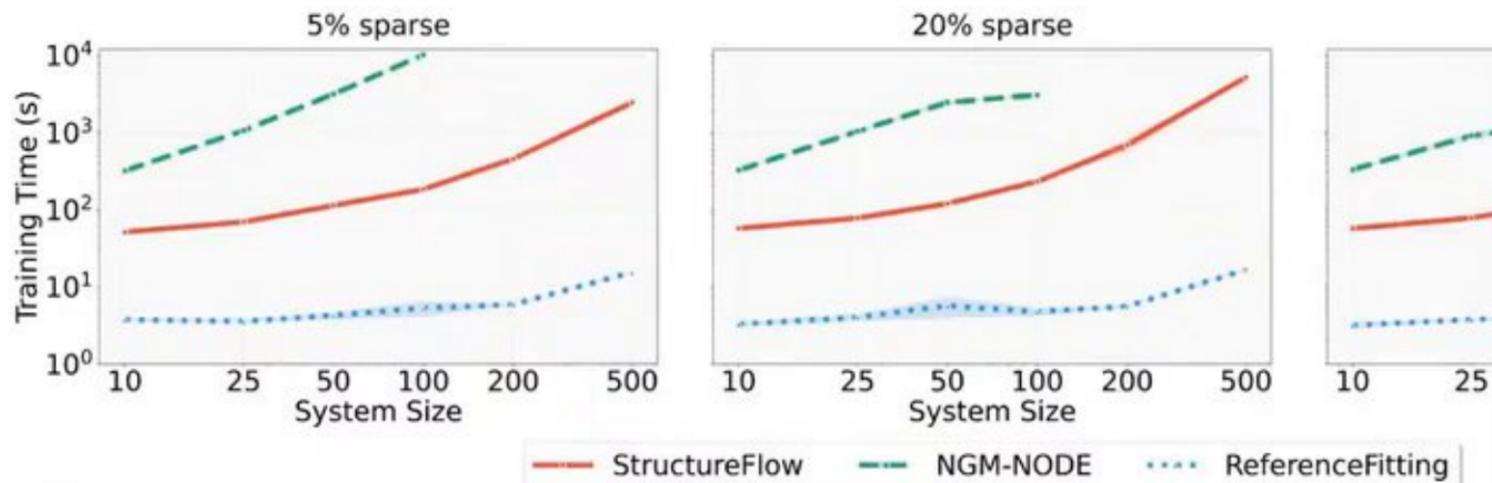
$$\mathcal{L}_{[\text{SF}]^2\text{M}}(\theta) = \mathbb{E}_{t, \mathbf{z}, \mathbf{x}} [\|\mathbf{v}_t^\theta(\mathbf{x}) - \mathbf{v}_t^\circ(\mathbf{x}|\mathbf{z})\|^2 + \lambda(t)\|\mathbf{s}_t^\theta(\mathbf{x}) - \nabla$$

How well does this work?

StructureFlow scales better to larger systems

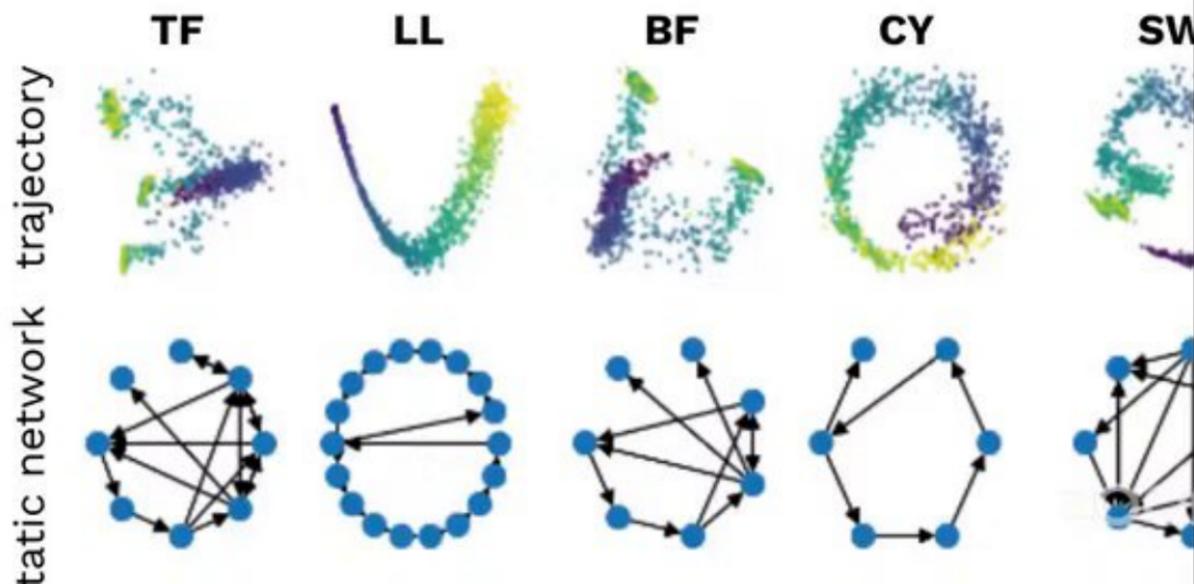


StructureFlow can be trained efficiently



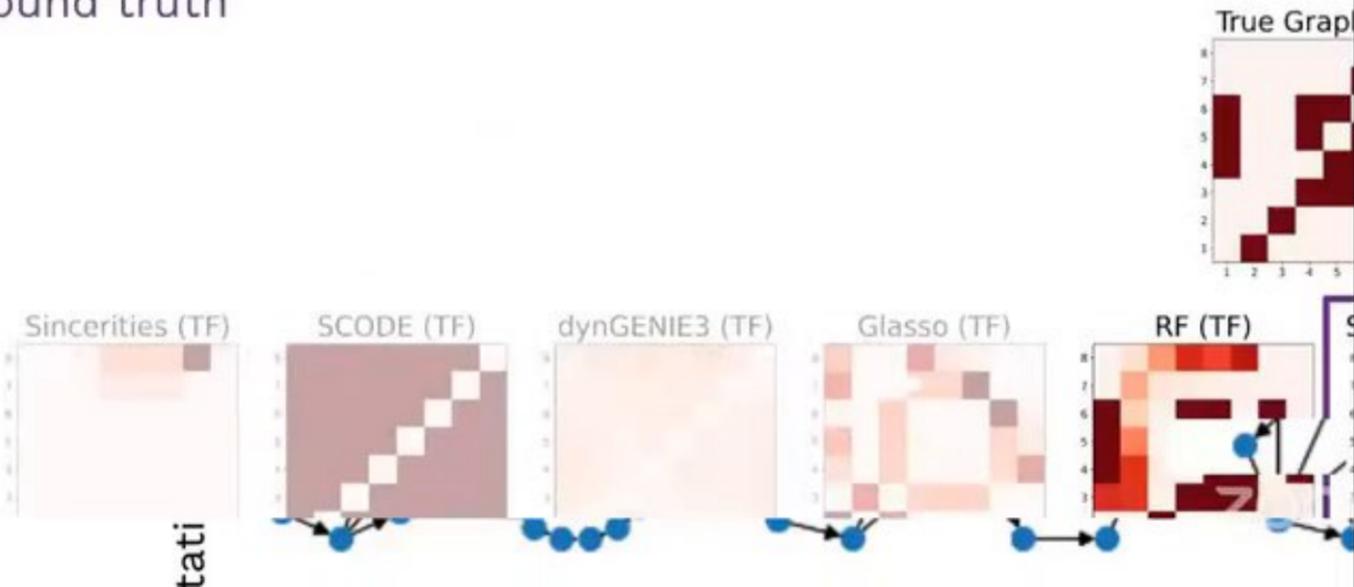
II - Simulated Biological Systems

Joint Trajectory Inference and GRN inference on small, curated, systems, each representing common gene regulatory mechanics
Data generation with BoolODE (Pratapa et al., 2020)

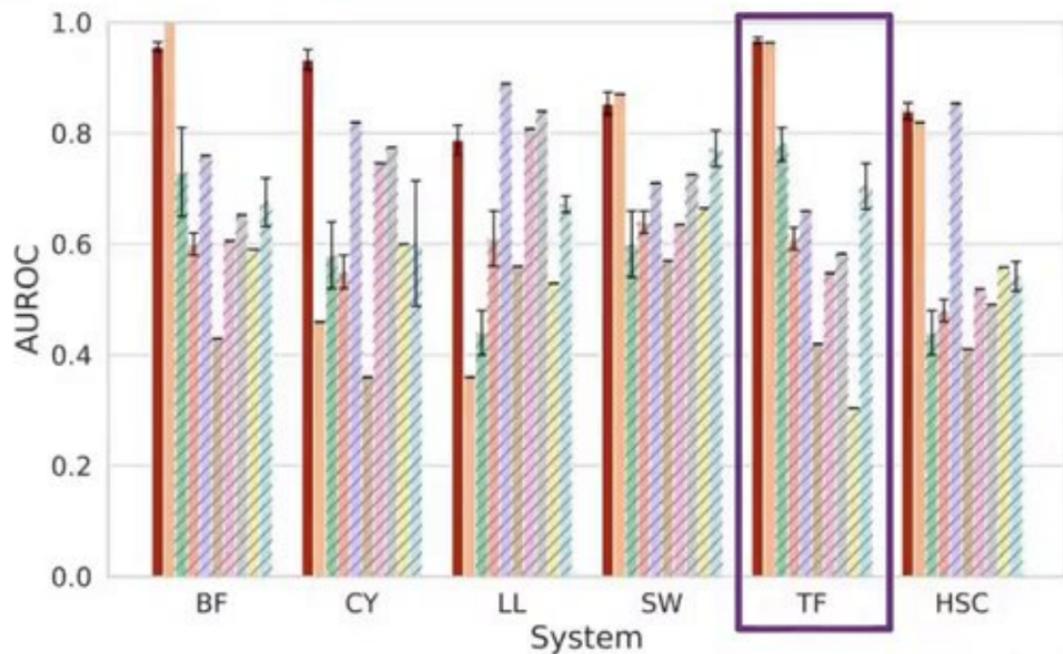


II - Simulated Biological Systems (*Structure Learning*)

Joint Trajectory Inference and GRN inference on small synthetic
→ GRN inference measures AP and AUROC of the recovered graph
ground truth



II - Simulated Biological Systems (Structure Learning)



II - Simulated Biological Systems (Trajectory inference)

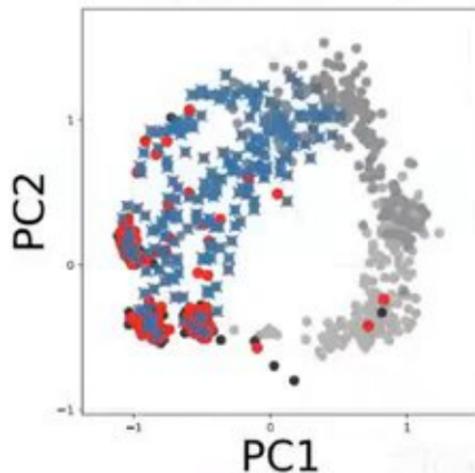
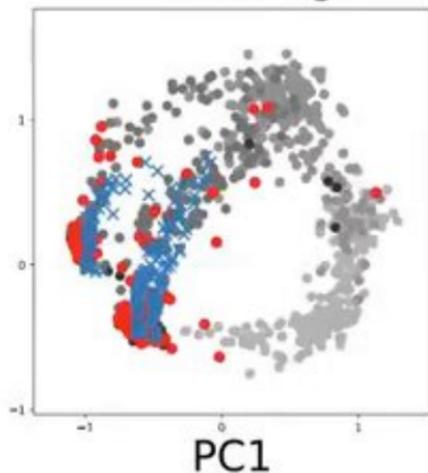
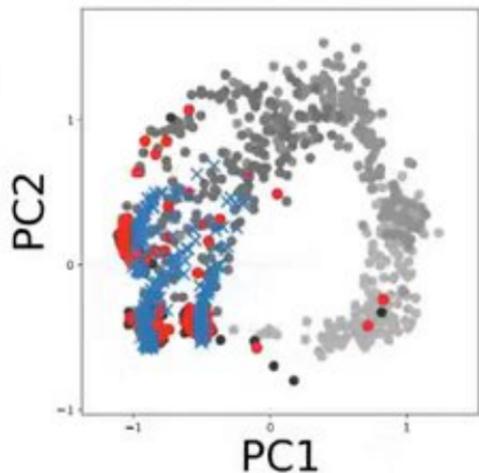
StructureFlow (ours)

Reference

Observational

Knockout g7

Observational



II - Simulated Biological Systems (Trajectory inference)

→ Trajectory inference measures W_2 and MMD, where one time hidden from the model during training

	TF		CY		LL		HSC (Curated)		$W_2 \downarrow$
	$W_2 \downarrow$	MMD \downarrow							
[SF] ² M (ODE)	0.76 ± 0.02	0.13 ± 0.01	0.52 ± 0.01	0.15 ± 0.00	0.85 ± 0.04	0.19 ± 0.01	0.67 ± 0.01	0.08 ± 0.00	0.38 ± 0.00
[SF] ² M (SDE)	0.79 ± 0.02	0.07 ± 0.01	0.53 ± 0.01	0.06 ± 0.00	1.11 ± 0.05	0.13 ± 0.01	0.64 ± 0.01	0.02 ± 0.00	0.45 ± 0.00
RF	1.03 ± 0.00	0.07 ± 0.00	0.76 ± 0.00	0.05 ± 0.00	1.19 ± 0.00	0.10 ± 0.00	0.72 ± 0.00	0.02 ± 0.00	0.81 ± 0.00
StructureFlow (ODE)	0.79 ± 0.01	0.13 ± 0.00	0.57 ± 0.01	0.14 ± 0.00	0.83 ± 0.02	0.19 ± 0.01	0.68 ± 0.01	0.07 ± 0.00	0.37 ± 0.00
StructureFlow (SDE)	0.82 ± 0.02	0.08 ± 0.01	0.59 ± 0.01	0.06 ± 0.00	1.05 ± 0.03	0.12 ± 0.01	0.66 ± 0.01	0.02 ± 0.00	0.42 ± 0.00

Top in **Bold**

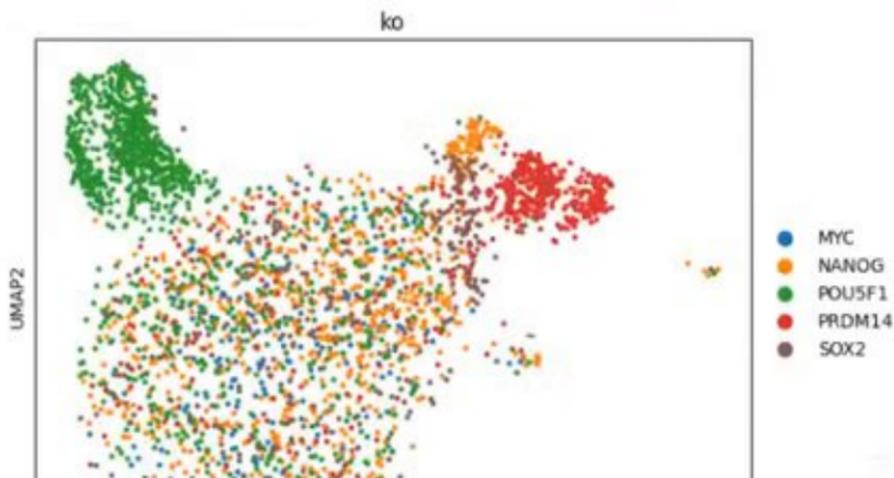
Second in **Blue**

II - Real Data (Renge, CRISPR KO)

Multi-marginal dataset (4 time-points) of iPSCs with CRISPR in (knockouts) for 23 transcription factors (Tokumasu et al., 2023)

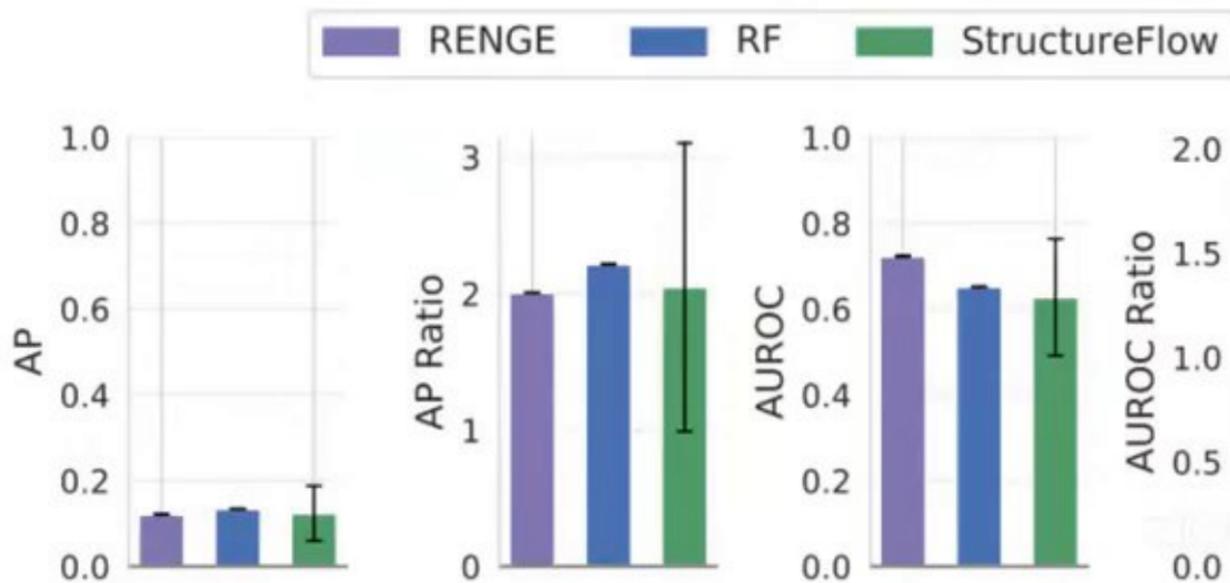
→ Higher dimensionality (103 genes)

→ Downside: ground truth network is a subset, experimental es



II - Real Data (Renge, CRISPR KO)

→ Gene Regulatory Network Inference Results



II - Real Data (Renge, CRISPR KO)

→ Trajectory Inference performance (prediction of unseen time)

	Average	
	$W_2 \downarrow$	MMD \downarrow
[SF] ² M (ODE)	5.75 ± 0.09	0.019 ± 0.004
[SF] ² M (SDE)	6.13 ± 0.11	0.018 ± 0.004
RF	6.53 ± 0.06	0.018 ± 0.004
STRUCTUREFLOW (ODE)	5.63 ± 0.04	0.020 ± 0.004
STRUCTUREFLOW (SDE)	5.94 ± 0.04	0.019 ± 0.004

Top in **Bold**
Second in **Blue**

Conclusion & Future Work

This is still a hard (and unsolved) problem!

We introduce an (efficient) method for jointly inferring network and dynamics (*with a single model and end-to-end optimization*)

- Still some limitations and things to improve ... future work?
- We don't think the model is learning *true* causal relationships
 - Still challenging to *systematically* select good hyperparameters

..



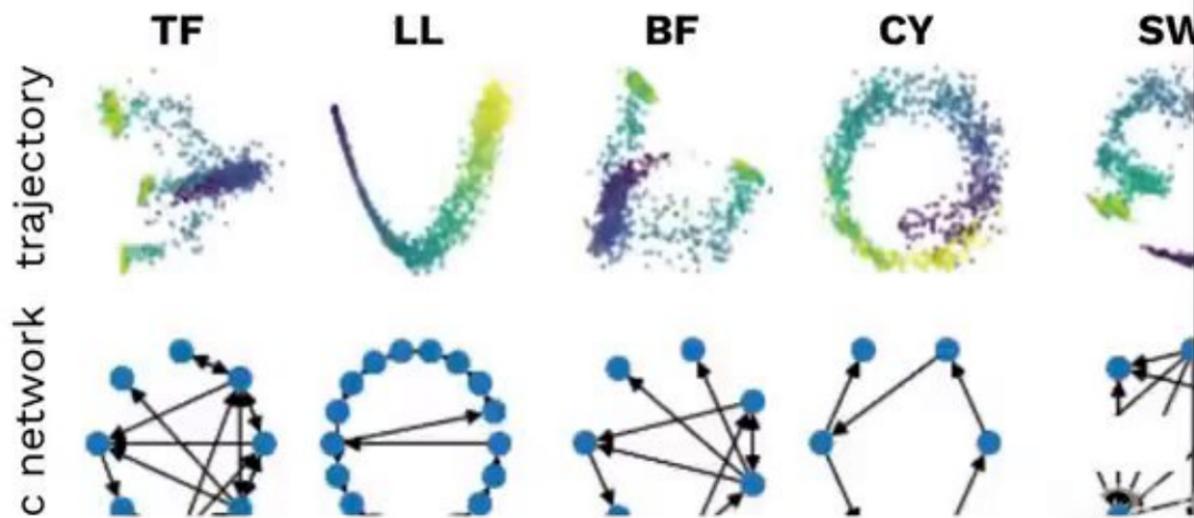
Acknowledgements



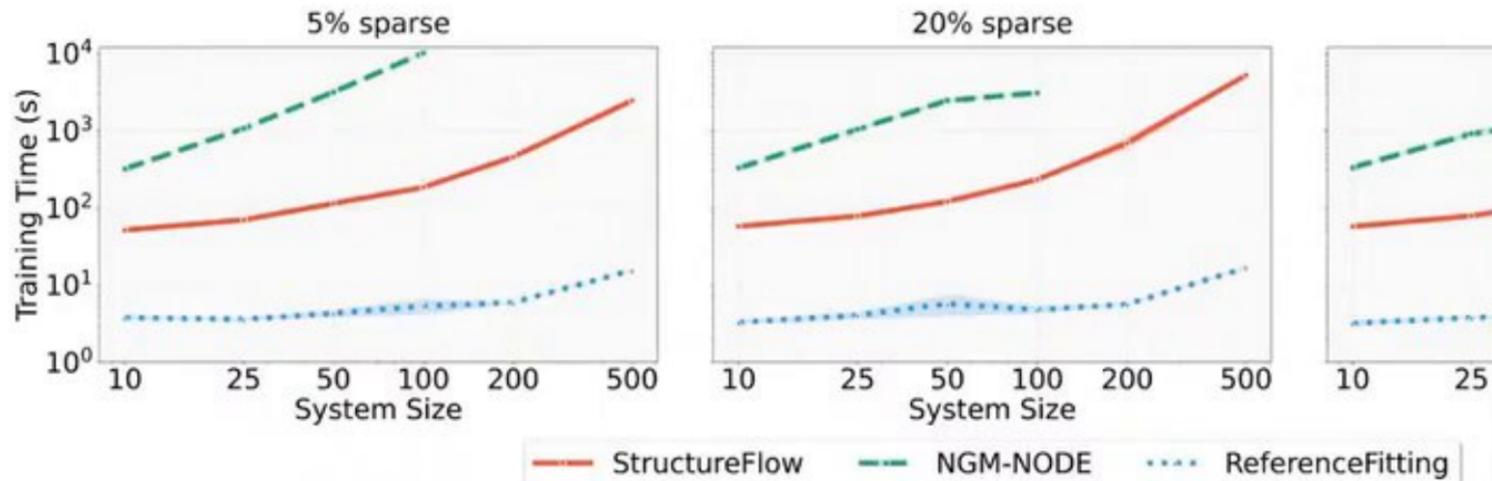
Questions?

II - Simulated Biological Systems

Joint Trajectory Inference and GRN inference on small, curated, systems, each representing common gene regulatory mechanics
Data generation with BoolODE (Pratapa et al., 2020)



StructureFlow can be trained efficiently





Dragon
data

mail.google.com/mail/u/0/#inboxFMfgzQcprwhpxDRZbpqVbpJJKpZdM

from: nsanjana@nygenome.org

148 emails found. All accounts : schedule a meeting to discuss AI

Neville Sanjana
Automatic reply Re: schedule a meeting to discuss AI ...
I'm out of office until september 12 with intermL... Sep 9

Automatic reply Re: AI company meeting?
I'm out of office until september 12 with inter... Aug 29

Neville, Gamze, Marcin, Dan
some personal news 5
agreed — you will be greatly missed but big c... Aug 20

July

Neville, Marcin, Rahul, Dan, Gamze

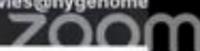
AI across biology: <http://www.csymposium-ai-for-life-science>

EasyChair MLC
slurm
ViewMail
More

SEP 11

daknowles@nygenome

Key Point 2:
can be anyth
gene expres





Dragon
data

Key Point 2:
can be anyth
gene expres

MLCB oral slides

MLCB oral slides

Name	Date Modified	Size	Kind
41_Linder Johannes Linder.pptx	Sep 9, 2025 at 12:00		werP...
42_Fannjiang Clara Fannjiang.key	Yesterday at 12:00		ynote
63_Shearer Courtney Shearer.pptx	Sep 9, 2025 at 12:00		werP...
114_Rocha Joao Felipe Rocha.pptx	Today at 9:06 AM		werP...
134_shaw Peter Shaw.pdf	Yesterday at 10:00		if Docu...
154_BLASSEL Luc Blassel.pdf	Yesterday at 11:00		if Docu...
2025 MLCB Keynote Jacob Schreiber.pptx	Today at 8:45 AM	28.3 MB	werP...
bee_MLCB_2025 Barbara.pptx	Yesterday at 1:17 PM	477.9 MB	werP...
MLCB (15 min) Alan Amin.key	Yesterday at 11:58 AM	4.5 MB	ynote
MLCB_2025_Battle Alexis Battle.pdf	Yesterday at 9:46 AM	6.2 MB	if Docu...
MLCB_2025_Battle Alexis Battle.pptx	Yesterday at 9:08 AM	37 MB	werP...
Perturbation_Benchmark_Hasanaj Euxhen Hasanaj.pdf	Today at 9:33 AM	1.1 MB	if Docu...

chedule a meeting to disc

zoom

The image shows a Zoom meeting window with a PowerPoint presentation. The slide is titled "STAGED: A Multi-Agent Neural Network for Learning Cellular Interaction Dynamics". The authors listed are João F. Rocha*, Ke Xu, Xingzhi Sun, Ananya Krishna, Dhananjay Bhaskar, Blanche Mongeon, Morgan Craig, Mark Gerstein, and Smita Krishnaswamy. The slide is associated with Yale School of Medicine and The Krishnaswamy Lab, which is represented by a circular logo featuring a stylized face. The Zoom interface includes a top toolbar with various editing and navigation tools, a left sidebar with a slide thumbnail list, and a bottom status bar showing "Slide 0 of 0" and "Good to go".

STAGED:
A Multi-Agent Neural Network for Learning
Cellular Interaction Dynamics

João F. Rocha*, Ke Xu, Xingzhi Sun, Ananya Krishna, Dhananjay Bhaskar,
Blanche Mongeon, Morgan Craig, Mark Gerstein, Smita Krishnaswamy

Yale SCHOOL OF MEDICINE

THE KRISHNASWAMY LAB

Slide 0 of 0 ... Good to go

STAGED:
A Multi-Agent Neural Network for Learning
Cellular Interaction Dynamics

João F. Rocha*, Ke Xu, Xingzhi Sun, Ananya Krishna, Dhananjay Bhaskar,
Blanche Mongeon, Morgan Craig, Mark Gerstein, Smita Krishnaswamy

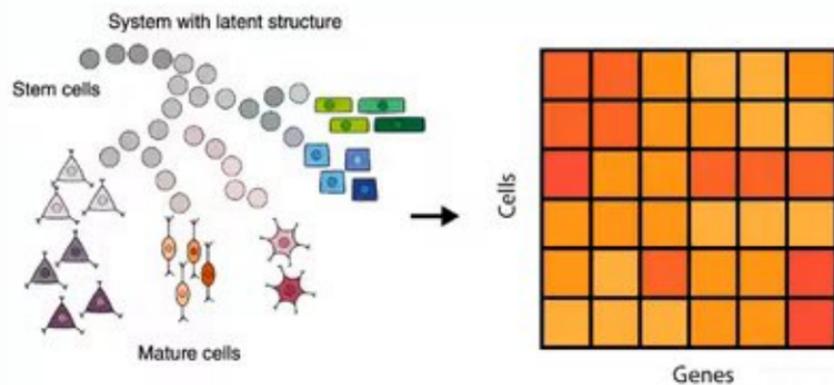
Yale SCHOOL OF MEDICINE



THE KRISHNASWAMY LAB 

Transcriptomics Data

Single Cell Data



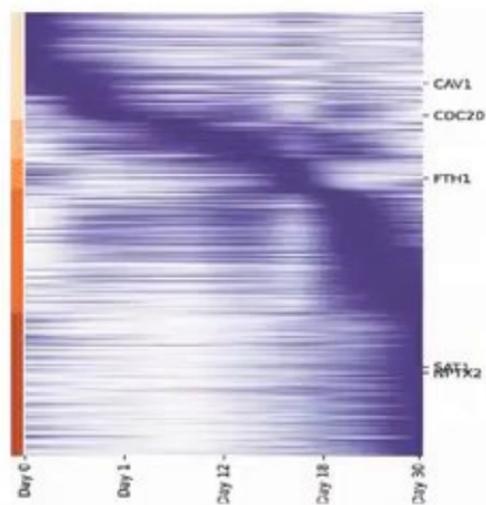
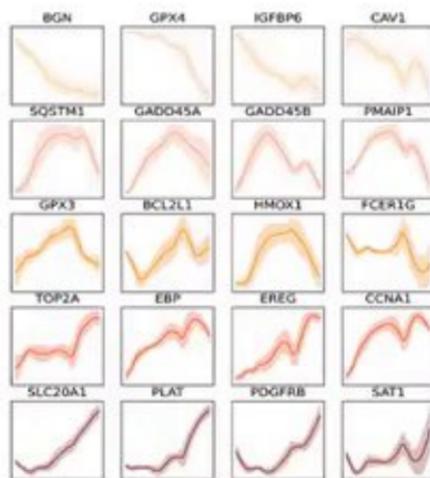
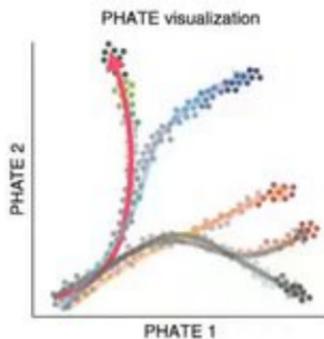
Trajectories in the data



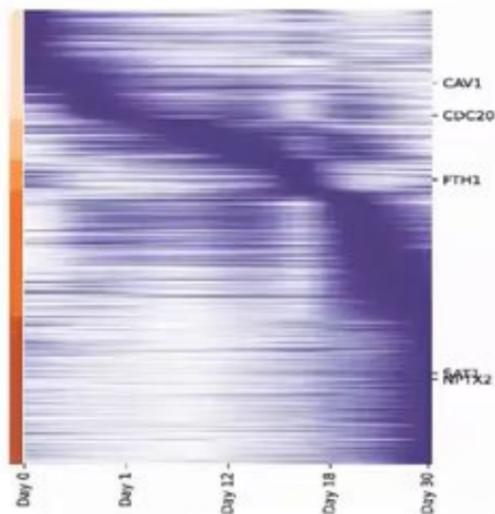
Huguet, G., et al. Manifold Interpolating Optimal-Transport Flows for Trajectory Inference. NEURIPS (2022).

Yale SCHOOL OF MEDICINE

Trajectories in the data

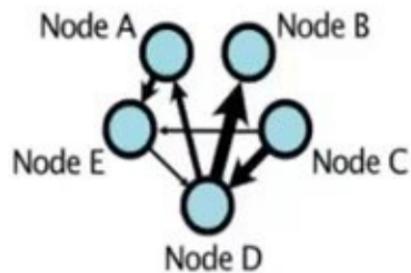


Gene Dynamics

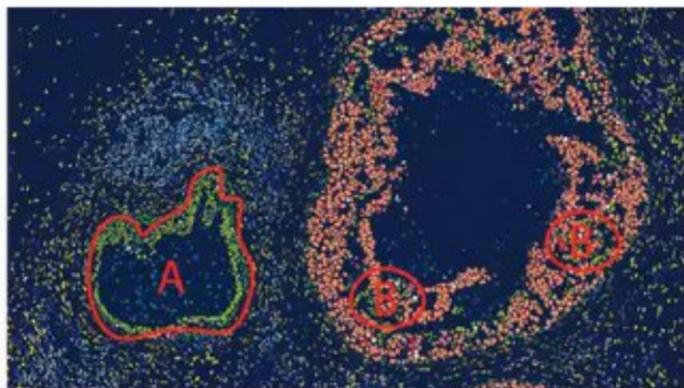


Granger
Causality

Gene Regulatory Networks

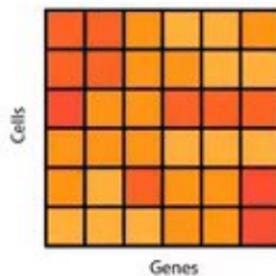


Spatial Transcriptomics Data



<https://ostz.ccr.cancer.gov/emerging-technologies/spatial-biology/xenium/> Accessed: Mar. 04, 2025

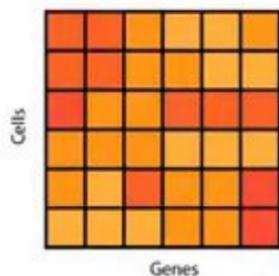
- **Spatial Transcriptomics:** Gives information on both the space and gene-level.



0.1	1.1
0.2	4.1
2.1	2.9
3.5	7.6
5.6	0.7
0.6	1.2

Position

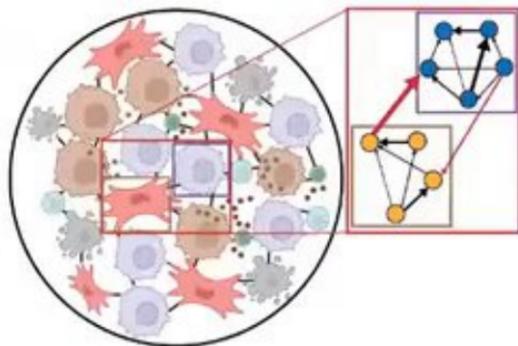
Spatial Transcriptomics Data: Hierarchical Graphs



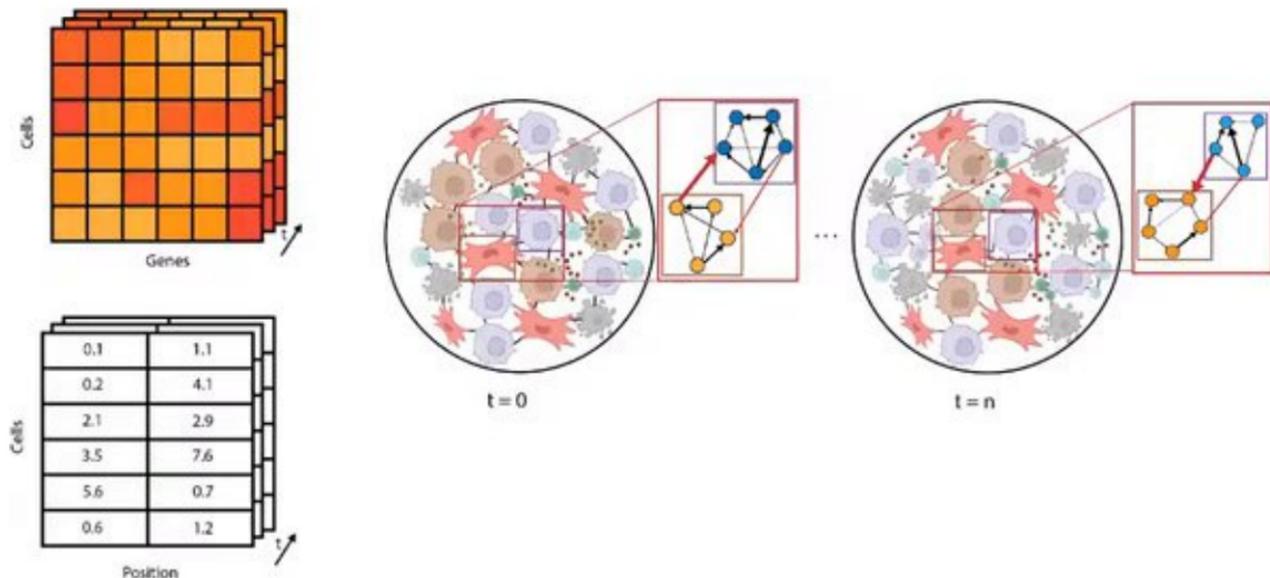
0.1	1.1
0.2	4.1
2.1	2.9
3.5	7.6
5.6	0.7
0.6	1.2

Cells

Position

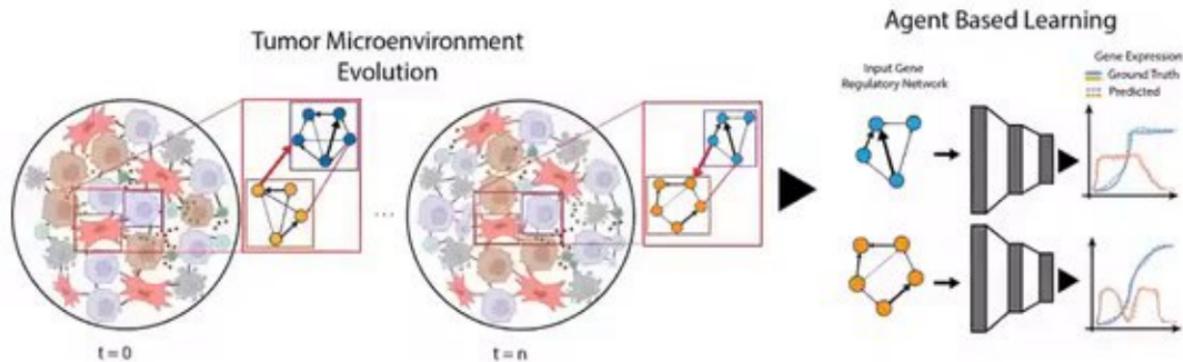


Time evolving Spatial Transcriptomics Data



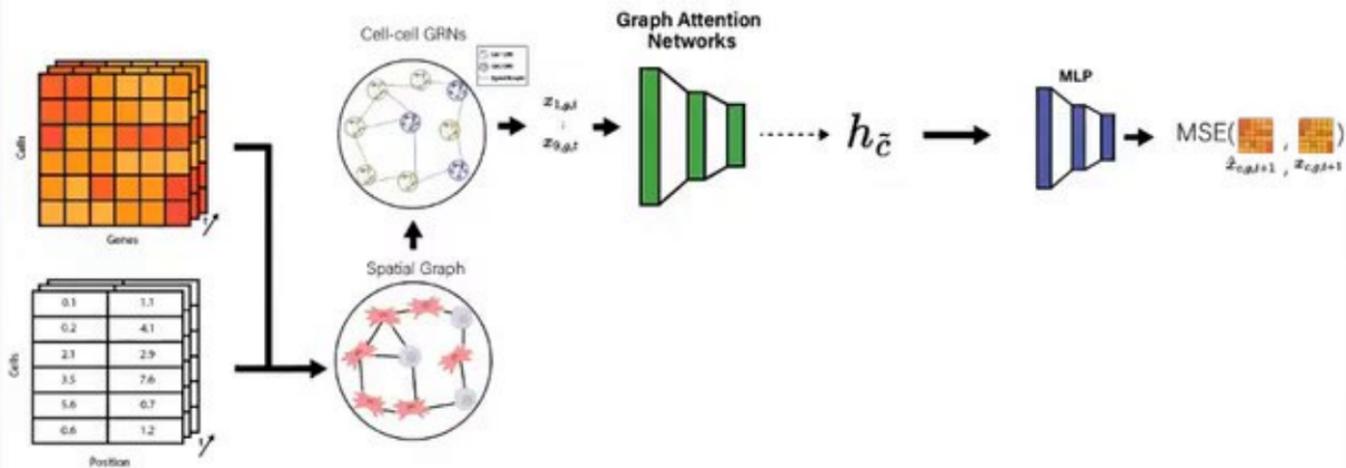
STAGED

Spatio-Temporal Agent-Based Graph Evolution Dynamics (STAGED)

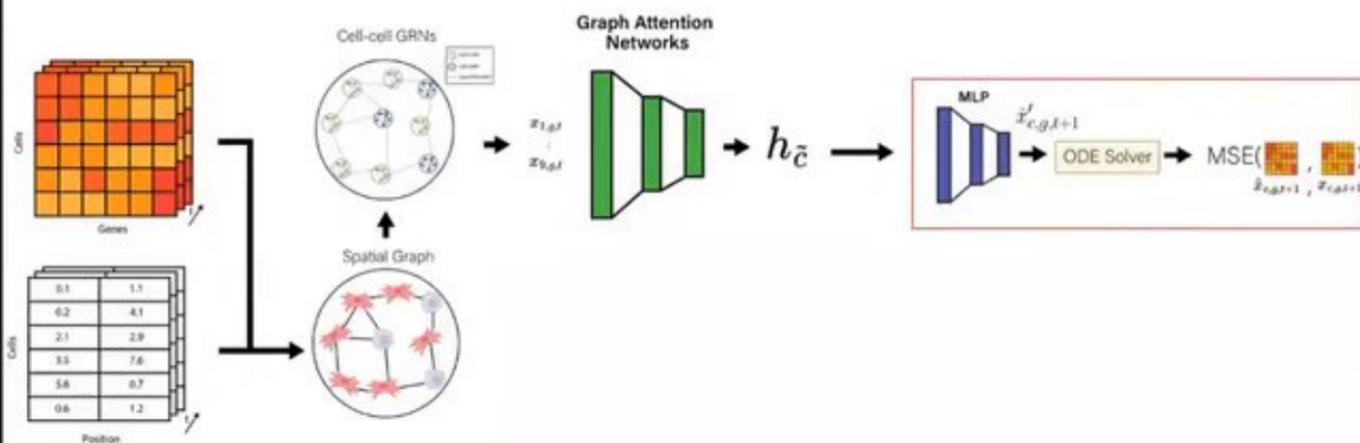


- **Level 1:** Cell-cell communication (agents talking)
- **Level 2:** Gene regulatory networks (within each cell)

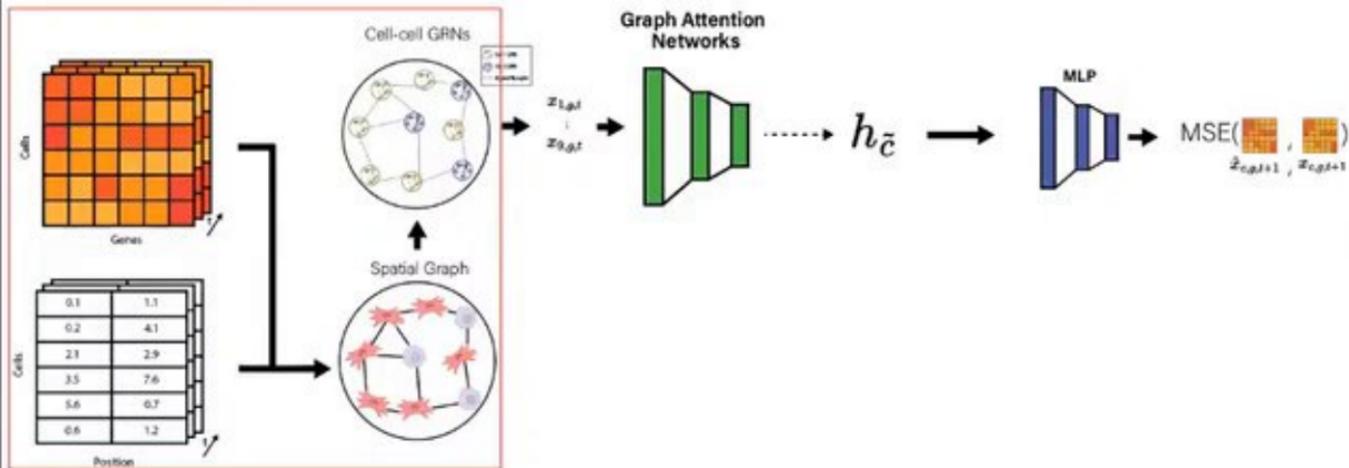
High-level Architecture



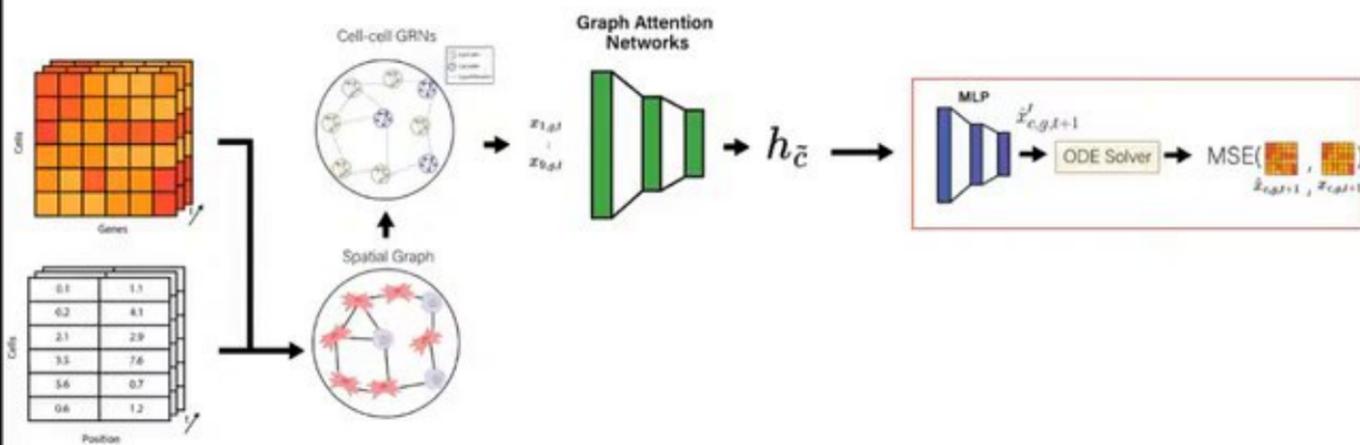
Better dynamics inference



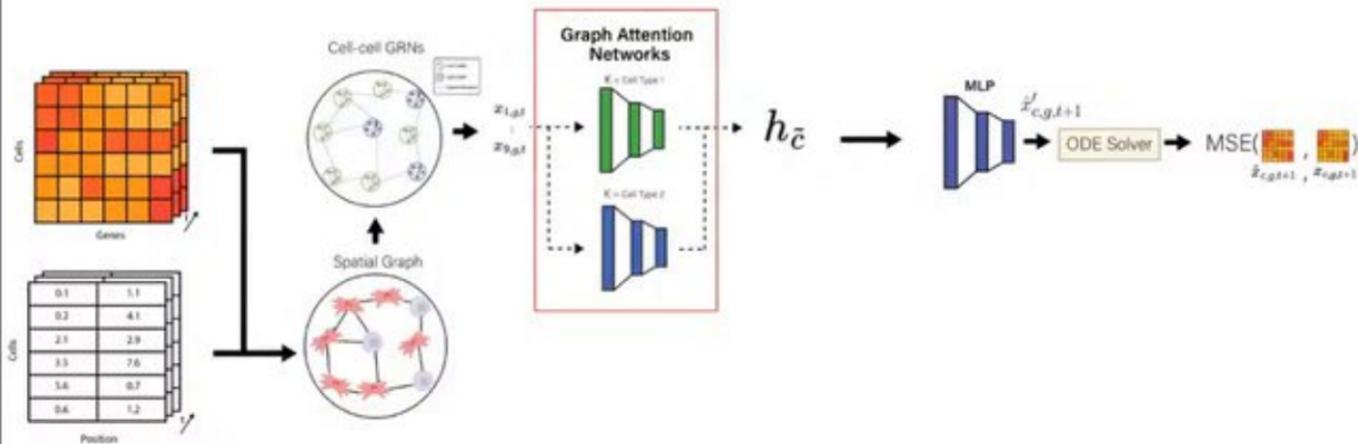
Biological Prior I: Ligand Receptor Pairs



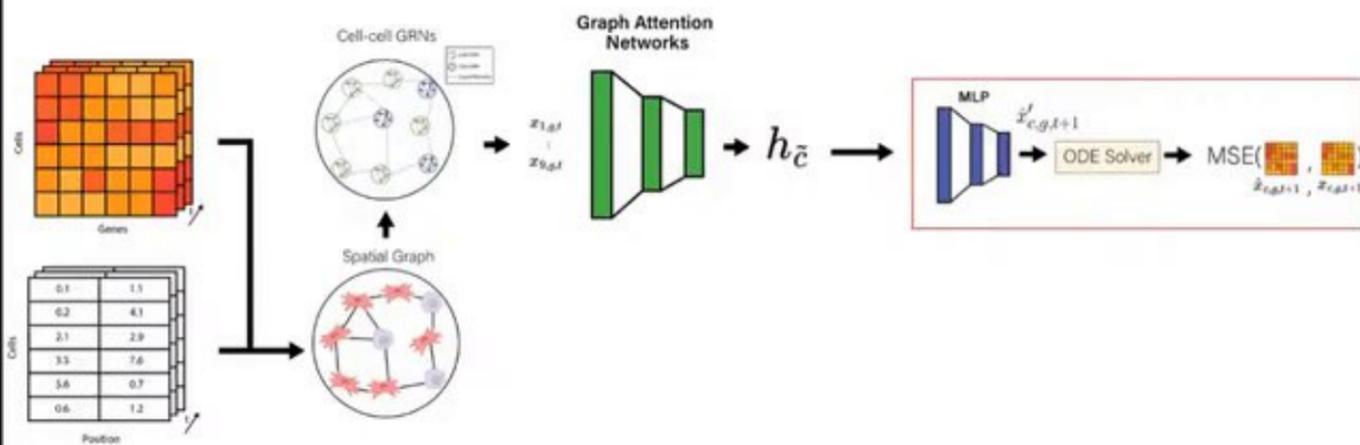
Better dynamics inference



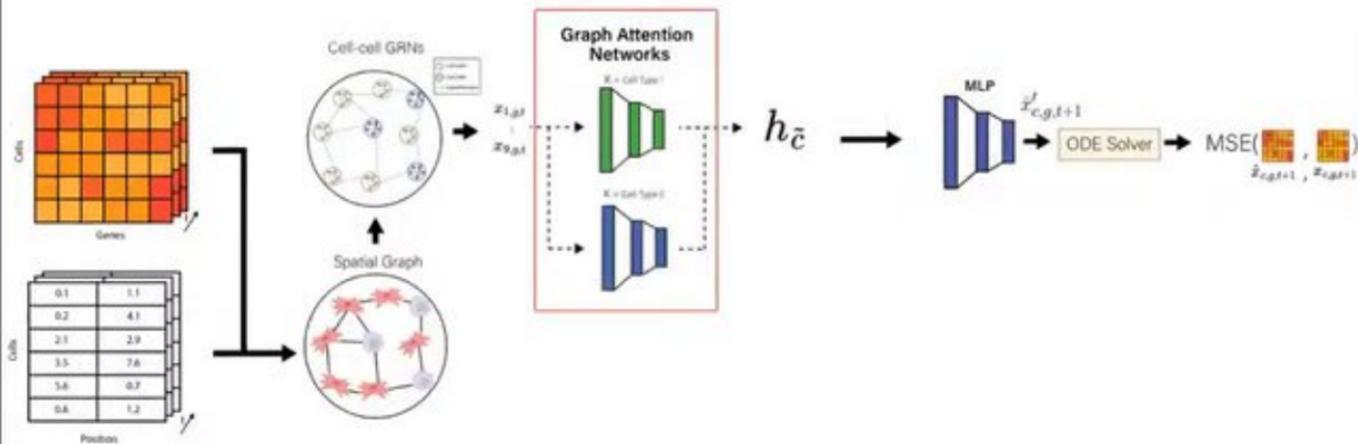
Biological Prior II: Agent-based learning



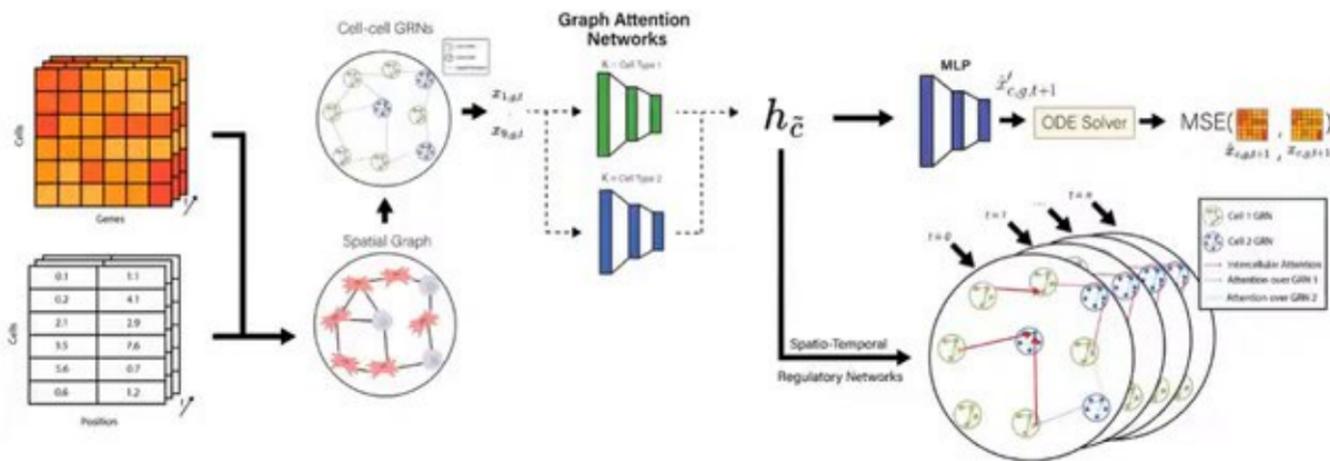
Better dynamics inference



Biological Prior II: Agent-based learning



Obtaining the cell-cell interactions from attention



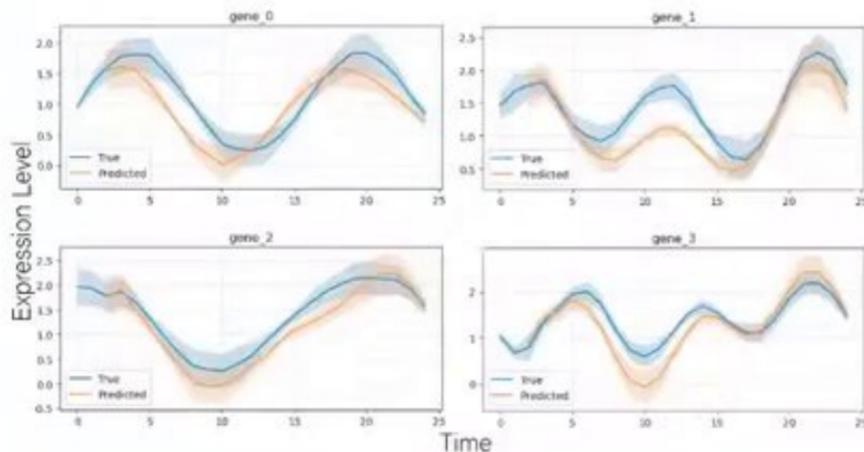
Spatial neighborhoods \rightarrow Attention mechanism \rightarrow Gene dynamics

Each agent in **STAGED** is a cell-type-specific model composed of a graph attention network (GAT) and a neural ODE.

Preliminary Results

Simulating Oscillatory Dynamics

Mimicking oscillatory dynamics

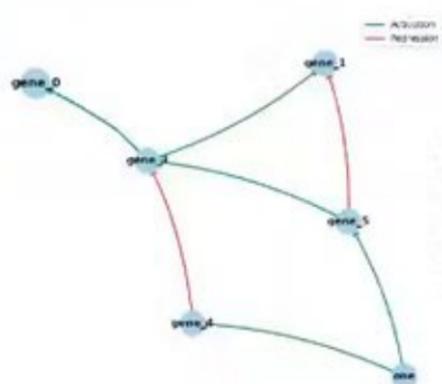


STAGED can accurately capture these dynamics, shown here one cell type for a subset of genes.

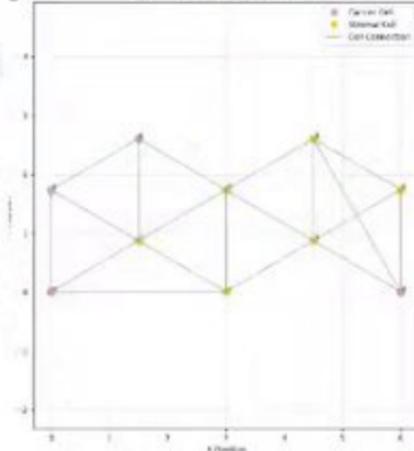
Results – Real Data

Simulations using simulated GRNs

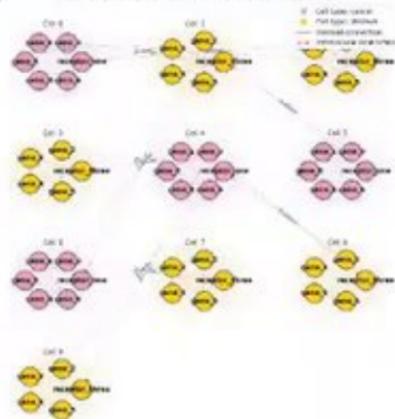
A Gene Regulatory Network



B Cell Positions with Neighbor Connections



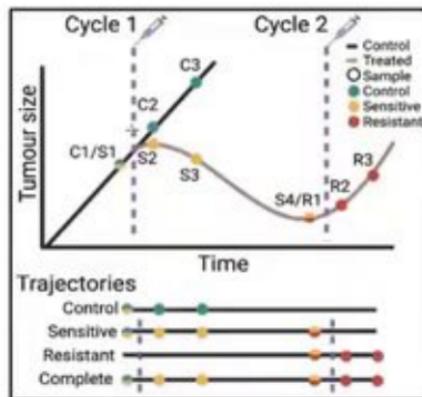
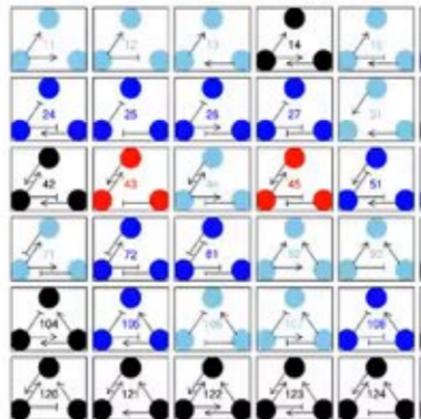
C Gene Interaction Network with Cell-Cell Interactions



Spatially resolved stochastic simulation of gene expression dynamics in a randomly constructed gene regulatory network (GRN) using GillesPy2.

Next Steps

- Use small simulation of GRNs
- Improve benchmarking and simulations
- Apply Multi-Agent RL techniques
- Application to new cancer data



Acknowledgements



Ke Xu



Xingzhi Sun



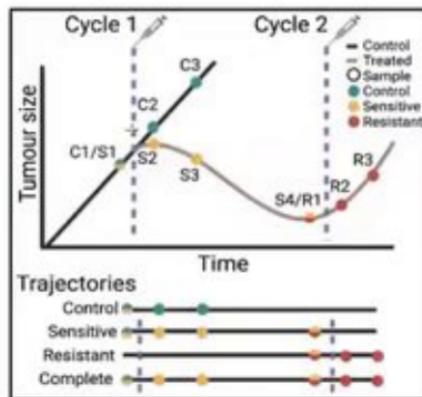
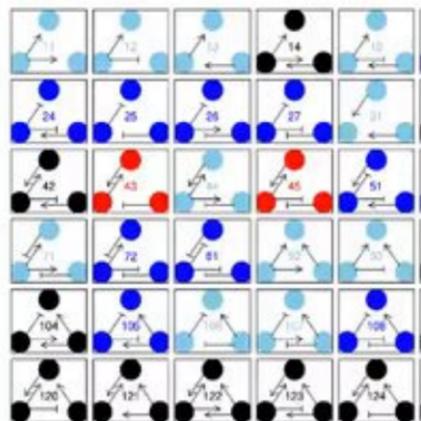
STAGED Authors

Ananya Krishna
Dhananjay Bhaskar
Blanche Mongeon
Morgan Craig
Mark Gerstein
Smita Krishnaswamy

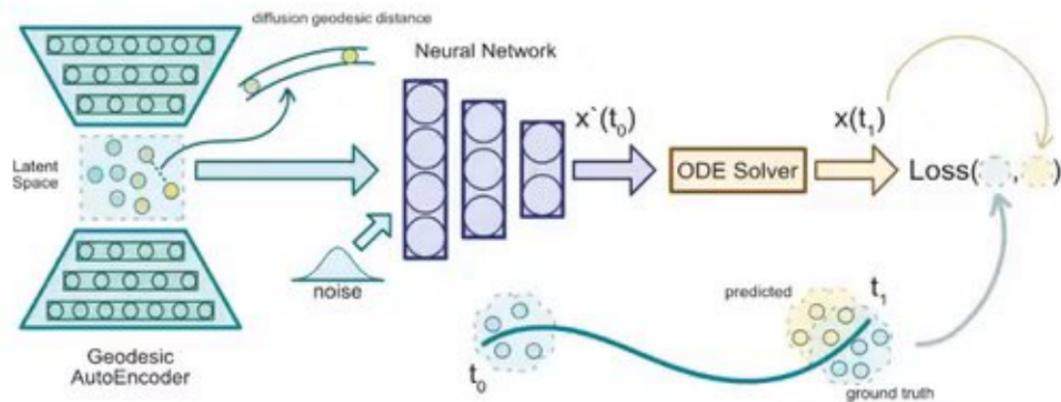


Next Steps

- Use small simulation of GRNs
- Improve benchmarking and simulations
- Apply Multi-Agent RL techniques
- Application to new cancer data



MIOFlow



MIOFlow models continuous population dynamics by learning optimal transport-driven flows in a manifold-aware latent space to interpolate between static data snapshots.

NYGC Events

Motivation

- Why perturbation response matters?
- **Current challenge:** FMs often don't outperform simple baselines.
- **Key gap:** we don't understand which part (embeddings, downstream models, datasets) drive performance.

Brief Communication | [Open access](#) | Published: 04 August 2025

Deep-learning-based gene perturbation effect prediction does not yet outperform simple linear baselines

Constantin Ahlmann-Eltze , Wolfgang Huber & Simon Anders

Nature Methods **22**, 1657–1661 (2025) | [Cite this article](#)

49k Accesses | 8 Citations | 181 Altmetric | [Metrics](#)

A systematic comparison of computational methods for expression forecasting

Eric Kenfeld, erik@stanford.edu, <https://orcid.org/0000-0002-2310-6191>

[†]Yanliao Yang, yyanliao@stanford.edu, <https://orcid.org/0000-0003-3429-3174>

[†]Joshua S. Weinstein, jweinst17@stanford.edu, <https://orcid.org/0000-0001-7013-1889>

^{††††}Alexis Balle, aballe@stanford.edu (corresponding author)

<https://orcid.org/0000-0002-5287-627X>

^{†††}Patrick Cahan, patrick.cahan@stanford.edu (corresponding author)

<https://orcid.org/0000-0003-3652-2540>

choice of metric, and especially for simple metrics like mean squared error, it is uncommon for expression forecasting methods to out-perform simple baselines. Our platform will serve as a



Multimodal Benchmarking of Foundation Model Representations for Cellular Perturbation Response Prediction

Euxhen Hasanaj¹, Elijah Cole¹, Shahin Mohammadi¹, Sohan Addagudi^{1,2}, Xingyi Zhang^{1,3}, Le Song^{1,3}, Eric Xing^{1,2,3}

¹ GenBio AI

² Carnegie Mellon University

³ Mohamed bin Zayed University of Artificial Intelligence

Presenter: Euxhen Hasanaj

zoom

Motivation

- Why perturbation response matters?
- **Current challenge:** FMs often don't outperform simple baselines.
- **Key gap:** we don't understand which part (embeddings, downstream models, datasets) drive performance.

Brief Communication | [Open access](#) | Published: 04 August 2025

Deep-learning-based gene perturbation effect prediction does not yet outperform simple linear baselines

[Constantin Ahlmann-Eltze](#)  [Wolfgang Huber](#) & [Simon Anders](#)

[Nature Methods](#) **22**, 1657–1661 (2025) | [Cite this article](#)

49k Accesses | 8 Citations | 181 Altmetric | [Metrics](#)

A systematic comparison of computational methods for expression forecasting

[Eric Kerfeld](#), [eric@hpi.edu](#), <https://orcid.org/0000-0002-2310-8191>

[Yanliao Yang](#), [yyang117@hpi.edu](#), <https://orcid.org/0000-0003-3429-3174>

[Joshua S. Wainstock](#), [jwainsto17@hpi.edu](#), <https://orcid.org/0000-0001-7013-1899>

¹[Alexis Battle](#), [abattle@hpi.edu](#) (corresponding author).

²<https://orcid.org/0000-0002-5287-627X>

³[Patrick Cahan](#), [patrick.cahan@hpi.edu](#) (corresponding author).

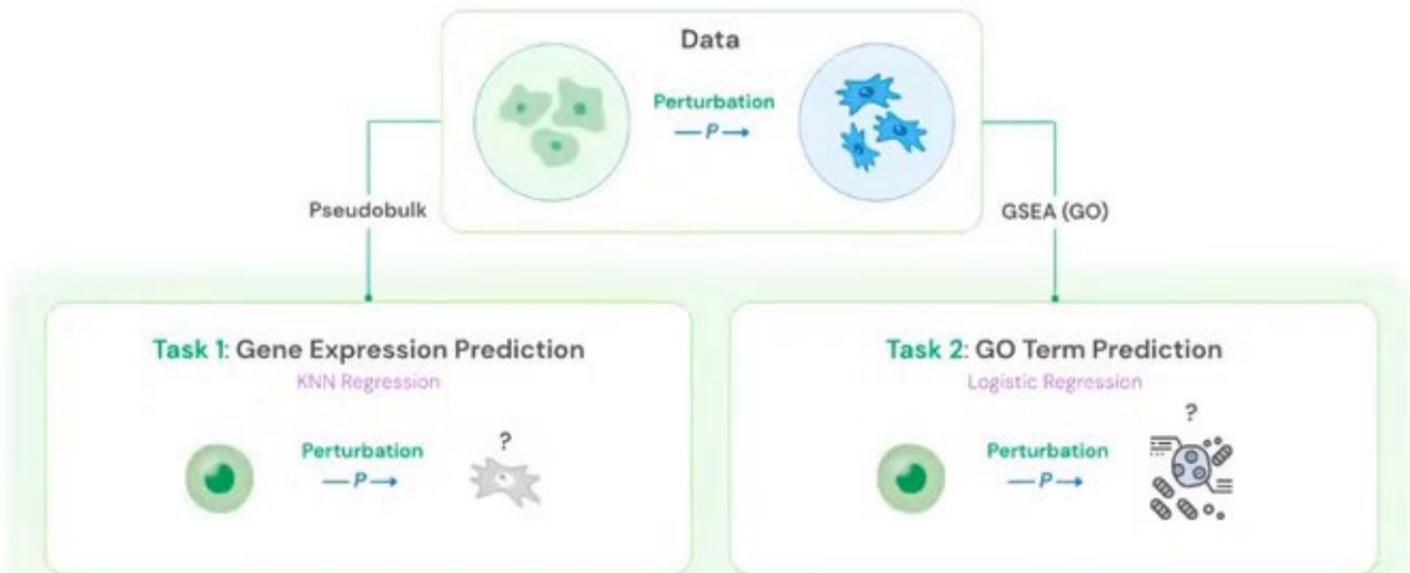
<https://orcid.org/0000-0003-3652-2540>

choice of metric, and especially for simple metrics like mean squared error, it is uncommon for expression forecasting methods to out-perform simple baselines. Our platform will serve as a

Goal

- What makes an embedding useful for perturbation response?
 1. We **fix the perturbation response model** to decouple embeddings from downstream predictors.
 2. We study a **diverse collection of modalities** and representation learning strategies.
 3. We utilize **simple metrics** (e.g., L2), but also more **biologically interpretable scoring methods**.

Methods: Problem Formulation



Methods: Embeddings Benchmarked

- Categories
 - **Baselines:** Random, No Change, PCA, Train Mean, Idealized
 - **No Change:** Predict the expression of control cells
 - **Idealized:** Embeddings obtained by fitting PCA on *all* pseudobulk data (train + test)
 - **Random:** Embeddings drawn at random from a Gaussian distribution
 - **Expression FMs:** AIDO.Cell (3M/10M/100M), scGPT, Geneformer, scPRINT, Transcriptformer
 - **Protein FMs:** AIDO.ProteinIF, ESM2, AIDO.StructureTokenizer, STRING-SPACE
 - **DNA FMs:** AIDO.DNA
 - **Network embeddings:** STRING (network, NBFNet, WaveGC)
 - **Prior Knowledge:** GenePT (GO BP/CC/MF/NCBI+UniProt), GenotypeVAE
- Stress Diversity
 - Expression, sequence, structure, prior knowledge

Datasets

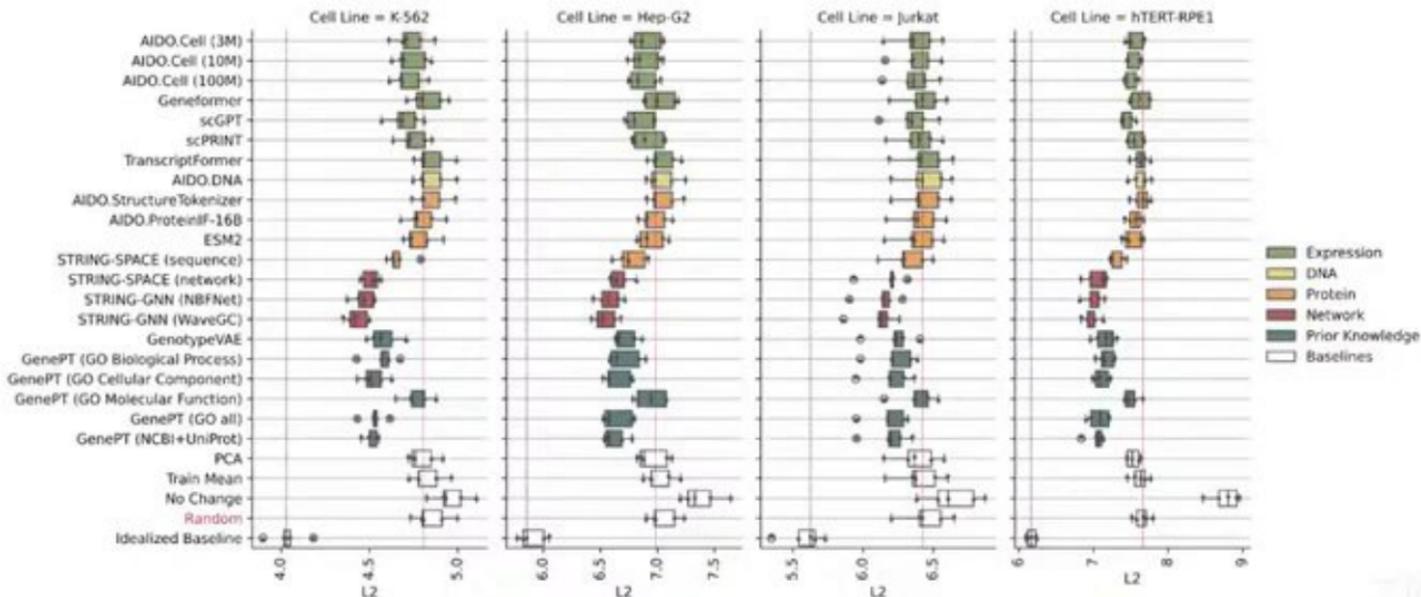
- **K-562** [1]
 - lymphoblast from individual with chronic myelogenous leukemia
 - 1534 perturbations
- **RPE1** [1]
 - retinal pigment epithelial cell from a healthy individual
 - 1752 perturbations
- **Jurkat** [2]
 - T cell from individual with leukemia
 - 1752 perturbations
- **Hep-G2** [2]
 - epithelial cell from an individual with hepatocellular carcinoma
 - 1752 perturbations

[1] Replogle, J. M. *et al.* Mapping information-rich genotype-phenotype landscapes with genome-scale Perturb-seq. *Cell* **185**, 2559-2575.e28 (2022).

[2] Nadig, A. *et al.* Transcriptome-wide analysis of differential expression in perturbation atlases. *Nat. Genet.* **57**, 1228–1237 (2025).

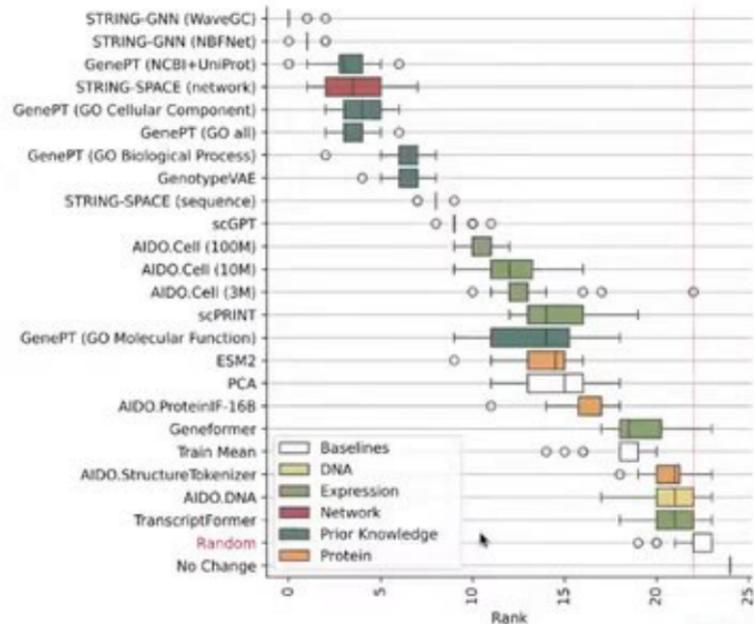
Results: Per-Cell Line Performance

Split-to-split variability: robust cross-validation is important!



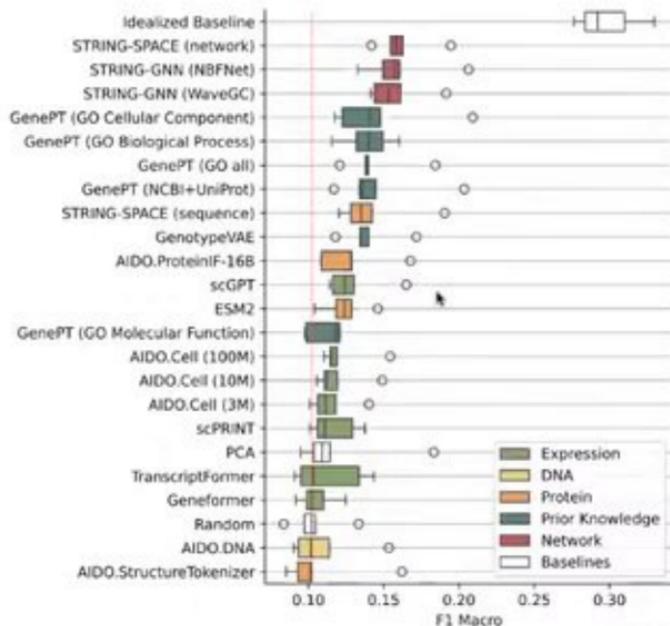
Results: Overall Ranking

1. **Knowledge-based and network embeddings** perform best.
2. Some cell **FMs outperform PCA**.
3. **Scaling helps** (AIDO.Cell 100M > 10M > 3M).
4. Protein FMs — **STRING-sequence best**.
5. GenePT using **"GO Cellular Component"** outperforms "GO Biological Process".



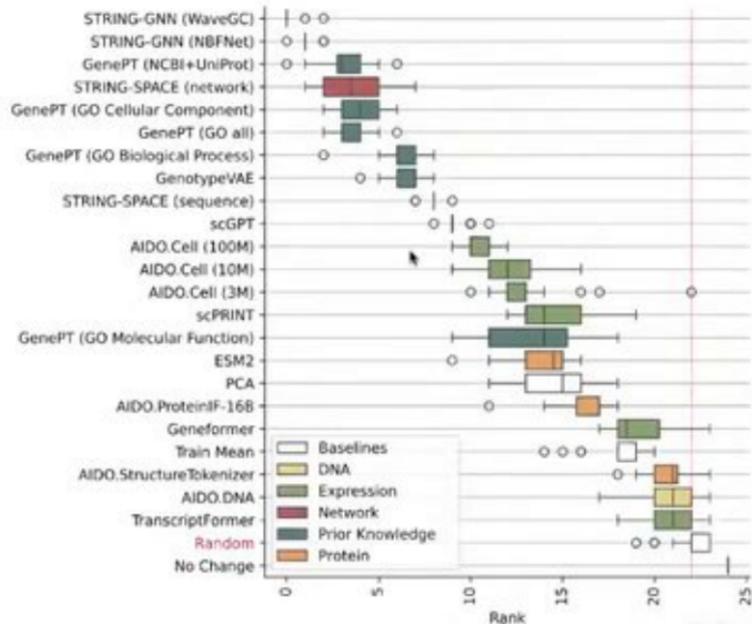
Results: Function Prediction

1. Models **ranked similarly** as for expression prediction.
2. Knowledge-driven embeddings again strongest.
3. **Protein embeddings slightly better** here vs expression.
4. There is still room for improvement (idealized baseline).



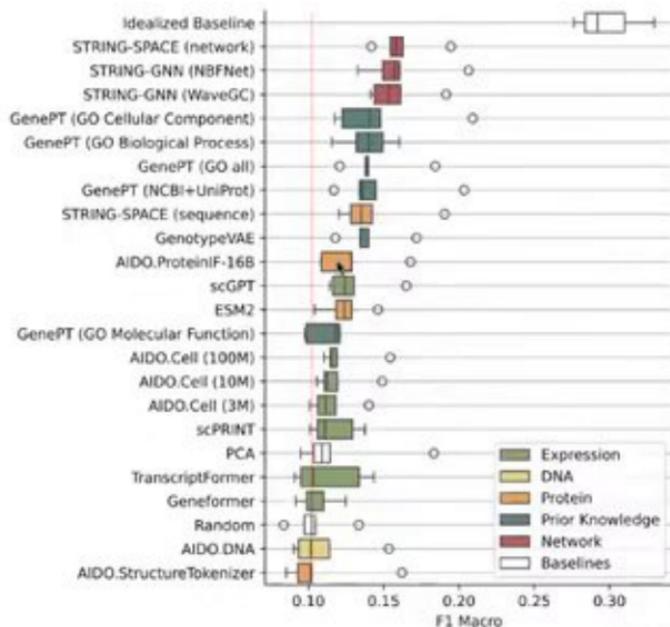
Results: Overall Ranking

1. **Knowledge-based and network embeddings** perform best.
2. Some cell FMs outperform PCA.
3. **Scaling helps** (AIDO.Cell 100M > 10M > 3M).
4. Protein FMs — **STRING-sequence** best.
5. GenePT using **"GO Cellular Component"** outperforms "GO Biological Process".



Results: Function Prediction

1. Models **ranked similarly** as for expression prediction.
2. Knowledge-driven embeddings again strongest.
3. **Protein embeddings slightly better** here vs expression.
4. There is still room for improvement (idealized baseline).



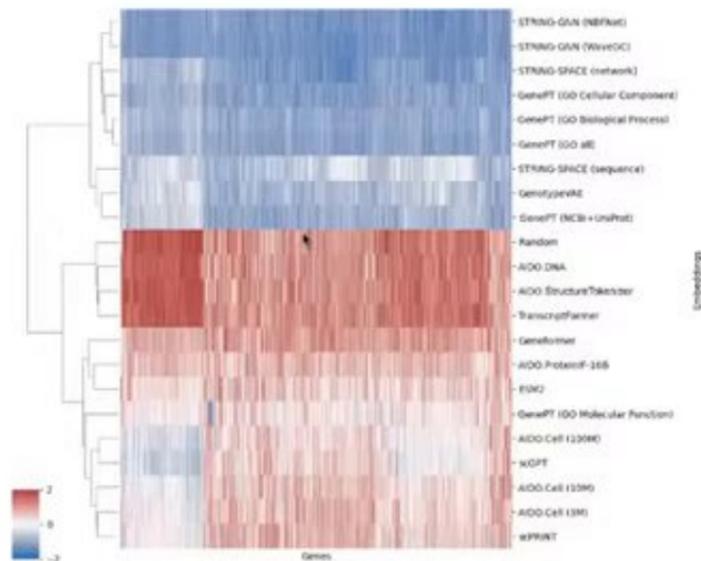
Limitations & Future Directions

- **Limitations**

- Only one model type (kNN/logistic)
- Pseudobulk, not single-cell
- Fixed embeddings (no fine-tuning)

- **Next steps**

- Fusion across scales / modalities
 - Extend beyond single-modality cell FMs by integrating diverse biological data types
- Multi cell-line generalization



Contributions

- Conducted a **systematic pan-modal study of perturbation embeddings**, and evaluated their predictive utility separately from downstream predictors.
- Evaluated embeddings from **20+ biological foundation models** and rigorously benchmarked them against simpler baselines.
- Introduced a **novel, biologically interpretable** formulation of perturbation response modeling based on predicting functional GO terms.
- Knowledge-based and network embeddings perform best.

Multimodal Benchmarking of Foundation Model Representations for Cellular Perturbation Response Prediction

Euxhen Hasanaj

Research Scientist

GenBio AI

Contact: euxhen.hasanaj@genbio.ai

Elijah Cole

Shahin Mohammadi

Sohan Addagudi

Xingyi Zhang

Le Song

Eric Xing



**Carnegie
Mellon
University**



bioRxiv

THE PREPRINT SERVER FOR BIOLOGY

Disclaimer: All authors are affiliated with GenBio AI.

NYGC Events



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley

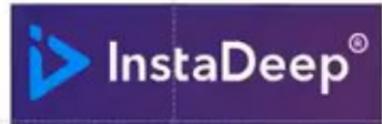
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 11.55am

mlcb.org for schedule



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley
 Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra
 Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott
 Adamson, Trevor Christensen, Yijie Kang

Restarting at 11.55am

mlcb.org for schedule

Click to add notes



cover slide — Saved to my Mac

Search (Cmd + Ctrl + U)

Record Comments Share

Home Insert Draw Design Transitions Animations **Slide Show** Record Review View

Play from Start Play from Current Slide Presenter View Custom Show Rehearse with Coach Set Up Slide Show Hide Slide Rehearse Timings Record Play Narrations Show Media Controls Always Use Subtitles Subtitle Settings

1



@



NEW YORK
GENOME CENTER®

Thank you to our sponsors







@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley
 Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 11.55am

mlcb.org for schedule

Click to add notes





@



Thank you to our sponsors



CORTEVA™
agriscience



InstaDeep®



Col

And helpers!

Sarah Curtiss , Kristen Weatherley
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das,
Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, La
Adamson, Trevor Christensen, Yijie Kang

NYGC Events



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley
 Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra
 Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott
 Adamson, Trevor Christensen, Yijie Kang

Restarting at 11.55am

mlcb.org for schedule

Click to add notes





@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley

Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 11.55am

mlcb.org for schedule

NYGC Events

Deep end-to-end likelihood-free inference of phylogenetic trees



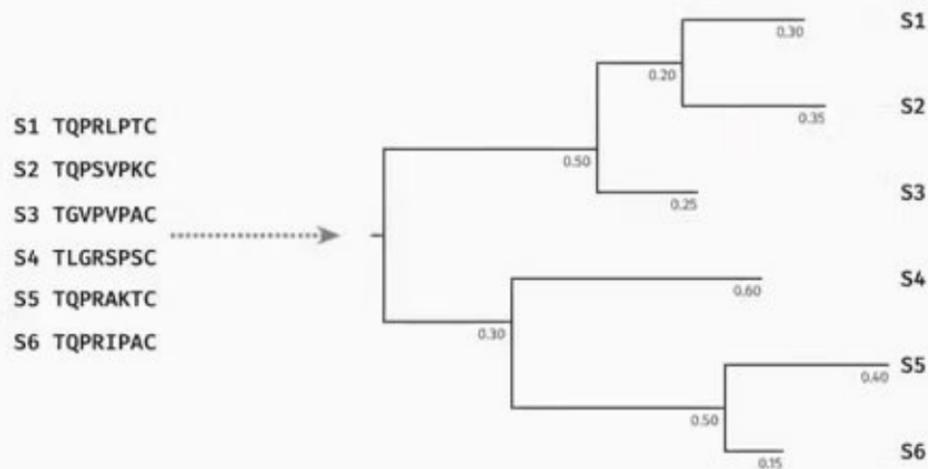
1

Luc Blassel, Nicolas Lartillot, Bastien Boussau, Laurent Jacob

MLCB - September 11th, 2025

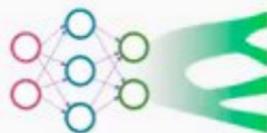


Context - Phylogenetic inference



Goal: describe **evolutionary-history** of MSA

Deep end-to-end likelihood-free inference of phylogenetic trees

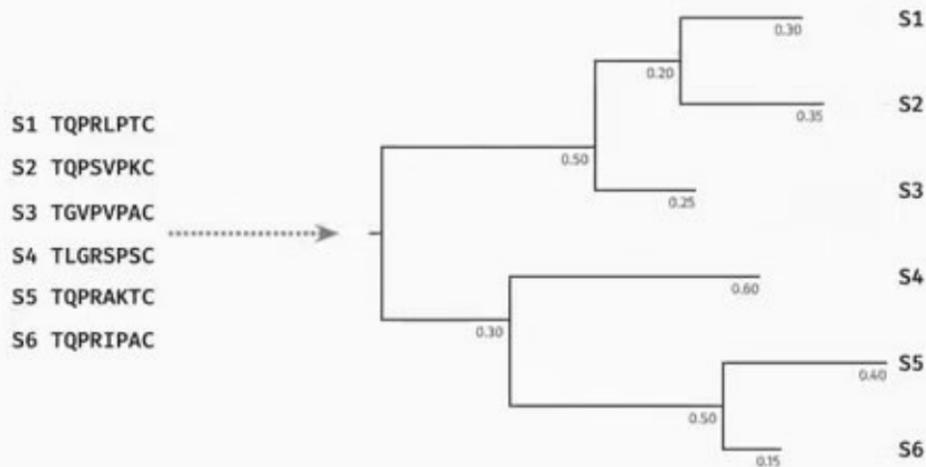


Luc Blassel, Nicolas Lartillot, Bastien Boussau, Laurent Jacob

MLCB - September 11th, 2025

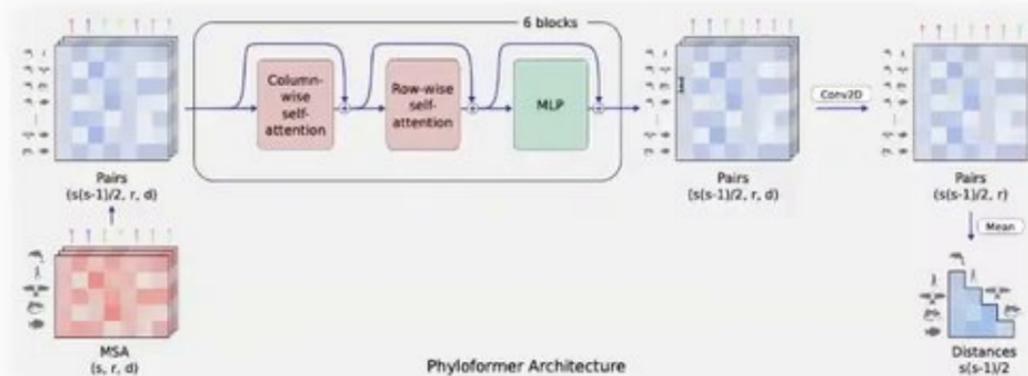


Context - Phylogenetic inference



Goal: describe **evolutionary-history** of MSA

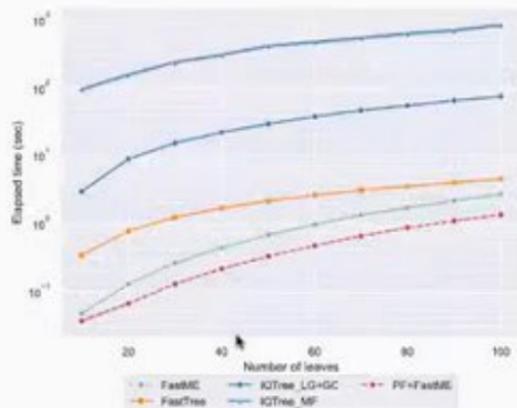
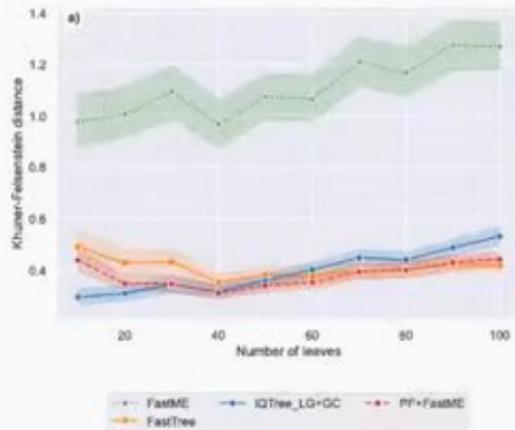
Related Work - Phyloformer, our first approach



- Input an **MSA**, get a **Distance matrix**
- Feed Distance matrix to **FastME** to get **tree**

Nesterenko et al. 2025; Lefort et al. 2015



Related Work - **Phyloformer** is good!

Tree inference accuracy (KF)

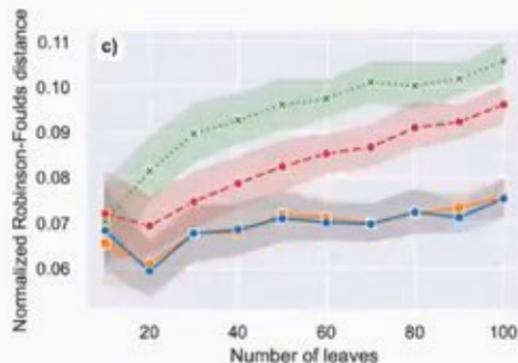
- Fairly **competitive** even on simple LG+GC model
- **Fast** because we use GPUs ¹

Runtime

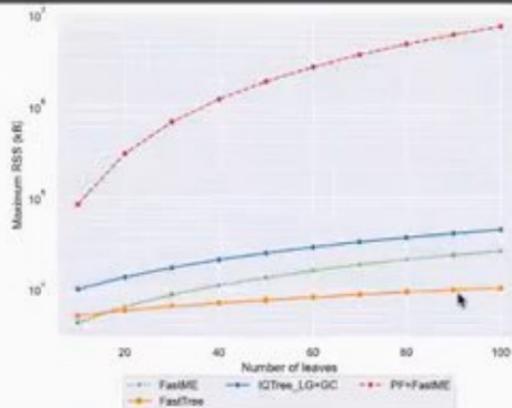
Nesterenko et al. 2025, ¹  Jean-Zay



Related Work - But also sometimes less good...



Topological accuracy (RF)

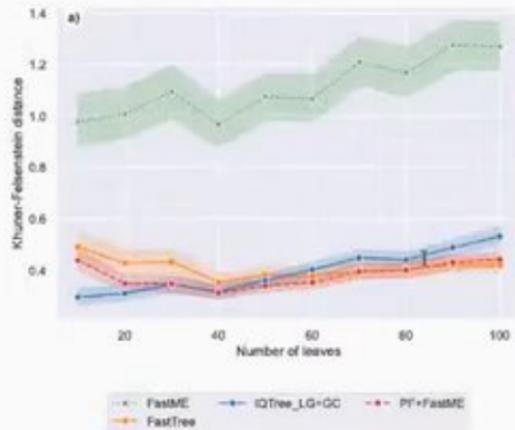


Memory usage

- **Gap** between PF and **ML** methods
- PF is **by far** the most **memory intensive**



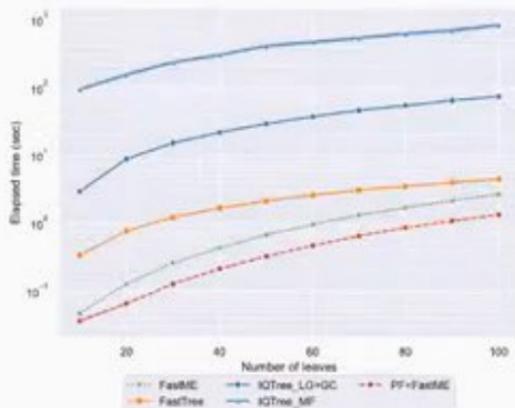
Related Work - **Phyloformer** is good!



Tree inference accuracy (KF)

- Fairly **competitive** even on simple LG+GC model
- **Fast** because we use GPUs ¹

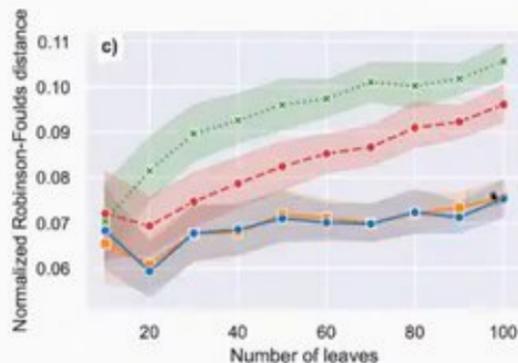
Nesterenko et al. 2025, ¹ 🏠 Jean-Zay



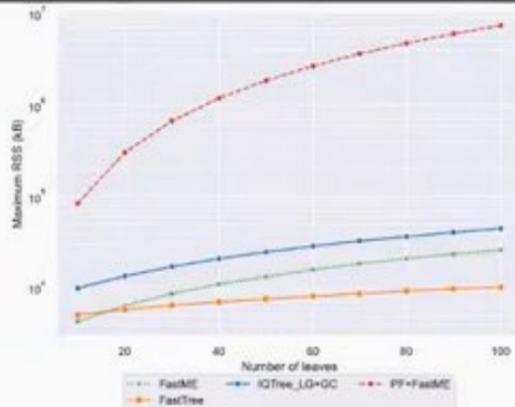
Runtime



Related Work - But also sometimes less good...



Topological accuracy (RF)



Memory usage

- **Gap** between PF and **ML** methods
- PF is **by far** the most **memory intensive**



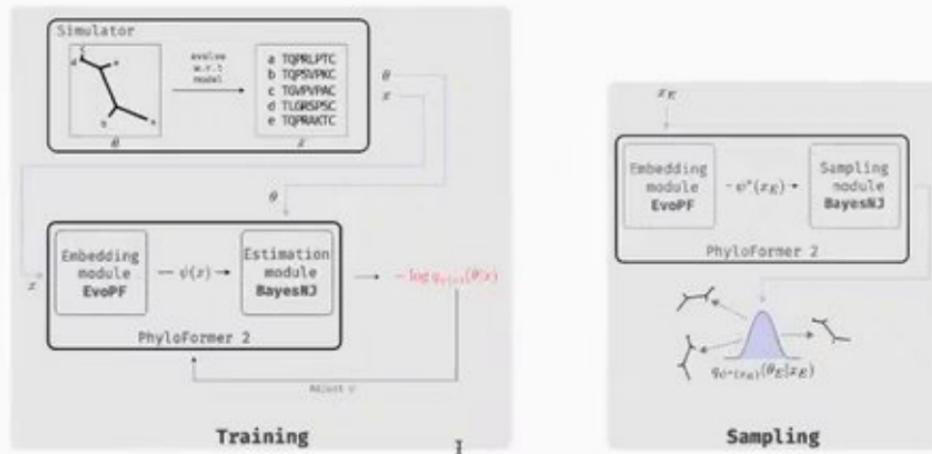
How to do phylogenetic inference end-to-end ?

Methods - Neural Posterior Estimation (NPE)

- Given a **probabilistic model** $p(x|\theta)$ with some prior $p(\theta)$
- We want to **estimate the posterior**: $p(\theta|x)$
- We build $q_\psi(\theta|x)$ a **family** of distributions **parametrized** by ψ (our NN)
- We find $q_{\psi^*} = \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{p(x)}[KL(q_\psi(\theta|x)||p(\theta|x))]$
- In practice we **maximize** $\mathbb{E}_{p(x,\theta)}[\log q_{\psi(x)}(\theta|x)]$ by **sampling** from $p(x, \theta)$

x : MSA, $\theta = (\tau, \ell)$: Phylogenetic tree, $\psi(x)$: NN applied to x

Methods - How do we do NPE?



- During **training** find $\psi^* = \underset{\psi}{\operatorname{argmin}} - \sum_i \log q_{\psi(x_i)}(\theta_i|x_i)$
- At **inference** time **sample** from: $q_{\psi^*(x_E)}(\theta_E|x_E)$ ¹

¹WIP: so for now only point-estimation

Methods - The EvoPF module, intro

the EvoPF module is an **adaptation** of the **EvoFormer** module from **AlphaFold2**. The tasks are **transpositions** of each other:

given input MSA ($n \times r$)

EvoFormer represent $r \times r$ relationships between sites

EvoPF represent $n \times n$ relationships between sequences

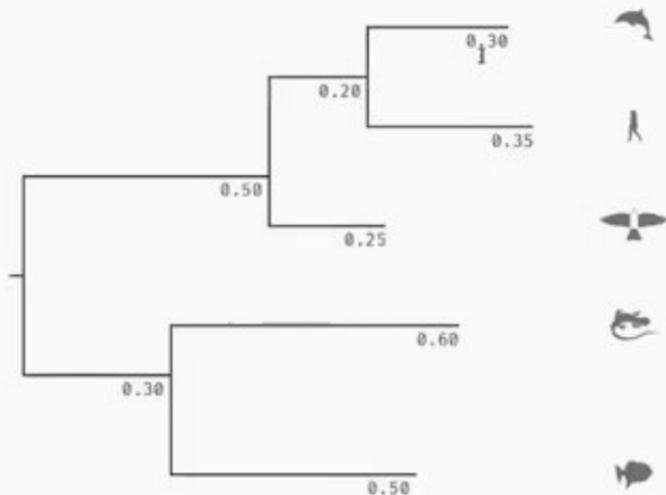
More expressive than MSA transformer

More lightweight than PF

Jumper et al. 2021; Rao et al. 2021

Methods - A tree is a series of merges

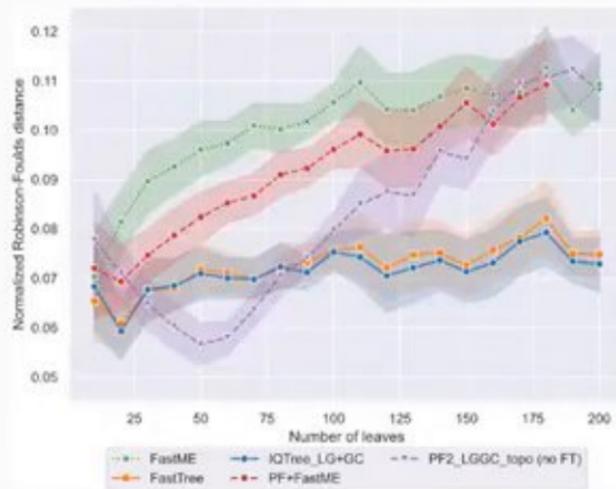
We want to describe the following tree:



Does it work ?

1

Results - Training topology only

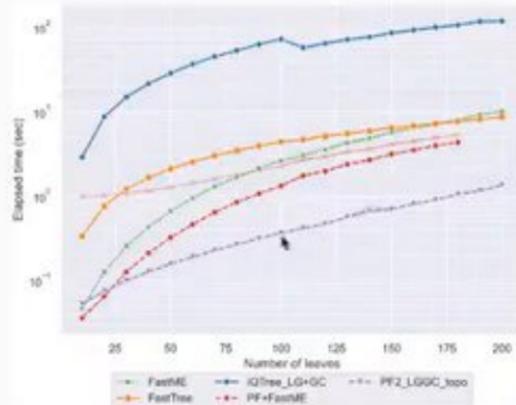


- **overfitting** on tree-size is an **issue**

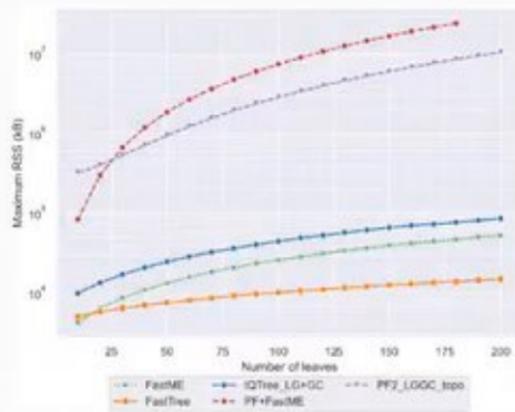
I

Same train set as PF1 paper: \approx 170k 50 seq LG+GC MSAs on rescaled BD trees

Results - Scalability



Execution time

Memory usage¹

¹ With 2x bigger sequence, and 4x bigger pair embeddings...

This is very much still a **work in progress...**

- Training on **more complex** data (*e.g. indels*) **increases** length-**overfitting**
- Learning **topology** and **branch-lengths** is also **challenging**
- How can we move **away** from **point-estimation** ?
- We might need to **adjust** our **priors** to compare with MCMC tools

Conclusion

Takeaways

- **Topologically** we manage to **beat** ML-methods¹ on LG
- While being **more scalable** than PF1
- Still needs some **work** for a fully **end-to-end** phylogenetic **inference** tool

What next ?

- Can we do **better** where computing $p(\theta|x)$ is **difficult** or **intractable**? (e.g. *Potts, epistasis, Selection, ...*)
- **Confident** this can **work** given our experience with **PF**

Prillo et al. 2023; Duchemin et al. 2023; Latrille et al. 2021

¹ Yay!

Thanks to:

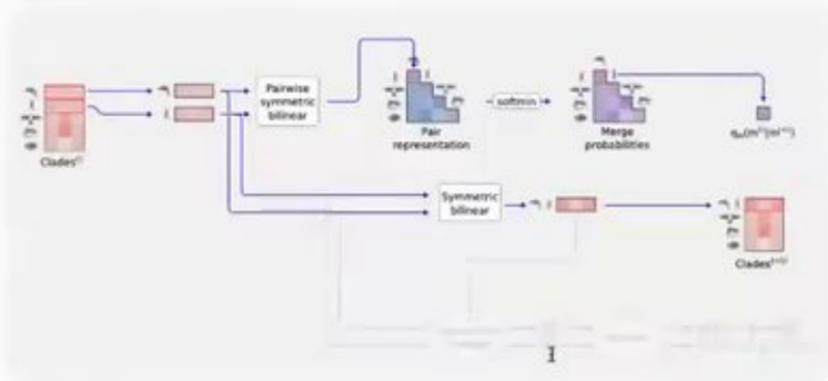
- **Luca Nesterenko**
- **Laurent Jacob**
- **Bastien Boussau**
- **Nicolas Lartillot**
- **Philippe Veber**
- **Vincent Garot**
- **Amélie Leroy**
- **Anybody that listened to me!**



Special thanks to Jean-Zay for all the GPUs!

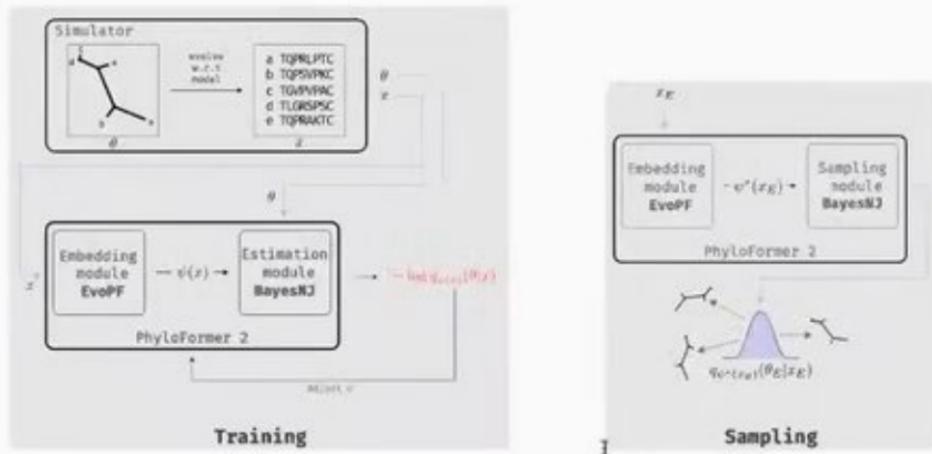
Methods - Bayesian, evaluating topological probabilities

$$m^{(i)} = (\sim, i) \quad (i^{(i)} = 0, \dots, i)$$



Update clade representation for next merge

Methods - How do we do NPE?



- During **training** find $\psi^* = \underset{\psi}{\text{argmin}} - \sum_i \log q_{\psi(x_i)}(\theta_i|x_i)$
- At **inference** time **sample** from: $q_{\psi^*(x_E)}(\theta_E|x_E)$ ¹

¹WIP: so for now only point-estimation

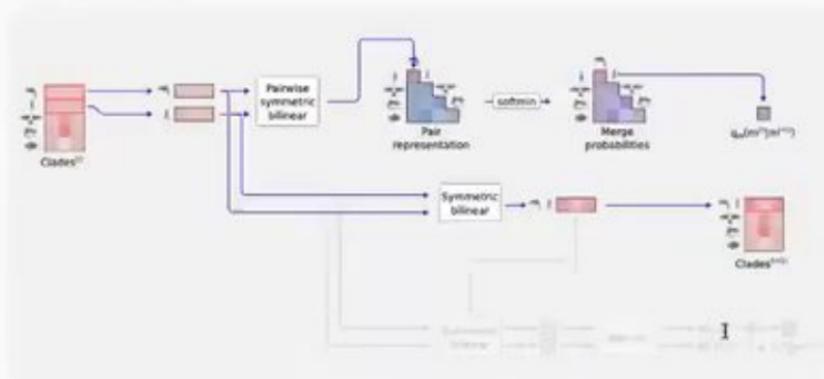
Methods - Neural Posterior Estimation (NPE)

- Given a **probabilistic model** $p(x|\theta)$ with some prior $p(\theta)$
- We want to **estimate the posterior**: $p(\theta|x)$
- We build $q_\psi(\theta|x)$ a **family** of distributions **parametrized** by ψ (our NN)
- We find $q_{\psi^*} = \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{p(x)}[KL(q_\psi(\theta|x)||p(\theta|x))]$
- In practice we **maximize** $\mathbb{E}_{p(x,\theta)}[\log q_{\psi(x)}(\theta|x)]$ by **sampling** from $p(x, \theta)$

x : MSA, $\theta = (\tau, \ell)$: Phylogenetic tree, $\psi(x)$: NN applied to x

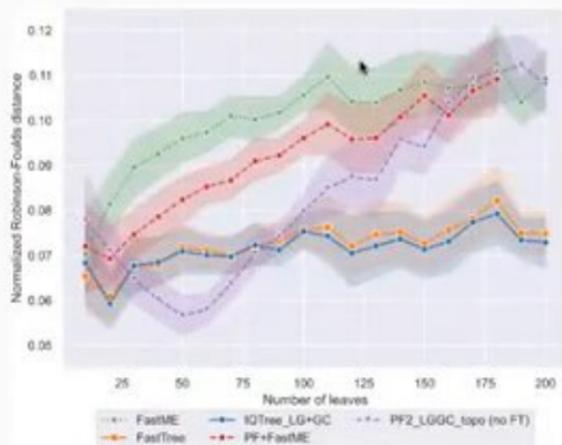
Methods - Bayesian, evaluating topological probabilities

$$m^{(i)} = (\tau, i) \quad (i^2 = 0, i)$$



Update clade representation for next merge

Results - Training topology only



- **overfitting on tree-size is an issue**

MLCB Schedule 2025

File Edit View Insert Format Data

Shubham Chaudhary (Heimholz)

	A	B
1	Slot	Poster
2	Spotlight day 1	1
3	Spotlight day 1	2
4	Spotlight day 1	3
5	Spotlight day 1	4
6	Spotlight day 1	5
7	Spotlight day 1	6
8	Spotlight day 1	7
9	Spotlight day 2	8
10	Spotlight day 2	9
11	Spotlight day 2	10
12	Spotlight day 2	11
13	Spotlight day 2	12
14	Spotlight day 2	13
15	Spotlight day 2	14
16	Spotlight day 2	15
17	Day 1	16
18	Day 1	17
19	Day 1	18
20	Day 1	19
21	Day 1	20
22	Day 1	21
23	Day 1	22
24	Day 1	23
25	Day 1	24
26	Day 1	25
27	Day 1	26
28	Day 1	27
29	Day 1	28
30	Day 1	29

MLCB oral slides

Name	Date Modified	Size	Kind
41_Linder Johannes Linder.pptx	Sep 9, 2025 at 10:21 PM	3.8 MB	PowerP... (.pptx)
42_Fanjiang Clara Fanjiang.key	Yesterday at 12:45 PM	16.1 MB	Keynote
43_Shearer Courtney Shearer.pptx	Sep 9, 2025 at 6:51 PM	2.8 MB	PowerP... (.pptx)
114_Rocha Jose Felipe Rocha.pptx	Today at 9:06 AM	12.7 MB	PowerP... (.pptx)
134_shaw Peter Shaw.pdf	Yesterday at 10:12 PM	1.2 MB	PDF Document
954_ILASSEL Luc Blassel.pdf	Yesterday at 11:55 PM	3.8 MB	PDF Document
2025_MLCB Keynote Jacob Schreiber.pptx	Today at 8:45 AM	28.3 MB	PowerP... (.pptx)
new_MLCB_2025 Barbara.pptx	Yesterday at 1:37 PM	477.9 MB	PowerP... (.pptx)
MLCB (15 min) Alan Amin.key	Yesterday at 11:58 AM	4.5 MB	Keynote
MLCB_2025_Battle Alexis Battle.pdf	Yesterday at 9:46 AM	6.2 MB	PDF Document
MLCB_2025_Battle Alexis Battle.pptx	Yesterday at 9:09 AM	37 MB	PowerP... (.pptx)
MLCB_2025_v2 Danielle Stevens.pptx	Today at 11:41 AM	92.3 MB	PowerP... (.pptx)
Perturbation_Benchmark_Hasana(Euxhan Hasana).pdf	Today at 9:33 AM	1.8 MB	PDF Document

Enhancing Bulk RNA-Seq Deconvolution Using Atlas-Level Deep Learning Embeddings

Decoding translation regulation at the 5' untranslated region of mRNA with augmented multi-species massively parallel reporter assay and active learning

Convert to table

AI Bookmarks

2:33 AM

applying

1:37 AM

boom

- 1 SPLISOSM: Kernel-based Multivariate Independence Testing for Mapping Spatial Transcript Diversity
- 2
- 3
- 4

SPLISOSM: Kernel-based Multivariate Independence Testing for Mapping Spatial Transcript Diversity

Jiayu Su
js5756@cumc.columbia.edu
Columbia Systems Biology & NYGC
MLCB 2025
Sept 11, 2025



Click to add notes

Restarting at 11.55am

mlcb.org for schedule

Zoom

SPLISOSM: Kernel-based Multivariate Independence Testing for Mapping Spatial Transcript Diversity

Jiayu Su

js5756@cumc.columbia.edu

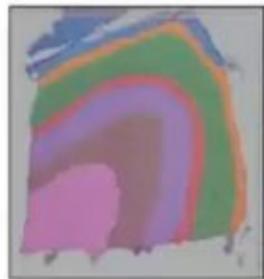
Columbia Systems Biology & NYGC

MLCB 2025

Sept 11, 2025



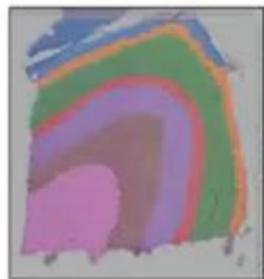
Identifying spatial patterns from high-dimensional data is hard



Human DLPFC



Identifying spatial patterns from high-dimensional data is hard



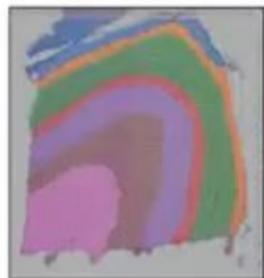
Human DLPFC



Human glioblastoma



Identifying spatial patterns from high-dimensional data is hard



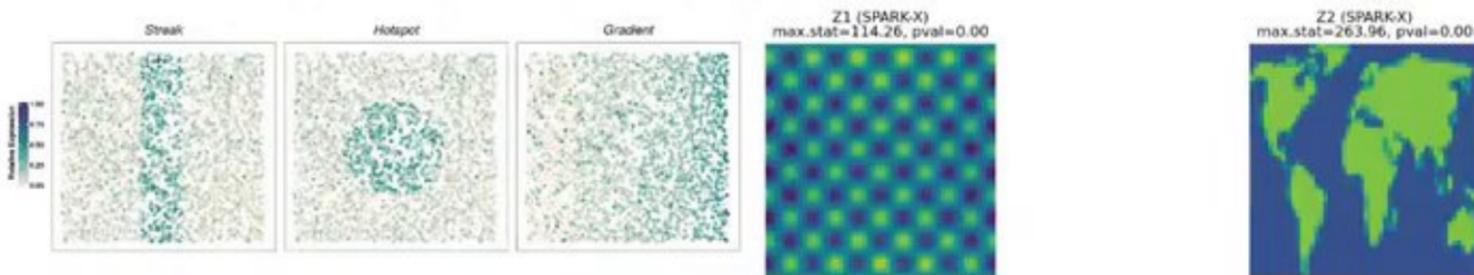
Human DLPFC



Human glioblastoma



Never blindly trust your lens: Pattern detection may not be robust to perturbation

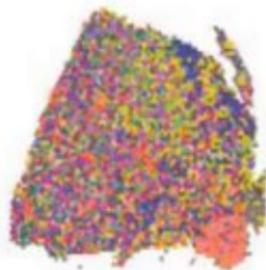


Zhu, Jiaqiang, Shiquan Sun, and Xiang Zhou. *Genome biology* 22.1 (2021): 184.

Identifying spatial patterns from high-dimensional data is hard



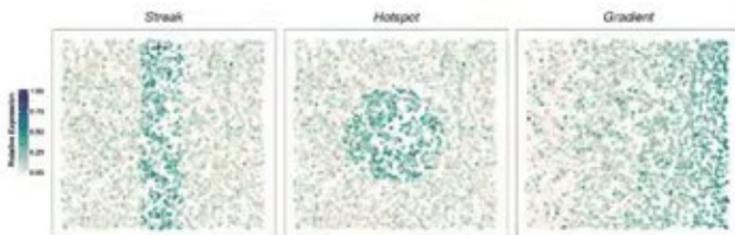
Human DLPFC



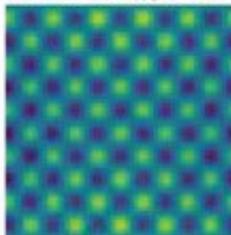
Human glioblastoma



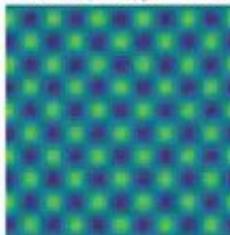
Never blindly trust your lens: Pattern detection may not be robust to perturbation



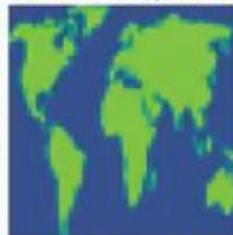
Z1 (SPARK-X)
max_stat=114.26, pval=0.00



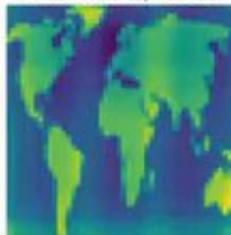
Z1-transformed (SPARK-X)
max_stat=0.00, pval=1.00



Z2 (SPARK-X)
max_stat=263.96, pval=0.00



Z2-transformed (SPARK-X)
max_stat=0.00, pval=1.00



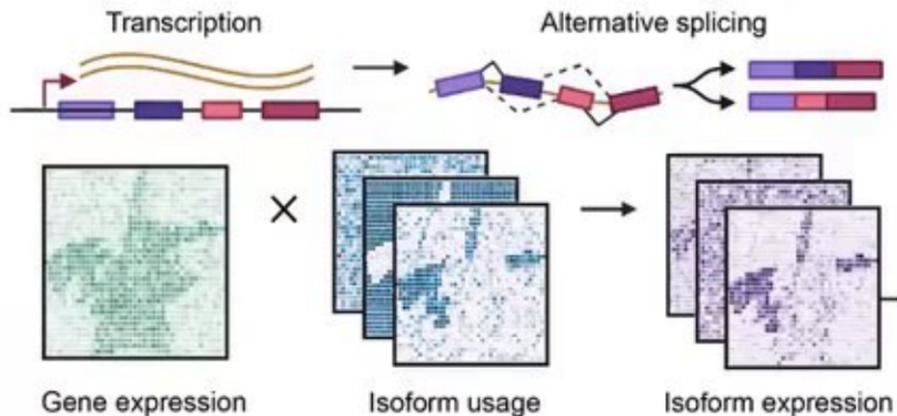
Zhu, Jiaqiang, Shiquan Sun, and Xiang Zhou. *Genome biology* 22.1 (2021): 184.

A new test to detect arbitrary spatial variability patterns

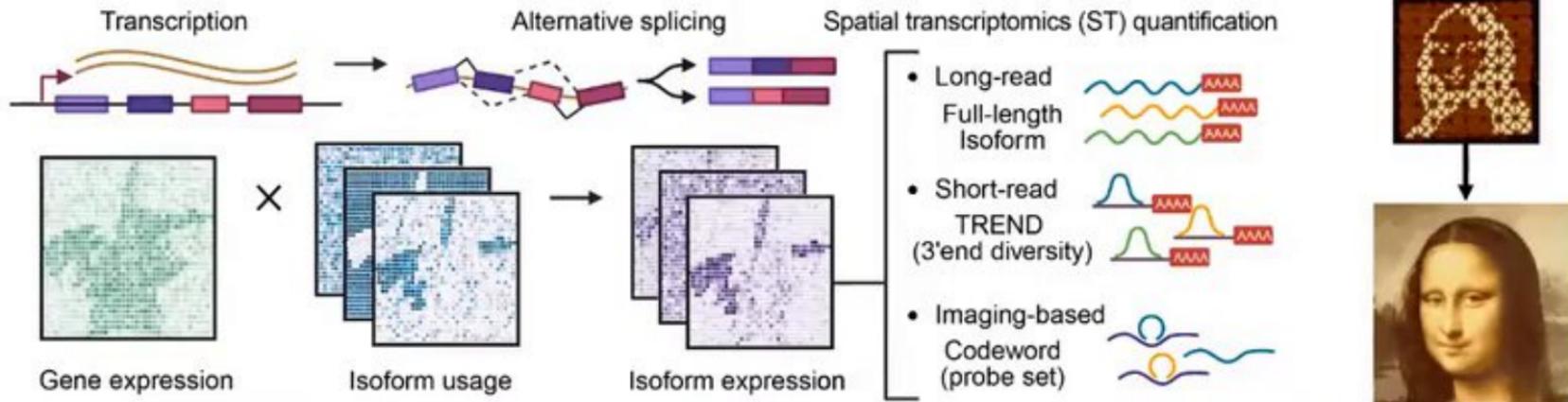
- **The conflict:** Space of spatial patterns is of dim $\#spots$, and we want all of them



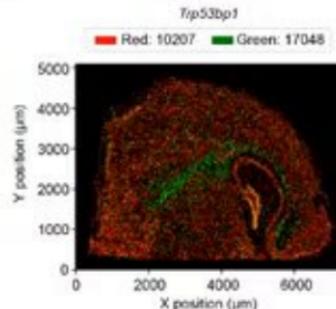
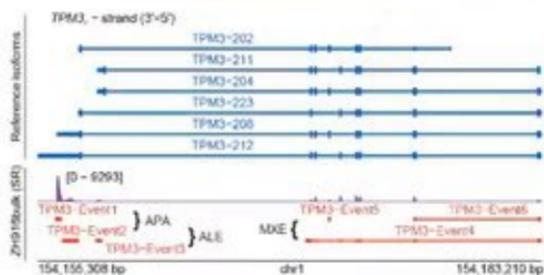
From univariate gene-centric to multivariate isoform-centric gold mines



From univariate gene-centric to multivariate isoform-centric gold mines



Through **reanalysis of existing datasets**, we revealed thousands of conserved spatial transcript diversity in healthy and malignant brain tissues and uncovered complex regulatory relationships

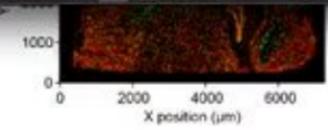
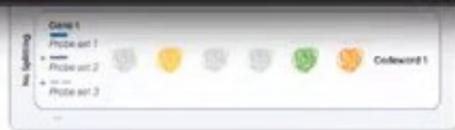
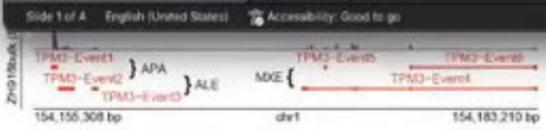


- 1 GeST: Towards Building A Generative Pretrained Transformer for Learning Cellular Spatial Context
- 2 GeST: Overview
- 3 GeST: Architecture
- 4 GeST: Results

GeST: Towards Building A Generative Pretrained Transformer for Learning Cellular Spatial Context

Contact: Minsheng Hao (hmsh653@gmail.com)

2025/9/11 1



mines



ty in
er #8

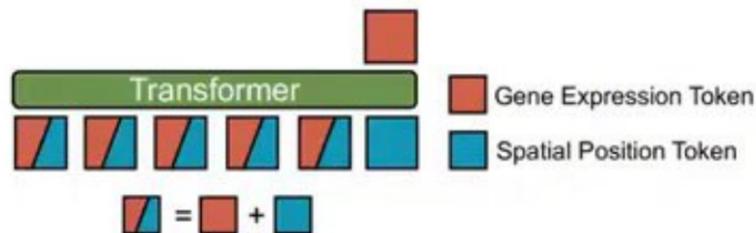


GeST: Towards Building A Generative Pretrained Transformer for Learning Cellular Spatial Context

Contact: Minsheng Hao (hmsh653@gmail.com)

Spatial-aware tokenization

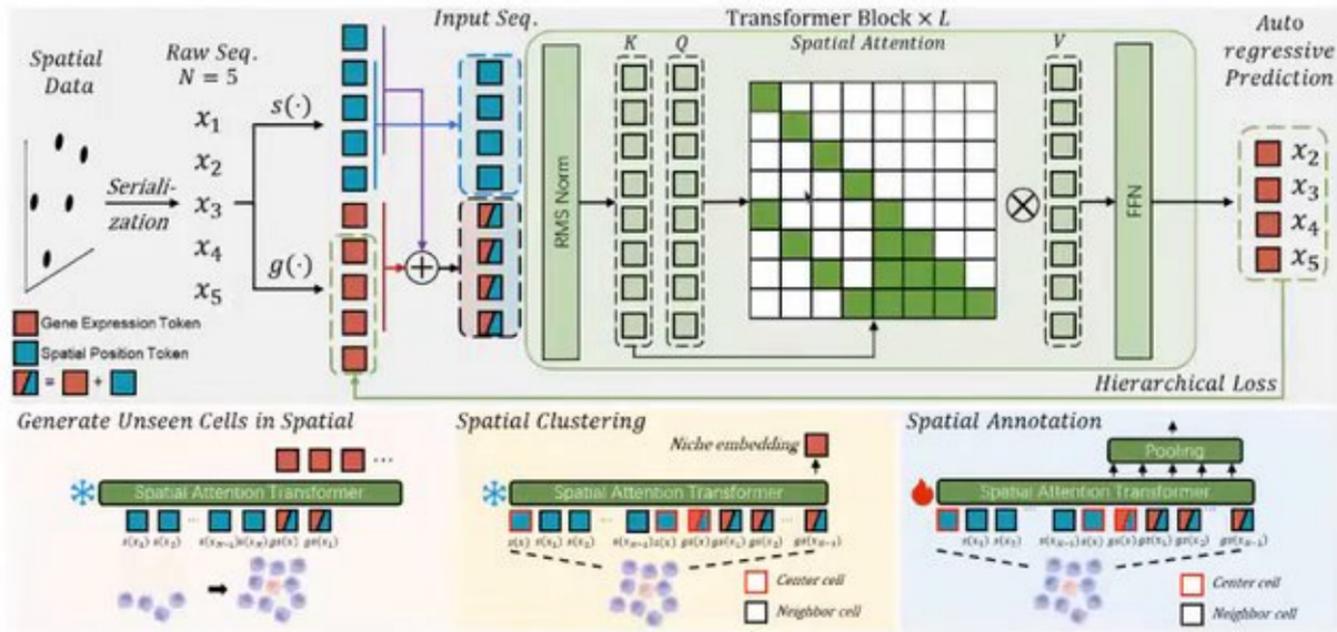
- Understanding spatial cell context is crucial, but current models is hard to model spatial dynamics
- GeST is a generative transformer that learns spatial context by predicting a cell from its neighbors.



Input = neighbors, target cell's spatial loc.

output = target cell's gene profile

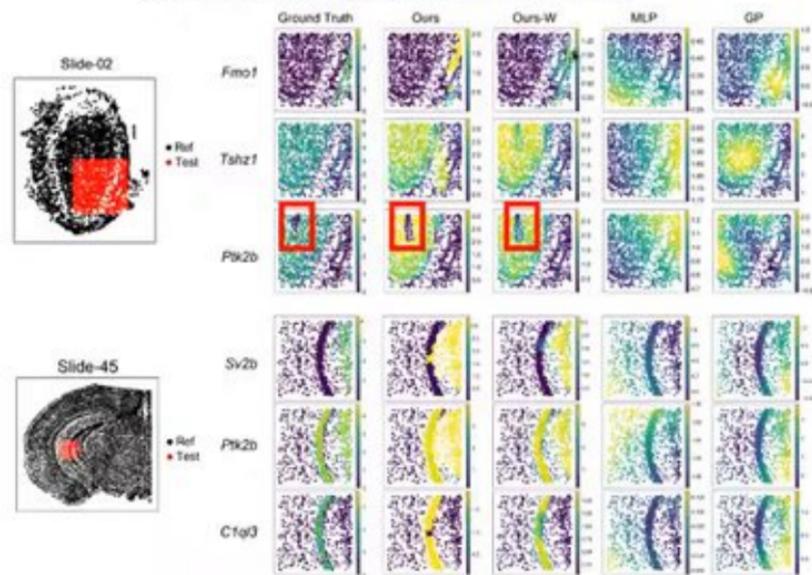
- Spatial serialization + Cell tokenization + Spatial attention
- Pretraining across multiple spatial transcriptomics datasets



- **Application:** Spatial generation + Spatial clustering + Spatial Annotation

- First **generative** pretrained model for spatial transcriptomics.
- Enables new applications like **in-silico perturbation analysis**

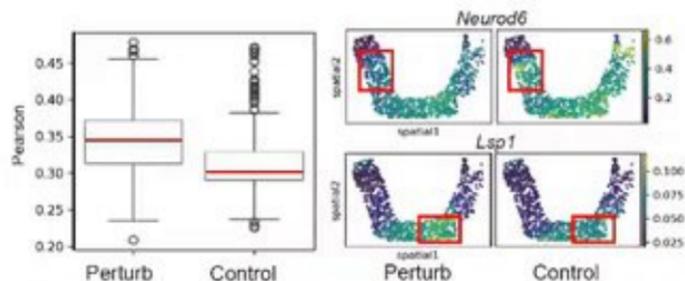
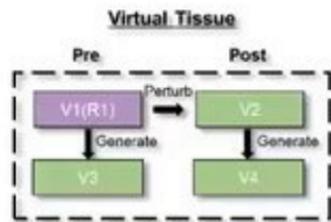
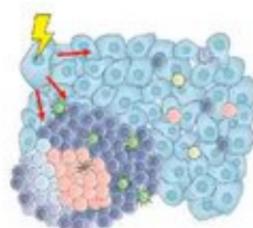
Generating unseen region (in red)



Method	Anterior Brain			MERFISH Mid Brain			Posterior Brain			Visium Primary Liver Cancer			Stereo-seq Sagittal Brain		
	Rall _i	R.200 _i	$\rho \uparrow$	Rall _i	R.200 _i	$\rho \uparrow$	Rall _i	R.200 _i	$\rho \uparrow$	Rall _i	R.200 _i	$\rho \uparrow$	Rall _i	R.200 _i	$\rho \uparrow$
Ours	1.374	1.296	0.301	1.379	1.319	0.265	1.386	1.340	0.242	1.303	1.126	0.540	1.405	1.357	0.326
Ours-W	1.352	1.244	0.340	1.367	1.288	0.302	1.376	1.322	0.275	1.320	1.150	0.499	1.399	1.340	0.323
GP	1.379	1.289	0.241	1.389	1.337	0.241	1.393	1.354	0.214	1.357	1.264	0.272	1.413	1.405	0.073
MLP	1.399	1.369	0.319	1.404	1.393	0.283	1.406	1.397	0.267	1.347	1.174	0.491	1.403	1.357	0.314

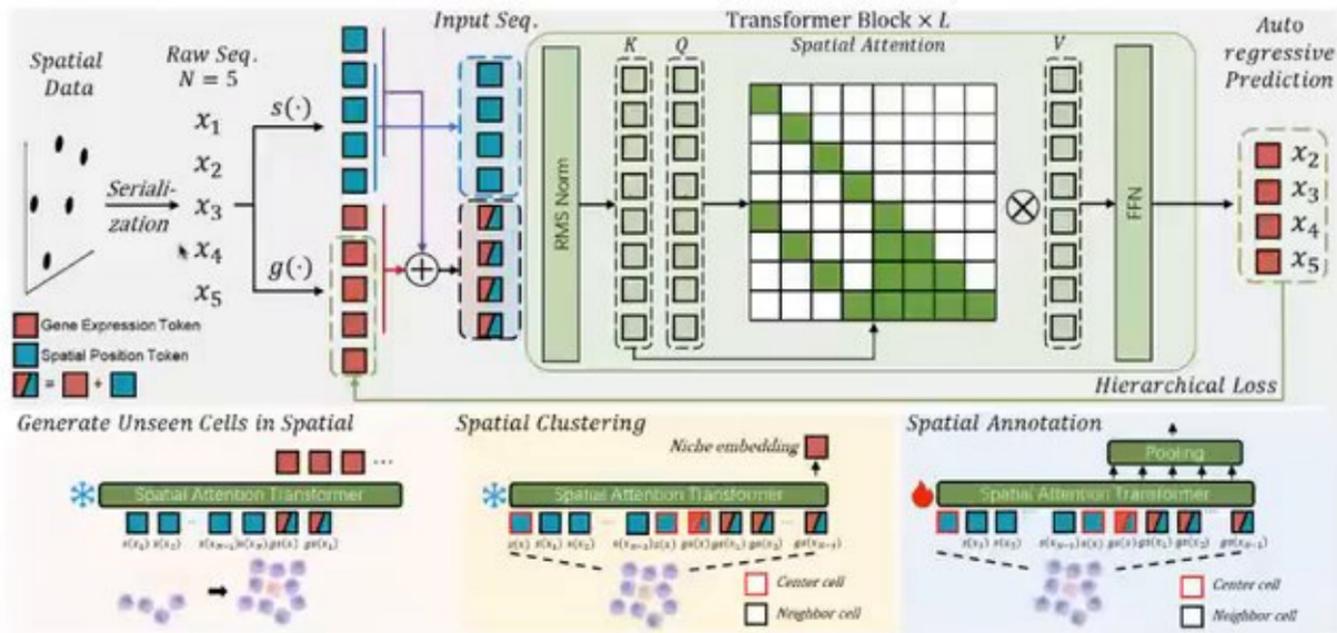
Spatial Perturbation

- Reference: A real experiment to study ischemic brain
- ✓ **Regenerate** the ischemic perturbation via GeST



- **Spatial serialization + Cell tokenization + Spatial attention**

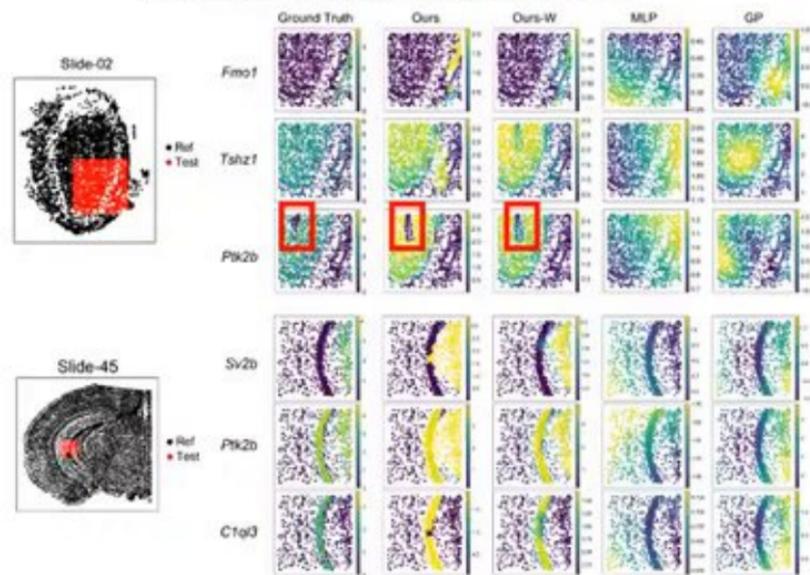
- Pretraining across multiple spatial transcriptomics datasets



- **Application: Spatial generation + Spatial clustering + Spatial Annotation**

- First **generative** pretrained model for spatial transcriptomics.
- Enables new applications like **in-silico perturbation analysis**

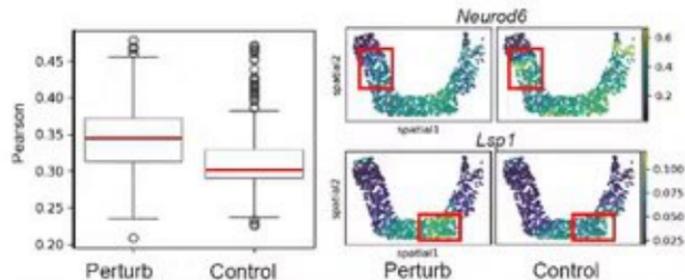
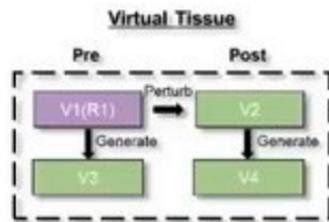
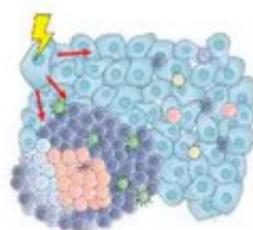
Generating unseen region (in red)



Method	Anterior Brain			MERFISH Mid Brain			Posterior Brain			Visium Primary Liver Cancer			Stereo-seq Sagittal Brain		
	Rall ₁	R.200 ₁	$\rho \uparrow$	Rall ₁	R.200 ₁	$\rho \uparrow$	Rall ₁	R.200 ₁	$\rho \uparrow$	Rall ₁	R.200 ₁	$\rho \uparrow$	Rall ₁	R.200 ₁	$\rho \uparrow$
Ours	1.374	1.296	0.301	1.379	1.319	0.265	1.386	1.340	0.242	1.303	1.126	0.540	1.405	1.357	0.326
Ours-W	1.352	1.244	0.340	1.367	1.288	0.302	1.376	1.322	0.275	1.320	1.150	0.499	1.399	1.340	0.323
GP	1.379	1.289	0.241	1.389	1.337	0.241	1.393	1.354	0.214	1.357	1.264	0.272	1.413	1.405	0.073
MLP	1.399	1.369	0.319	1.404	1.393	0.283	1.406	1.397	0.267	1.347	1.174	0.491	1.403	1.357	0.314

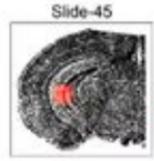
Spatial Perturbation

- Reference: A real experiment to study ischemic brain
- ✓ **Regenerate** the ischemic perturbation via GeST



• First generative pretrained model for spatial transcriptomics.
 • Enables new applications like in-silico perturbation analysis

Generat



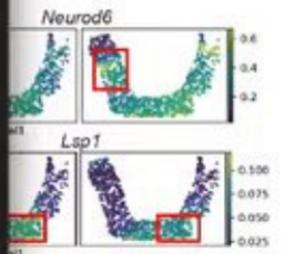
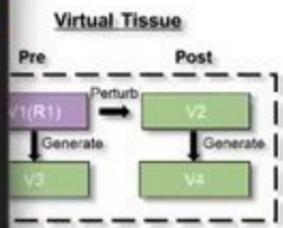
Method	Anterior Brain				Rall ₁	R ₂₀₀	ρ ↑	Rall ₁	R ₂₀₀	ρ ↑	Rall ₁	R ₂₀₀	ρ ↑	Rall ₁	R ₂₀₀	ρ ↑	
	Rall ₁	R ₂₀₀	ρ ↑	Rall ₁													R ₂₀₀
Ours	1.374	1.296	0.301	1.379													
Ours-W	1.352	1.244	0.340	1.367	1.288	0.302	1.376	1.322	0.275	1.320	1.150	0.499	1.399	1.340	0.323		
GP	1.379	1.289	0.241	1.389	1.337	0.241	1.393	1.354	0.214	1.357	1.264	0.272	1.413	1.405	0.073		
MLP	1.399	1.369	0.319	1.404	1.393	0.283	1.406	1.397	0.267	1.347	1.174	0.491	1.403	1.357	0.314		

MLCB2025 spotlight slides

Name	Date Modified	Size	Kind
1_Schreiber Jacob Schreiber.pptx	Sep 9, 2025 at 3:10 PM	1.7 MB	PowerP...n (.pptx)
2_Seitz Evan Seitz.pdf	Sep 9, 2025 at 8:47 PM	2.6 MB	PDF Document
3_Nair Surag Nair.pdf	Sep 9, 2025 at 2:44 PM	2.2 MB	PDF Document
4_Talk_mahbuba Mahbuba Tasmin.pptx	Sep 9, 2025 at 5:00 PM	4.2 MB	PowerP...n (.pptx)
5_Ayub Shanza Ayub.pptx	Sep 9, 2025 at 11:03 AM	9.8 MB	PowerP...n (.pptx)
6_aideryev Pavel Aideryev.pdf	Yesterday at 12:23 AM	586 KB	PDF Document
7_Liu Simon Liu.pptx	Sep 9, 2025 at 9:08 PM	9 MB	PowerP...n (.pptx)
8_Su ZY Su.pptx	Today at 1:04 AM	17.2 MB	PowerP...n (.pptx)
9_Minsheng&Tianyu Tianyu Liu.pptx	Sep 8, 2025 at 12:09 AM	1.7 MB	PowerP...n (.pptx)
10_Lam Kevin Lam.pdf	Sep 9, 2025 at 6:19 PM	1.6 MB	PDF Document
11_Chaudhary Shubham Chaudhary.pptx	Sep 9, 2025 at 4:22 PM	21 MB	PowerP...n (.pptx)
12_Demathelin Antoine de Mathelin.pdf	Sep 9, 2025 at 2:15 PM	3.7 MB	PDF Document
14_Sergio Mares.key	Yesterday at 12:50 AM	1.7 MB	Keynote
16_Datal Saykhoom Dalal.key	Sep 9, 2025 at 8:48 PM	5.1 MB	Keynote
MLCB2025_Mares Sergio Mares.pptx	Yesterday at 12:36 AM	356 KB	PowerP...n (.pptx)
MLCB2025_Mares_Final Sergio Mares.pptx	Yesterday at 12:55 AM	1.8 MB	PowerP...n (.pptx)

on

study ischemic brain
 perturbation via GeST



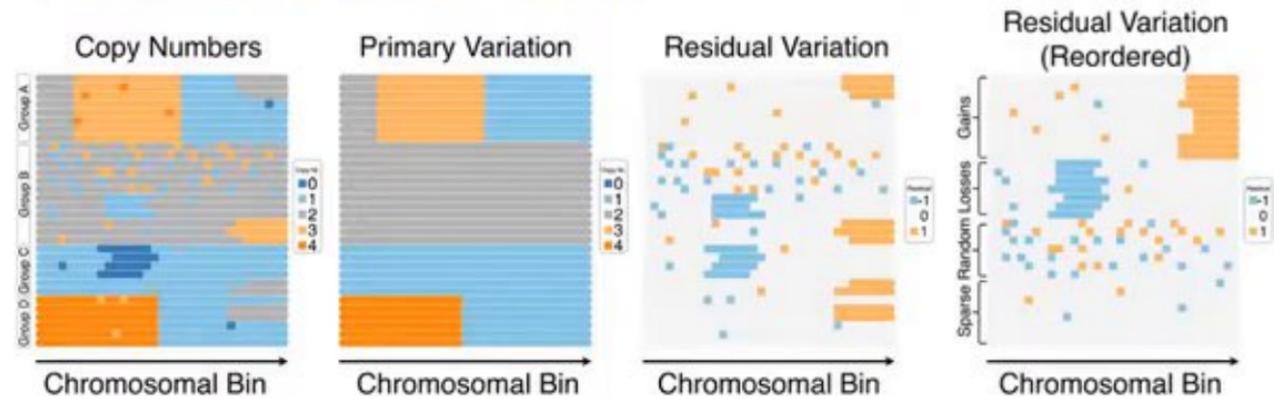
Perturb Control Perturb Control



CN-SBM: Categorical Block Modelling For Primary and Residual Copy Number Variation

Kevin Lam^{1,2}, William Daniels², J Maxwell Douglas², Daniel Lai², Samuel Aparicio², Benjamin Bloem-Reddy¹, Yongjin Park^{1,2}

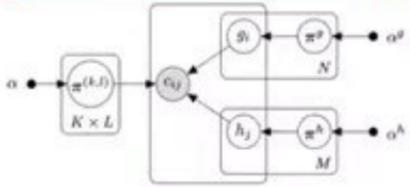
¹Department of Statistics, University of British Columbia; ²Department of Molecular Oncology, BC Cancer Research Centre



Copy Number Variation (CNV): alterations in the number of copies of DNA segments (amplifications and deletions)

- **Primary (main) variation:** large-scale, recurrent alterations
- **Residual variation:** finer sample-specific deviations

CN-SBM Model

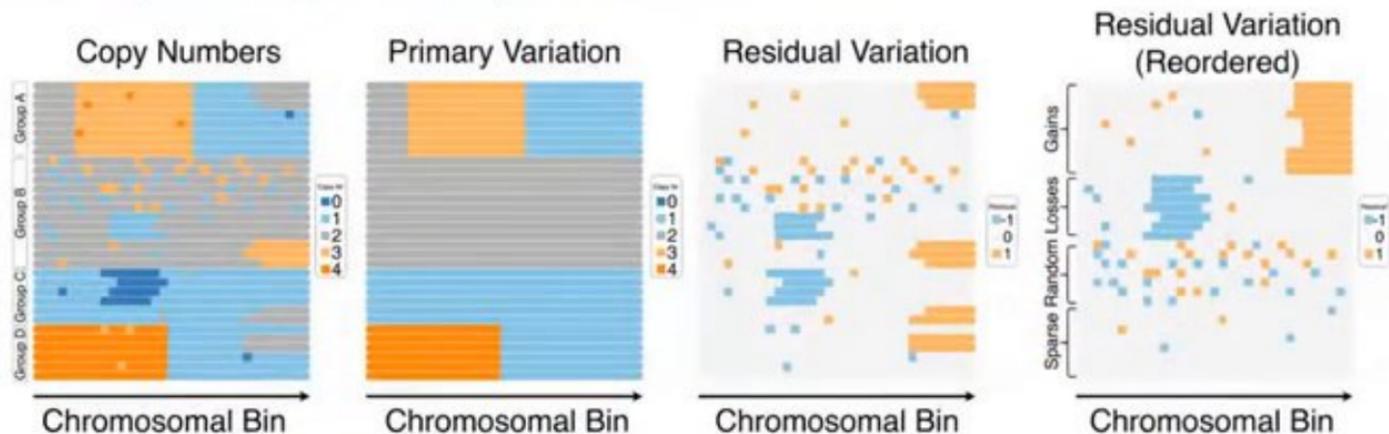


- Explicitly models categorical CN states ($c_{ij} \in \{0, 1, \dots, 10, \geq 11\}$)
- Jointly clusters samples and genomic bins

CN-SBM: Categorical Block Modelling For Primary and Residual Copy Number Variation

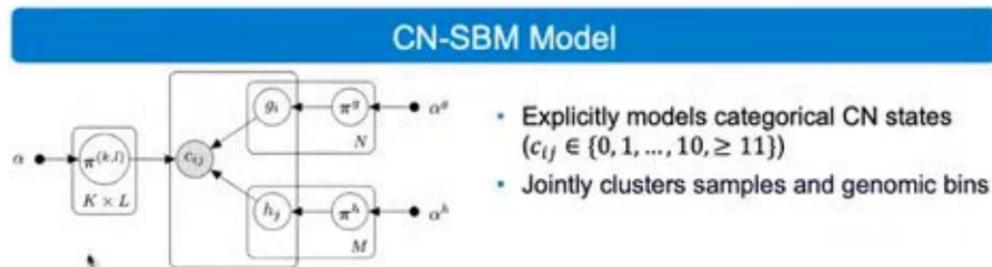
Kevin Lam^{1,2}, William Daniels², J Maxwell Douglas², Daniel Lai², Samuel Aparicio², Benjamin Bloem-Reddy¹, Yongjin Park^{1,2}

¹ Department of Statistics, University of British Columbia; ² Department of Molecular Oncology, BC Cancer Research Centre



Copy Number Variation (CNV): alterations in the number of copies of DNA segments (amplifications and deletions)

- **Primary (main) variation:** large-scale, recurrent alterations
- **Residual variation:** finer sample-specific deviations



CN-SDM Detects Clinically Relevant Clusters in Low-Grade Glioma

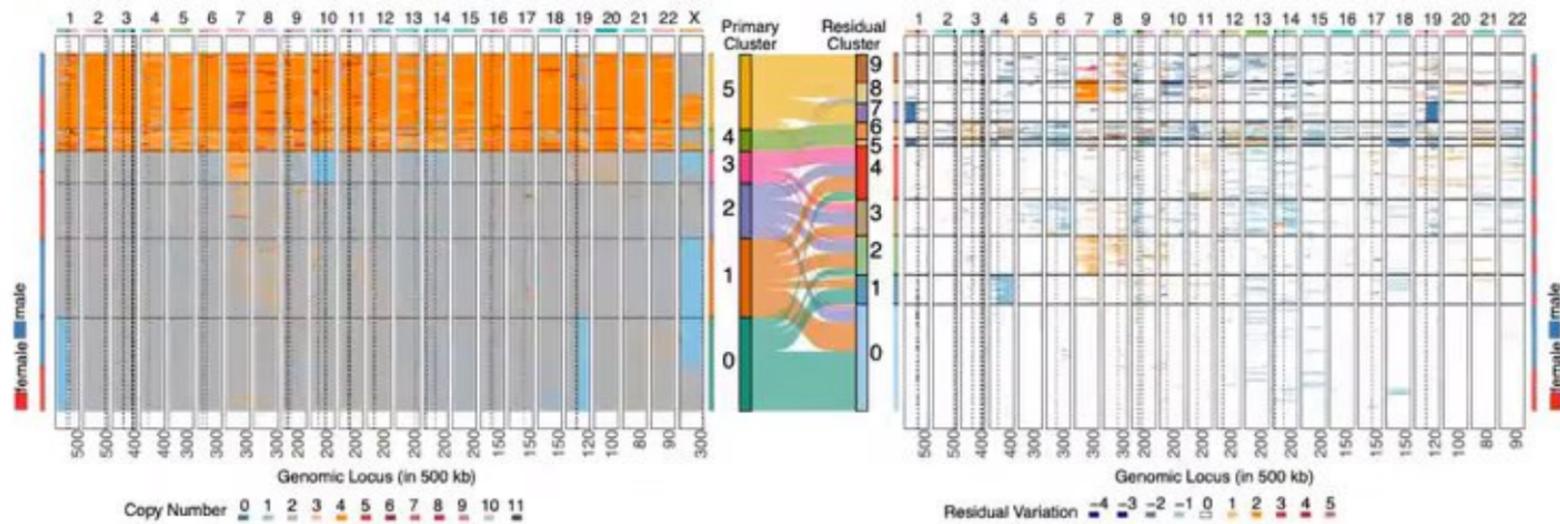


Figure: Primary and residual CNVs in low-grade glioma TCGA-cohort ($N = 490$)

- **Primary clusters:** arm- or whole-chromosome alterations (chr. 1, 7, 10, 19, X)
- **Residual clusters:** sparse, fine alterations

- 1
- 2
- 3
- 4
- 5
- 6

AI-based histopathology phenotyping reveals germline loci shaping breast cancer morphology

Shubham Chaudhary
*Institute of AI for Health
 Helmholtz Munich & TUM*

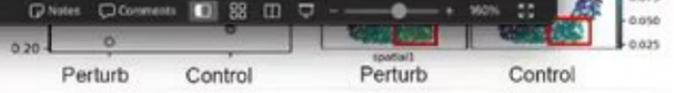
MLCB
 11th Sep 2025

**HELMHOLTZ
 MUNICH**

Click to add notes

Side 1 of 7 English (United States) Accessibility: Investigate

	Rall	R.2000	ρ												
Ours	1.374	1.296	0.301	1.379	1.319	0.265	1.386	1.340	0.242	1.303	1.126	0.540	1.405	1.357	0.326
Ours-W	1.352	1.244	0.340	1.367	1.288	0.302	1.376	1.322	0.275	1.320	1.150	0.499	1.399	1.340	0.323
GP	1.379	1.289	0.241	1.389	1.337	0.241	1.393	1.354	0.214	1.357	1.264	0.272	1.413	1.405	0.073
MLP	1.399	1.369	0.319	1.404	1.393	0.283	1.406	1.397	0.267	1.347	1.174	0.491	1.403	1.357	0.314



ic brain
 ST
 Post
 V2
 Generate
 V4
 0.6
 0.4
 0.2
 0.050
 0.075
 0.100
 0.125
 0.150
 0.175
 0.200
 0.225
 0.250
 0.275
 0.300
 0.325
 0.350
 0.375
 0.400
 0.425
 0.450
 0.475
 0.500
 0.525
 0.550
 0.575
 0.600

AI-based histopathology phenotyping reveals germline loci shaping breast cancer morphology

Shubham Chaudhary
*Institute of AI for Health
Helmholtz Munich & TUM*

MLCB
11th Sep 2025

Background: Genome-wide association studies (GWAS)

GWAS is a method used to identify genetic variants that are associated with a specific trait or disease.

- The study compares the genetic variations across large groups of individuals with the presence of a trait / disease
- A statistical testing procedure is used to identify specific genetic variants (typically SNPs) associated with the trait or disease
- These SNPs are then used as markers to identify other genetic variations that may be contributing to the trait or disease

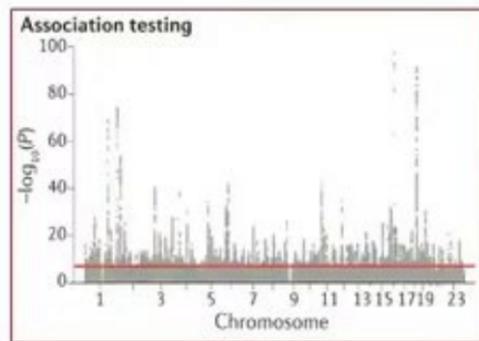
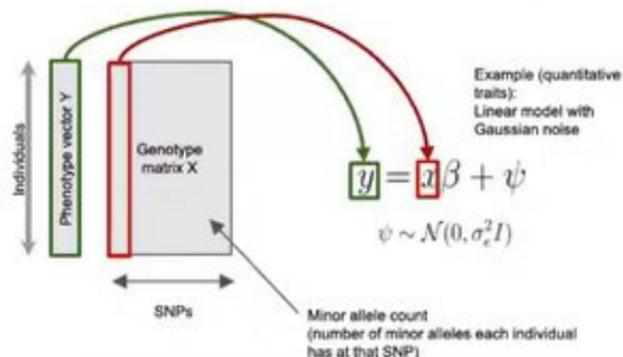


Figure adapted from Uffelmann et al, Nat Rev Methods Primers, 2021

Genetics

Expression

Disease

```

%CTGAAACTGGGGGACTGACGTGCAACG
%CTGCAACTGGGGGACTGACGTGCAACG
%CTGCAACTGGGGGACTGACGTGCAACG
%CTGAAACTGGGGGATTGACGTGCAACG
%CTGCAACTGGGGATTGACGTGCAACG
%CTGCAACTGGGGATTGACGTGCAACG
    
```



- GWAS have identified thousands of genetic variants associated with human complex traits and disease
- Challenge: We still lack understanding of biological processes and pathways affected by these variants
- Genetic analyses of intermediate phenotypes such as gene expression can help mitigate this by identifying tissues, cell types and genes affected by disease loci

Genetics

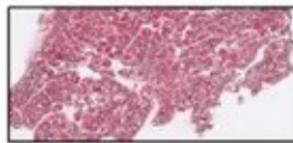
```

4CCTGAAACTGGGGGACTGACGTGGAACG
4CCTGCAACTGGGGGACTGACGTGCAACGC
4CCTGCAACTGGGGGACTGACGTGCAACGC
4CCTGAAACTGGGGGATTGACGTGGAACG
4CCTGCAACTGGGGGATTGACGTGCAACGC
4CCTGCAACTGGGGGATTGACGTGCAACGC
    
```

Expression



Medical imaging



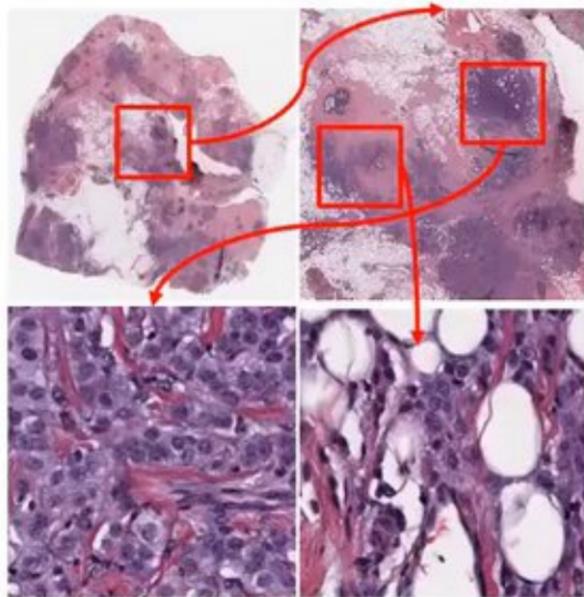
Disease

→ Can histological images reveal how germline variants shape tumor morphology?

H&E-stained whole slide images in cancer

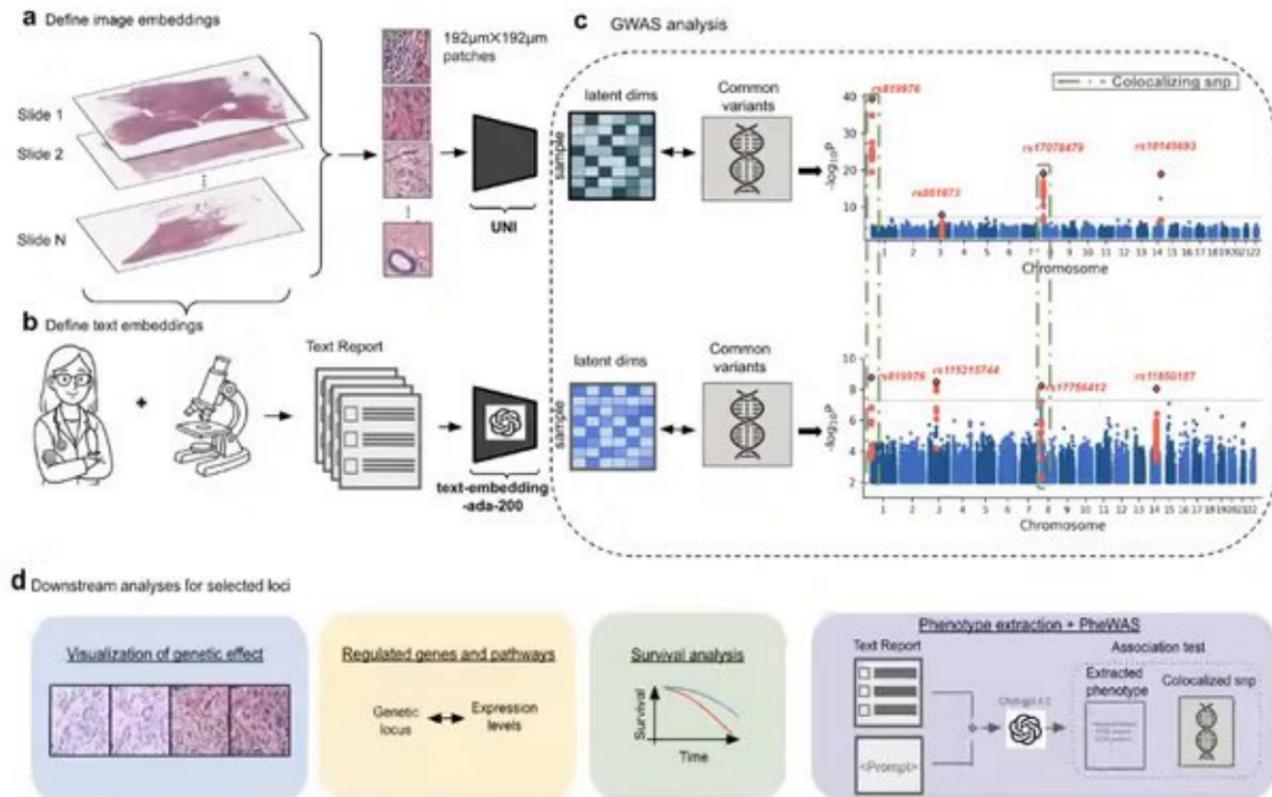
- Are microscopy images of tissue sections
- Used in the clinic for disease diagnosis (e.g., cancer, steatohepatitis)
- Showcase detailed cell and tissue phenotypes related to disease

Breast slide from TCGA



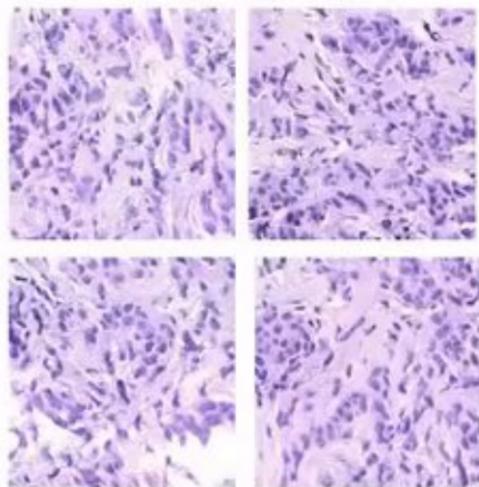
Our approach: Multimodal analysis

Our approach: Multimodal analysis

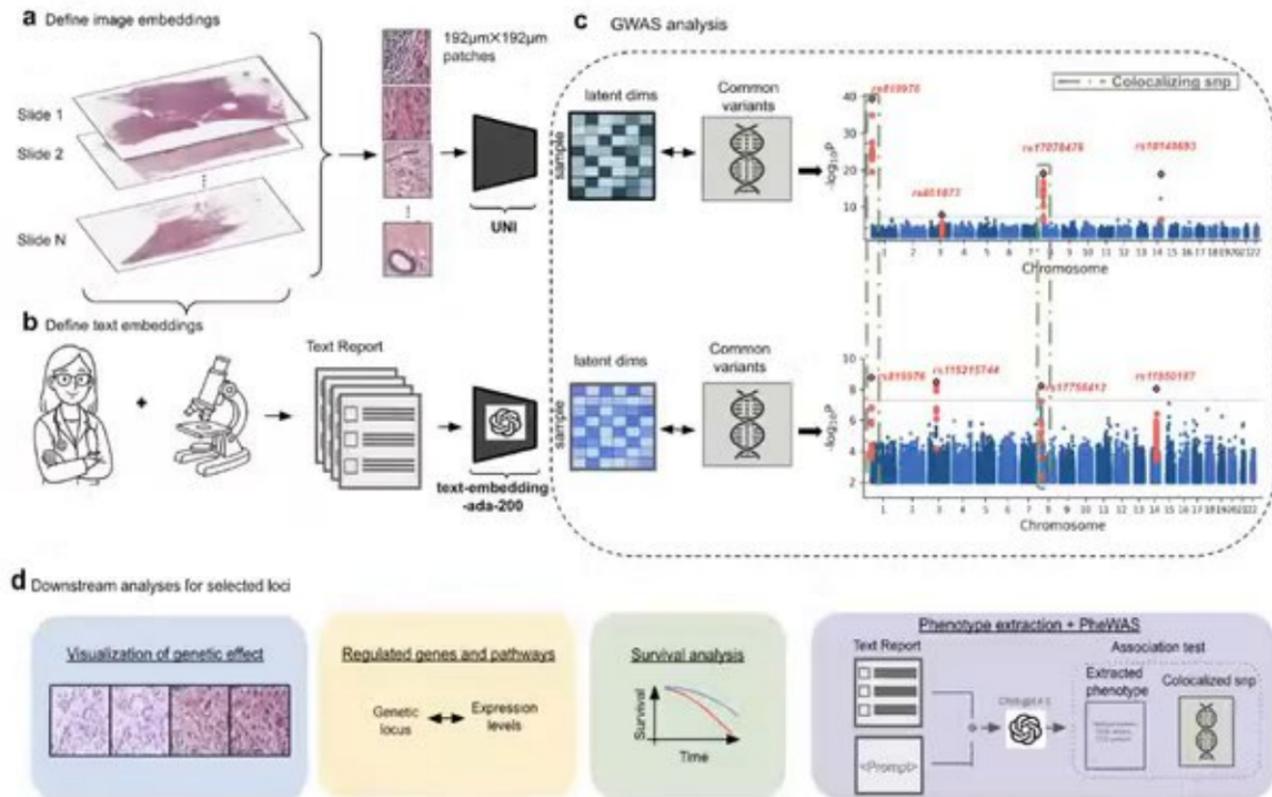


Characterization of genetic effect on tissue Histology

Genetic effect in tissue

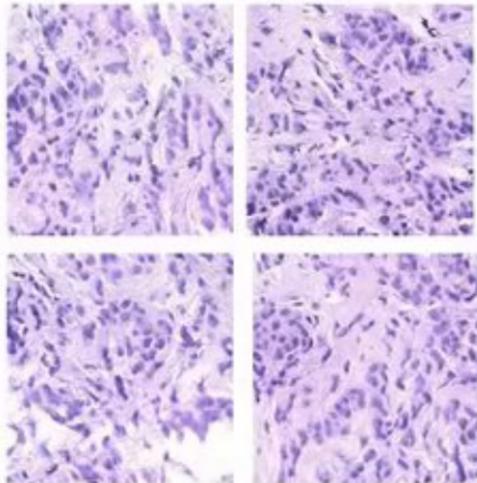


Our approach: Multimodal analysis



Characterization of genetic effect on tissue Histology

Genetic effect in tissue



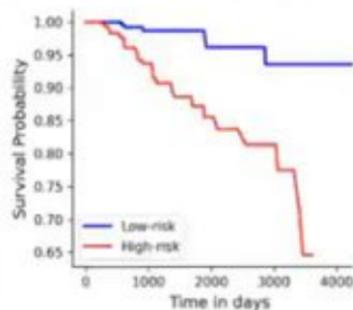
Regulated genes and pathways

MSigDB Term	P value	Adjusted P-value	genes
G2-M Checkpoint	1.1e-4	0.013	PLK1,E2F1,CDC6
Mitotic Spindle	9.4e-4	0.05	PLK1,SAC3D1
E2F Targets	1.5e-3	0.06	PLK1,TK1
Myc Targets V2	3.0e-2	0.102	PLK1
DNA Repair	4.5e-2	0.108	SAC3D1

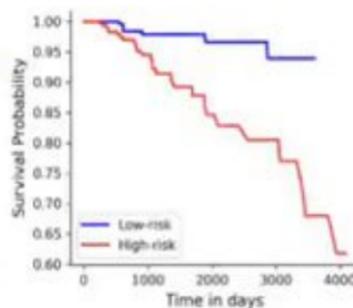
Genetic locus ↔ Expression levels

Survival analysis

(i) KM curve (Gene:MIR205HG)



(ii) KM curve (Gene:C4BPA)

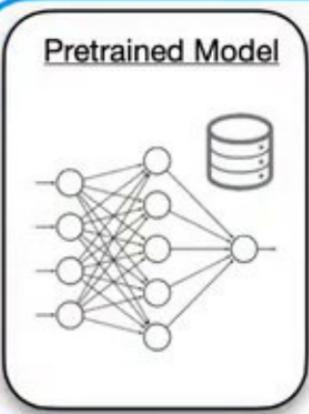


Addressing the Cold-Start Problem for Personalized Combination Drug Screening

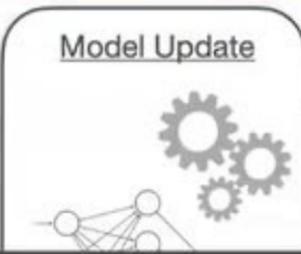
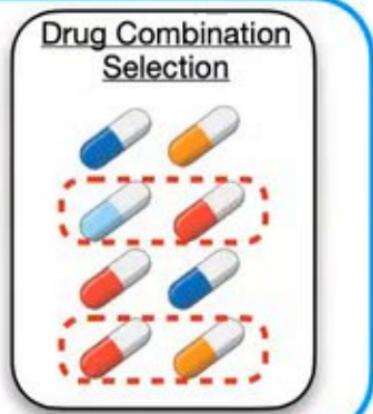


Memorial Sloan Kettering Cancer Center

Antoine de Mathelin, Christopher Tosh and Wesley Tansey



Cold Start Problem



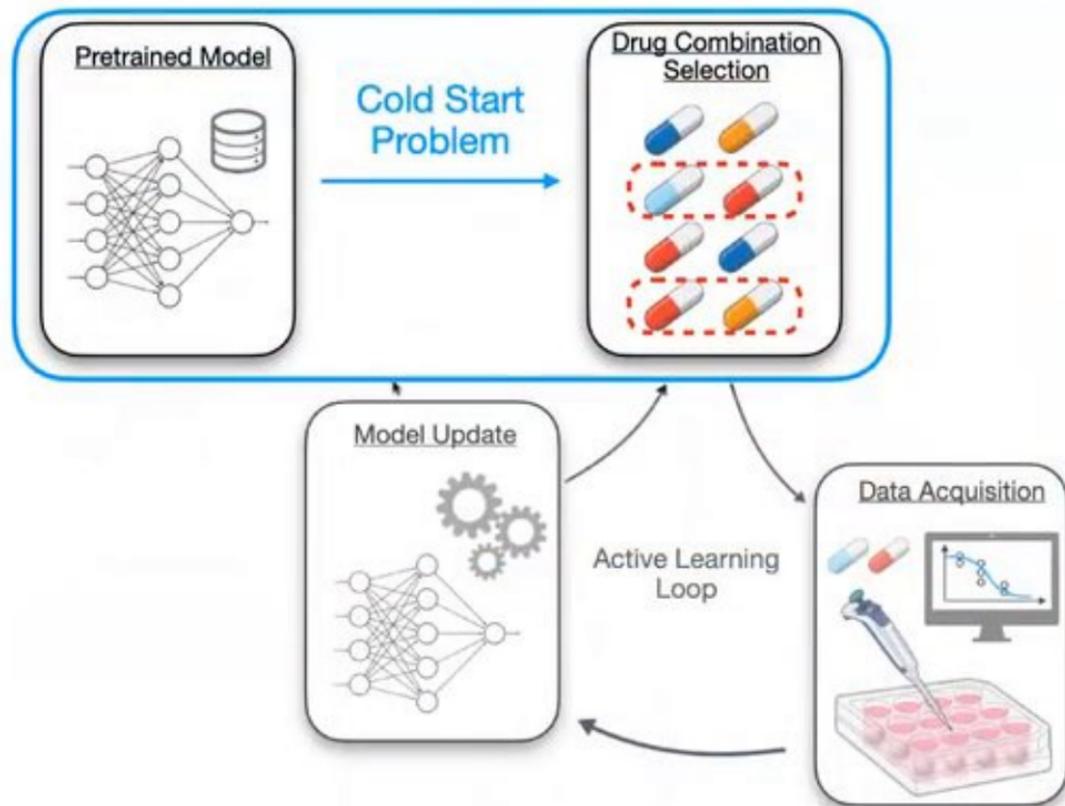
Active Learning Loop



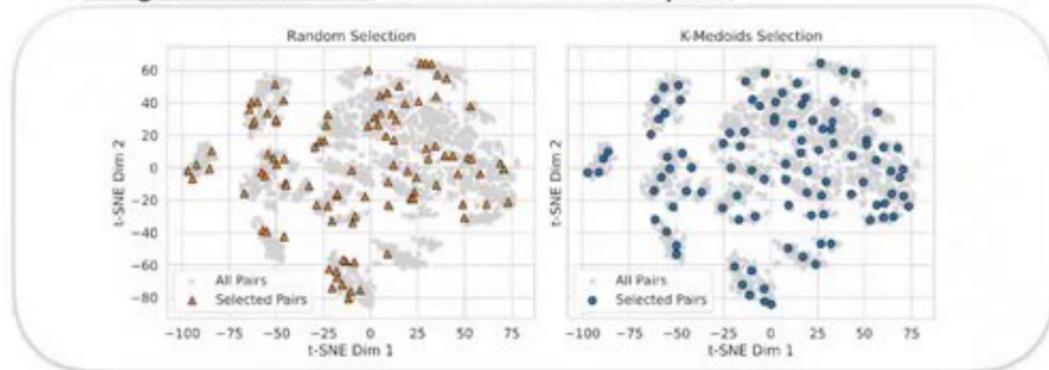
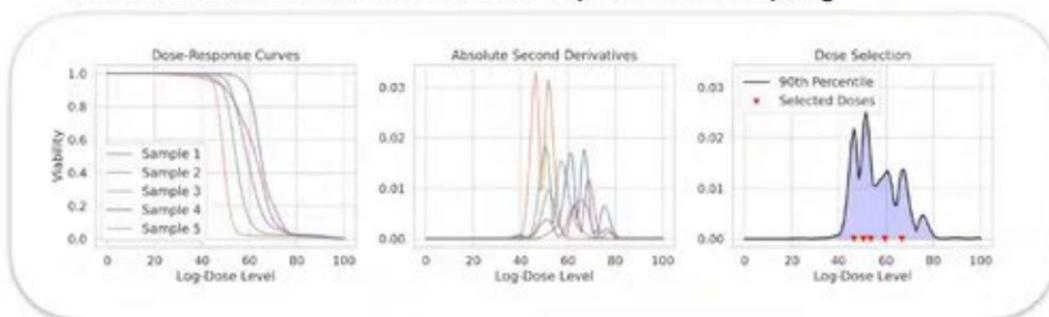
Addressing the Cold-Start Problem for Personalized Combination Drug Screening

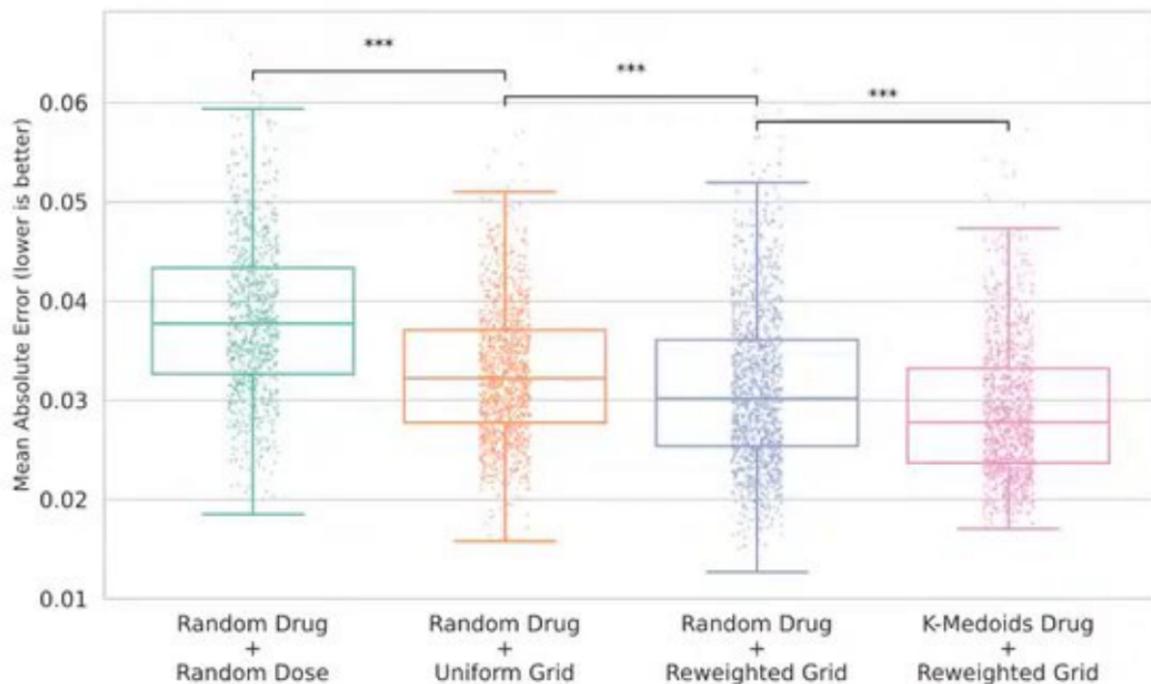


Antoine de Mathelin, Christopher Tosh and Wesley Tansey



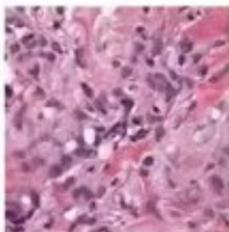
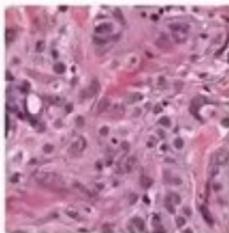
Antoine de Mathelin, Christopher Tosh and Wesley Tansey

Drug Pair Selection: K-medoids in AUC Space**Dose Selection: Curvature-based Importance Sampling**

Antoine de Mathelin, Christopher Tosh and Wesley Tansey

Characterization of genetic effect on tissue Histology

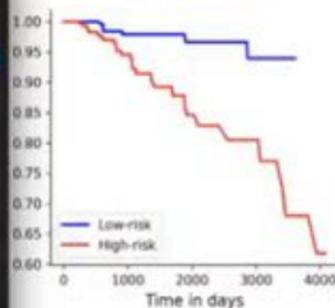
Genetic effect



MLCB2025 spotlight slides

Name	Date Modified	Size	Kind
1_Schreiber Jacob Schreiber.pptx	Sep 8, 2025 at 3:10 PM	1.7 MB	PowerP... (pptx)
2_Seltz Evan Seltz.pdf	Sep 8, 2025 at 8:47 PM	2.5 MB	PDF Document
3_Nair Surag Nair.pdf	Sep 8, 2025 at 2:44 PM	2.2 MB	PDF Document
4_Talk_mahbuba Mahbuba Tasmin.pptx	Sep 8, 2025 at 5:00 PM	4.2 MB	PowerP... (pptx)
5_Ayub Shanza Ayub.pptx	Sep 8, 2025 at 11:03 AM	9.8 MB	PowerP... (pptx)
6_avdiyev Pavel Avdiyev.pdf	Yesterday at 12:23 AM	586 KB	PDF Document
7_Liu Simon Liu.pptx	Sep 8, 2025 at 9:08 PM	9 MB	PowerP... (pptx)
8_Su ZY Su.pptx	Today at 1:04 AM	17.2 MB	PowerP... (pptx)
9_Minsheng&Tianyu Tianyu Liu.pptx	Sep 8, 2025 at 12:09 AM	1.7 MB	PowerP... (pptx)
10_Lam Kevin Lam.pdf	Sep 8, 2025 at 6:19 PM	1.5 MB	PDF Document
11_Choudhary Shubham Choudhary.pptx	Sep 8, 2025 at 4:22 PM	21 MB	PowerP... (pptx)
12_Demathelin Antoine de Mathelin.pdf	Sep 8, 2025 at 3:15 PM	3.7 MB	PDF Document
14_Sergio Mares.key	Yesterday at 12:50 AM	1.7 MB	Keynote
16_Dalal Saykhoom Dalal.key	Sep 8, 2025 at 8:48 PM	5.1 MB	Keynote
MLCB2025_Mares Sergio Mares.pptx	Yesterday at 12:30 AM	356 KB	PowerP... (pptx)
MLCB2025_Mares_Final Sergio Mares.pptx	Yesterday at 12:55 AM	1.8 MB	PowerP... (pptx)

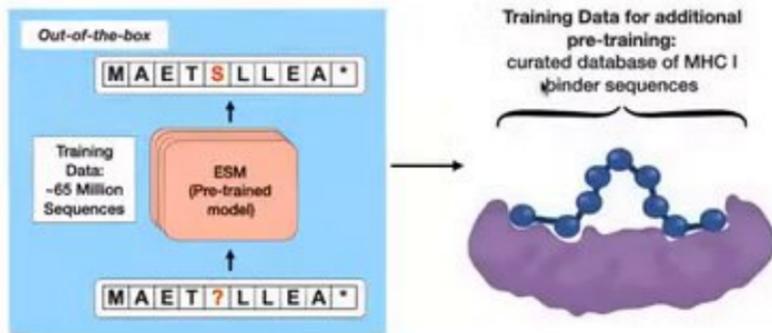
(ii) KM curve (Gene:C4BPA)



"If you're i

Characterization of genetic effect on tissue Histology

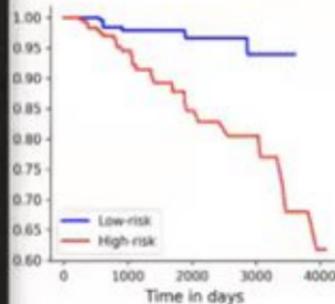
We apply additional domain pre-training to test whether this improves predictions



1. Liu, Z. et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* 376, 1129–1133 (2021).

Berlin

(ii) KM curve (Gene:C4BPA)



"If you're i

**It's infeasible to test all neoantigens
experimentally – we need accurate predictions
for prioritization**

Continued domain-specific pre-training of protein language models for pMHC-I binding prediction



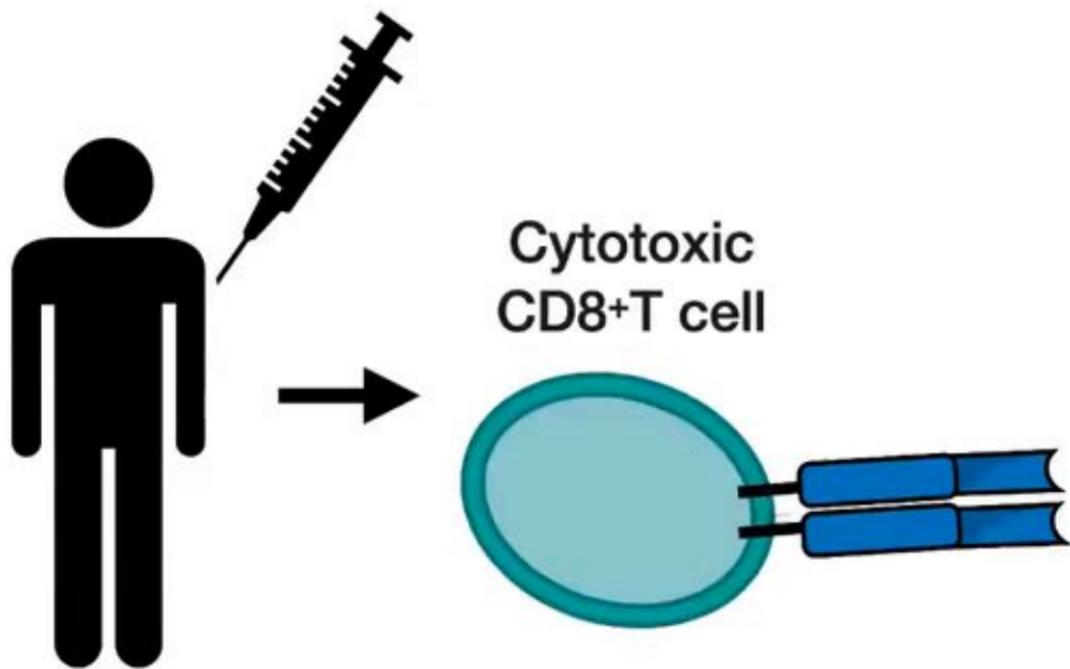
Sergio E. Mares

PhD Student

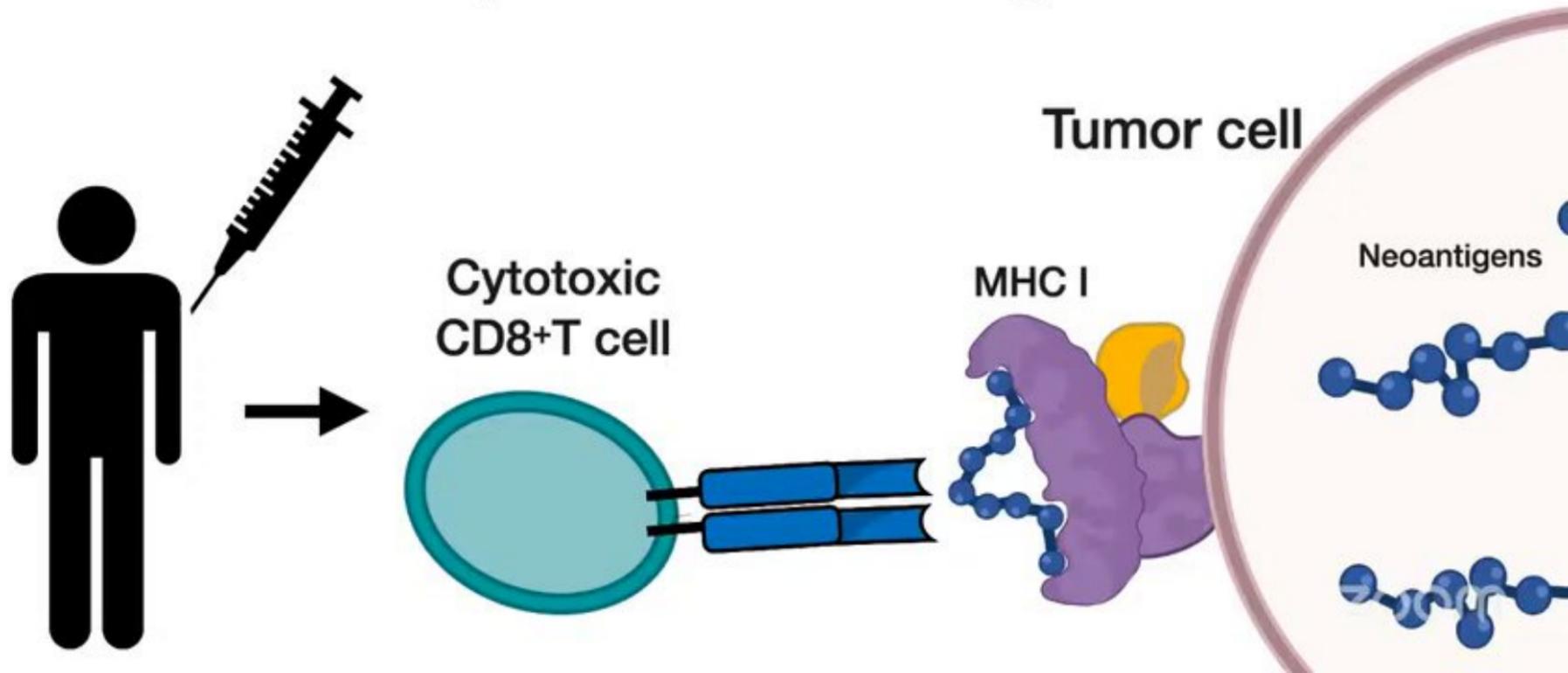
Center for Computational Biology

Vaccines stimulate the immune system to respond to neoantigens

Vaccines stimulate the immune system to respond to neoantigens

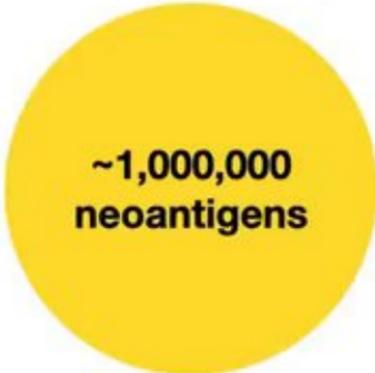


Vaccines stimulate the immune system to respond to neoantigens



It's infeasible to test all neoantigens experimentally – we need accurate predictions for prioritization

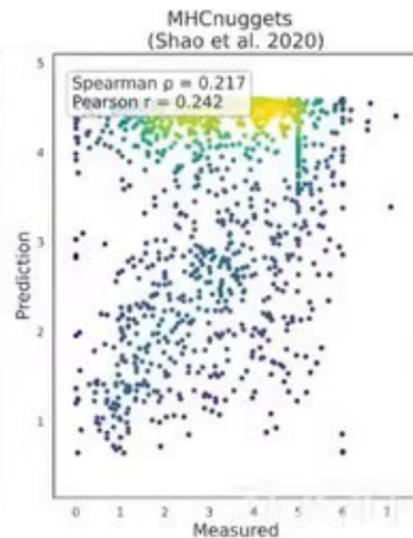
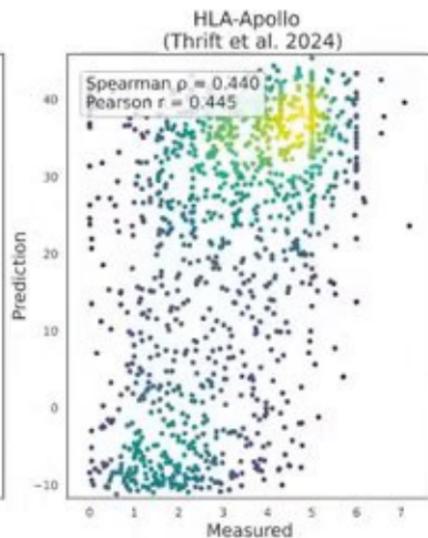
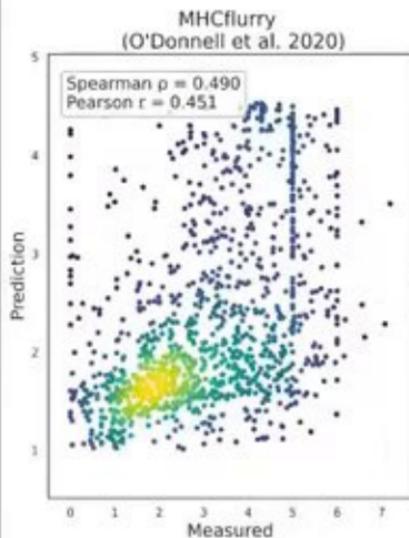
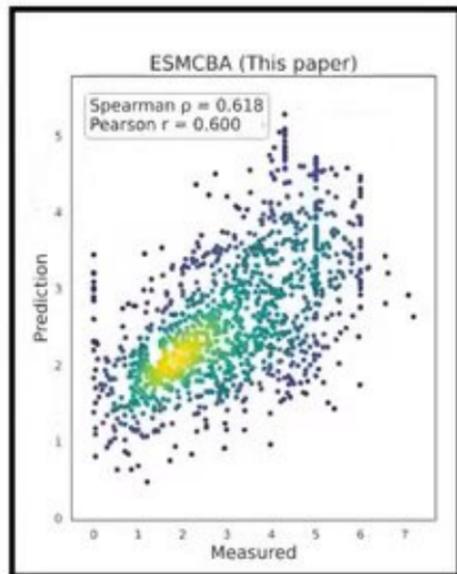
The Cancer
Genome Atlas



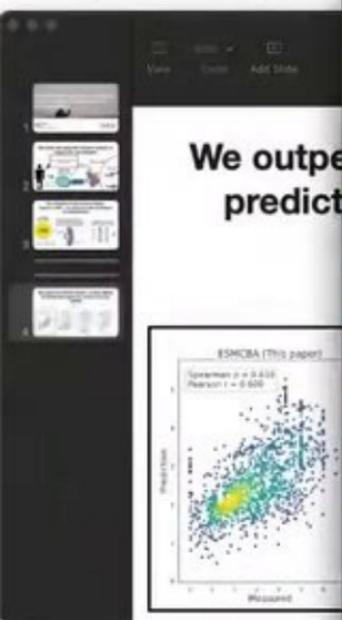
~1,000,000
neoantigens

We outperform SOTA models in binding affinity prediction leveraging pre-trained language models

Poster #14



Characterization of genetic effect on tissue Histology

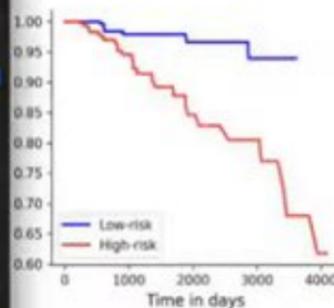


"If you're i

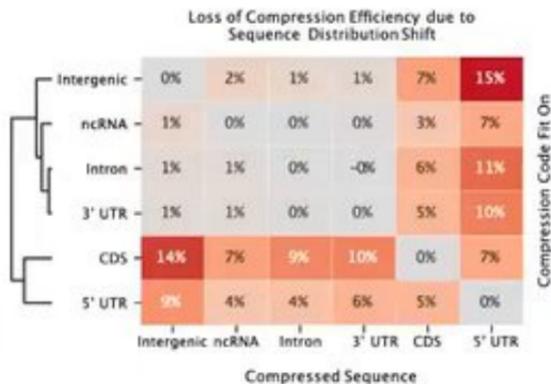
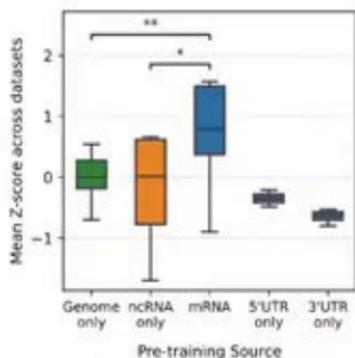
MLCB2025 spotlight slides

Name	Date Modified	Size	Kind
1_Schreiber Jacob Schreiber.pptx	Sep 8, 2025 at 5:10 PM	1.7 MB	PowerP...n (.pptx)
2_Seltz Evan Seltz.pdf	Sep 8, 2025 at 8:47 PM	2.5 MB	PDF Document
3_Nair Surag Nair.pdf	Sep 8, 2025 at 2:44 PM	2.2 MB	PDF Document
4_Talk_mahbuba Mahbuba Tasmin.pptx	Sep 9, 2025 at 5:00 PM	4.2 MB	PowerP...n (.pptx)
5_Ayub Shanza Ayub.pptx	Sep 9, 2025 at 11:03 AM	9.9 MB	PowerP...n (.pptx)
6_aidreyev Pavel Ardeev.pdf	Yesterday at 12:23 AM	586 KB	PDF Document
7_Liu Simon Liu.pptx	Sep 8, 2025 at 9:08 PM	8 MB	PowerP...n (.pptx)
8_Su ZY Su.pptx	Today at 1:04 AM	17.2 MB	PowerP...n (.pptx)
9_Minsheng&Tianyu Tianyu Liu.pptx	Sep 8, 2025 at 12:09 AM	1.7 MB	PowerP...n (.pptx)
10_Liam Kevin Lam.pdf	Sep 9, 2025 at 6:39 PM	1.5 MB	PDF Document
11_Chaudhary Shubham Chaudhary.pptx	Sep 8, 2025 at 4:22 PM	21 MB	PowerP...n (.pptx)
12_Demathelin Antoine de Mathelin.pdf	Sep 8, 2025 at 2:16 PM	1.7 MB	PDF Document
M_Sergio Mares.key	Yesterday at 12:50 AM	1.7 MB	Keynote
15_Datal Saythoom Dalat.key	Sep 8, 2025 at 9:48 PM	3.1 MB	Keynote
MLCB2025_Mares Sergio Mares.pptx	Yesterday at 12:38 AM	356 KB	PowerP...n (.pptx)
MLCB2025_Mares_Final Sergio Mares.pptx	Yesterday at 12:55 AM	1.8 MB	PowerP...n (.pptx)

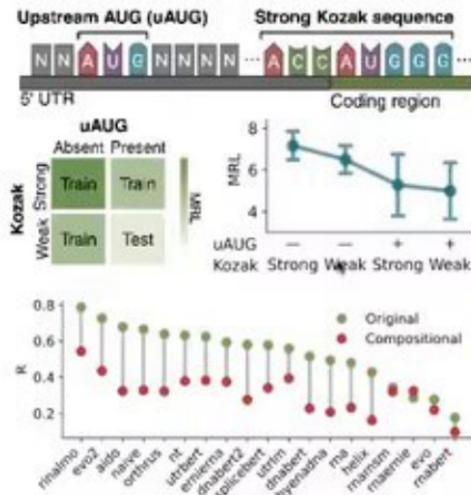
(ii) KM curve (Gene:C4BPA)



What else did we do?



→ Quantified why DNA / ncRNA models perform poorly at mRNA tasks



→ Assessed if nucleotide FMs can generalize to novel combinations of familiar sequence features

mRNABench: A curated benchmark for mature mRNA property and function prediction

Ruian Shi*, **Taykhoom Dalal***, **Philip Fradkin***, Divya Koyyalagunta, Simran Chhabria,
Andrew Jung, Cyrus Tam, Defne Ceyhan, Jessica Lin, Kaitlin U. Lavery, Ilyes Baali, Bo
Wang, Quaid Morris

Machine Learning in Computational Biology 2025



Memorial Sloan Kettering
Cancer Center

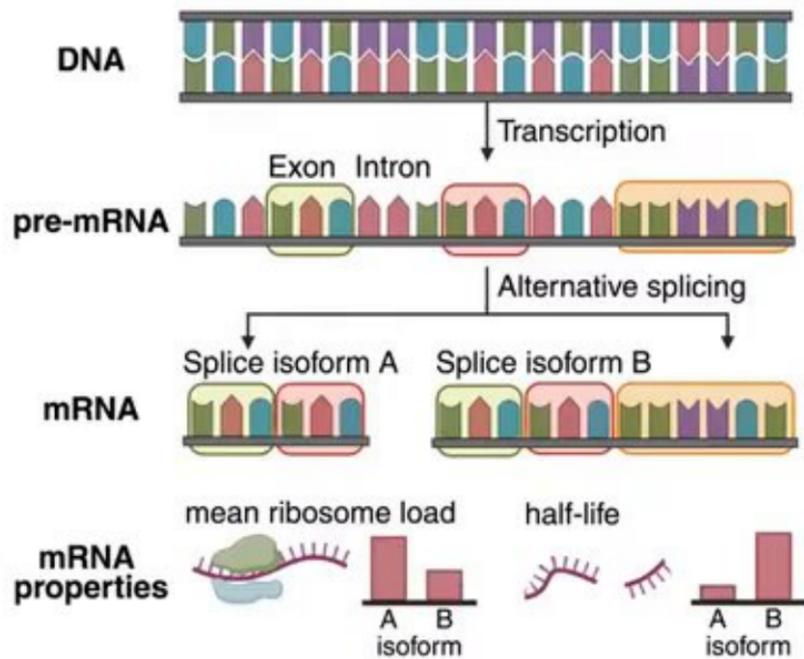


UNIVERSITY OF
TORONTO



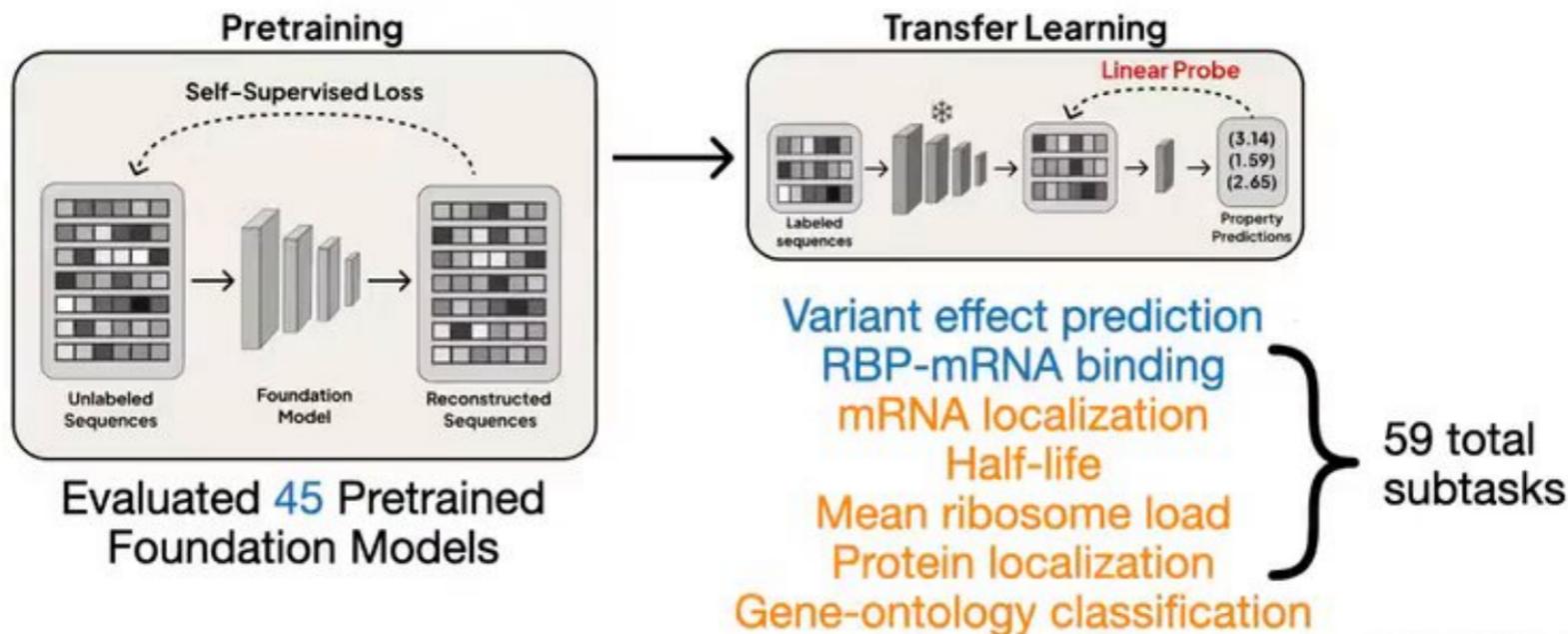
VECTOR
INSTITUTE

Motivation

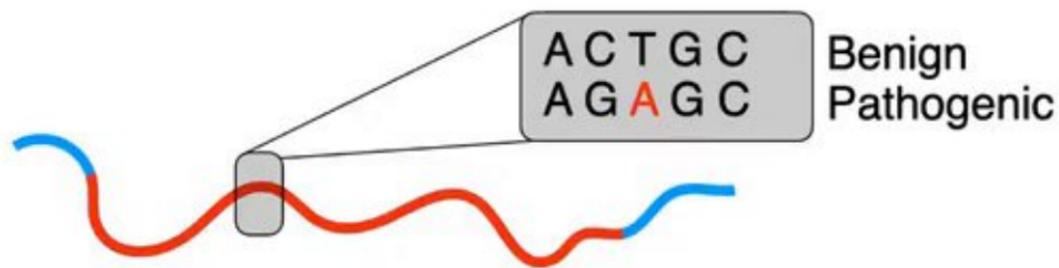


- Alternative splicing can generate multiple splice isoforms with distinct properties and functions
- Self-supervised nucleotide FMs can be used to predict these properties but there are no existing mRNA-focused benchmarks

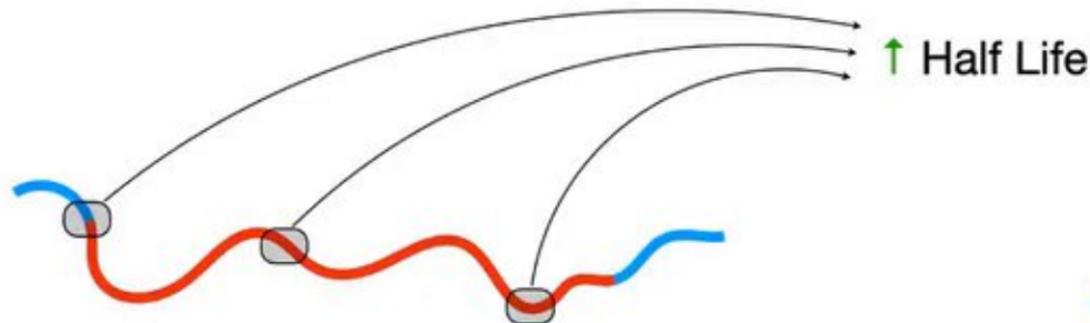
mRNABench is the first mRNA-specific nucleotide FM benchmarking suite



We categorized tasks in our benchmark as either being “**local**” or “**global**”



Local
Variant effect prediction
RBP-mRNA binding
Mean Ribosome Load (MPRA)

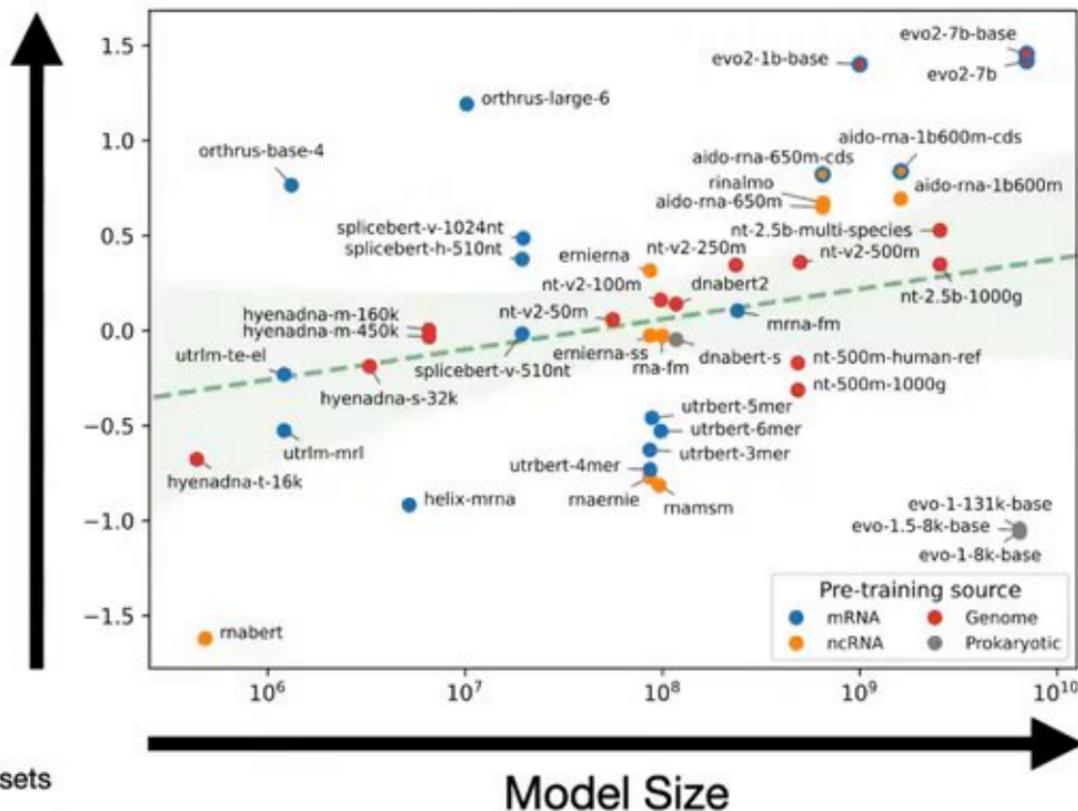


Global
mRNA localization
Half-life
Mean ribosome load
Protein localization
Gene-ontology classification

* UTR CDS

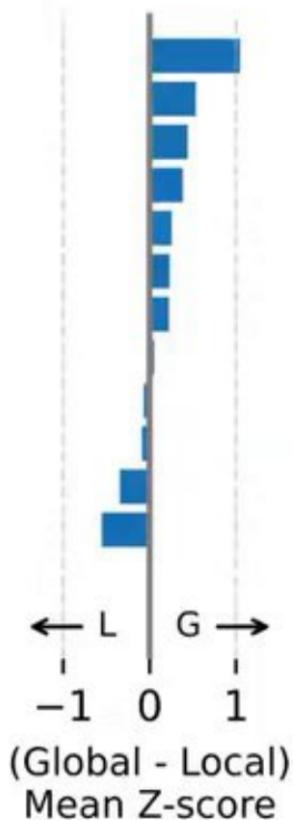
The top performing models are **Evo2** and **Orthrus**

*Average Performance Across Tasks



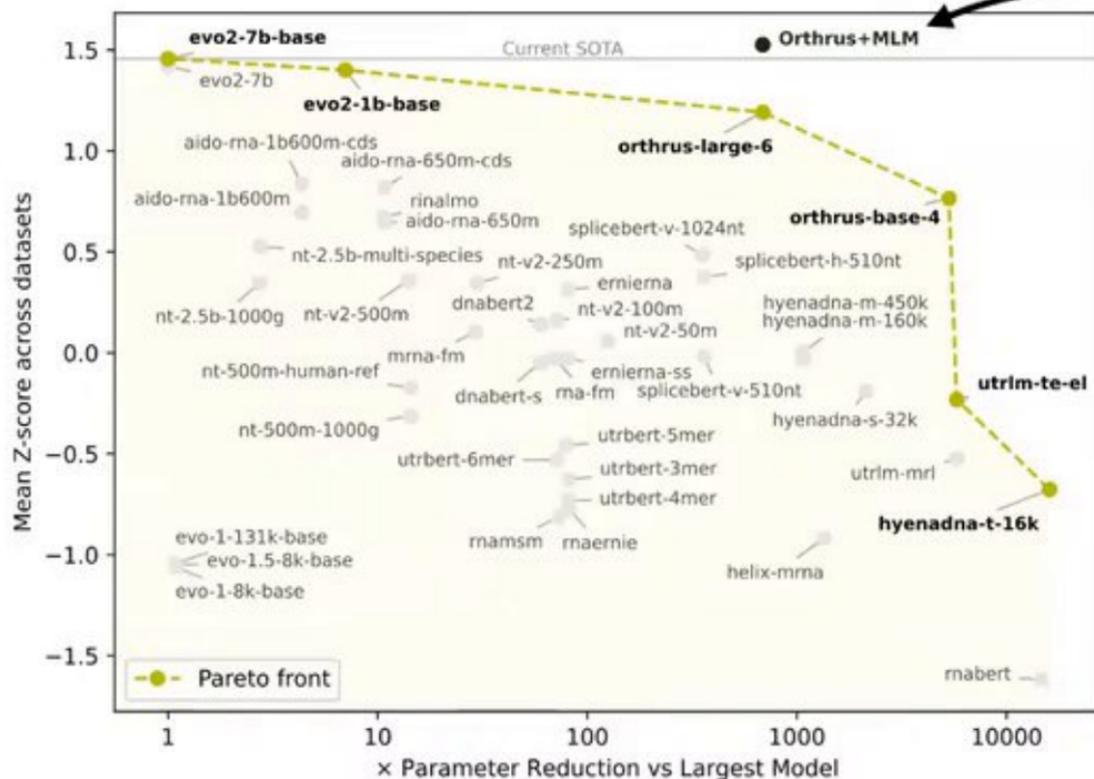
*Mean Z-score across datasets

Orthrus
 Evo2
 SpliceBERT
 RNA-FM
 AIDO.RNA
 RNA-MSM
 RNAErnie
 RNABERT
 UTR-LM
 ERNIE-RNA
 RiNALMo
 3UTRBERT



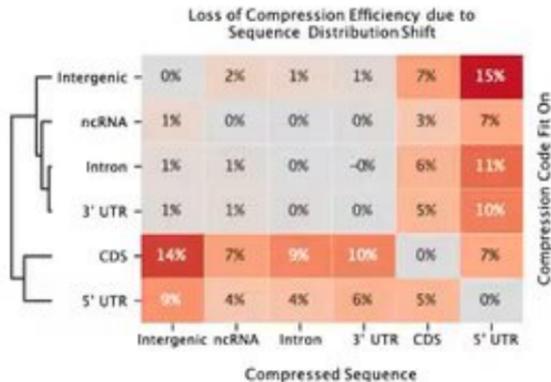
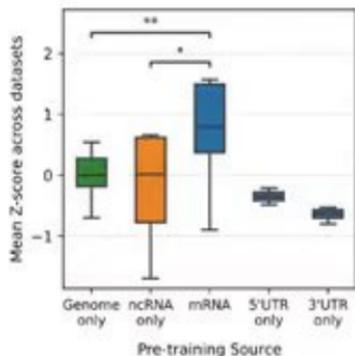
Orthrus does much better on “**global**” tasks than “**local**” ones due to its **contrastive learning objective**

What did we do with this information?

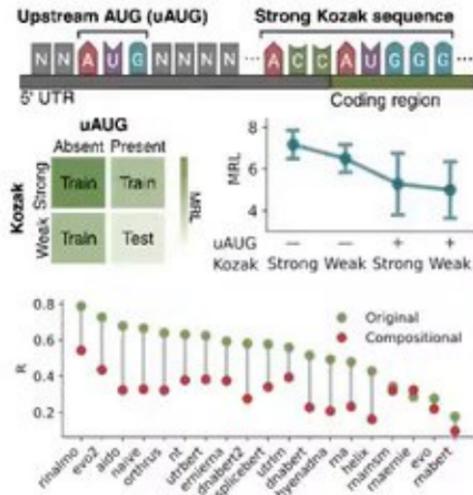


Trained a SOTA Orthrus variant using a joint **contrastive learning** and **masked language modeling** objective

What else did we do?



→ Quantified why DNA / ncRNA models perform poorly at mRNA tasks



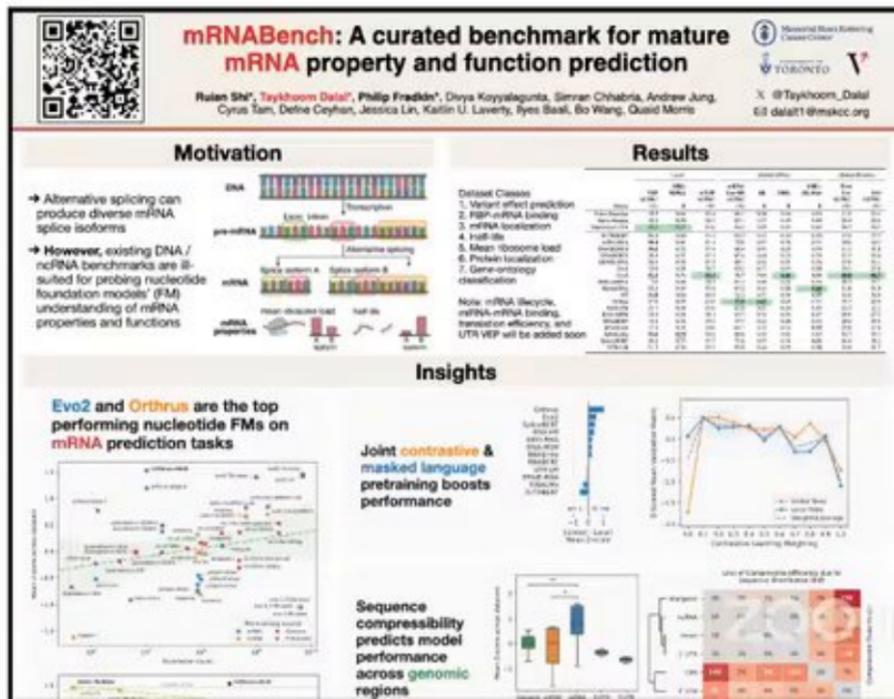
→ Assessed if nucleotide FMs can generalize to novel combinations of familiar sequence features

If any of this sounds interesting...

Preprint



Poster 15



NYGC Events

MLCB @ NYGC



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley

Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 1.30pm ET

mlcb.org for schedule

NYGC Events



"zoom.us" is requesting to bypass the system private window picker and directly access your screen and audio.

This will allow zoom.us to record your screen and system audio, including personal or sensitive information that may be visible or audible.

[Allow](#)

[Open System Settings](#)

[Deny](#)



Macintosh HD



Screenshot
2025-0...6.88.png



Click **Open zoom.us** on the dialog shown by your browser
If you don't see a dialog, click **Join from Zoom Workplace app** below

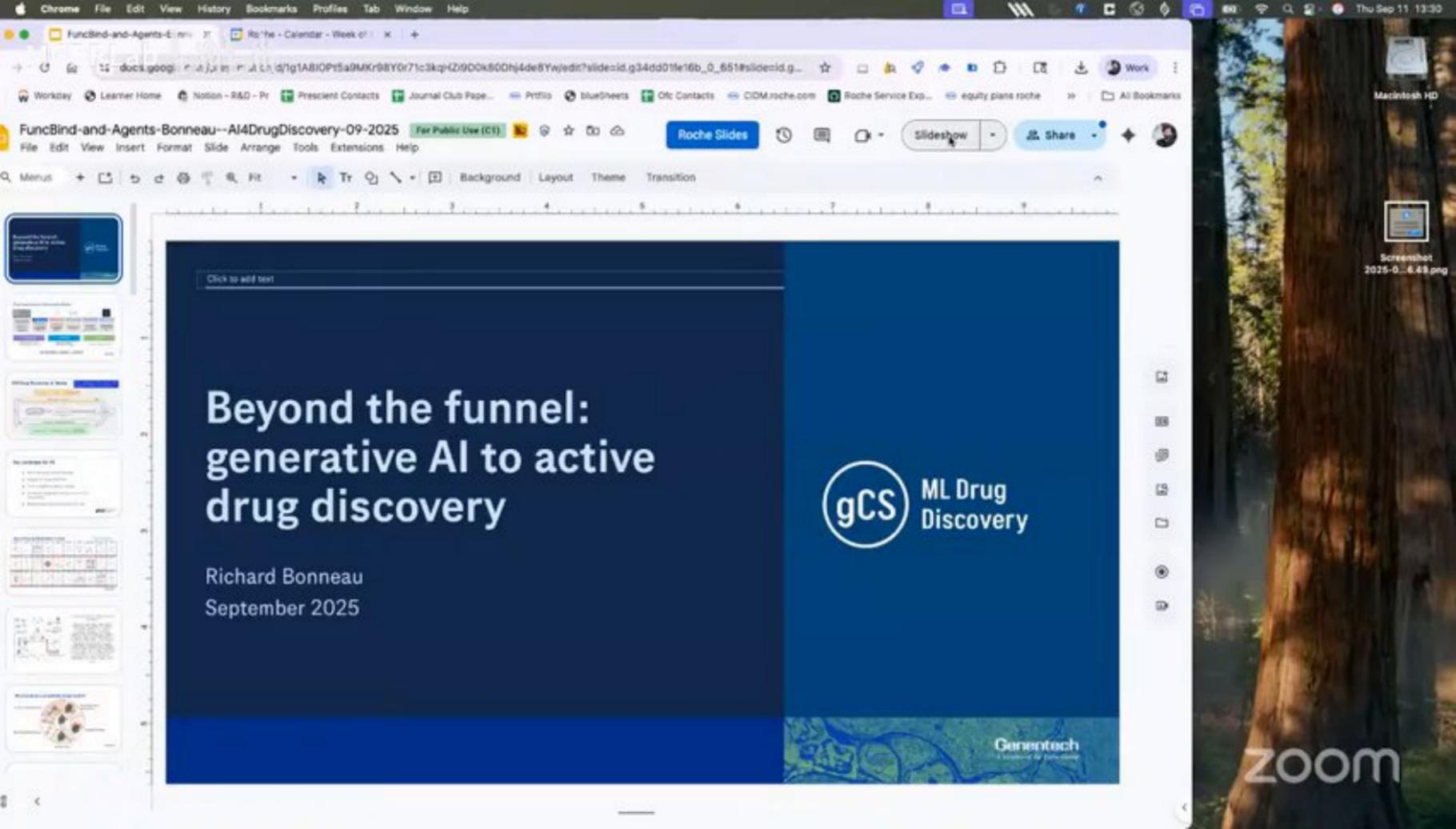
By joining a meeting, you agree to our [Terms of Service](#) and [Privacy Statement](#)

Join from Zoom Workplace app

Did not open Zoom Workplace app? x
Please [download](#) and install the app and click Join from Zoom Workplace app again.

Don't have the Zoom Workplace app installed? [Download Now](#)





Click to add text

Beyond the funnel: generative AI to active drug discovery

Richard Bonneau
September 2025



Screenshot
2025-0...-6.48.png

Macintosh HD

Beyond the funnel: generative AI to active drug discovery

Richard Bonneau
September 2025



MCBRLab bilibili

Stack of desktop widgets including a calendar, a notes board, and a Zoom meeting control panel.

Screen Recording

 "zoom.us.app" would like to record this computer's screen and audio.

Grant access to this application in Privacy & Security settings, located in System Settings.

[Open System Settings](#) [Deny](#)

Macintosh HD

Screenshot 2025-0...-6.48.png

zoom

Click to add text

Beyond the funnel: generative AI to active drug discovery

Richard Bonneau
September 2025

gCS ML Drug Discovery

Genentech
A member of the Roche Group



Loading...

Beyond the funnel: generative AI to active drug discovery

Richard Bonneau
September 2025



ML Drug
Discovery

AI for Drug Discovery @ Genentech & Roche

Richard
Bonneau
ML Drug
Discovery



Computational drug discovery, driving new approaches and synergies between small and large molecule drug discovery

Small Mol & Peptide Drug Discovery

computational and machine learning tools for **small molecule and peptide design** and optimization

Large Molecule Drug Discovery

computational and machine learning tools for **large molecule design** and optimization

Frontier Research

foundational machine learning methods and approaches

ML Engineering

machine learning engineering for enabling efforts across all MLDD pillars

Prescient LM

internal development of **large language models** across gRED & Roche

Technology Dev

new capabilities for unmet needs and exploratory work

Platforms

engineering, integration, extensibility, and scaling

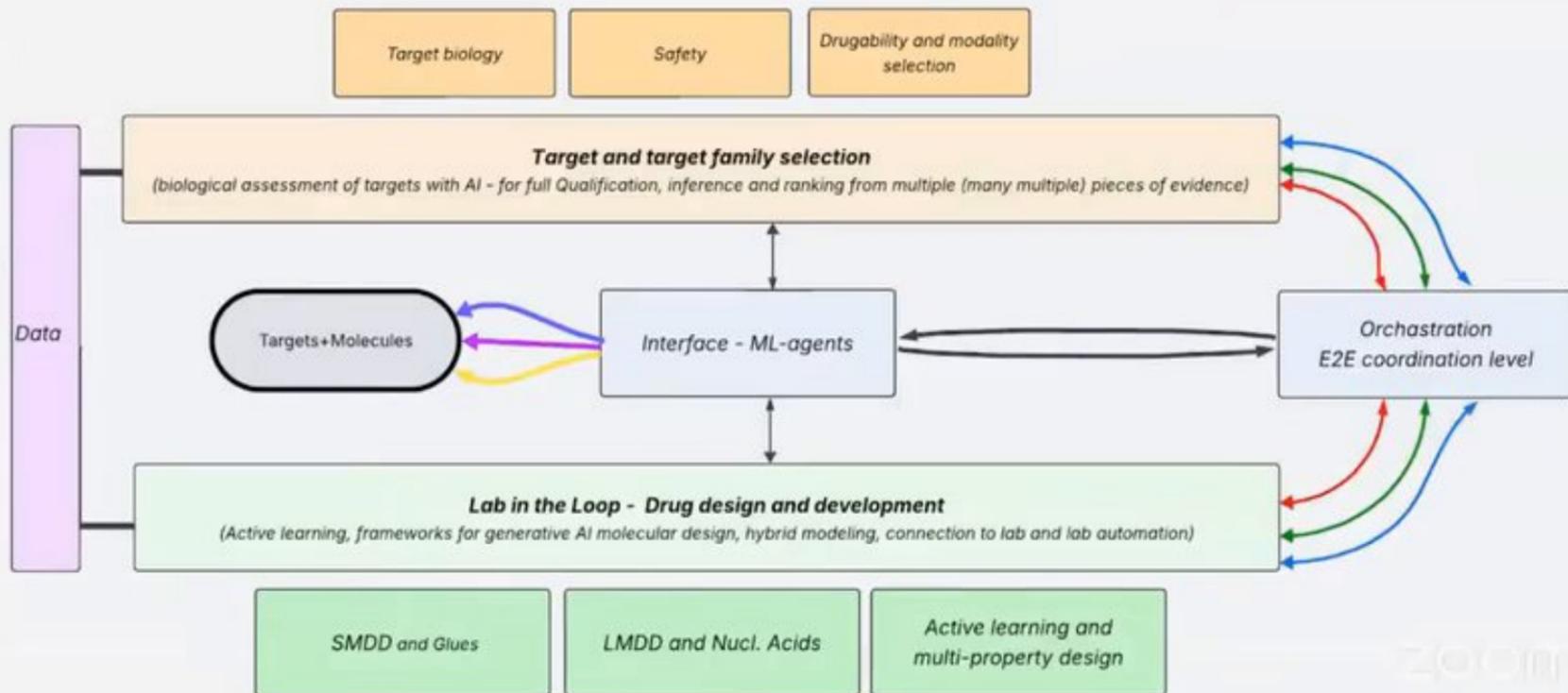
Portfolio

application and program impact

new capabilities, platforms, & portfolio

E2E Drug Discovery @ Roche

w/ Sara Mostafavi, Fabian Birzele, Anna Craig, Ryan Copping, John Marioni

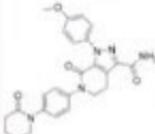
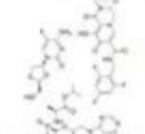
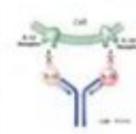
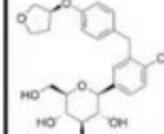
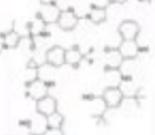
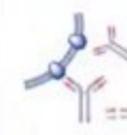
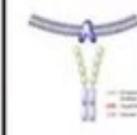
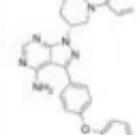
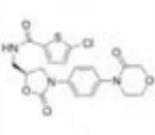
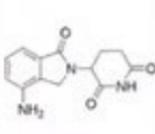
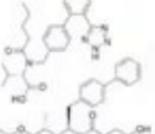
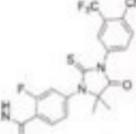
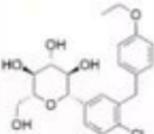
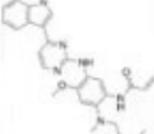


Key challenges for DD

- Lab in the Loop / Active learning
- Integration across timelines
- Cross modality workflows needed
- Contrast in depth/size/runtime of workflow components
- Motivating process optimization with ML

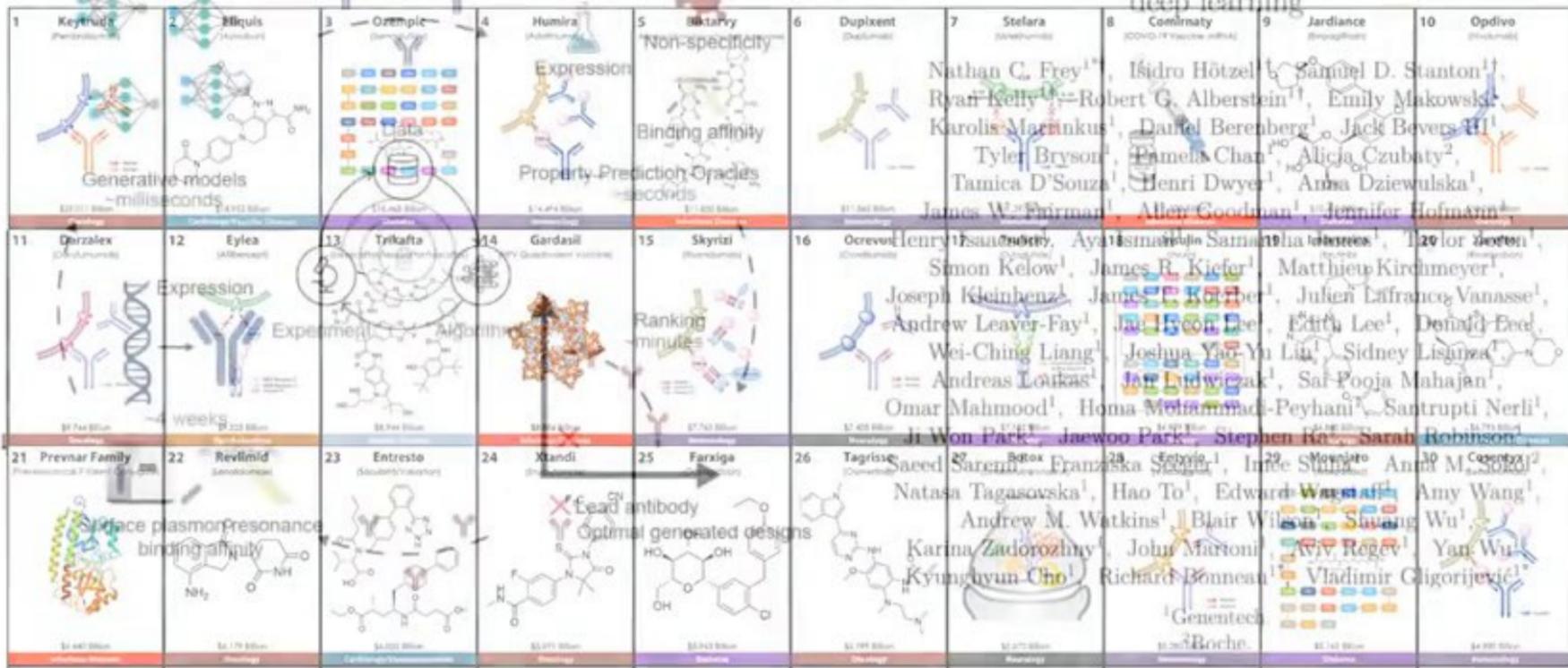
Top 30 Drug by Retail Sales in 2023

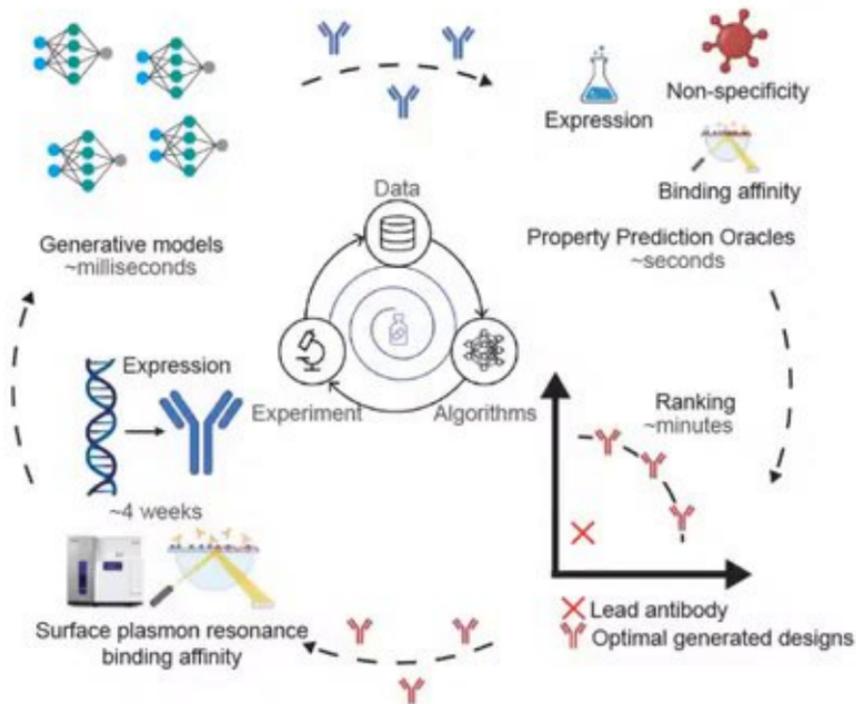
[Njardarson group](#)

<p>1 Keytruda (Pembrolizumab)</p>  <p>\$23.071 Billion <i>Roche/Genentech</i></p>	<p>2 Eliquis (Apixiban)</p>  <p>\$18.923 Billion <i>Cardinal/Novartis</i></p>	<p>3 Ozempic (Semaglutide)</p>  <p>\$16.468 Billion <i>Cardinal/Novartis</i></p>	<p>4 Humira (Adalimumab)</p>  <p>\$14.456 Billion <i>Roche/Genentech</i></p>	<p>5 Biktarvy (Bictegravir/Emsavirenat/Vemurafenon)</p>  <p>\$11.900 Billion <i>Roche/Genentech</i></p>	<p>6 Dupixent (Dupilumab)</p>  <p>\$11.565 Billion <i>Roche/Genentech</i></p>	<p>7 Stelara (Stelezumab)</p>  <p>\$11.287 Billion <i>Roche/Genentech</i></p>	<p>8 Comirnaty (COVID-19 Vaccine, mRNA)</p>  <p>\$11.220 Billion <i>Roche/Genentech</i></p>	<p>9 Jardiance (Empagliflozin)</p>  <p>\$10.407 Billion <i>Roche/Genentech</i></p>	<p>10 Opdivo (Nivolumab)</p>  <p>\$10.220 Billion <i>Roche/Genentech</i></p>
<p>11 Darzalex (Daratumumab)</p>  <p>\$9.766 Billion <i>Roche/Genentech</i></p>	<p>12 Eylea (Aflibercept)</p>  <p>\$9.222 Billion <i>Roche/Genentech</i></p>	<p>13 Trikafta (Elexcaftor/Tenaxcaftor/Proxacaftor)</p>  <p>\$8.766 Billion <i>Roche/Genentech</i></p>	<p>14 Gardasil (HPV Quadrivalent Vaccine)</p>  <p>\$8.566 Billion <i>Roche/Genentech</i></p>	<p>15 Skyrizi (Risankumab)</p>  <p>\$7.763 Billion <i>Roche/Genentech</i></p>	<p>16 Ocrevus (Ocrelizumab)</p>  <p>\$7.428 Billion <i>Roche/Genentech</i></p>	<p>17 Trulicity (Dulaglutide)</p>  <p>\$7.122 Billion <i>Roche/Genentech</i></p>	<p>18 Insulin (Insulin)</p>  <p>\$6.879 Billion <i>Roche/Genentech</i></p>	<p>19 Imbruvica (Ibrutinib)</p>  <p>\$6.640 Billion <i>Roche/Genentech</i></p>	<p>20 Xarelto (Edoxaban)</p>  <p>\$6.793 Billion <i>Cardinal/Novartis</i></p>
<p>21 Prevnar Family (Pneumococcal Polysaccharide Conjugate)</p>  <p>\$6.460 Billion <i>Roche/Genentech</i></p>	<p>22 Revlimid (Lenalidomide)</p>  <p>\$6.179 Billion <i>Roche/Genentech</i></p>	<p>23 Entresto (Sacubitril/Valsartan)</p>  <p>\$6.020 Billion <i>Cardinal/Novartis</i></p>	<p>24 Xtandi (Enzalutamide)</p>  <p>\$5.971 Billion <i>Roche/Genentech</i></p>	<p>25 Farxiga (Dapagliflozin)</p>  <p>\$5.640 Billion <i>Roche/Genentech</i></p>	<p>26 Tagrisso (Osimertinib)</p>  <p>\$5.779 Billion <i>Roche/Genentech</i></p>	<p>27 Botox (OnabotulinumtoxinA)</p>  <p>\$5.470 Billion <i>Roche/Genentech</i></p>	<p>28 Entyvio (Vedolizumab)</p>  <p>\$5.380 Billion <i>Roche/Genentech</i></p>	<p>29 Mounjaro (Tirzepatide)</p>  <p>\$5.143 Billion <i>Roche/Genentech</i></p>	<p>30 Cosentyx (Secukinumab)</p>  <p>\$4.980 Billion <i>Roche/Genentech</i></p>

Top 30 Drug by Retail Sales in 2023

Lab-in-the-loop therapeutic Nanobody design with deep learning





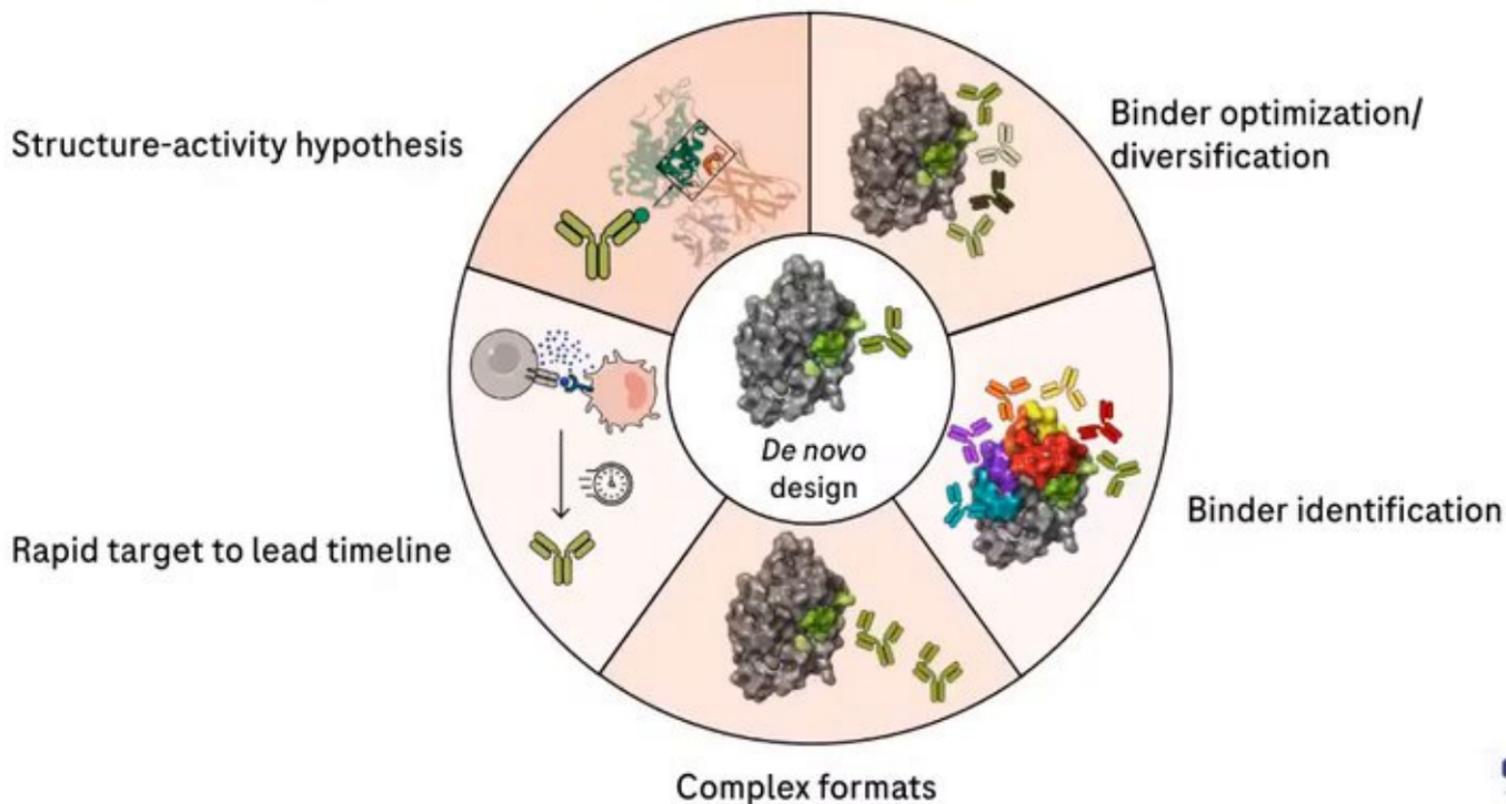
Lab-in-the-loop therapeutic antibody design with deep learning

Nathan C. Frey^{1†}, Isidro Hötzel^{1†}, Samuel D. Stanton^{1†},
 Ryan Kelly^{1†}, Robert G. Alberstein^{1†}, Emily Makowski¹,
 Karolis Martinkus¹, Daniel Berenberg¹, Jack Bevers III¹,
 Tyler Bryson¹, Pamela Chan¹, Alicja Czuby²,
 Tamica D'Souza¹, Henri Dwyer¹, Anna Dziewulska¹,
 James W. Fairman¹, Allen Goodman¹, Jennifer Hofmann¹,
 Henry Isaacson¹, Aya Ismail¹, Samantha James¹, Taylor Joren¹,
 Simon Kelow¹, James R. Kiefer¹, Matthieu Kirchmeyer¹,
 Joseph Kleinhenz¹, James T. Koerber¹, Julien Lafrance-Vanasse¹,
 Andrew Leaver-Fay¹, Jae Hyeon Lee¹, Edith Lee¹, Donald Lee¹,
 Wei-Ching Liang¹, Joshua Yao-Yu Lin¹, Sidney Lisanza¹,
 Andreas Loukas¹, Jan Ludwiczak¹, Sai Pooja Mahajan¹,
 Omar Mahmood¹, Homa Mohammadi-Peyhani¹, Santrupti Nerli¹,
 Ji Won Park¹, Jaewoo Park¹, Stephen Ra¹, Sarah Robinson¹,
 Saeed Saremi¹, Franziska Seeger¹, Imee Sinha¹, Anna M. Sokol²,
 Natasa Tagasovska¹, Hao To¹, Edward Wagstaff¹, Amy Wang¹,
 Andrew M. Watkins¹, Blair Wilson¹, Shuang Wu¹,
 Karina Zadorozhny¹, John Marioni¹, Aviv Regev¹, Yan Wu¹,
 Kyunghyun Cho¹, Richard Bonneau^{1*}, Vladimir Gligorijević^{1*}

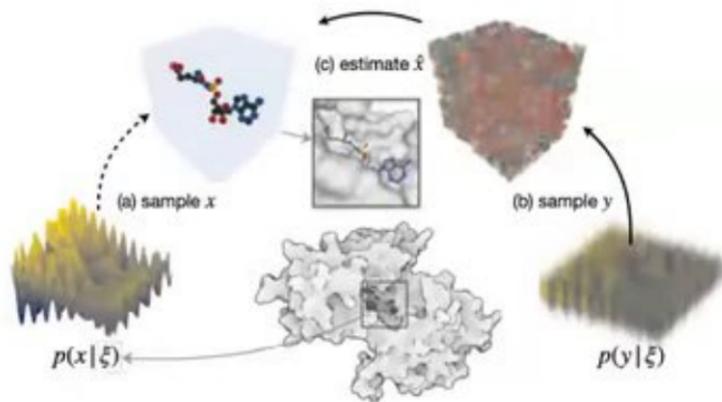
¹Genentech.

²Roche.

What would *de novo* antibody design enable?

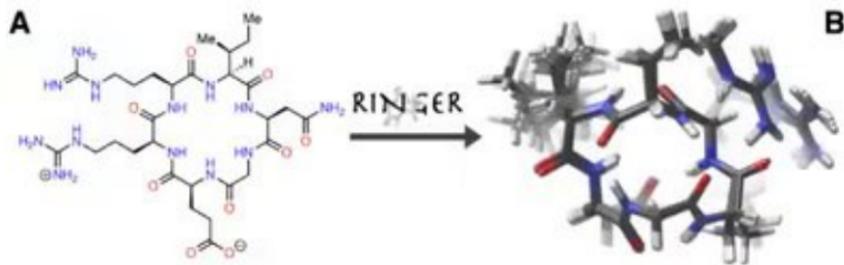


Molecular Generation and Conformational Sampling



Oliveira Pinheiro et al. "Structure-based drug design by denoising voxel grids." *ICML 2024*

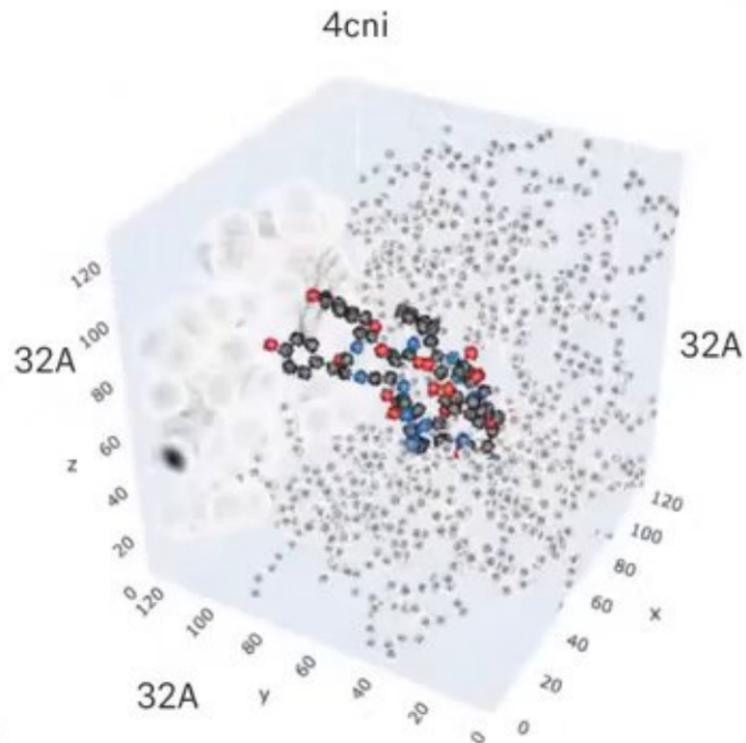
pocket-conditioned 3D molecule generation
with **walk jump sampling**



Grambow et al. "Accurate and Efficient Structural Ensemble Generation of Macrocyclic Peptides using Internal Coordinate Diffusion". arXiv:2305.19800v2 2024

efficient conformer generation of macrocyclic peptides and proteins with **walk-jump sampling** and **diffusion** models

Representing molecules as atomic density fields



De novo & X-modality : FuncBind

Unconditional generation

Target conditioned generation

Voxels

VoxMol [Pinheiro et al., NeurIPS 2023](#)

Modality: small molecule

VoxBind [Pinheiro et al., ICML 2024](#)

Modality: small molecule

Neural Fields

FuncMol [Kirchmeyer et al., NeurIPS 2024](#)

Modality: small molecule, macrocyclic peptide

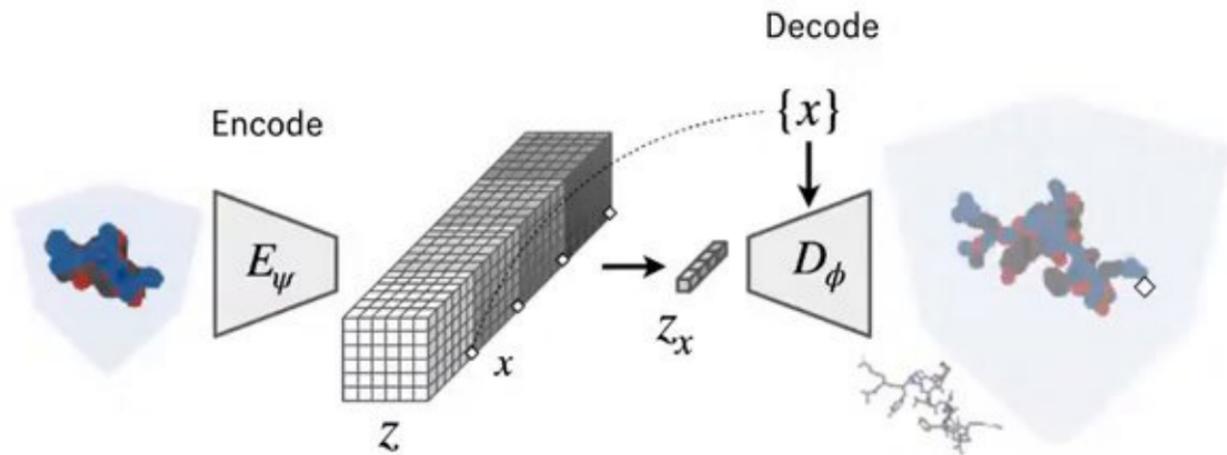
FuncBind (under review)

Modality: unified model for small molecules, macrocyclic peptides, antibody CDR loops

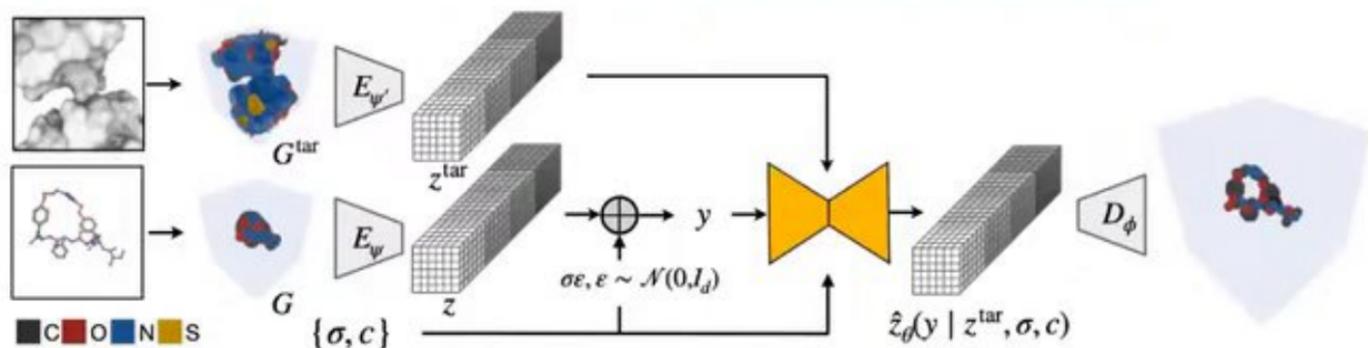
Improved scalability and speed going from voxels to neural fields

Several technical improvements from FuncMol to FuncBind

Embedding molecules into a latent space using neural fields

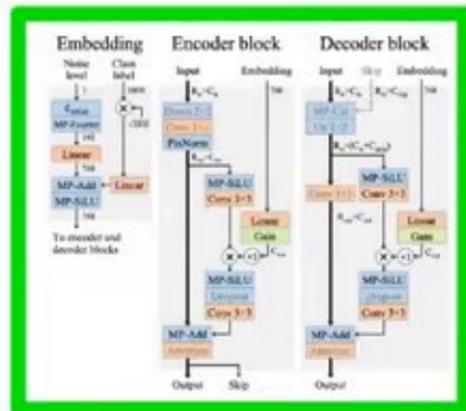


FuncBind's workhorse: a powerful denoiser architecture



Large model leveraging all relevant interface data

Scale via neural field representation



Sampling via score-based approaches

e.g. Walk-Jump Sampling [[Saremi & Hyvärinen, 2019 JMLR](#)]



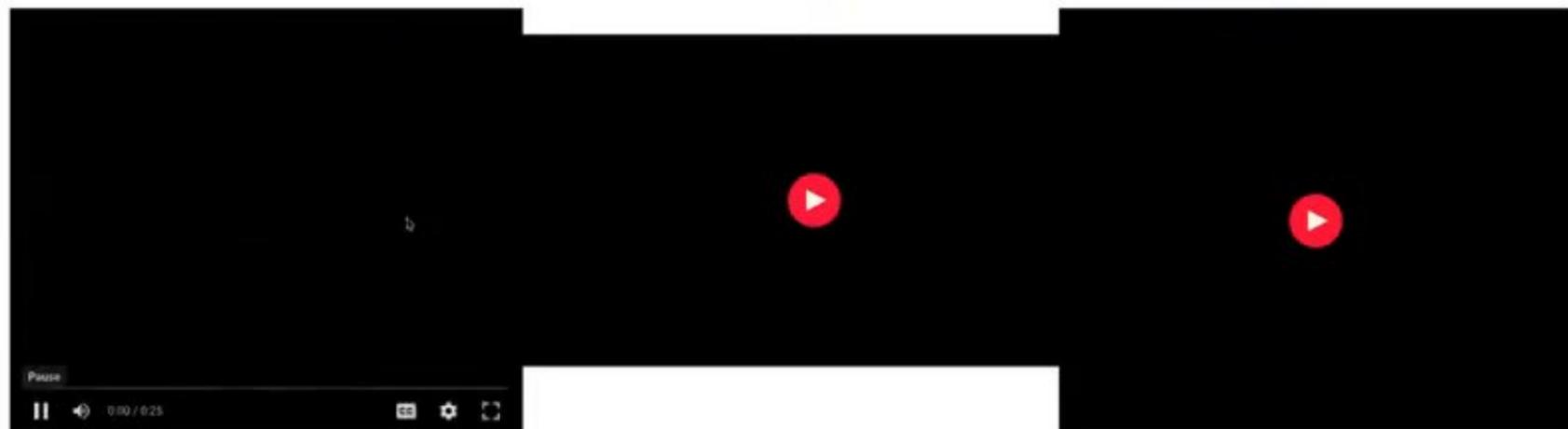
Sample

Pocket conditioned de-novo x-modal generation

small molecule

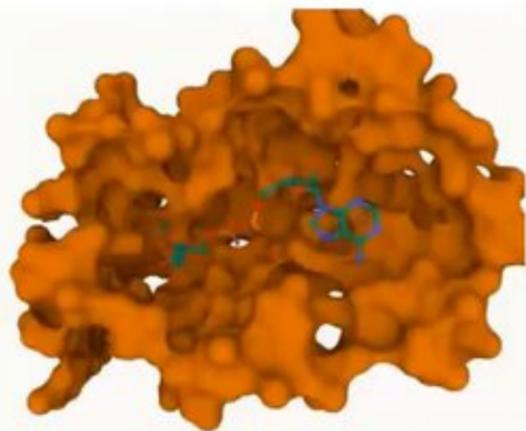
antibody

peptide

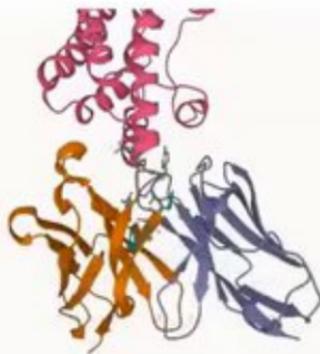


Pocket conditioned de-novo x-modal generation

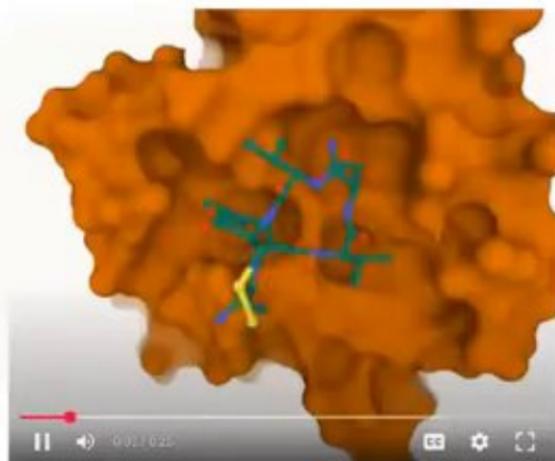
small molecule



antibody



peptide



FuncBind generates new non canonical amino acids

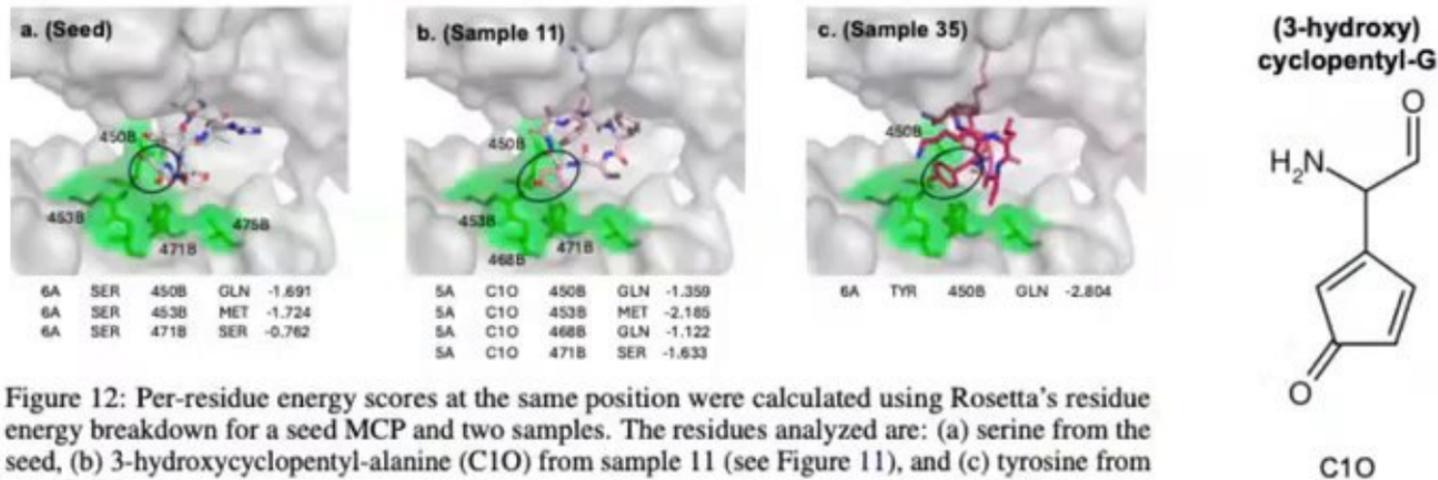


Figure 12: Per-residue energy scores at the same position were calculated using Rosetta's residue energy breakdown for a seed MCP and two samples. The residues analyzed are: (a) serine from the seed, (b) 3-hydroxycyclopentyl-alanine (C1O) from sample 11 (see Figure 11), and (c) tyrosine from sample 35.

C1O (b) is a new non canonical amino acid that interacts with pocket residues that neither the seed (a) nor a chemically similar AA at the same position (c) engage

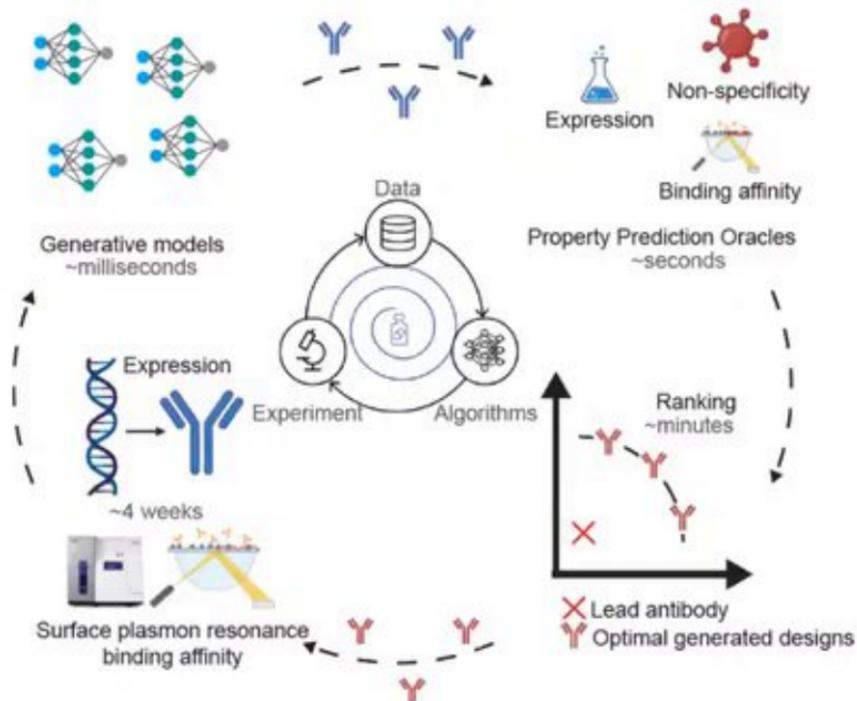
Evaluation

In-silico evaluation on public benchmarks

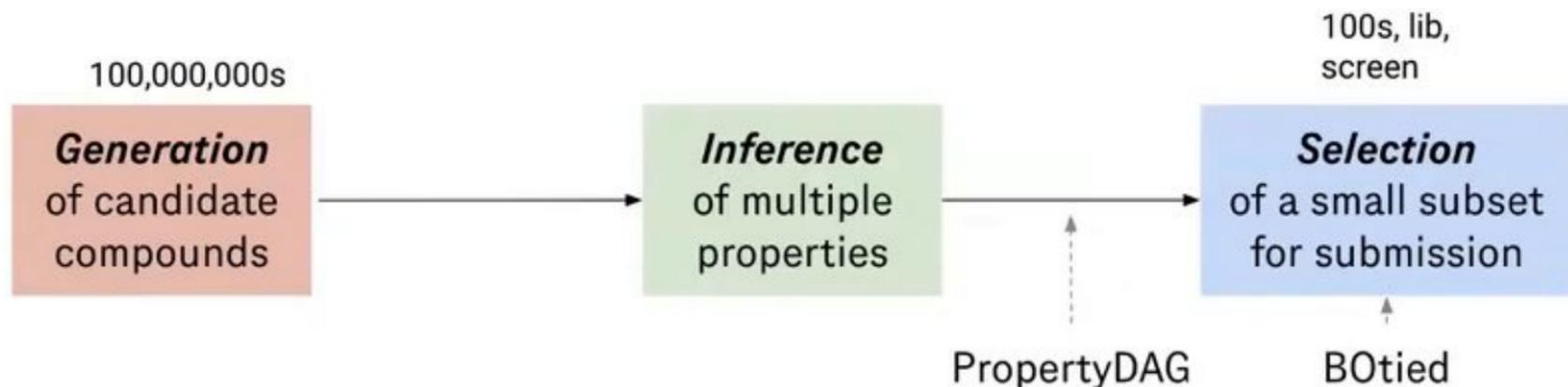
- state-of-the-art on pocket conditioned small molecule generation
- outperforms the state-of-art on epitope conditioned CDR redesign
- promising results on our new (to be released) benchmark for pocket conditioned macro cyclic peptide generation

Wet lab validation

- For CDR H3 redesign, we achieved a 42% binding rate on our confident samples and found a 5X better binder than approved drugs for targets.



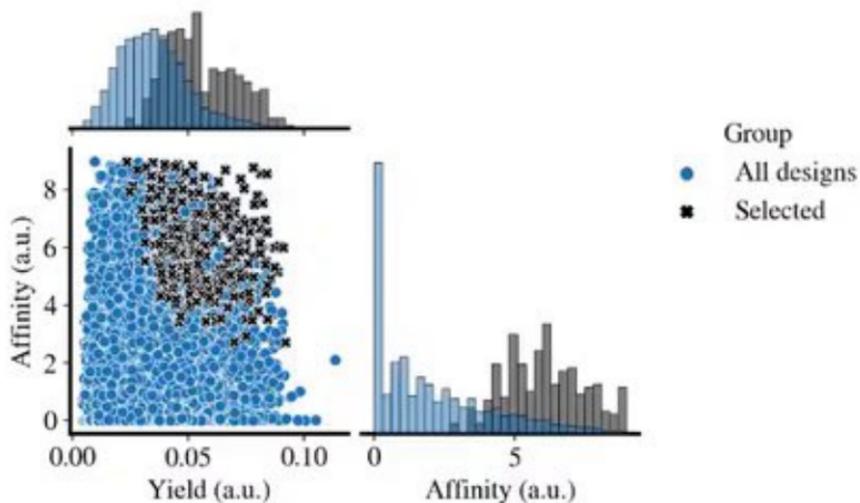
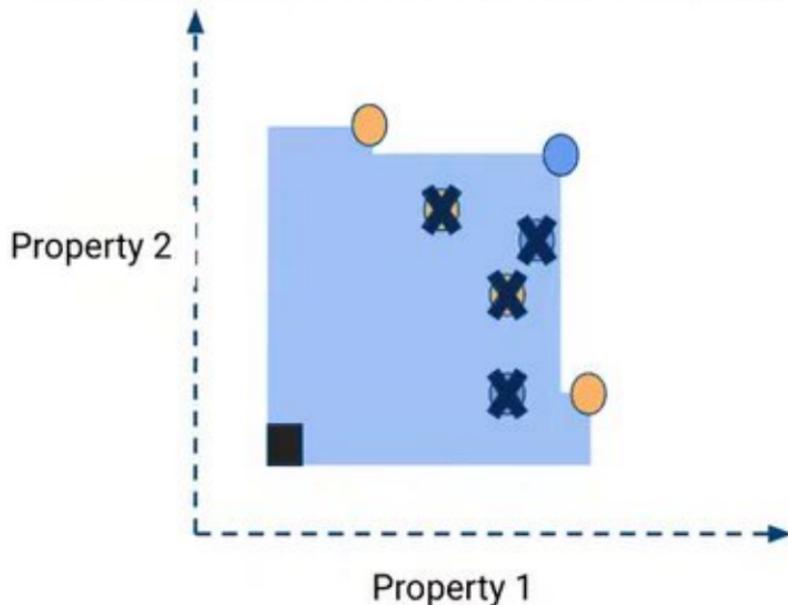
Bayesian optimization for design selection: PropertyDAG and BOTied



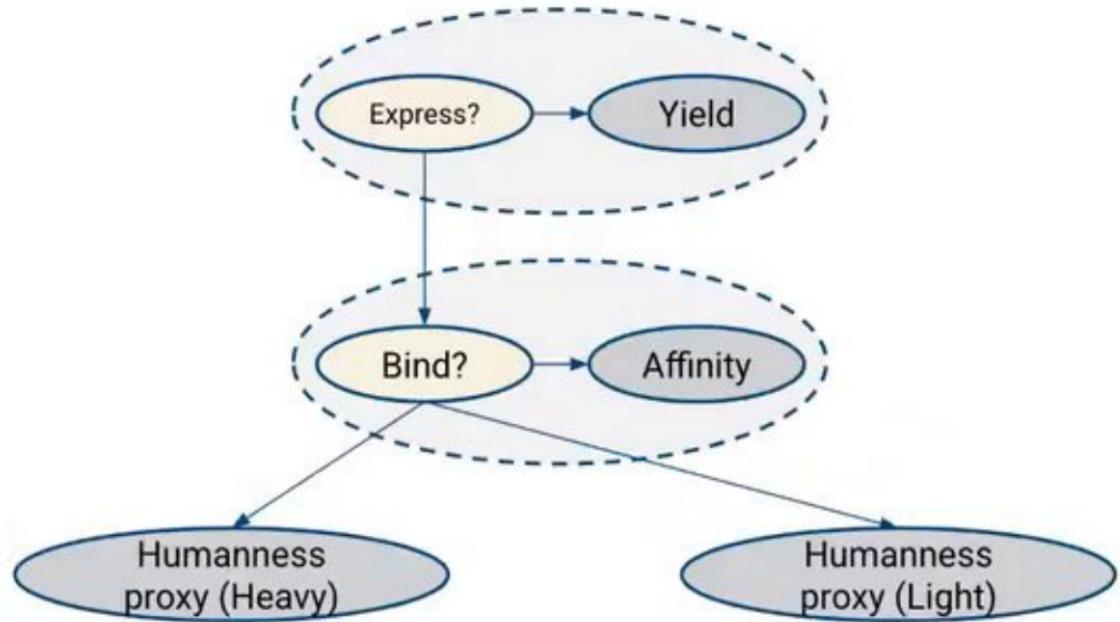
Park, Ji Won, et al. "PropertyDAG: Multi-objective Bayesian optimization of partially ordered, mixed-variable properties for biological sequence design" (2022) NeurIPS, AI4Science

Active learning: multi-objective sample selection through expanding the Pareto frontier

PropertyDAG: Expected hypervolume improvement (EHVI) guides exploration vs. exploitation balance towards expanding the Pareto frontier over multiple developability properties.

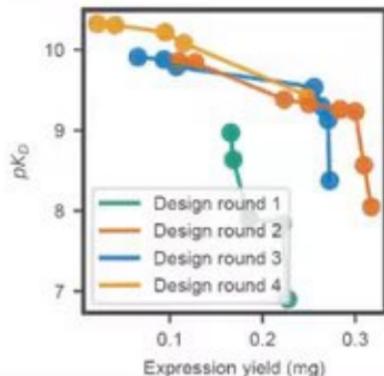
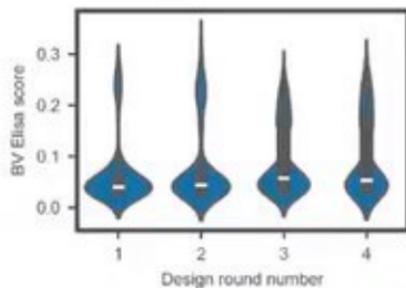
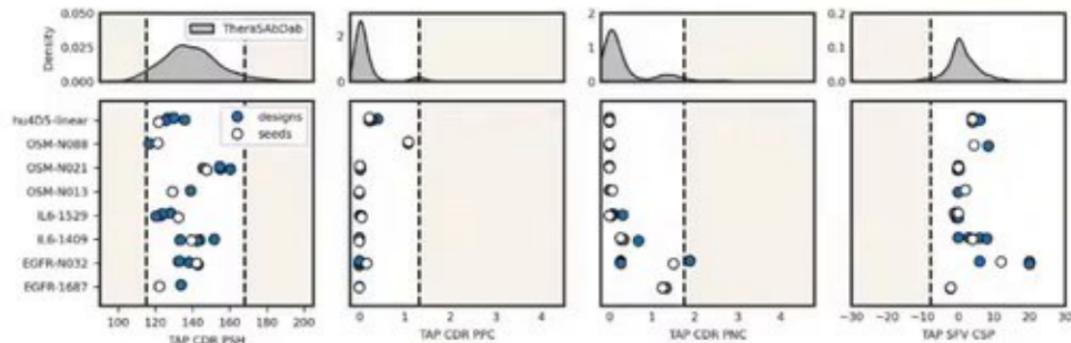


PropertyDAG: antibodies



Multi-property optimization enables therapeutic antibody design

Improving expression yield and in-silico developability properties of lead candidate antibodies





Generates large amounts of experimental data that is currently underutilized

Large scale screens

- HTS Screens (High throughput screening)
- DEL Screens (DNA Encoded Libraries)
- Fragment Screens

...

Biochemical and cell assays

- SAR potency assay data (Structure Activity Relationships)
- ADME assays (Absorption, Distribution, Metabolism, Excretion)
- Toxicology assay data

...

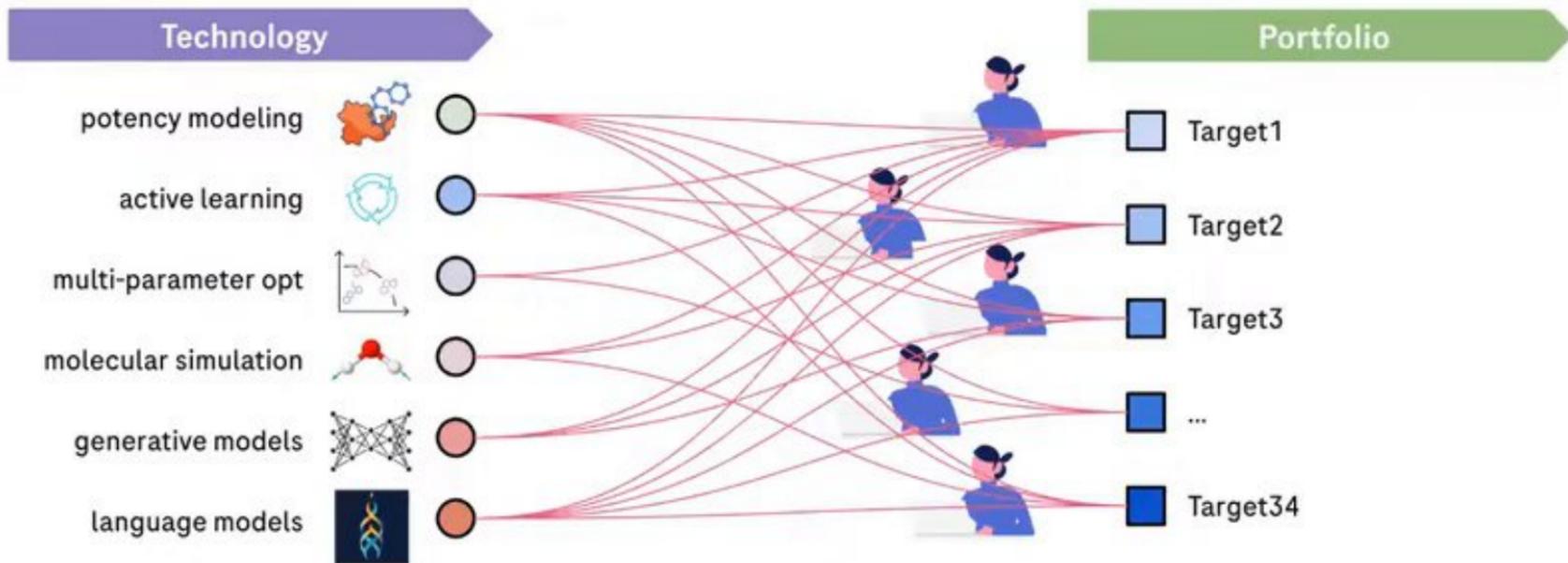
Animal studies

- Toxicology
- PK/PD (Pharmacodynamics, Pharmacokinetics)
- ADME (Absorption, Distribution, Metabolism, Excretion)

...

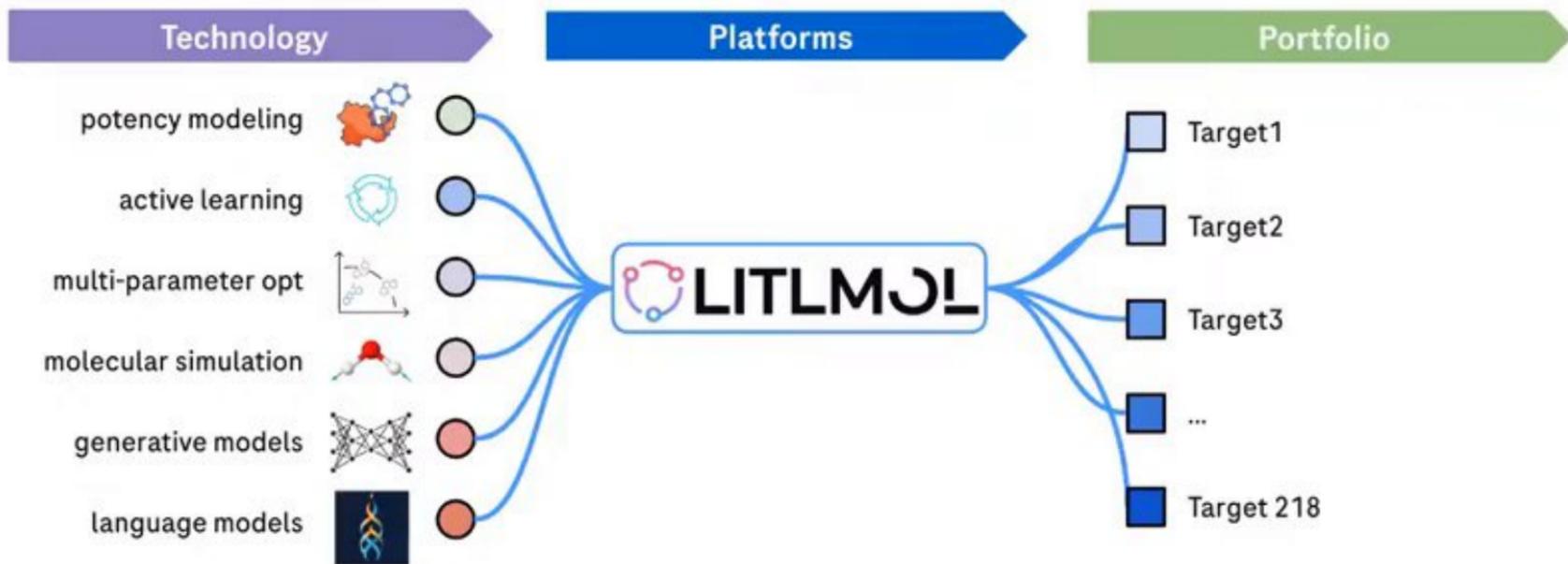


Small Molecule Drug Discovery without a common platform



Ad hoc adoption requires **significant** human resources and manual implementation
Does not scale - $O(N)$

Small Molecule Drug Discovery at Prescient



Instead, we invest in **building systems and platforms** to deliver on the pipeline

We are Clueless about Pretrained Models

- Publicly accessible commercial and open models were pretrained with a massive amount of data of unknown type, origin, composition and content.
- Biases, noise and intentional poisoning in such unknown data will inevitably sway these models' predictions, and we will never know why these models make such predictions.
- For mission-critical scenarios, such as in healthcare, security and others, we must ask ourselves whether we should and will allow such opacity.

Source: Kyunghyun Cho, CHIL '24 Keynote

Foundational models: Train and maintain **Genie-base** and **Genie-large**

Genie-base: 5-15B parameters, supports managed FT, RAG; used for prototyping and running locally

Genie-large: 30-70B parameters, used for complex applications



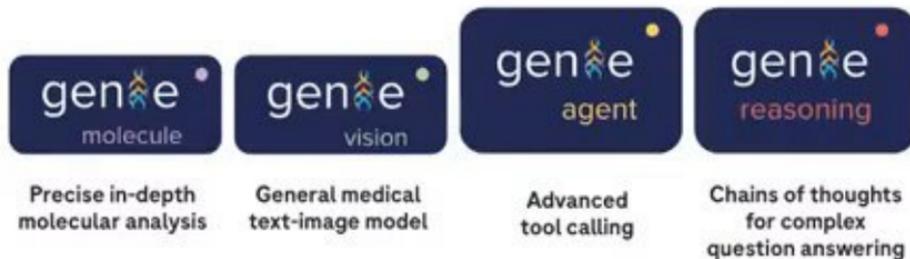
Improve **drug discovery-related** model capabilities

Tool calling to facilitate the lab-in-the-loop discovery

Reasoning: enable complex QA, multi-hop and better CoT

Multi-modality: train the model on protein sequences, small molecule structures, images

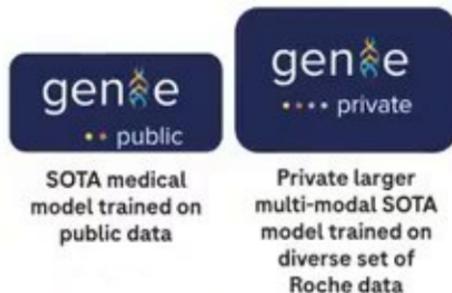
Chat: improve chat interface, instruction following for everyday tasks



Achieve **state-of-the-art** and make a **public version**

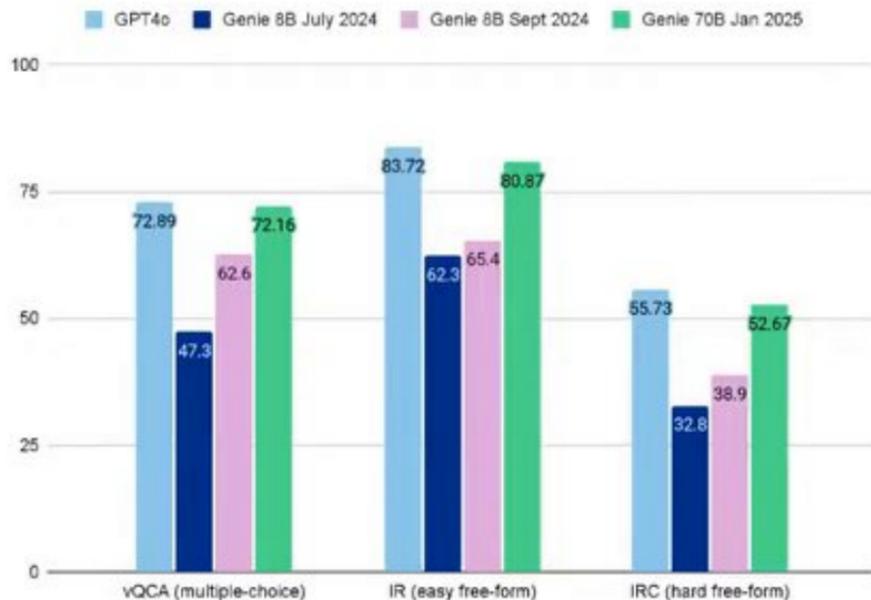
Public version: smaller version trained on a public dataset with agentic and reasoning improvements

Private version: our best state-of-the-art LLM for drug discovery



Early results: Genie for Target Safety Assessment

Accuracies in %



Internal models for core inference

- Safety
- Target assessment and data integration
- Digital twins
- ...

Early results: Genie for Tool Calling

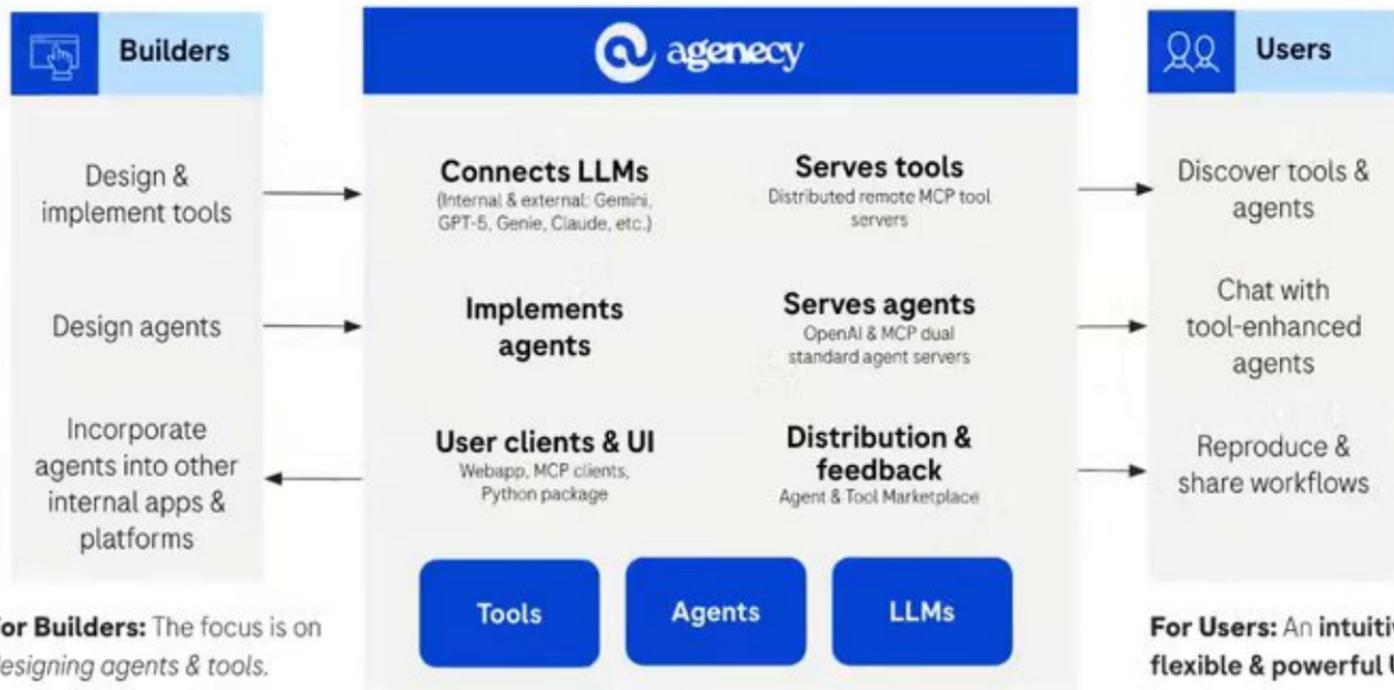
Accuracies in %



Models for Agents and Orchestration

- Tool calling
- Methods encapsulation and interface
- Automated and semi-automated prototyping and design
- Automated experimental design
- ...

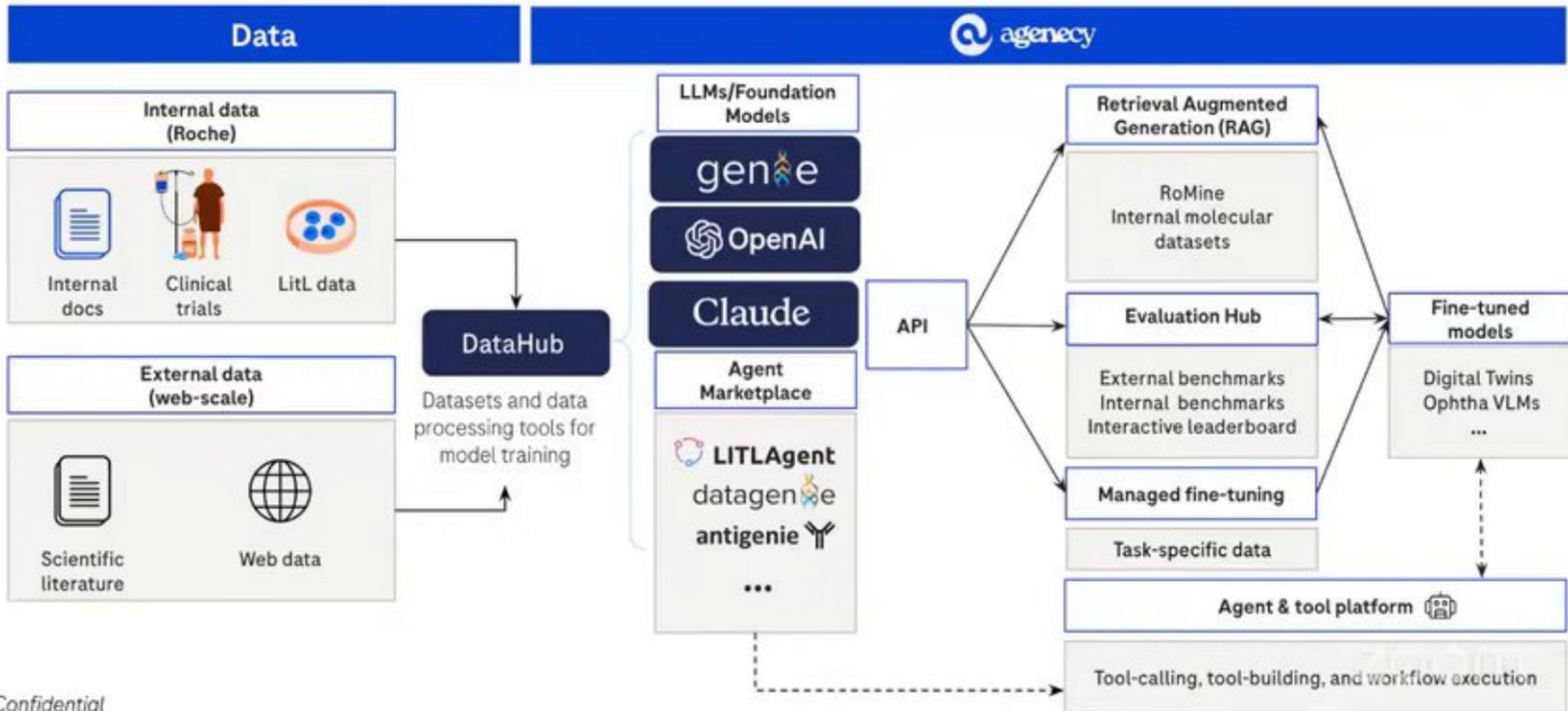
Agency: A unified platform for building & using agents



For Builders: The focus is on *designing agents & tools*. Complex engineering is **seamlessly abstracted away**.

For Users: An **intuitive, flexible & powerful UI** to use and chat with different LLMs, agents, and tools.

Agency: the Roche agentic ecosystem





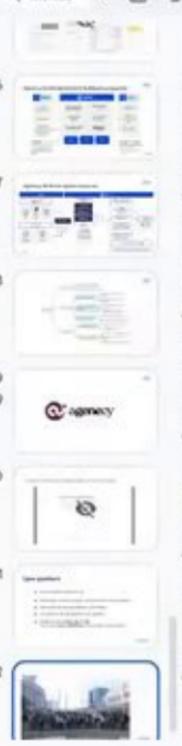


Open questions

- Core models and training
- Data layer, model registry, tool/method encapsulation
- GenAI+tools encapsulation and testing
- Integration, interoperability and uptime
- Embrace the **Long Tail ? OR**
Focus on **main workflows** that power value chain?

Thanks!





Click to add text

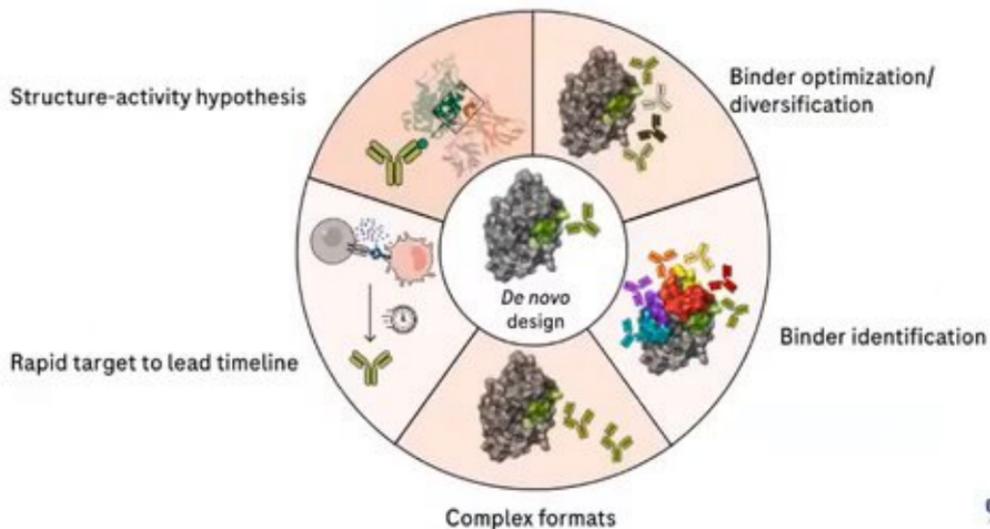
Thanks!

Genentech
A Division of the Roche Group

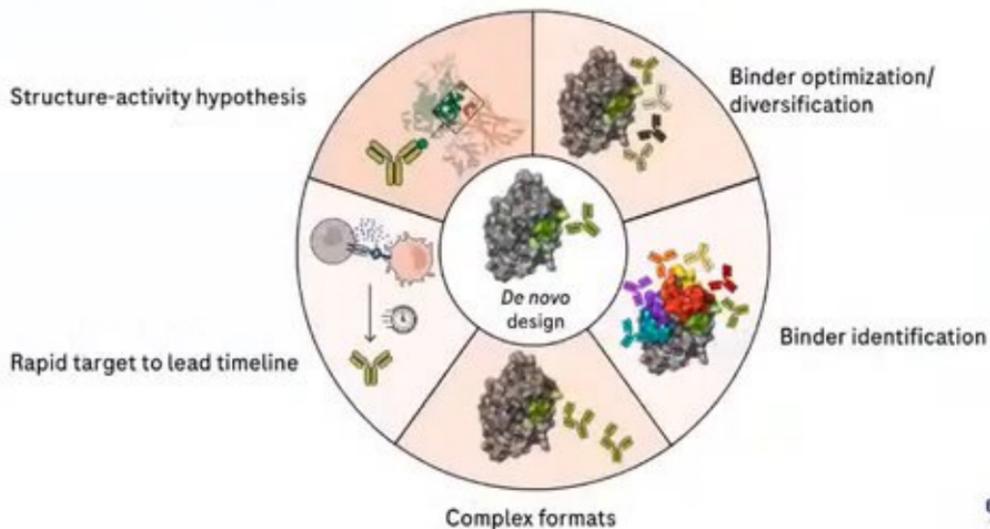


zoom

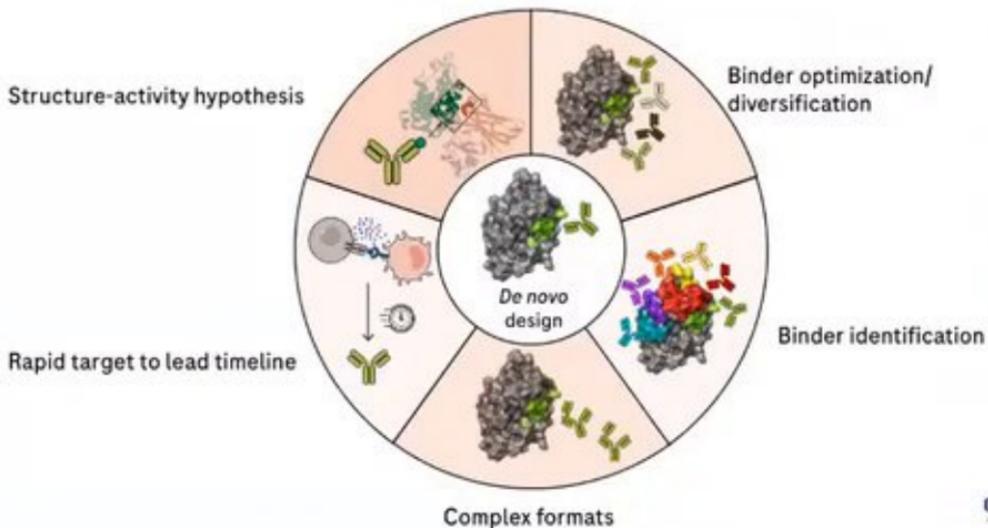
What would *de novo* antibody design enable?



What would *de novo* antibody design enable?



What would *de novo* antibody design enable?

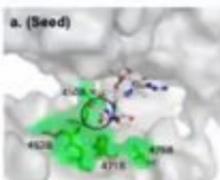


zoom

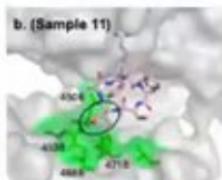
Click to add text

Prescient
Design
A Genentech Accelerator

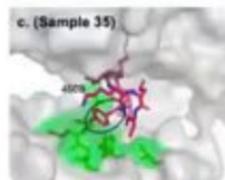
FuncBind generates new non canonical amino acids



6A	SER	430S	GLN	-1.891
6A	SER	433S	MET	-1.724
6A	SER	471S	SER	-0.762

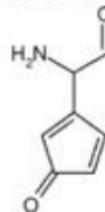


5A	C10	430S	GLN	-1.208
5A	C10	433S	MET	-2.196
5A	C10	466S	GLN	-1.122
5A	C10	471S	SER	-1.633



6A	TYR	430S	GLN	-2.804
----	-----	------	-----	--------

(3-hydroxy)
cyclopentyl-G



C10

Figure 12: Per-residue energy scores at the same position were calculated using Rosetta's residue energy breakdown for a seed MCP and two samples. The residues analyzed are: (a) serine from the seed, (b) 3-hydroxycyclopentyl-alanine (C10) from sample 11 (see Figure 11), and (c) tyrosine from sample 35.

C10 (b) is a new non canonical amino acid that interacts with pocket residues that neither the seed (a) nor a chemically similar AA at the same position (c) engage

18

Genentech
A Member of the Roche Group

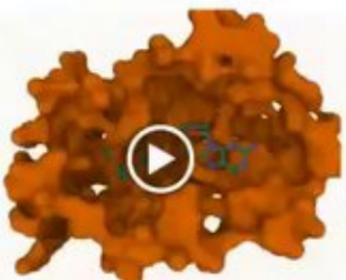
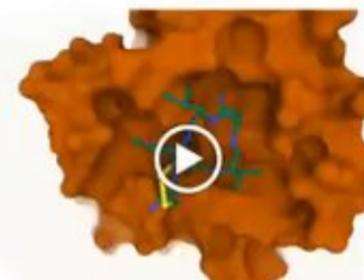
zoom



Click to add text

Prescient Design
A Genentech Acquisition

Pocket conditioned de-novo x-modal generation

small molecule	antibody	peptide
		

Genentech
A Member of the Roche Group



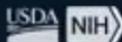
Rich Bonneau

mamp-mi:

A deep learning approach to predicting epitope immunogenicity in plants

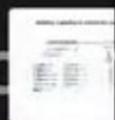
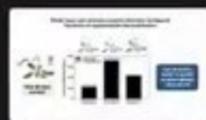
Danielle M. Stevens¹, David Yang², Tetiana J. Liang¹, Tianran Li¹, Brandon Vega¹, Gitika L. Cooker¹, Kazuma Krasileva^{1,2}

¹ Plant & Microbial Biology, UC Berkeley
² Center for Computational Biology, UC Berkeley
³ Plant Pathology, UC Davis



MCSB 0225

biksdj



Berkeley
UNIVERSITY OF CALIFORNIA

mamp-mi:

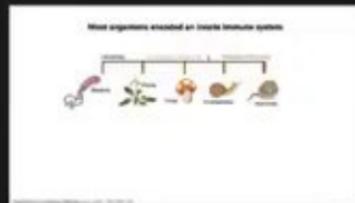
A deep learning approach to predicting epitope immunogenicity in plants

Daniella M. Stevens¹, David Yang², Tatiana J. Liang¹, Tianrun Li¹, Brandon Vega¹, Gitta L. Couker¹, Kazuo Krasakawa^{1,2}

¹ Plant & Microbial Biology, UC Berkeley
² Center for Computational Biology, UC Berkeley
³ Plant Pathology, UC Davis



MLCB 2023



bfkscjdb



Slide 1 of 23

I



1 2 3 4 5 6 7 8 9 10

1

mamp-ml

A deep learning approach to predicting optimal immunogenicity in plants

Danielle M. Browne¹, David Yang¹, Yuhua J. Liang¹, Brandon Vogel¹, Gabe L. Condon¹, Kaitia Kraibit¹

¹Plant & Molecular Biophysics, UC Berkeley

USDA NIH



2

Most organisms evolved an innate immune system



Plant immune system

Plant immune system

3

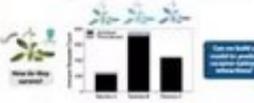
Plant immune system for disease breakthrough pathogen proliferation system



Plant immune system

4

Plants have also evolved adaptive immunity for rapid the spread of agricultural production

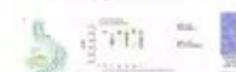


How to this system?

Can we help it control for genetic diversity?

5

Alphavirus like to produce consensus when an underlying structure is not defined



mamp-ml

A deep learning approach to predicting immunogenicity in plants

Danielle M. Browne¹, David Yang¹, Yuhua J. Liang¹, Brandon Vogel¹, Gabe L. Condon¹, Kaitia Kraibit¹

¹Plant & Molecular Biophysics, UC Berkeley

²Center for Computational Biology, UC Berkeley

³Plant Pathology, UC Davis

USDA NIH

bioRxiv

NYGC Events

mamp-ml:

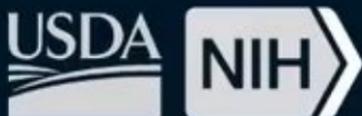
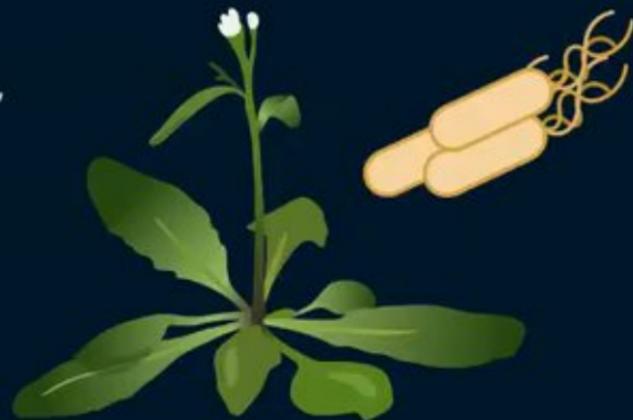
A deep learning approach to predicting epitope immunogenicity in plants

Danielle M. Stevens¹, David Yang², Tatiana J. Liang¹, Tianrun Li³,
Brandon Vega¹, Gitta L. Coaker³, Ksenia Krasileva^{1,2}

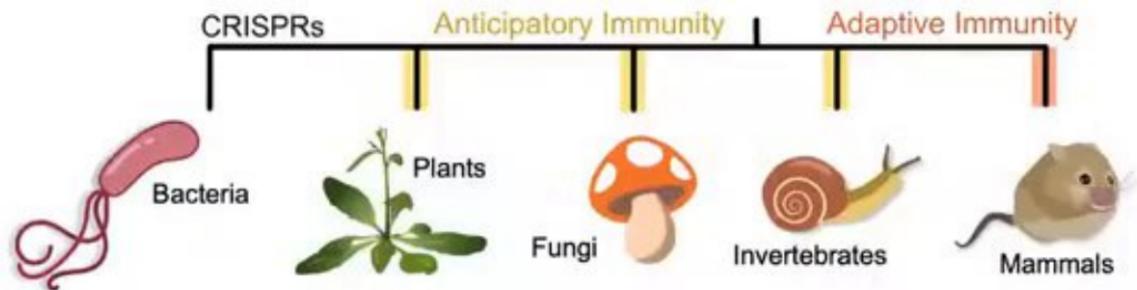
¹ Plant & Microbial Biology, UC Berkeley

² Center for Computational Biology, UC Berkeley

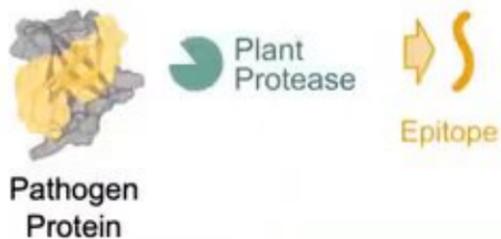
³ Plant Pathology, UC Davis



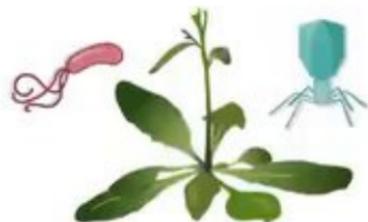
Most organisms encoded an innate immune system



Plant receptors survey for diverse threats through pathogen proteinaceous epitopes



Plants have vast immune receptor diversity far beyond the limits of experimental characterization

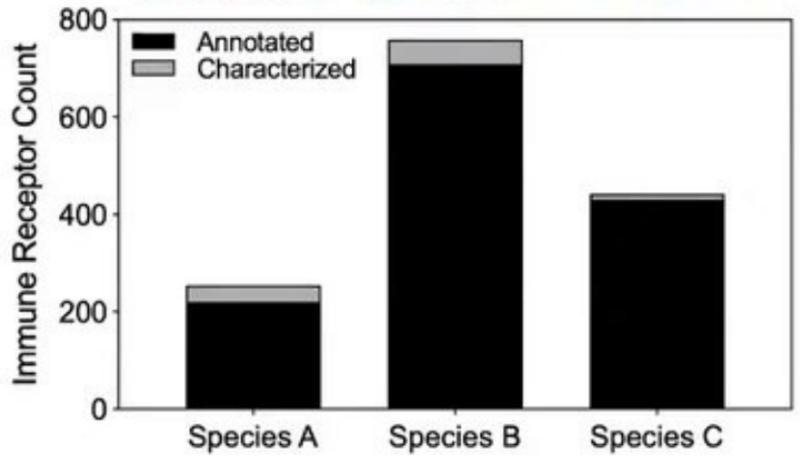


How do they survive?

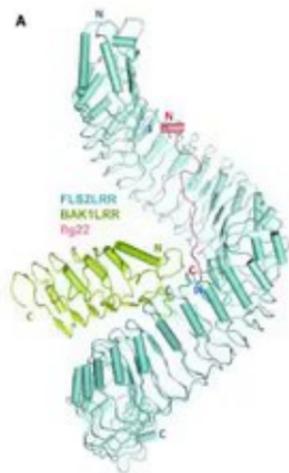
Plants have vast immune receptor diversity far beyond the limits of experimental characterization



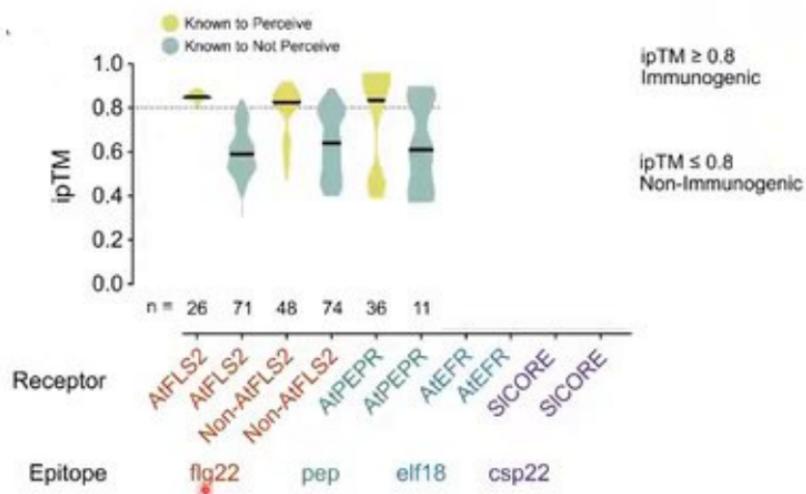
How do they survive?



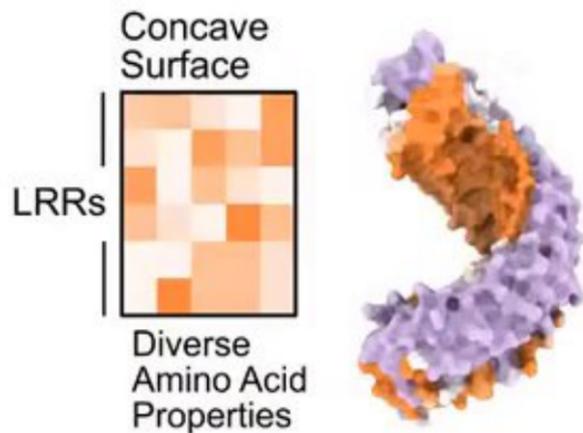
AlphaFold3 fails to predict outcomes when an underlying structure is not solved



AtFLS2-flg22



Goal: Fine-tune pLM for epitope-receptor immune outcomes



Generating the data needed for fine-tuning a language model

Scraping Data

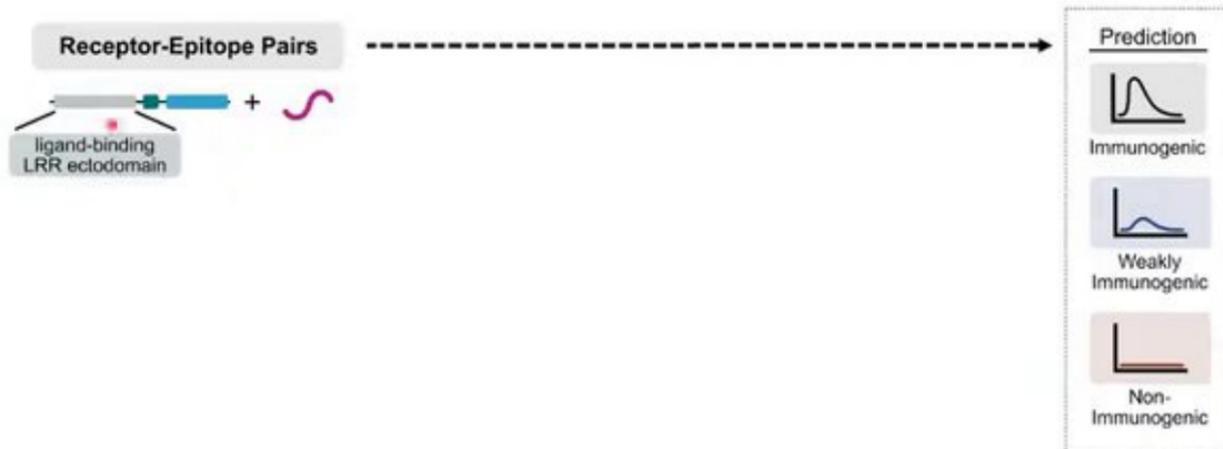
74 previous published studies from 1999 - 2025

91 plant species

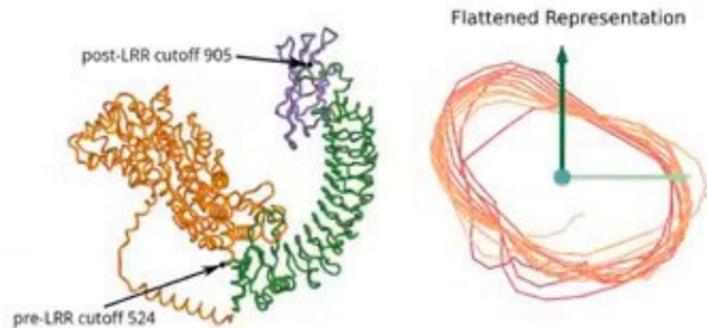
n = 1337 interactions

Epitope		Receptor		
bacteria	flg22	FLS2	RLKs	
	flgII-28	FLS3		
	csp22	CORE		
	elf18	EFR		
oomycete	pep-25	PERU		
plant	scoop	MIK2		
	screw	NUT		
herbivore	In11	INR		RLPs
	plant	crip21		
fungal	pg	RPL42		
multi-kingdom	nlp	RPL23		

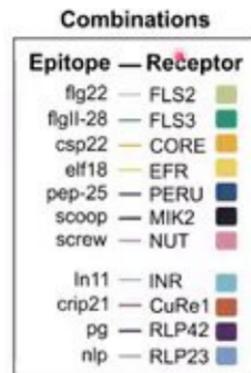
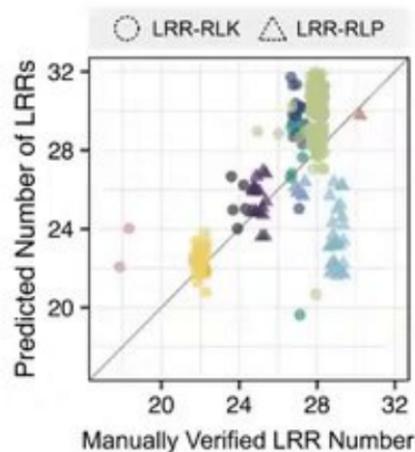
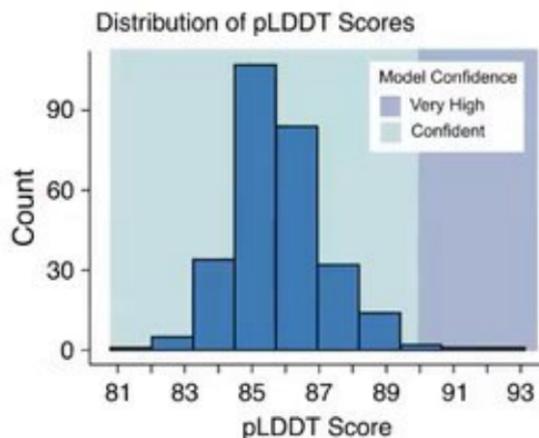
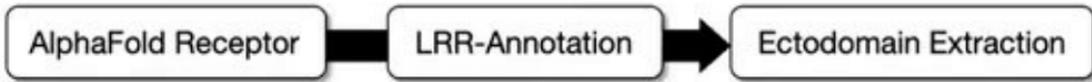
Building a pipeline to extract the epitope-binding LRR Ectodomain



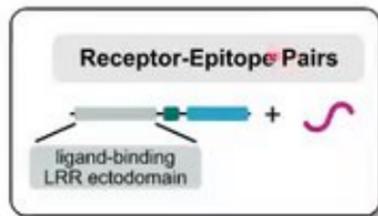
Using LRR-Annotation to extract LRR ectodomain and track surface residues



AlphaFold and LRR-Annotation extracts epitope-binding LRR Domain

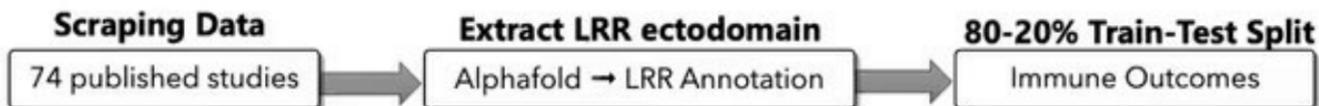


mamp-ml, a fine-tuned pLM enables prediction of immunogenic outcomes



In Feed-Forward NN:
 Loss Function = Cross-Entropy
 Dropout – 0.2, L2 Regularization to limit overfitting

mamp-ml outperforms embeddings only model

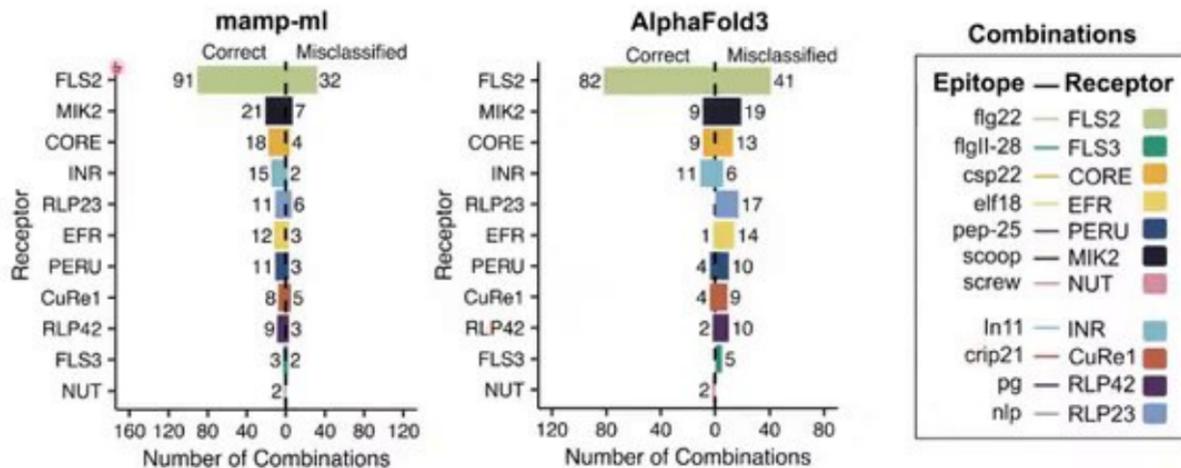


SeqOnly

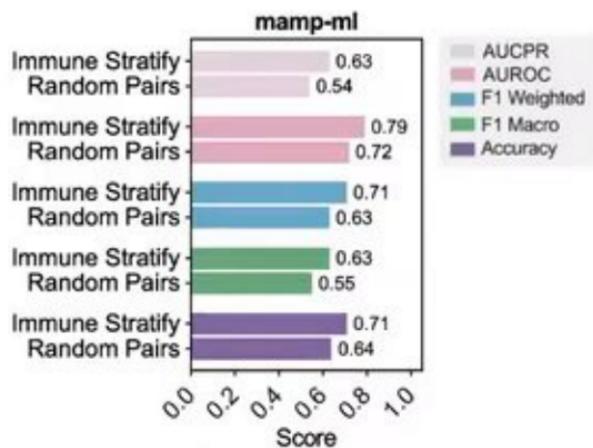
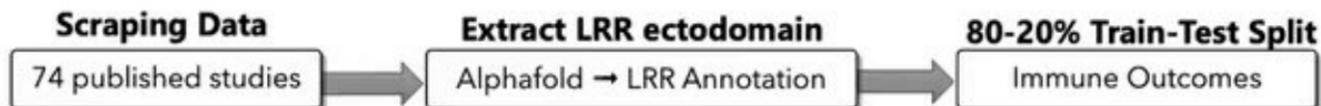
ESM2: pretrained & frozen



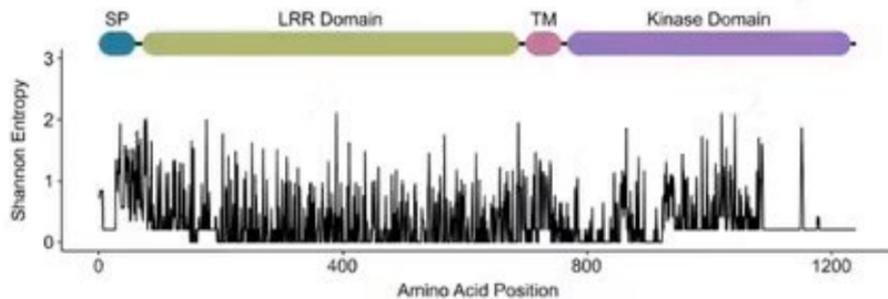
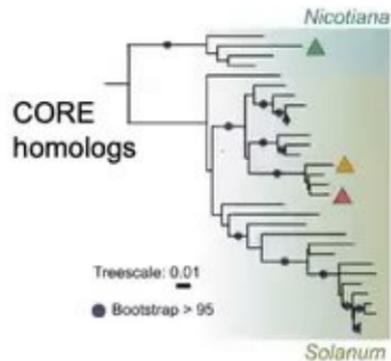
Mamp-ml predicts receptor-epitope immunogenicity without structural context



Validating mamp-ml for receptor-epitope interactions



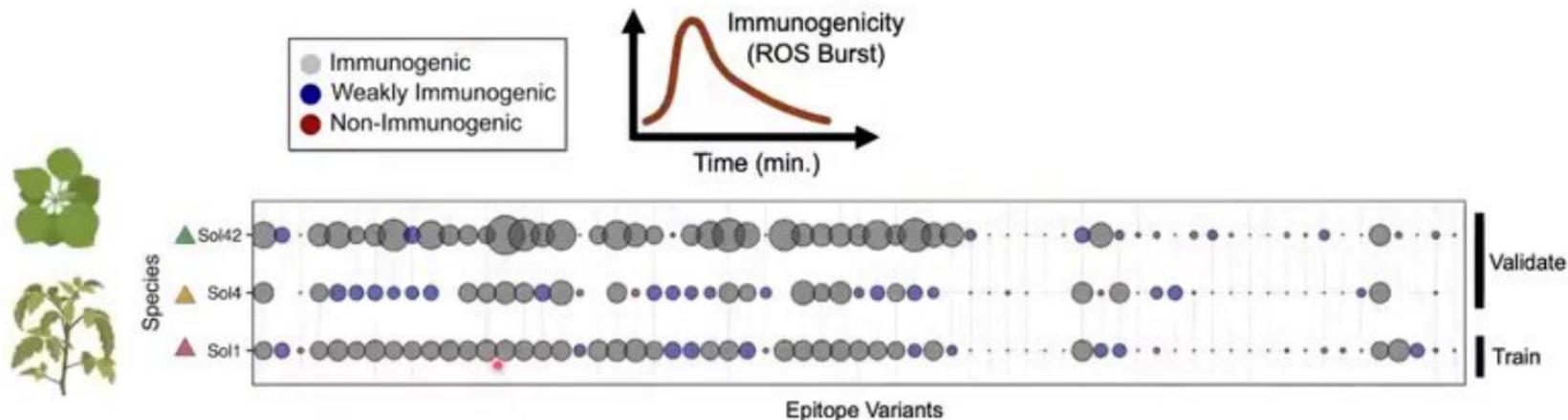
Using CORE-csp22 as a receptor-epitope validation dataset



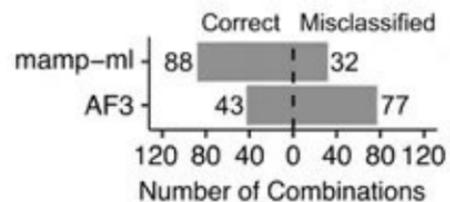
Collecting experimental validation data



Tatiana Liang,
UC Berkeley Undergraduate



Mamp-ml has improved prediction power over AlphaFold3



Conclusions and Future Directions

Fine-tuned models are more accurate than generalized approaches



Mamp-ml performs poorly to completely new receptor-epitope interactions and thus, needs to be generalized in the future



bioRxiv
THE PREPRINT SERVER FOR BIOLOGY



Available in
GitHub +
Google Colab

Acknowledgements

Krasileva Lab

- **Ksenia Krasileva**
- China Lunde Shaw
- Wei Wei
- Grace Stark
- Jude Edwards
- Nicole Dubs
- Kyungyong Seong
- Chandler Sutherland
- **Tatiana Liang**
- **Brandon Vega**



Collaborators

- **Gitta Coaker, UC Davis**
- **Jerry Li, UC Davis**
- Georg Felix, University of Tubingen
- **David Yang (Now at UCSF in Kortemme Lab)**
- Zack Lippman, Cold Spring Harbor
- Patrick Shih, UC Berkeley

UC Berkeley, Center for
Computational Biology



Berkeley
UNIVERSITY OF CALIFORNIA



Innovative
Genomics
Institute



Conclusions and Future Directions

Fine-tuned models are more accurate than generalized approaches

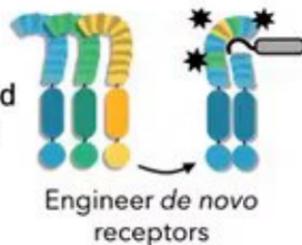


Mamp-ml performs poorly to completely new receptor-epitope interactions and thus, needs to be generalized in the future

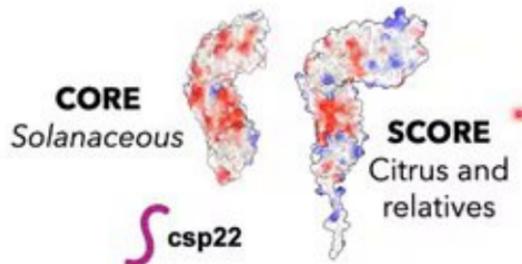
Mamp-ml accelerates *in silico* screening of new receptor homologs & epitope variants



Integration with generative module could accelerate engineering improved receptors



Mamp-ml may enable discovery of convergent evolved receptors to a shared target epitope



bioRxiv

THE PREPRINT SERVER FOR BIOLOGY



Available in
GitHub +
Google Colab

NYGC Events



- 1 Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT
- 2 Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT
- 3
- 4
- 5
- 6
- 7
- 8



Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT

Palash Sethi, Juannan Zhou
Department of Biology, University of Florida

MLCB 2025
September 11, 2025

zoom

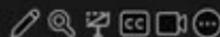
Hi.



Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT

Palash Sethi, Juannan Zhou
Department of Biology, University of Florida

MLCB 2025
September 11, 2025



Slide 1 of 33



Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT

Palash Sethi, Juannan Zhou
Department of Biology, University of Florida

MLCB 2025
September 11, 2025

Hi.



- 1 Learning gene interactions and functional landscapes from entire bacterial proteomes - BacPT
- 2 Learning gene interactions and functional landscapes from entire bacterial proteomes - BacPT
- 3
- 4
- 6
- 6
- 7
- 8



Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT

Palash Sethi, Juannan Zhou
Department of Biology, University of Florida

MLCB 2025
September 11, 2025



NYGC Events

Learning gene interactions and functional landscape from entire bacterial proteomes - BacPT

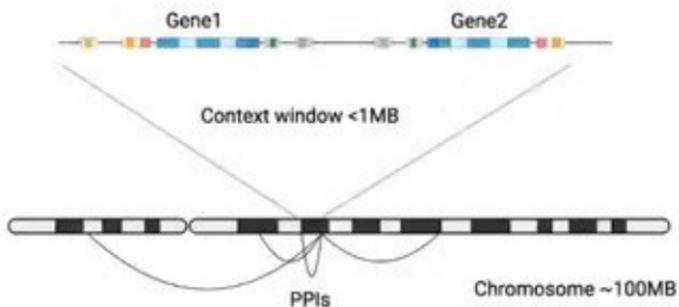
Palash Sethi, Juannan Zhou
Department of Biology, University of Florida

Motivation

Genome foundation models provide useful local contextual embeddings

Motivation

Contextualized whole-genome embeddings are underexplored



Long-range gene
interaction

Whole proteome model for bacteria

Bacterial genomes

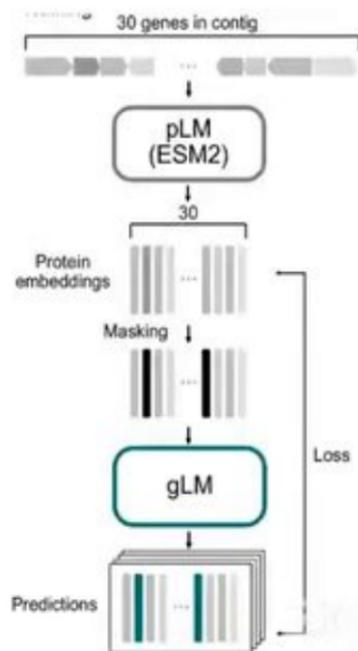
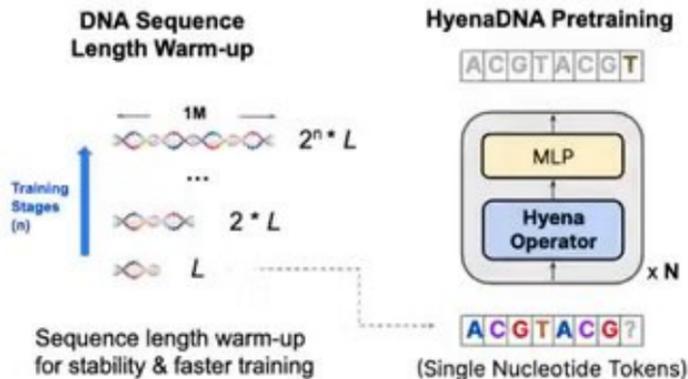
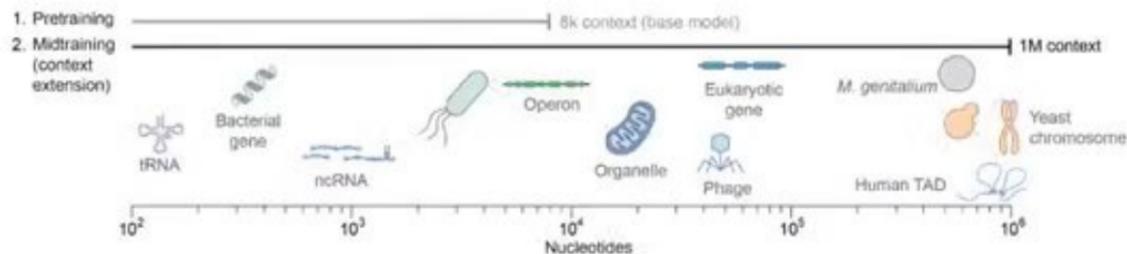
- Small (a few megabases, thousands of genes)
- Minimal intergenic regions (6-14%)
- Large sequence data availability (Refseq, MGnify)
- High-throughput phenotyping (e.g. BacDive database)

Applications of bacterial proteome foundation models

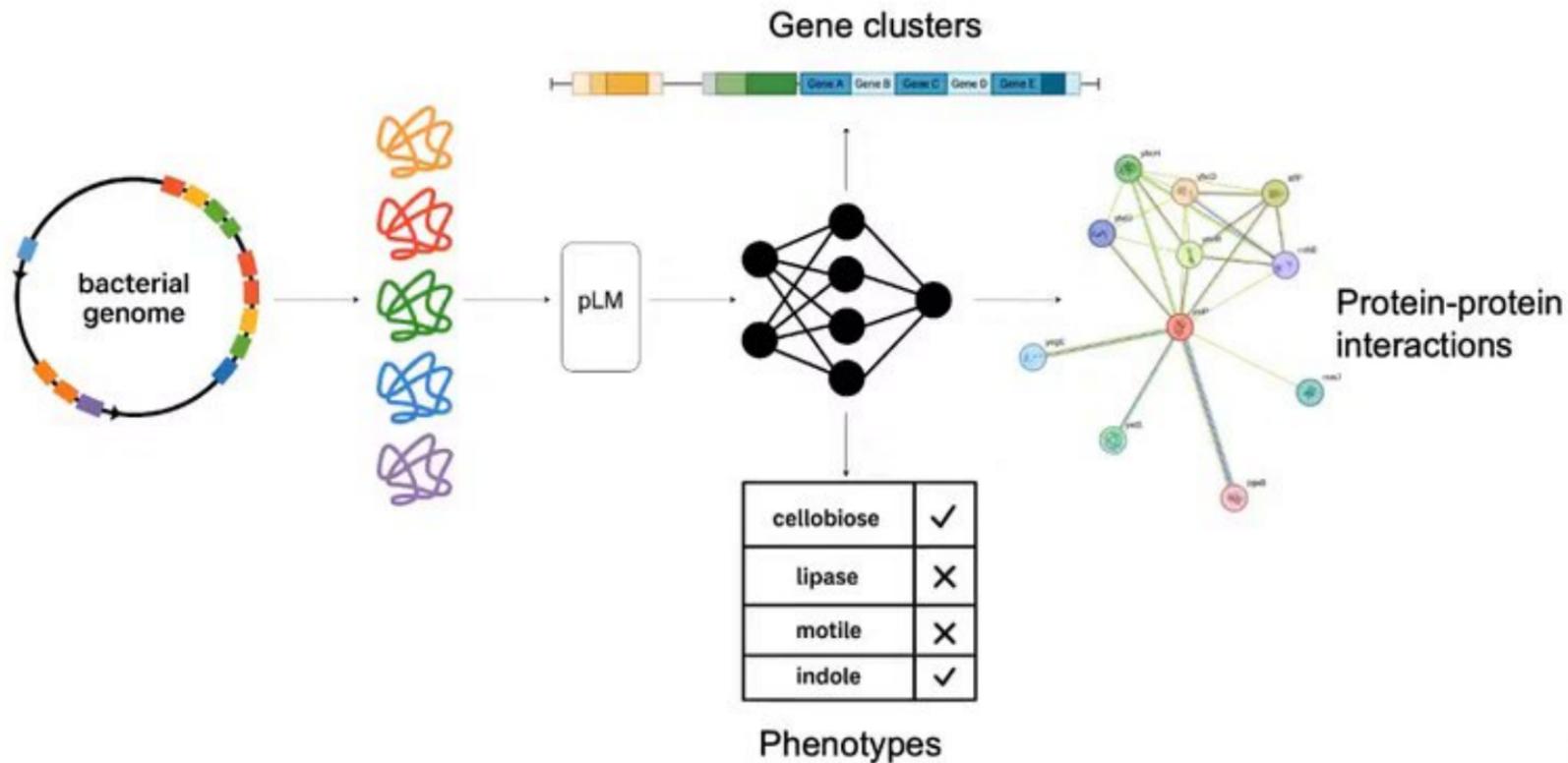
- Genome annotation
- Antibiotic resistance prediction
- Natural product discovery
- Microbiome studies

Motivation

Currency models are limited in their context lengths

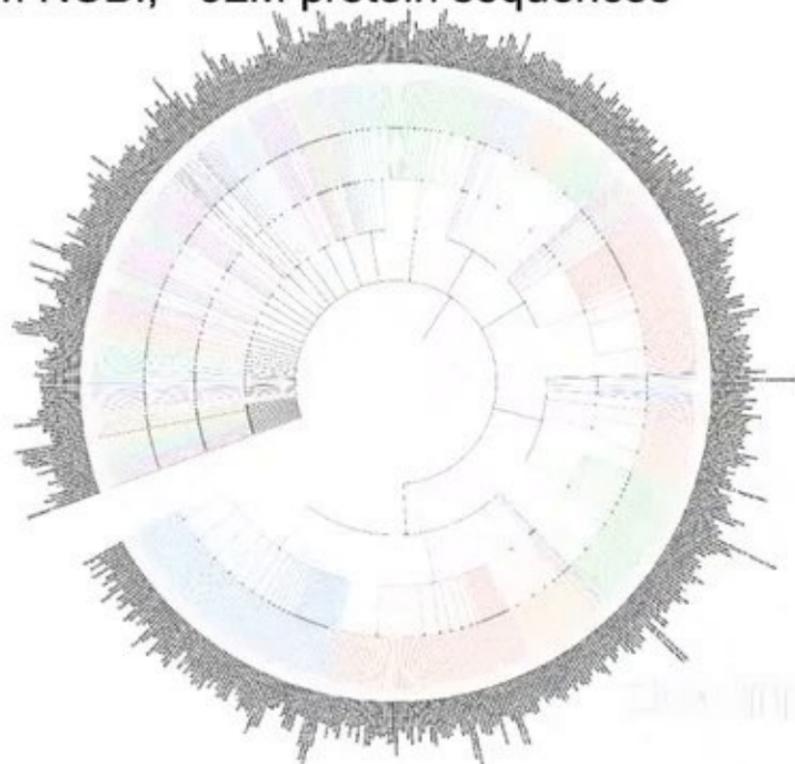
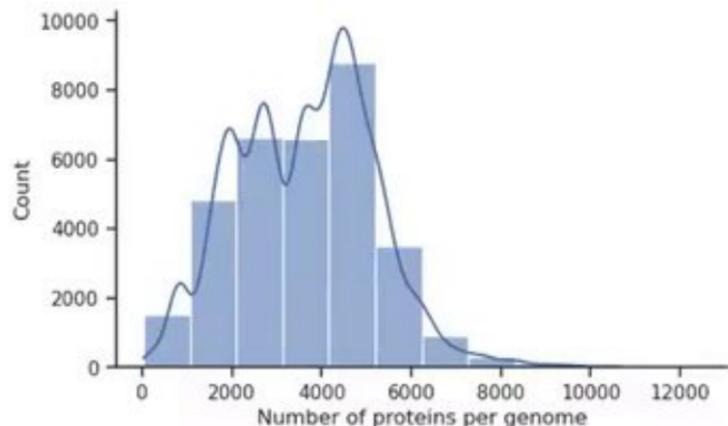


BacF1

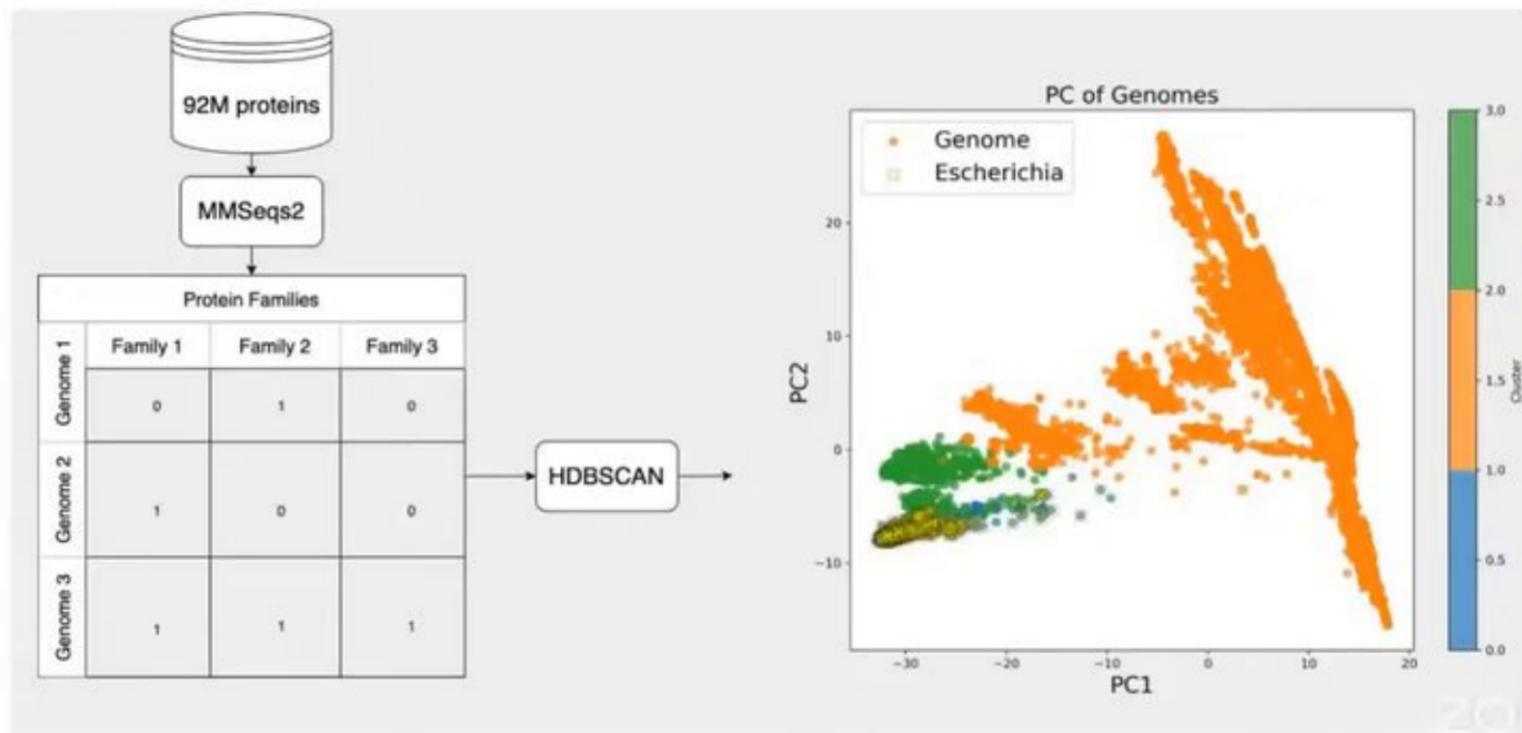


Methods - data

Trained on 33,140 bacterial proteomes from NCBI, ~92M protein sequences

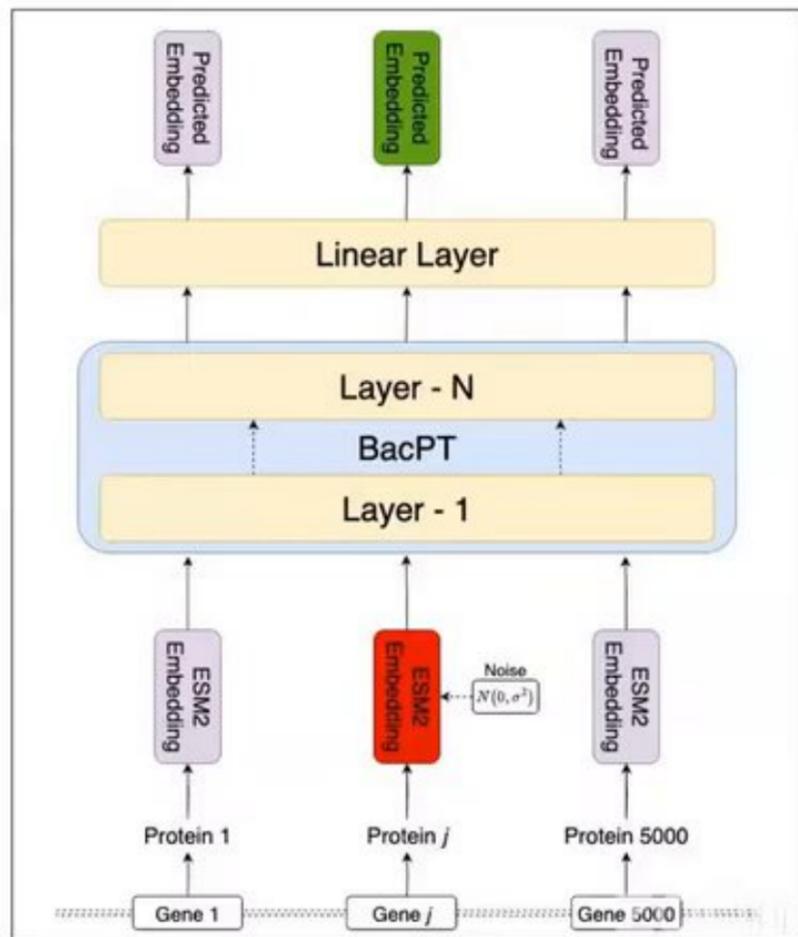


Methods - train/test split



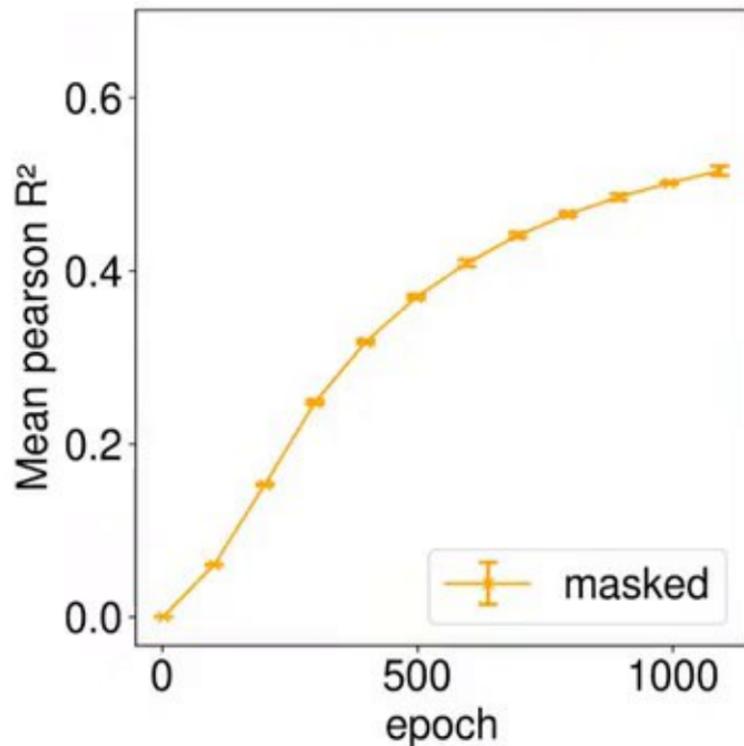
Methods - model architecture

- RoBERTa with rotary position embeddings
- 10 layers, 5 attention heads
- ESM2 - 35M model for input
- 5000 context length
- Trained via progressive noising with **MSE** loss to reconstruct input ESM embeddings
- 1100 epochs over 12 days on 16 NVIDIA A100 GPUs



Performance metrics - 1

- Correlation between predicted protein embeddings and ESM embeddings on test genomes

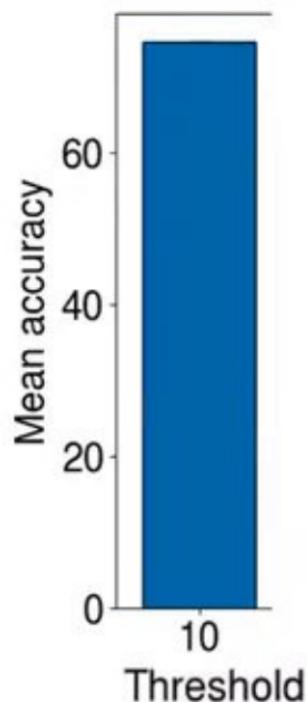


Performance metrics - 2

Metrics	BacPT	gLM
Pseudo Accuracy (%)	63.4	59.2
Absolute Accuracy (%)	78	71.9

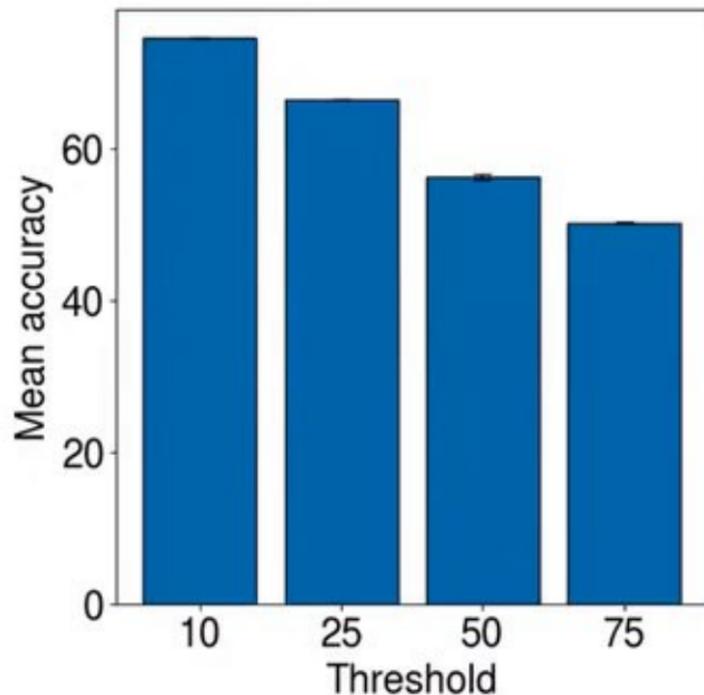
Performance metrics - 3

- Mean accuracy of BacPT at predicting the correct gene cluster for masked proteins, based on 10,000 gene cluster identified in training data



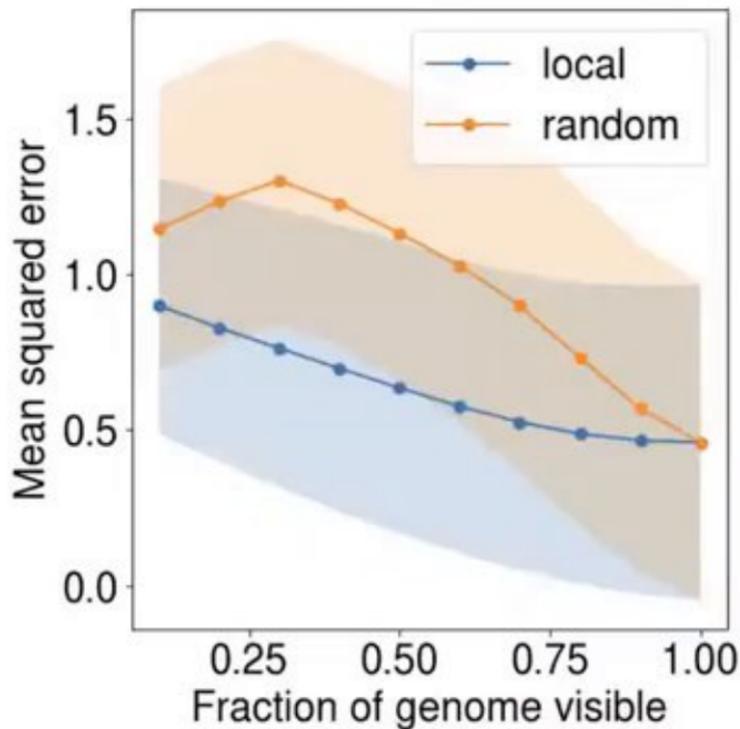
Performance metrics - 3

- Mean accuracy of BacPT at predicting the correct gene cluster for masked proteins, based on 10,000 gene cluster identified in training data



Performance metrics - 4

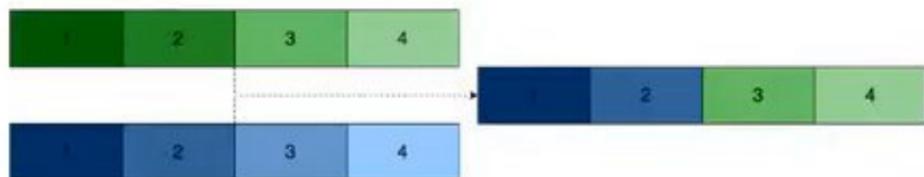
- Predictions for masked proteins improve when larger proportions of genomes are visible



Performance metric - 5

Can BacPT capture genome integrity?

Yes, linear probes can classify natural versus chimeric genomes

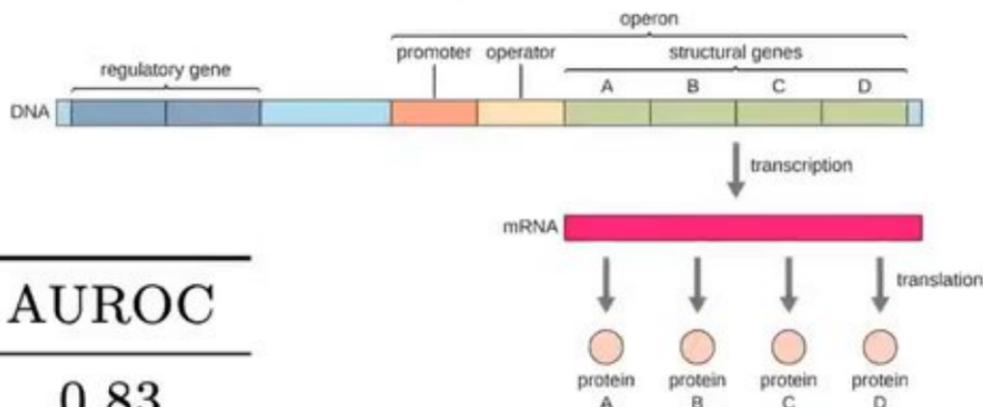


Model	Accuracy	AUROC
BacPT	0.77	0.82
ESM2	0.47	0.53

Applications

1. Operons
2. Protein-protein interactions
3. Whole-proteome representations

BacPT can identify operons

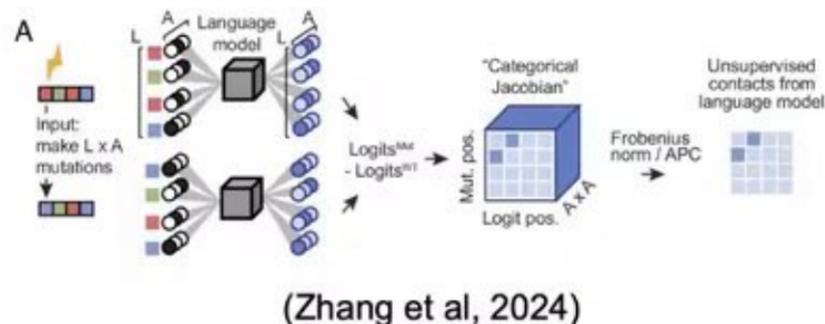
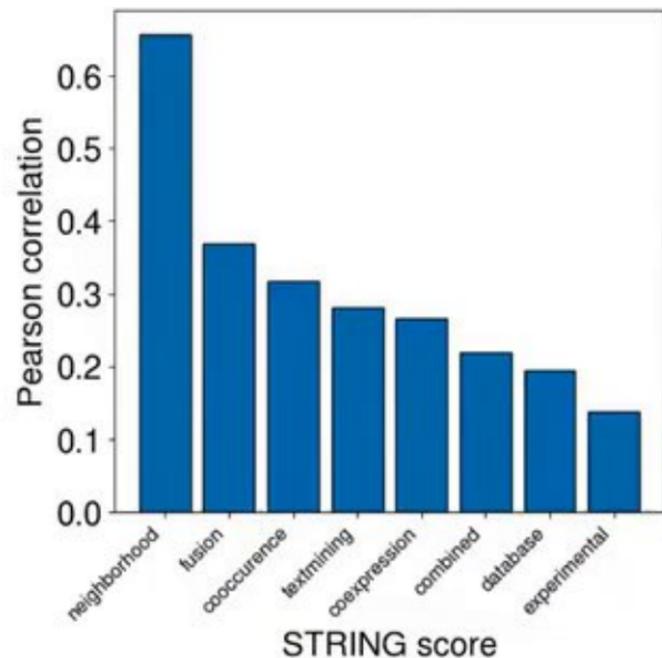


MODEL	ACCURACY	AUROC
BACPT	0.75	0.83
ESM2	0.66	0.75

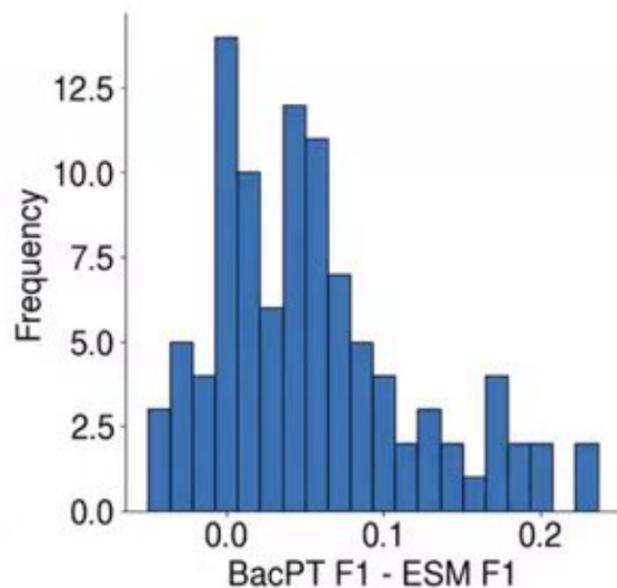
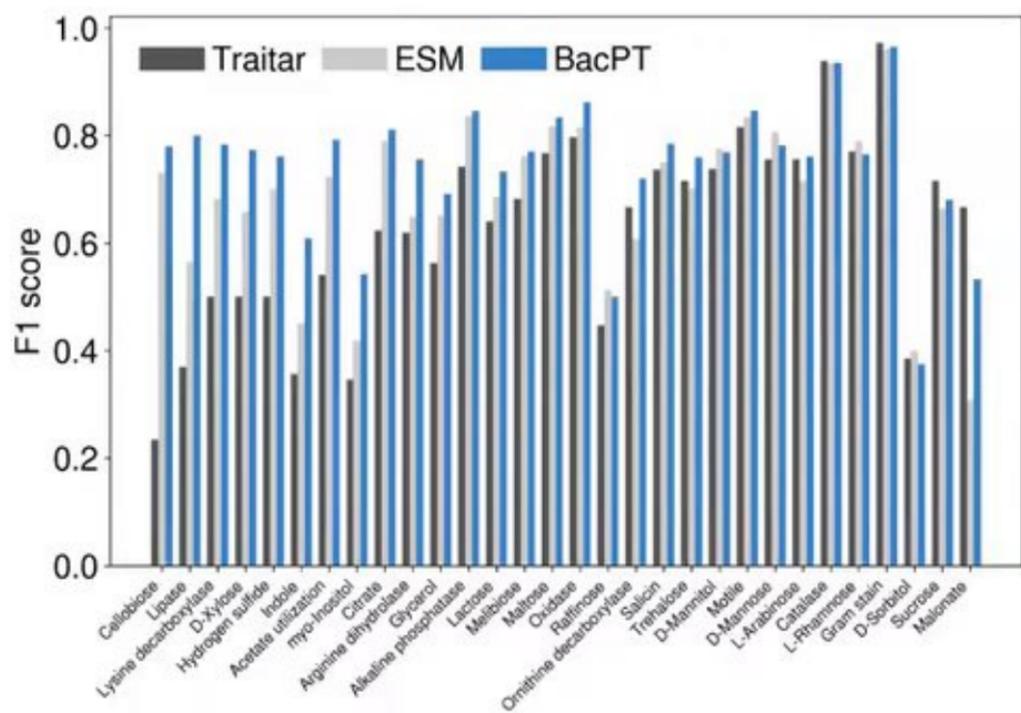
Linear probing of operon embeddings from BacPT outperform ESM2

BacPT learns protein-protein and other gene interactions

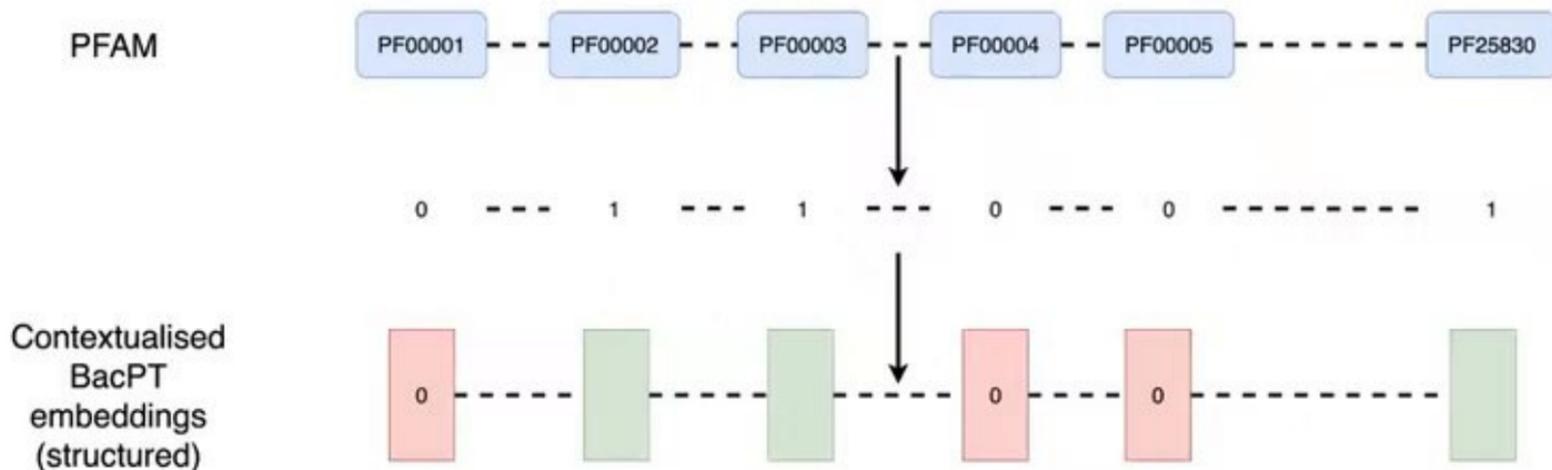
Using jacobian



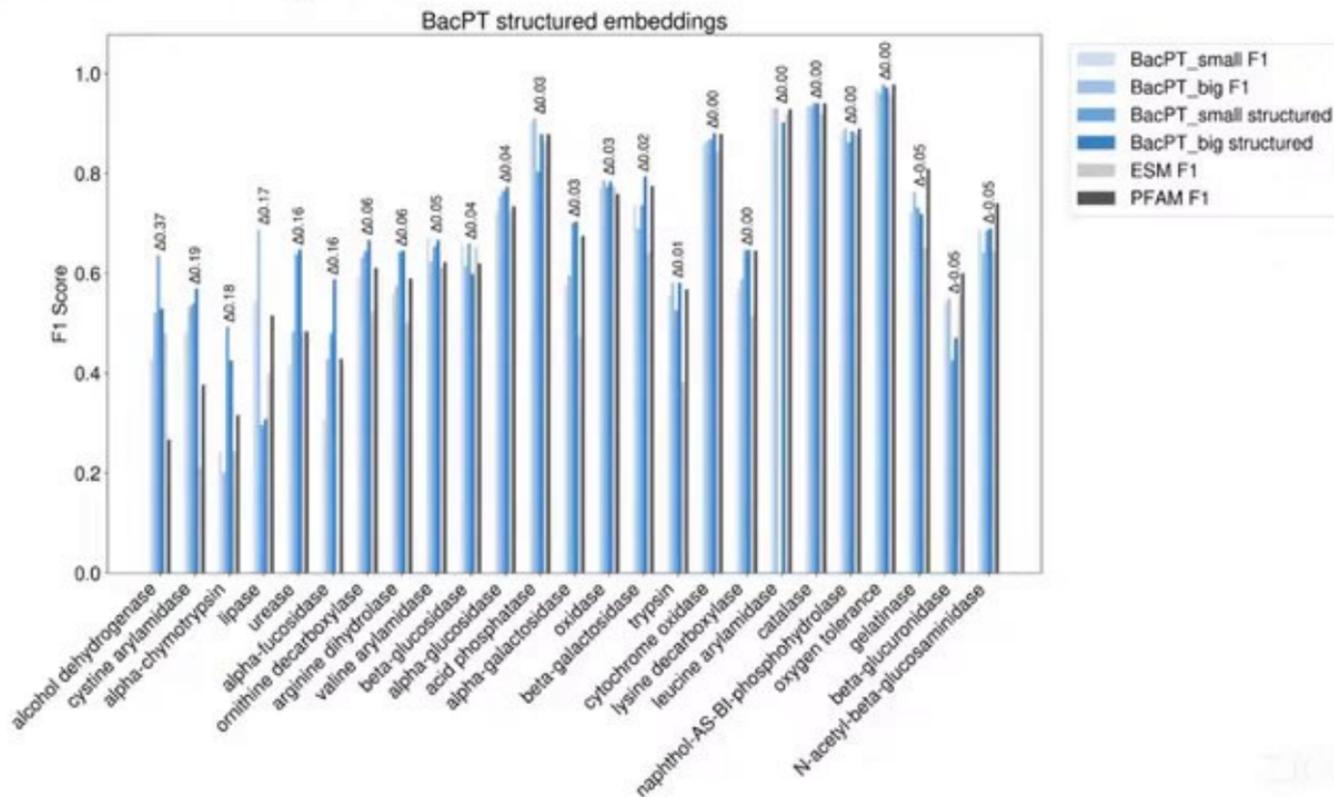
BacPT predicts organism-level traits



BacPT structural embeddings



BacPT predicts organism-level traits



Conclusion

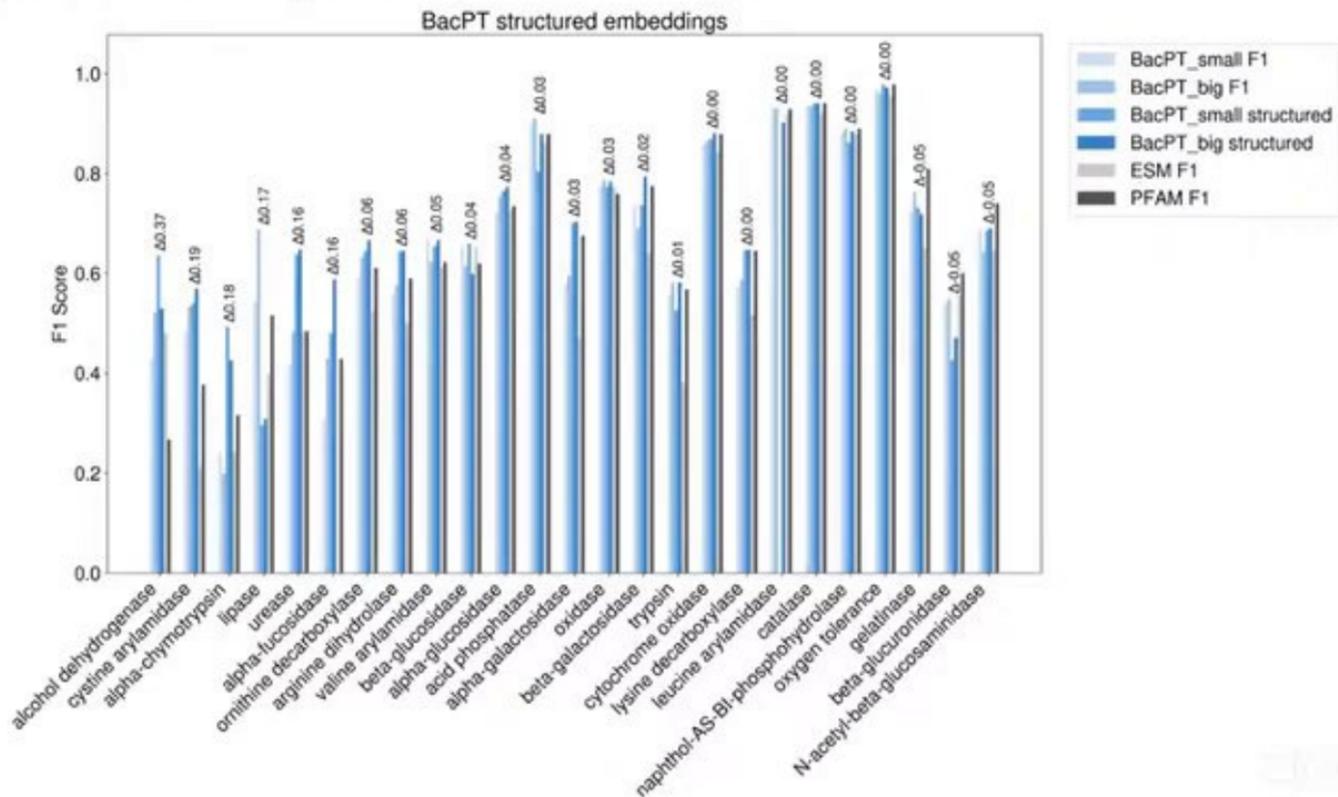
Whole proteome bacterial embeddings via BacPT

- Captures whole-proteome context length
- Identifies higher order genomic structures such as operons
- Learns protein-protein interaction
- Predicts organism-level complex traits

Future Work

- Gene clusters
- Microbial interaction
- Attribution maps for complex traits

BacPT predicts organism-level traits



Conclusion

Whole proteome bacterial embeddings via BacPT

- Captures whole-proteome context length
- Identifies higher order genomic structures such as operons
- Learns protein-protein interaction
- Predicts organism-level complex traits

Future Work

- Gene clusters
- Microbial interaction
- Attribution maps for complex traits

Acknowledgements



Zhou lab @ UF

- Marc G Chevrette @ UW Madison



National Institutes
of Health



National Institute
of General Medical
Sciences

UF | Biodiversity Institute
UNIVERSITY of FLORIDA

Questions?

Learning gene interactions and functional landscapes from entire bacterial proteomes

Palash Sethi¹, Marc G Chevrette², and Juannan Zhou*¹

¹Department of Biology, University of Florida, Gainesville, FL, 32611

²Department of Microbiology & Cell Science, University of Florida, Gainesville, FL, USA



palash.sethi@ufl.edu

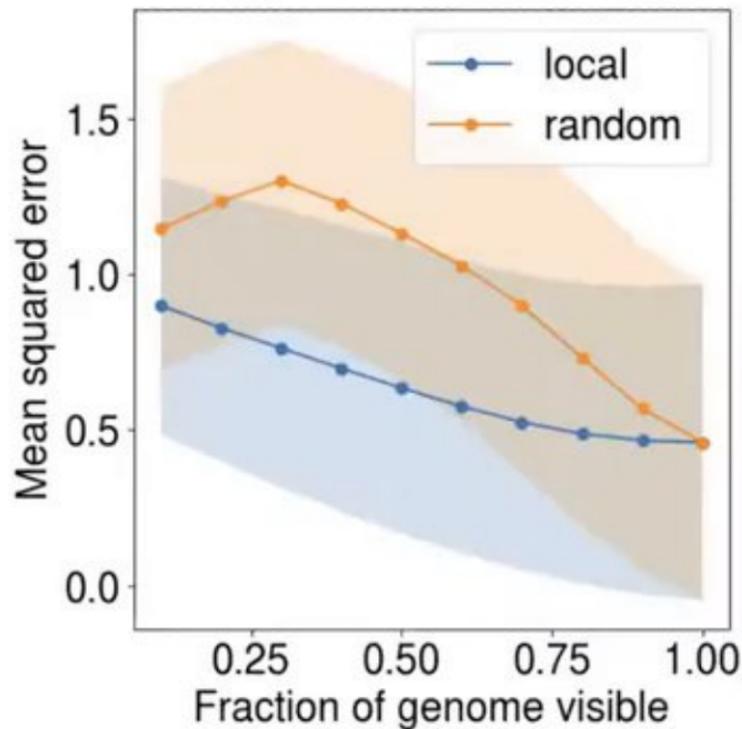


@PalashSethi6



Performance metrics - 4

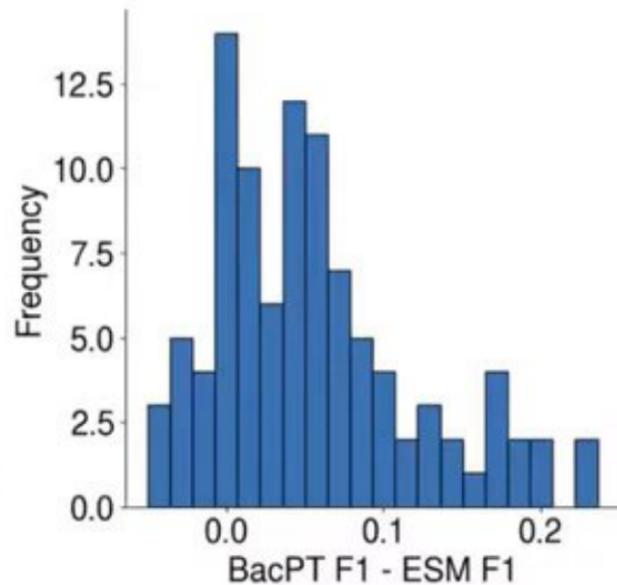
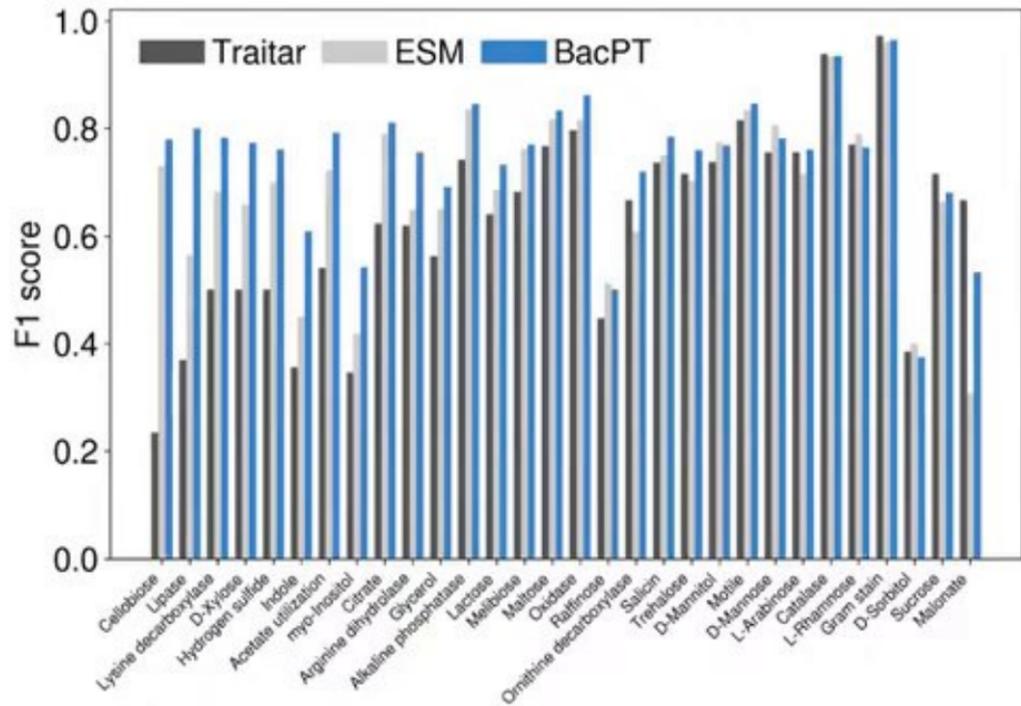
- Predictions for masked proteins improve when larger proportions of genomes are visible



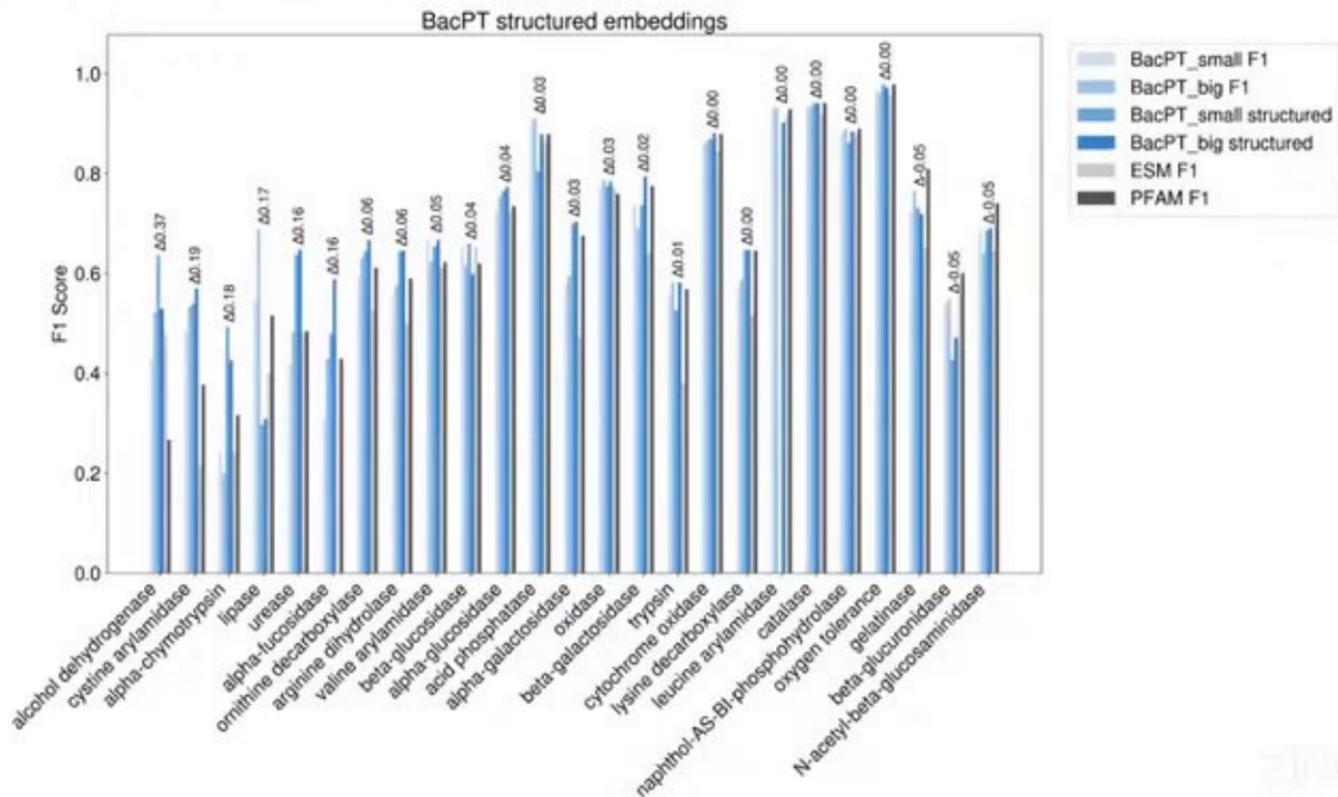
Applications

1. Operons
2. Protein-protein interactions
3. Whole-proteome representations

BacPT predicts organism-level traits



BacPT predicts organism-level traits



NYGC Events



Sliding Window Interaction Grammar (SWING): a generalized interaction language model for peptide and protein interactions

- 1 Sliding Window Interaction Grammar (SWING) is a generalized interaction language model for peptide and protein interactions.
- 2 The computational systems immunology lab
- 3
- 4

Jishnu Das
Assistant Professor, Center for Systems Immunology, Departments of Immunology and Computational & Systems Biology
Director, Computational Immunogenomics Core
University of Pittsburgh School of Medicine



MCLB 2025
September 11, 2025

www.jishnulab.org

Twitter: @jishnu1729, LinkedIn: www.linkedin.com/in/jishnu-das/

Click to add notes

Restarting at 1.30pm ET

mlcb.org for schedule



Sliding Window Interaction Grammar (SWING): a generalized interaction language model for peptide and protein interactions

Jishnu Das

Assistant Professor, Center for Systems Immunology, Departments of Immunology and Computational & Systems Biology
Director, Computational Immunogenomics Core
University of Pittsburgh School of Medicine



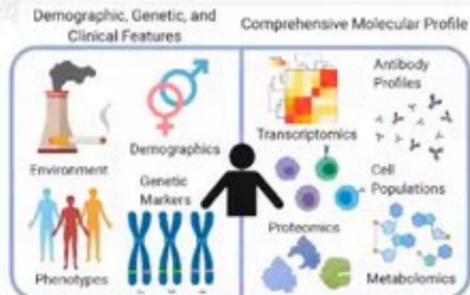
MCLB 2025
September 11, 2025



www.jishnulab.org

Twitter: [@jishnu1729](https://twitter.com/jishnu1729), LinkedIn: www.linkedin.com/in/jishnu-das/

Das computational systems immunology lab



Multi-omic integration to uncover correlates and mechanisms of immune regulation using predictive and interpretable machine learning

Network-based integration of genomic, epigenomic and transcriptomic datasets to uncover molecular phenotypes of immune disorders

*=co-first, ^=corresponding

Representative ML publications

Das et al *Bioinformatics* 2012

Ackerman, Das et al *Nature Medicine* 2018

Suscovich*, Fallon*, Das* et al *Science Translational Medicine* 2020

Das et al *PLoS Pathogens* 2020

Das et al *Med (Cell Press)* 2021

Wu, ... Das[^], Sperry[^], Billiar[^] *Annals of Surgery* 2022

Bing, ... Das[^] *Patterns (Cell Press)* 2022

Pedireddy, ... Das[^] *Cell Reports* 2022

Rahimikollu, Xiao, ... Das[^] *Nature Methods* 2024

Saha, Chakraborty, ... Das[^], Sarkar[^], *Science Translational Medicine* 2024

Chenchenko, Chhibbar, ... Das[^] *Nature Methods* (in press)

Representative networks publications

Wang*, Wei*, Thijssen*, Das* et al *Nature Biotechnology* 2012

Das et al *BMC Systems Biology* 2012

Das et al *Science Signaling* 2013

Wei*, Das* et al *PLoS Genetics* 2014

Das et al *Human Mutation* 2014

Vo*, Das* et al *Cell* 2016

Fragoza*, Das* et al *Nature Communications* 2019

Ningappa, ... Das[^] *Cell Reports Medicine* 2022

Berkowitz, ... Das[^] *Arthritis & Rheumatology* 2023

Chhibbar, Guha Roy ... Das[^] *Cell Reports* 2024

Rosen, ... Das[^], Torok *JCI Insight* 2025

Perturbations to cellular networks

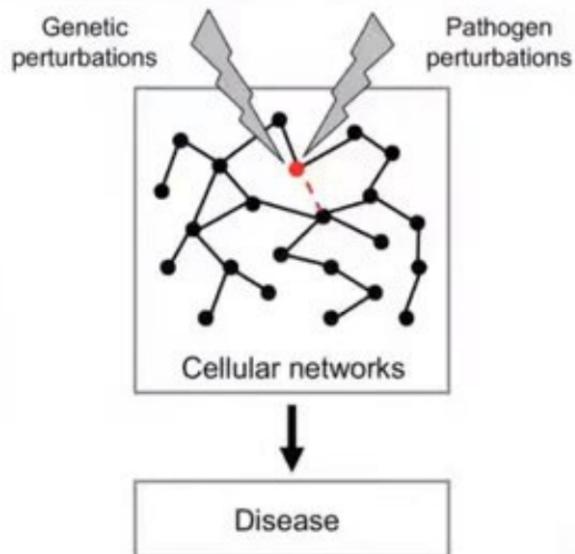
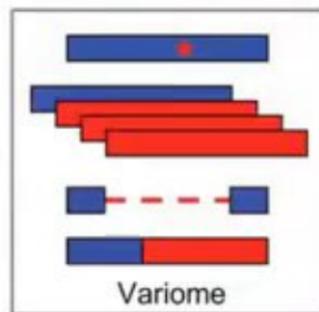
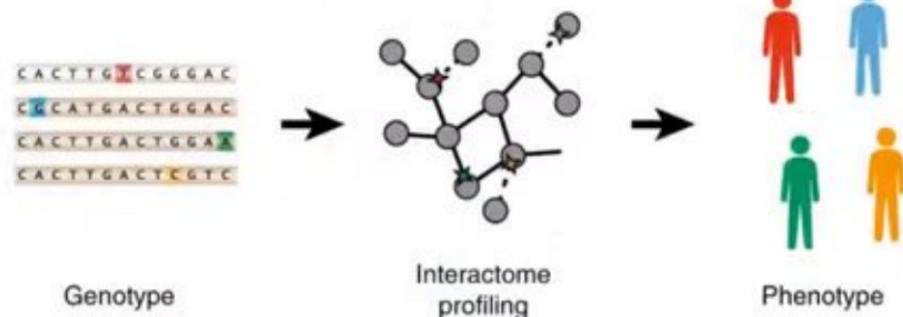


Current Opinion in Genetics & Development

Volume 23, Issue 6, December 2013, Pages 649-657

Edgotype: a fundamental link between genotype and phenotype

[Nidhi Sahni](#)^{1,2}, [Song Yi](#)^{1,2}, [Quan Zhong](#)^{1,2}, [Noor Jailkhani](#)^{1,2}, [Benoit Charlotteaux](#)^{1,2},
[Michael E. Cusick](#)^{1,2}, [Marc Vidal](#)^{1,2} ✉



Sliding Window Interaction Grammar (SWING): a generalized interaction language model for peptide and protein interactions

nature methods

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature methods](#) > [articles](#) > [article](#)

Article | [Open access](#) | Published: 28 July 2025

Sliding Window Interaction Grammar (SWING): a generalized interaction language model for peptide and protein interactions

[Jane C. Siwek](#), [Alisa A. Omelchenko](#), [Prabal Chhibbar](#), [Sanya Arshad](#), [AnnaElaine Rosengart](#), [Iliyan Nazarali](#), [Akash Patel](#), [Kiran Nazarali](#), [Javad Rahimikollu](#), [Jeremy S. Tilstra](#), [Mark J. Shlomchik](#), [David R. Koes](#), [Alok V. Joglekar](#)  & [Jishnu Das](#) 

[Nature Methods](#) **22**, 1707–1719 (2025) | [Cite this article](#)

15k Accesses | **1** Citations | **21** Altmetric | [Metrics](#)

J Siwek*, A Omelchenko*, P Chhibbar*... [J Das](#)[^] *Nature Methods* 2025

Key contributors



Alok Joglekar



Alisa Omelchenko



Jane Siwek



Prabal Chhibbar

Interaction language models to decode peptide & protein interactions

nature biotechnology

[Explore content](#) [About the journal](#) [Publish with us](#)

[nature](#) > [nature biotechnology](#) > [articles](#) > [article](#)

Article | [Published: 26 January 2023](#)

Large language models generate functional protein sequences across diverse families

nature biotechnology

[Explore content](#) [About the journal](#) [Publish with us](#)

[nature](#) > [nature biotechnology](#) > [articles](#) > [article](#)

Article | [Open access](#) | [Published: 24 April 2023](#)

Efficient evolution of human antibodies from general protein language models

nature machine intelligence

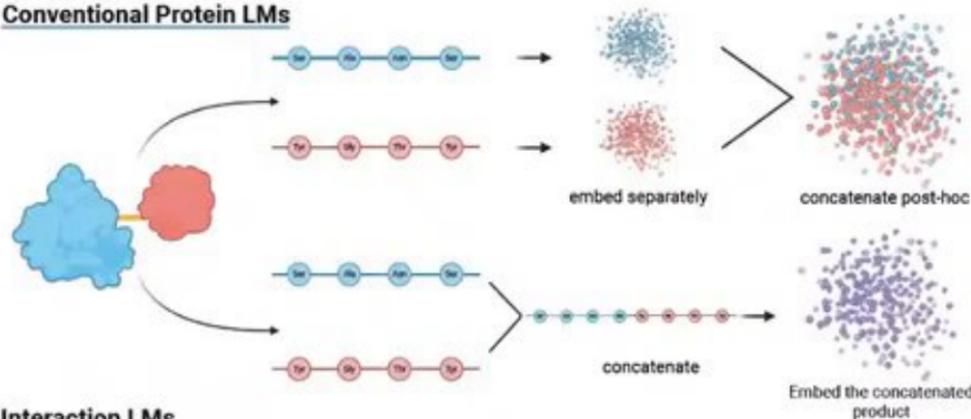
[Explore content](#) [About the journal](#) [Publish with us](#)

[nature](#) > [nature machine intelligence](#) > [analyses](#) > [article](#)

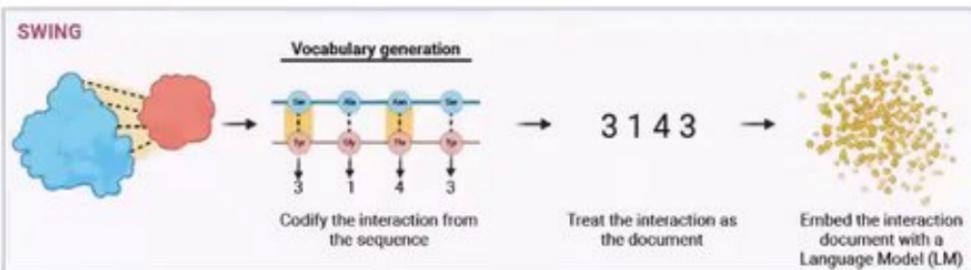
Analysis | [Published: 21 March 2022](#)

Learning functional properties of proteins with language models

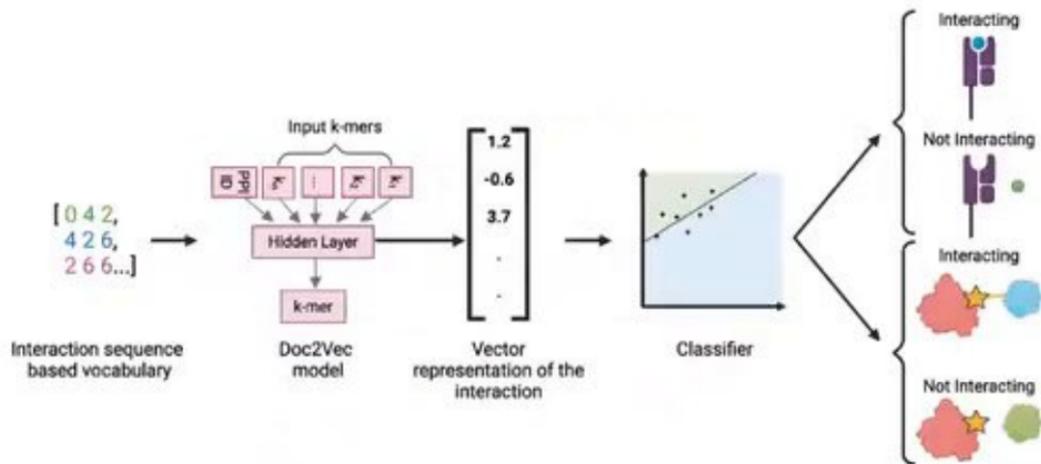
Conventional Protein LMs



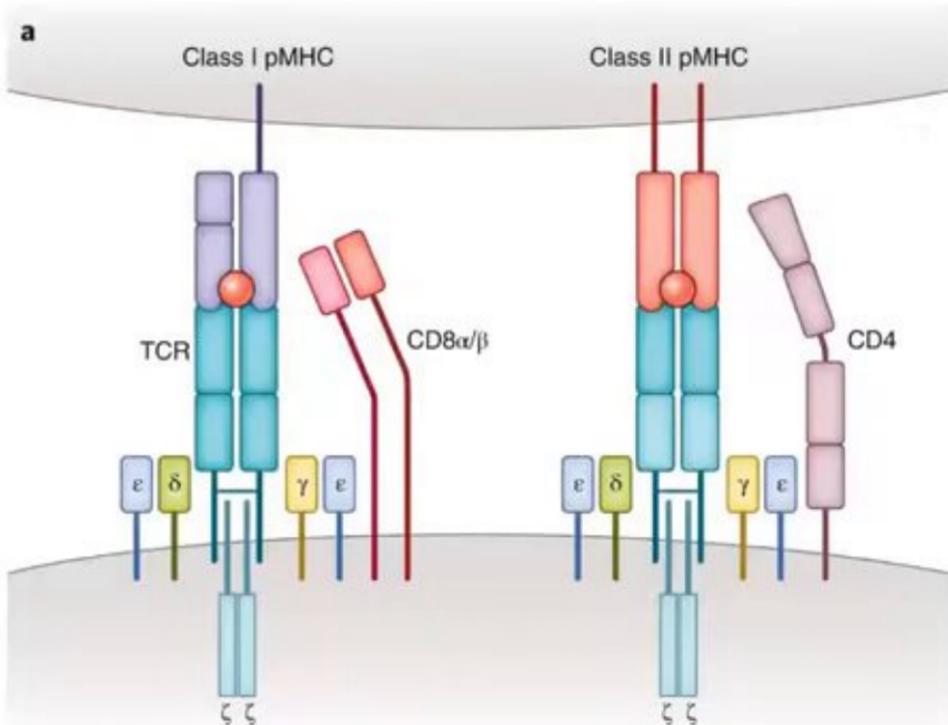
Interaction LMs



SWING: Sliding Window Interaction Grammar – learning the language of peptide and protein interactions



Predicting pMHC binding is a well-studied and important problem



Joglekar and Li *Nature Methods* 2021

[Immunogenetics. 2013 Sep; 65\(9\): 655-665.](#)

Published online 2013 Jun 18. doi: [10.1007/s00251-013-0714-9](#)

MHCcluster, a method for functional clustering of MHC molecules

[Martin Thomsen](#), [Claus Lundegaard](#), [Søren Buus](#), [Ole Lund](#), and [Morten Nielsen](#)[✉]

JOURNAL ARTICLE

NetMHCpan-4.1 and NetMHCIIpan-4.0: improved predictions of MHC antigen presentation by concurrent motif deconvolution and integration of MS MHC eluted ligand data

[Birkir Reynisson](#), [Bruno Alvarez](#), [Sinu Paul](#), [Bjoern Peters](#), [Morten Nielsen](#)

Author Notes

Nucleic Acids Research, Volume 48, Issue W1, 02 July 2020, Pages W449-W454,
<https://doi.org/10.1093/nar/gkaa379>

nature biotechnology

[Explore content](#) [About the journal](#) [Publish with us](#) [Subscribe](#)

[nature](#) > [nature biotechnology](#) > [brief communications](#) > [article](#)

Brief Communication | Published: 14 October 2019

Robust prediction of HLA class II epitopes by deep motif deconvolution of immunopeptidomes

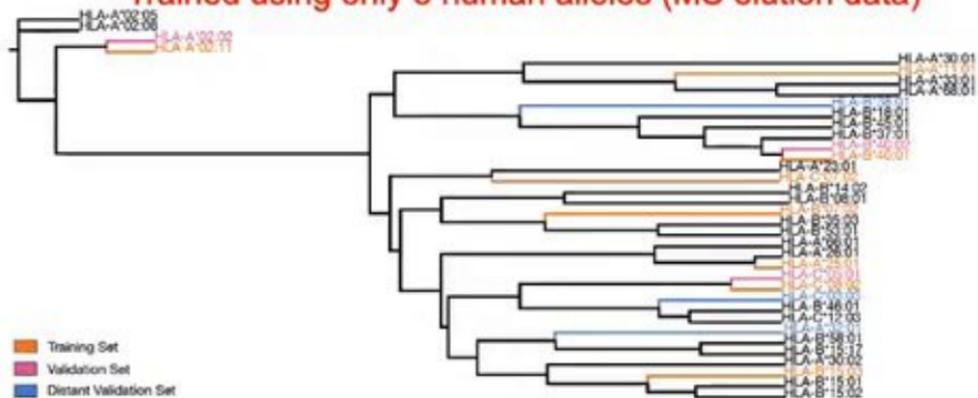
[Julien Bacle](#), [Justine Michaux](#), [Georg Alexander Rocknoer](#), [Marion Arnaud](#), [Sara Bobisse](#), [Chloe Chonq](#), [Philippe Guillaume](#), [Georges Couxos](#), [Alexandre Harari](#), [Camilla Jannus](#), [Michal Bassani-Sternberg](#) & [David Glezer](#)

Nature Biotechnology **37**, 1283-1286 (2019) | [Cite this article](#)

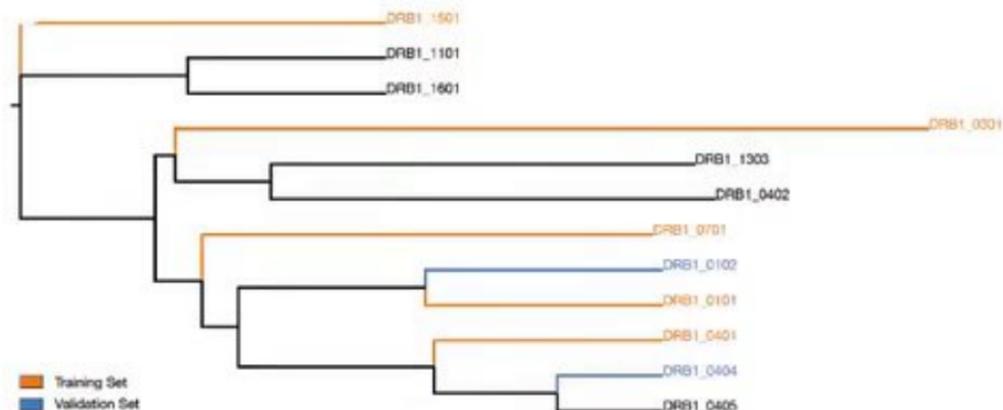
14k Accesses | 156 Citations | 118 Altmetric | [Metrics](#)

Learning the language of Class I pMHC interactions

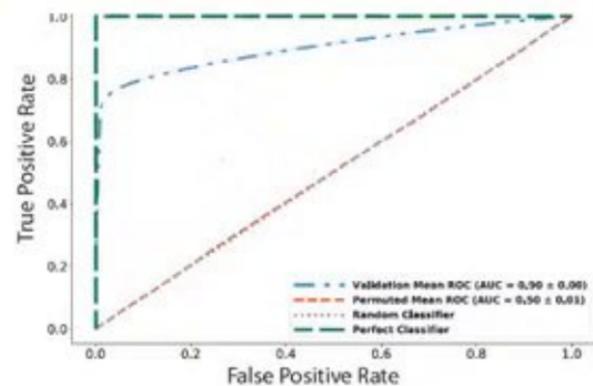
Trained using only 8 human alleles (MS elution data)



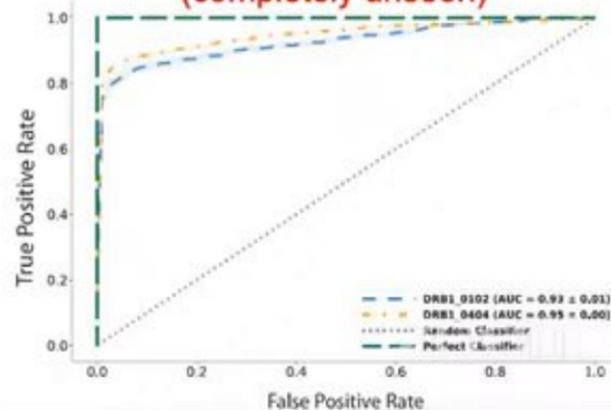
Learning the language of Class II pMHC interactions



K-fold cross-validation on the training set



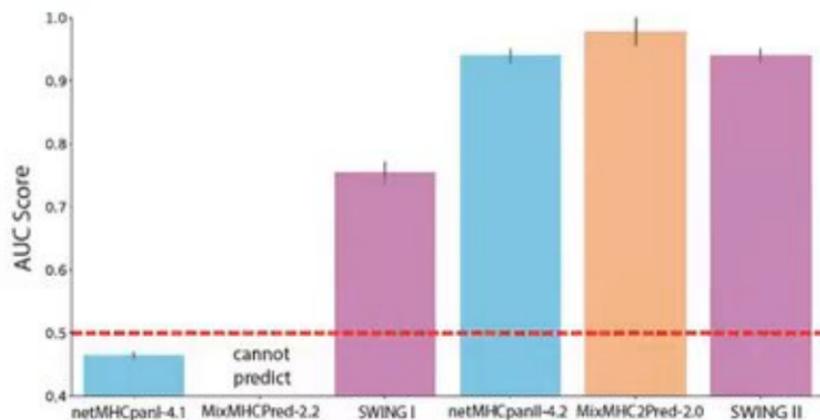
Cross prediction on the validation set (completely unseen)



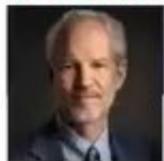
Generalizing the language of all pMHC interactions

Generalizing the language of all pMHC interactions

No existing method can go from Class I -> Class II



SWING can make de-novo predictions of pMHC interactions in a murine model of SLE

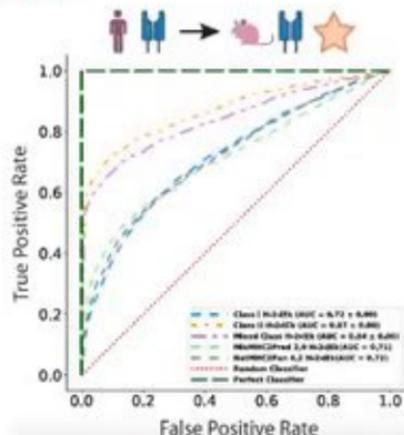
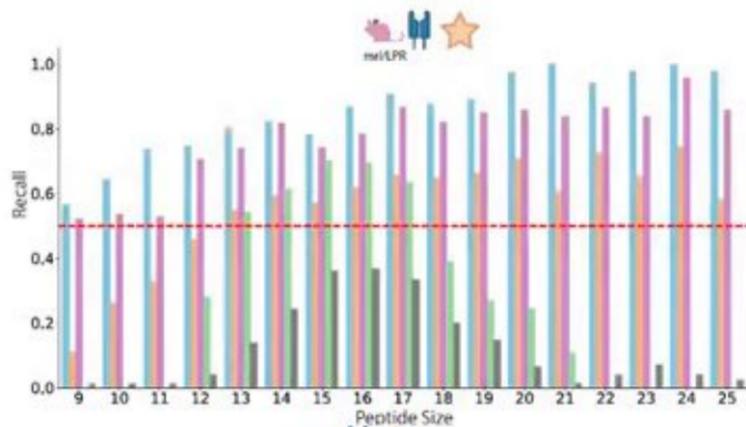


Few-shot
learning –
prediction of de-
novo pMHC
binding across
organisms

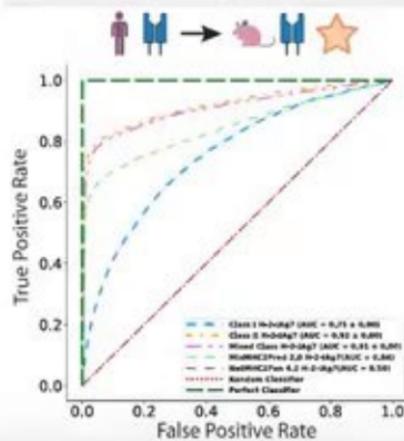
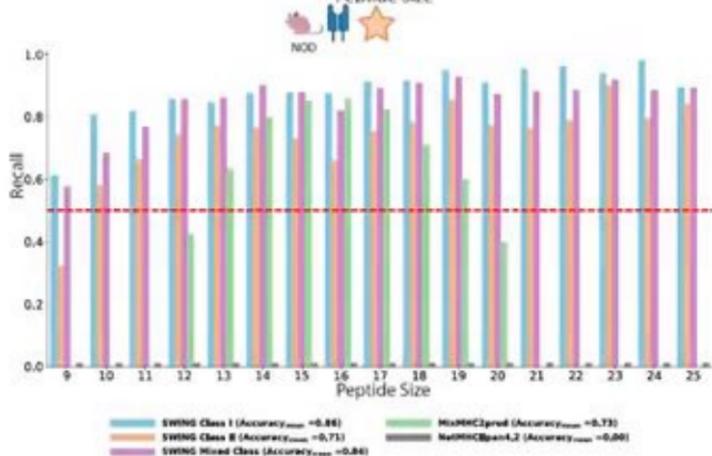
SWING can make de-novo predictions of pMHC interactions in a murine model of SLE



SWING model learnt using human data highly predictive of H-2-IEk binding in a murine model of SLE (Mrl/lpr)



SWING model learnt using human data highly predictive of IAg7 binding in a murine model of T1D (NOD)



Few-shot learning – prediction of de-novo pMHC binding across organisms

■ SWING Class 1 (Accuracy_{Test} = 0.84) ■ MixMHC2pred (Accuracy_{Test} = 0.73)
■ SWING Class 2 (Accuracy_{Test} = 0.71) ■ NetMHCpan4.2 (Accuracy_{Test} = 0.60)
■ SWING Mixed Class (Accuracy_{Test} = 0.84)

False Positive Rate

Modern variant effect predictors learn complex relationships but not the nuanced language of protein interactions

nature

Explore content ▾ About the journal ▾ Publish with us ▾ Subscribe

nature > articles > article

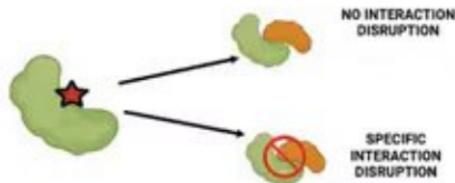
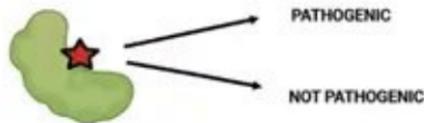
Article | Published: 27 October 2021

Disease variant prediction with deep generative models of evolutionary data

Jonathan Frank, Pascal Notin, Mubinda Dias, Allan Gomes, Joseph K. Min, Kelly Brock, Yarin Gal¹ & Dorothea S. Marks²

Nature 999, 91–95 (2021) | [Cite this article](#)

67k Accesses | 109 Citations | 487 Altmetric | [Metrics](#)



nature genetics

Explore content ▾ About the journal ▾ Publish with us ▾

nature > nature genetics > articles > article

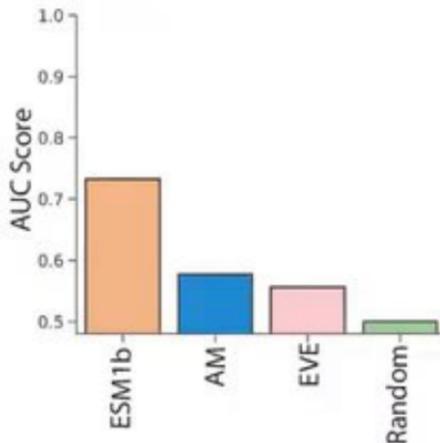
Article | [Open access](#) | Published: 10 August 2023

Genome-wide prediction of disease variant effects with a deep protein language model

Nadav Brandes, Grant Goldman, Charlotte H. Wang, Chun Jimmie Ye¹ & Yaeliz Miralles¹

Nature Genetics 55, 1812–1822 (2023) | [Cite this article](#)

58k Accesses | 28 Citations | 180 Altmetric | [Metrics](#)



Science

Current issue First release papers Archive View ▾

Search manuscript

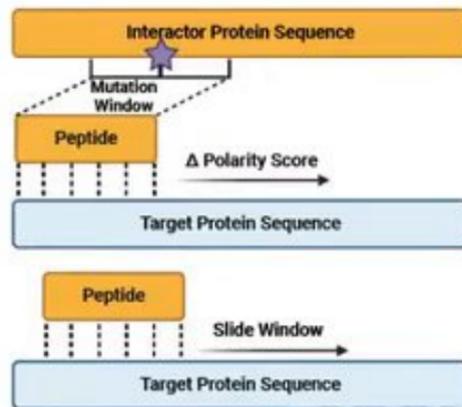
Science > Science articles > Science articles

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

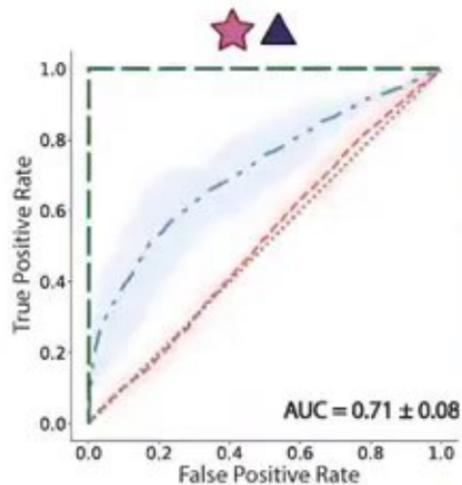
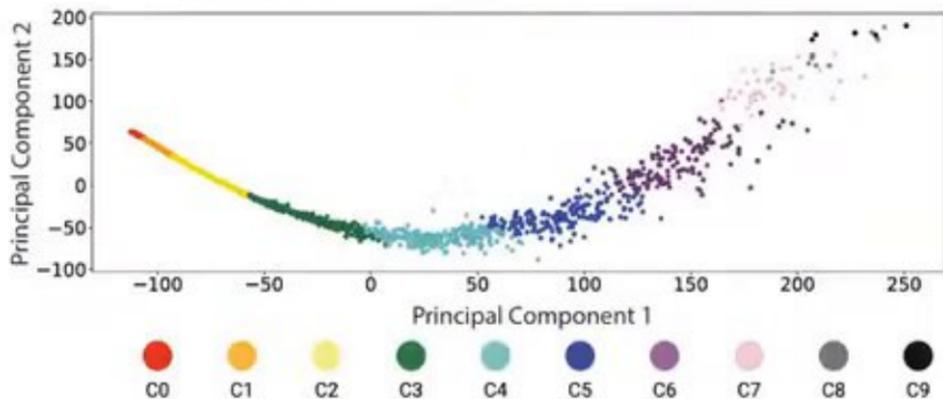
Accurate proteome-wide missense variant effect prediction with AlphaMissense

Jonathan Frank, Pascal Notin, Mubinda Dias, Allan Gomes, Joseph K. Min, Kelly Brock, Yarin Gal¹ & Dorothea S. Marks²

Science 378, 108–113 (2022) | [Cite this article](#)

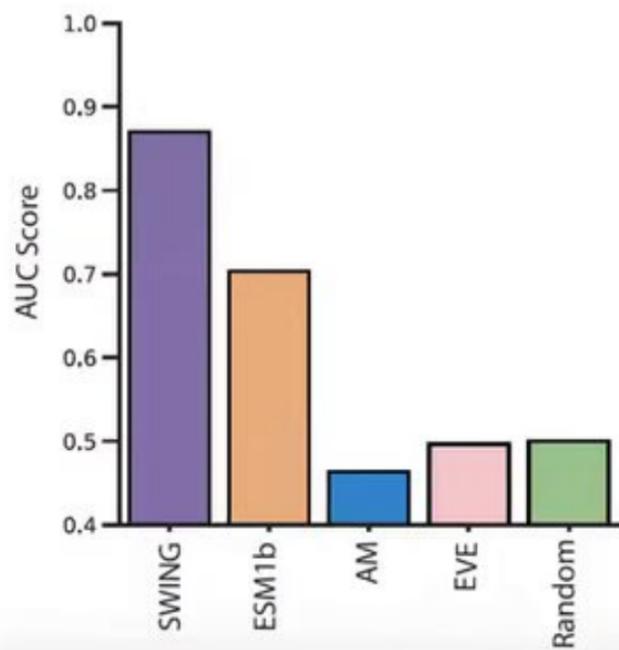


SWING learns relationships well beyond sequence similarity

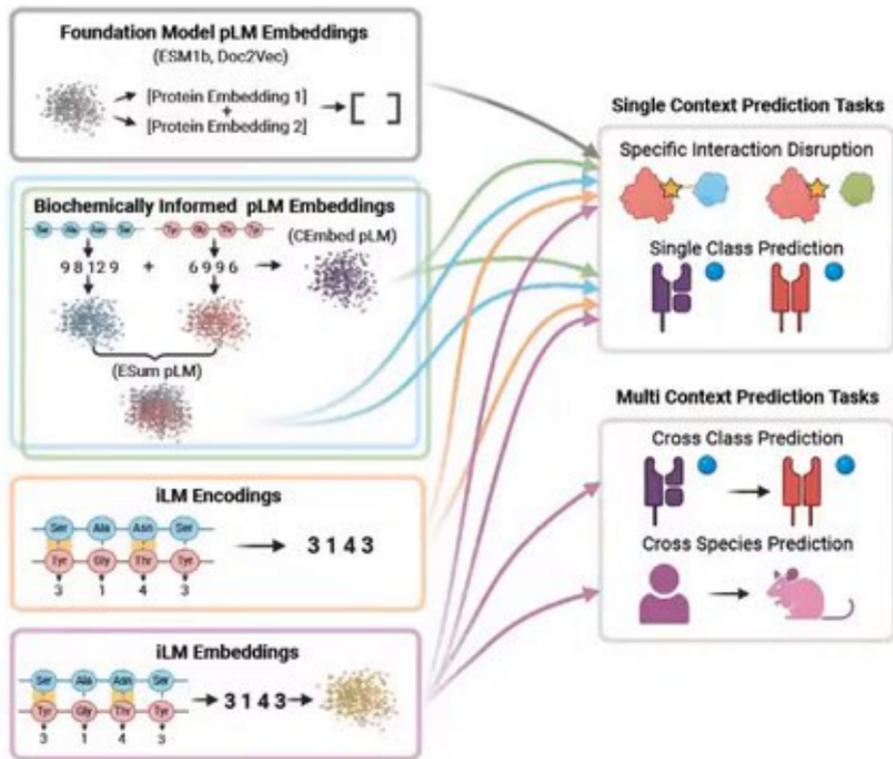


	Validation Clusters	Dropped Clusters	AUC Score
Fold 1	C0, C1	C2, C3	0.67
Fold 2	C2, C3	C1, C4	0.70
Fold 3	C4, C5	C3, C6	0.72
Fold 4	C6, C7	C5, C8	0.78

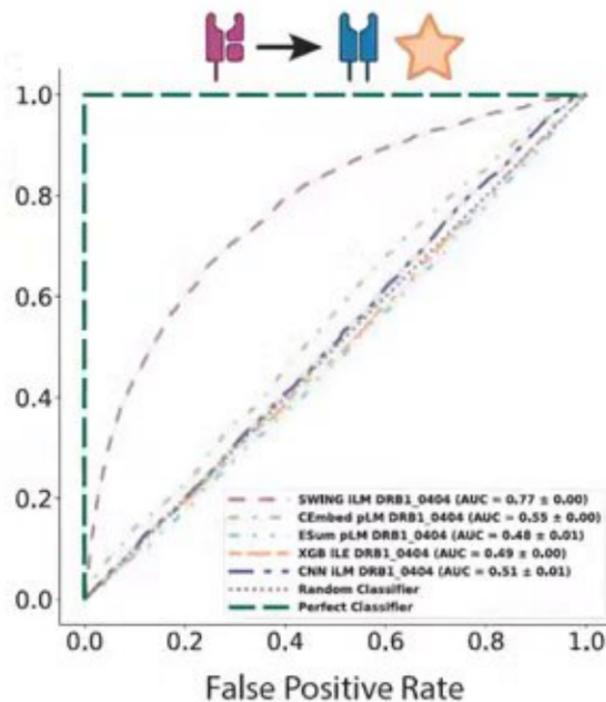
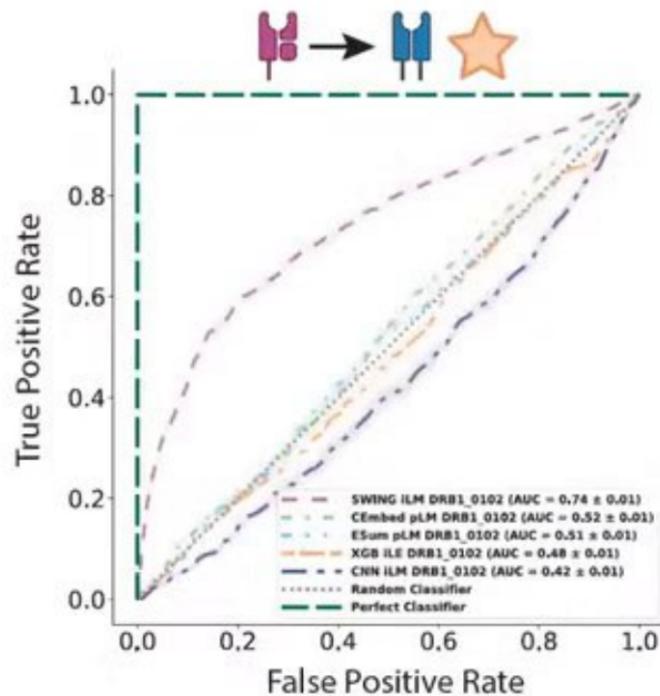
SWING outperforms state-of-the-art VEP methods and uncovers molecular phenotypes missed by these approaches



SWING performs better than/as well as related architectures for single context prediction tasks



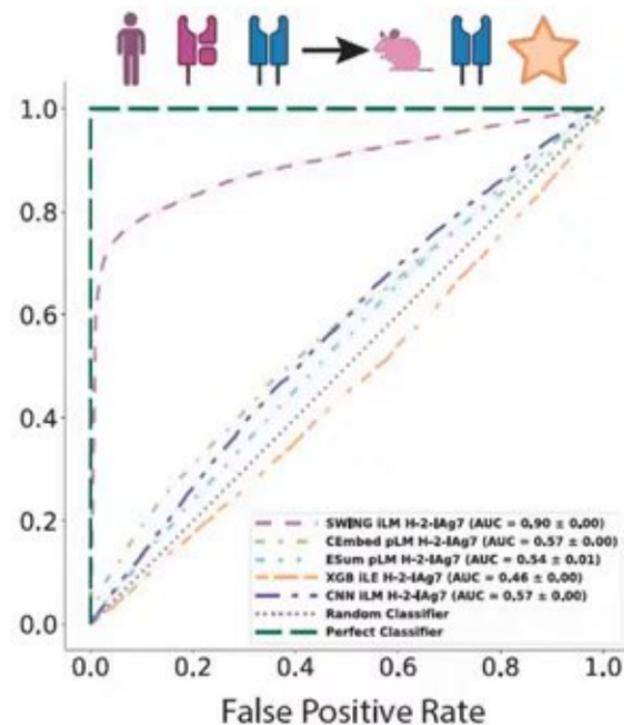
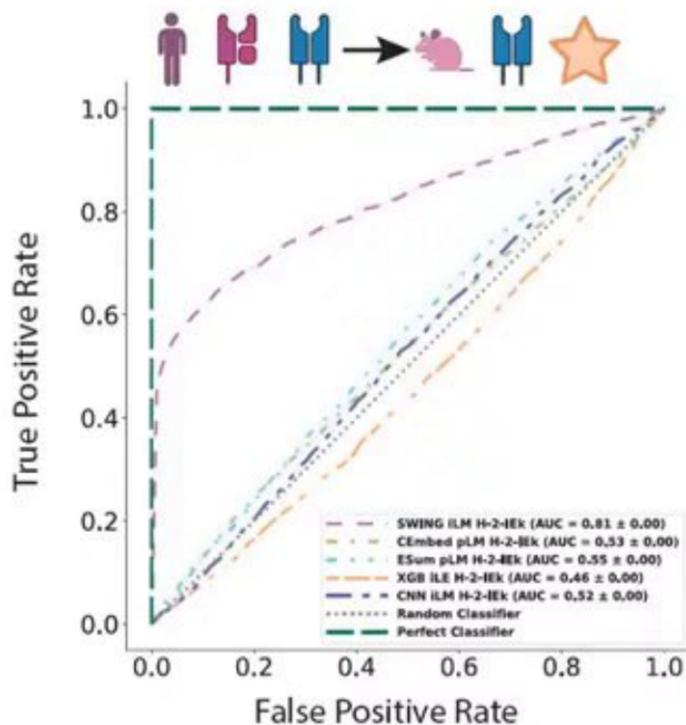
SWING significantly outperforms related architectures for multi-context prediction tasks and unseen alleles - I



 Cross Prediction  Standard Cross Validation  MHC-I  MHC-II



SWING significantly outperforms related architectures for multi-context prediction tasks and unseen alleles - II



Conclusions

	SWING	NetMHCpan	MixMHCpred	AF	ESM1b	EVE
Embeds the interaction	✓	✗	✗	✗	✗	✗
Works across peptide-peptide interaction contexts	✓	✗	✗	✗	✗	✗
Predicts Class I pMHC interactions	✓	✓	✓			
Predicts Class II pMHC interactions	✓	✓	✓			
Predicts Class I and Class II interactions with the same model	✓	✗	✗			
Predicts across MHC Classes	✓	✗	✗			
Generalizes across species while trained on one species	✓	✗	✗			
Predicts variant effects on the phenotype	✓			✓	✗	✓
Predicts edgotype effects on the phenotype	✓			✗	✓	✗



Systems Immunology Faculty Position at Pitt



Interested candidates should send a CV and statement of research interests to **Rachel Gottschalk** (rachel.gottschalk@pitt.edu) and **Jishnu Das** (jishnu@pitt.edu).



The Center for Systems Immunology (CSI) and the Department of Immunology at the University of Pittsburgh are seeking to appoint a tenure-track Assistant Professor. Candidates with strong research training and accomplishments in systems and computational immunology and/or synthetic biology and immune engineering will be considered. The position offers a competitive start-up package and dedicated space within a new, highly interactive research building on the University campus that houses faculty in the CSI, the Departments of Immunology and the Hillman Cancer Center.

The candidate will be expected to develop a leading externally funded research group; mentor graduate students including those in the CMU-Pitt Computational Biology program and contribute to graduate teaching within such programs affiliated with the School of Medicine.

The University of Pittsburgh has excellence in Immunology and Computational Biology made possible by a highly collaborative faculty, a diverse and strong trainee pool, and world class core facilities and infrastructure.

Interested candidates should send a CV and statement of research interests to Rachel Gottschalk (rachel.gottschalk@pitt.edu) and Jishnu Das (jishnu@pitt.edu).

Our Funders

4 R01/R01-equivalents grants as PI/MPI

NIAID New Innovators DP2 Award



DP2AI164325 (Role: PI, Active)

National Institute of Allergy and Infectious Diseases

NHGRI IGVF Consortium U01



U01HG012041 (Role: MPI, Active)

National Human Genome Research Institute

NIAID Flu Systems Vaccinology R01



R01AI170108 (Role: MPI, Active)

National Institute of Allergy and Infectious Diseases

NIAID Flu Systems Vaccinology U01



U01AI179514 (Role: MPI, Active)

National Institute of Allergy and Infectious Diseases

jishnu@pitt.edu



www.jishnulab.org

Twitter: @jishnu1729



Rainin Innovator Award



1. DoD Idea Development Award
2. DoD SL20118

Other Grants



National Institute of Allergy and Infectious Diseases

1. NIAID R01AI167711
2. NIAID U19AI181984



National Eye Institute

NEI U01EY034711



NIAMS P50 CORT
P50AR080612-01



National Heart, Lung, and Blood Institute

NIAID R35HL166219



National Institute of Mental Health

NIMH R01MH136952



NCI R01CA2774731

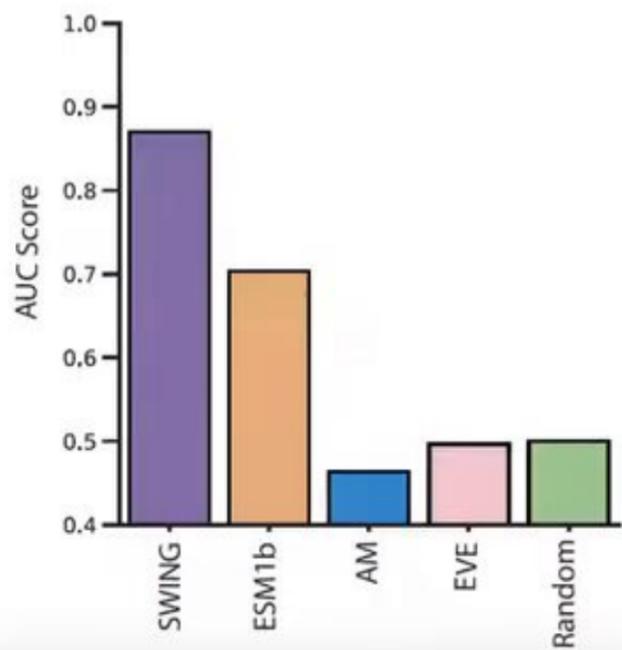


CIHR IRSC

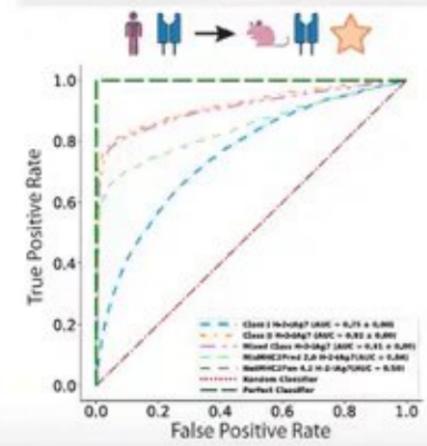
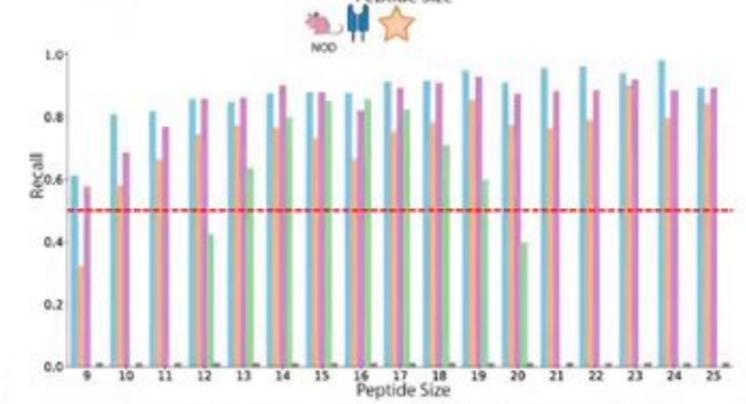
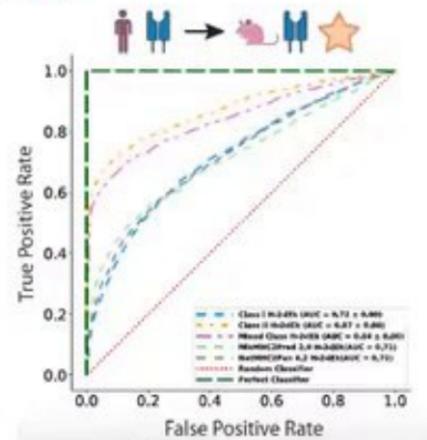
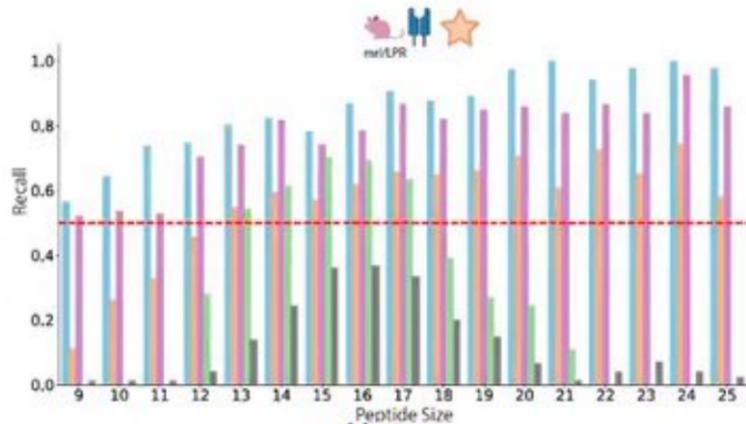
Canadian Institutes of Health Research / Instituts de recherche en santé du Canada

CIHR
DCO150GP

SWING outperforms state-of-the-art VEP methods and uncovers molecular phenotypes missed by these approaches



SWING can make de-novo predictions of pMHC interactions in a murine model of SLE



SWING model learnt using human data highly predictive of H-2-IEk binding in a murine model of SLE (Mrl/lpr)

SWING model learnt using human data highly predictive of IAg7 binding in a murine model of T1D (NOD)

Few-shot learning – prediction of de-novo pMHC binding across organisms

■ SWING Class 1 (Accuracy_{Test} = 0.84) ■ MiMHC2pred (Accuracy_{Test} = 0.73)
■ SWING Class 2 (Accuracy_{Test} = 0.71) ■ NetMHCpan4.2 (Accuracy_{Test} = 0.60)
■ SWING Mixed Class (Accuracy_{Test} = 0.84)

NYGC Events



@



Thank you to our sponsors

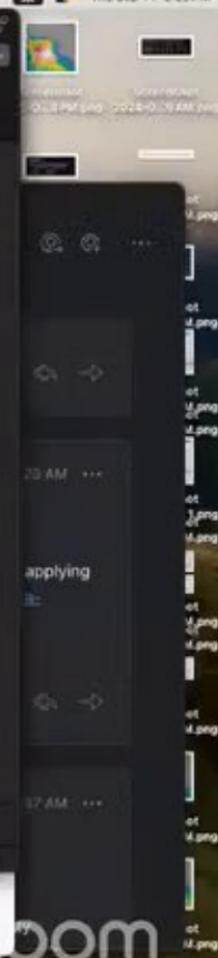


And helpers!

Sarah Curtiss , Kristen Weatherley
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 3.50pm ET

mlcb.org for schedule





@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley

Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 3.50pm ET

mlcb.org for schedule



@



Thank you to our sponsors

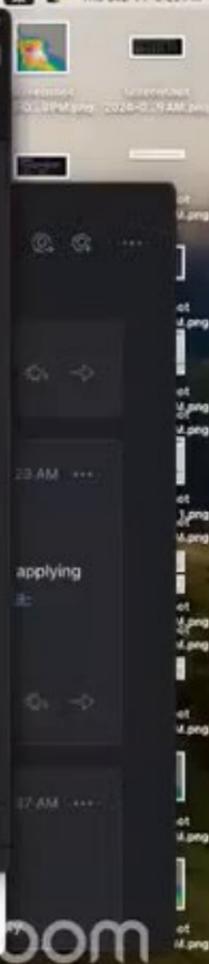


And helpers!

Sarah Curtiss , Kristen Weatherley
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 3.50pm ET

mlcb.org for schedule





Home Insert Draw Design Trans

MCBRLab logo

Click to add notes

Slide 1 of 1 English (United States) Accessibility: Investigate

MLCB oral slides

Name	Date Modified	Size	Kind
29_SWINO_MLCB2025_ID Jishnu Das.pptx	Today at 12:28 PM	21.2 MB	PowerP... (pptx)
41_Linder Johannes Linder.pptx	Sep 9, 2025 at 10:21 PM	3.8 MB	PowerP... (pptx)
42_Fanjiang Clara Fanjiang.key	Yesterday at 12:41 PM	16.1 MB	Keynote
83_Shearer Courtney Shearer.pptx	Sep 9, 2025 at 8:51 PM	2.8 MB	PowerP... (pptx)
102_Sethi Palash Sethi.pptx	Today at 1:13 PM	20.8 MB	PowerP... (pptx)
114_Rocha Joao Felipe Rocha.pptx	Today at 9:06 AM	12.7 MB	PowerP... (pptx)
134_shaw Peter Shaw.pdf	Yesterday at 10:12 PM	1.2 MB	PDF Document
154_BLASSSEL Luc Blasssel.pdf	Yesterday at 11:56 PM	3.6 MB	PDF Document
2025 MLCB Keynote Jacob Schneber.pptx	Today at 8:45 AM	28.3 MB	PowerP... (pptx)
2025-09-11 MLCB RocketSHIP Samuel Sledzieski.pptx	Today at 12:58 PM	179.8 MB	PowerP... (pptx)
bee_MLCB_2025 Barbara.pptx	Yesterday at 1:17 PM	477.8 MB	PowerP... (pptx)
Michael_Broddiaconi_Look_Mom Michael Broddiaconi.pdf	Today at 1:05 PM	28.7 MB	PDF Document
MLCB (75 min) Alan Amin.key	Yesterday at 11:58 AM	4.5 MB	Keynote
MLCB_2025_Battle Alexis Battle.pdf	Yesterday at 9:46 AM	6.2 MB	PDF Document
MLCB_2025_Battle Alexis Battle.pptx	Yesterday at 9:08 AM	37 MB	PowerP... (pptx)
MLCB_2025_v2 Danielle Stevens.pptx	Today at 11:41 AM	82.3 MB	PowerP... (pptx)
Perturbation_Benchmark_Hasana) Eshana Hasana).pdf	Today at 9:33 AM	1.8 MB	PDF Document

Restarting at 3.50pm ET

mlcb.org for schedule

Timeline of events:

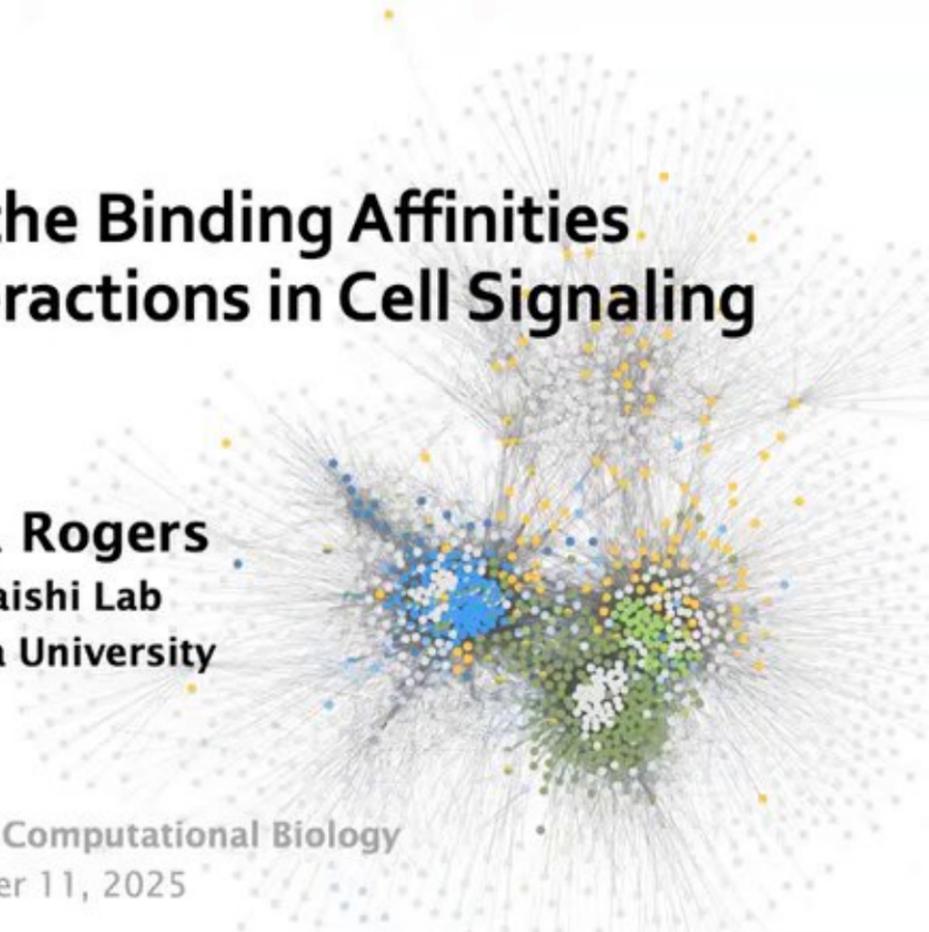
- 23 AM ...
- applying
- 37 AM ...

NYGC Events

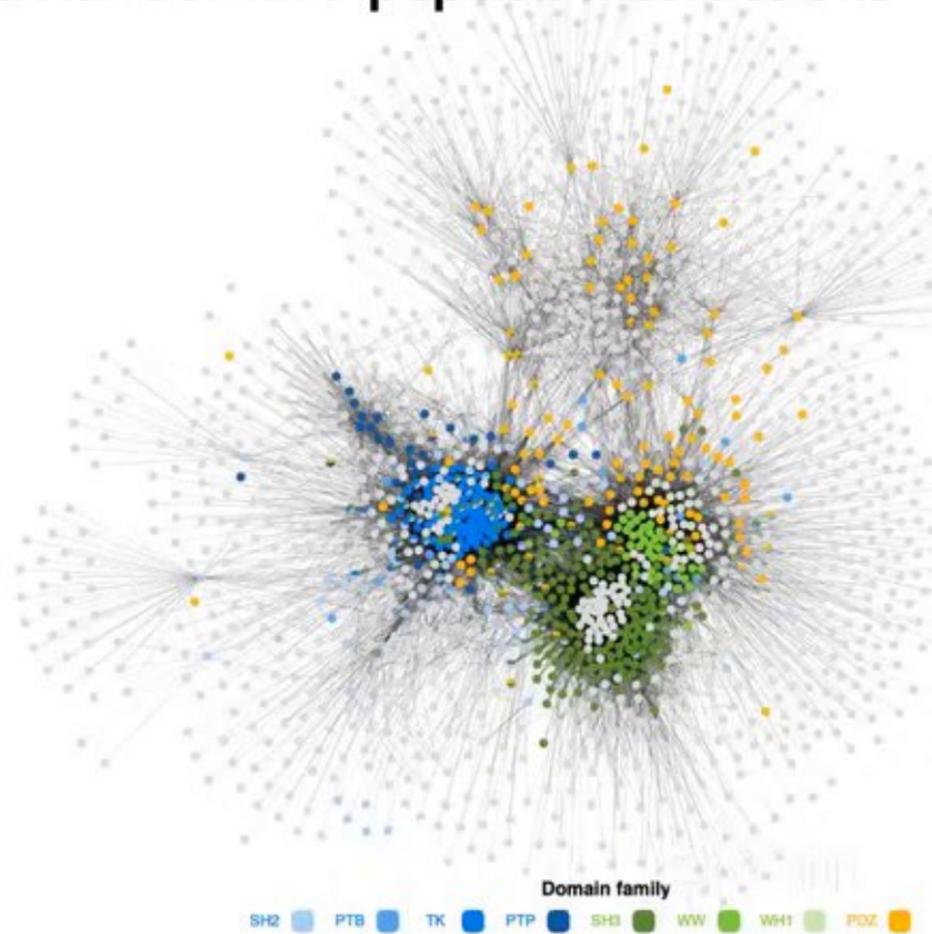
Machine Learning the Binding Affinities of Protein–Peptide Interactions in Cell Signaling

Julia R Rogers
 AlQuraishi Lab
 Columbia University

Machine Learning in Computational Biology
 September 11, 2025

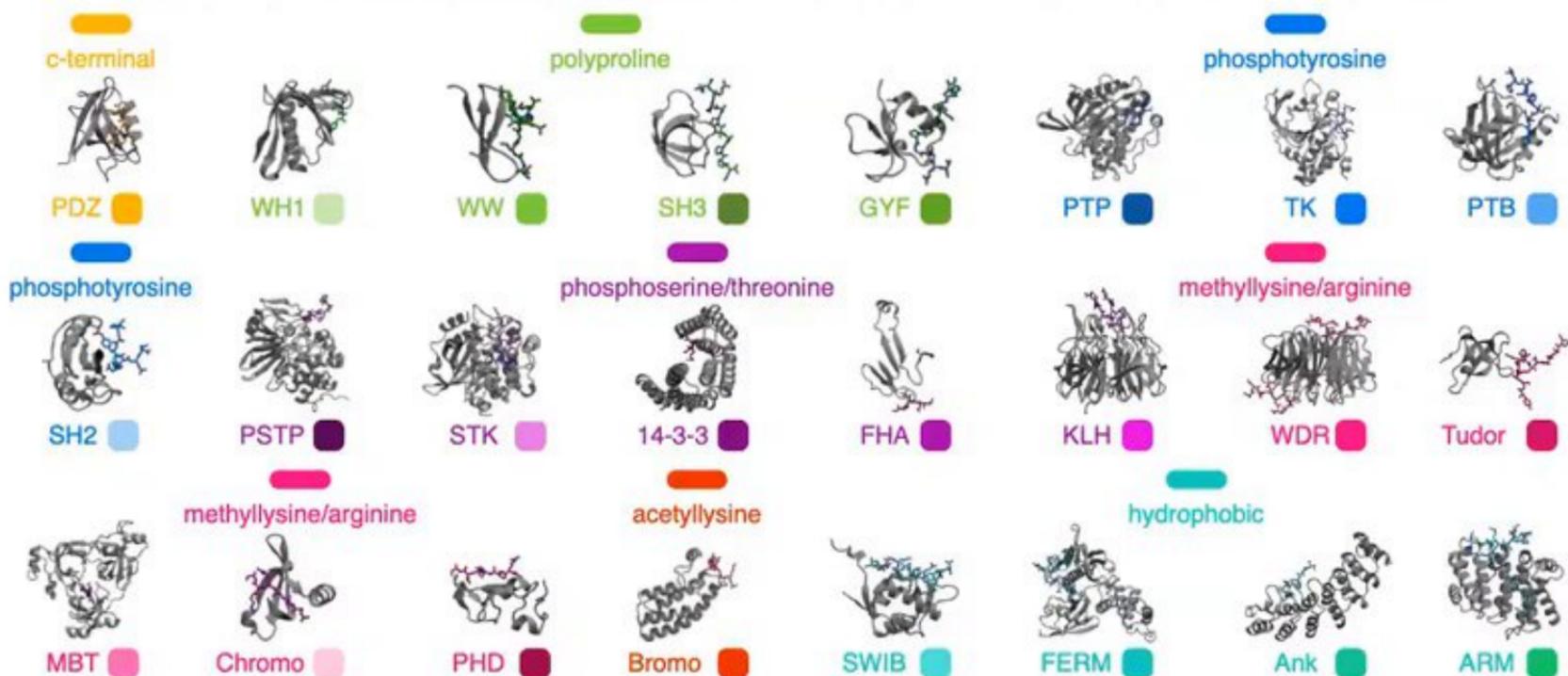


Cellular signals processed by modular domain-peptide interactions

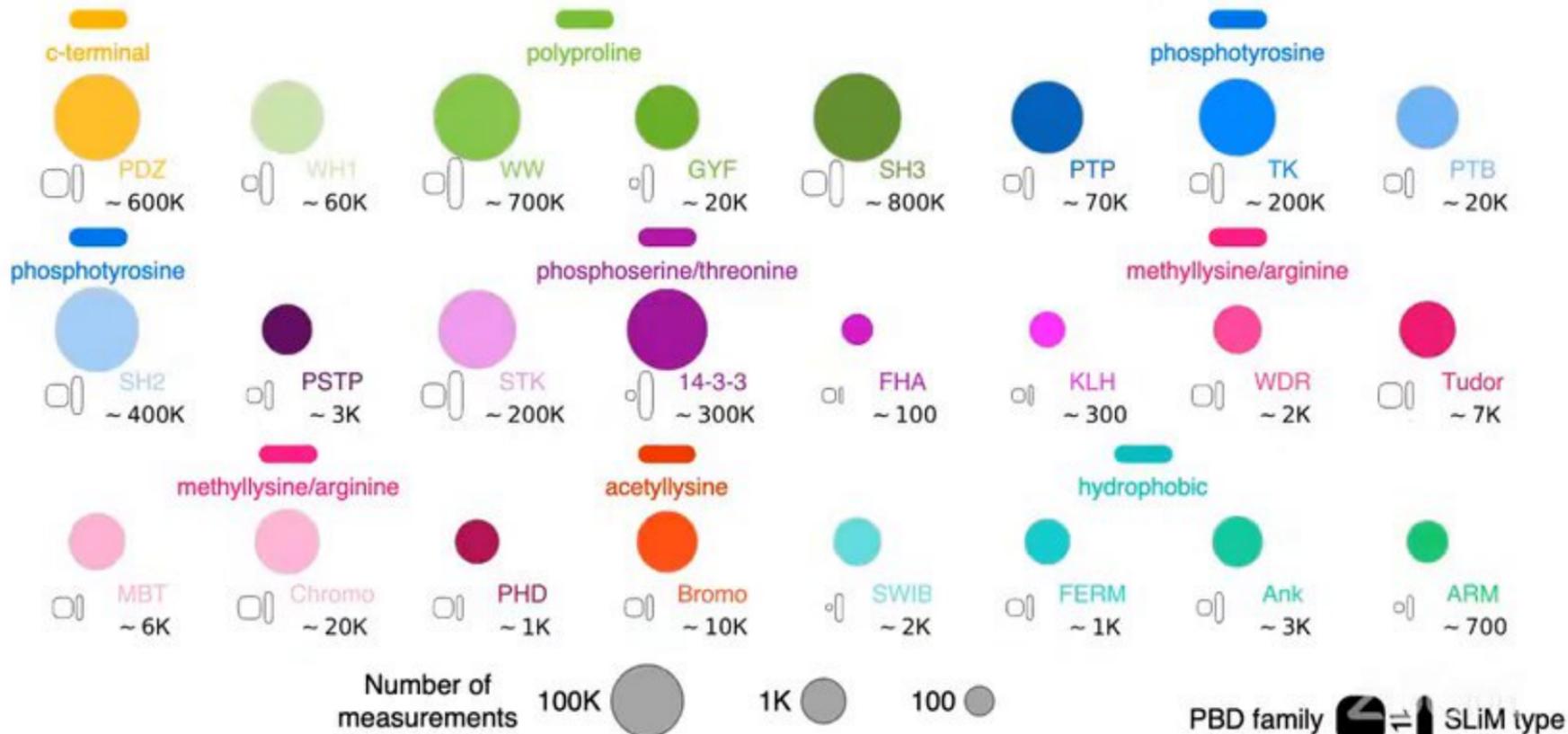


Diversity of modular domain-peptide interactions

(out of ~100-200 human PBD families and ~7 binding chemistry)

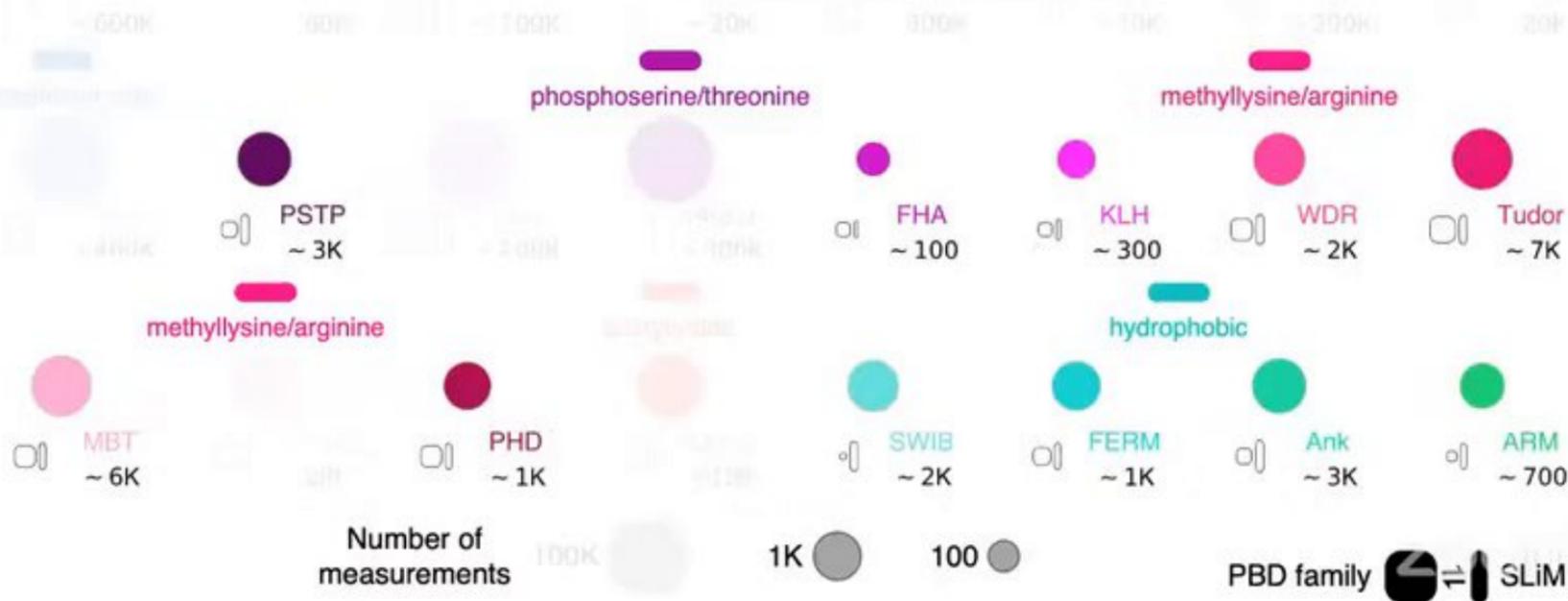


Diversity of modular domain-peptide interactions and disparate experimental coverage

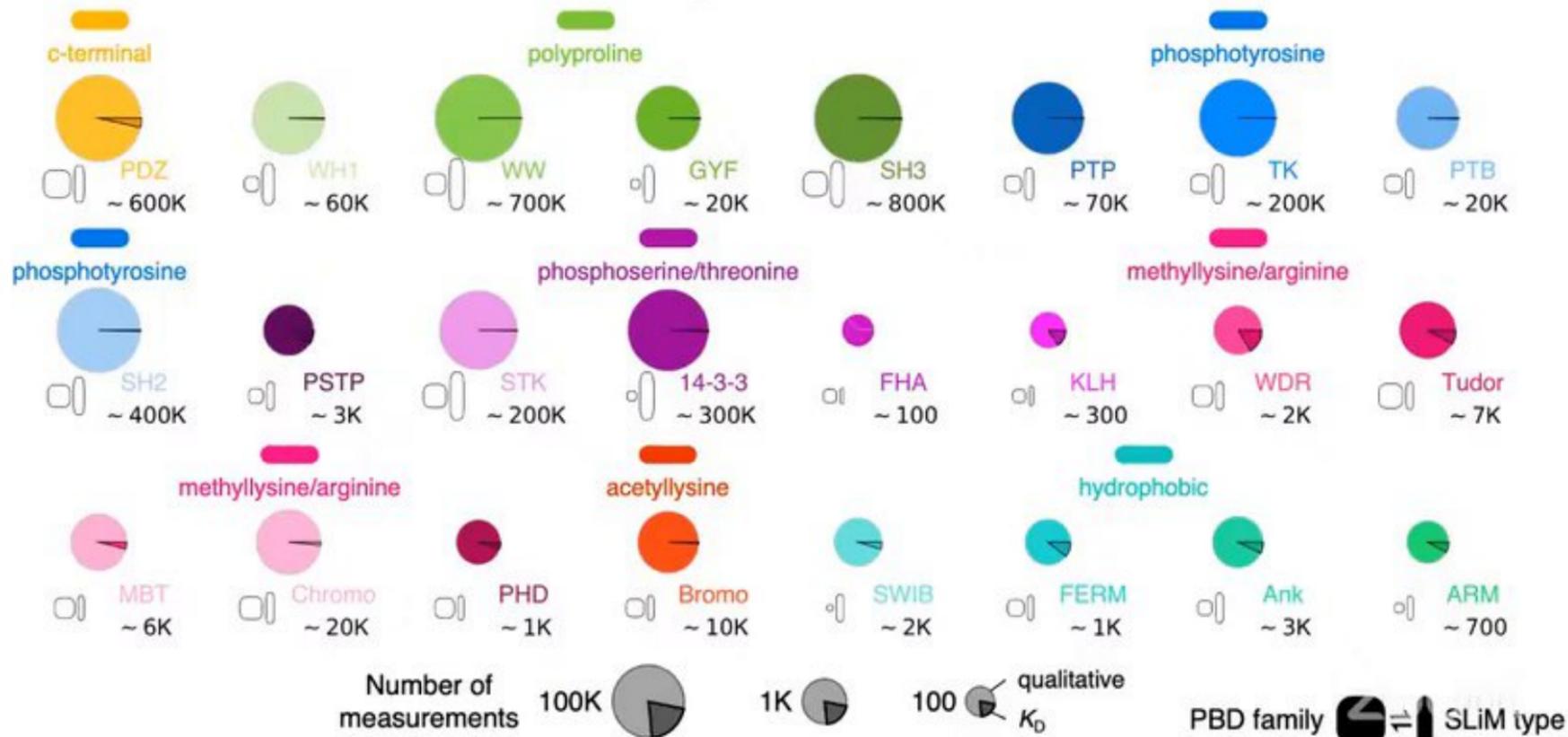


Diversity of modular domain-peptide interactions and disparate experimental coverage

Poorly characterized
~100–1,000 measured interactions

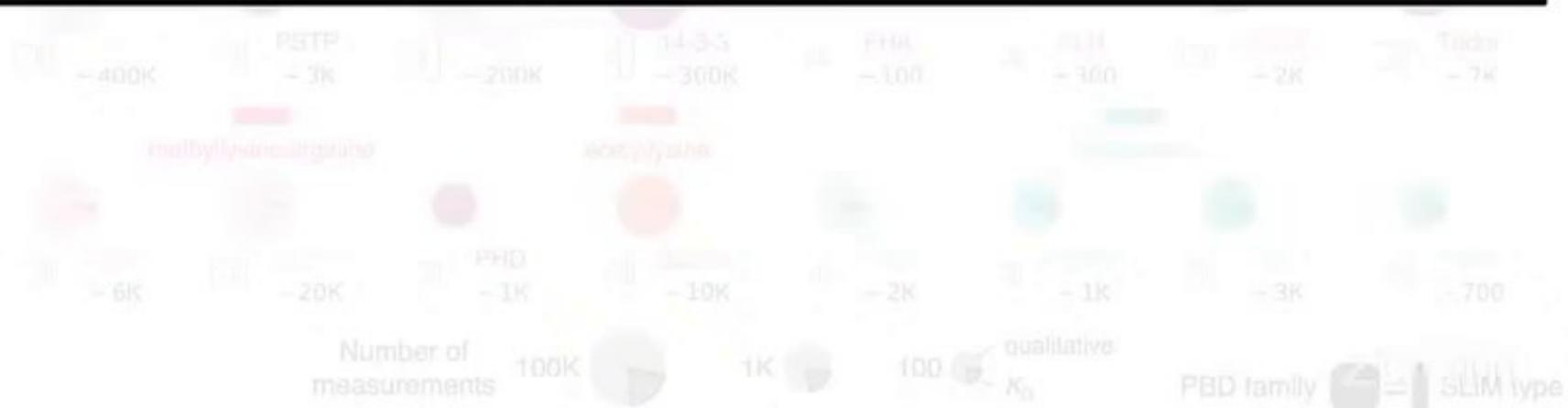


Diversity of modular domain-peptide interactions and limited knowledge of quantitative affinities



Diversity of modular domain-peptide interactions and disparate experimental coverage

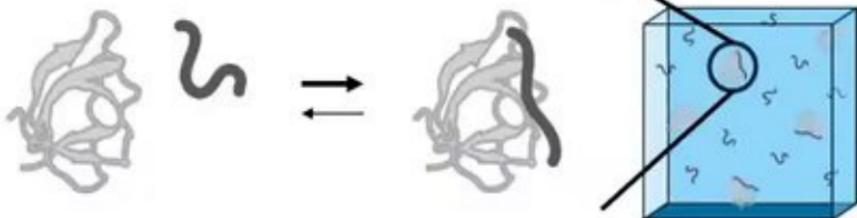
Development of a machine learning approach for universally predicting domain-peptide interaction affinity



Statistical mechanical description of domain-peptide interactions

Unbound

Bound



Machine learning a quantitative statistical mechanical model for ΔF to predict interaction affinity directly from amino acid sequence

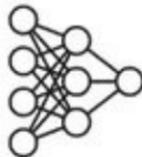
Peptide sequence:

GTTPPPYTV



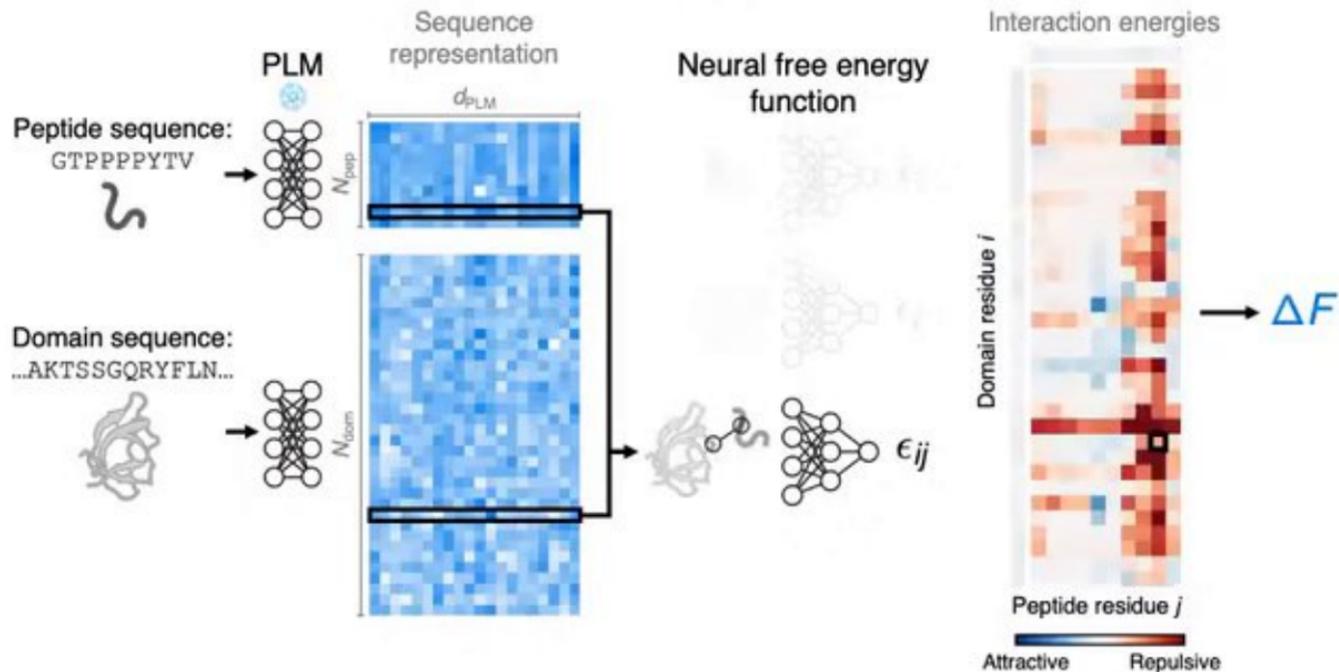
Domain sequence:

...AKTSSGQRYFLN...



ΔF

Machine learning a quantitative statistical mechanical model for ΔF parameterized by representations shared across domain families



$$\Delta F = \sum_i^{N_{\text{dom}}} \epsilon_i^{(\text{dom})} + \sum_j^{N_{\text{pep}}} \epsilon_j^{(\text{pep})} + \sum_{ij}^{N_{\text{dom}} \times N_{\text{pep}}} \epsilon_{ij}^{(\text{int})} + \Delta S$$

Quantitative statistical mechanical model (QSM)

Machine learning a quantitative statistical mechanical model for ΔF

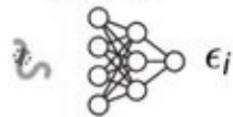
How to infer a free energy function from data?

Peptide sequence:

GTPPPPYTV



Neural free energy function



Domain sequence:

...AKTSSGQRYFLN...



Interaction energies



ΔF

Attractive Repulsive

$$\Delta F = \sum_i^{N_{\text{dom}}} \epsilon_i^{(\text{dom})} + \sum_j^{N_{\text{pep}}} \epsilon_j^{(\text{pep})} + \sum_{ij}^{N_{\text{dom}} \times N_{\text{pep}}} \epsilon_{ij}^{(\text{int})} + \Delta S$$

Quantitative statistical mechanical model (QSM)

Machine learning a quantitative statistical mechanical model for ΔF using abundant qualitative interaction data

Peptide sequence
STDPHPTV



Domain sequence:
AKTSRGGRYKEL



Neural free energy function



Interaction matrix



$\rightarrow \Delta F$

Attractive Repulsive

Qualitative data



Binding probability at peptide concentration c

$$P_{\text{bind}} = \frac{(c/c_0)e^{-\beta\Delta F}}{1 + (c/c_0)e^{-\beta\Delta F}}$$

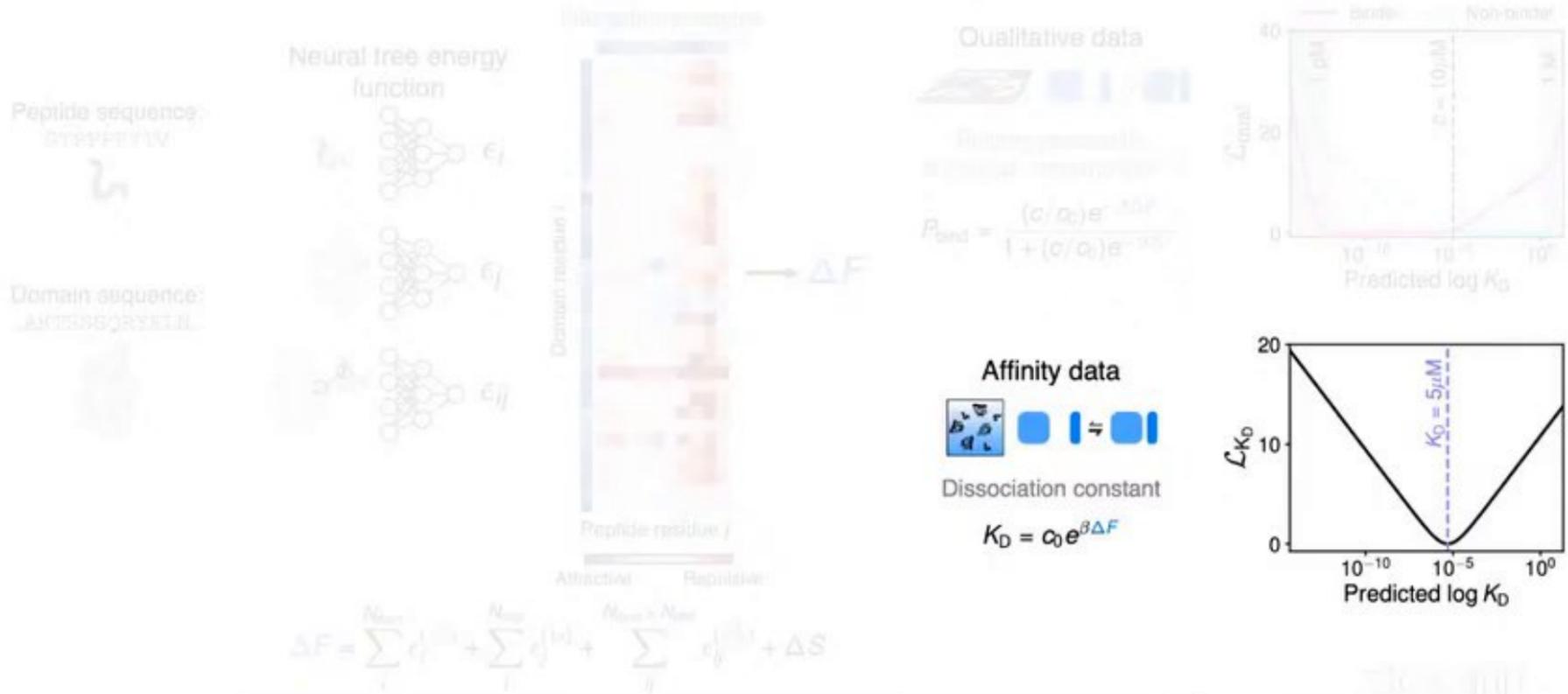
$$\Delta F = \sum_i \epsilon_i^{(i)} + \sum_j \epsilon_j^{(j)} + \sum_{i,j} \epsilon_{ij}^{(ij)} + \Delta S$$

Quantitative statistical mechanical model (QSM)

Predicted observable



Machine learning a quantitative statistical mechanical model for ΔF using quantitative affinity data



Quantitative statistical mechanical model (QSM)

Predicted observable

Optimization objective

Machine learning a quantitative statistical mechanical model for ΔF using a multitask learning objective to leverage all available interaction data

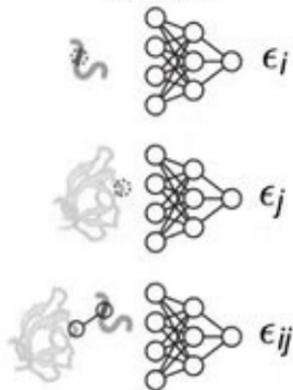
Peptide sequence:
GTPPPPYTV



Domain sequence:
...AKTSSGQRYFLN...



Neural free energy function



Interaction energies



ΔF

Attractive Repulsive

$$\Delta F = \sum_i^{N_{\text{dom}}} \epsilon_i^{(L)} + \sum_j^{N_{\text{pep}}} \epsilon_j^{(L)} + \sum_{ij}^{N_{\text{dom}} \times N_{\text{pep}}} \epsilon_{ij}^{(L)} + \Delta S$$

Qualitative data



Binding probability at peptide concentration c

$$P_{\text{bind}} = \frac{(c/c_0)e^{-\beta\Delta F}}{1 + (c/c_0)e^{-\beta\Delta F}}$$

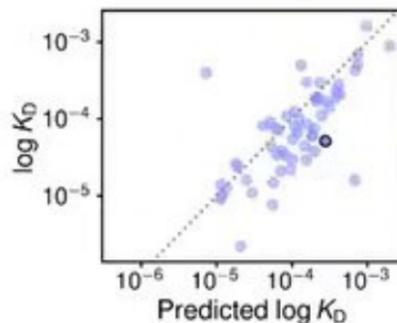
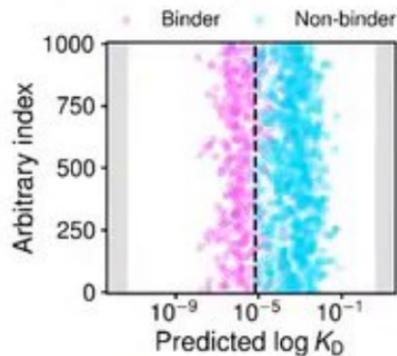
Affinity data



Dissociation constant

$$K_D = c_0 e^{\beta\Delta F}$$

$$\mathcal{L} = \mathcal{L}_{\text{qual}}(\text{blue square vs blue bar}) + \alpha \mathcal{L}_{K_D}(\text{blue square vs blue bar with equilibrium symbol})$$



Quantitative statistical mechanical model (QSM)

Predicted observable

Optimization objective

Machine learning a quantitative statistical mechanical model for ΔF

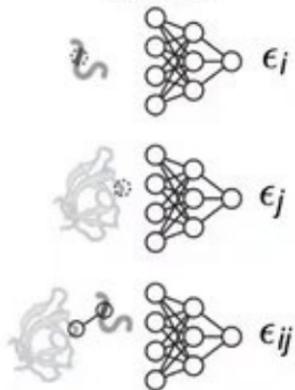
Peptide sequence:
GTPPPPYTV



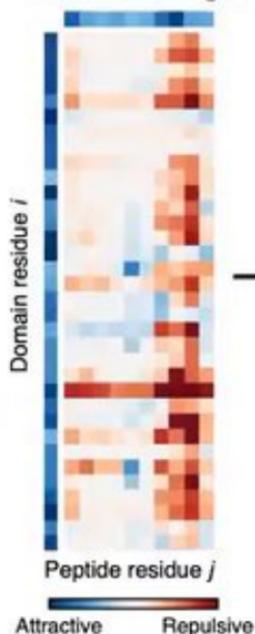
Domain sequence:
...AKTSSGQRYFLN...



Neural free energy function



Interaction energies



ΔF

$$\Delta F = \sum_{i,j} \epsilon_{ij} + \sum_{i,j} \epsilon_j + \sum_{i,j} \epsilon_j + \Delta S$$

Qualitative data



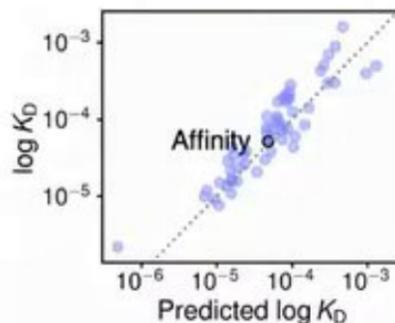
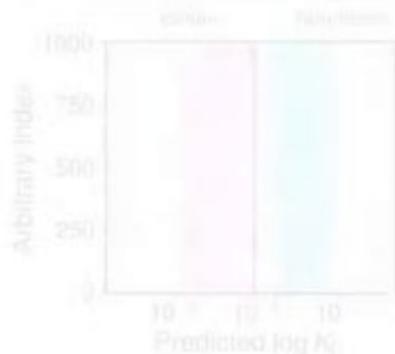
Empirical model + regularization

$$P_{best} = \frac{(c/c_0)e^{-\beta \Delta F}}{1 + (c/c_0)e^{-\beta \Delta F}}$$

Affinity data



$$K_D = c_0 e^{-\beta \Delta F}$$



$$\mathcal{L} = \mathcal{L}_{qual}(\dots) + \alpha \mathcal{L}_{K_D}(\dots)$$

Quantitative statistical mechanical model (QSM)

Predicted observable

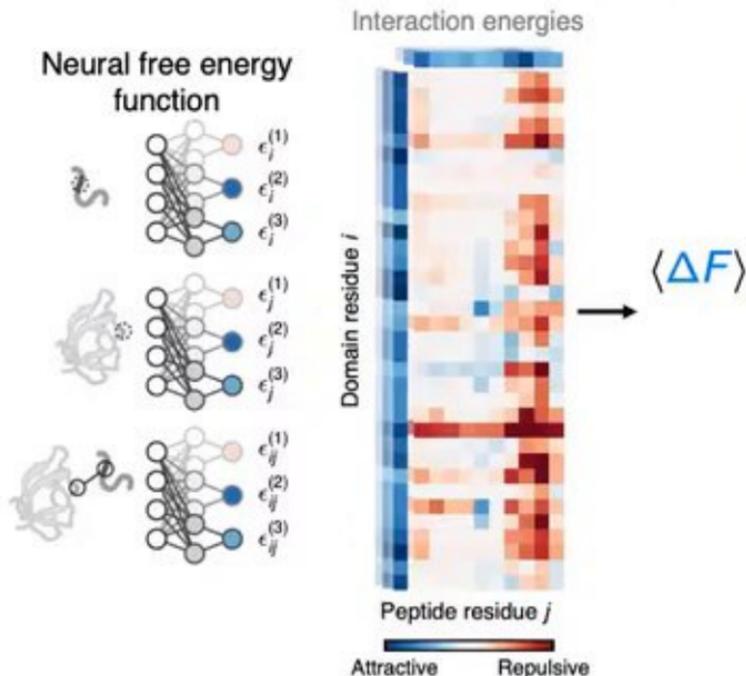
Optimization objective

Machine learning a quantitative statistical mechanical model for ΔF that also estimates predictive error

Peptide sequence:
GTPPPPYTV



Domain sequence:
...AKTSSGQRYFLN...



Ensemble neural free energy function

\Rightarrow Estimates affinity + absolute error

Scientific application

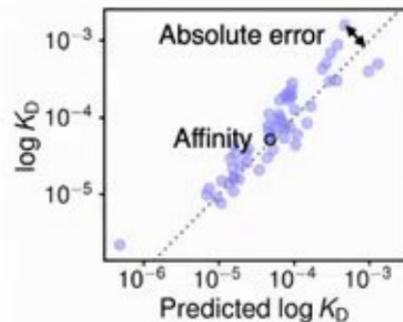
Ensemble mean := Predicted binding free energy

Affinity data



Yeast two-hybrid interaction

$$K_D = \alpha_0 e^{-\beta \Delta F}$$



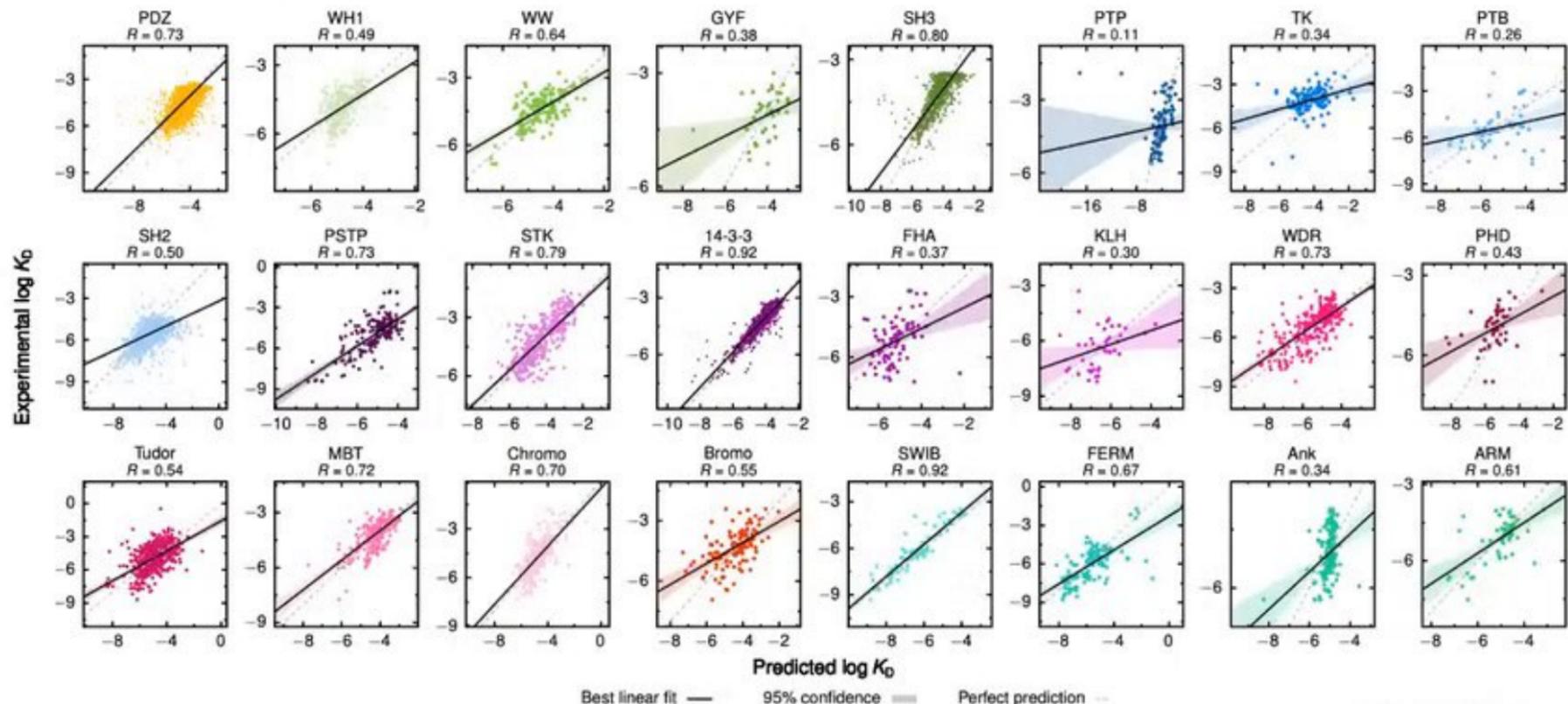
Quantitative statistical mechanical model (QSM)

Predicted observable

Optimization objective

Do QSM models predict quantitative affinities within PBD families?

QSM predicts quantitative affinities within PBD families



OSM predicts quantitative affinities within PBD families



Leverage statistical mechanics to construct a multitask objective for training a quantitative predictor of interaction affinity



Predicted log K_d

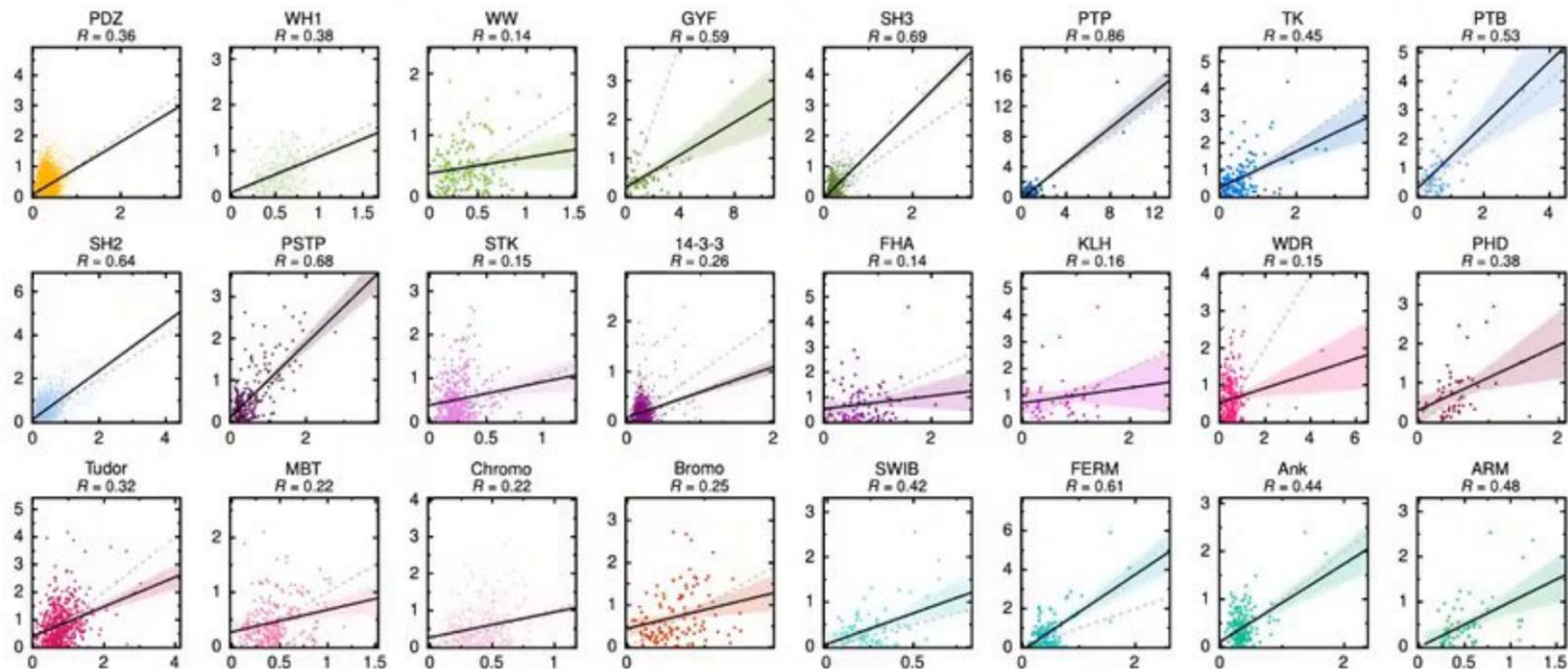
Experimental log K_d

$$\Delta F(\text{WT}; \text{WT})$$

$$\Delta F(\text{WT}; \text{WT})$$

How trustworthy are QSM affinity predictions?

QSM estimates uncertainty by predicting its own absolute error

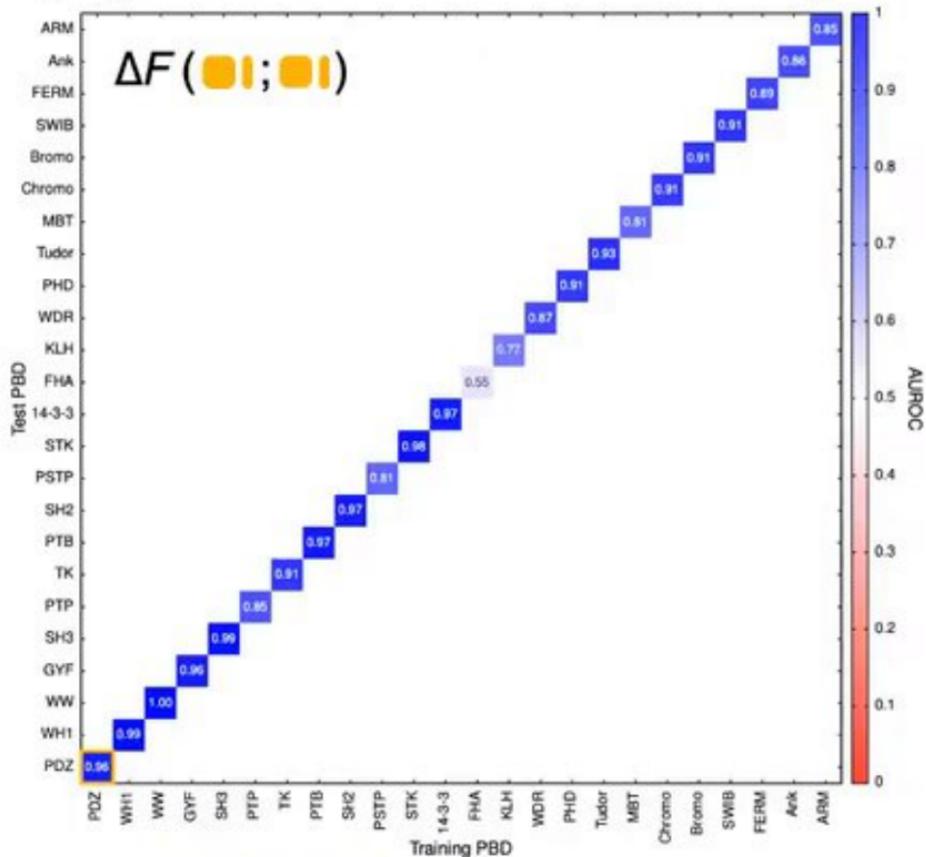
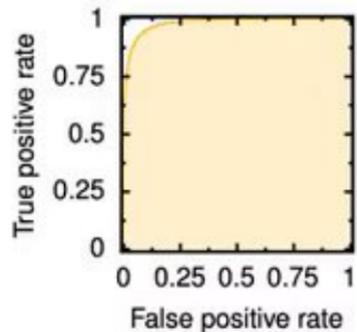


Best linear fit — 95% confidence — Perfect prediction - - -

Are QSM free energy functions universal?

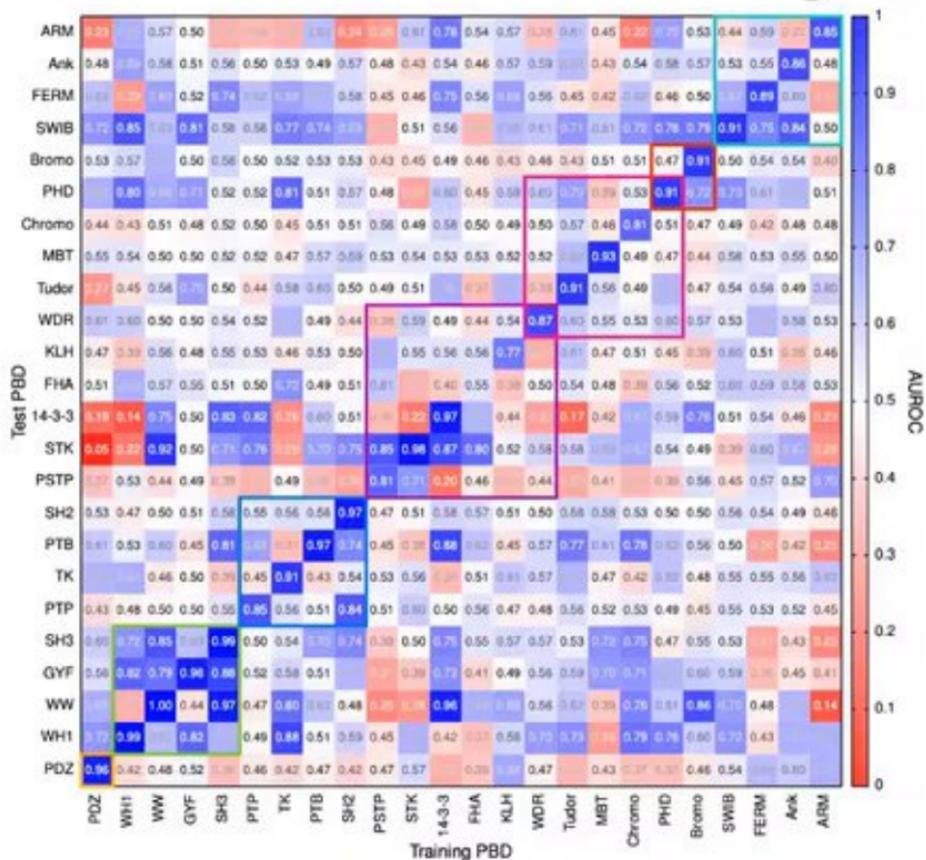
QSM free energy functions identify interactions within PBD families

Evaluated for identifying 'interaction' vs 'non-interacting' domain-peptide pairs



OSM free energy functions transfer across PBD families

Evaluated as classifier of 'interaction' vs 'non-interacting' domain-peptide pairs



PROTEIN QSM free energy functions transfer across PBD families
Evaluated as classifier of 'interaction' vs 'non-interacting' domain-peptide pairs

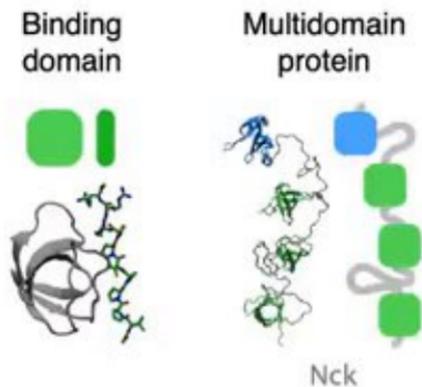


**Leverage a protein language model
to construct a universal predictor of interaction affinity**

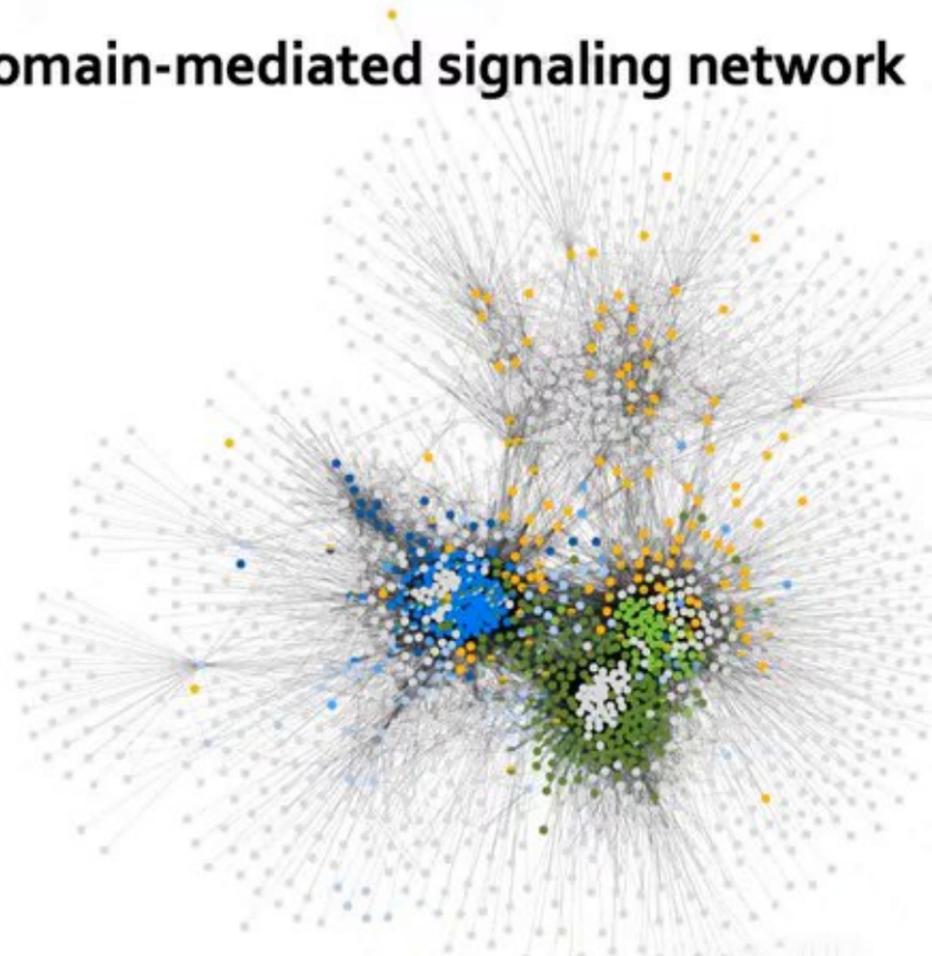


$\Delta F(\text{blue}; \text{green})$

Can QSM discover multivalent protein-protein interactions mediated by domain-peptide interactions?



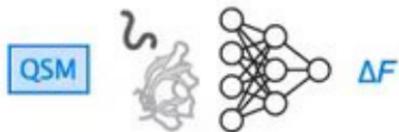
Constructing a proteome-scale domain-mediated signaling network



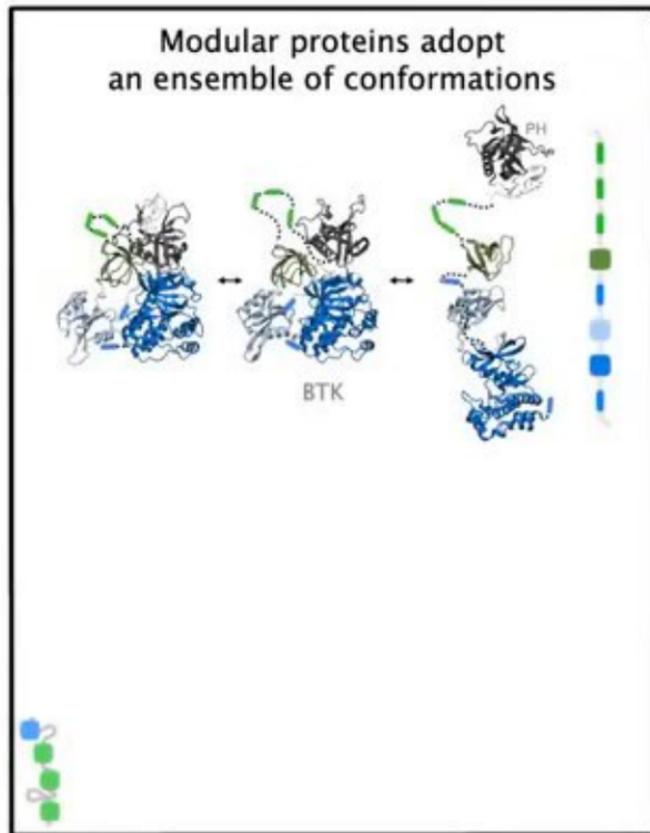
Domain family

SH2 PTB TK PTP SH3 WW WH1 PDZ

Constructing a proteome-scale domain-mediated signaling network



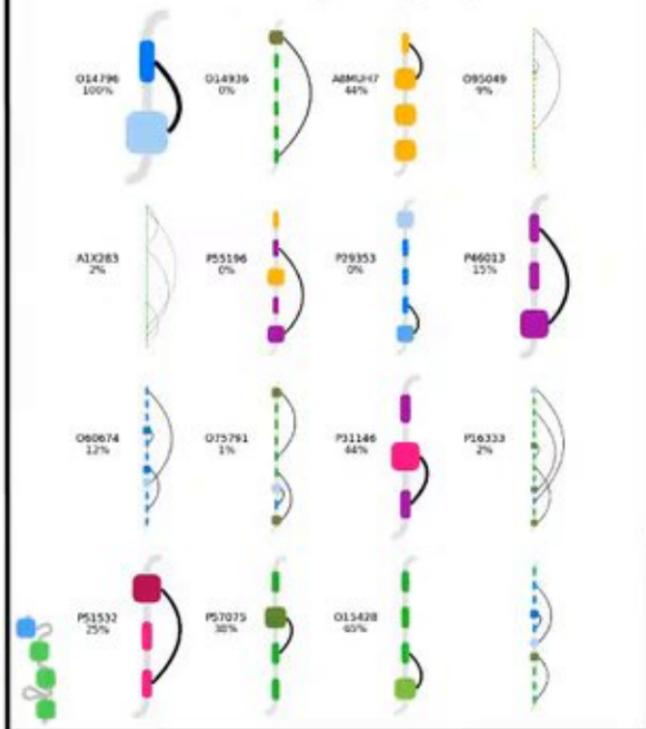
Quantitative models of cell signaling enabled by machine learning to decode the biophysical logic of cellular information processing



Quantitative models of cell signaling enabled by machine learning

to decode the *biophysical logic* of cellular information processing

How is protein *conformational plasticity* tuned for signaling logic?



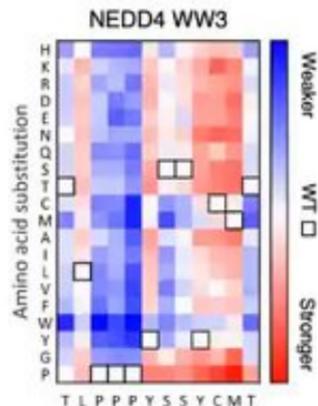
PBD-mediated protein-protein interaction network



Quantitative models of cell signaling enabled by machine learning

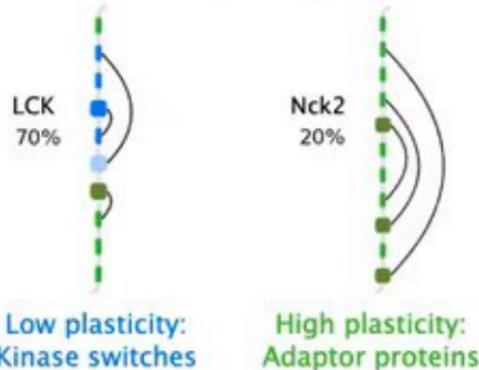
to decode the *biophysical logic* of cellular information processing

How does sequence variation tune PBD-SLiM *binding affinity*?



In silico
deep mutational scanning

How is protein *conformational plasticity* tuned for signaling logic?



What *multi-input computations* are performed by modular proteins?



PBD-mediated protein-protein interaction network



Acknowledgements

AIQuraishi Lab



Mohammed
AIQuraishi



Data curation

Ethan Eickmann
Undergraduate



Allen Na
Masters



Sev Ihnat
Masters



Bridget Liu
Undergraduate



Ashley Hsu
Undergraduate



Teemu Rönkkö
Visiting PhD



Lindsey Yang
Undergraduate



Henry Low
MD/PhD



Sumaiyah Rahman
Undergraduate



Angelina Yan
Undergraduate



Martin Culka
Postdoc



Model benchmarking

Shriya Mahakala
Undergraduate



Emily Duniec
Masters



Machine learning a quantitative statistical mechanical model for ΔF

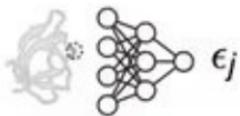
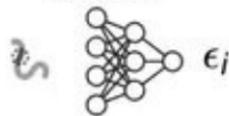
How to infer a free energy function from data?

Peptide sequence:

GTTPPPYTV



Neural free energy function



Domain sequence:

...AKTSSGQRYFLN...



Interaction energies



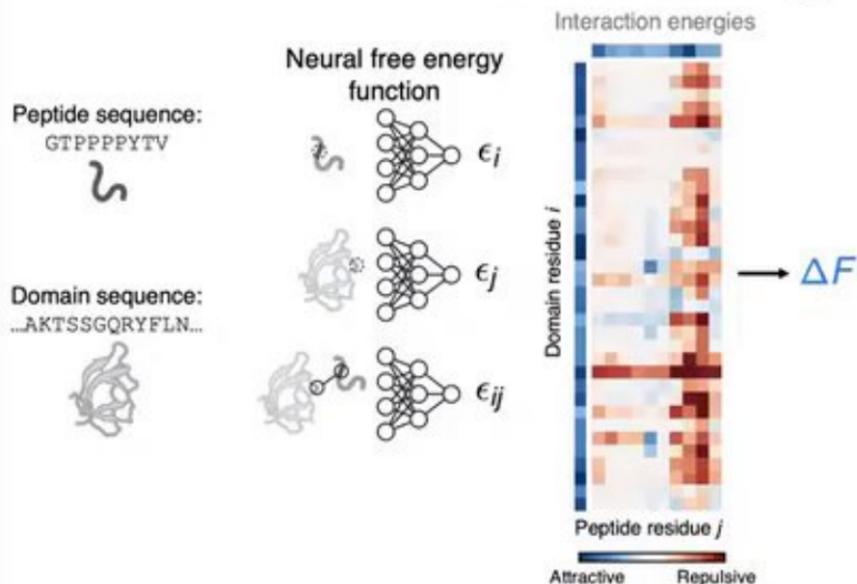
Attractive Repulsive

$$\Delta F = \sum_i^{N_{\text{dom}}} \epsilon_i^{(\text{dom})} + \sum_j^{N_{\text{pep}}} \epsilon_j^{(\text{pep})} + \sum_{ij}^{N_{\text{dom}} \times N_{\text{pep}}} \epsilon_{ij}^{(\text{int})} + \Delta S$$

Quantitative statistical mechanical model (QSM)

Machine learning a quantitative statistical mechanical model for ΔF

How to infer a free energy function from data?

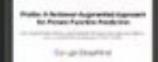


$$\Delta F = \sum_i^{N_{\text{dom}}} \epsilon_i^{(1)} + \sum_j^{N_{\text{pep}}} \epsilon_j^{(2)} + \sum_{ij}^{N_{\text{dom}} \times N_{\text{pep}}} \epsilon_{ij}^{(3)} + \Delta S$$

Quantitative statistical mechanical model (QSM)

NYGC Events

134_show Peter Shaw.pdf



1



2



3



4



5



47

48

49

50

ProtEx: A Retrieval-Augmented Approach for Protein Function Prediction

Peter Shaw, Bhaskar Gurram, David Belanger, Andreea Gane, Maxwell L. Bileschi, Lucy J. Colwell, Kristina Toutanova, Ankur P. Parikh



Google DeepMind

Zoom Meeting Controls

- Join
- Chat
- Share
- Comments
- Bookmarks
- File List:
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB Keynote
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB PDF Document
 - JB PDF Document
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB PDF Document
 - JB Keynote
 - JB PDF Document
 - JB PowerP...n (.pptx)
 - JB PowerP...n (.pptx)
 - JB PDF Document

Zoom Meeting Interface

- Zoom Meeting
- Screen Share
- 2:23 PM
- 2024-09-19 AM
- 2:23 AM
- 2:37 AM
- Zoom Meeting

134_show Peter Sheard

- Hide Toolbar
- Thumbnails
- Table of Contents
- Highlights and Notes
- Bookmarks
- Contact Sheet
- Continuous Scroll
- Single Page
- Two Pages
- Soft Proof with Profile
- Actual Size
- Zoom In
- Zoom Out
- Show Markup Toolbar
- Hide Toolbar
- Customize Toolbar...
- Slideshow**
- Enter Full Screen

: A Retrieval-Augmented Approach for Protein Function Prediction

Bhaskar Gurram, David Belanger, Andreea Gane, Maxwell L. Bileschi, Lucy J. Colwell, Kristina Toutanova, Ankur P. Parikh

Google DeepMind

Google DeepMind

Windows taskbar and application windows:

- System tray: Thu Sep 11 4:17 PM
- Taskbar: Comments, Share
- File Explorer: 134_show Peter Sheard
- Microsoft Word: 134_show Peter Sheard
- Microsoft Teams: 23 AM, 37 AM
- Zoom: boom

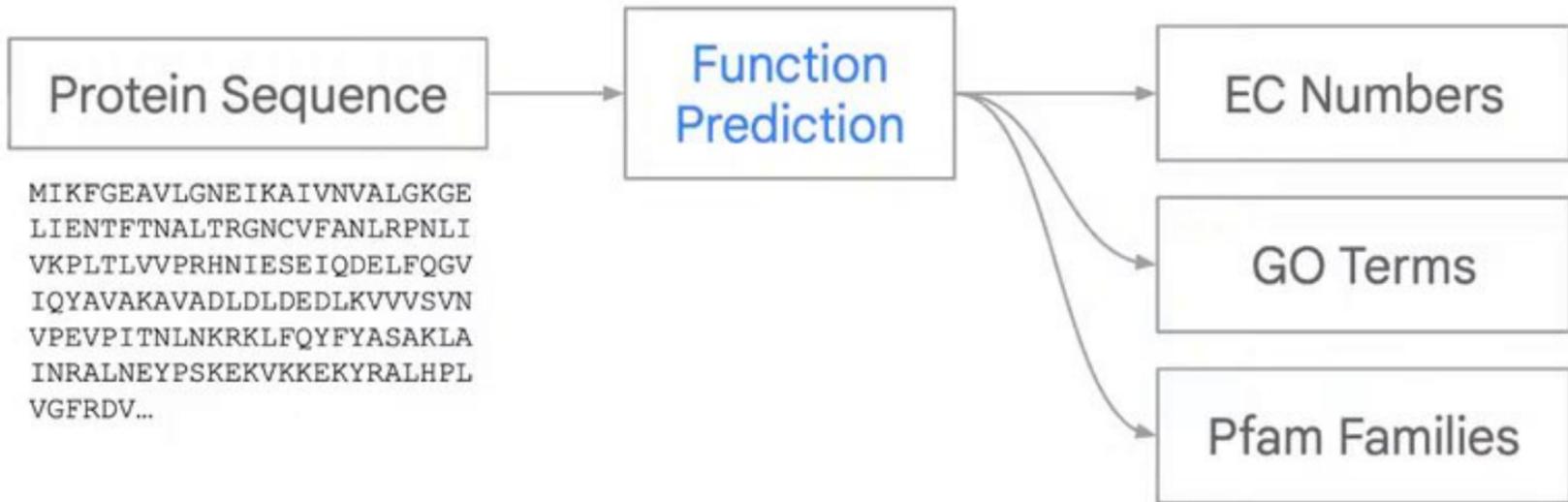
Calendar application interface:

- Days: 47, 48, 49, 50
- Buttons: Schedule, Spotlights & Posters

ProtEx: A Retrieval-Augmented Approach for Protein Function Prediction

Peter Shaw, Bhaskar Gurram, David Belanger, Andreea Gane, Maxwell L. Bileschi, Lucy J. Colwell, Kristina Toutanova, Ankur P. Parikh

 DeepMind



Protein Sequence

Function Prediction

EC Numbers

GO Terms

Pfam Families

MIKFGEAVLGNEIKAIVNVALGKGE
LIENTETNALTRGNCVEANLRPNLT
VKP
IQY
VPE
INR
VGF

News > [All EMBL-EBI news](#) > Technology and innovation

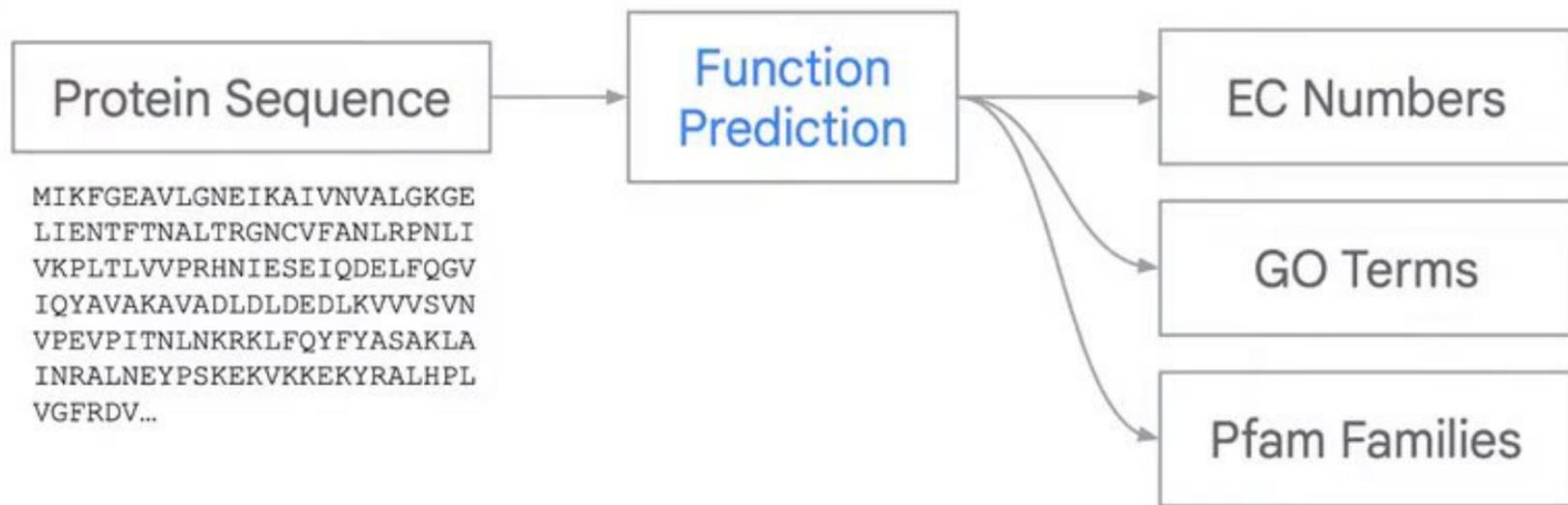
Natural language processing for rapid protein annotation

Over 40 million protein annotations have been added to the UniProt database using a Google Research natural language processing model

*Gane et al. (2022)
in collaboration with*

EMBL-EBI

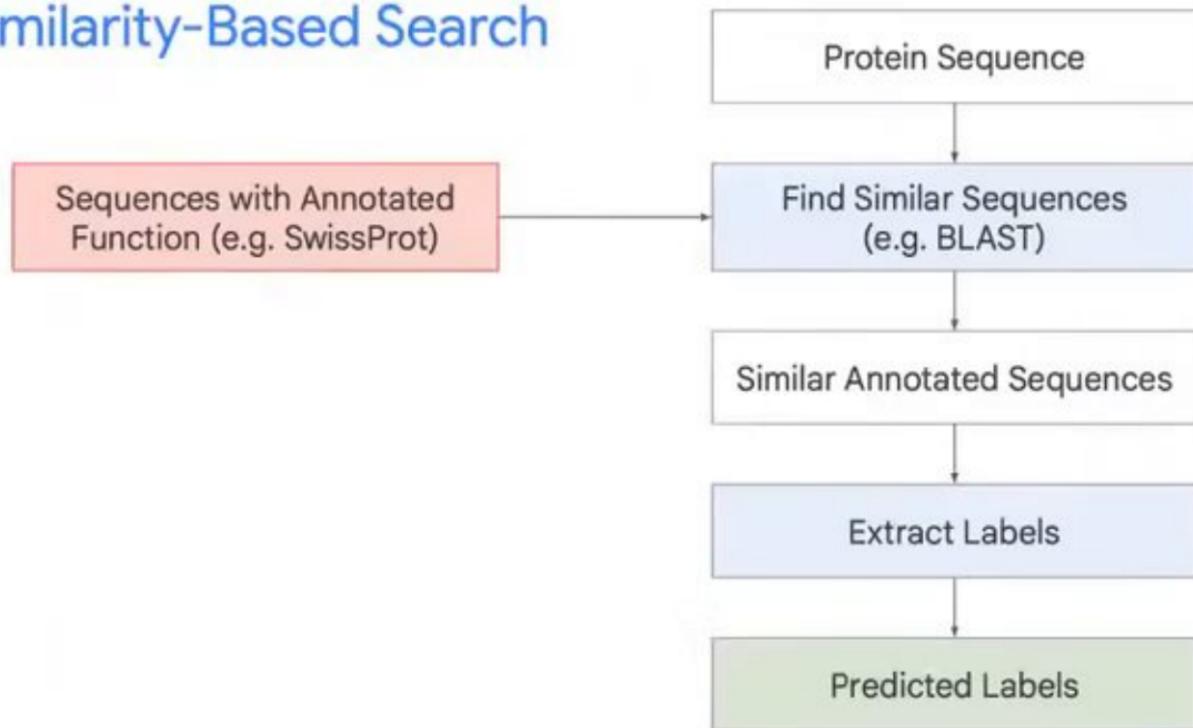




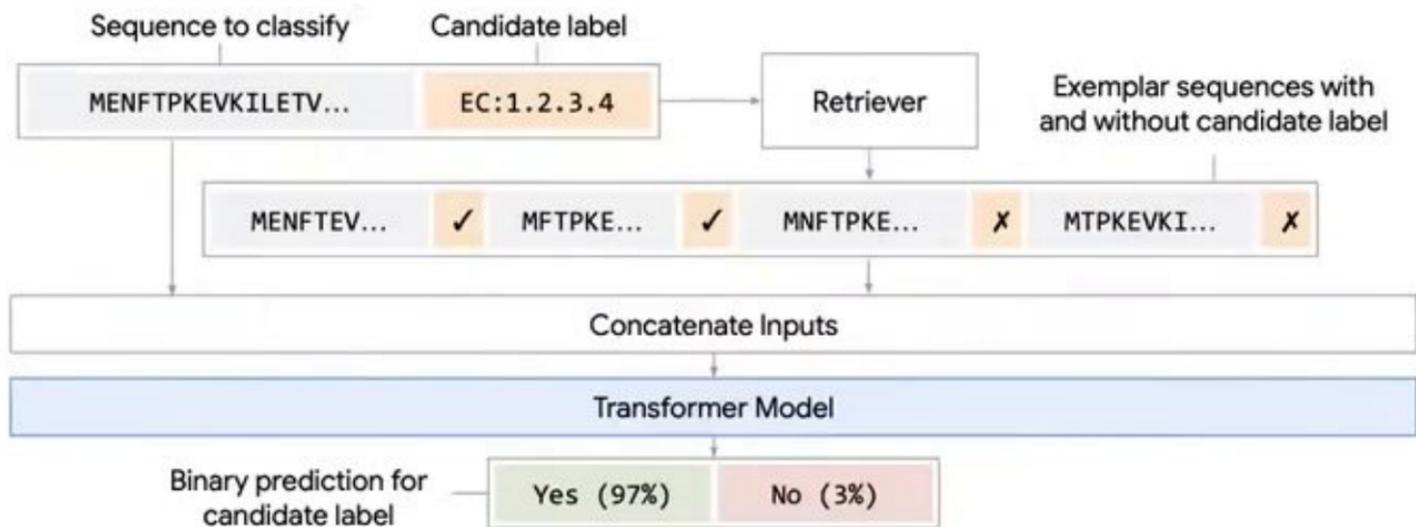
Key machine learning challenges:

- *Protein sequences that are dissimilar from annotated proteins*
- *Rare functional labels*

Similarity-Based Search



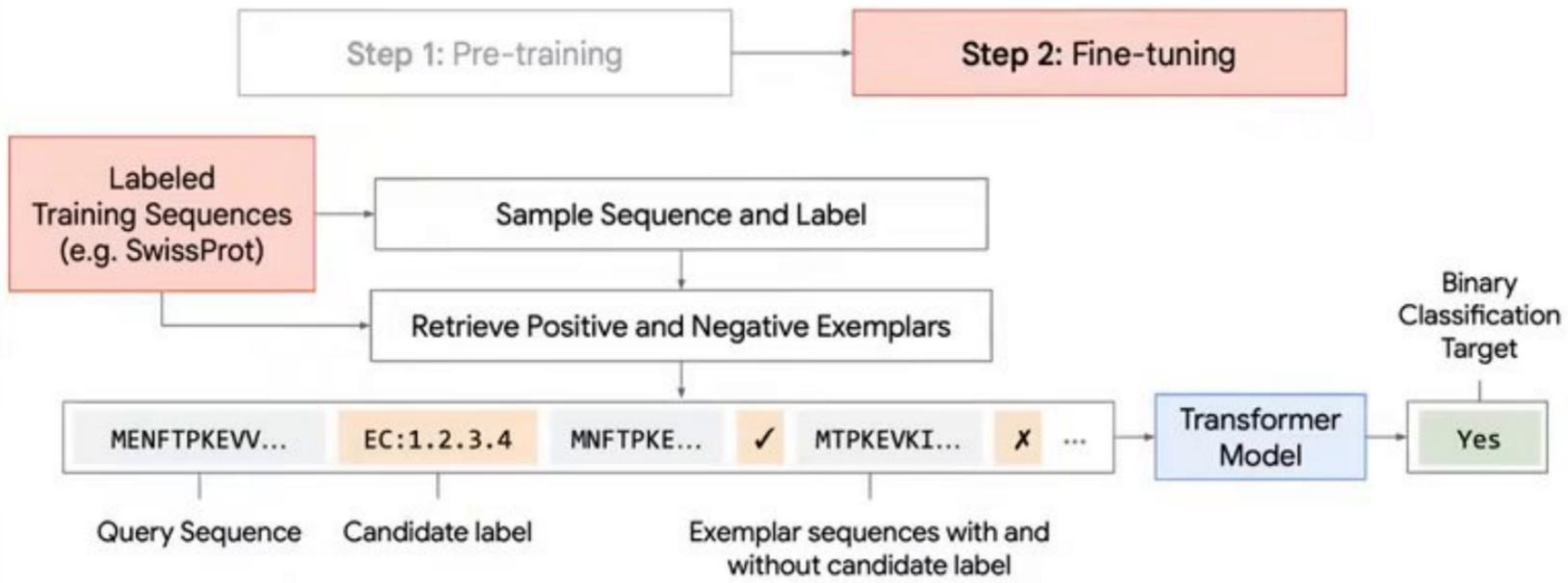
ProtEx Overview



ProtEx Training Pipeline

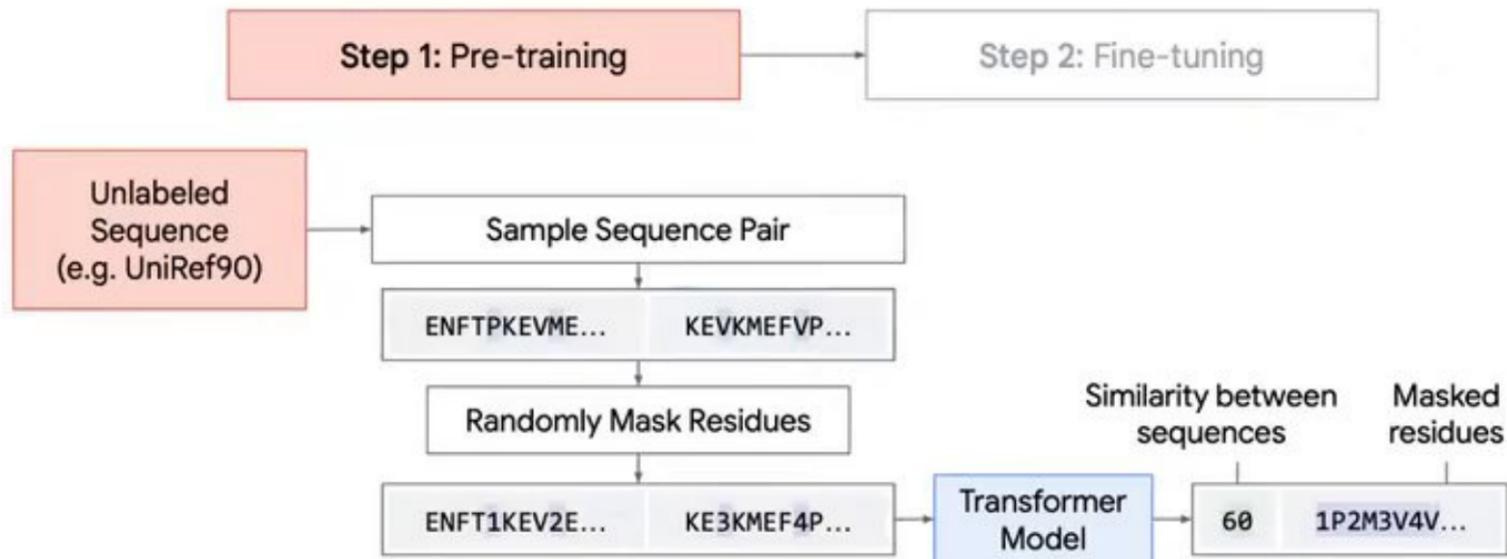


ProtEx Training Pipeline: Fine-tuning



Bias training time sampling towards less similar exemplars to improve generalization

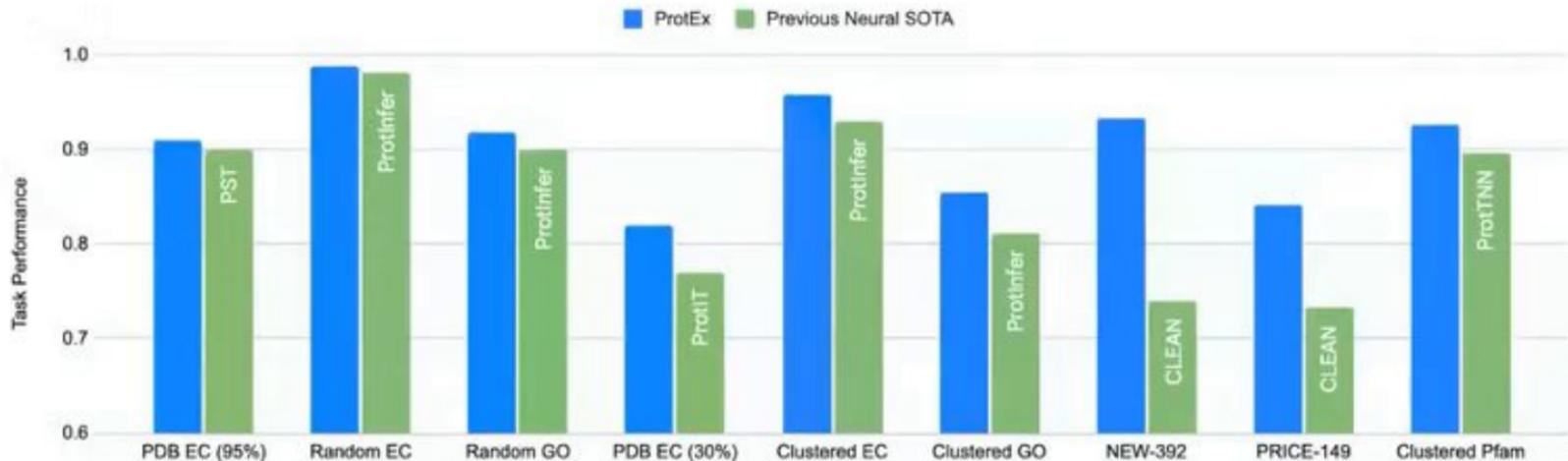
ProtEx Training Pipeline: Pre-training



Pre-training encourages model to implicitly align and compare multiple sequences

Experiments and Results

Task Performance on EC, GO, and Pfam Benchmarks

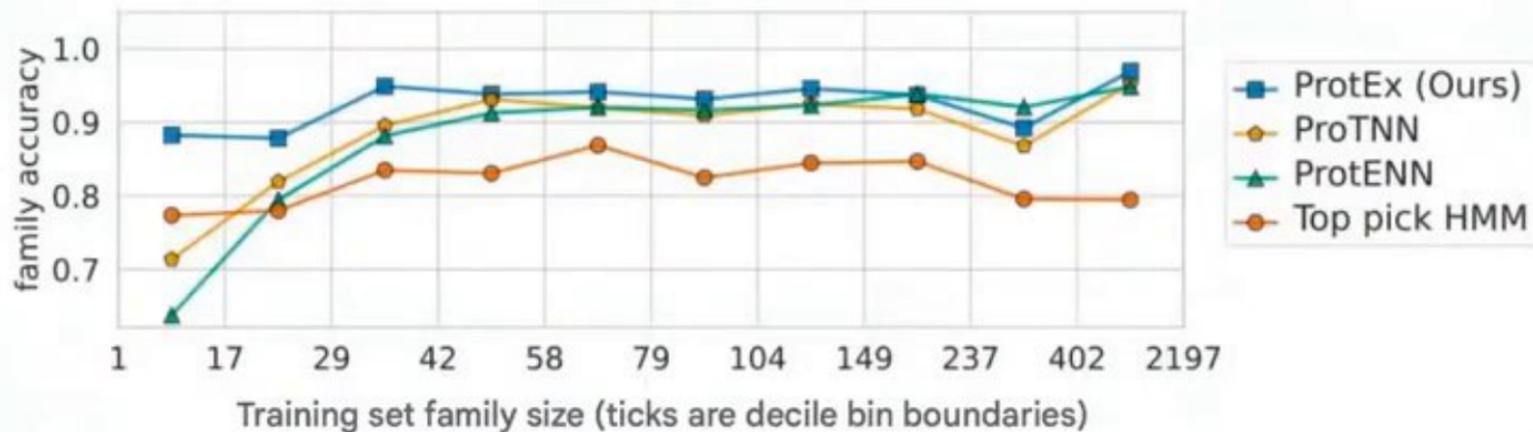


ProtEx outperforms strong neural methods from prior work across EC, GO, and Pfam baseline tasks

We report Max F1 for the random and clustered EC and GO splits of Sanderson et al. (2023), Weighted AUC for NEW-392 and PRICE-149 splits of Yu et al. (2023), and Family Accuracy for the Pfam clustered split of Bileschi et al. (2022). Previous SOTA neural methods are from PST (Chen et al., 2024), ProtInfer (Sanderson et al., 2023), ProtIR (Zhang et al., 2024), CLEAN (Yu et al., 2023), and ProtTNN (Dohan et al., 2021).

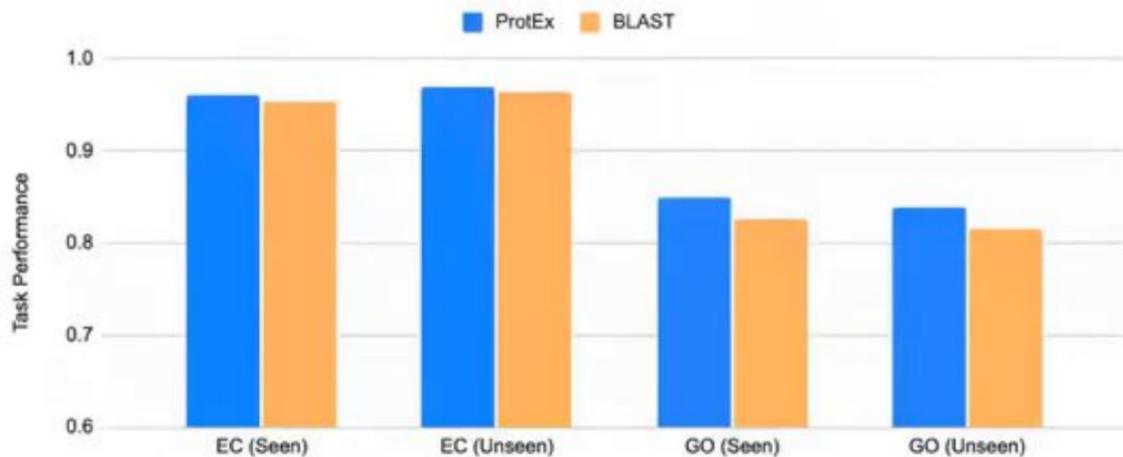
Performance for rare functional labels

ProtEx performance on clustered Pfam task stratified by number of annotated proteins per family



ProtEx's improvements are largest when fewer training examples are available for a given functional label

Performance for unseen functional labels



ProtEx generalizes well even if we randomly remove a subset of functional labels during fine-tuning

Conclusions and Future Work

- Retrieving and conditioning on relevant annotated proteins can improve function prediction accuracy
- “Agents” with tools are an increasingly popular instance of inference-time retrieval
- What other types of information could be useful to retrieve at inference time for protein function prediction?

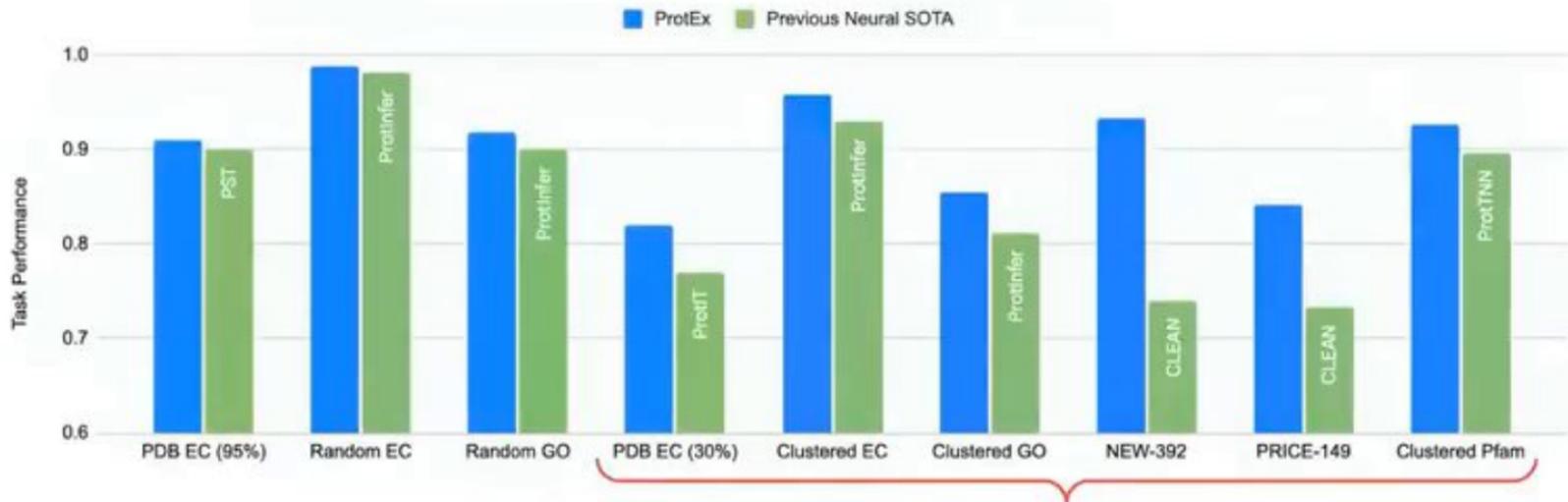
Thank you!

<https://github.com/google-deepmind/protex>

Peter Shaw, Bhaskar Gurram, David Belanger, Andreea Gane, Maxwell L. Bileschi,
Lucy J. Colwell, Kristina Toutanova, Ankur P. Parikh

Google DeepMind

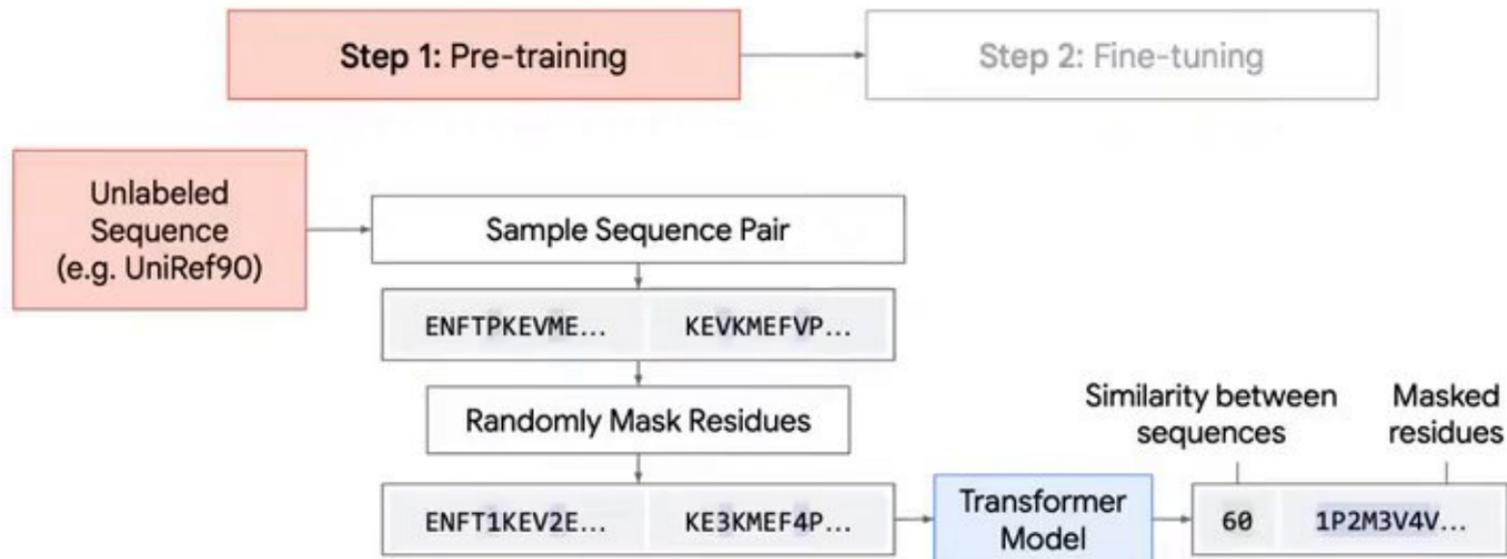
Task Performance on EC, GO, and Pfam Benchmarks



Especially for sequences dissimilar to the training set

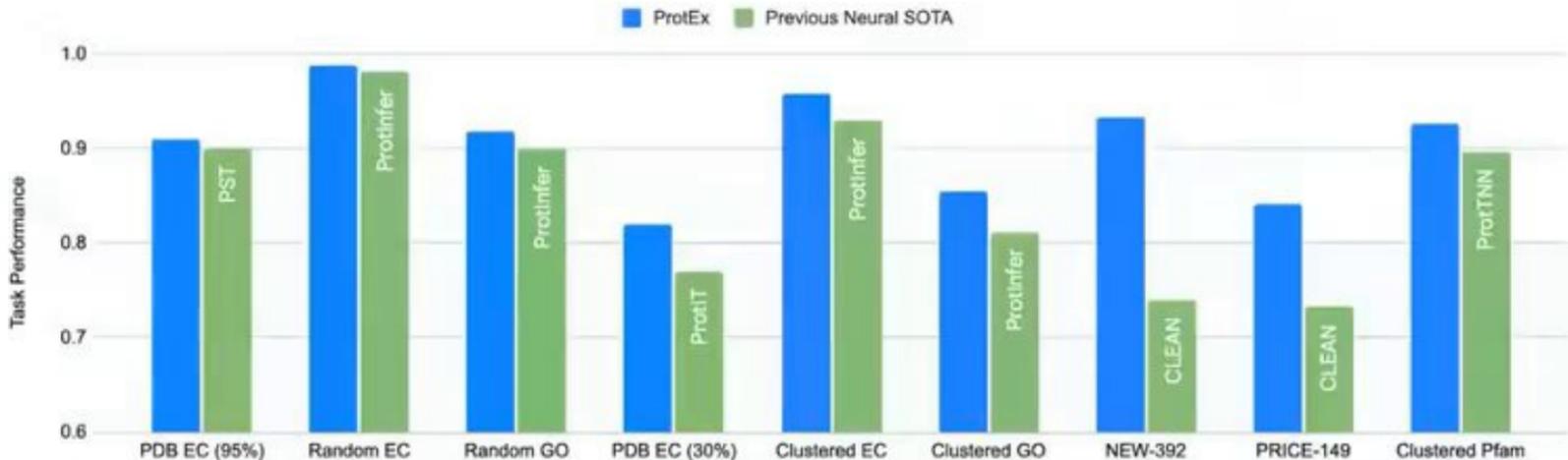
We report Max F1 for the random and clustered EC and GO splits of Sanderson et al. (2023), Weighted AUC for NEW-392 and PRICE-149 splits of Yu et al. (2023), and Family Accuracy for the Pfam clustered split of Bileschi et al. (2022). Previous SOTA neural methods are from PST (Chen et al., 2024), ProtInfer (Sanderson et al., 2023), ProtIR (Zhang et al., 2024), CLEAN (Yu et al., 2023), and ProtTNN (Dohan et al., 2021).

ProtEx Training Pipeline: Pre-training



Pre-training encourages model to implicitly align and compare multiple sequences

Task Performance on EC, GO, and Pfam Benchmarks



ProtEx outperforms strong neural methods from prior work across EC, GO, and Pfam baseline tasks

We report Max F1 for the random and clustered EC and GO splits of Sanderson et al. (2023), Weighted AUC for NEW-392 and PRICE-149 splits of Yu et al. (2023), and Family Accuracy for the Pfam clustered split of Bileschi et al. (2022). Previous SOTA neural methods are from PST (Chen et al., 2024), ProtInfer (Sanderson et al., 2023), ProtIR (Zhang et al., 2024), CLEAN (Yu et al., 2023), and ProtTNN (Dohan et al., 2021).

Task Performance on EC, GO, and Pfam Benchmarks



Ablations highlight that conditioning on exemplars is key to ProtEx's performance

We report Max F1 for the random and clustered EC and GO splits of Sanderson et al. (2023), Weighted AUC for NEW-392 and PRICE-149 splits of Yu et al. (2023), and Family Accuracy for the Pfam clustered split of Bileschi et al. (2022).

Conclusions and Future Work

- Retrieving and conditioning on relevant annotated proteins can improve function prediction accuracy
- “Agents” with tools are an increasingly popular instance of inference-time retrieval
- What other types of information could be useful to retrieve at inference time for protein function prediction?

Thank you!

<https://github.com/google-deepmind/protex>

Peter Shaw, Bhaskar Gurram, David Belanger, Andreea Gane, Maxwell L. Bileschi,
Lucy J. Colwell, Kristina Toutanova, Ankur P. Parikh

Google DeepMind

NYGC Events

Learning Protein Dynamics at Proteome Scale

Samuel Sledzieski
Machine Learning in Computational Biology
New York Genome Center
11 September 2025

How do we learn about protein function?

Transport?

Structure?

Catalytic Activity?

Signaling?

Immunity?

Gene Regulation?

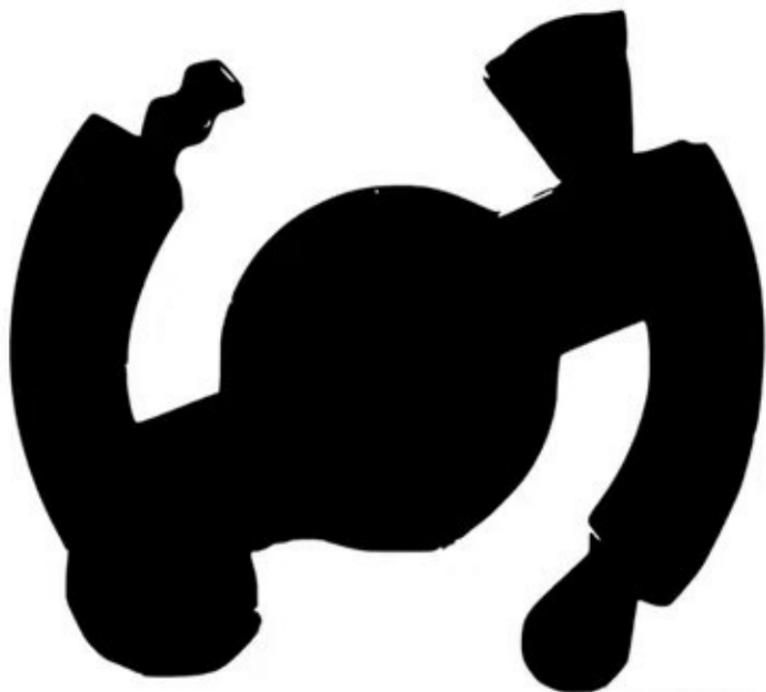
MLPGLALLLLAAWTARALEVPTD
GNAGLLAEPQIAMFCGRINMNMN
VQNGKWSDSPSGTKTCIDTKEGI
LQYCQEVYPELQITNVVEANQPV
TIQNWCKRGRKQCKTHPHFVIPI
RCLVGEFVSDALLVPDKCKFLHQ
ERMDVCETHLHWHTVAKETCSEK
STNLHDYGMLLPCGIDKFRGVEF
VCCPLAEESDNVDSADAEEEDSD
VWWGGADTDYADGSEDKVVEVAE
EEEVAEVEEEEADDEDEDGDE
VEEEAEEPYEATERTTSIATTT
TTTTESVEEVVREVCSEQAETGP
CRAMISRWFYFDVTEGKCAPFFYG
GCGGNRNNFDTEEYCMVCGSAM
SQSLLKTTQEPLARDPVKLPTTA
ASTPDAVDKYLETPGDENEHAHF
QKAKERLEAKHRERMSQVMREWE
FAERQAKNLRKADKKAVLQHEQF

How do we learn about protein function?

MLPGLALLLLAAWTARALEVPTDG
NAGLLAEPQIAMFCGRLNMHMNVQ
NGKWSDPSGKTCIDTKEGILQY
CQEVYPELQITNVVEANQPVTIQN
WCKRGRKQCKTHPHFVIPYRCLVG
EFVSDALLVPDKCKFLHQERMDVC
ETHLHWHTVAKETCSEKSTNLHDY
GMLLPCGIDKFRGVEFVCCPLAEE
SDNVDSADAEEDSDVWVGADTD
YADGSEDKVVEVAEEEEVAEVEEE
EADDEDEDGDEVEEEAEPEYEE
ATERTTSIATTTTTTTESVEEVVR
EVCSEQAETGPCRAMISRWFVDVT
EGKCAPFFYGGCGNRNFDTEEY
CMAVCGSAMSQSLLKTTQEPLARD
PVKLPTTAASTPDAVDKYLETPGD
ENEHAHFQKAKERLEAKHRERMSQ
VMREWEEAERQAKNLPKADKKA
VIQHFQEKVESLEQE

How do we learn about protein function?

MLPGLALLLLAAWTARALEVPTDG
NAGLLAEPQIAMFCGRLNMHMNVQ
NGKWSDPSGKTKCIDTKEGILQY
CQEVYPELQITNVVEANQPVTIQN
WCKRGRKQCKTHPHFVIPYRCLVG
EFVSDALLVPDKCKFLHQERMDVC
ETHLHWHTVAKETCSEKSTNLHDY
GMLLPCGIDKFRGVEFVCCPLAEE
SDNVDSADAEEDSDVWGGADTD
YADGSEDKVVEVAEEEEVAEVEEE
EADDEDEDDEGDEVEEEAEPEYEE
ATERTTSIATTTTTTTEVVEEVVR
EVCSEQAETGPCRAMISRWFVDVT
EGKCAPFFYGGCGGNRNNFDTEEY
CMAVCGSAMSQSLLKTTQEPLARD
PVKLPTTAASTPDAVDKYLETPGD
ENEHAHFQKAKERLEAKHRERMSQ
VMREWEAAERQAKNLPKADKKA VI
QHFQEKVESLEQE



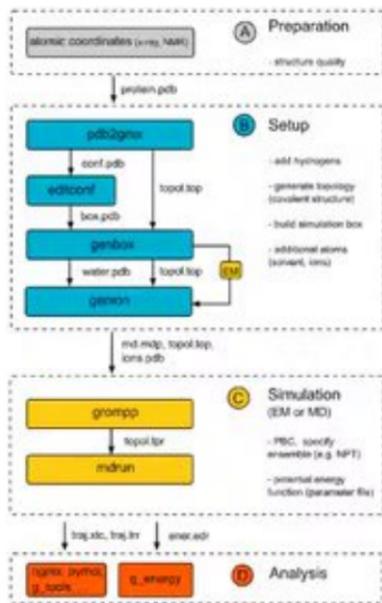
How do we learn about protein function?

- Allosteric regulation requires conformational change.
- Receptor activation and signaling requires helical motion.
- Natively disordered regions are functional on binding.
- Catalytic sites may be dynamically accessible.

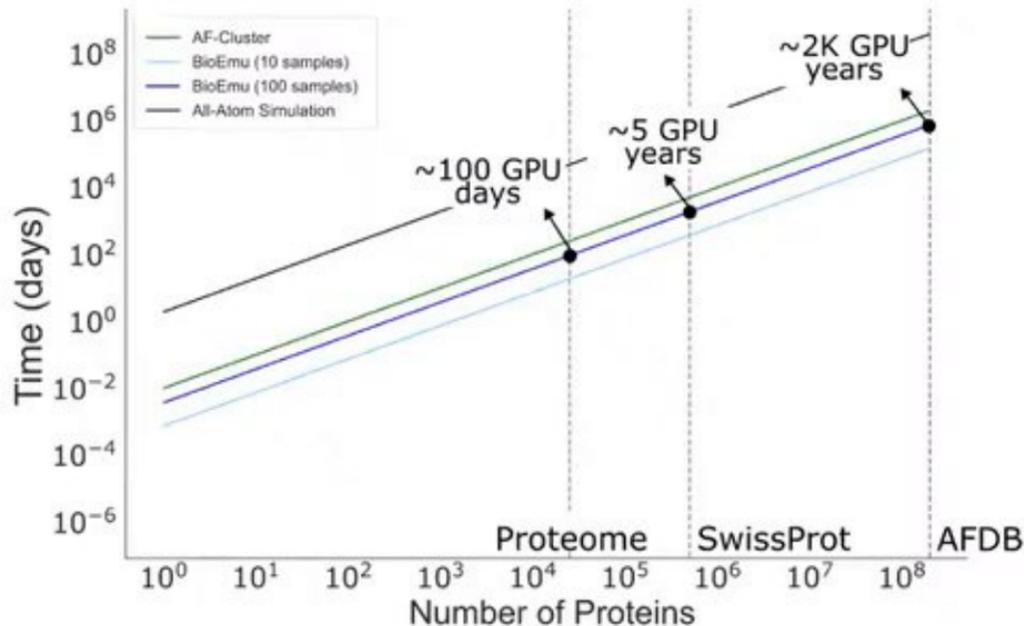


Current approaches to modeling dynamics

Simulation

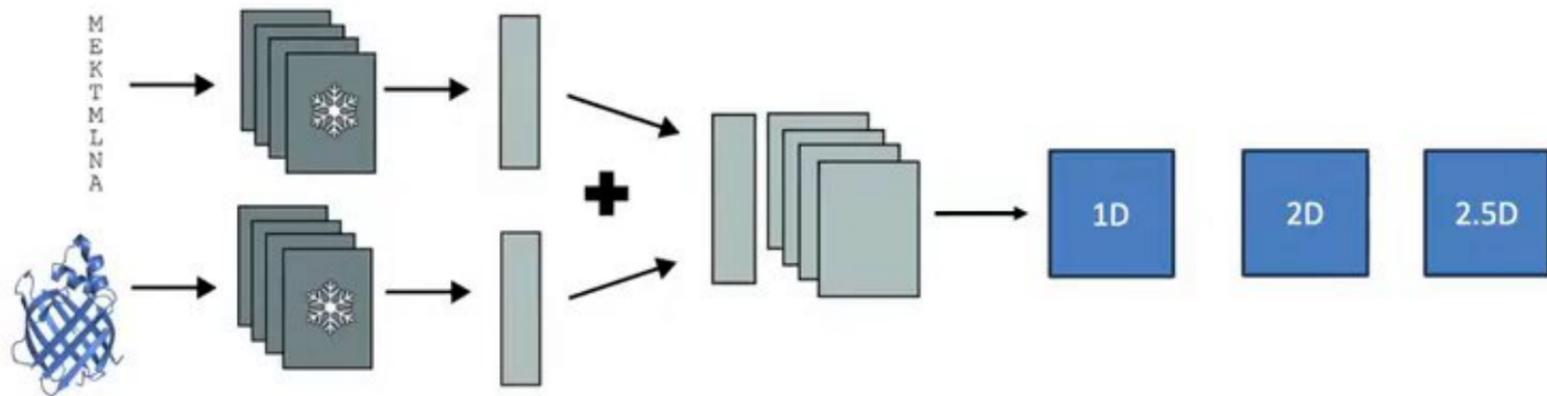


Limitations – genome scale \times variant effect

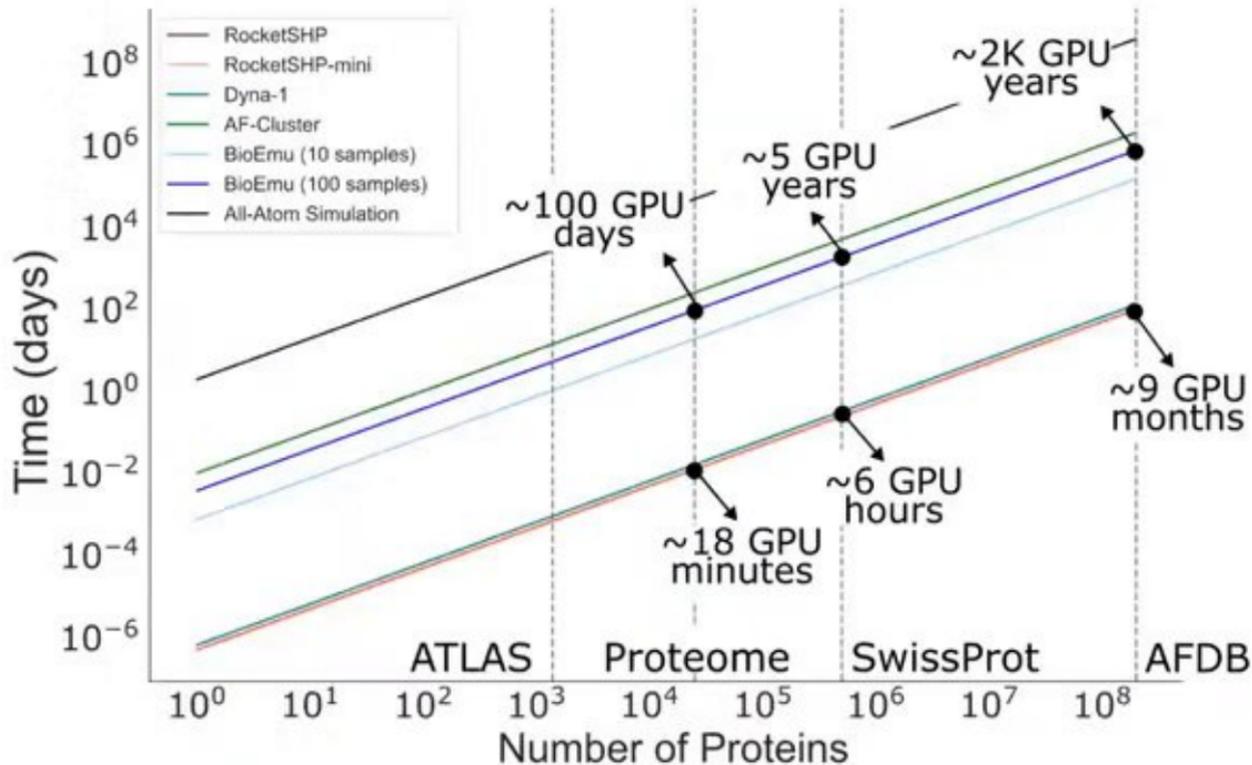


- We aim for high throughput *de novo* prediction of properties of protein dynamics.
- Differential splicing of isoforms, missense variants, somatic mutations all contribute to differential dynamics (>1 million in ClinVar as of 2021)
- Sequence space is vast: we want to functionally characterize non-human and engineered proteins

RocketSHP: fast Structure Heterogeneity Profiles

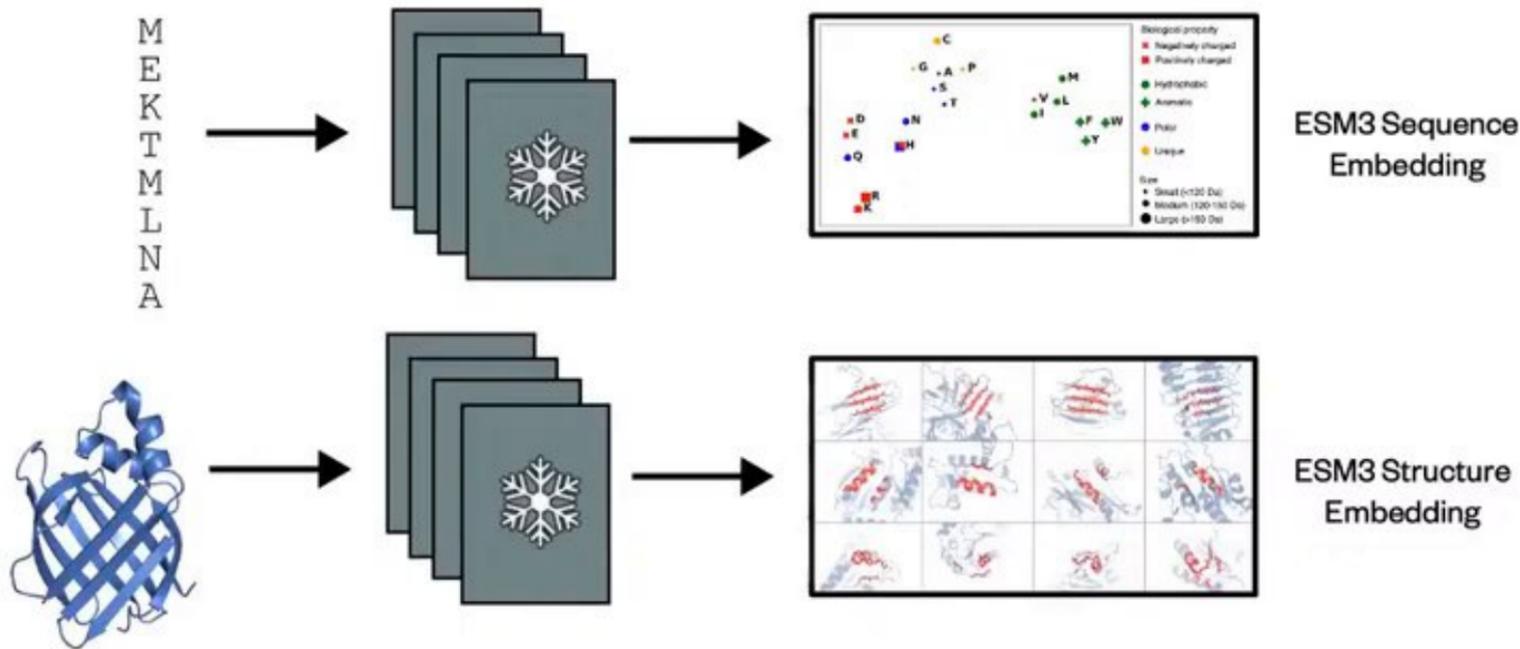


RocketSHP is highly scalable



~0.1s per protein
>1,000x speedup

Enabled by rich pre-trained features...



...and large-scale MD data sets



ATLAS: protein flexibility description from atomistic molecular dynamics simulations

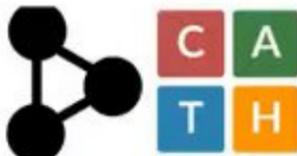
Yann Vander Meersche¹, Gabriel Cretin², Aria Gheeraert³, Jean-Christophe Gelly^{1,4*} and Tatiana Galochkina^{1,4*}

¹Université Paris Cité and Université des Antilles and Université de la Réunion, INSERM, B1GR, F-75014 Paris, France

²To whom correspondence should be addressed. Tel: +33 1 81 72 43 30; Email: tatiana.galochkina@u-paris.fr

³Correspondence may also be addressed to Jean-Christophe Gelly. Tel: +33 1 81 72 43 22; Email: jean-christophe.gelly@u-paris.fr

- 1390 proteins with diverse ECOD domains
- 3 replicates at 300K
- 2fs time step, coordinates saved every 10ps for 100ns
- GROMACS w/ CHARMM36m
- Clustered into 1,039 Foldseek clusters for train/validation/test splits



mdCATH: A Large-Scale MD Dataset for Data-Driven Computational Biophysics

Antonio Mirarchi^{1,†}, Toni Giorgino^{2,1,†}, and Gianni De Fabritis^{1,3,4,*}

¹Computational Science Laboratory, Universitat Pompeu Fabra, Barcelona Biomedical Research Park (PRBB), Carrer Dr. Aiguader 88, Barcelona, 08003, Spain.

²Biophysics Institute, National Research Council (CNR-IBF), Via Celoria 26, Milan, 20133, Italy

³Institut Catalana de Recerca i Estudis Avançats (ICREA), Passeig Lluís Companys 23, Barcelona, 08010, Spain.

⁴Acellera Labs, Doctor Trueta 183, Barcelona, 08005, Spain.

*corresponding author(s): Gianni De Fabritis (g.defabritis@gmail.com), Toni Giorgino (toni.giorgino@cnr.it)

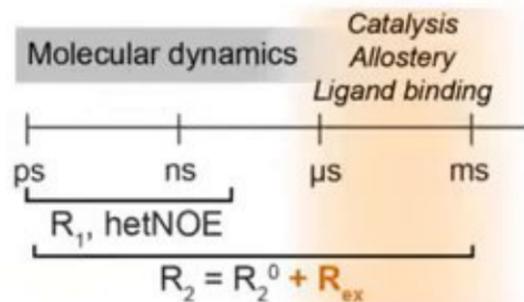
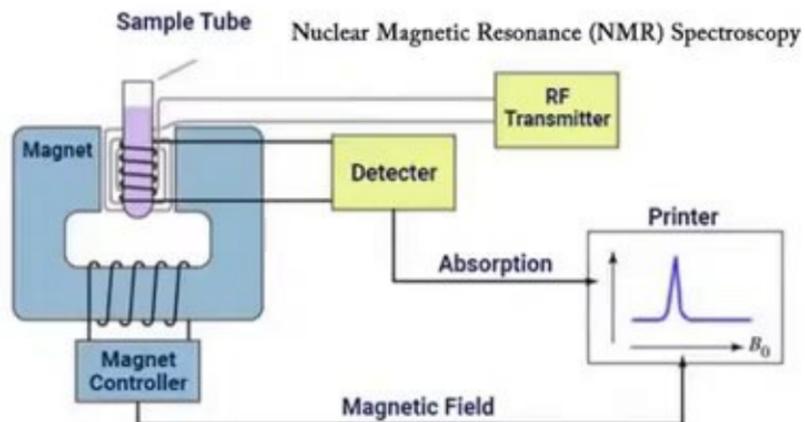
[†]A. M. and T. G. contributed equally to this work

- mdCATH: 5,398 CATH domains
- 5 replicates at 5 temperatures (320K, 348K, 379K, 413K, 450K)
- 4fs time step, coordinates saved every 1ns for ~464ns
- ACEMD on GPUGRID.net w/ CHARMM22
- Clustered into 3,178 Foldseek clusters for train/validation/test splits

RocketSHP efficiently predicts descriptors of protein dynamics

Model	RMSF		GCC-LMI		SHP	Model Information		
	RMSE	Spear. ρ (\uparrow)	GDD	IMSD	KL-Div	Seq.	Struct.	Params.
RocketSHP-mini	0.104	0.680	1.023	2.532	1.089	✓	✗	1.5M
RocketSHP-seq	0.090	0.719	0.944	2.387	1.225	✓	✗	29.7M
RocketSHP	0.083	0.789	0.859	1.898	1.551	✓	✓	30.2M
BioEmu (10 samples)	0.192	0.715	2.550	3.357	1.923	✓	✗	31.2M
BioEmu (100 samples)	0.202	0.748	1.215	1.871	1.778	✓	✗	31.2M
Dyna-1	0.361	0.267	-	-	-	✓	✓	188.9M
Dyna-1 (calibrated)	0.191	0.267	-	-	-	✓	✓	188.9M

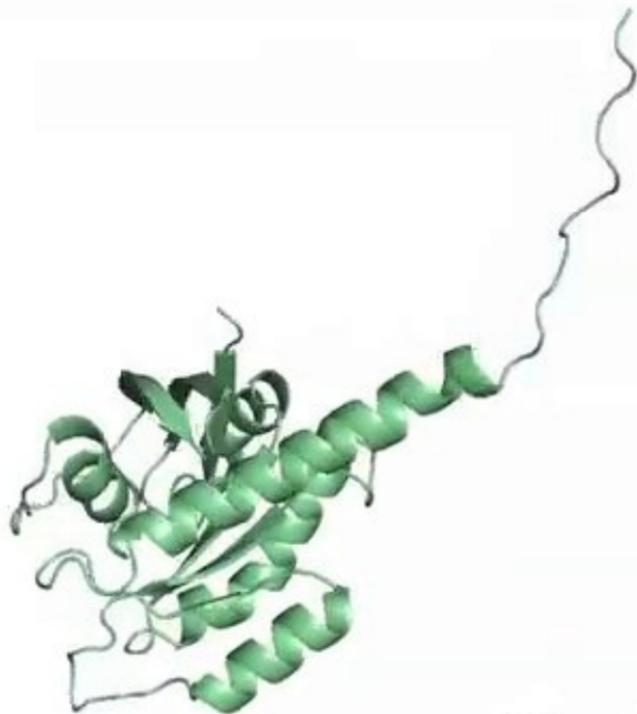
Predicting NMR hetNOE with RocketSHP



RocketSHP predicts temperature-dependent behavior

$$B_{obs} = B_0 e^{kt} \quad \longrightarrow \quad B_{T2} = B_{T1} e^{k(T2 - T1)}, k = 0.0045 K^{-1}$$

Case Study: Allosteric Networks in KRAS

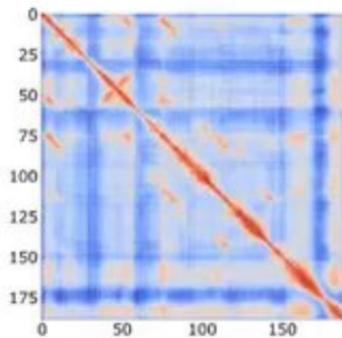


100ns all-atom MD

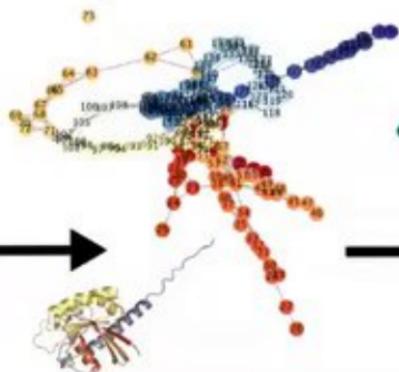
- KRAS is a frequently-mutated oncogene implicated in diverse cancer types
- Popular but challenging therapeutic target due to allosteric behavior
- “Molecular switch” relays signals only on binding by GTP
- Can we model the “allosteric network” of KRAS?

Case Study: Allosteric Networks in KRAS

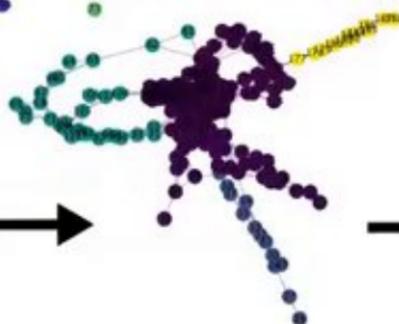
Predicted Correlations



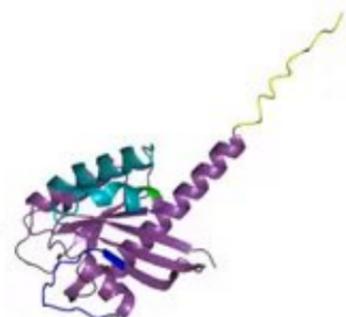
By Residue Index



Clustered
(Girvan-Newman)

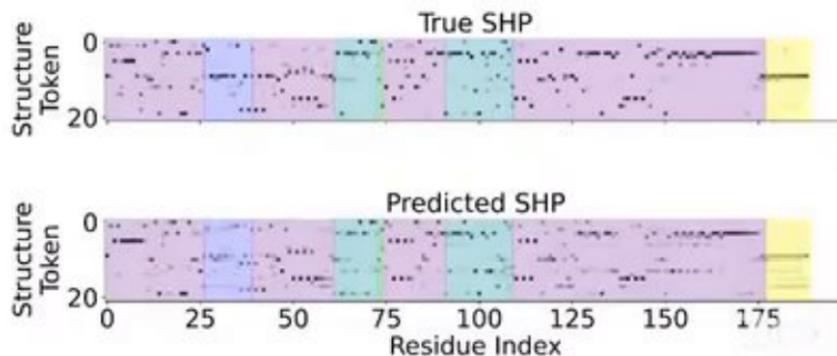
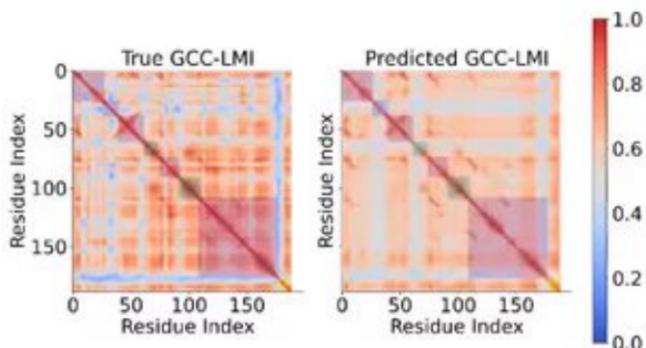
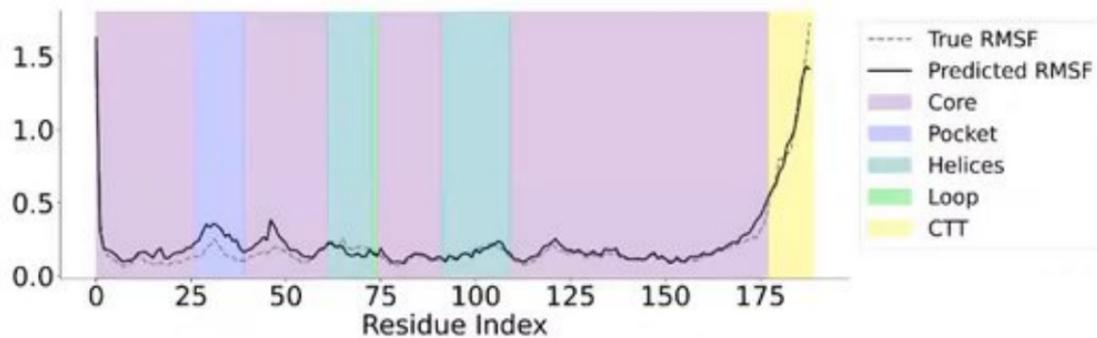
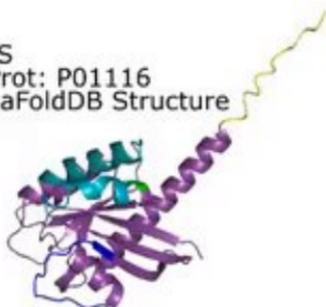


Protein Sub-structures

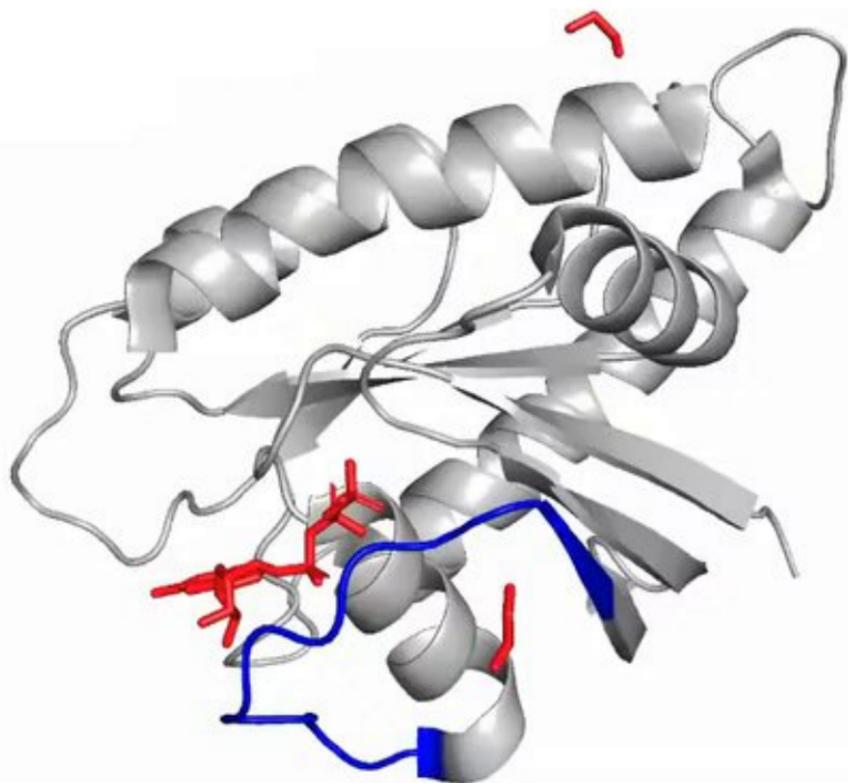


Case Study: Allosteric Networks in KRAS

KRAS
UniProt: P01116
AlphaFoldDB Structure



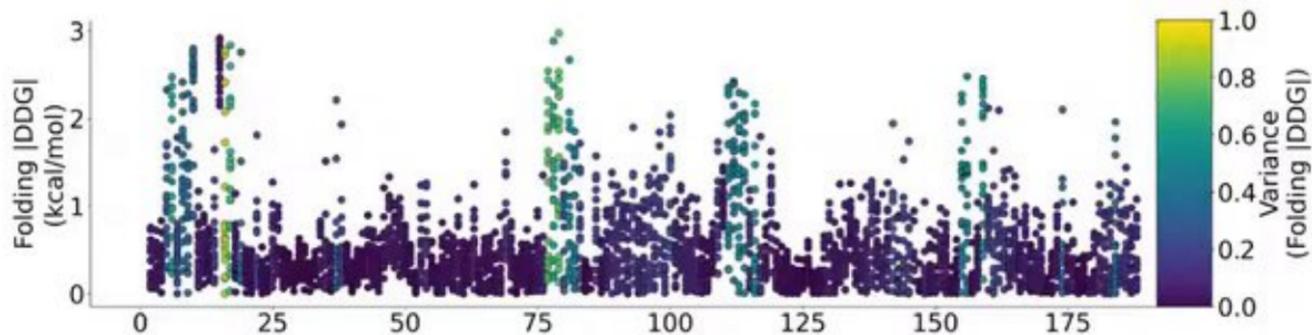
Case Study: Allosteric Networks in KRAS



PDB: 4DSN

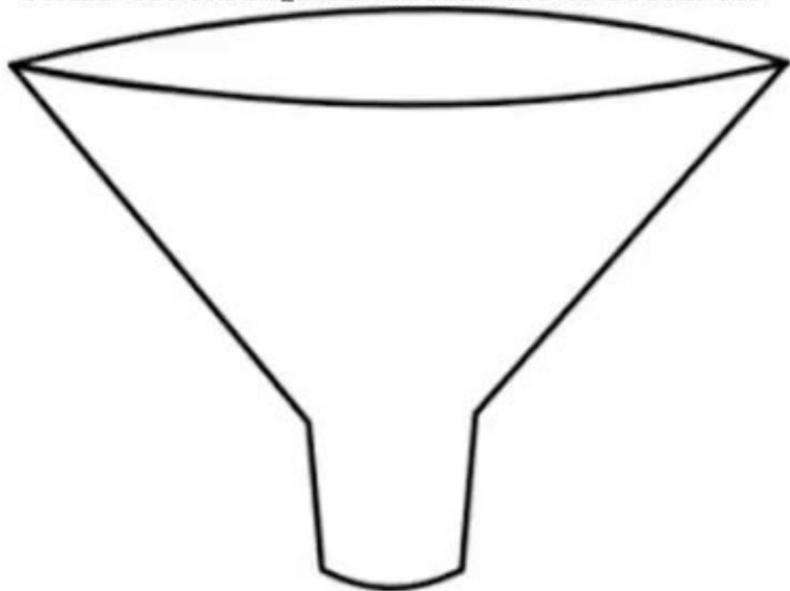
KRAS w/ bound ligand

Case Study: Allosteric Networks in KRAS



Outlook and Limitations

MLPGLALLLLAAWTARALEVPTDGNAGLLAEPQIAMFCGRL
NMHMNVQNGKWSDPSGTTKTCIDTKEGILQYCQEVYPELQI
TNVVEANQPVTIQNWCKRGRKQCKTHPHFVIPYRCLVGEFV
SDALLVPDKCKFLHQERMDVCETHLHWHTVAKETCSEKSTN



Acknowledgements



Sonya Hanson

Structural and Molecular
Biophysics

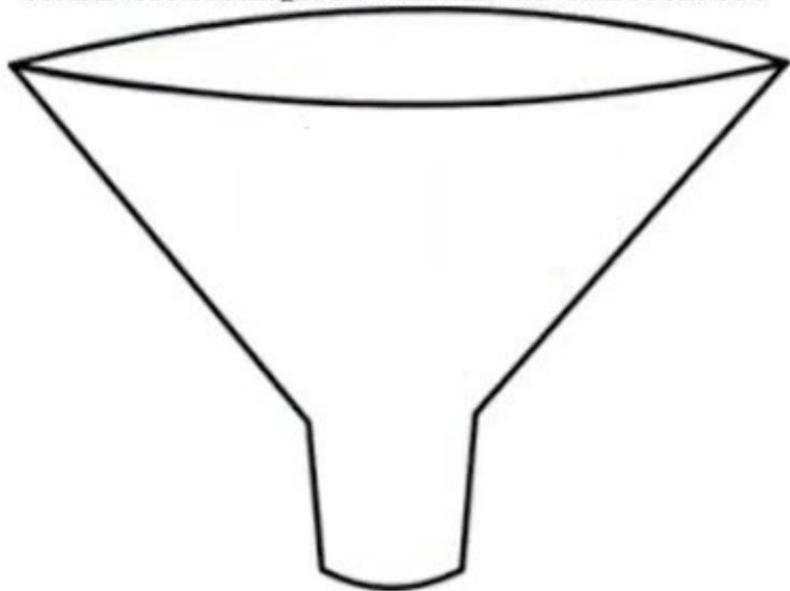
SIM NS
FOUNDATION

Pilar Cossio, Olga Troyanskaya, Miro Astore, Mihir Bafna, Siavash Golkar,
Bowen Jing, Mahsa Mofidi, Anand Ojha, Abhilash Sahoo, Lane Votapka



Outlook and Limitations

MLPGLALLLLAAWTARALEVPTDGNAGLLAEPQIAMFCGRL
NMHMNVQNGKWSDPSGTKTCIDTKEGILQYCQEVYPELQI
TNVVEANQPVTIQNWCKRGRKQCKTHPHFVIPYRCLVGEFV
SDALLVPDKCKFLHQERMDVCETHLHWHTVAKETCSEKSTN



RocketSHP efficiently predicts descriptors of protein dynamics

Model	RMSF		GCC-LMI		SHP	Model Information		
	RMSE	Spear. ρ (\uparrow)	GDD	IMSD	KL-Div	Seq.	Struct.	Params.
RocketSHP-mini	0.104	0.680	1.023	2.532	1.089	✓	✗	1.5M
RocketSHP-seq	0.090	0.719	0.944	2.387	1.225	✓	✗	29.7M
RocketSHP	0.083	0.789	0.859	1.898	1.551	✓	✓	30.2M
BioEmu (10 samples)	0.192	0.715	2.550	3.357	1.923	✓	✗	31.2M
BioEmu (100 samples)	0.202	0.748	1.215	1.871	1.778	✓	✗	31.2M
Dyna-1	0.361	0.267	-	-	-	✓	✓	188.9M
Dyna-1 (calibrated)	0.191	0.267	-	-	-	✓	✓	188.9M

NYGC Events

Michael_Brocidiacono_...



1



2



2



2



3



3



Look mom, no experimental data! Learning to score protein-ligand interactions from simulations

Michael Brocidiacono
University of North Carolina at Chapel Hill

September 11, 2025

MCBRLab **bilibili**

Michael_Brocidiacono

- Hide Tab Bar
- Hide Toolbar
- ✓ Thumbnails
- Table of Contents
- Highlights and Notes
- Bookmarks
- Contact Sheet
- ✓ Continuous Scroll
- Single Page
- Two Pages
- Soft Proof with Profile
- Show Document Background
- Adjust Size
- Zoom In
- Zoom Out
- Zoom to Selection
- Show Markup Toolbar
- Hide Toolbar
- Customize Toolbar...
- Slideshow**
- Enter Full Screen

1

2

3

Look_Mom Michael Brocidiacono.pdf

Michael_Brocidiacono_Look_Mom Michael Brocidiacono.pdf

Look mom, no experimental data! Learning to infer protein-ligand interactions from simulations

Michael Brocidiacono
University of North Carolina at Chapel Hill

September 11, 2025

1 / 18

The universe of drug targets is expanding

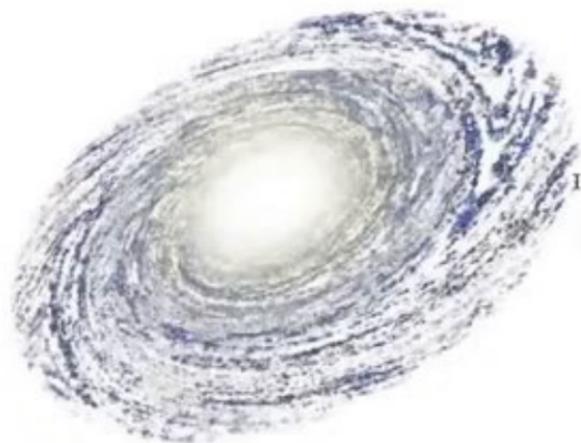
Look mom, no experimental data! Learning to score protein-ligand interactions from simulations

Michael Brocidiacono

University of North Carolina at Chapel Hill

September 11, 2025

The universe of drug targets is expanding



But the target to drug pipeline is broken

1

But the target to drug pipeline is broken



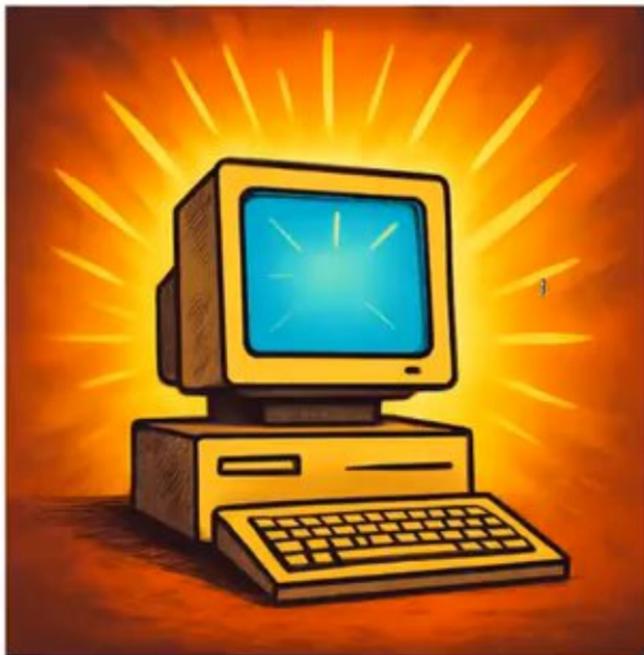
Finding small-molecule binders is hard

- High-throughput screening (HTS) is slow and expensive (and hits require a lot of optimization)



Can we computationally identify binders?

- Can pure AI solve our problems?



What's going on with free energy calculations?



RL



R



L

- Free energy is very spicy to calculate (10K dimensional integral).

I

$$\Delta G_{\text{binding}} = kT \log \frac{Z_{RL}}{Z_R Z_L}$$

$$Z = \int_{-\infty}^{\infty} e^{\frac{-U(x)}{kT}} dx$$

Combining ML and physics



I

Force matching

If we have a dataset of forces from MD simulations (pretty cheap), we can use a neural network to remove degrees of freedom from the system.

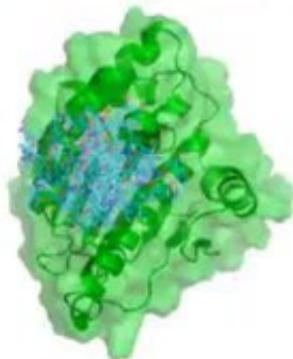
$$\text{Minimize} \left(-\frac{df_{NN}(x_S)}{dx_S} - \langle \mathbf{F}_{x_S} \rangle \right)^2$$

1

Ligand force matching (LFM) overview

a)

30-60K random
small molecules
(ZINC database)



i

Is it working?

1. Choose targets with many known binders in ChEMBL (w/ and w/o cocrystal poses)

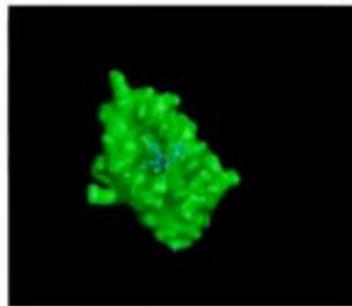
I

Benchmarking targets

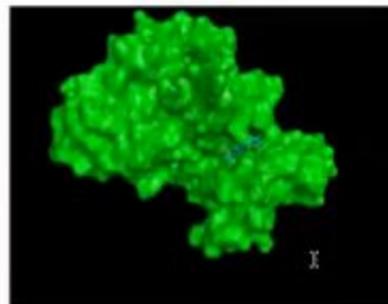
MCL1



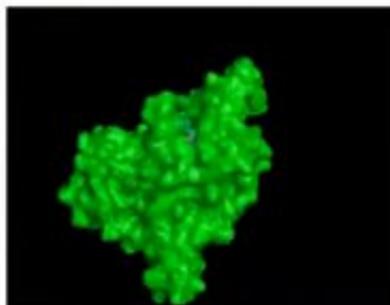
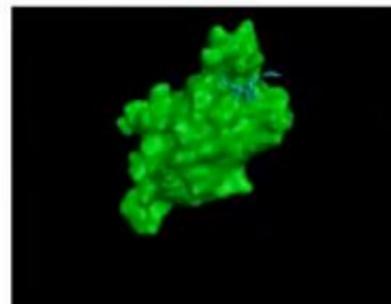
48,284 training datapoints



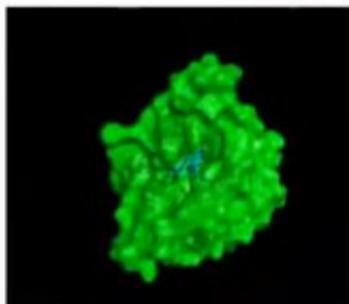
CDK2



BRD4



29,275 training datapoints



19,842 training datapoints

Results (crystal poses)

Model	EF_{\max}^B	AUC	% < 2 Å
-------	---------------	-----	---------

i

What's going on?



i

DD actives + UD decoys

Model	EF_{\max}^B	AUC
Vina	1.0 [1.0, 2.6]	0.23 ± 0.02
GNINA	7.9 [4.9, 17]	0.65 ± 0.02
LFM	13 [7.9, 31]	0.55 ± 0.04

i

Moving forward

- LFM is achieving promising results, though poor docking is holding us back



Acknowledgments



Grants

NIH R01GM140154

Tropsha Lab

Dr. Alexander Tropsha

Dr. Eugene Muratov

Enes Kelestemur

Dr. Holli-Joi Martin

Dr. James Wellnitz

Dr. Jon-Michael Beasley

Kelvin Idanwekhai

Kushal Koirala

Dr. Matthew Hart

Nyssa Tucker

Dr. Marcello DeLuca

Ricardo Teighi

Popov Lab

Dr. Konstantin Popov

Dr. James Wellnitz

Brandon Novoy

Rishabh Dey

Koes Lab

Dr. David Koes

Rishal Aggarwal

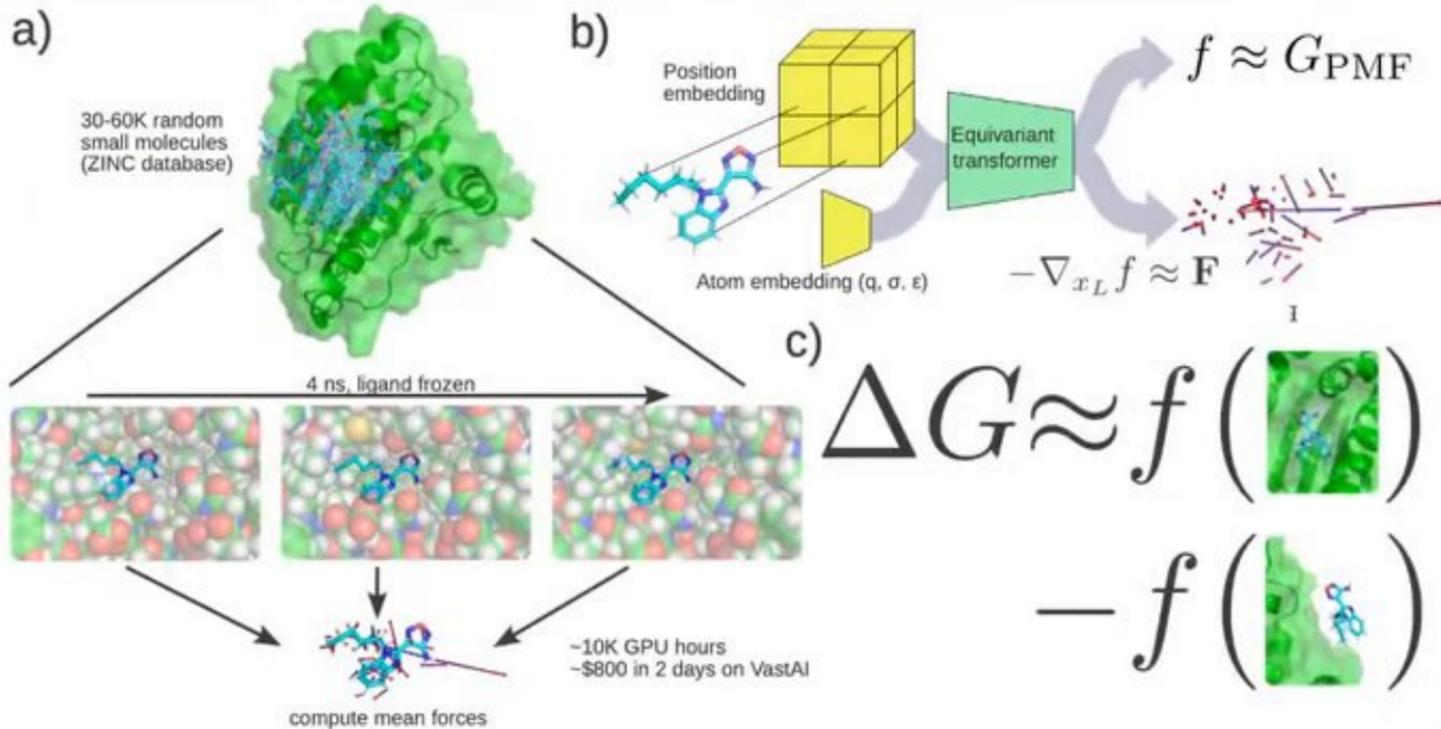
Dr. Paul Francoeur

Results (crystal poses)

Model	EF_{\max}^B	AUC	% < 2 Å
Vina (UD)	3.7 [1.0, 13]	0.36 ± 0.09	13 ± 5.0 %

I

Ligand force matching (LFM) overview



RocketSHP efficiently predicts descriptors of protein dy



Model

RocketSHP-m
RocketSHP-se
RocketSHP

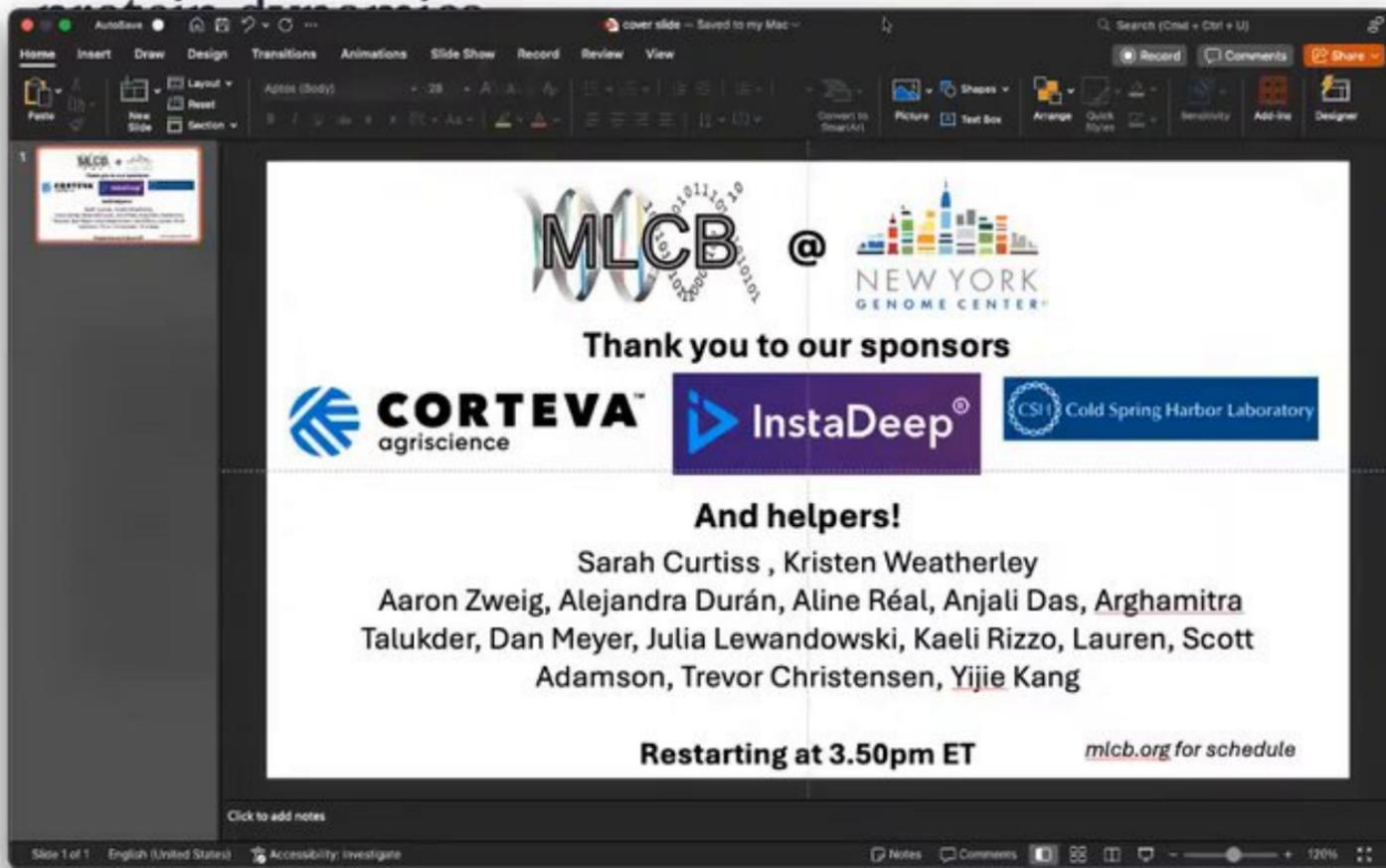
BioEmu (10 s
BioEmu (100 s

Dyna-1
Dyna-1 (calibr

Information
Params.
1.5M
29.7M
30.2M
31.2M
31.2M
188.9M
188.9M

2023.09.11

RocketSHP efficiently predicts descriptors of



cover slide -- Saved to my Mac --

Home Insert Draw Design Transitions Animations Slide Show Record Review View

Record Comments Share

Pages New Slide Section

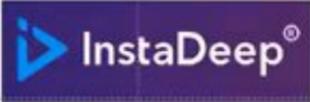
Address (Body) 28

Convert to BeamerArt Picture Text Box Arrange Quick Styles Sensitivity Add-ons Designer

1

MLCB @ NEW YORK GENOME CENTER

Thank you to our sponsors

And helpers!

Sarah Curtiss , Kristen Weatherley
Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 3.50pm ET mlcb.org for schedule

Click to add notes

Slide 1 of 1 English (United States) Accessibility: Investigate

Notes Comments

120%

Information

ct.	Params.
-----	---------

X	1.5M
---	------

X	29.7M
---	-------

✓	30.2M
---	-------

X	31.2M
---	-------

X	31.2M
---	-------

✓	188.9M
---	--------

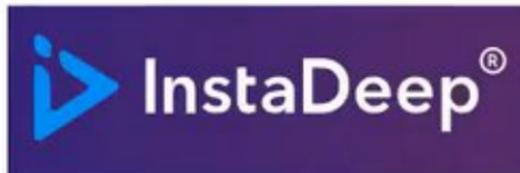
✓	188.9M
---	--------



@



Thank you to our sponsors



And helpers!

Sarah Curtiss , Kristen Weatherley

Aaron Zweig, Alejandra Durán, Aline Réal, Anjali Das, Arghamitra Talukder, Dan Meyer, Julia Lewandowski, Kaeli Rizzo, Lauren, Scott Adamson, Trevor Christensen, Yijie Kang

Restarting at 3.50pm ET

mlcb.org for schedule