

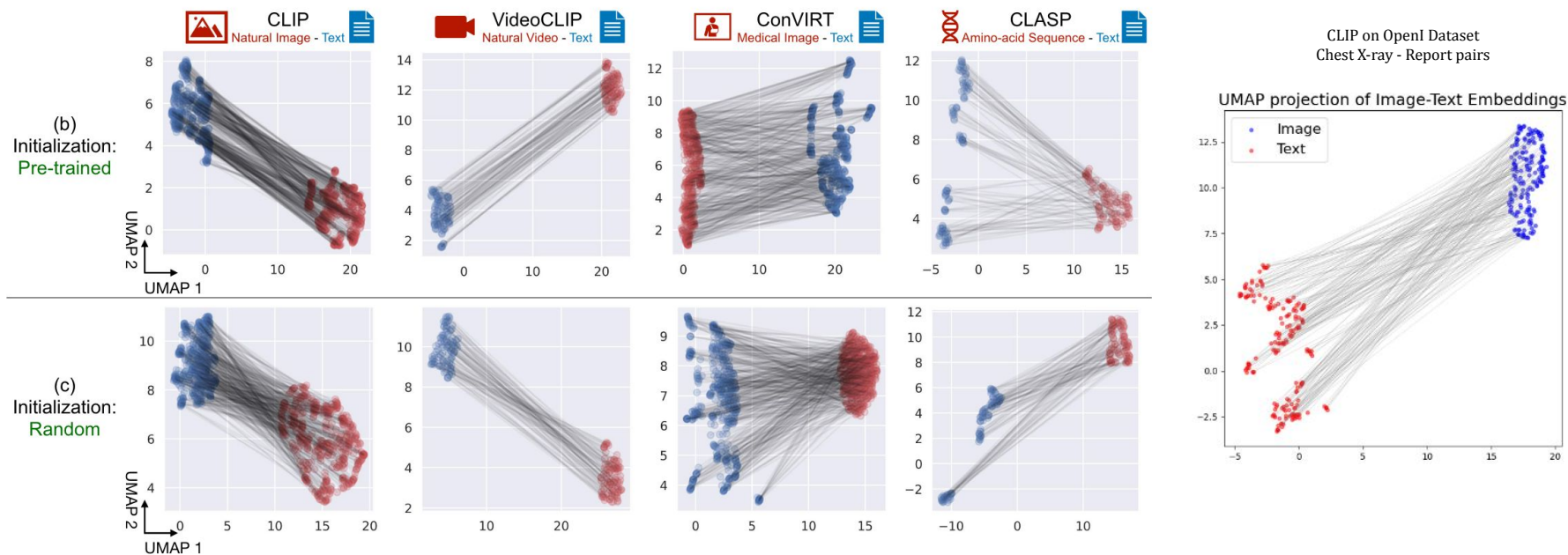
# Improving Medical Multi-modal Contrastive Learning with Expert Annotations

Yogesh Kumar  and Pekka Marttinen 

Department of Computer Science, Aalto University, Finland

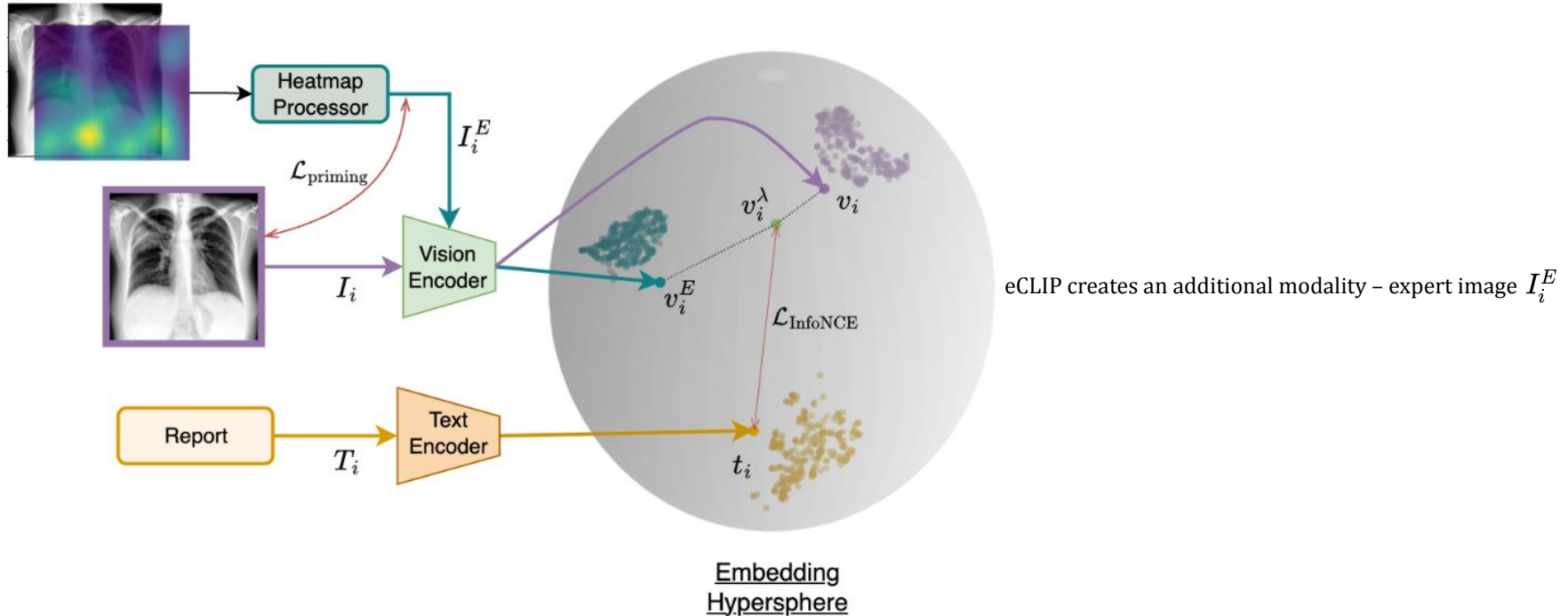
# Modality Gap in CLIP Embeddings

UMAP projection reveals that CLIP embeddings form distinct clusters that pertain to their respective modalities (e.g., images and text)



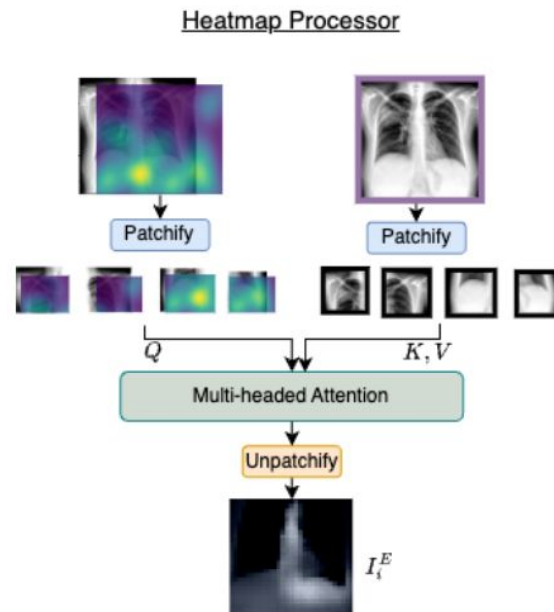
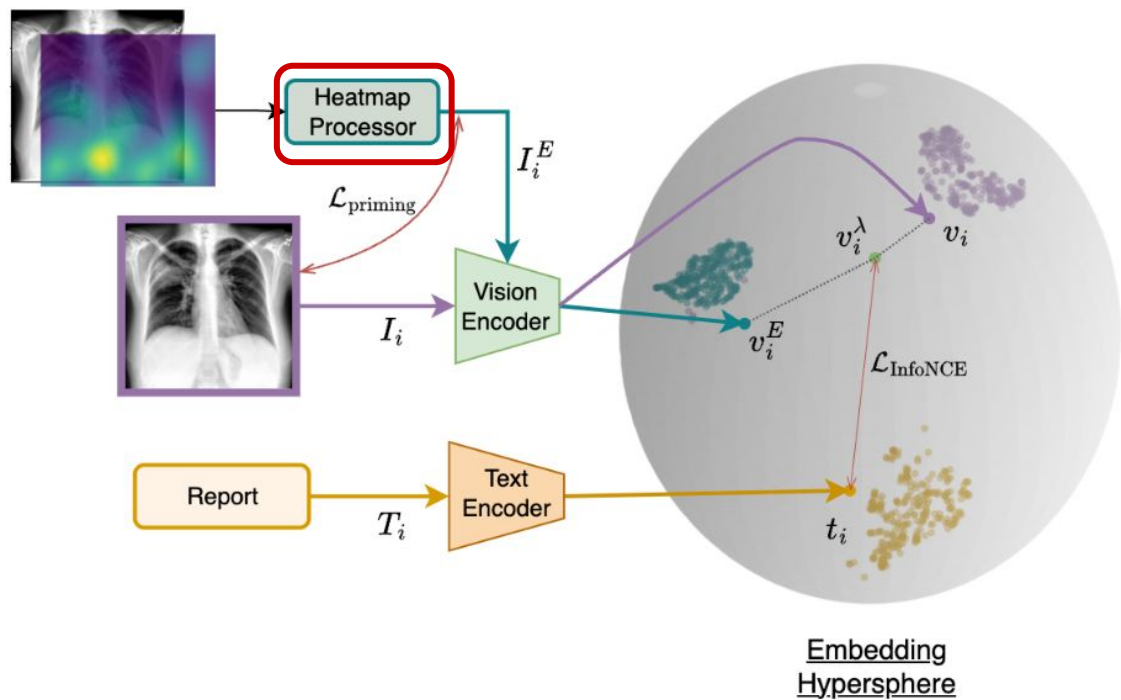
# Expert Annotated CLIP - eCLIP

- For medical data, pairs of chest X-rays and their corresponding radiologist report form the image-text pairs used for CLIP pretraining
- Additionally, we can leverage the heatmaps of radiologist's gaze patterns obtained through eye-tracking



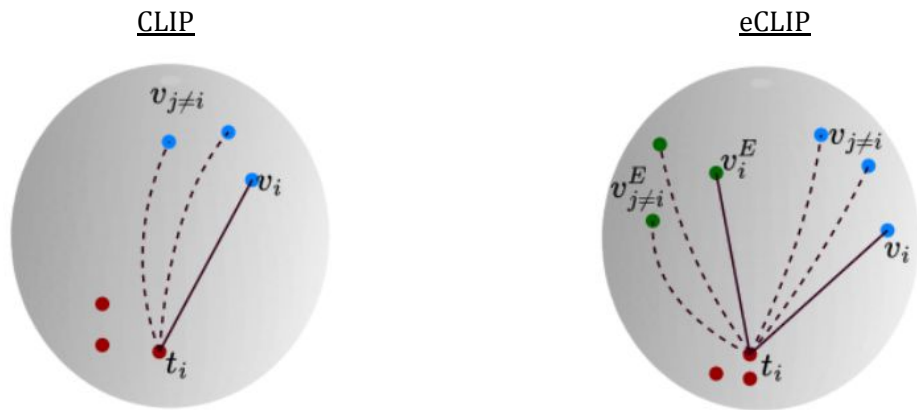
# Heatmap Processor

Heatmap processor is a single layer of multi-headed attention used to fuse the eye-gaze heatmaps with X-ray image, forming the “expert image”  $I_i^E$



# Additional Image-Text pairs for Contrastive Loss

- eCLIP exploits the modality gap to create additional pairs for contrastive loss using the expert image embeddings  $v_i^E$
- Heatmap processor is crucial to ensure that image embeddings  $v_i$  and the expert image embeddings  $v_i^E$  are not too similar



$$\mathcal{L}_{\text{text}} = \mathbb{E}_{(t_i, v_i) \sim \text{pos}} \left[ -\log \frac{\exp(\text{sim}(t_i, v_i)/\tau)}{\exp(\text{sim}(t_i, v_i)/\tau) + \sum_{j \neq i} \exp(\text{sim}(t_i, v_j)/\tau)} \right]$$

# Experiments

## Zero Shot Image Classification (F1-Score)

Model	Dataset			
	Chexpert 5x200	MIMIC 5x200	RSNA	CXR 14x100
CLIP <sub>ViT</sub> Base	0.540 $\pm$ .017	0.465 $\pm$ .004	0.805 $\pm$ .001	0.183 $\pm$ .011
+naive	0.506 $\pm$ .011	0.426 $\pm$ .006	0.805 $\pm$ .004	0.151 $\pm$ .009
+DAKL	0.474 $\pm$ .007	0.400 $\pm$ .002	0.759 $\pm$ .001	0.106 $\pm$ .003
+ $m^3$ -mix	0.542 $\pm$ .021	0.465 $\pm$ .013	0.798 $\pm$ .004	0.183 $\pm$ .020
+expert (ours)	<b>0.563<math>\pm</math>.021</b>	<b>0.477<math>\pm</math>.004</b>	<b>0.814<math>\pm</math>.003</b>	<b>0.193<math>\pm</math>.017</b>

## Pre-training Sample Efficiency

