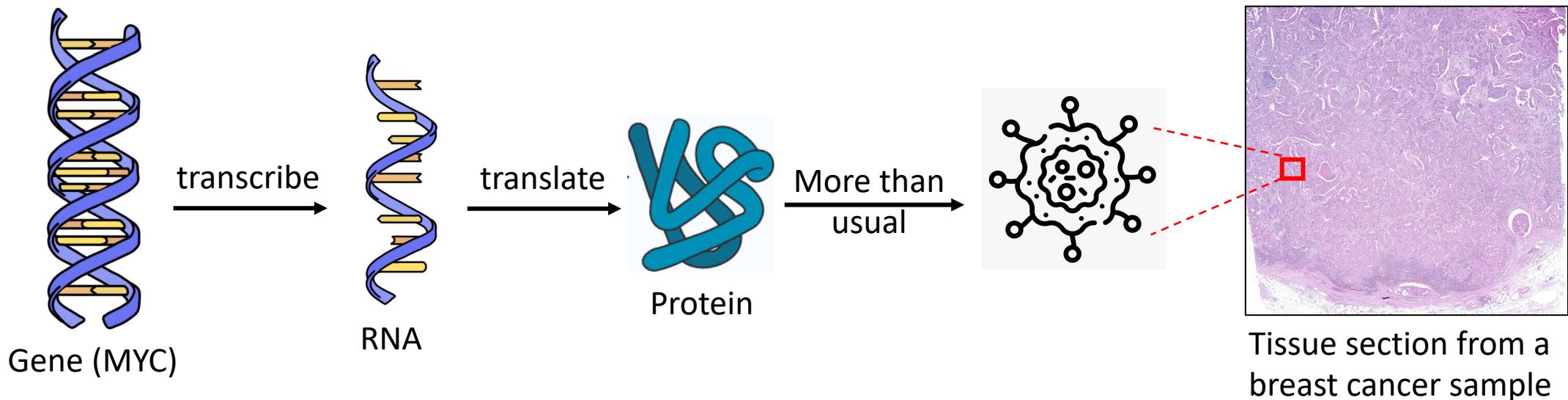


Predicting Spatial Transcriptomics from Histology Images via Biologically Informed Flow Matching

Tinglin Huang

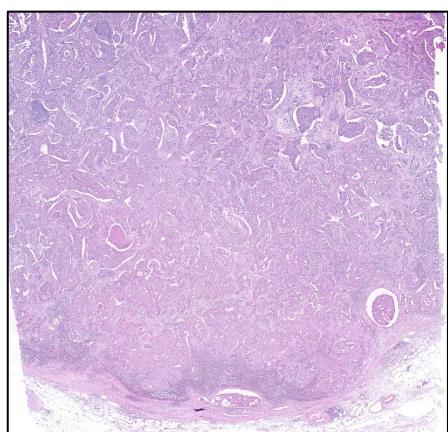
Background

- Gene expression regulates the cellular processes by central dogma
 - E.g., the overexpression of MYC gene can lead to cancer cells
 - Detecting the gene expression level can monitor the disease progression and treatment response

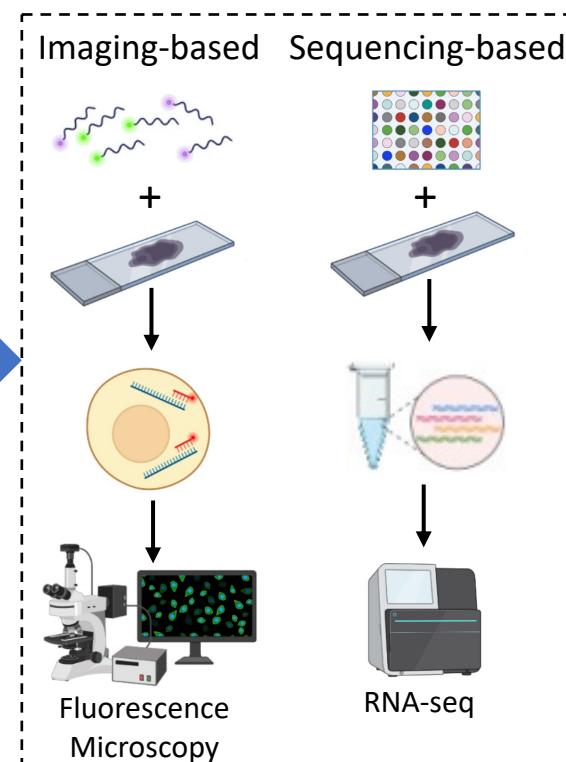


Spatial Transcriptomics (1)

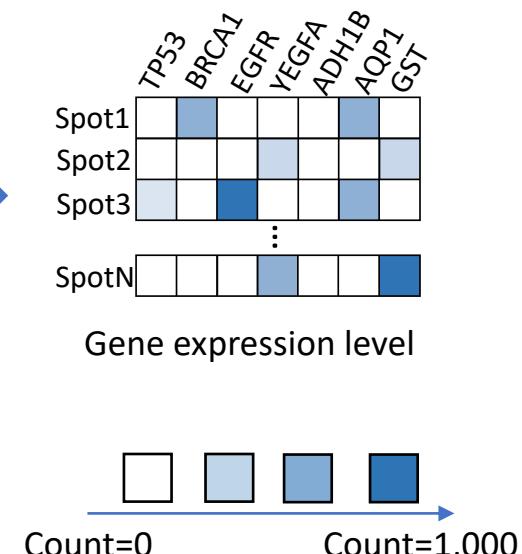
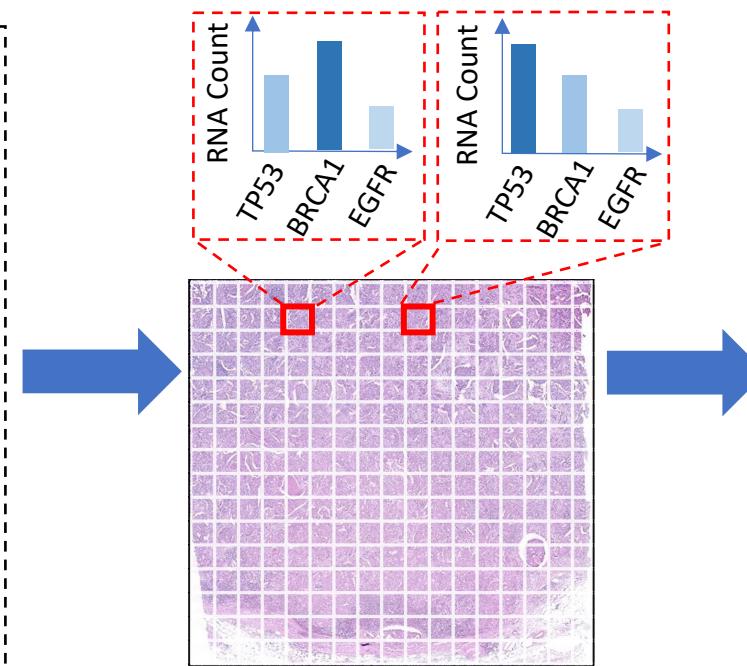
- A technology that can detect spatial distribution of RNA transcripts
 - The number of RNA transcripts represents the gene expression level



Tissue section



ST technology



Count=0 Count=1,000

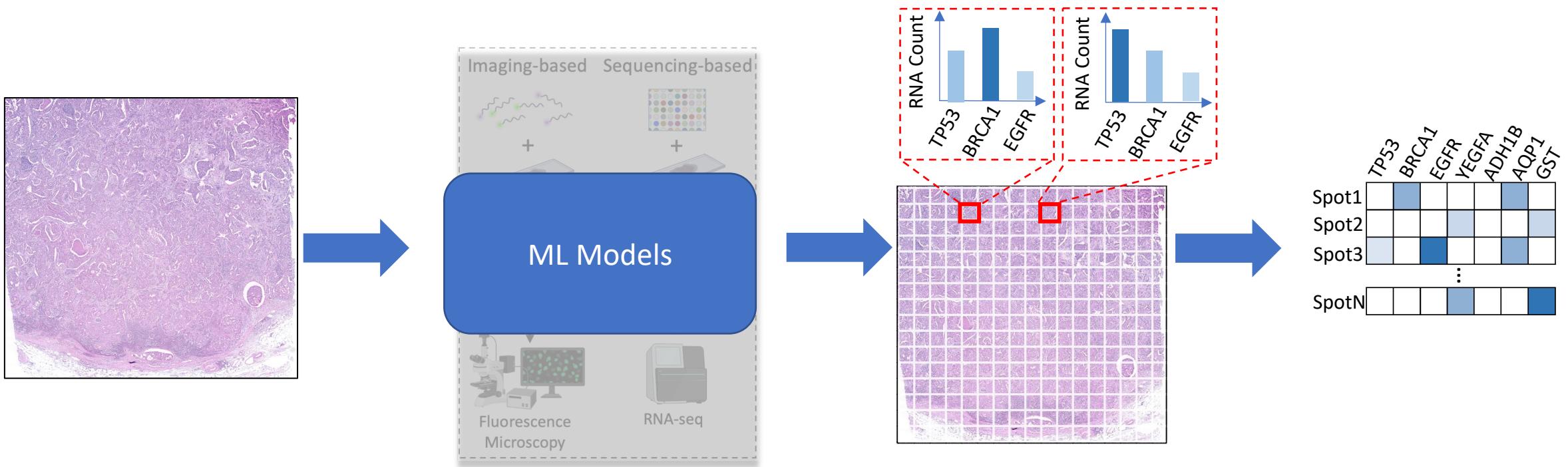
Spatial Transcriptomics (2)

- ST methods are low throughput and rely on specialized equipment
 - A single slide costs approximately \$1,000-1,500 USD and 7-10 days



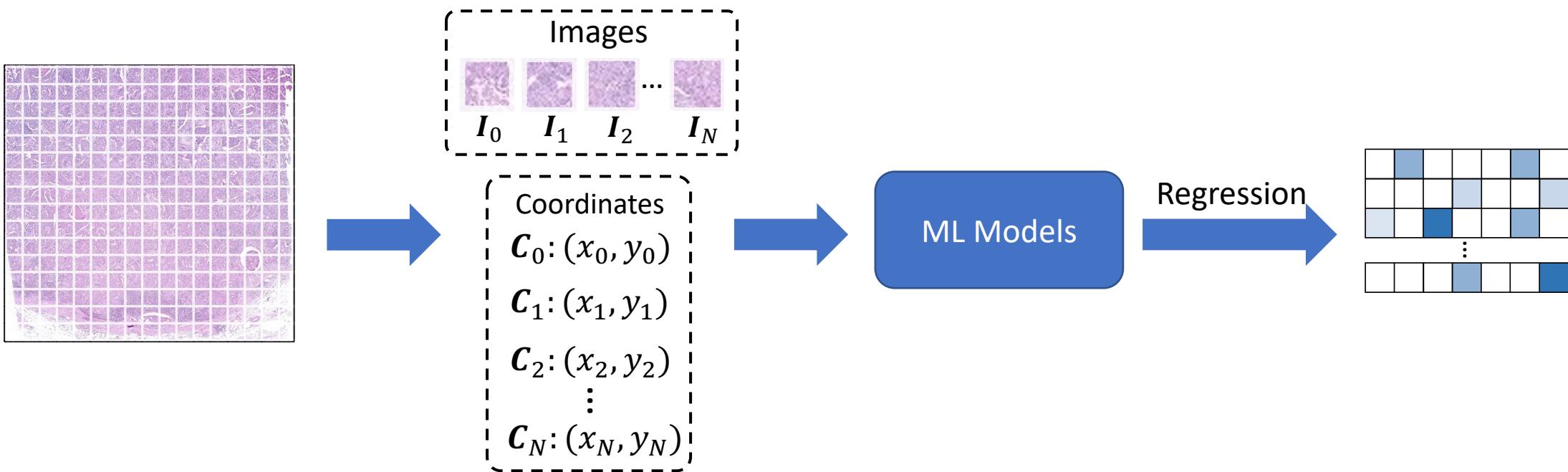
Spatial Transcriptomics (3)

- ML is one of the promising alternatives



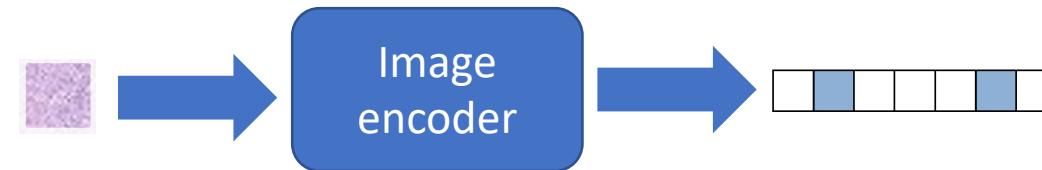
Problem formulation

- Predict the gene expression for each spot (regression task)
 - Input: spot images $\{I_1, \dots, I_N\}$, coordinates $\{C_1, \dots, C_N\}$
 - Output: gene expression $\{Y_1, \dots, Y_N\}$

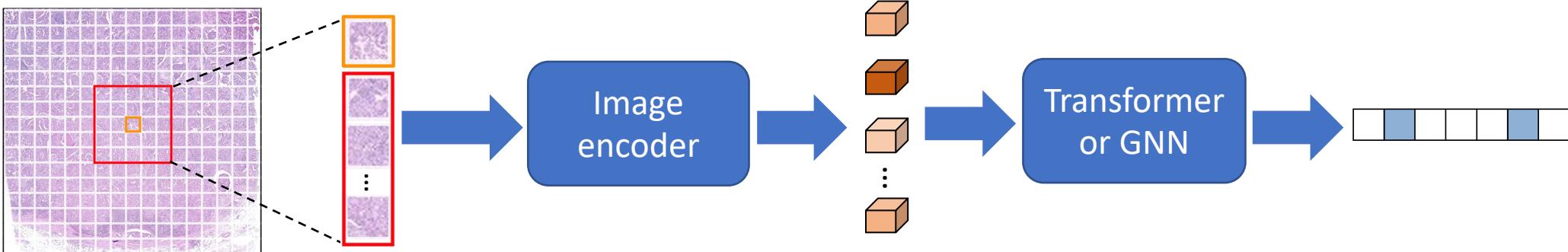


Modeling strategies

- Spot-level methods solely utilize images: $p(Y_i|I_i)$
 - Pathology foundation model shows great performance

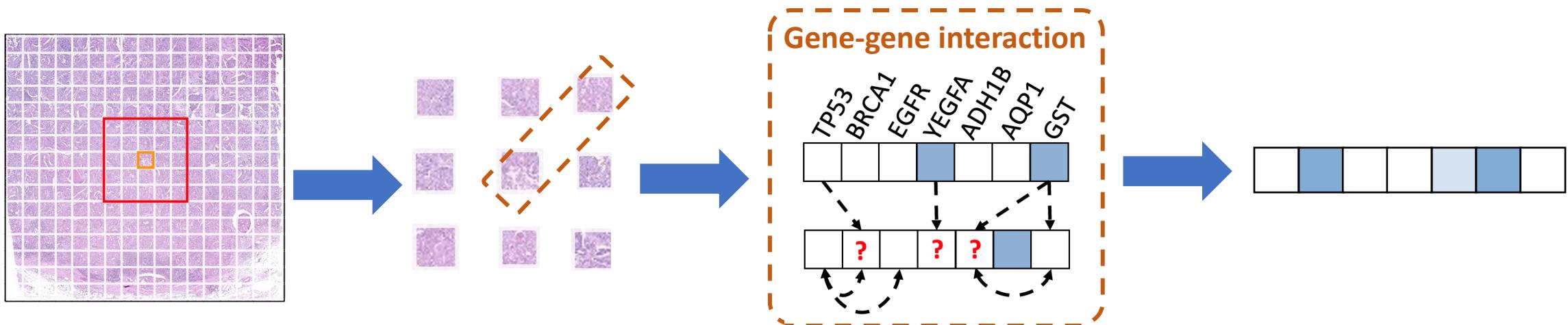


- Slide-level methods incorporate the slide context: $p(Y_i|I_1, \dots, I_N)$
 - Aggregate the neighboring spots with Transformer/GNN



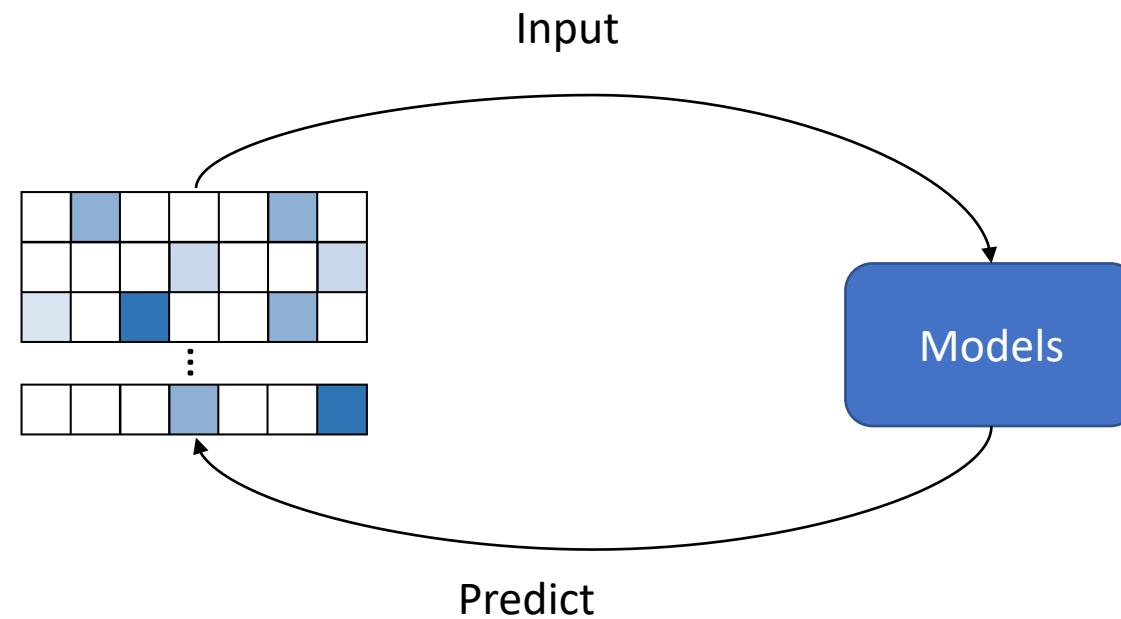
Gene-Gene Interaction

- The interaction between genes regulates the expression significantly
 - Gene regulation pathway: gene1→gene2→gene3
 - But it has been overlooked so far



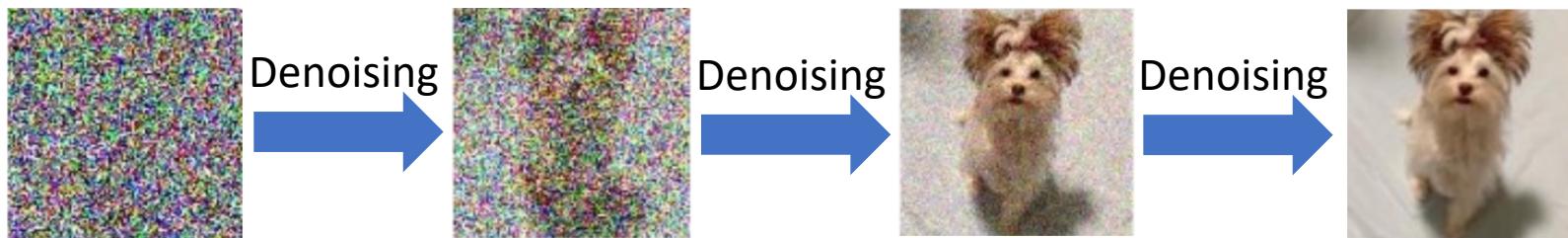
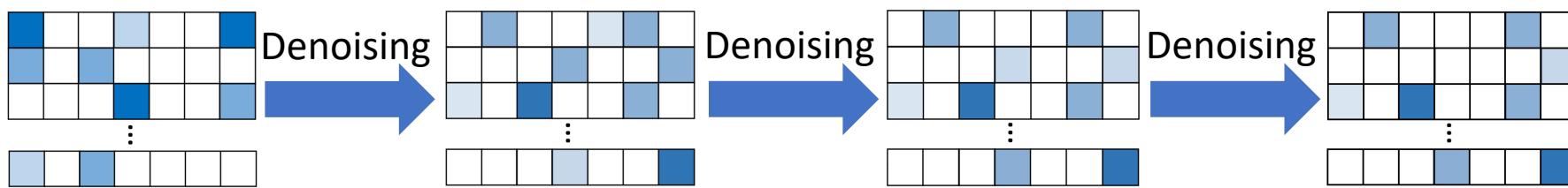
Chicken-and-Egg Problem

- Gene-gene interaction requires gene expression as context, yet gene expression itself is the target of prediction



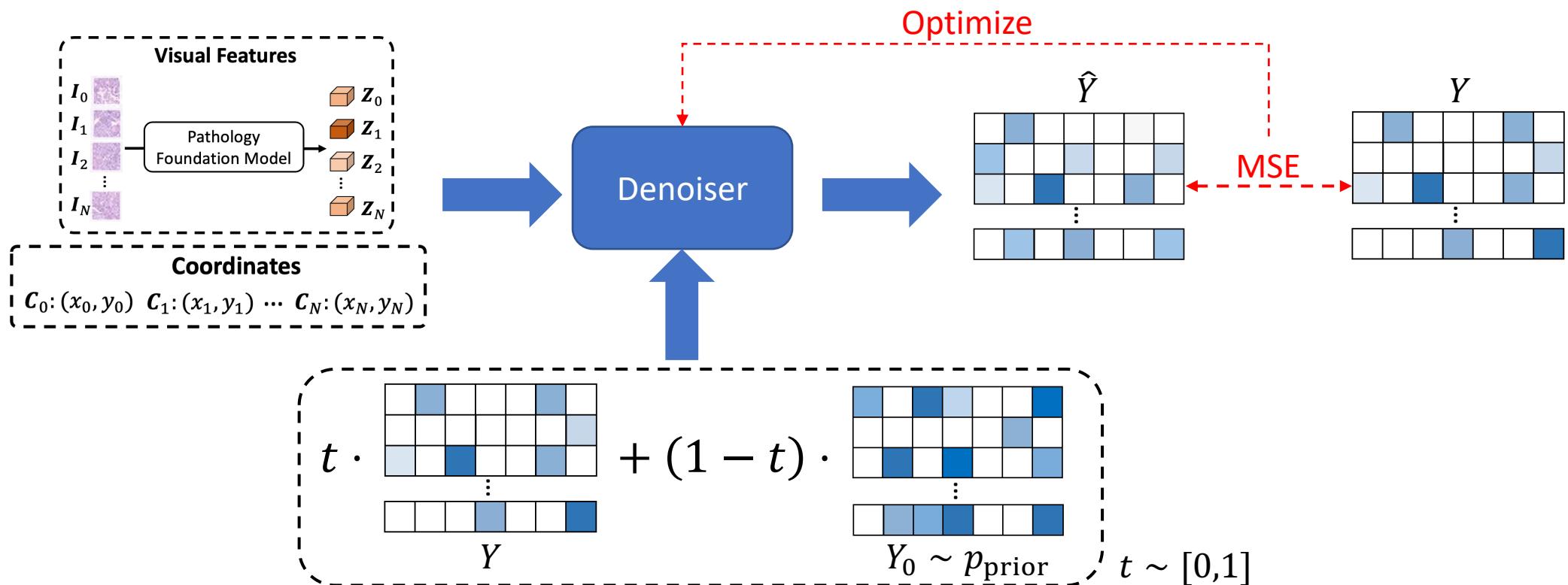
Generative Gene Expression Prediction

- Generative model provides a recipe to receive gene expression as input
 - Iteratively refine the gene expression instead of one-step regression



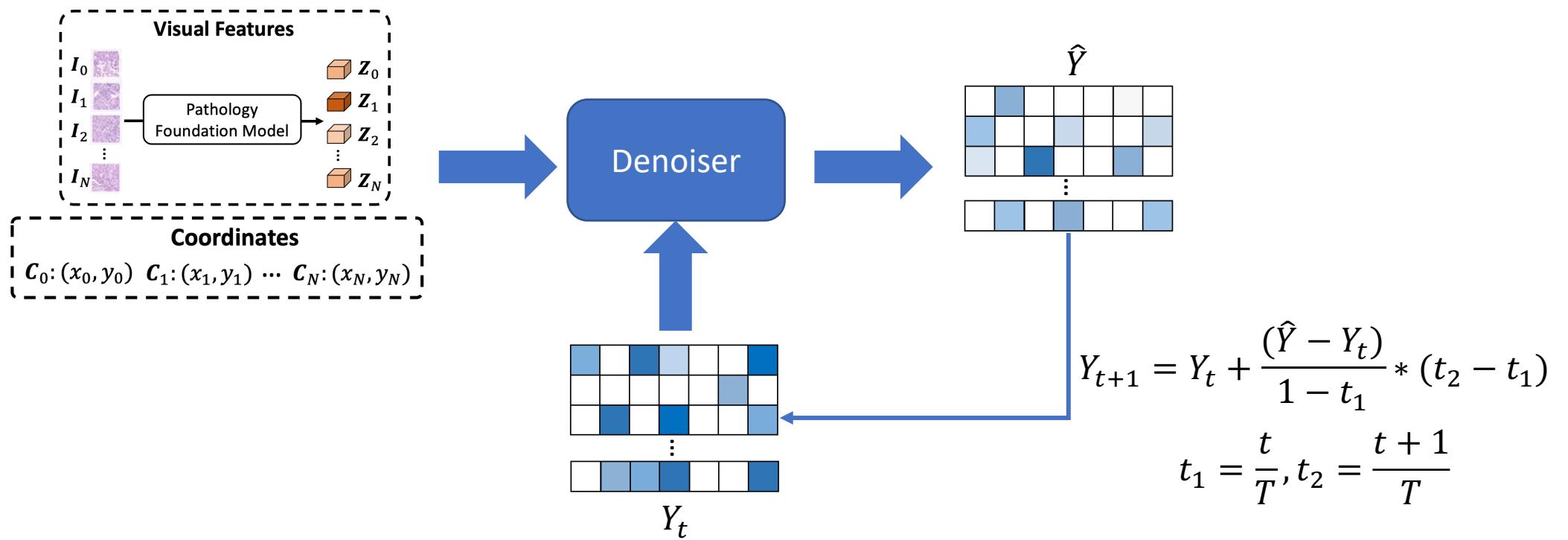
Flow Matching: Train

- Optimized by recovering gene expression from the interpolation between ground truth and a sampled noise



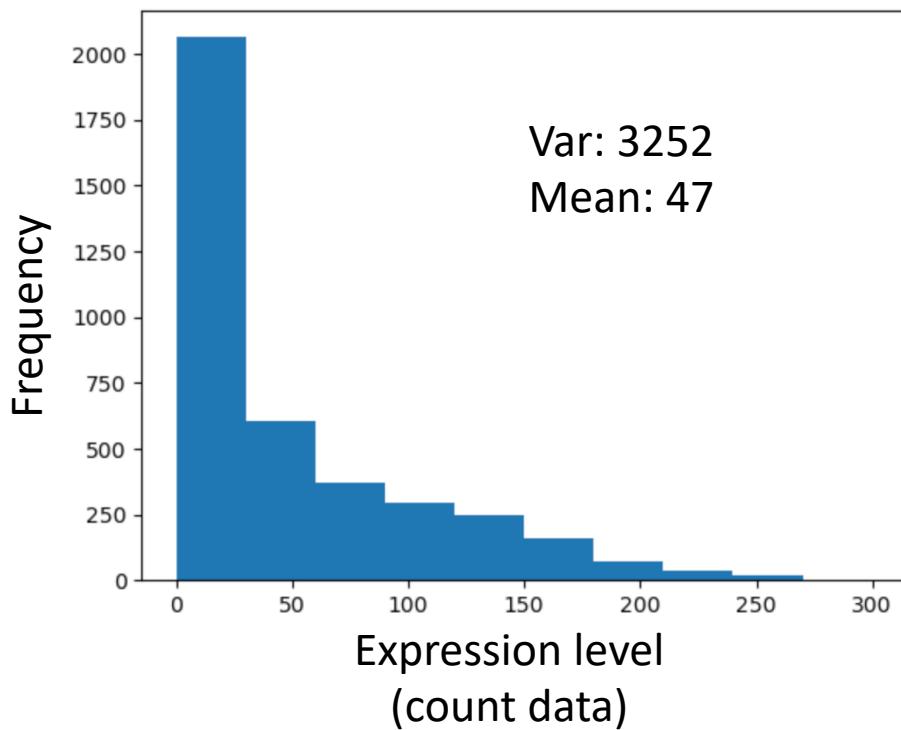
Flow Matching: Inference

- Iteratively refines the gene expression by interpolating between the noisy input Y_t and the predicted expression \hat{Y}



ZINB distribution as prior

- Gene expression is typically sparse and exhibits overdispersion
 - Overdispersion: variance > mean



- Negative binomial distribution is suitable for representing overdispersed data
- It is further extended to **zero-inflated negative binomial** (ZINB) distribution for modeling the sparse characteristic
 - The sample can randomly be dropped

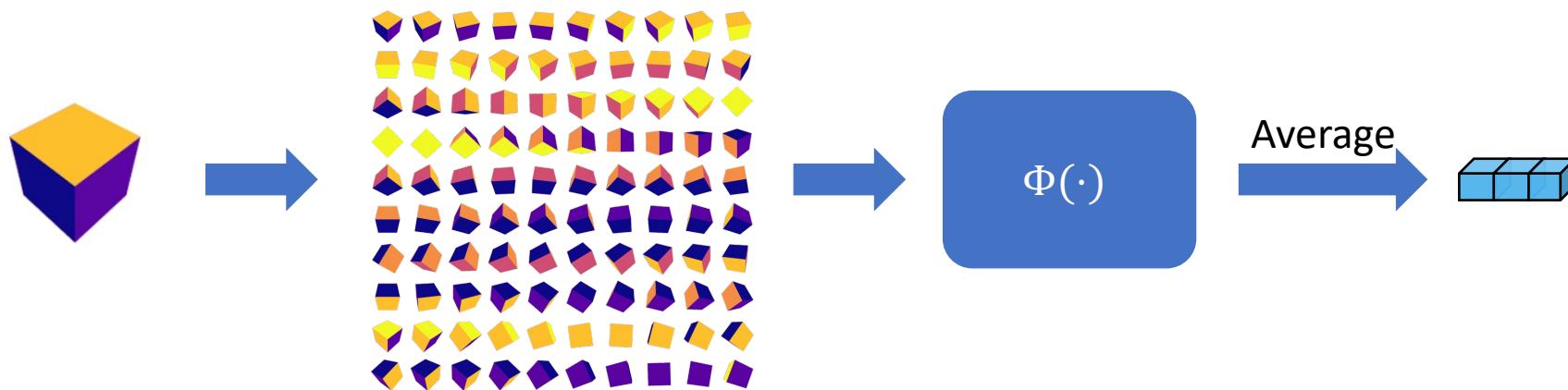
$$p(y | \mu, \phi, \pi) = \begin{cases} \pi + (1 - \pi) \left(\frac{\Gamma(y+\phi)}{\Gamma(\phi) y!} \right) \left(\frac{\phi}{\phi+\mu} \right)^{\phi} \left(\frac{\mu}{\phi+\mu} \right)^y & \text{if } y = 0, \\ (1 - \pi) \left(\frac{\Gamma(y+\phi)}{\Gamma(\phi) y!} \right) \left(\frac{\phi}{\phi+\mu} \right)^{\phi} \left(\frac{\mu}{\phi+\mu} \right)^y & \text{if } y > 0, \end{cases}$$

Negative binomial distribution

y : count, π : dropout probability, μ : mean, ϕ : number of failures

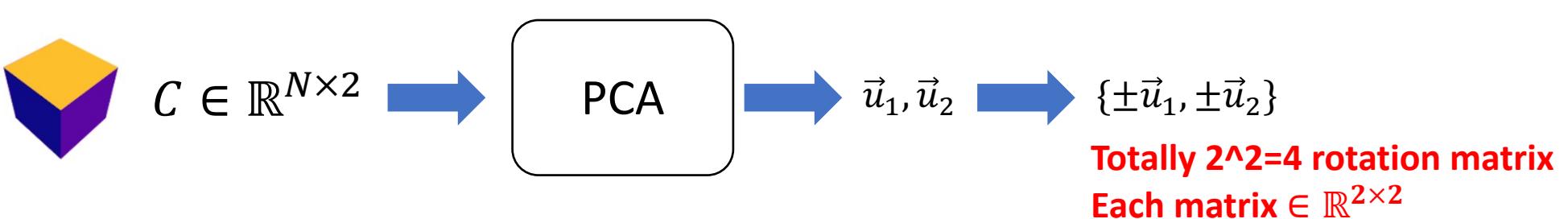
Architecture of Denoiser (1)

- Input: coordinates, visual features, gene expression
 - $C \in \mathbb{R}^{N \times 2}, Z \in \mathbb{R}^{N \times D}, Y \in \mathbb{R}^{N \times G}$
- Extend Frame Averaging Transformer[1]
 - Goal: invariantly embed the spot with rotation/translation on coordinates
 - Any model $\Phi(\cdot)$ is invariant if it considers all kinds of transformation t



Architecture of Denoiser (2)

- Frame averaging (FA) prove the group average can be conducted over **a carefully selected subset** $\mathcal{F}(X) \subset G_{\text{all}}$
 - Frame \mathcal{F} is a mapping from coordinate space into a group
- Such a subset can be calculated by PCA
 - Two principal components \vec{u}_1, \vec{u}_2
 - $\vec{u}_i \in \mathbb{R}^{2 \times 1}$



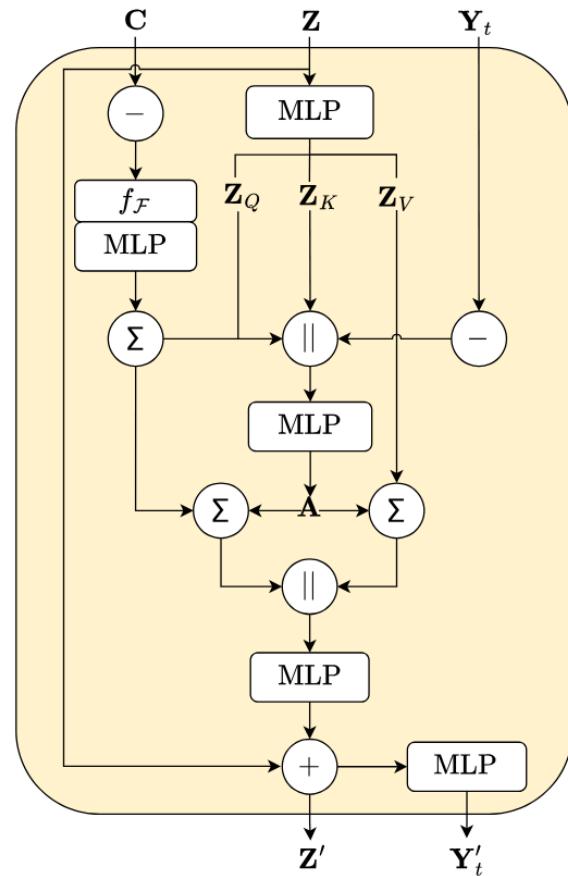
Architecture of Denoiser (3)

- Extend FAFormer[1] to incorporate all these information
 - $\mathcal{N}(i)$: k-nearest neighbors of i -th spot
 - Coordinate is embedded with frame averaging
- Attention and aggregation

$$\{c'_{i \rightarrow j}\}_{j \in \mathcal{N}(i)} = \text{FAMLP}(c_i, \{c_j\}_{j \in \mathcal{N}(i)})$$

$$a_{ij} = \text{Softmax} \left(\text{MLP}(z_i^Q || z_i^K || c'_{i \rightarrow j} || (y_i - y_j)) \right)$$

$$z = \text{MLP} \left(\sum_{j \in \mathcal{N}(i)} a_{ij} z_j^K || \sum_{j \in \mathcal{N}(i)} a_{ij} c'_{i \rightarrow j} \right) + z$$



Benchmark

- 10 datasets covering 48 patients and 74 slides
 - Metric: Pearson correlation between prediction and groundtruth
- Baselines
 - Pathology foundation models: Ciga, UNI, Gigapath
 - Slide-level approaches: STNet, BLEEP, Hist2ST, HisToGene, TRIPLEX

	Spot-based					Gigapath-slide	Slide-based			Ciga	STFlow UNI	Gigapath
	Ciga	UNI	Gigapath	STNet DenseNet121	BLEEP ResNet50		Hist2ST ViT	HisToGene ViT	TRIPLEX Ciga			
IDC	0.423 _{.002}	0.502 _{.050}	0.514 _{.064}	0.380 _{.048}	0.346 _{.094}	OOD	0.052 _{.032}	0.350 _{.063}	0.492 _{.042}	0.460 _{.028}	0.589 _{.063}	0.565 _{.055}
PRAD	0.343 _{.001}	0.357 _{.000}	0.386 _{.008}	0.346 _{.006}	0.303 _{.004}	0.386 _{.006}	0.065 _{.038}	0.253 _{.005}	0.351 _{.023}	0.380 _{.001}	0.420 _{.005}	0.415 _{.013}
PAAD	0.406 _{.008}	0.424 _{.060}	0.436 _{.054}	0.370 _{.047}	0.347 _{.059}	0.394 _{.041}	0.111 _{.004}	0.303 _{.007}	0.429 _{.045}	0.440 _{.047}	0.506 _{.078}	0.513 _{.063}
SKCM	0.492 _{.003}	0.613 _{.020}	0.578 _{.001}	0.385 _{.054}	0.407 _{.130}	0.543 _{.014}	0.195 _{.010}	0.321 _{.028}	0.576 _{.091}	0.608 _{.072}	0.707 _{.028}	0.651 _{.089}
COAD	0.275 _{.054}	0.287 _{.005}	0.287 _{.008}	0.249 _{.063}	0.172 _{.014}	OOD	0.071 _{.006}	0.266 _{.015}	0.305 _{.004}	0.344 _{.023}	0.328 _{.013}	0.325 _{.023}
READ	0.051 _{.005}	0.162 _{.080}	0.151 _{.081}	0.116 _{.032}	0.098 _{.063}	0.188 _{.048}	0.034 _{.025}	-0.006 _{.013}	0.129 _{.062}	0.137 _{.075}	0.243 _{.002}	0.260 _{.023}
CCRCC	0.136 _{.005}	0.186 _{.050}	0.187 _{.062}	0.213 _{.071}	0.107 _{.023}	0.183 _{.052}	0.100 _{.053}	0.112 _{.036}	0.229 _{.036}	0.250 _{.054}	0.335 _{.070}	0.326 _{.065}
HCC	0.042 _{.001}	0.051 _{.000}	0.054 _{.002}	0.078 _{.034}	0.066 _{.021}	0.026 _{.005}	0.015 _{.001}	0.028 _{.015}	0.044 _{.022}	0.105 _{.030}	0.128 _{.017}	0.125 _{.019}
LUNG	0.544 _{.001}	0.511 _{.030}	0.568 _{.038}	0.526 _{.025}	0.476 _{.021}	0.530 _{.025}	0.302 _{.063}	0.477 _{.057}	0.563 _{.036}	0.584 _{.027}	0.608 _{.021}	0.602 _{.013}
LYMPH	0.235 _{.006}	0.234 _{.050}	0.275 _{.049}	0.237 _{.063}	0.204 _{.061}	0.284 _{.042}	0.096 _{.079}	0.238 _{.062}	0.286 _{.055}	0.307 _{.052}	0.305 _{.056}	0.305 _{.053}
Average	0.305	0.347	0.344	0.290	0.252	/	0.104	0.234	0.340	0.361	0.419	0.409

Summary

- Task: gene expression prediction based on histology image
- Repurpose one-step regression task as a generative task
 - Apply flow matching with a novel transformer architecture as denoiser
 - Introduce ZINB distribution as prior, which is biologically informed
- Achieve best performance over all the baselines on 10 datasets
 - STFlow is foundation model-agnostic, which can consistently achieve better performance