

EPFL

EVENT

Week 11: Context-aware learning

■ École
polytechnique
fédérale
de Lausanne

VENUE & DATE

Course presentation
11/26/2025

PRESENTED BY

Jiying Zhang
First-year PHD
Patrick's Lab



- Zhang et al., [EquiPocket: an E\(3\)-Equivariant Geometric Graph Neural Network for Ligand Binding Site Prediction](#), ICML 2024
- Li et al., [Contextual AI models for single-cell protein biology](#), Nature Methods, 2024

- Zhang et al., EquiPocket: an E(3)-Equivariant Geometric Graph Neural Network for Ligand Binding Site Prediction, ICML 2024
- Li et al., Contextual AI models for single-cell protein biology, Nature Methods, 2024

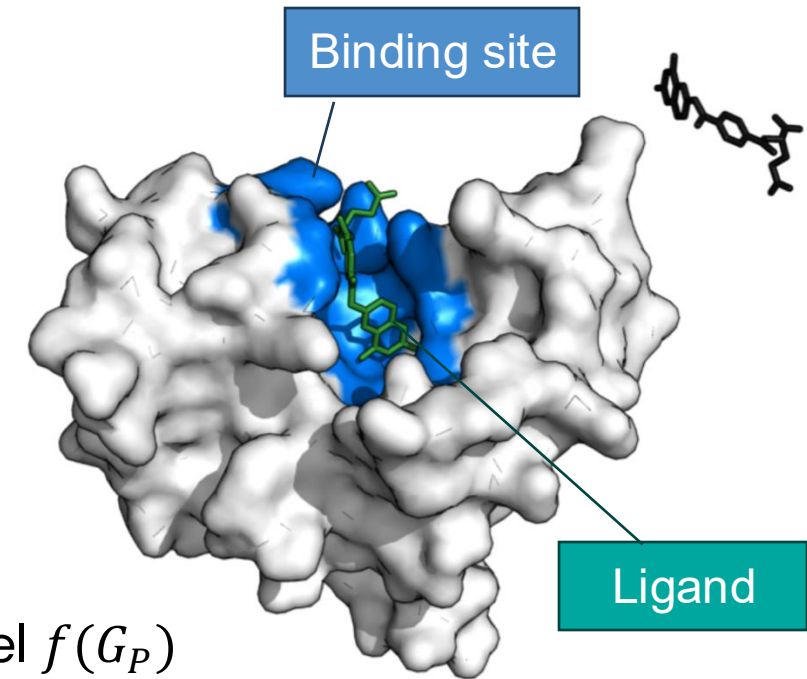
The Core Problem

Predicting the binding sites of target proteins.

It is a fundamental step in modeling drug discovery

- Protein-ligand docking tasks.
- Structure-based molecular generation
- Understanding biological interactions in living systems

Problem Statement. Given a protein $G_P = (\mathcal{V}_p, E_c, E_D)$, the model $f(G_P)$ is required to predict the atoms of **binding sites**: $\mathcal{V}_B \in \mathcal{V}_p$.



Geometry-Based

Examples: Fpocket, LigSite

Exploit hand-crafted geometric algorithms to find hollow spaces. Limited expressivity

Machine Learning

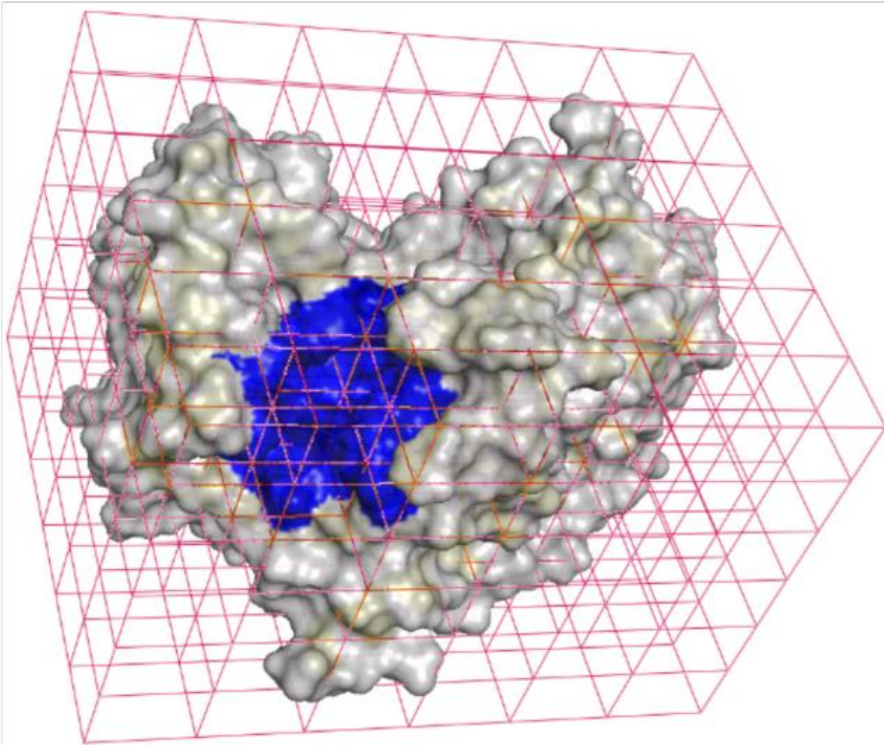
Examples: P2Rank

Use Random Forest with surface geometry features. Still relies on hand-crafted descriptors

Deep Learning (CNN)

Examples: DeepSite, DeepSurf

Treats protein as a 3D image (voxelization). Currently the dominant approach but has flaws



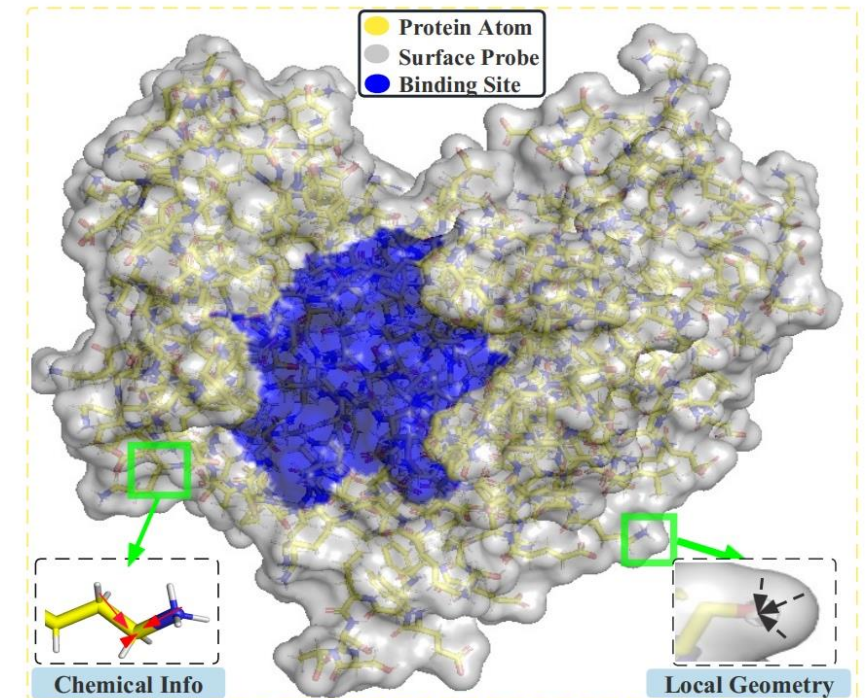
(a) CNN-based methods

- **Irregular Structures:** Regular voxels struggle to model the naturally shape of proteins. (fixed box)
- **Rotation Sensitivity:** Rotation changes the input, confusing the model.
- **Surface Coarseness:** Voxel grids are too coarse to capture the fine-grained geometry of surface atoms.
- **Size Shift Unawareness:** Fixed grid sizes often cut off large proteins or waste space for small ones.

Key Innovations

EquilPocket is an $E(3)$ -equivariant GNN that addresses previous limitations

- **Graph-based:** naturally handles irregular protein
- **Equivariant:** Insensitive to rotation and translation
- **Surface-centric:** Specifically models surface geometry using probes
- **Adaptive:** Uses dense attention to handle variable protein sizes



Protein Graphs

- Protein Graph $G_P = (V_P, E_C, E_D)$: Nodes V_P are atoms, edges E_C are chemical bonds and E_D are spatial neighbors

Surface Probe Set

- A set of 3D coordinates generated by **rolling a solvent probe** (e.g., water) across the protein surface (alg. **MSMS**)
 - $\mathbf{x}_i \in \mathbb{R}^3$: Its 3D coordinate.
 - $p_i \in V_P$: The index of its **Nearest Protein Atom**.

Surface Graph

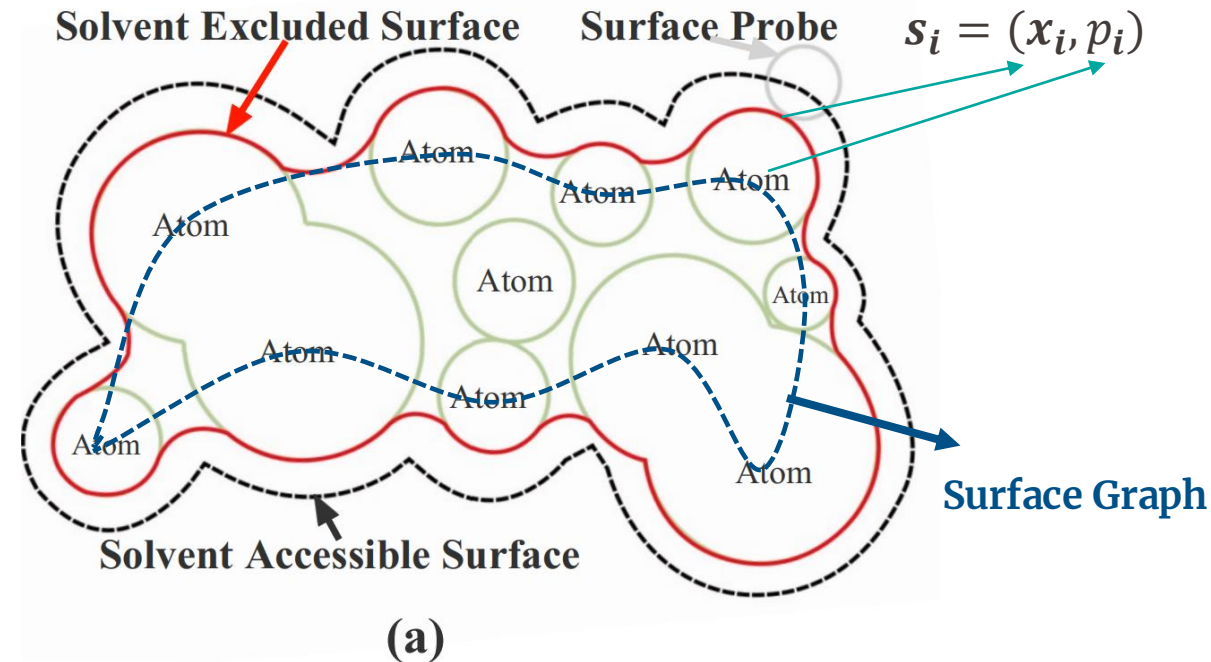
- $G_S = (V_S, E_S)$: $V_S = \{p_i\}$ (Nearest Protein Atoms).

E(3)-Equivariance

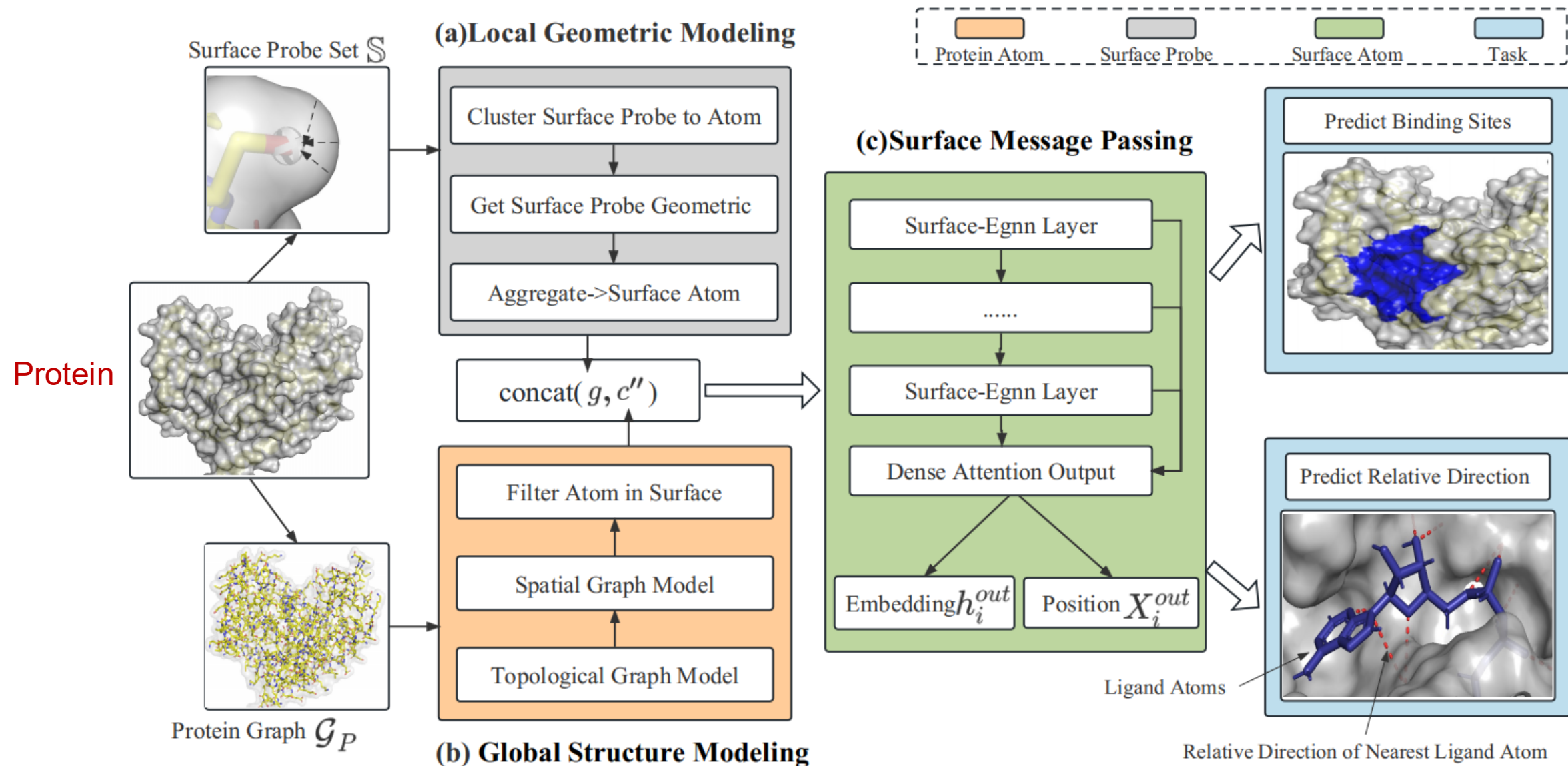
- A function f is equivariant if transforming the input transforms the output similarly:

$$f(g \cdot X) = g \cdot f(X)$$

- Invariant: $f(g \cdot X) = f(X)$



Three modules: Local Geometric Modeling, Global Structure, Surface Message Passing Modeling



Extracting Surface Features (context)

- Collect surrounding probes $S_i = \{s_j = (\mathbf{x}_j, p_j) \in S \mid p_j = i\}$ for each surface atom $i \in V_S$.

Geometric descriptor g_i

- For each surrounding surface probe $s_j \in S_i$, we first search its two nearest surface probes from S as s_{j_1} and s_{j_2} , $\bar{\mathbf{x}}_i = \text{mean}(S_i)$

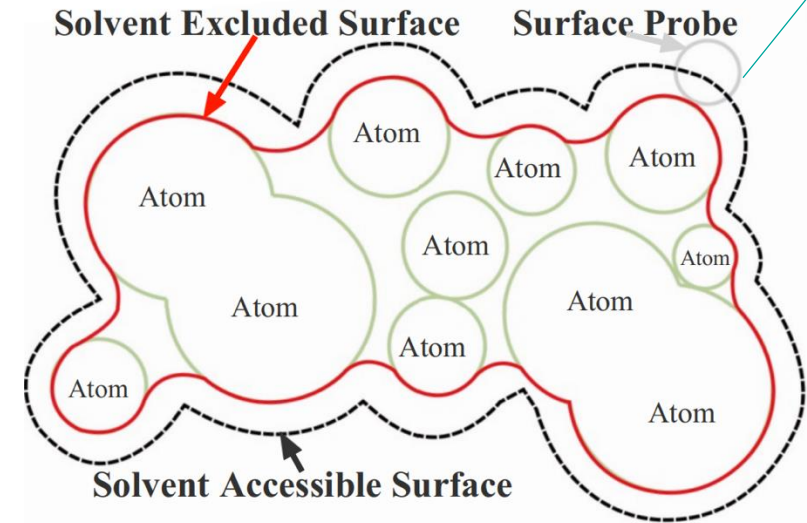
$$\mathbf{x}_{jj_1} = \mathbf{x}_j - \mathbf{x}_{j_1}, \mathbf{x}_{jj_2} = \mathbf{x}_j - \mathbf{x}_{j_2}, \mathbf{x}_{j,\text{center}} = \mathbf{x}_j - \bar{\mathbf{x}}_i,$$

$$\mathbf{x}_{j,\text{protein}} = \mathbf{x}_j - \mathbf{x}_i, \mathbf{x}_{\text{center},\text{protein}} = \bar{\mathbf{x}}_i - \mathbf{x}_i.$$

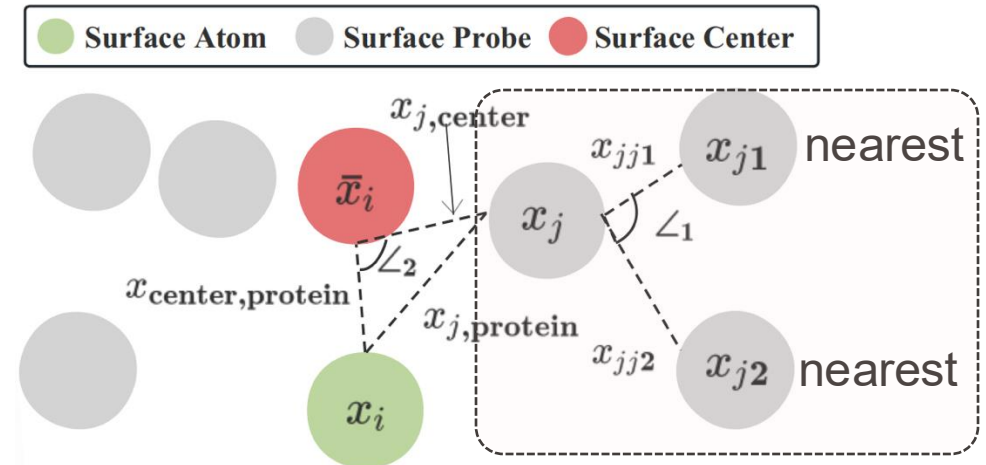
$$g(s_j) := [\|\mathbf{x}_{jj_1}\|_2, \|\mathbf{x}_{jj_2}\|_2, \angle_1, \|\mathbf{x}_{j,\text{center}}\|_2, \|\mathbf{x}_{j,\text{protein}}\|_2, \|\mathbf{x}_{\text{center},\text{protein}}\|_2, \angle_2],$$

$$\angle_1 = \frac{\mathbf{x}_{jj_1} \cdot \mathbf{x}_{jj_2}}{\|\mathbf{x}_{jj_1}\|_2 \|\mathbf{x}_{jj_2}\|_2} \quad \angle_2 = \frac{\mathbf{x}_{j,\text{center}} \cdot \mathbf{x}_{\text{center},\text{protein}}}{\|\mathbf{x}_{j,\text{center}}\|_2 \|\mathbf{x}_{\text{center},\text{protein}}\|_2}$$

$$g_i = [\text{Pooling}(\{\text{MLP}(g(s_j))\}_{s_j \in S_i}), \text{MLP}(\text{Pooling}(\{g(s_j)\}_{s_j \in S_i}))].$$



(a)



(b)

Compute relative distances and angles between probes

Chemical-Graph Modeling

- Uses GNN to process chemical features (atom types) and chemical bonds

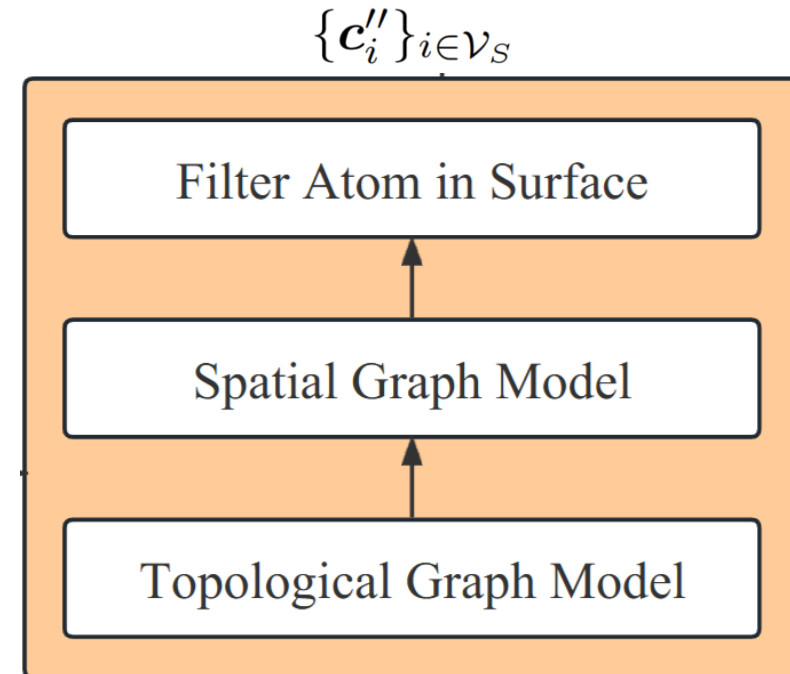
$$\{c'_i\}_{i \in \mathcal{V}_P} = \text{GNN}(\{c_i\}_{i \in \mathcal{V}_P}, \mathcal{E}_C),$$

Spatial-Graph Modeling

- Uses EGNN to incorporate 3D coordinates x_i and spatial edges.

$$\{c''_i\}_{i \in \mathcal{V}_P} = \text{EGNN}(\{x_i, c'_i\}_{i \in \mathcal{V}_P}, \mathcal{E}_D).$$

- This captures the global structural context of the protein.



We perform equivariant message passing exclusively on the **Surface Graph** to refine features.

Surface-EGNN Layer

- Input: Invariant features $h^{(0)} = [c'', g]$ and coordinates $\mathbf{X}_i(0) = [\mathbf{x}_i, \bar{\mathbf{x}}_i]$.

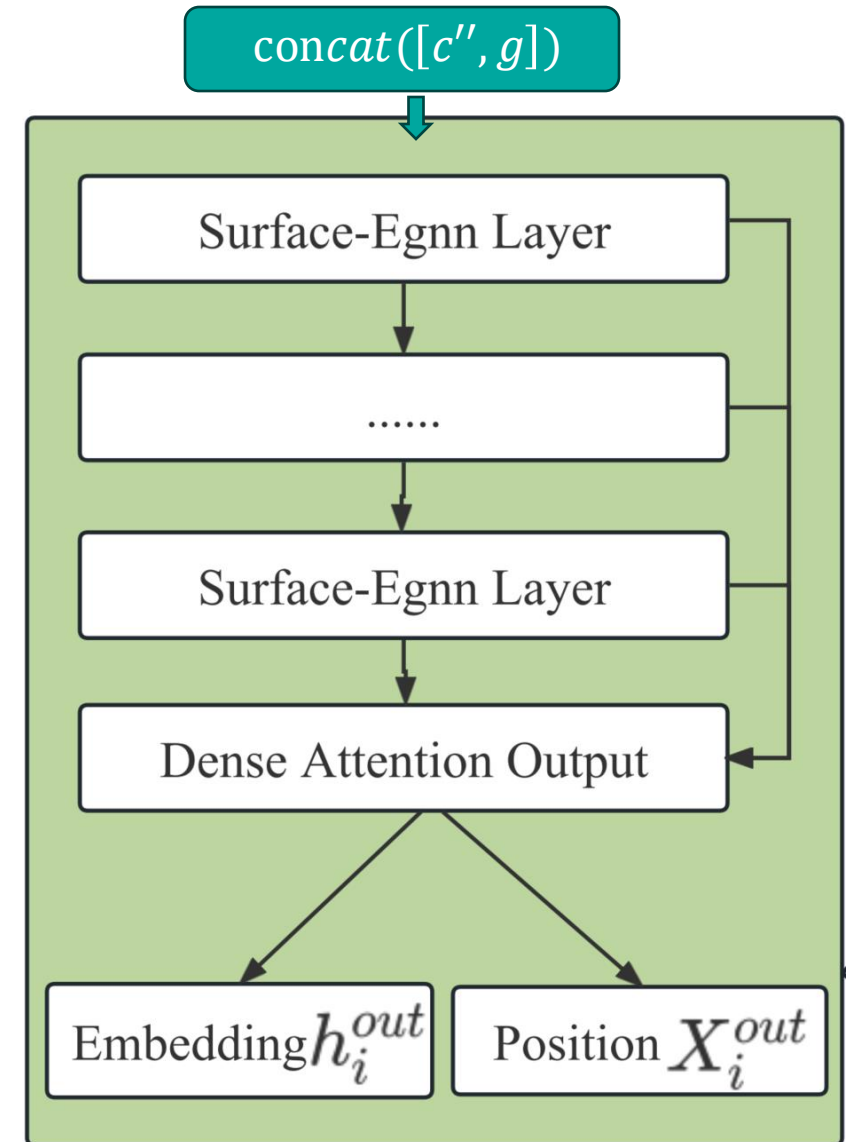
$$\mathbf{m}_{ij} = \phi_m \left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)}, f_x(\mathbf{X}_i^{(l)}, \mathbf{X}_j^{(l)}), e_{ij} \right),$$

$$\mathbf{h}_i^{(l+1)} = \phi_h \left(\mathbf{h}_i^{(l)}, \sum_{j \in \mathcal{N}(i)} \mathbf{m}_{ij} \right),$$

$$\mathbf{X}_i^{(l+1)} = \mathbf{X}_i^{(l)} + \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} (\mathbf{x}_{i,1}^{(l)} - \mathbf{x}_{j,1}^{(l)}) \phi_x(\mathbf{m}_{ij}),$$

where ϕ_m, ϕ_h, ϕ_x are MLPs, and

$$f_x(\mathbf{X}_i, \mathbf{X}_j) = [\|\mathbf{x}_{ij}\|_2, \|\mathbf{x}_{ci}\|_2, \|\mathbf{x}_{cj}\|_2, \angle_{ci,ij}, \angle_{cj,ij}, \angle_{ci,cj}],$$



Challenge: A fixed receptive field causes over-smoothing in small proteins and insufficient coverage in large ones.

Dense Attention Output Layer.

1. Calculate the density of neighbors at different hop distances n_i

$$n_i^{(l)} = \frac{|\{j \in \mathcal{V}_P \mid 0 \leq \|\mathbf{x}_i - \mathbf{x}_j\|_2 < l\theta\}|}{N_P},$$

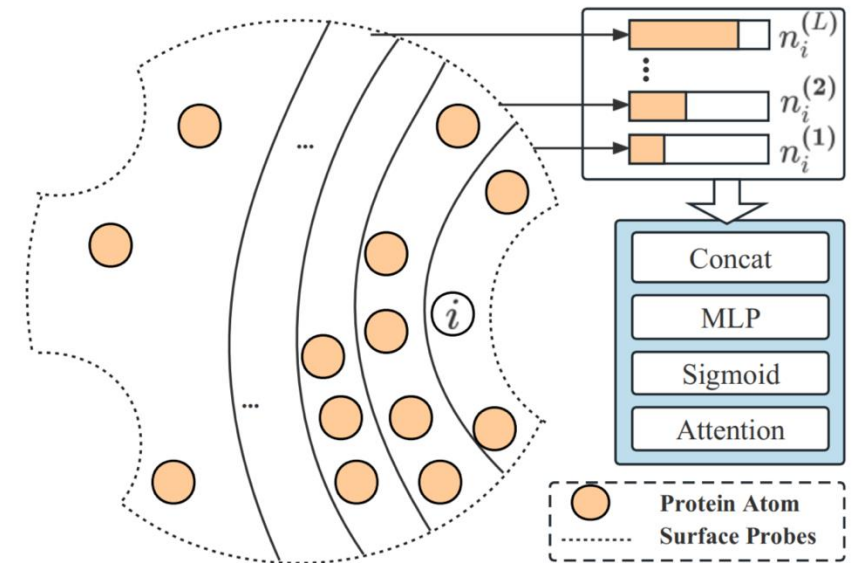
$$\mathbf{n}_i = [n_i^{(0)}, n_i^{(1)}, \dots, n_i^{(L)}, N_P] \in \mathbb{R}^{L+2}$$

2. An MLP generates attention weights a_i to adaptively balance information from different GNN layers.

$$\mathbf{a}_i = \text{Sigmoid}(\phi_a(\mathbf{n}_i))$$

$$\mathbf{h}_i^{\text{out}} = \text{Concat}(a_{i0}\mathbf{h}_i^{(0)}, \dots, a_{iL}\mathbf{h}_i^{(L)}),$$

$$\mathbf{X}_i^{\text{out}} = \frac{1}{L+1} \sum_{l=0}^L \mathbf{X}_i^{(l)},$$



Binding Site Prediction

- **Binary Classification:** label $y = 1$: surface atom i is within 4Å of a ligand.
$$\hat{y}_i = \mathbf{Sigmoid}(\text{MLP}(\mathbf{h}_i^{\text{out}}))$$
- **Loss:** Dice Loss (L_b) handles class imbalance effectively.

Relative Direction

- **Auxiliary Task:** Predict the unit vector direction d_i from atom $\mathbf{x}_i \in G_p$ to the nearest ligand atom m_i .
$$\mathbf{d}_i = \frac{\mathbf{m}_i - \mathbf{x}_i}{\|\mathbf{m}_i - \mathbf{x}_i\|_2}, \quad \hat{\mathbf{d}}_i = \frac{\mathbf{x}_i^{\text{out}} - \mathbf{x}_i}{\|\mathbf{x}_i^{\text{out}} - \mathbf{x}_i\|_2}.$$
- **Loss:** Cosine Similarity Loss (L_d)

Total Loss

$$L = L_b + L_d$$

- Combines semantic segmentation with geometric regression

Datasets

Dataset	Role	Number of Binding Sites	Description
scPDB	Training	1	~16k entries, curated binding sites
COACH420	Testing	1,2,3,4,5...	Standard benchmark (mli g subset)
HOLO4K	Testing	1,2,3,4,5...	Larger, more complex proteins

Evaluation Metrics

- **DCC**: Distance between the predicted and true binding site center. (Success if $< 4\text{\AA}$)
- **DCA**: Shortest distance from predicted center to any ligand atom

Failure Rate: % of proteins without any binding site center

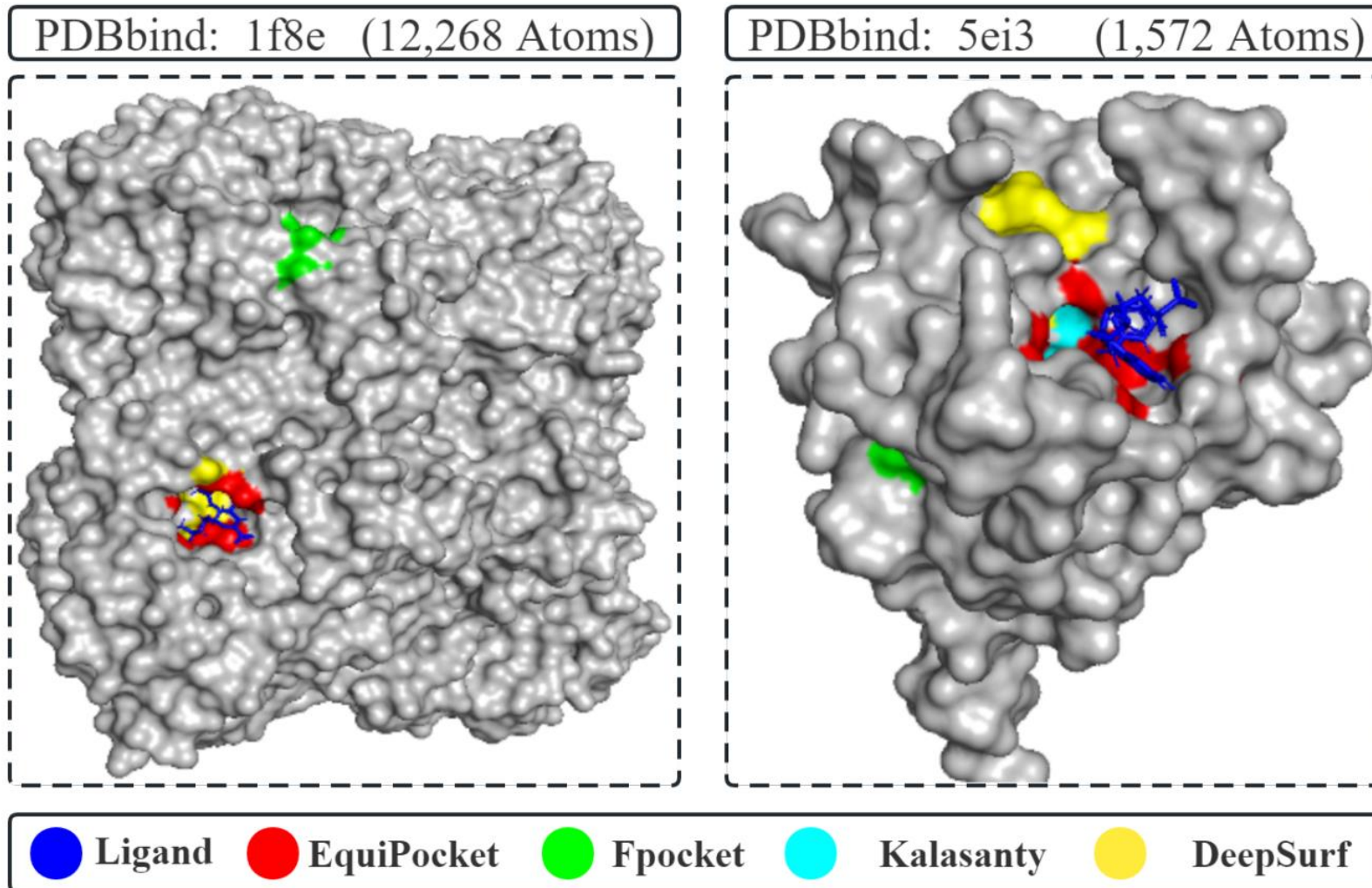
$$\text{Success Rate(DCC/DCA)} = \frac{|\{\text{Predicted sites} \mid \text{DCC/DCA} < \overset{4\text{\AA}}{\text{threshold}}\}|}{|\{\text{True sites}\}|}$$

$$\text{Failure Rate} = \frac{|\{\text{Protein} \mid \text{predicted binding center} = \emptyset\}|}{|\{\text{Protein}\}|}$$

EquilPocket outperforms geometry-based, ML-based, and CNN-based methods

Methods	Type	Param (M)	Failure Rate ↓	COACH420		HOLO4K		PDBbind2020	
				DCC↑	DCA↑	DCC↑	DCA↑	DCC↑	DCA↑
Fpocket ^b	Geometric-based	\	0.000	0.228	0.444	0.192	0.457	0.253	0.371
P2rank ^b	Machine-learning	\	0.000	0.366	0.628	0.314	0.621	0.503	0.677
DeepSite ^b	CNN-based	1.00	\	\	0.564	\	0.456	\	\
Kalasanty ^b		70.6	0.120	0.335	0.636	0.244	0.515	0.416	0.625
DeepSurf ^b		33.1	0.054	0.386	0.658	0.289	0.635	0.510	0.708
RecurPocket ^b		21.2	0.075	0.354	0.593	0.277	0.616	0.492	0.663
GAT	Topological Graph	0.03	0.110	0.039(0.005)	0.130(0.009)	0.036(0.003)	0.110(0.010)	0.032(0.001)	0.088(0.011)
GCN		0.06	0.163	0.049(0.001)	0.139(0.010)	0.044(0.003)	0.174(0.003)	0.018(0.001)	0.070(0.002)
GCN2		0.11	0.466	0.042(0.098)	0.131(0.017)	0.051(0.004)	0.163(0.008)	0.023(0.007)	0.089(0.013)
SchNet	Spatial	0.49	0.140	0.168(0.019)	0.444(0.020)	0.192(0.005)	0.501(0.004)	0.263(0.003)	0.457(0.004)
Egnn	Graph	0.41	0.270	0.156(0.017)	0.361(0.020)	0.127(0.005)	0.406(0.004)	0.143(0.007)	0.302(0.006)
EquiPocket-L	Ours	0.15	0.552	0.070(0.009)	0.171(0.008)	0.044(0.004)	0.138(0.006)	0.051(0.003)	0.132(0.009)
EquiPocket-G		0.42	0.292	0.159(0.016)	0.373(0.021)	0.129(0.005)	0.411(0.005)	0.145(0.007)	0.311(0.007)
EquiPocket-LG		0.50	0.220	0.212(0.016)	0.443(0.011)	0.183(0.004)	0.502(0.008)	0.274(0.004)	0.462(0.005)
EquiPocket		1.70	0.051	0.423(0.014)	0.656(0.007)	0.337(0.006)	0.662(0.007)	0.545(0.010)	0.721(0.004)

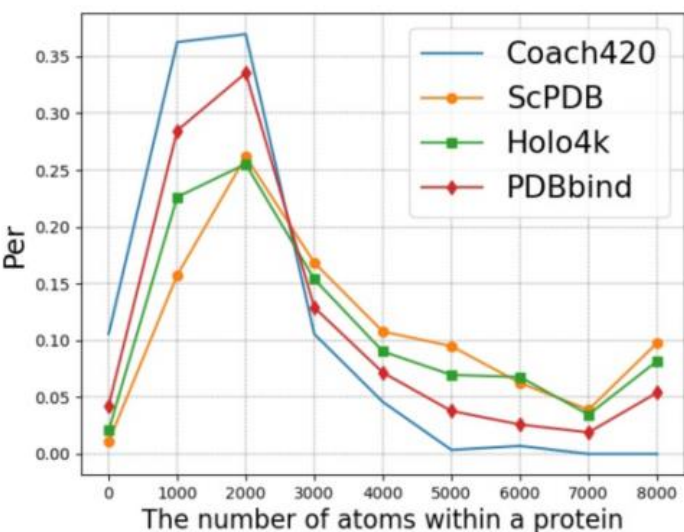
EquilPocket outperforms geometry-based, ML-based, and CNN-based methods



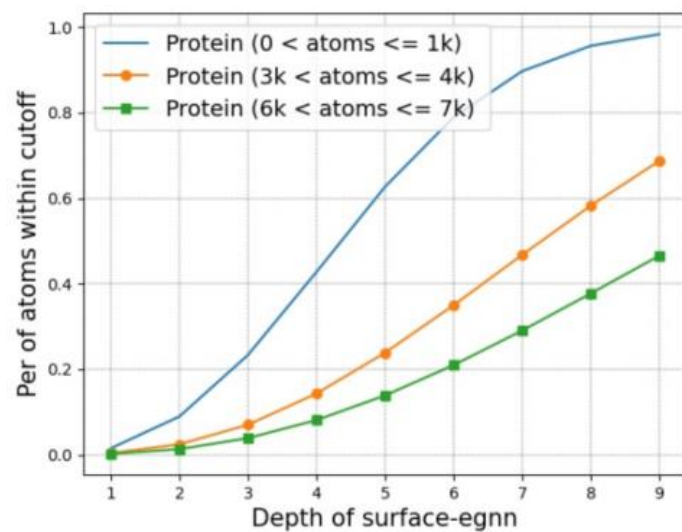
EquiPocket outperforms geometry-based, ML-based, and CNN-based methods

Methods	Type	Param (M)	Failure Rate ↓	COACH420		HOLO4K		PDBbind2020	
				DCC↑	DCA↑	DCC↑	DCA↑	DCC↑	DCA↑
EquiPocket-L	Ours	0.15	0.552	0.070(0.009)	0.171(0.008)	0.044(0.004)	0.138(0.006)	0.051(0.003)	0.132(0.009)
EquiPocket-G		0.42	0.292	0.159(0.016)	0.373(0.021)	0.129(0.005)	0.411(0.005)	0.145(0.007)	0.311(0.007)
EquiPocket-LG		0.50	0.220	0.212(0.016)	0.443(0.011)	0.183(0.004)	0.502(0.008)	0.274(0.004)	0.462(0.005)
EquiPocket		1.70	0.051	0.423(0.014)	0.656(0.007)	0.337(0.006)	0.662(0.007)	0.545(0.010)	0.721(0.004)

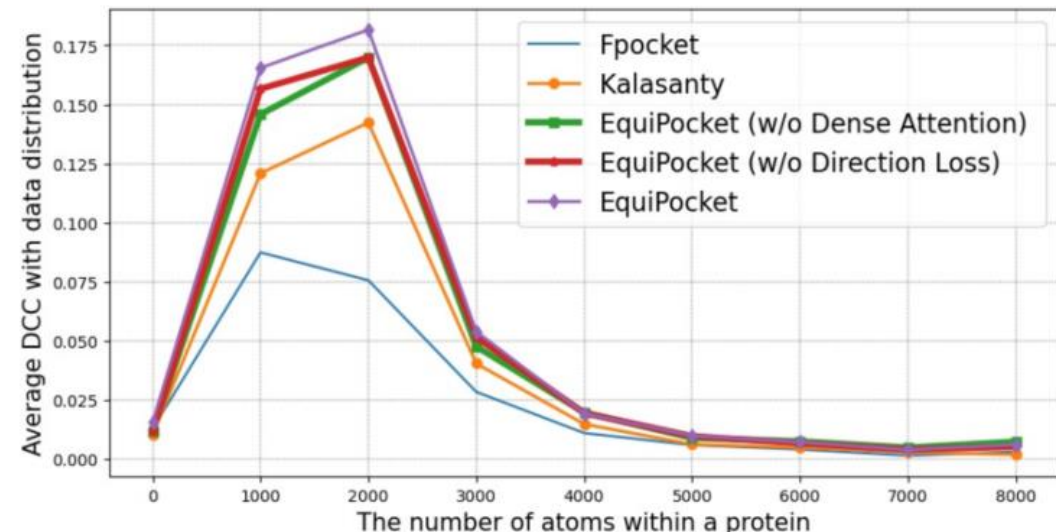
The protein size shift and model performances for proteins of various sizes



(a)



(b)



(c)

EquiPocket is an **E(3)-Equivariant** Graph Neural Network (GNN) that leverages a **dual-graph architecture** (atomic and surface) to integrate **geometric context** for rotation-invariant **ligand binding site prediction**.

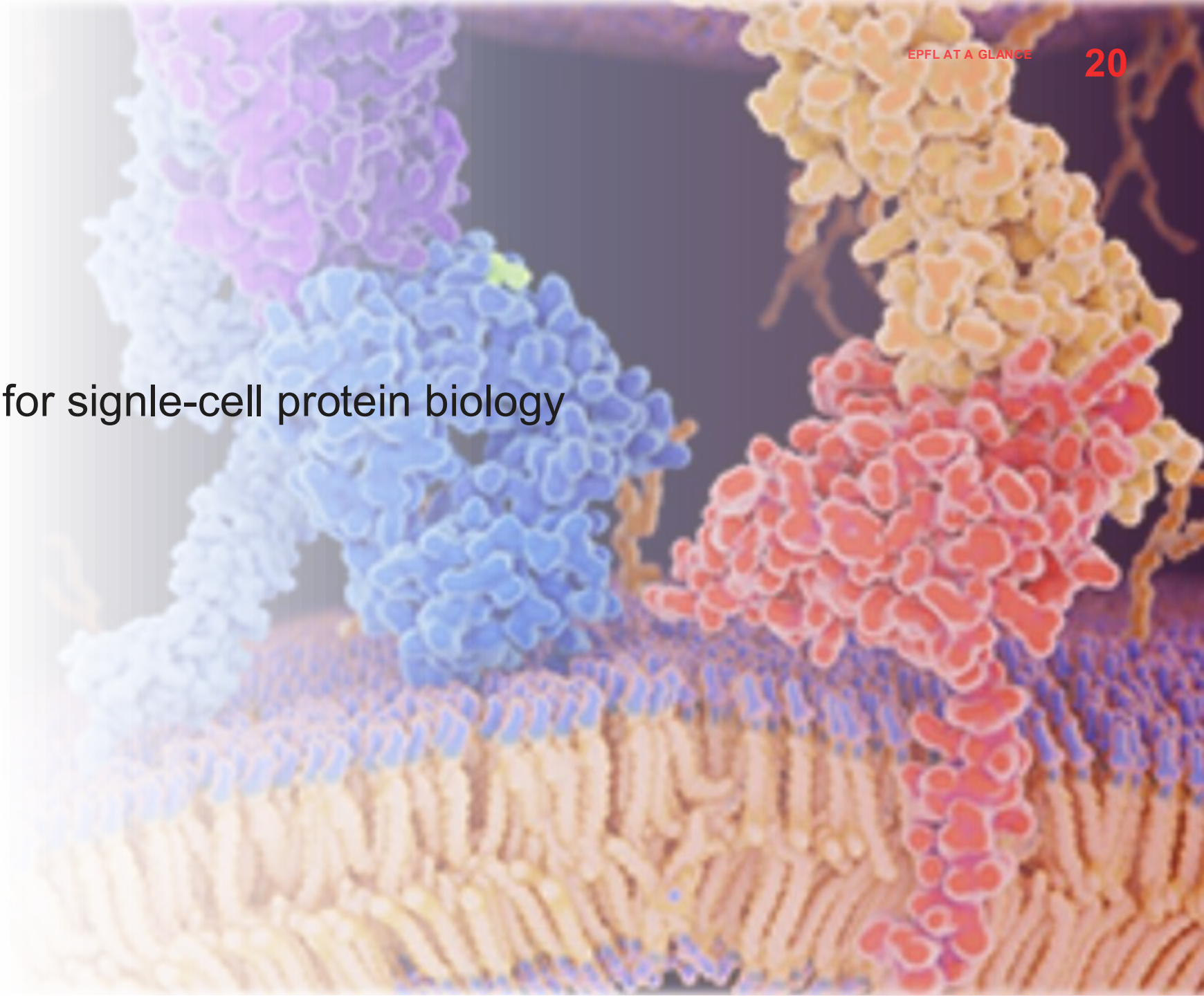
Positives:

- **E(3) Equivariance:** Guarantees rotation/translation invariance for the final prediction and consistency of feature learning, improving generalization.
- **Geometric Context:** Effectively models the protein surface geometry to provide crucial context for binding site prediction.

Limitations:

- **Scalability Issue:** Performance can deteriorate or become inefficient for very large proteins due to the sheer size of the surface probe set and complex graph message passing.

- EquiPocket:
- Contextual AI models for single-cell protein biology

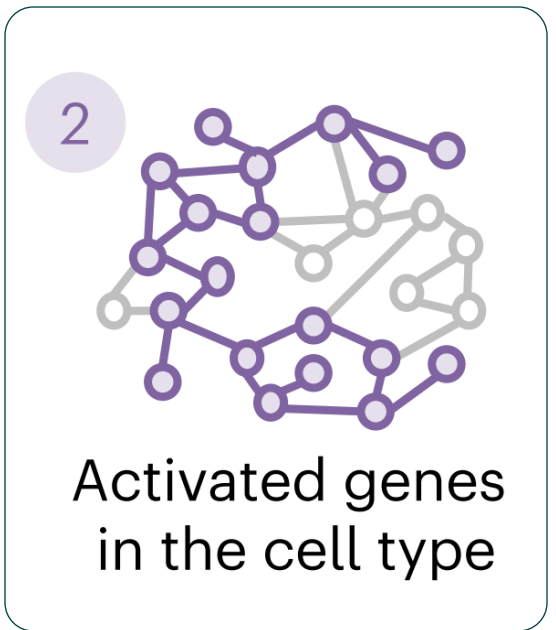
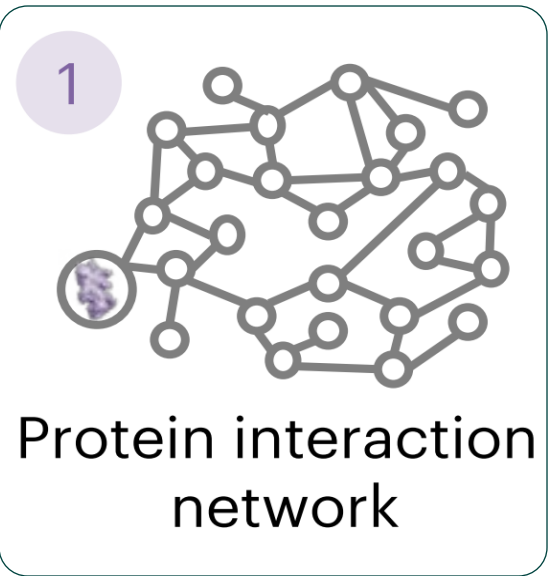
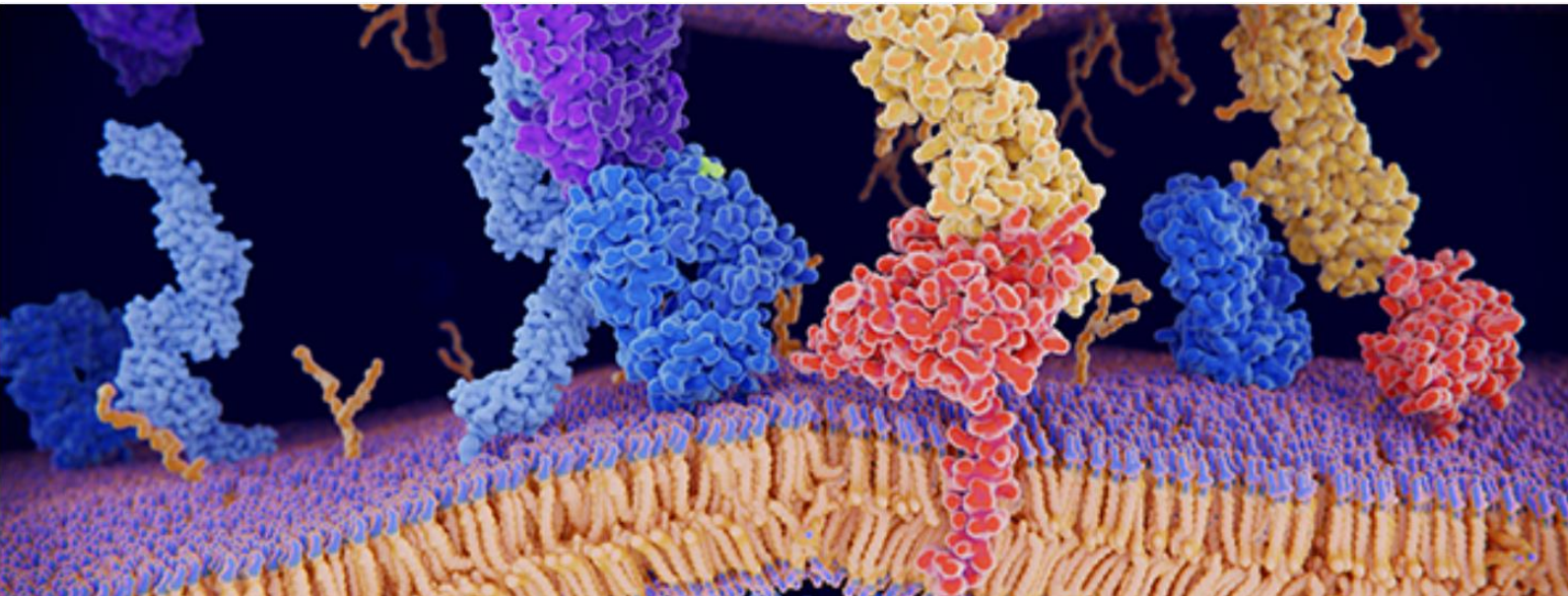


EPFL Background—protein in the biological contexts²¹

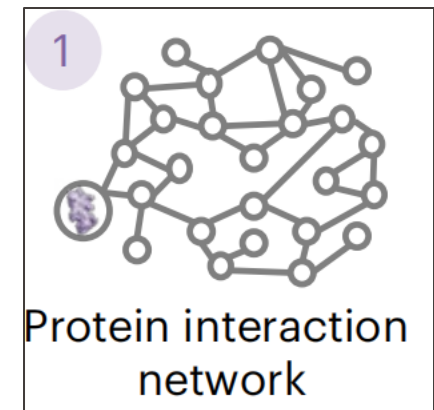
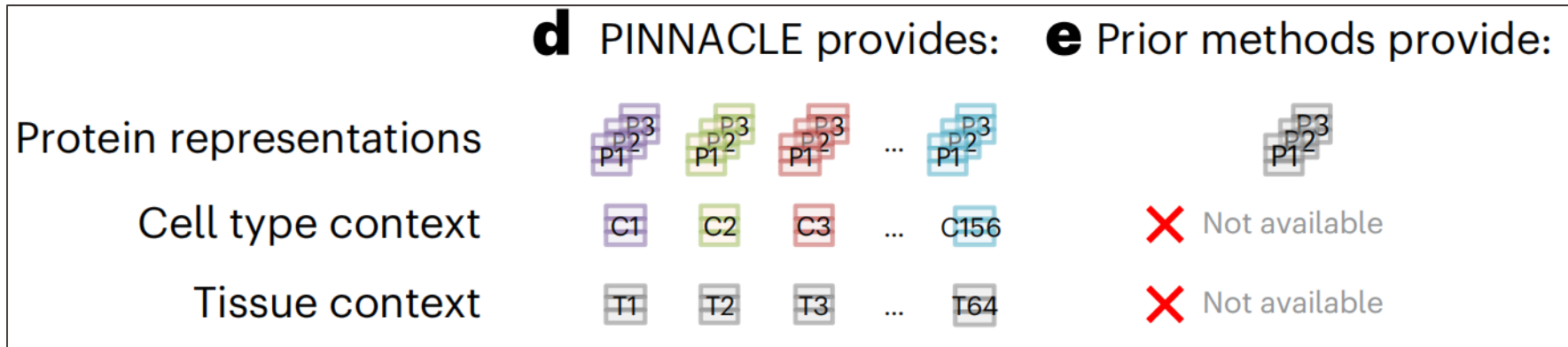
Importance of biological contexts (e.g. cellular and tissue)

Understanding protein function and developing molecular therapies require deciphering the **cell types** in which proteins act as well as the interactions between proteins.

- The expression of genes and the function of proteins encoded by these genes depend on **cellular and tissue contexts**
- Protein functions and interactions will be different in various biological contexts.
- Modeling protein interactions across biological contexts remains challenging for existing algorithms



- **Existing works:** Generate one representation for each protein, not tailored to specific biological contexts (cell types and disease states) (e.g. GAT)
- **Weakness:** These representations cannot identify protein functions that vary across different cell types, which in turn hamper the prediction of protein roles in a cell type-specific manner.



Single-Cell Atlases

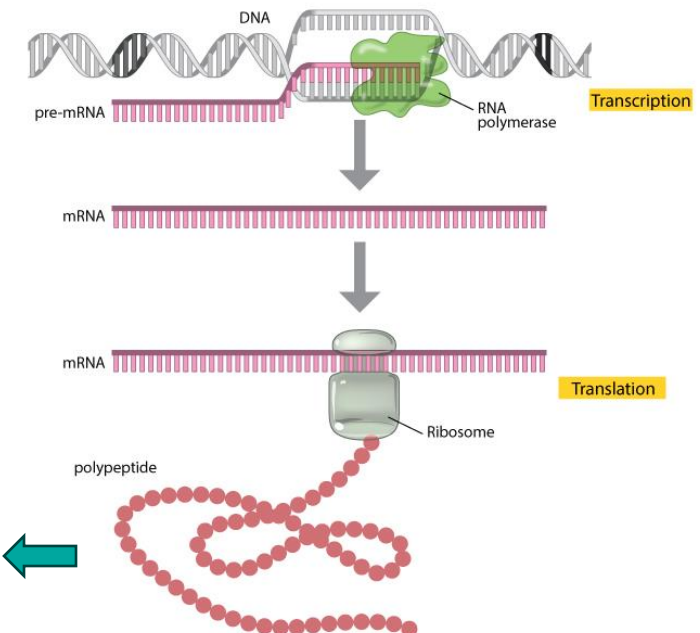
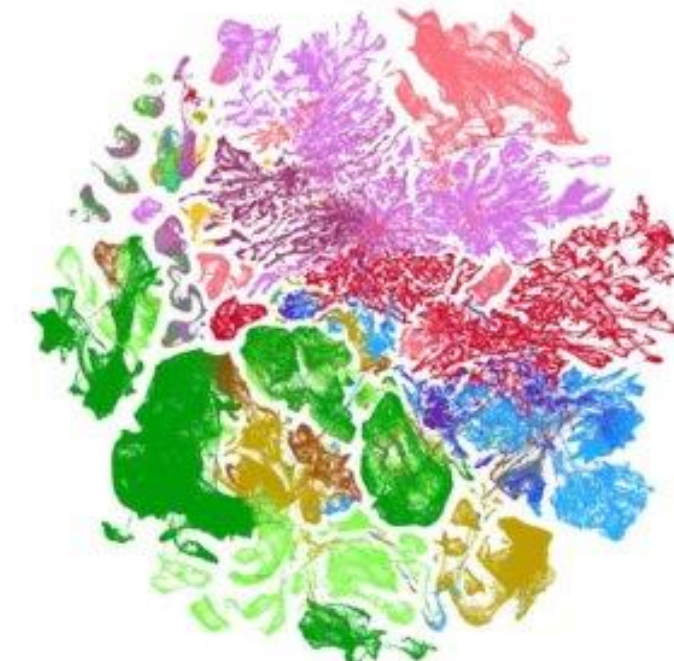
We are in the era of **single-cell transcriptomics**. Resources like *Tabula Sapiens* provide a high-resolution "parts list" of activated genes across hundreds of unique cell types and tissues.

The Missing Link

While we have the **schematic** (protein networks) and the **inventory** (single-cell atlases), they have historically remained separate.

The **challenge** lies in effectively merging these massive, disparate datasets to inform protein models.

Single-cell atlas



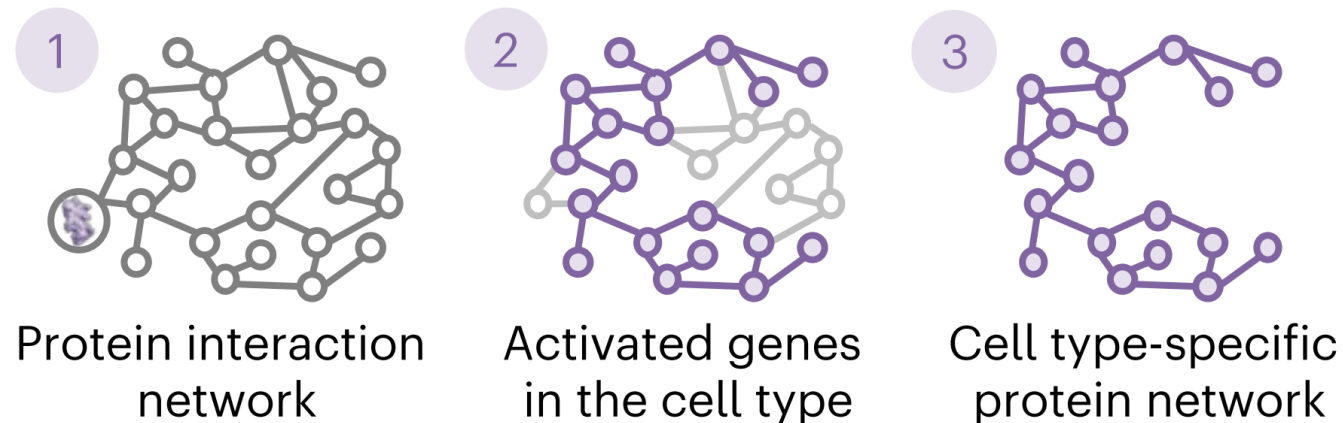
PINNACLE (Protein Network-based Algorithm for Contextual Learning)

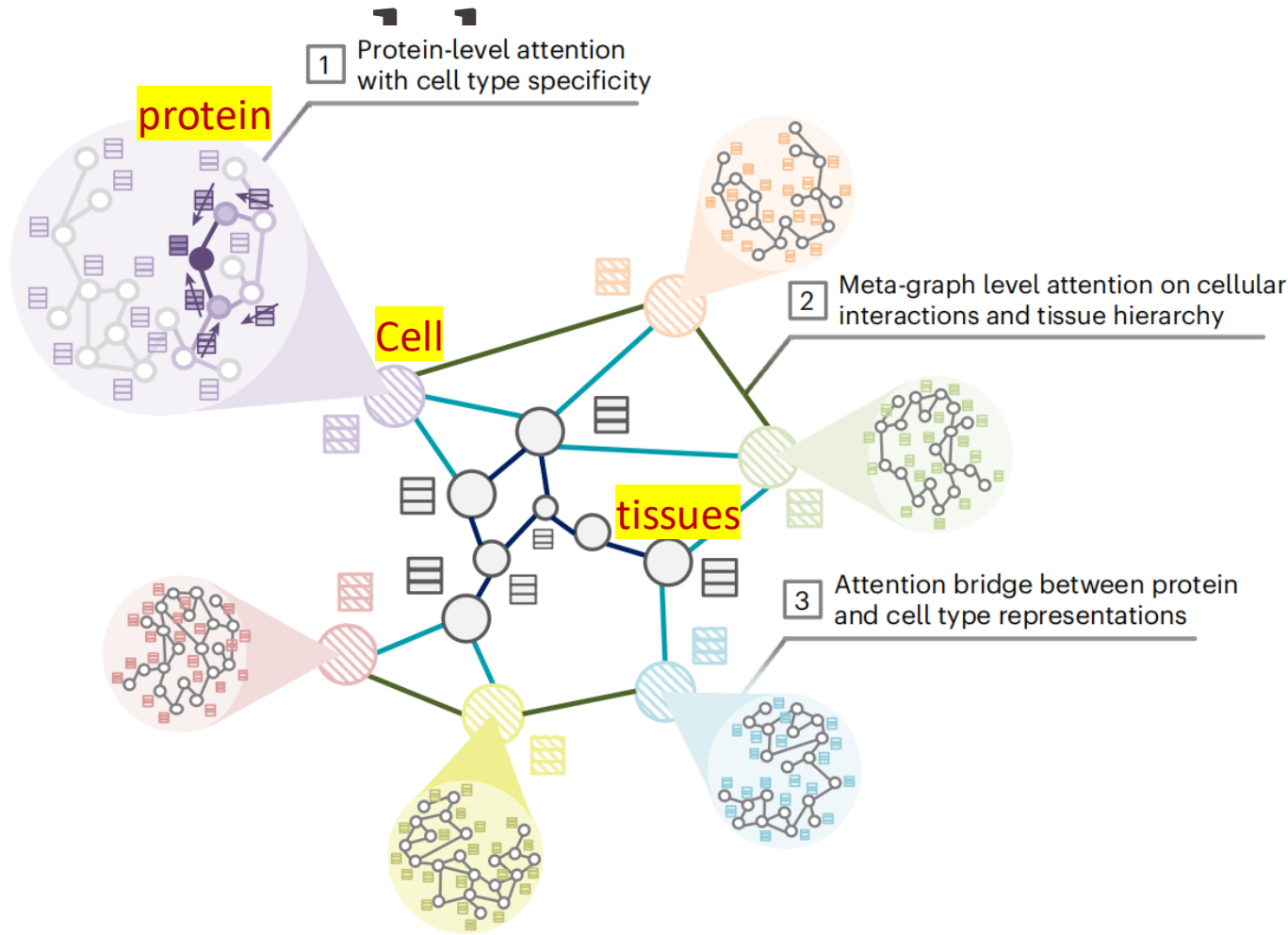
Geometric Deep Learning

PINNACLE integrates single-cell transcriptomics, protein-protein interaction networks, and tissue hierarchies. It moves beyond static models to generate high-resolution protein representations.

Distinct Cell-Type Specificity

Instead of one single vector **per protein**, PINNACLE generates **394,760 distinct representations**—one for every cell type in which a protein is active. This allows for precise, cell-type-specific therapeutic targeting.





➤ **PPI**

Cell type-specific PPI (Protein-protein interaction) networks $\{G_i\}$

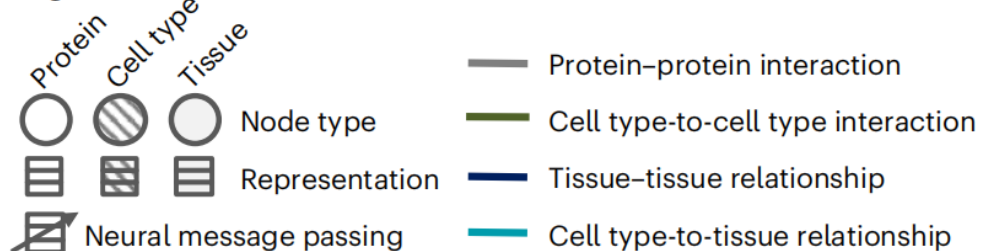
- Edges: Protein-Protein interaction

➤ **Metagraph:**

Cell types and tissues form a network.

- Edges
 - cell type-cell type interactions
 - cell type-tissue associations
 - Tissue-tissue relationships

Legend



1. Protein-level attention with cell type specificity.

$$\mathbf{h}_u^{\text{PP}} \leftarrow \text{AGG} \left(\sigma \left(\sum_{v \in \mathcal{N}_u} \alpha_{u,v} W^{\text{PP}} \mathbf{h}_v^{\text{PP}} \right) \right) \quad \alpha_{u,v} = \frac{\exp(\sigma(\mathbf{a}^T \cdot [\mathbf{h}_u \| \mathbf{h}_v]))}{\sum_{v \in \mathcal{N}_u} \exp(\sigma(\mathbf{a}^T \cdot [\mathbf{h}_u \| \mathbf{h}_v]))}$$

2. Metagraph-level attention on cellular interactions and tissue hierarchy.

$$\mathbf{h}_{c_i}^{\text{CC}} \leftarrow \text{AGG} \left(\sigma \left(\sum_{c \in \mathcal{N}_{c_i}} \alpha_{c_i,c} W^{\text{CC}} \mathbf{h}_c^{\text{CC}} \right) \right) \quad (2) \quad \mathbf{h}_{t_i}^{\text{TT}} \leftarrow \text{AGG} \left(\sigma \left(\sum_{t \in \mathcal{N}_{t_i}} \alpha_{t_i,t} W^{\text{TT}} \mathbf{h}_t^{\text{TT}} \right) \right)$$

$$\mathbf{h}_{c_i}^{\text{CT}} \leftarrow \text{AGG} \left(\sigma \left(\sum_{t \in \mathcal{N}_{c_i}} \alpha_{c_i,t} W^{\text{CT}} \mathbf{h}_t^{\text{CT}} \right) \right) \quad (3) \quad \mathbf{h}_{t_i}^{\text{TC}} \leftarrow \text{AGG} \left(\sigma \left(\sum_{c \in \mathcal{N}_{t_i}} \alpha_{t_i,c} W^{\text{TC}} \mathbf{h}_c^{\text{TC}} \right) \right)$$

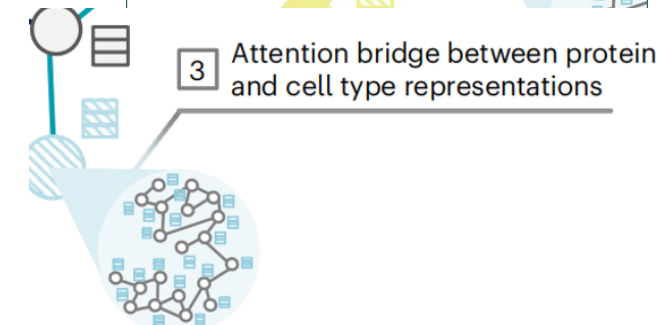
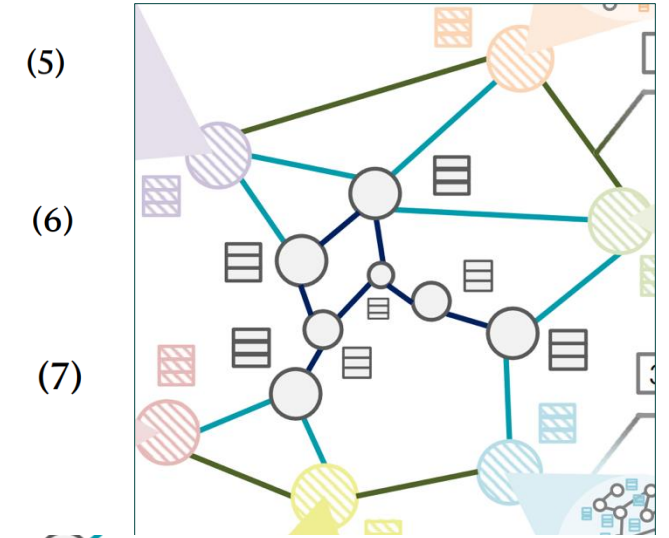
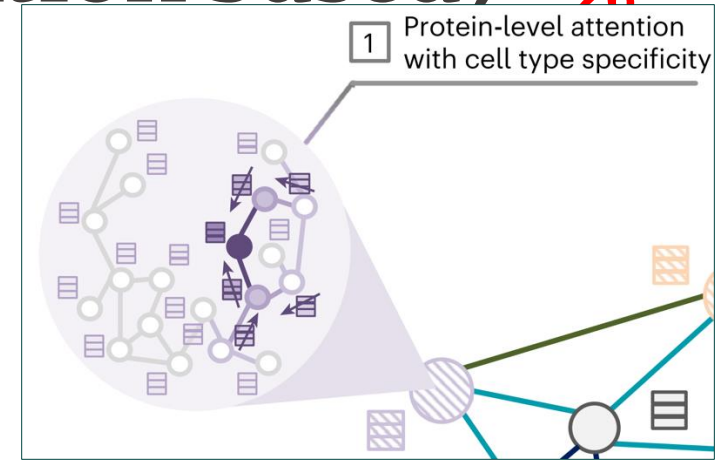
$$\mathbf{h}_{c_i} = \beta^{\text{CC}} \mathbf{h}_{c_i}^{\text{CC}} + \beta^{\text{CT}} \mathbf{h}_{c_i}^{\text{CT}} \quad (4) \quad \mathbf{h}_{t_i} = \beta^{\text{TT}} \mathbf{h}_{t_i}^{\text{TT}} + \beta^{\text{TC}} \mathbf{h}_{t_i}^{\text{TC}}$$

learned attention weights $\beta^r, r = CC, CT, TC, TT, \quad \beta^r = \frac{\exp(m_r)}{\sum_{r \in R} \exp(m_r)}$

$$m_r = \sum_{u \in V_q} \mathbf{s}^T \cdot \tanh(M \cdot \mathbf{h}_u^r + \mathbf{b})$$

3. Bridge between protein and cell type embeddings.

$$\mathbf{h}_{c_i} \leftarrow \mathbf{h}_{c_i} + \text{AGG} \left(\sigma \left(\sum_{u \in V_{c_i}} \gamma_{c_i,u} \mathbf{h}_u \right) \right). \quad \mathbf{h}_u \leftarrow \mathbf{h}_u + \gamma_{c_i,u} \mathbf{h}_{c_i}.$$



Self-supervised learning

$$\mathcal{L} = \mathcal{L}_{\text{protein}} + (1 - \theta)(\mathcal{L}_{\text{celltype}} + \mathcal{L}_{\text{tissue}}),$$

- Each cell type-specific PPI network (link prediction) and cell type identity

$$\mathcal{L}_{\text{protein}} = \theta \mathcal{L}_{\text{ppi}} + \lambda \mathcal{L}_{\text{celltypeid}}$$

$$\mathcal{L}_{\text{ppi}} = \sum_{c_i \in \mathcal{C}} \sum_{u, v \in V_{c_i}} y_{u, v} \log(\hat{y}_{u, v}) + (1 - y_{u, v}) \log(1 - \hat{y}_{u, v})$$

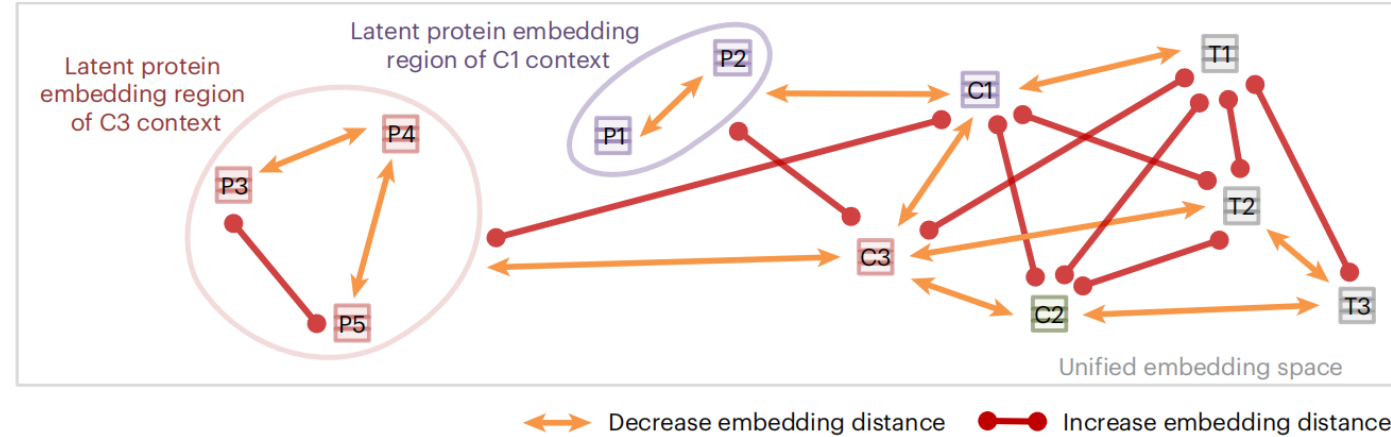
$$\mathcal{L}_{\text{celltypeid}} = \sum_{c_i \in \mathcal{C}} \sum_{u \in V_{c_i}} \|\mathbf{z}_u - \mathbf{z}_{c_i}\|_2^2 \quad (\text{cell type identity})$$

- Cell type level (link prediction) BCE(C-C + C-T)

$$\mathcal{L}_{\text{celltype}} = \mathcal{L}_{\text{celltype}}^{\text{CC}} + \mathcal{L}_{\text{celltype}}^{\text{CT}}$$

$$\mathcal{L}_{\text{celltype}}^{\text{CC}} = \sum_{c_i, c_j \in \mathcal{C}} y_{c_i, c_j} \log(\hat{y}_{c_i, c_j}) + (1 - y_{c_i, c_j}) \log(1 - \hat{y}_{c_i, c_j})$$

$$\mathcal{L}_{\text{celltype}}^{\text{CT}} = \sum_{c_i \in \mathcal{C}} \sum_{t_k \in \mathcal{T}} y_{c_i, t_k} \log(\hat{y}_{c_i, t_k}) + (1 - y_{c_i, t_k}) \log(1 - \hat{y}_{c_i, t_k}).$$



- Tissue level (link prediction) BCE(T-T + T-C)

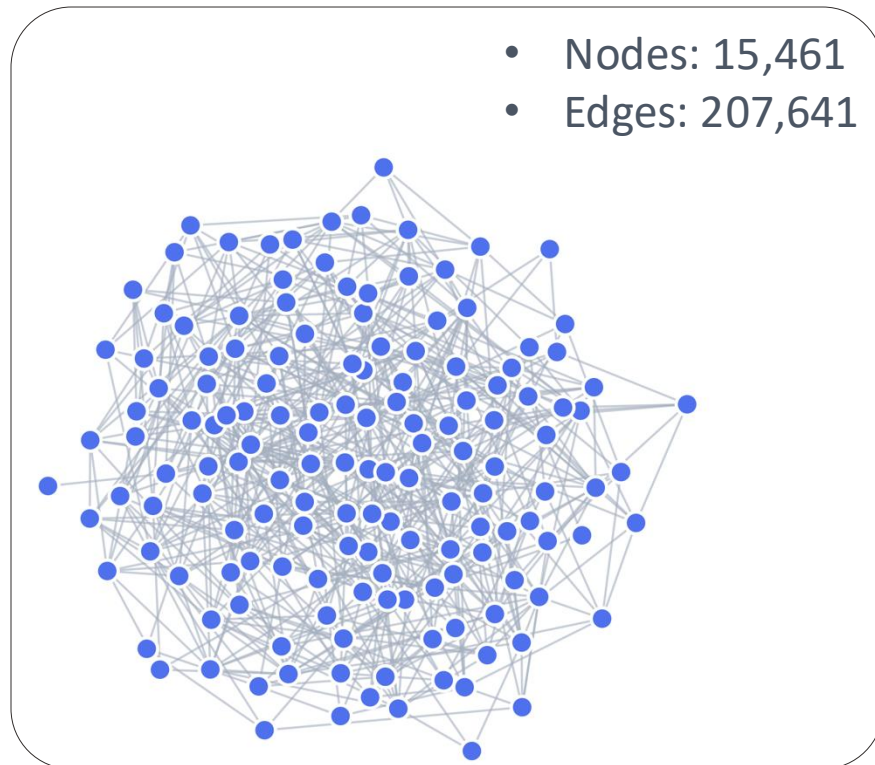
$$\mathcal{L}_{\text{tissue}} = \mathcal{L}_{\text{tissue}}^{\text{TT}} + \mathcal{L}_{\text{tissue}}^{\text{TC}}$$

$$\mathcal{L}_{\text{tissue}}^{\text{TT}} = \sum_{t_k, t_q \in \mathcal{T}} y_{t_k, t_q} \log(\hat{y}_{t_k, t_q}) + (1 - y_{t_k, t_q}) \log(1 - \hat{y}_{t_k, t_q})$$

$$\mathcal{L}_{\text{tissue}}^{\text{TC}} = \sum_{t_k \in \mathcal{T}} \sum_{c_i \in \mathcal{C}} y_{t_k, c_i} \log(\hat{y}_{t_k, c_i}) + (1 - y_{t_k, c_i}) \log(1 - \hat{y}_{t_k, c_i}).$$

Reference human physical PPI network

BioGRID, HuRI and Menche et al^[1] with 15,461 nodes and 207,641 edges.



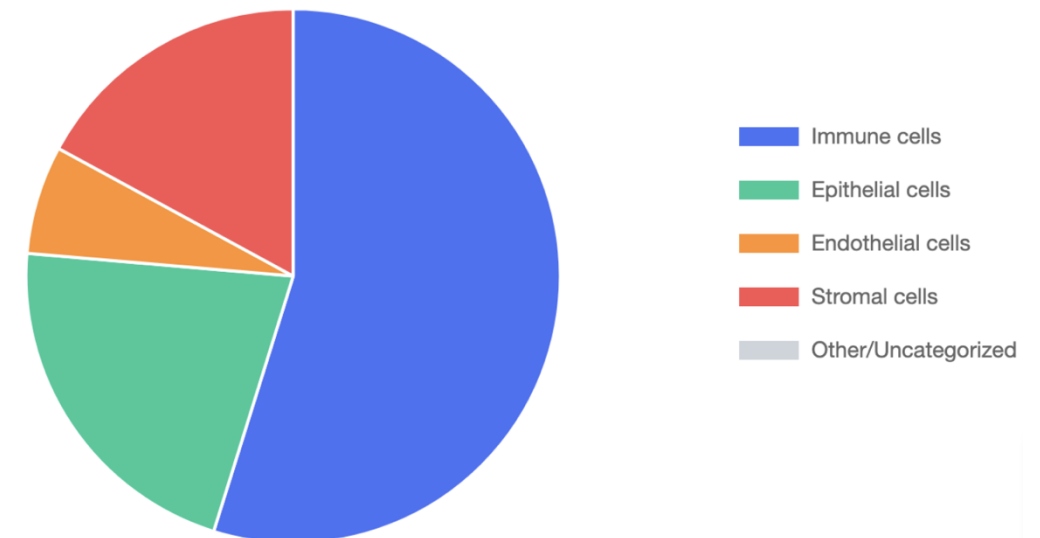
Multiorgan, single-cell transcriptomic atlas of humans.

Total Cells: 483,152 (from 15 Donors, 59 Specimens) *Tabula Sapiens*

Data Summary:

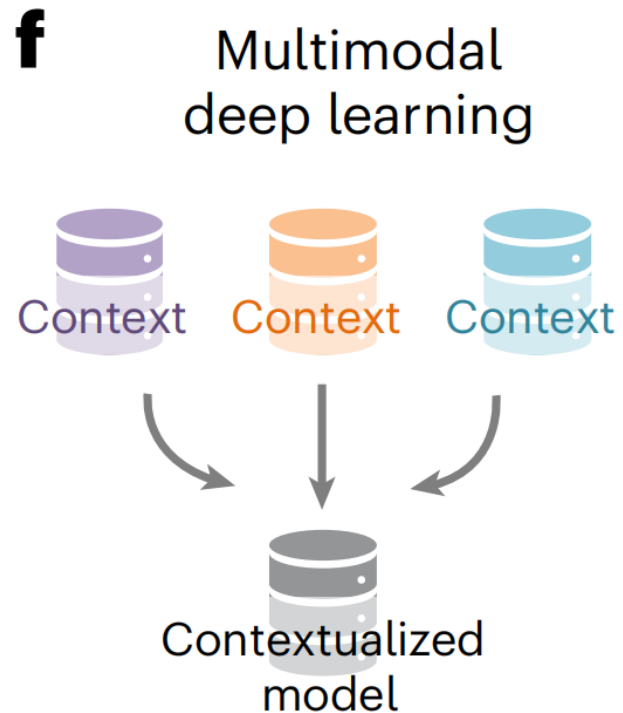
- Immune cells: **264,824** (54.81%)
- Epithelial cells: **104,148** (21.56%)
- Endothelial cells: **31,691** (6.56%)
- Stromal cells: **82,478** (17.07%)

Proportional Makeup of Cell Compartments

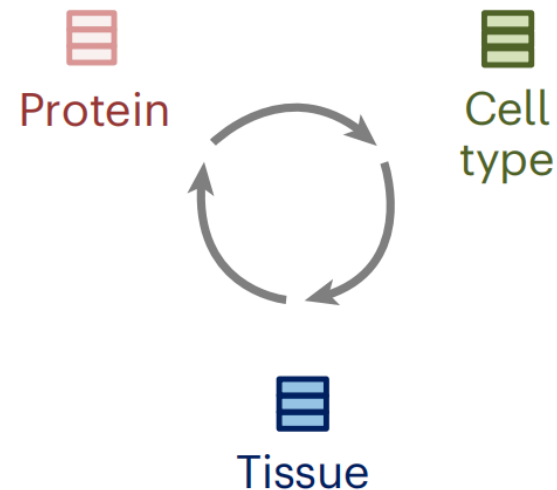


Tasks:

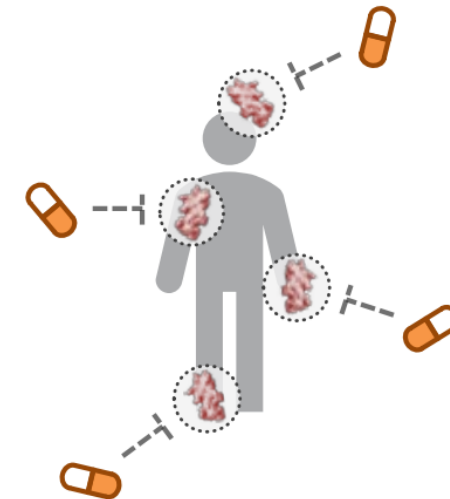
- enhance three-dimensional (3D) structural protein representations,
- analyze the effects of drugs across cell type contexts, nominate therapeutic targets in a cell type-specific manner
- retrieve tissue hierarchy in a zero-shot manner
- perform context-specific transfer learning



g Context-specific transfer learning



h Contextualized prediction



Results: captures cellular and tissue organization

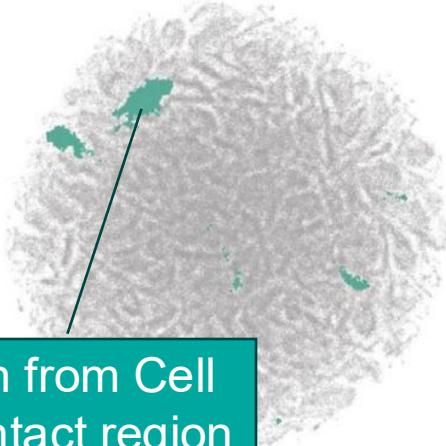
Two dimensional UMAP(uniform manifold approximate projection) plots of contextualized **protein representations** generated by PINNACLE from six different cell type context.

Quantified by neighborhood enrichment scores (NES), the higher is better.

Protein from Cell text contact region

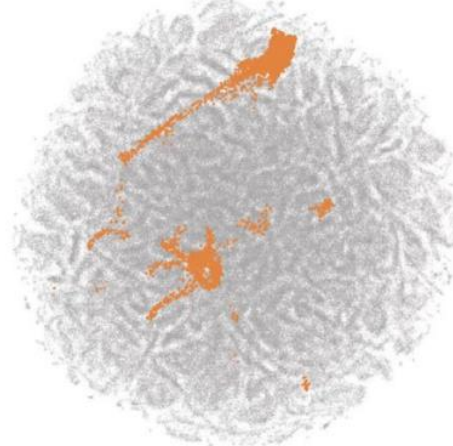
proteins from other cell types

a Protein embedding region of cell type context: medullary thymic epithelial cell



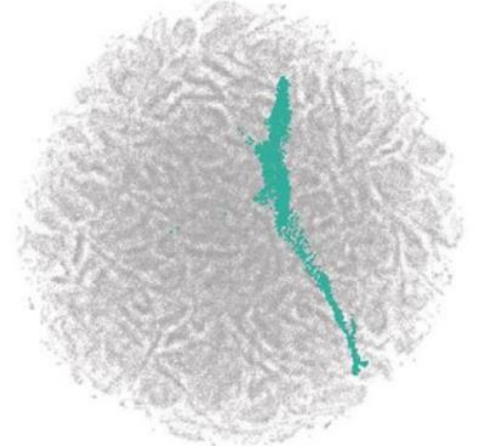
Mean SAFE NES = 53.01 ± 7.24
 Max SAFE NES = 58.93
 SAFE enriched neighborhoods = 182

b Protein embedding region of cell type context: bronchial vessel endothelial cell



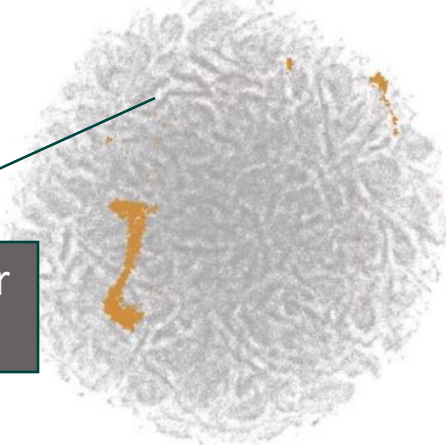
Mean SAFE NES = 23.43 ± 12.20
 Max SAFE NES = 36.99
 SAFE enriched neighborhoods = 264

c Protein embedding region of cell type context: mesenchymal stem cell



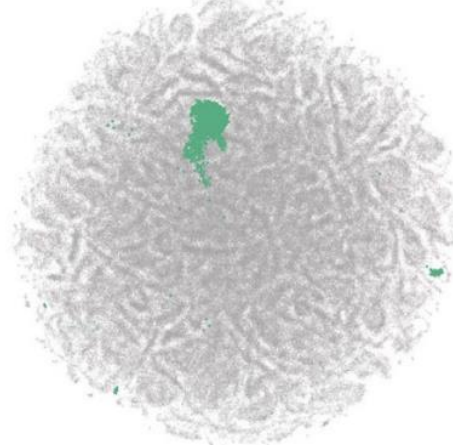
Mean SAFE NES = 15.12 ± 20.25
 Max SAFE NES = 63.86
 SAFE enriched neighborhoods = 144

d Protein embedding region of cell type context: lung microvascular endothelial cell



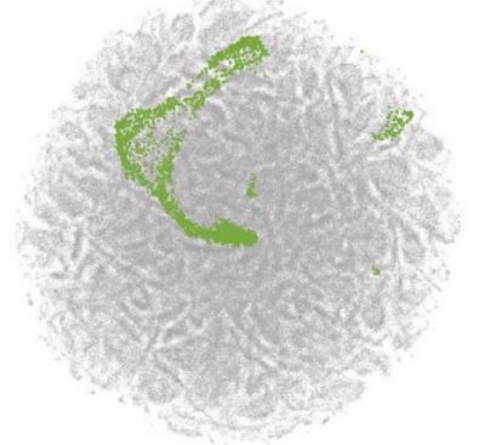
Mean SAFE NES = 37.75 ± 4.87
 Max SAFE NES = 44.10
 SAFE enriched neighborhoods = 355

e Protein embedding region of cell type context: kidney epithelial cell



Mean SAFE NES = 52.90 ± 12.95
 Max SAFE NES = 62.78
 SAFE enriched neighborhoods = 137

f Protein embedding region of cell type context: fibroblast of breast

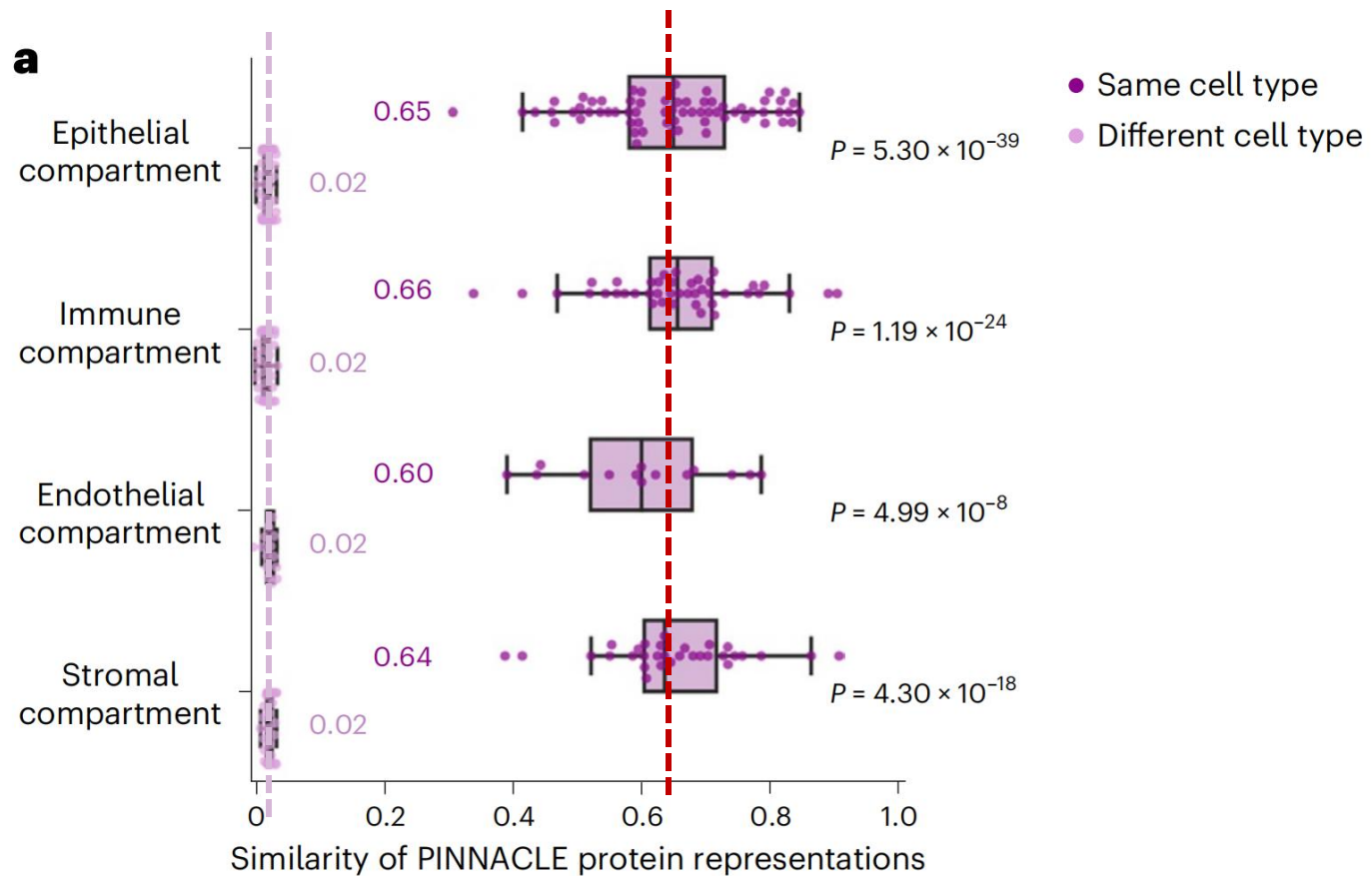


Mean SAFE NES = 49.67 ± 4.12
 Max SAFE NES = 65.83
 SAFE enriched neighborhoods = 235

Embedding similarities using PINNACLE's protein representations

Same cell type v.s. Different Cell type

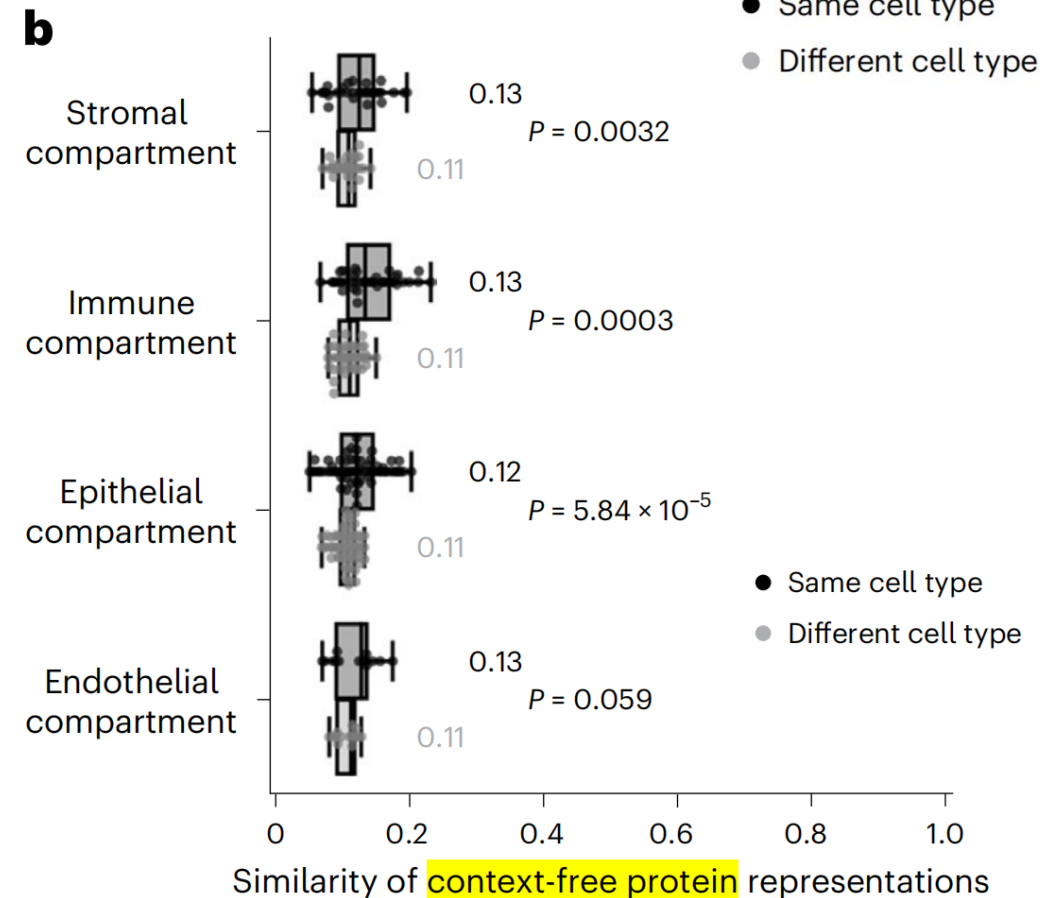
Protein presentations from the same cell type show high similarity



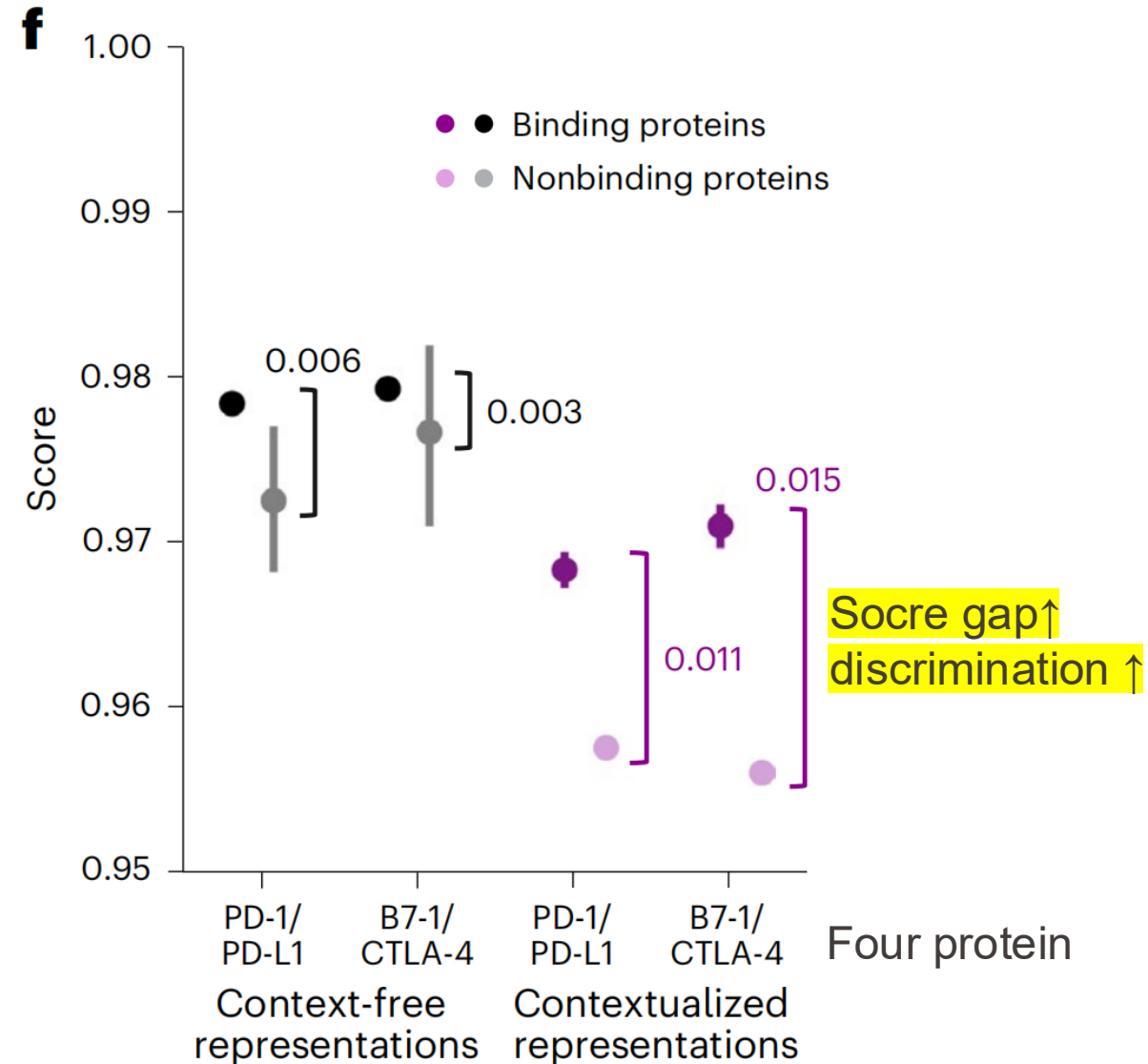
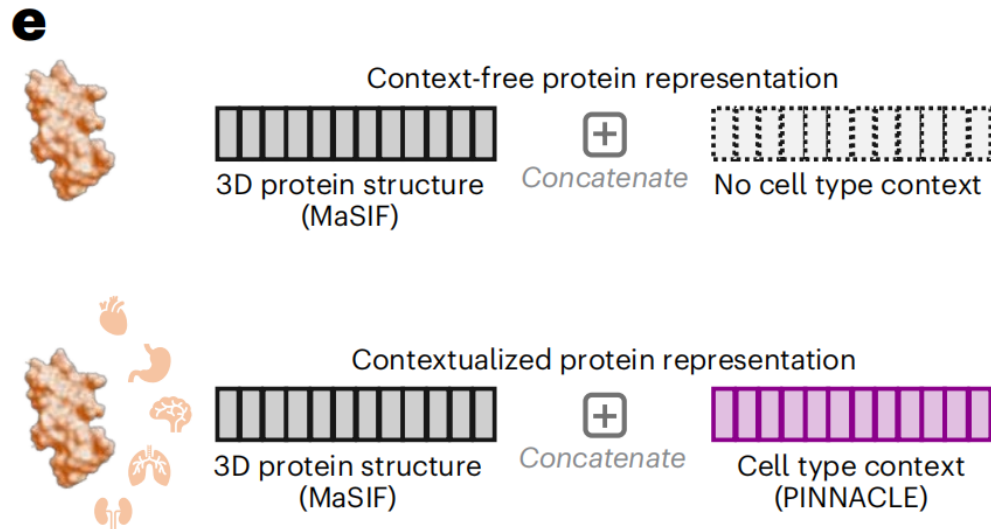
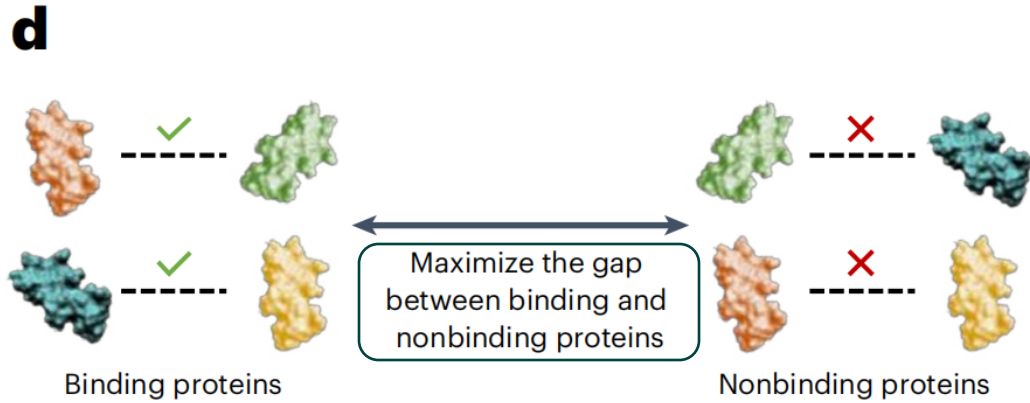
KS two-sample KS test P-value, lower is better (from two distribution)

Context specific v.s. Context-free

Context-free PINNACLE: remove cell type and tissue network and all cell type- and tissue-related components of PINNACLE's architecture and objective function



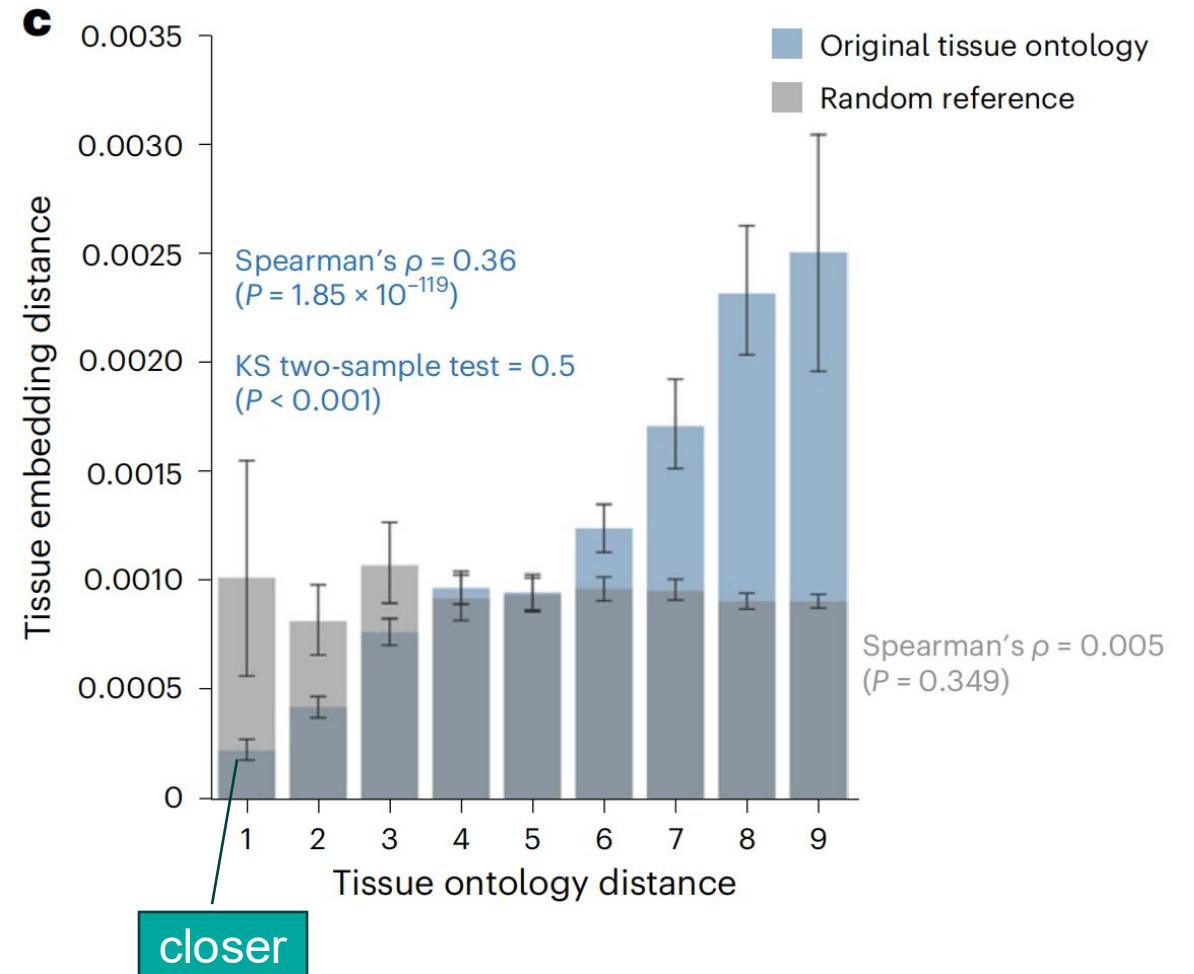
PINNACLE enhances 3D structural representations of PPIs



Context information increases the discrimination of the protein representations

Tissues discrimination ↑

This validates PINNACLE's ability **to learn and encode the underlying biological structure of tissues**. The closer the tissues are in the real world (ontology), the closer they are in the model's embedding space.

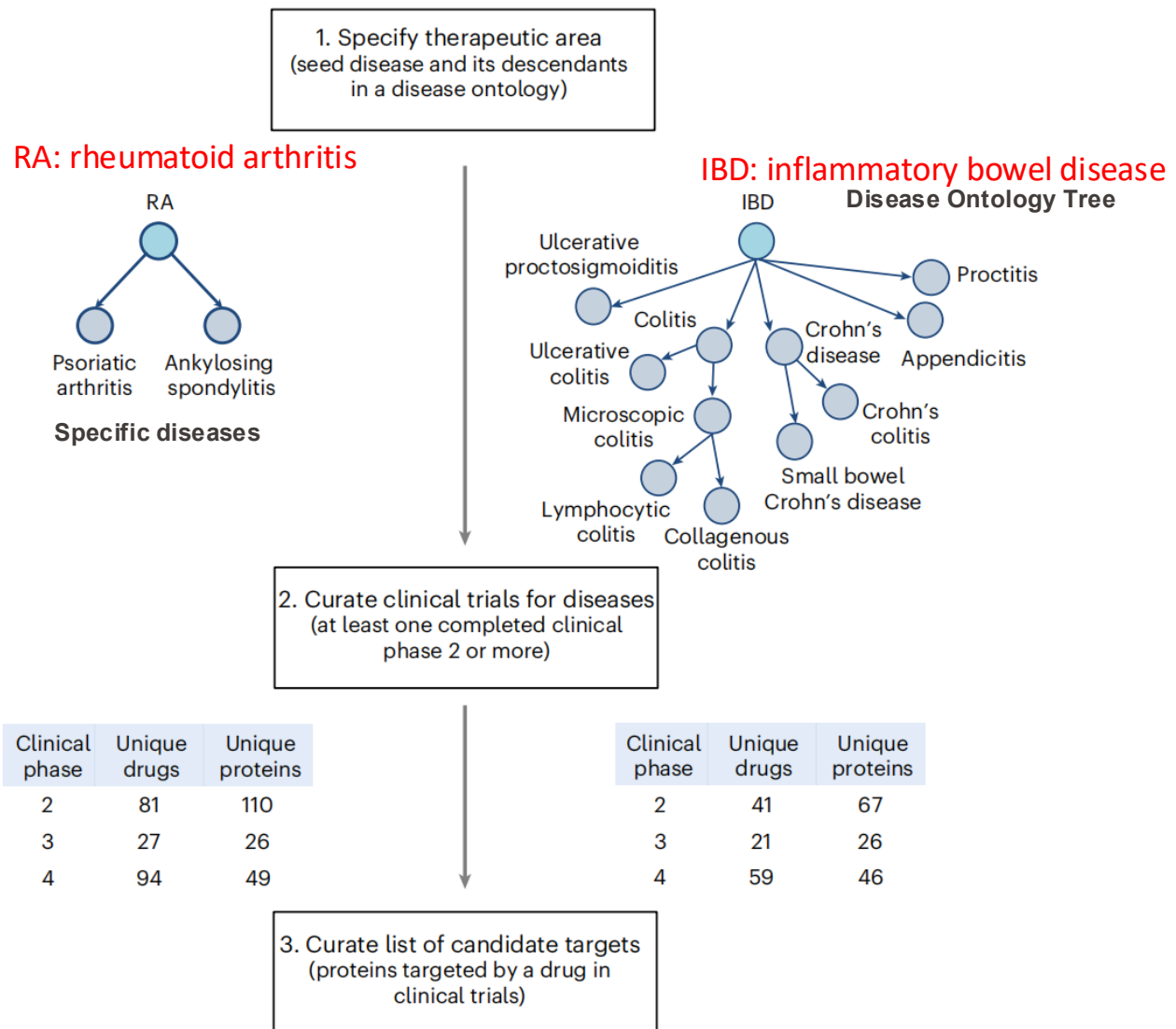


As the biological distance (X-axis) increases, the computational distance (Y-axis) also increases.

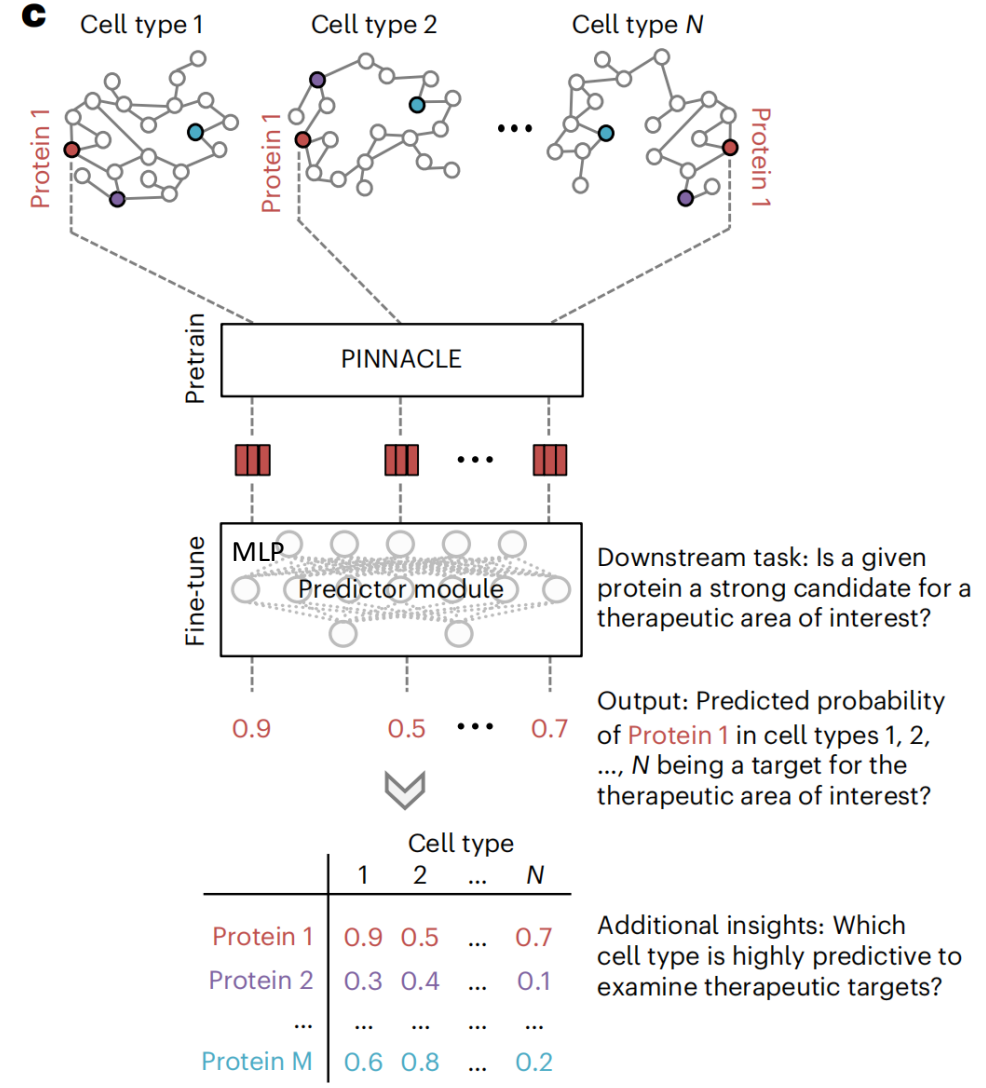
Therapeutic target prioritization (via fine-tuning contextualized protein representations)

The predictor module (that is, MLP) fine-tunes the (pretrained) contextualized protein representations for predicting whether a given protein is a strong candidate for the therapeutic area of interest. (RA: rheumatoid arthritis + IBD: inflammatory bowel disease)

a



c

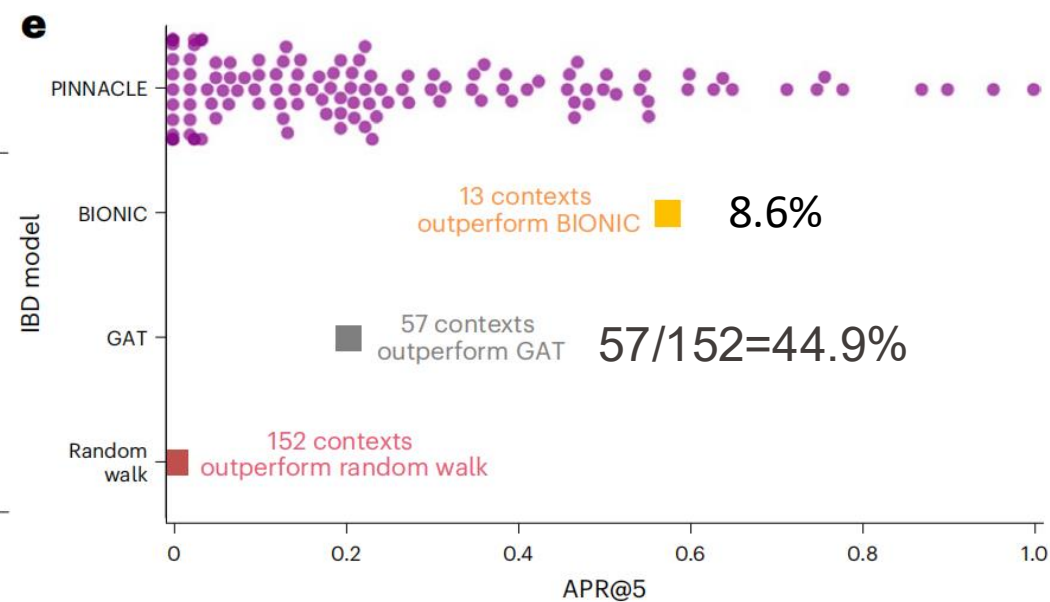
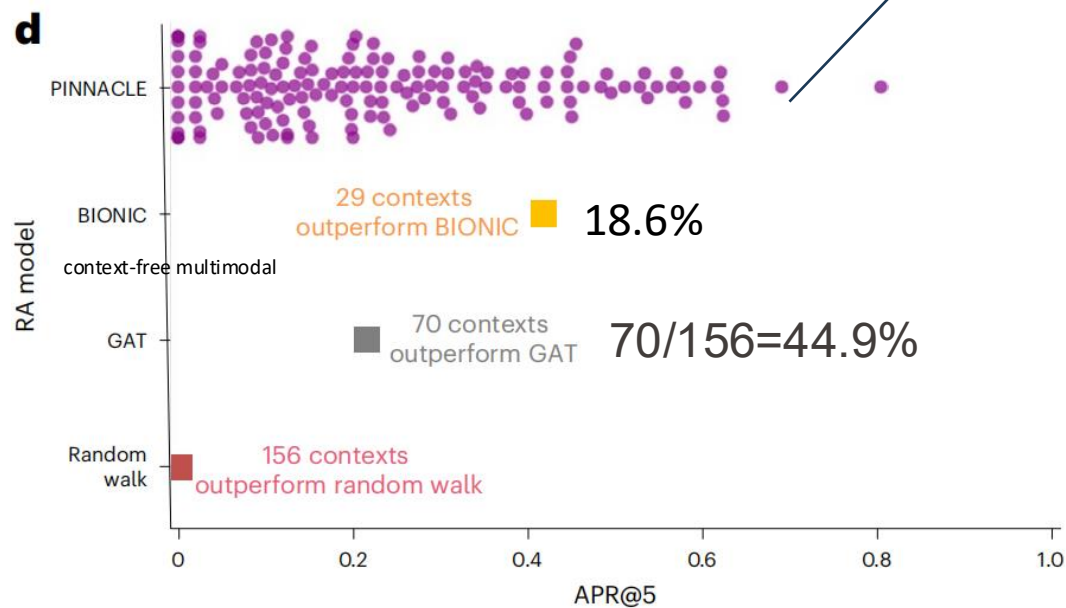


EPFL Contextual models outperform context-free target prediction 35

Benchmarking of context-aware and context-free approaches for RA (d) and IBD (e) therapeutic areas.

Baselines: Random walk, GAT, BIONIC (context free model)

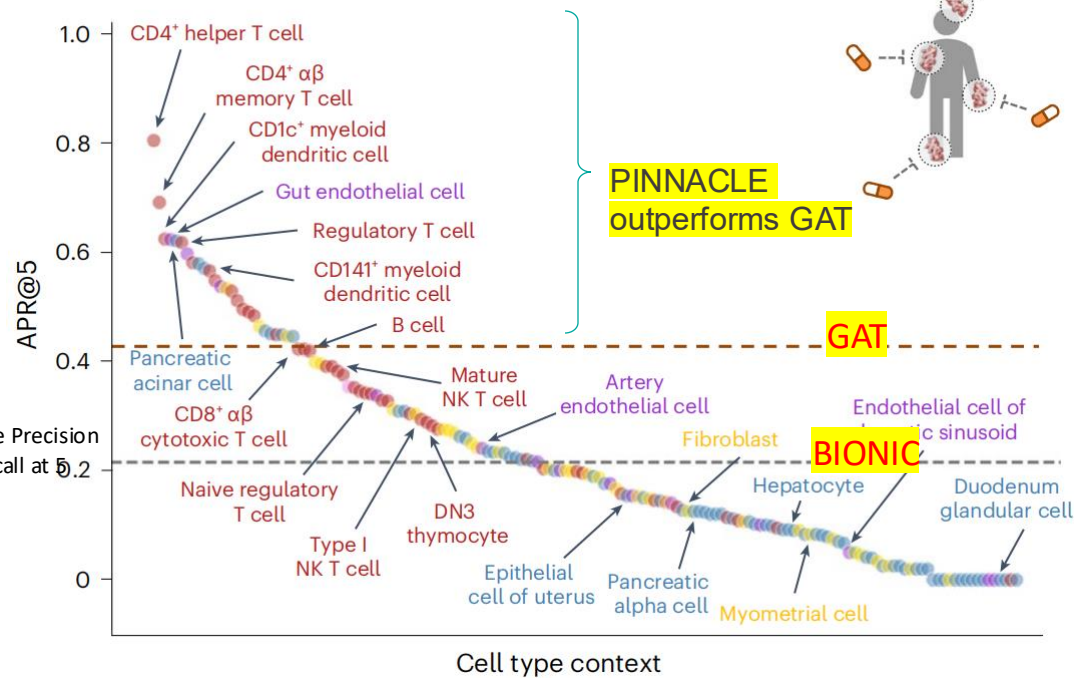
Dot: performance (averaged across ten random seeds) of protein representations from a given context



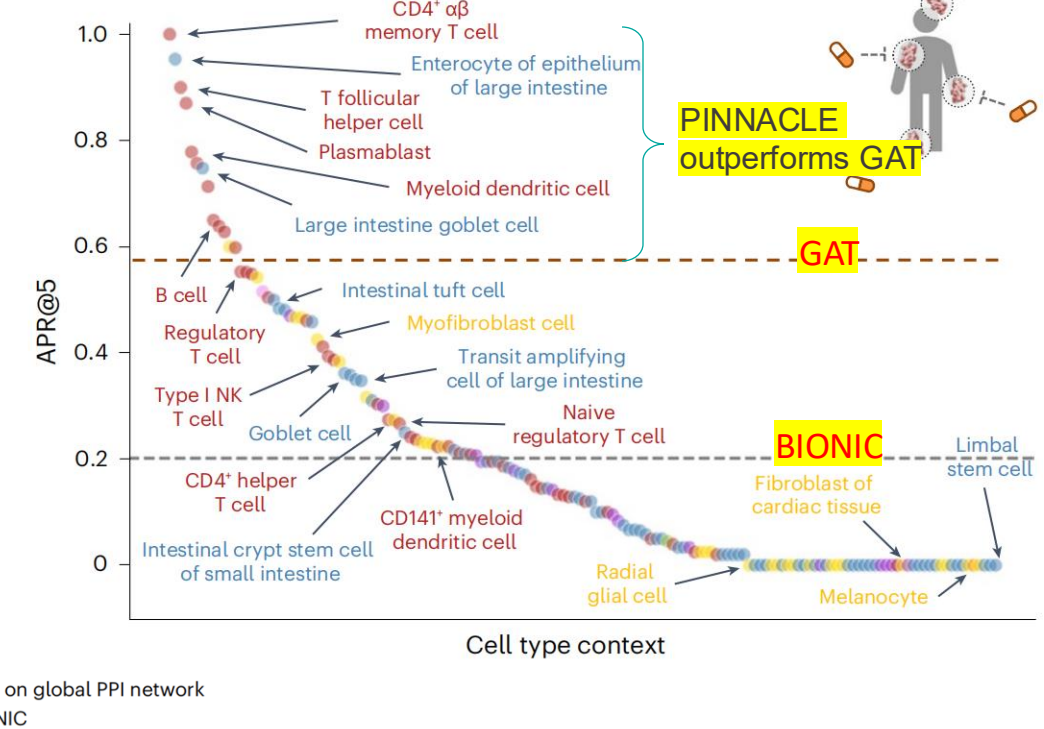
Performance of contextualized target prioritization for RA and IBD therapeutic areas

PINNACLE can nominate targets across cell type contexts

a RA therapeutic area

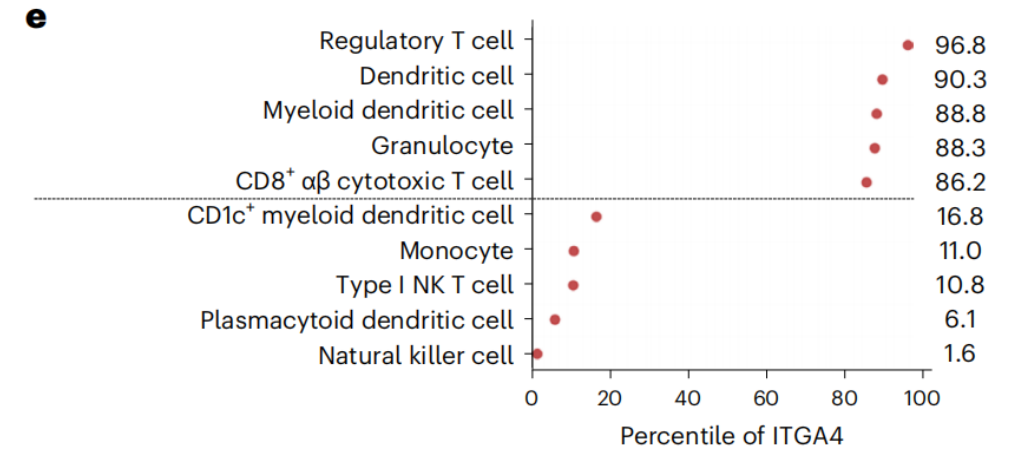
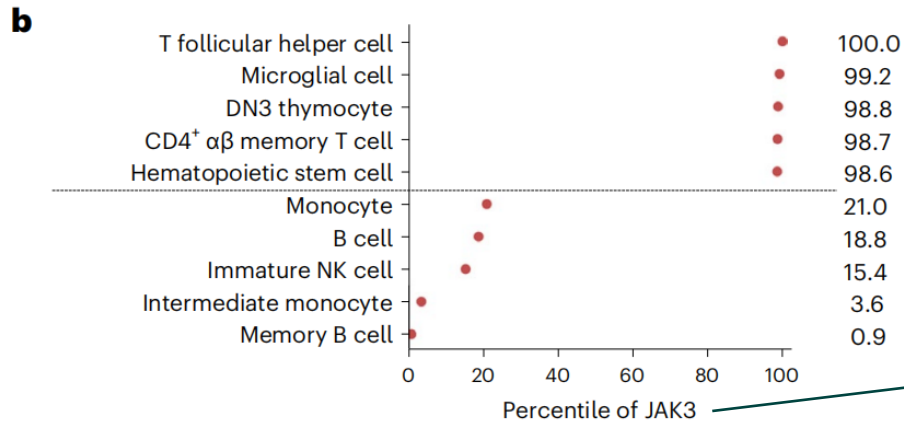


d IBD therapeutic area

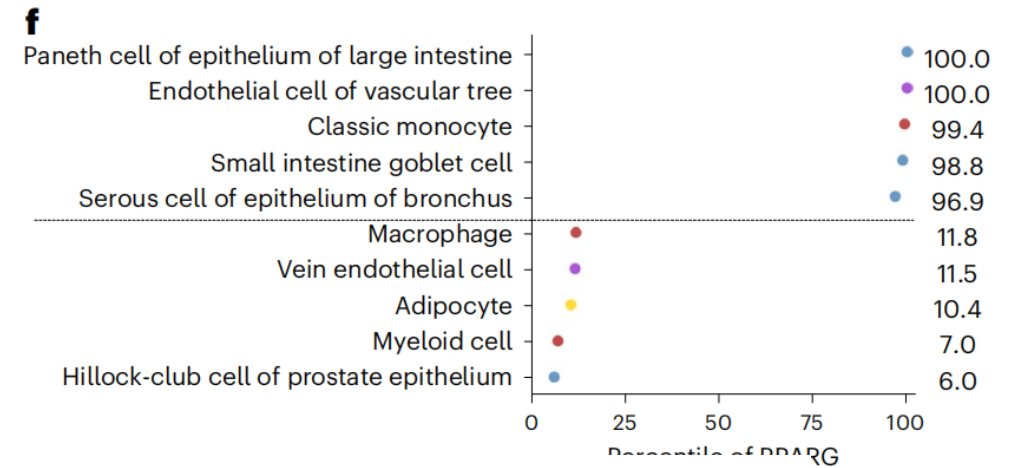
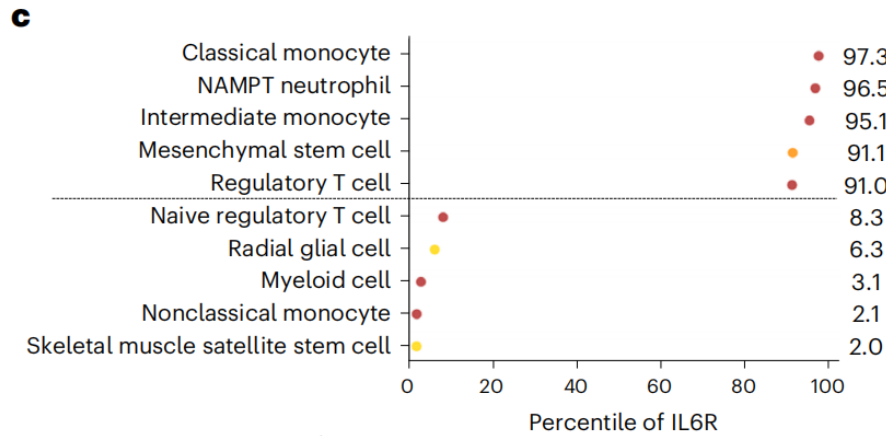


▪ Predictive cell type contexts reflect MoAs in RA therapies

▪ Predictive cell type contexts elucidate MoAs in IBD therapies



protein



PINNACLE

- Epithelial context
- Endothelial context
- Immune-stromal context
- Germ line context
- GAT on global PPI network
- Stromal context
- Immune context
- Stromal-epithelial context
- BIONIC

compartment

PINNACLE is a protein representation learning method that use the context learning for enhancing the context-specific property of the presentations

Positives:

- **Sophisticated Modeling:** Utilizes GNNs to capture complex topological interactions within the static PPI network backbone.
- **High-Resolution Targeting:** Achieves cell-type specificity in protein representations, moving beyond "average" network views using *single cell atlases* context.

Limitations:

- **Inter-Cell Signaling:** Long-range connections (protein interactions between different cell types) are currently only modeled implicitly through the Metagraph structure, simplifying direct molecular dialogue.

EPFL

EVENT

**Thank you
for your
attention !**

■ École
polytechnique
fédérale
de Lausanne

PRESENTED BY

Jiying Zhang