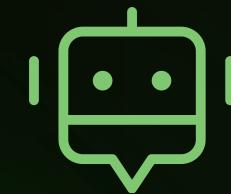


# DSL501



## Region-Aware MedCLIP for Medical Image–Text Understanding

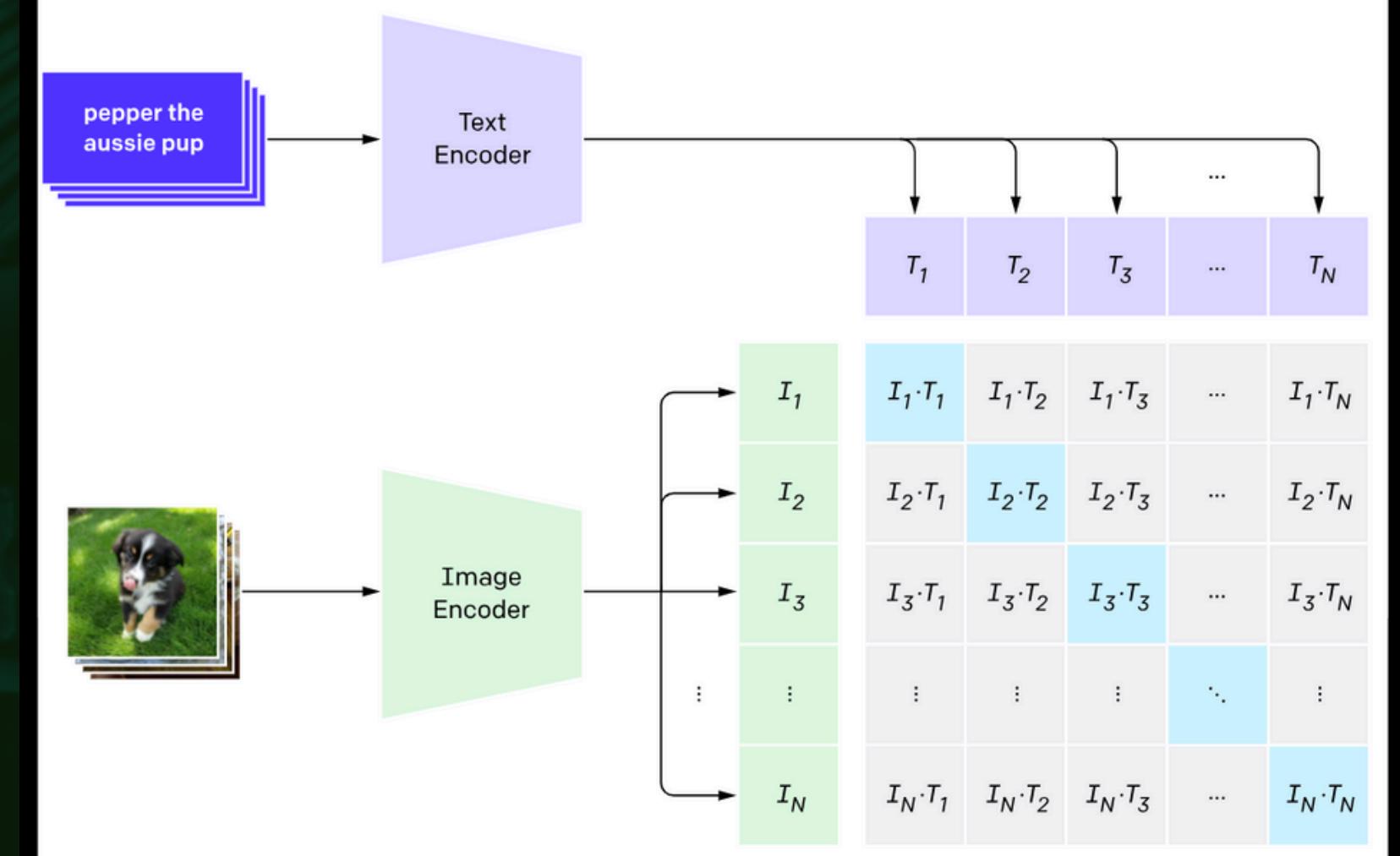
BY S H A S H A N K Y A D A V  
3 R D - Y E A R B . T E C H D S A I

# Introduction

- What is Contrastive Language-Image Pre-Training?
- Why CLIP is good (large-scale, general knowledge)?



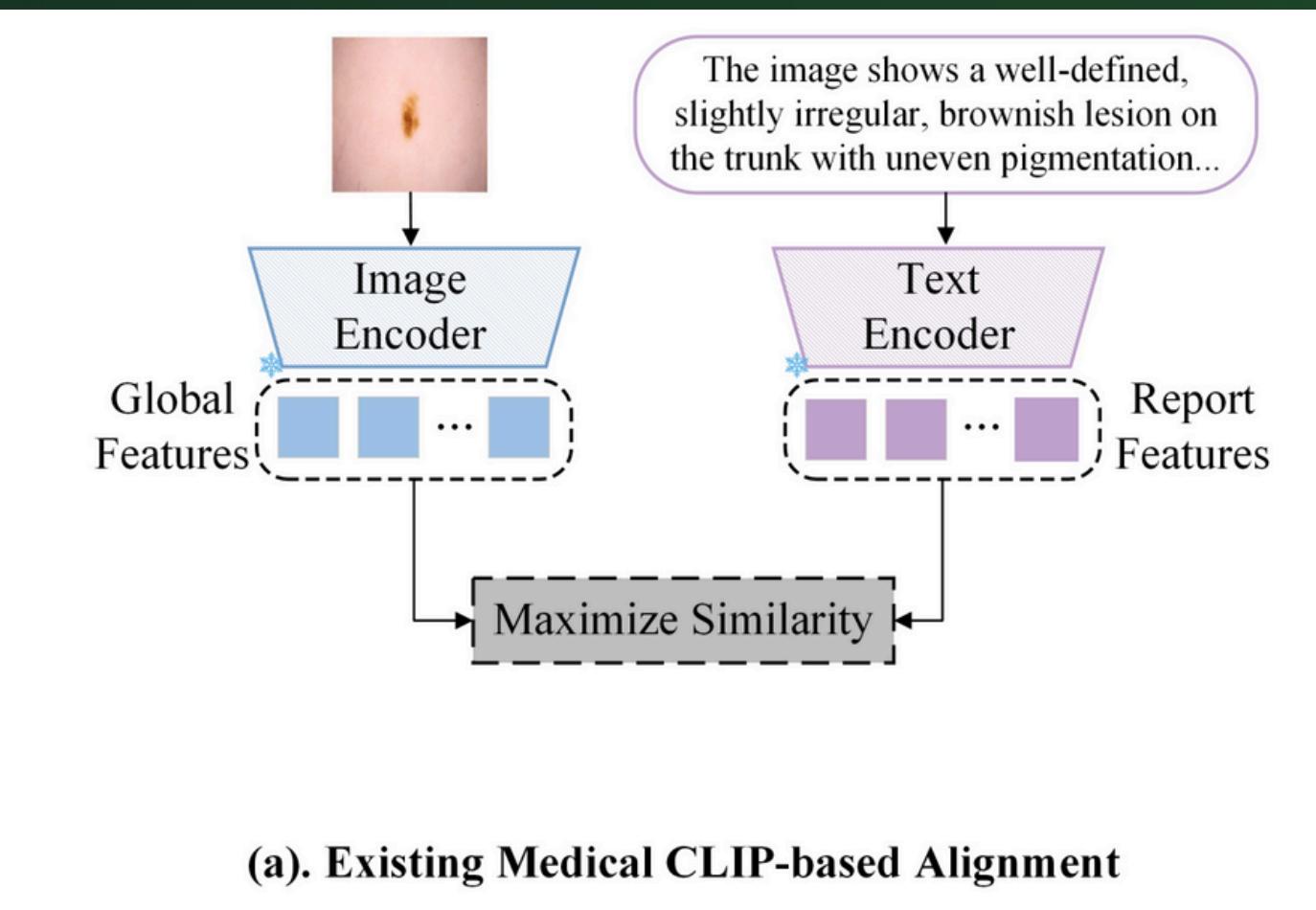
## 1. Contrastive pre-training



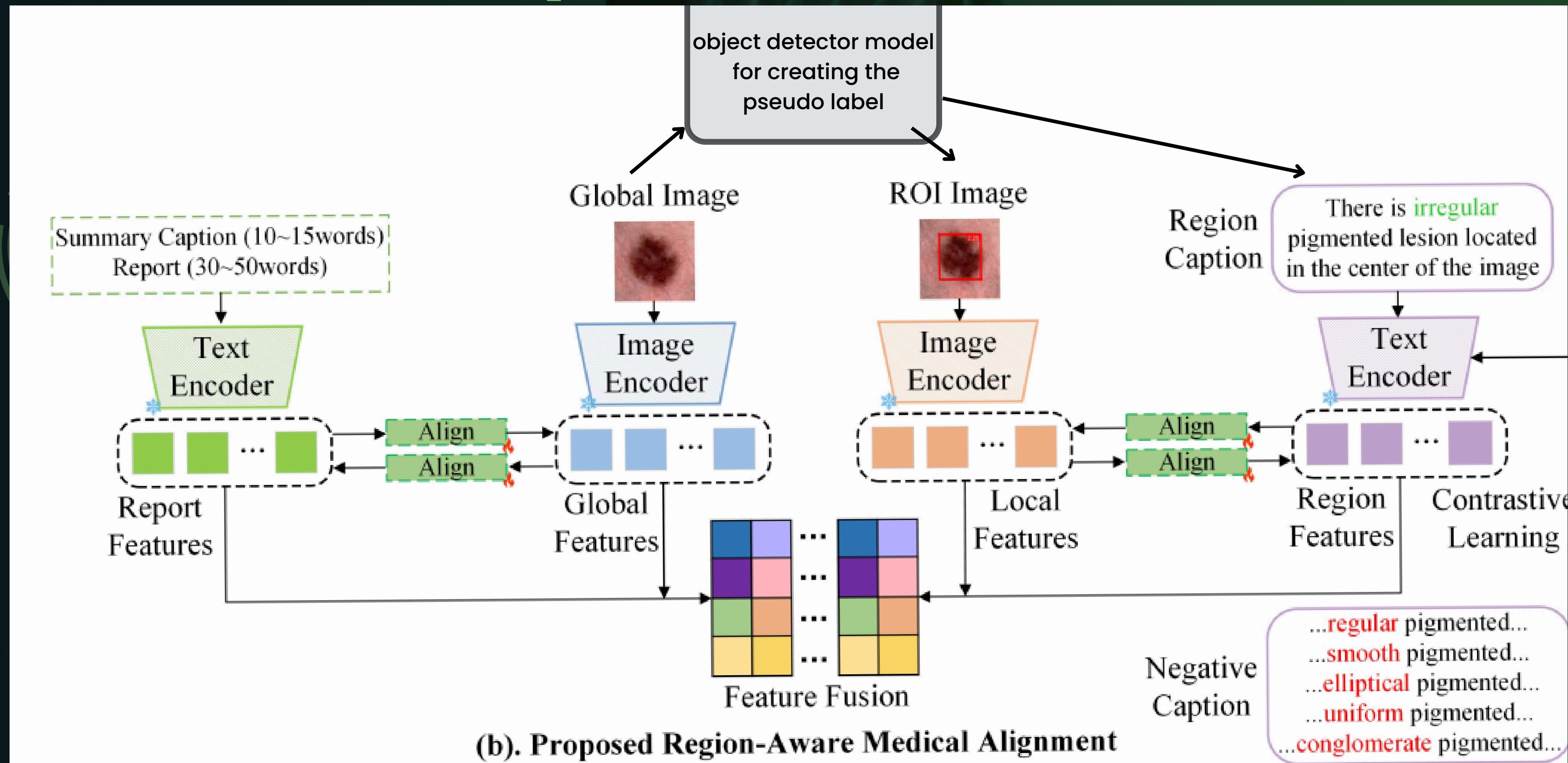
# Problem

- Why CLIP is limited in medical imaging
- Medical images require fine-grained reasoning.
- Regions matter (lesions, nodules, abnormalities).

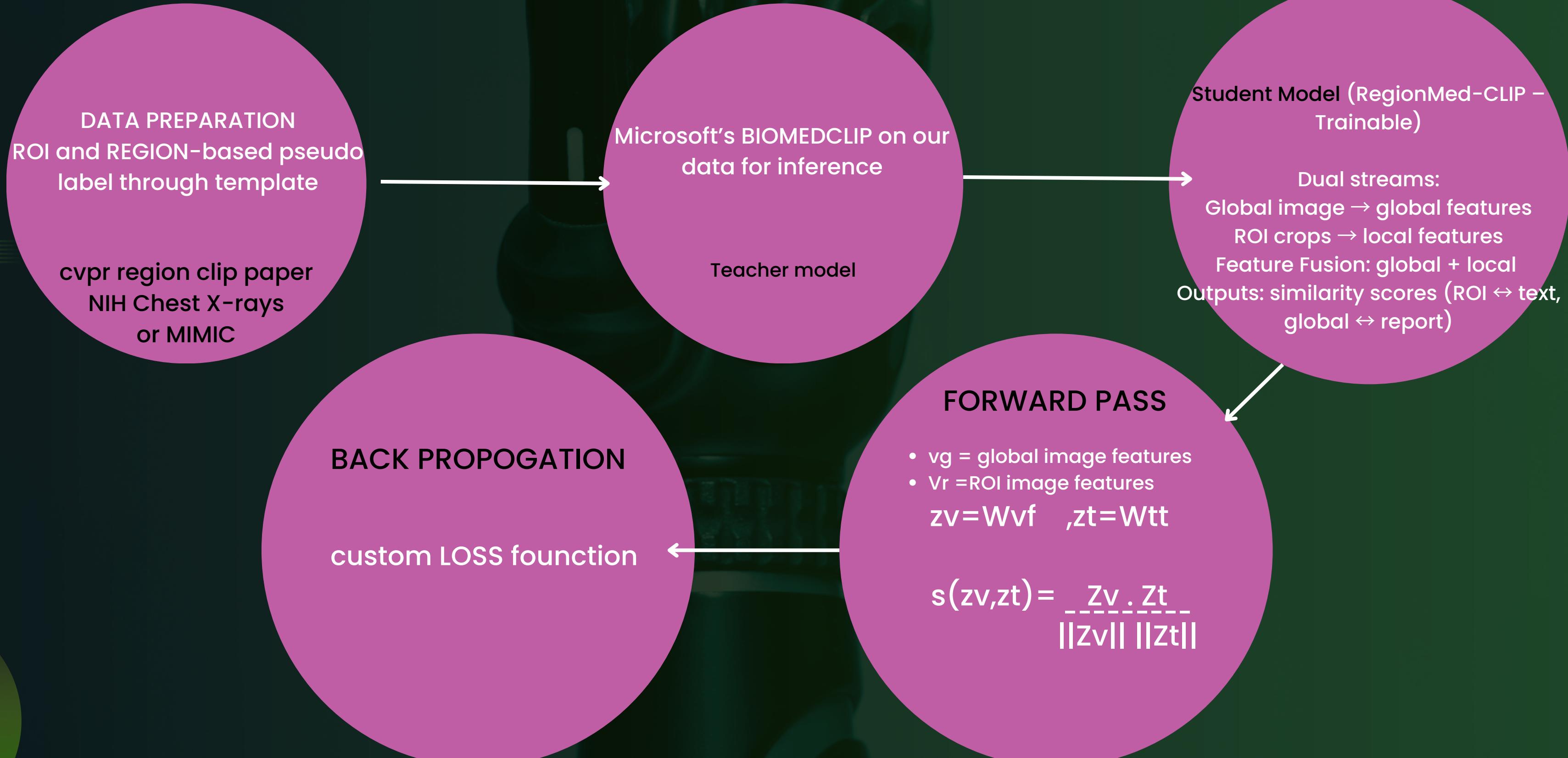
?



# Proposed workflow



# TRAINING



# LOSS Fn

**Global Contrastive Loss (InfoNCE style):**

$$L_{global} = -\log \frac{\exp(s(z_g, z_t^g)/\tau)}{\sum_j \exp(s(z_g, z_t^j)/\tau)}$$

**Distillation (KL divergence):**

$$L_{distill} = \sum_i KL(p_{teacher}(T|I_r^i) \parallel p_{student}(T|I_r^i))$$

**ROI Contrastive Loss:**

For each ROI–caption pair:

$$L_{ROI} = -\sum_i \log \frac{\exp(s(z_r^i, z_t^i)/\tau)}{\sum_j \exp(s(z_r^i, z_t^j)/\tau)}$$

**Negative Caption Loss:**

Push ROI embeddings away from mismatched captions:

$$L_{neg} = -\sum_i \log \frac{\exp(s(z_r^i, z_t^i)/\tau)}{\exp(s(z_r^i, z_t^i)/\tau) + \sum_{t_{neg}} \exp(s(z_r^i, z_{t_{neg}})/\tau)}$$

$$L_{total} = \alpha L_{global} + \beta L_{ROI} + \gamma L_{neg} + \lambda L_{distill}$$

where  $\alpha, \beta, \gamma, \lambda$  control contributions (often set = 1 except distill ~0.3–0.5).

# Expected outcome

Zero-shot disease classification

- Higher AUROC & accuracy vs. CLIP/MedCLIP.

Image  $\leftrightarrow$  Report retrieval

Region  $\leftrightarrow$  Caption matching

Fine-grained lesion description retrieval.

# Thank You!