



南京航空航天大学

# Discovering Fine-Grained Visual-Concept Relations by Disentangled Optimal Transport Concept Bottleneck Models

Yan Xie<sup>1\*</sup> Zequn Zeng<sup>1\*</sup> Hao Zhang<sup>1†</sup> Yucheng Ding<sup>1</sup> Yi Wang<sup>1</sup> Zhengjue Wang<sup>2</sup> Bo Chen<sup>1</sup>  
Hongwei Liu<sup>1</sup>

<sup>1</sup>National Key Laboratory of Radar Signal Processing, Xidian University, Xi'an, 710071, China

<sup>2</sup>State Key Laboratory of Integrated Service Networks, Xidian University, Xi'an, 710071, China

CVPR 2025

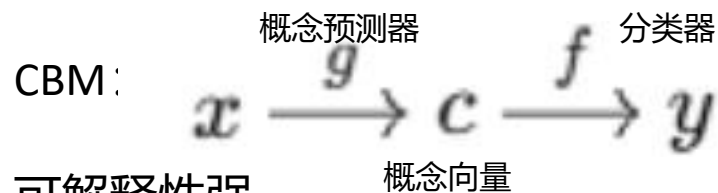
# Introduction



南京航空航天大学

普通CNN:

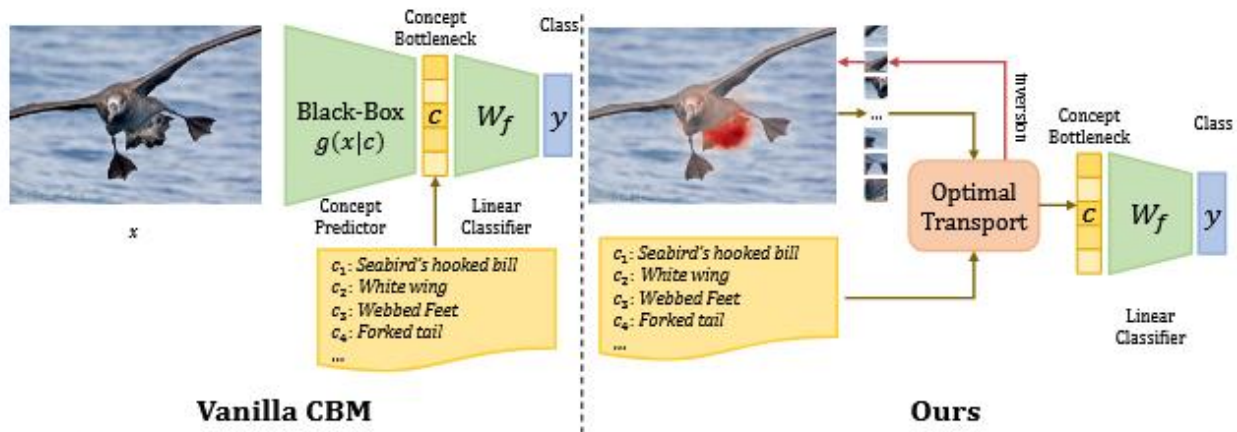
图片  $\rightarrow$  卷积层  $\rightarrow$  分类结果



可解释性强

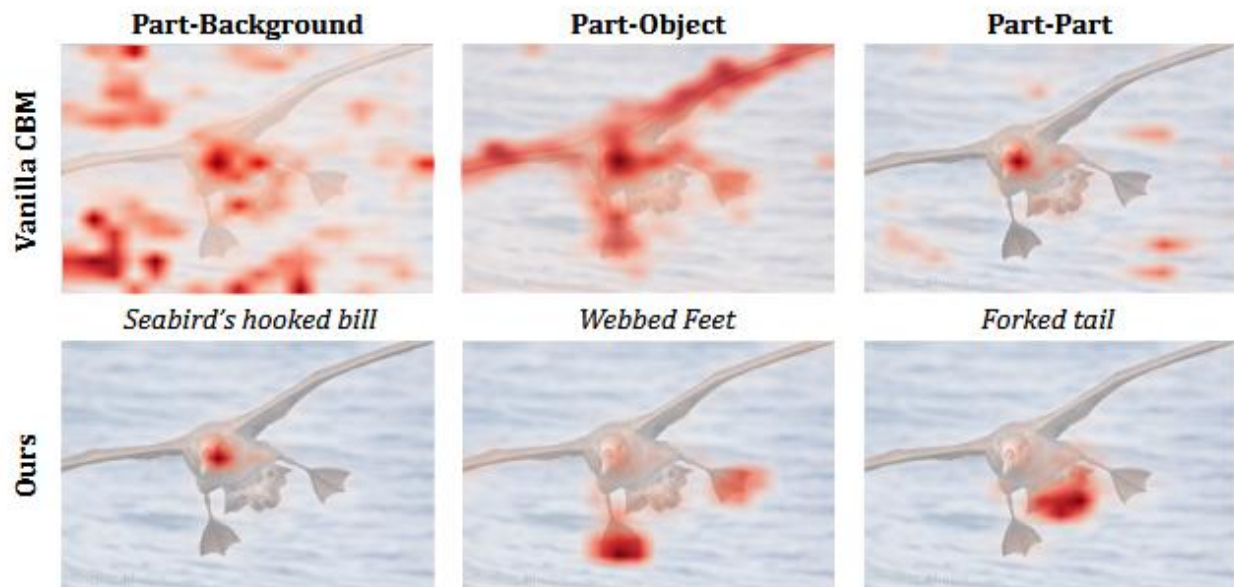
问题 1: 概念标注成本极高

问题 2: 概念-输入不对齐



(a) Brief comparison about the coarse-grained design of vanilla CBM and fine-grained design of our proposed DOT-CBM.

Figure 1. Comparison between vanilla CBM and our proposed DOT-CBM, including architecture design and inversion heatmap for concept predictions. (a) Due to the black-box mapping from images to concepts, Vanilla CBM needs Grad-CAM [37] techniques to locate concept predictions back to image space while DOT-CBM can provide an explicit inversion heatmap visualization. (b) Due to lack of fine-grained alignment, vanilla CBM produces spurious correlations in three levels of granularity (part-background, part-object and part-part level) that mislocalizes the local concept to the background, the whole object, and incorrect local region.



(b) Inversion heatmap visualization.

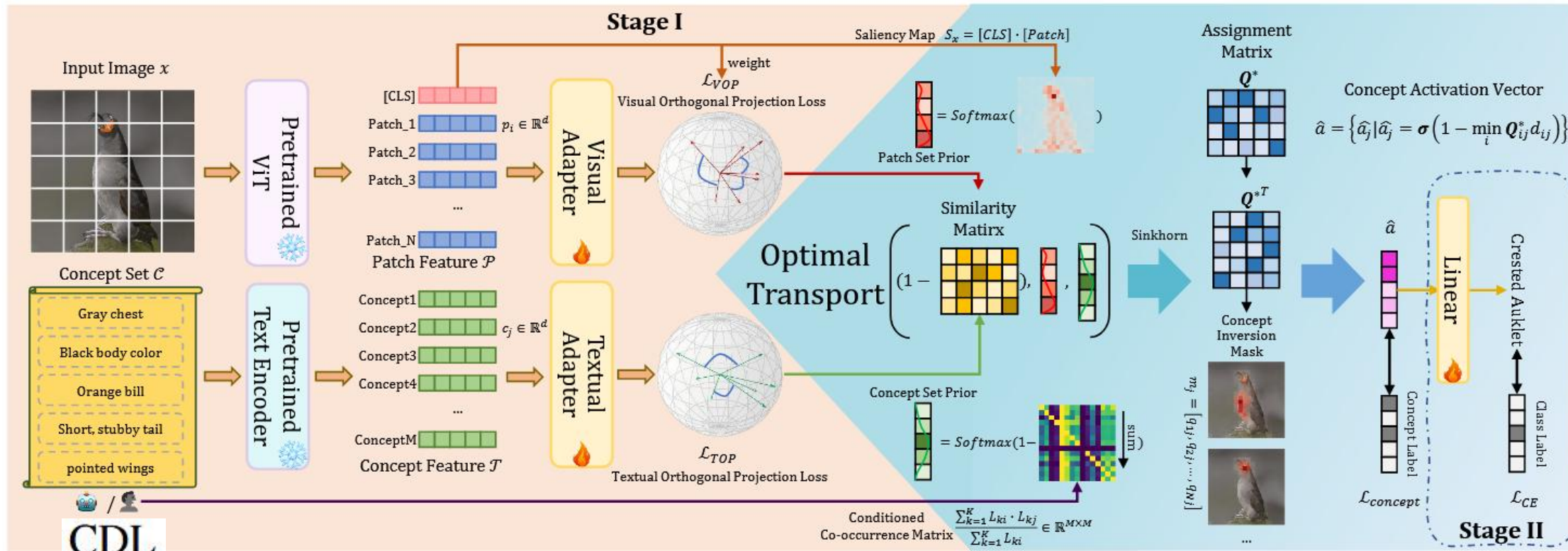


Figure 2. Overview of our proposed DOT-CBM. The overall CBM framework proceeds from left to right. In the first stage, the model transforms input images into concept activation vectors. In forward pass, local embeddings from a pre-trained Vision Transformer (ViT) and a text encoder are processed through learnable adapters to generate two feature sets (Sec. 3.2). Two orthogonal projection losses are applied to both feature sets to constrain Adapter training. An Optimal Transport (OT) algorithm optimizes the Assignment Matrix, representing the explicit correlation between concepts and image patches. Concept activation values are derived by combining this matrix with a cost matrix, and the loss is supervised by concept labels for training (Sec. 3.3). To address data bias, we use the Saliency Map from the pre-trained ViT for rough foreground-background separation as the visual prior in OT. A Conditioned Co-occurrence Matrix, based on frequency statistics, serves as the prior for concepts, enhancing the model’s ability to distinguish co-occurring concepts (Sec. 3.4). In the second stage, consistent with the vanilla CBM framework, concept activation values are used to predict final class labels via a linear classification network, which is trained using class labels (Sec. 3.5).

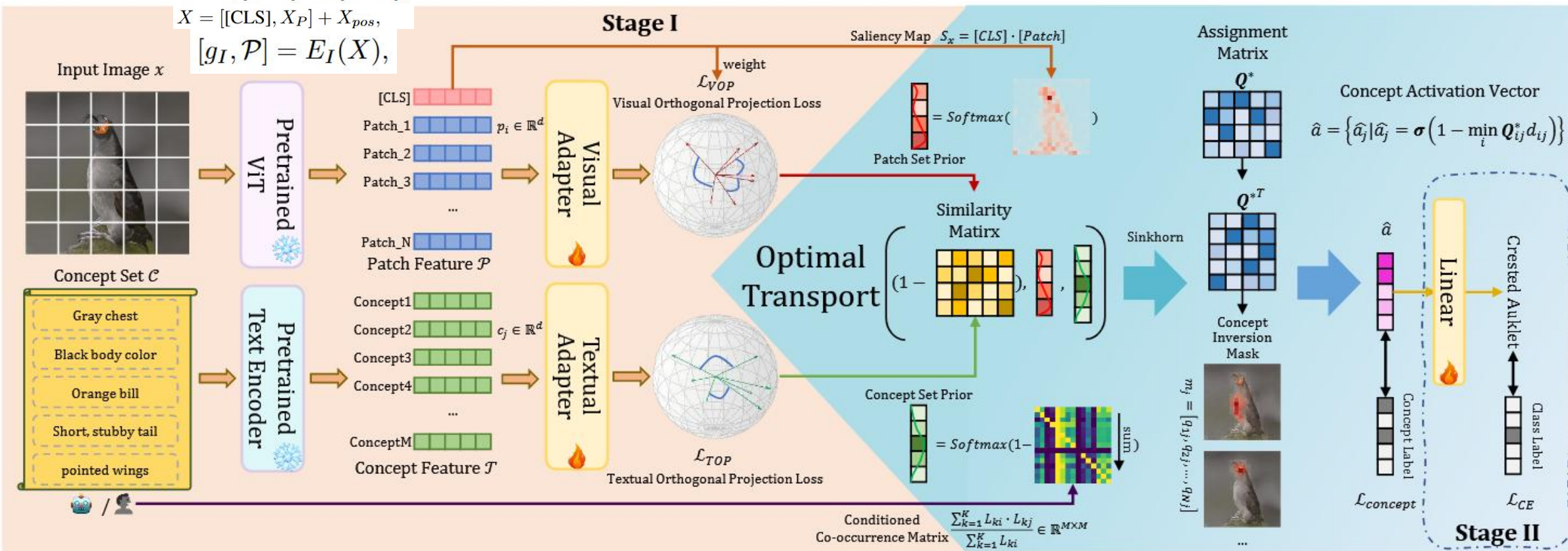
# Method Disentanglement on Local representations



$$X_P = [I_{p1} \cdot W_p, \dots, I_{pN} \cdot W_p],$$

$$X = [[CLS], X_P] + X_{pos},$$

$$[g_I, \mathcal{P}] = E_I(X),$$



$$S_I = g_I \cdot P,$$

$$\langle A_v(p_i), A_v(p_j) \rangle = \frac{A_v(p_i) \cdot A_v(p_j)}{\|A_v(p_i)\| \|A_v(p_j)\|},$$

$$\mathcal{L}_{VOP} = \mathbb{E}_{(p_i, p_j) \sim \mathcal{P}} [S_I(i) \cdot S_I(j) \cdot \langle A_v(p_i), A_v(p_j) \rangle]$$

$$\mathcal{T} = \{t_j = E_T(c_j)\}_{j=1}^M$$

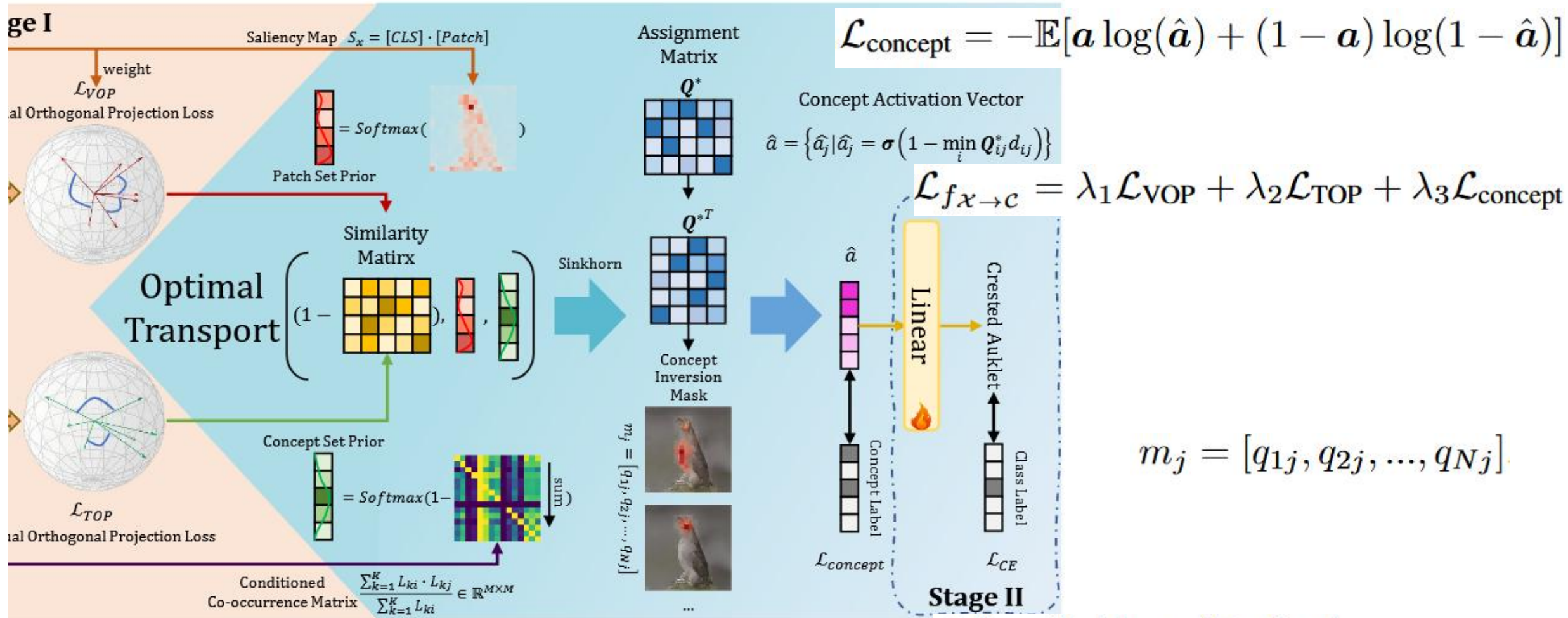
$$\mathcal{L}_{TOP} = \mathbb{E}_{(t_i, t_j) \sim \mathcal{T}} [\langle A_t(t_i), A_t(t_j) \rangle],$$

$$\text{Conditioned Co-occurrence Matrix } \frac{\sum_{k=1}^K L_{ki} \cdot L_{kj}}{\sum_{k=1}^K L_{ki}} \in \mathbb{R}^{M \times M}$$

$$m_j = [q_{1j}, q_{2j}, \dots, q_{Nj}]$$

$$\hat{a} = \{\hat{a}_j | \hat{a}_j = \sigma(1 - \min_i Q_{ij}^* d_{ij})\}$$

# Method Patch-concept alignment via optimal transport



$$\Theta = \sum_{i=1}^N \theta_i \delta_{p_i}$$

$$\Gamma = \sum_{j=1}^M \gamma_j \delta_{c_j}$$

$$Q = [q_{ij}]_{N \times M} \quad d_{ij} = d(p_i, t_j) = (1 - \langle p_i, t_j \rangle) \geq 0$$

$$\min_{Q \in \mathcal{Q}} \sum_{i,j} q_{ij} d_{ij} - \varepsilon H(Q)$$

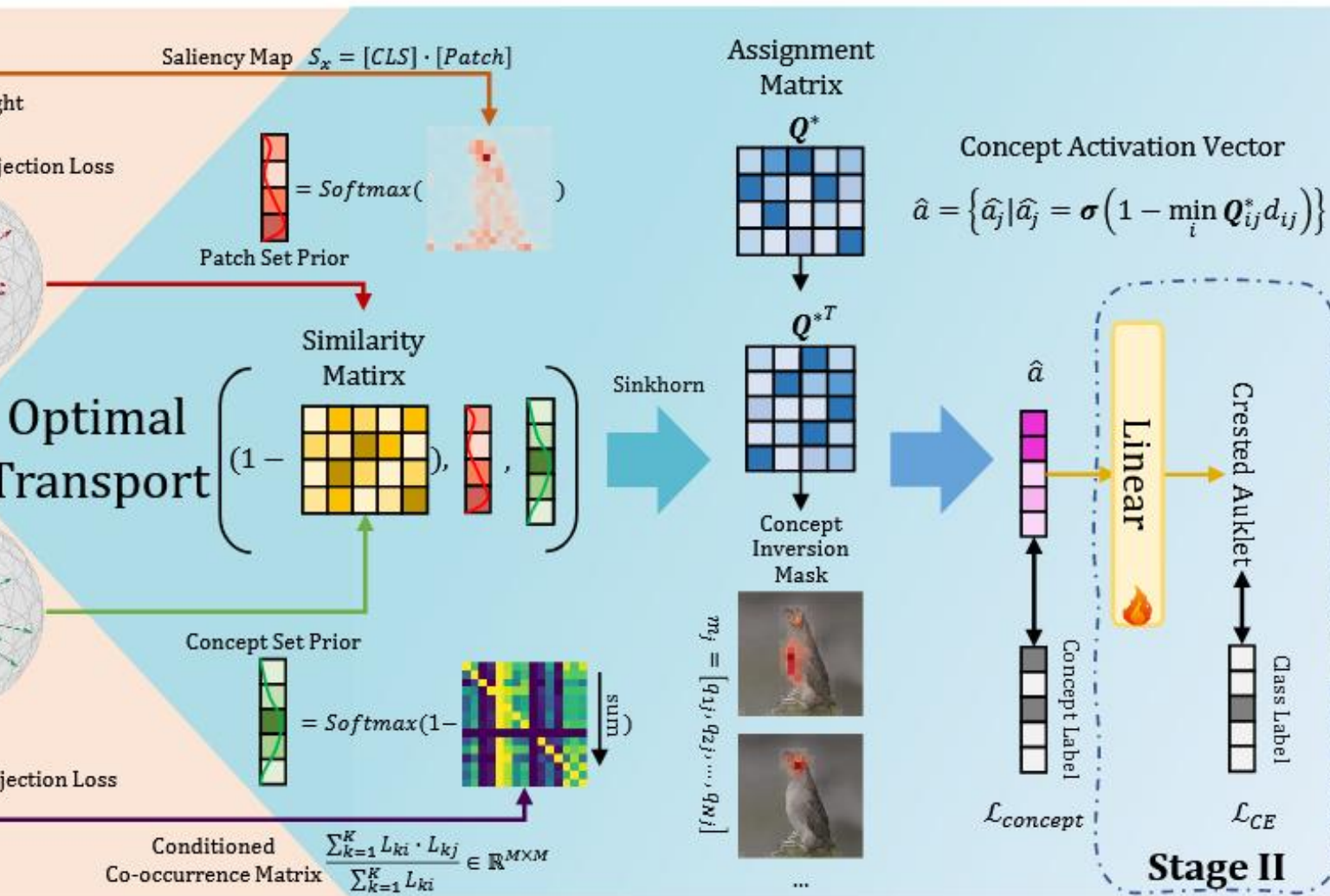
$$\text{s.t. } \mathcal{Q} = \{Q \in \mathbb{R}_+^{N \times M} \mid Q \mathbf{1}_M = \theta, Q^T \mathbf{1}_N = \gamma\}$$

$$Q^* = \text{Sinkhorn}(D, \theta, \gamma)$$

$$d_{OT} = \sum_{i,j} q_{ij} d_{ij}$$

$$\hat{a} = \{\hat{a}_j | \hat{a}_j = \sigma(1 - d_{OT_j})\}_{j=1}^M$$

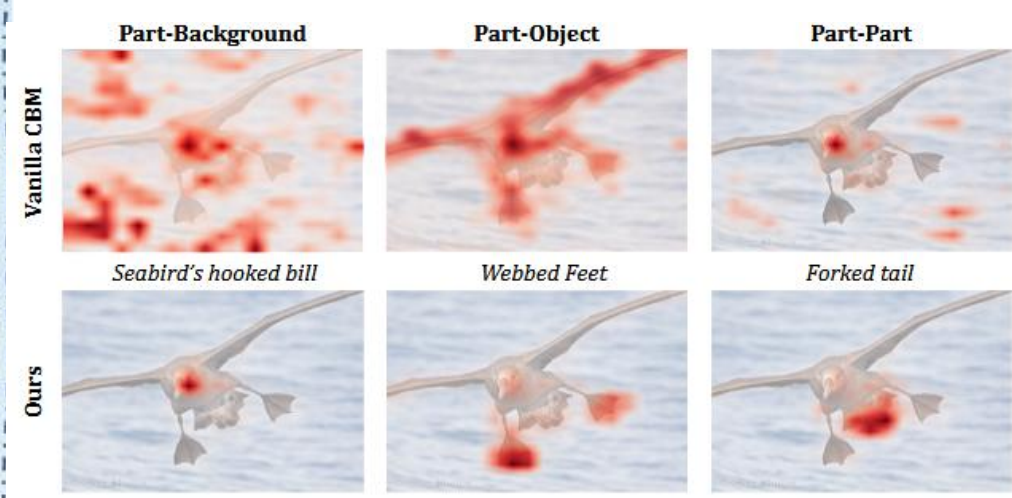
# Method Prior distributions Modeling



$$\Theta = \sum_{i=1}^N \theta_i \delta_{p_i}, \quad \Gamma = \sum_{j=1}^M \gamma_j \delta_{c_j},$$

$$S_I = g_I \cdot P,$$

$\theta = \text{Softmax}(S_I)$ , Patch set prior  $\theta$



(b) Inversion heatmap visualization.

$$r_i = \sum_{j=1}^M \left( \frac{\sum_{k=1}^K L_{ki} \cdot L_{kj}}{\sum_{k=1}^K L_{ki}} \right) \gamma = 1 - \text{Softmax}(r_i)$$

class-concept label form a matrix  $L = [l_{i,j}]_{K \times M}$

Method	Classification Accuracy ( $\uparrow$ )				Part Detection $mAP_{0.5}$ ( $\uparrow$ )			
	ImageNet	CUB	CIFAR100	AWA2	PartImageNet	CUB	RIVAL	PASCAL-Parts
Vanilla-CBM [17]	79.17	78.32	80.04	93.15	26.94	27.02	26.76	17.32
CEM [48]	81.29	80.47	81.23	95.92	29.19	30.86	30.83	20.85
LaBo [45]	82.93	81.30	84.10	96.92	38.27	39.28	40.34	30.17
SparseCBM [38]	82.85	82.07	84.75	95.56	39.84	40.23	40.68	33.71
CoopCBM [39]	82.73	82.10	84.66	<b>97.08</b>	41.17	45.16	43.82	35.75
<b>DOT-CBM</b>	<b>83.84</b>	<b>85.39</b>	<b>85.83</b>	96.83	<b>50.12</b>	<b>53.47</b>	<b>50.93</b>	<b>44.18</b>
	(+0.91)	(+3.29)	(+1.08)	(-0.25)	(+8.95)	(+8.31)	(+7.11)	(+8.43)

Table 1. Performance comparison of Classification Accuracy (%) and Part Detection ( $mAP_{0.5}$ ).

Method	CUB		Dogs	
	ID	OOD	ID	OOD
Vanilla-CBM [17]	86.9	27.7	87.3	29.4
CEM [48]	84.0	34.0	83.8	36.5
LaBo [45]	85.3	39.4	85.9	41.7
Sparse-CBM [38]	84.5	35.8	85.8	38.6
Coop-CBM [39]	85.8	36.2	86.1	40.3
<b>DOT-CBM w/o. prior</b>	86.5	36.9	87.0	40.8
<b>DOT-CBM w. prior</b>	<b>87.2</b>	<b>49.7</b>	<b>89.1</b>	<b>52.4</b>
	(+0.3)	(+10.3)	(+1.8)	(+10.7)

Table 2. Performance comparison of OOD generalization.

OT	Prior	LrD	Classification Accuracy ( $\uparrow$ )				Part Detection $mIOU$ ( $\uparrow$ )			
			ImageNet	CUB	CIFAR100	AWA2	PartImageNet	CUB	RIVAL	PASCAL-Parts
$\times$	$\times$	$\times$	81.60	80.92	82.73	96.07	0.35	0.41	0.37	0.28
$\checkmark$	$\times$	$\times$	82.03	81.75	82.91	96.24	0.42	0.49	0.42	0.34
$\checkmark$	$\checkmark$	$\times$	83.22	84.92	85.49	96.47	0.50	0.61	0.53	0.46
$\times$	$\times$	$\checkmark$	82.73	82.10	84.66	97.08	0.47	0.58	0.49	0.41
$\checkmark$	$\checkmark$	$\checkmark$	83.84	85.39	85.83	96.83	0.52	0.66	0.54	0.49

Table 3. Ablation studies on image classification and part detection. OT, Prior, and LrD denote the optimal transport design, prior distribution modeling, and local representation disentanglement.

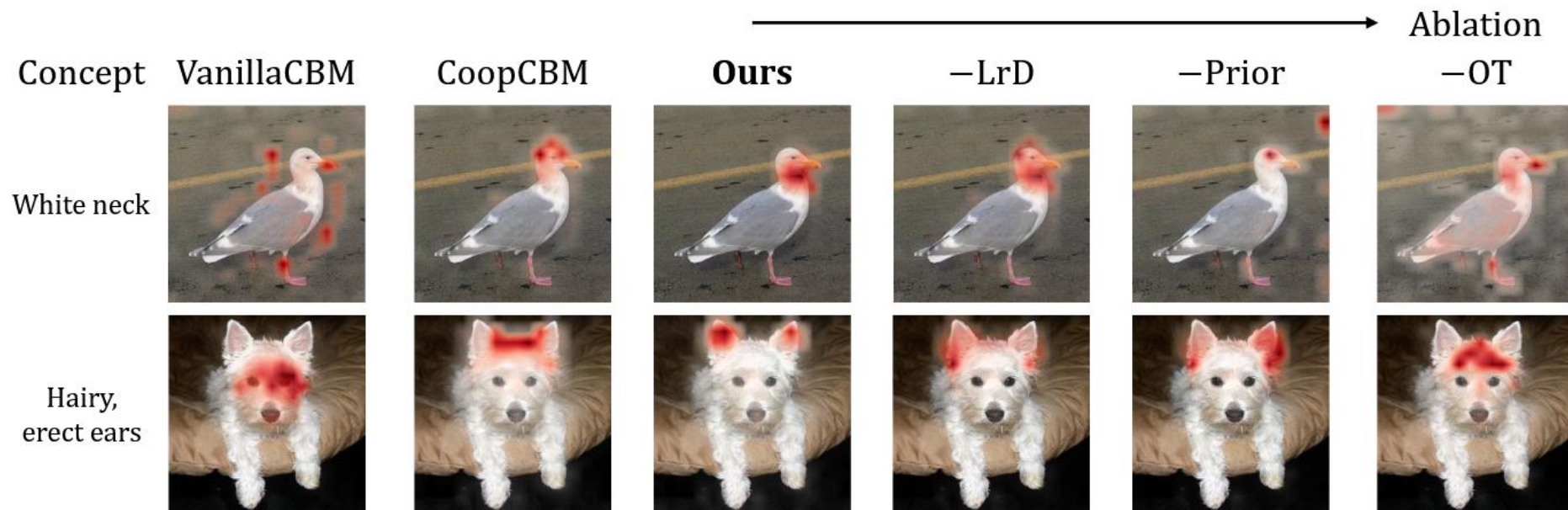


Figure 4. Qualitative results for the concept inversion. VanillaCBM and CoopCBM adopt Grad-CAM to provide a concept inversion heatmap, while our proposed DOT-CBM can generate the inversion mask on its own.



Thanks