

Learning Relative Gene Expression Trends from Pathology Images in Spatial Transcriptomics

Kazuya Nishimura^{1,2}, Haruka Hirose², Ryoma Bise³,
Kaito Shiku³, Yasuhiro Kojima²

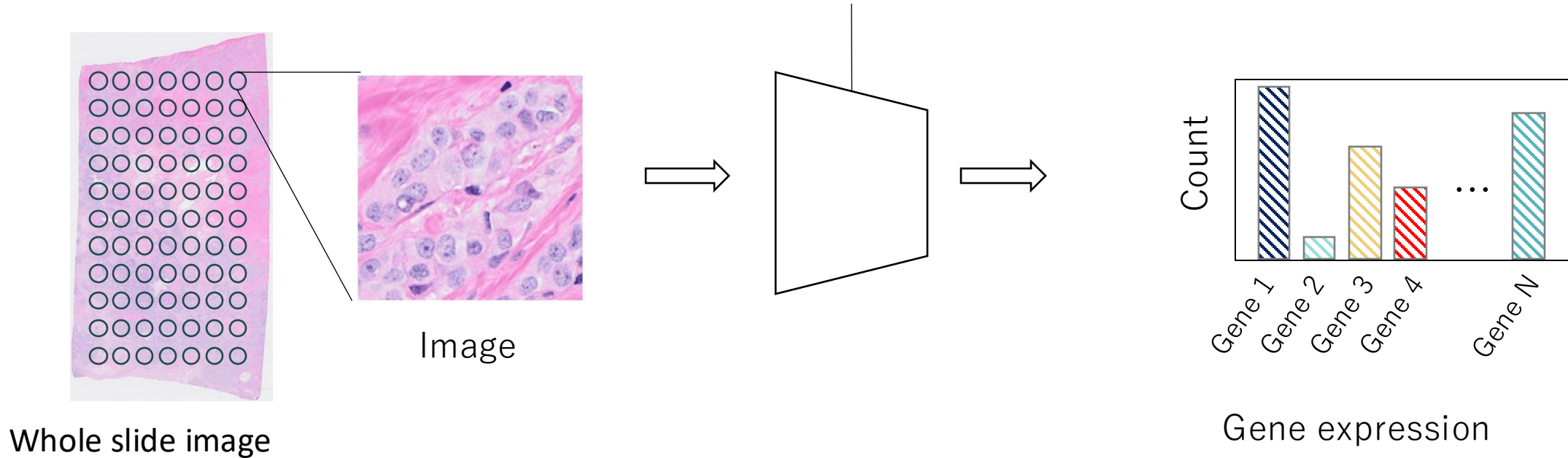
1 Osaka University,

2 National Cancer Center Japan,

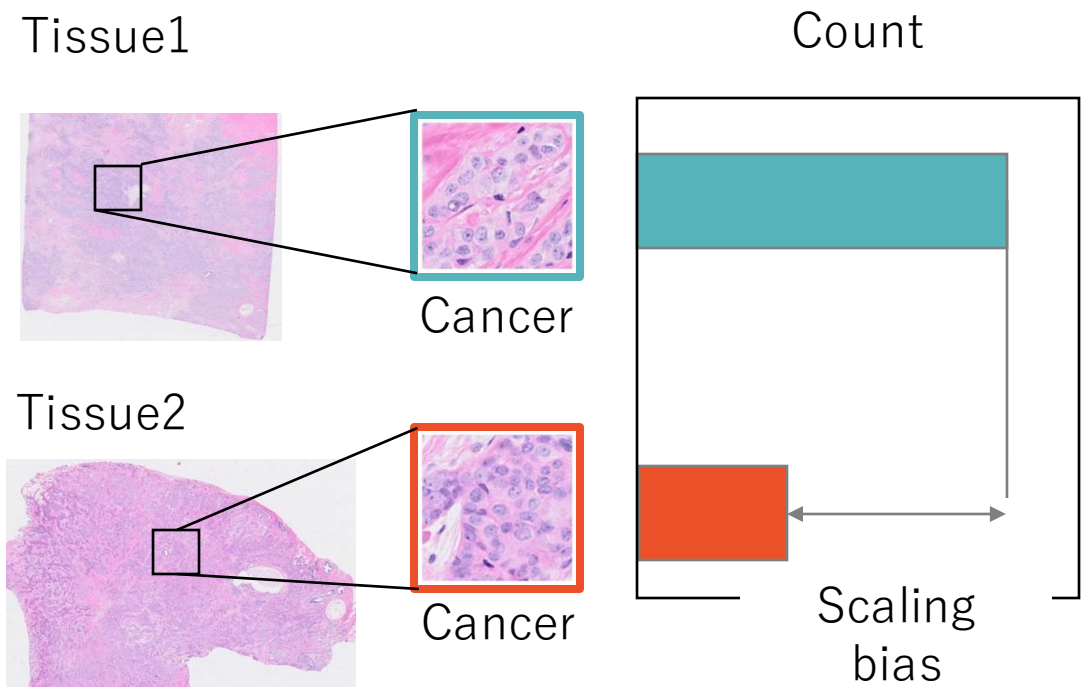
3 Kyushu University



CNN, GNN, Transformer, etc.

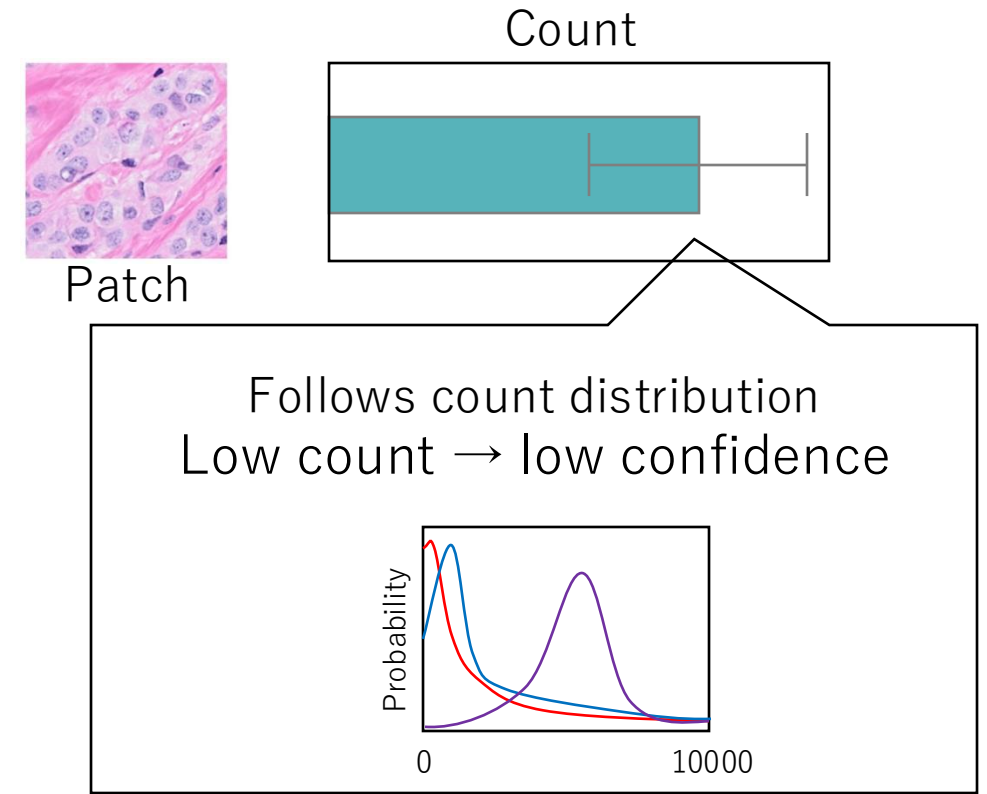


Scaling bias (Batch effect)

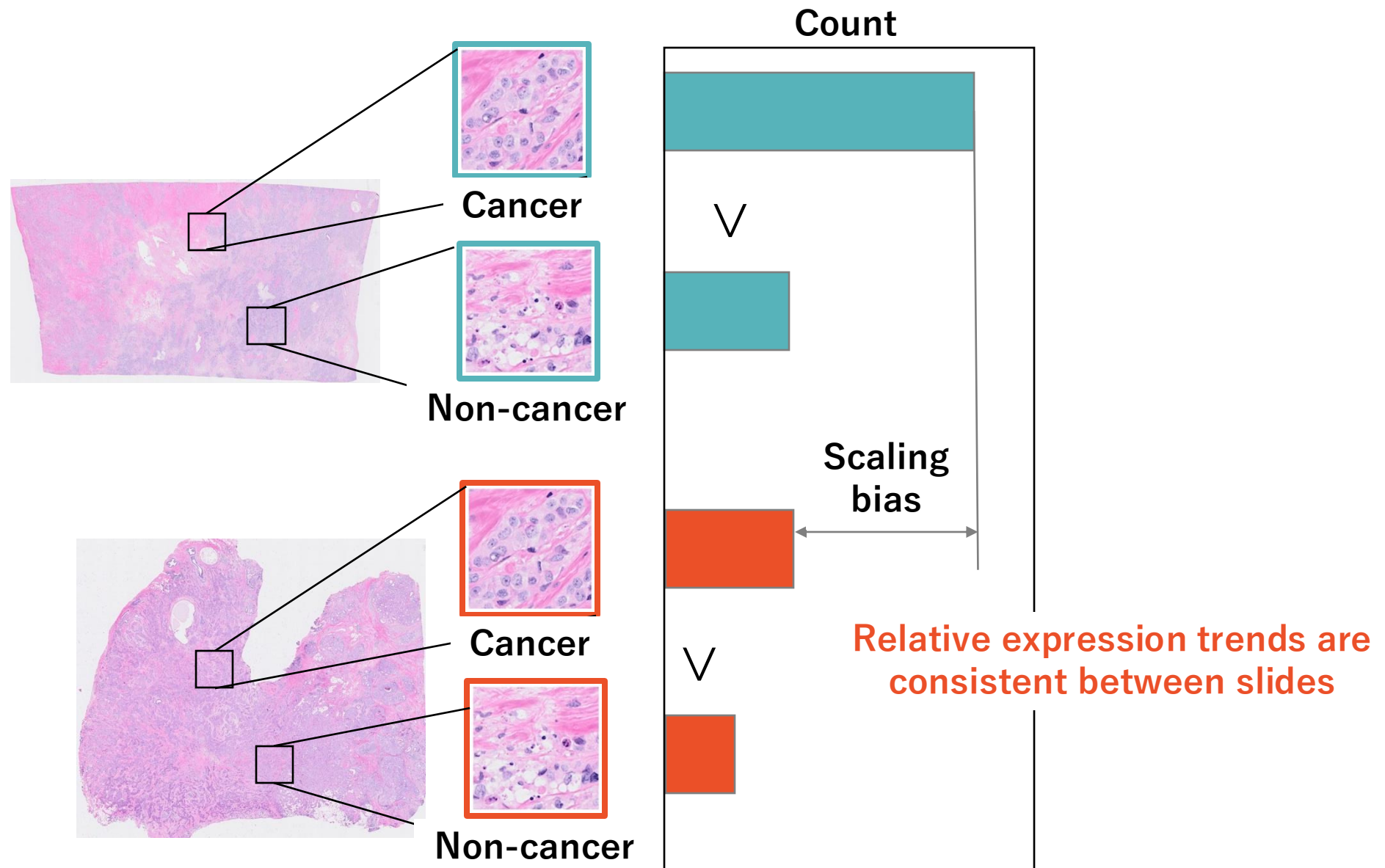


Stochastic noise

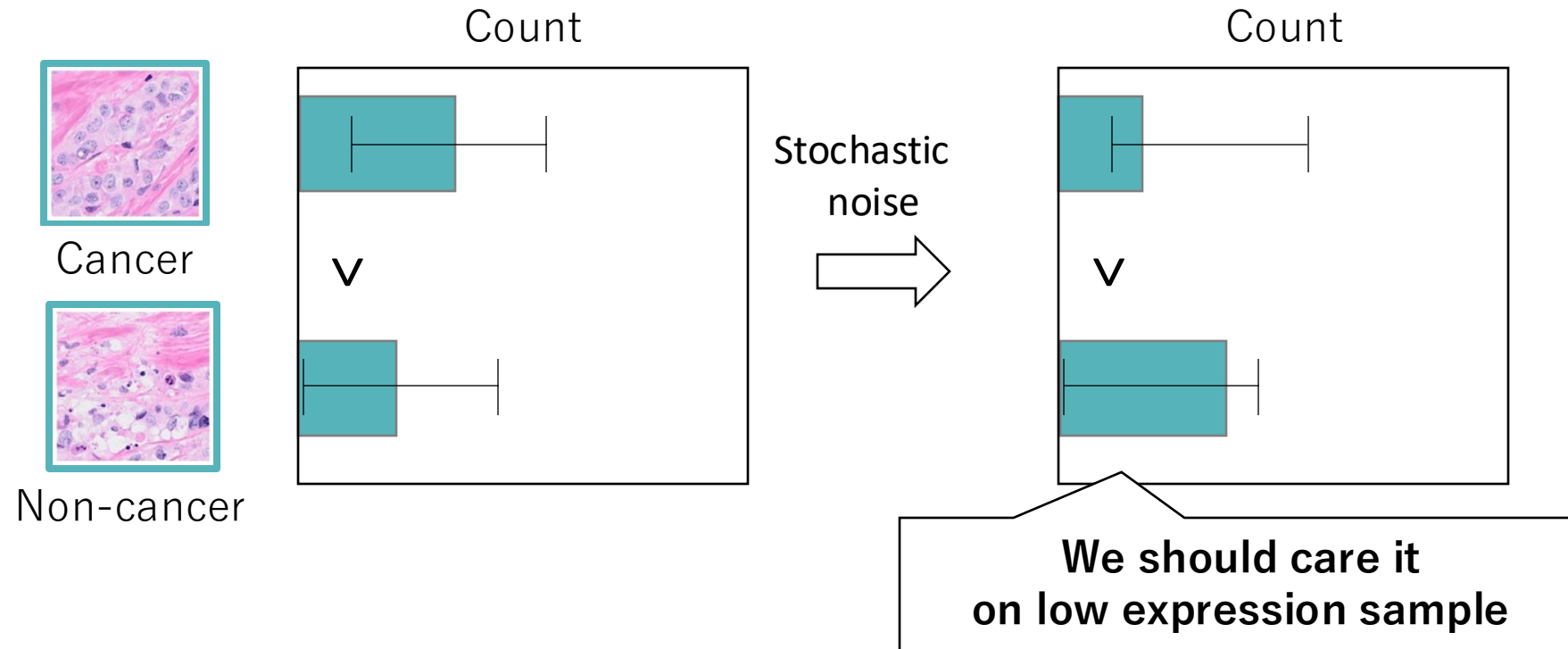
- Temporal dynamics
- Heterogeneity of cell
- Follows count dist.



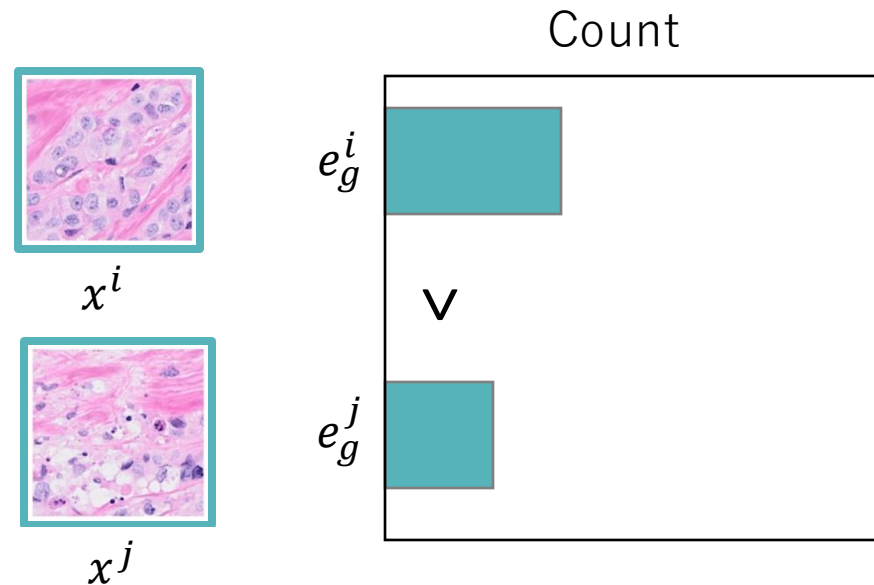
Learning to relative expression trends



Stochastic noise may flip relation



Models gene expression as a discrete probabilistic distribution!



$$P(e^i | x^i, x^j, e_g^i + e_g^j) = \prod_{g=1}^{N_g} P(e_g^i | x^i, x^j, e_g^i + e_g^j)$$

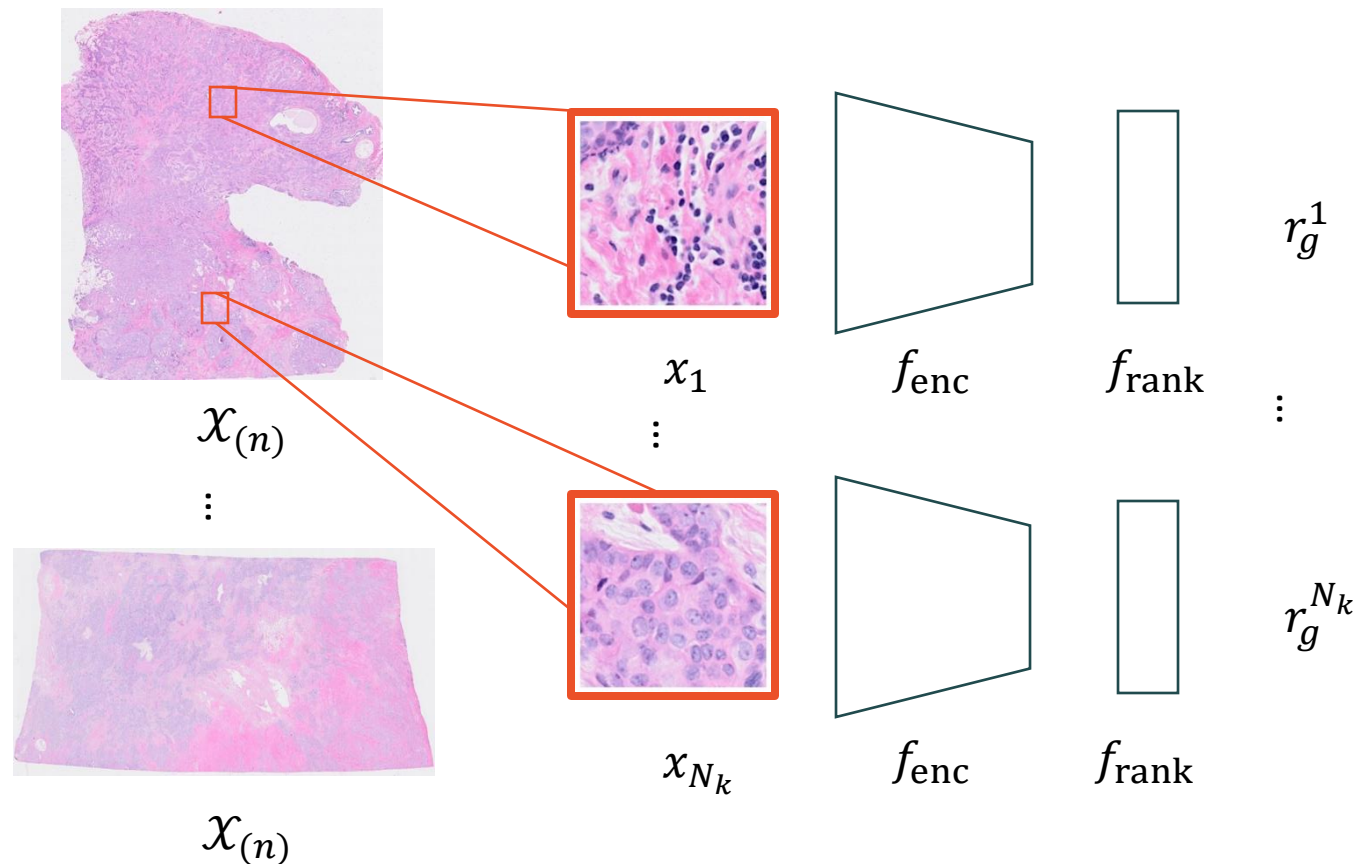
$$P(e_g^i | x^i, x^j, e_g^i + e_g^j) = \text{Binomial}(e_g^i + e_g^j, p_g^{i|j})$$

Assume Binomial distribution by
Viewed as consequence of the counting process given relative frequencies

Learning to relative expression with modeling

Aim: estimate rank score r_g^i which reflect expression relation

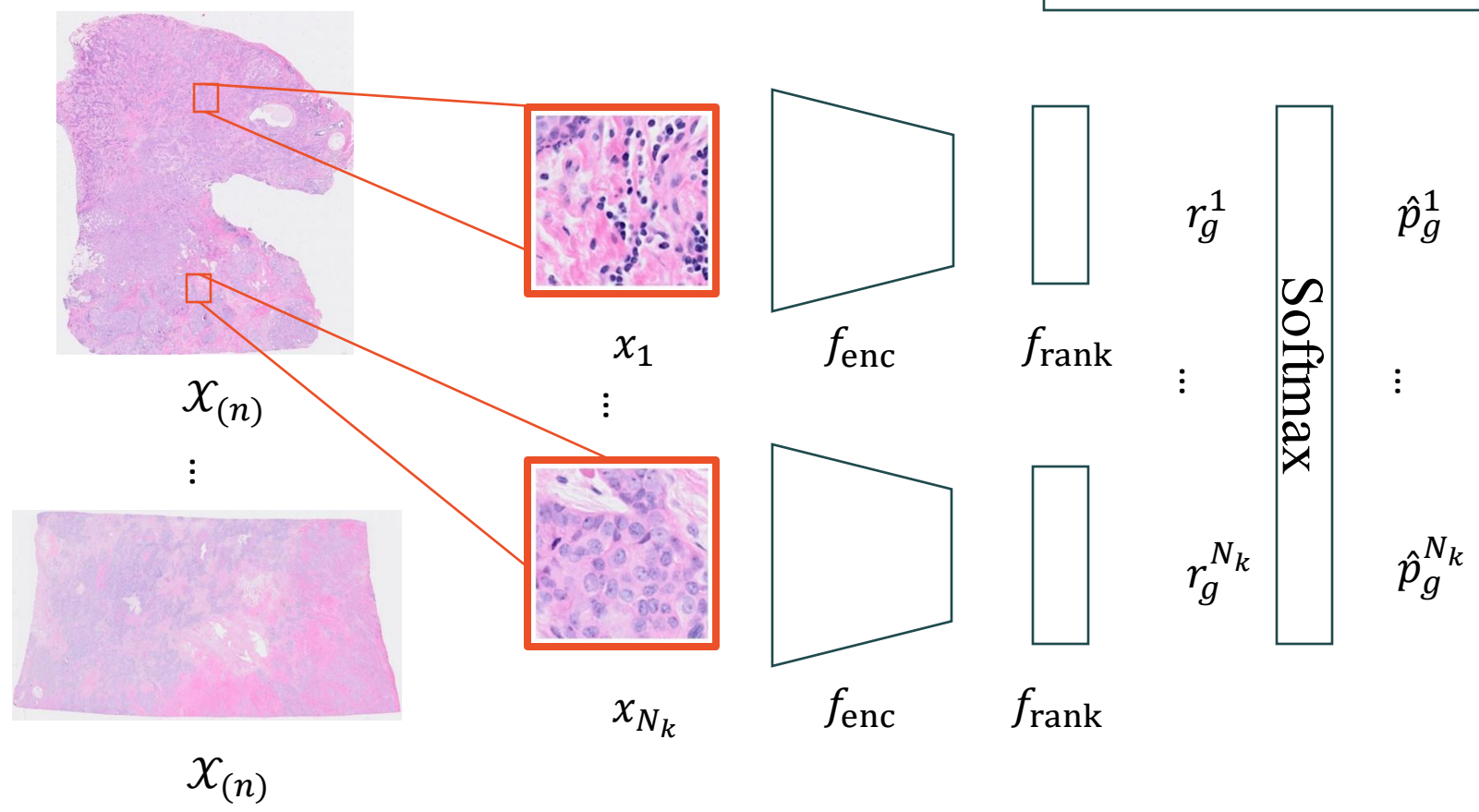
If $e_g^i > e_g^j$, r_g^i should $r_g^i > r_g^j$



Learning to relative expression with modeling

$$\text{Loss} = -\log P(e^i | x^i, x^j, e_g^i + e_g^j) = \sum_{g=1}^{N_g} P(e_g^i | x^i, x^j, e_g^i + e_g^j)$$
$$= -\sum_{g=1}^{N_g} (e_g^i \log \hat{p}_g^i + e_g^j \log \hat{p}_g^j)$$

$P(e_g^i | x^i, x^j, e_g^i + e_g^j) = \text{Binomial}(e_g^i + e_g^j, \hat{p}_g^i)$



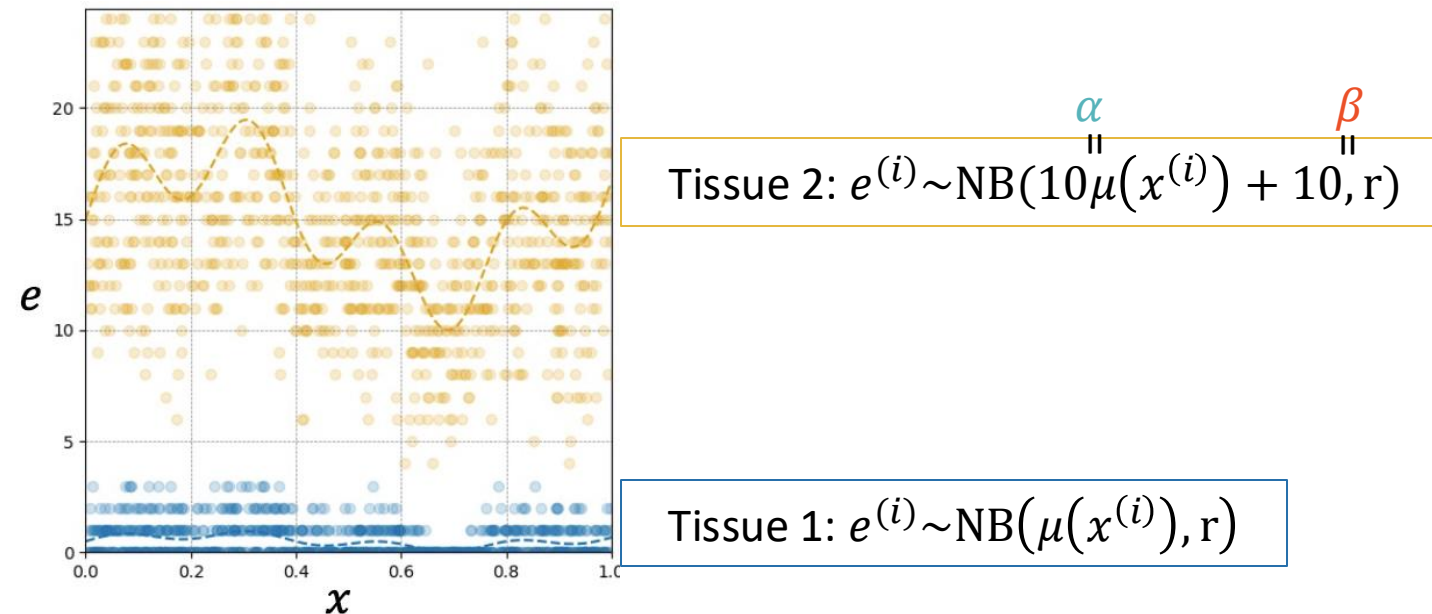
Experiments on synthetic 1d data to understand characteristic of loss function

Gene expression follows negative binomial distribution

$$e^{(i)} \sim \text{NB}(\mu(x^{(i)}), r)$$
$$\mathcal{D} = \{x^{(i)}, e^{(i)}\}$$

Due to batch effect the signal will be change:

Scaling factor $\alpha \mu(x^{(i)}) + \beta$ additive factor



Experiments on synthetic 1d data to understand characteristic of loss function

	Loss	Uniform	Imbalanced	
Point	MSE	0.748	0.583	<u>Not consider relative expression</u>
	Po	0.777	0.603	
	NB	0.788	0.601	
Pair	Rank	0.835	0.738	<u>Not consider noise</u>
	PairSTrank	<i>0.907</i>	<i>0.818</i>	
List	PCC	0.858	0.560	
	ListSTrank	0.945	0.828	

Experiments on real dataset (HEST 1k)

11

Evaluated on 7 types of organ

	Loss	IDC	PRAD	PAAD	COAD	READ	ccRCC	IDC-L	Ave.
Point	MSE	0.393	0.484	0.307	0.556	0.140	0.093	0.168	0.306
	Po	0.314	<i>0.485</i>	0.336	0.524	0.172	0.091	0.134	0.293
	NB	0.199	0.491	0.119	0.538	<i>0.160</i>	0.075	0.126	0.244
Pair	Rank	0.317	0.317	0.181	0.566	0.047	0.059	0.110	0.228
	PairSTrank	<i>0.494</i>	0.458	0.346	<i>0.613</i>	0.136	0.127	<i>0.228</i>	<i>0.343</i>
List	PCC	0.472	0.459	0.307	0.640	0.105	0.102	0.198	0.326
	ListSTrank	0.510	0.459	<i>0.343</i>	0.597	0.140	<i>0.125</i>	0.238	0.345

Task. Patch-level gene expression estimation

Difficulty. batch effect and stochastic noise of observed data

Idea.

1. Batch effect
 - > learning to relative expression
2. Stochastic noise
 - > Model relative expression with discrete distribution