

# Boosting Medical Visual Understanding from Multi-Granular Language Learning (ICLR 2026)

Zihan Li\*<sup>1</sup> Yiqing Wang\*<sup>2</sup> Sina Farsiu<sup>†2</sup> Paul Kinahan<sup>†1</sup>  
<sup>1</sup>University of Washington <sup>2</sup>Duke University

\*Equal Contribution <sup>†</sup>Co-corresponding Author

# The Problem: CLIP's Single-Granularity Limitation

## Standard CLIP

### Single-label, single-granularity

Image → "Chest X-ray"

#### Problems:

- Ignores multi-label nature
- Loses hierarchical structure
- Poor feature transferability
- Clinically meaningless attention



## MGLL (Ours)

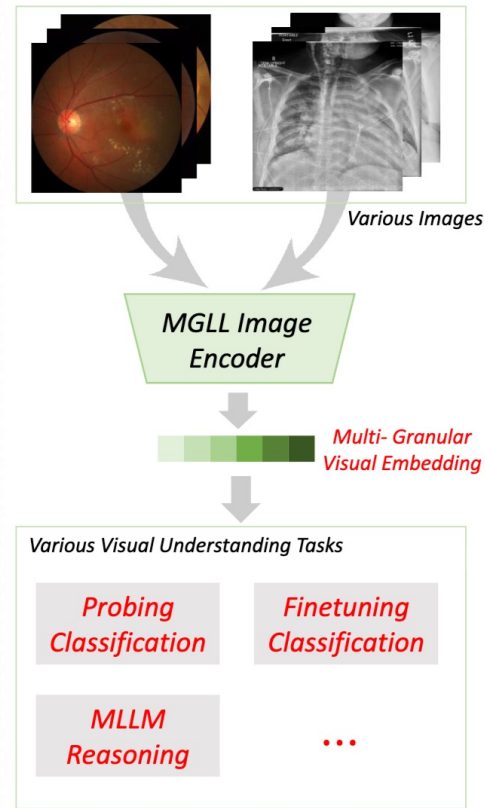
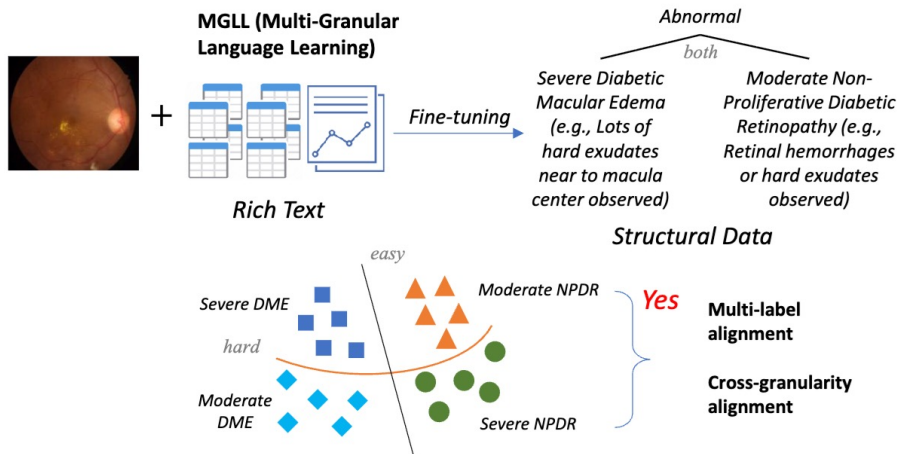
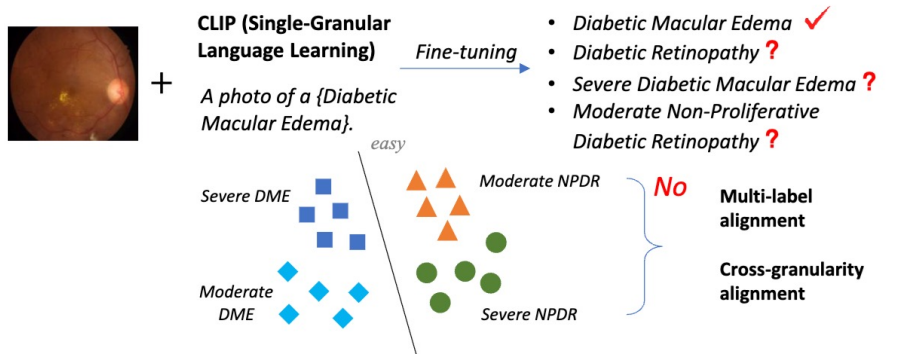
### Multi-label, multi-granularity

Granularity 1: "DX" (modality)  
Granularity 2: "Chest X-ray" (study)  
Granularity 3: "PA View" (series)

#### Advantages:

- Hierarchical alignment
- Plug-and-play, ~ 3% overhead

# The Problem: CLIP's Single-Granularity Limitation



# Multi-Granularity in Medical Imaging

## Retinal Fundus (3 Granularities)

### G1: Binary Status

*"Abnormal"*

### G2: Disease Category

*" Diabetic Retinopathy "*

### G3: Clinical Explanation

*"Retinal hemorrhages and hard exudates observed in macula"*

## Chest X-ray (3 Granularities)

### G1: Modality

*"DX" (Digital Radiography)*

### G2: Study Description

*"Chest X-ray"*

### G3: Series Description

*"PA View" (posteroanterior)*

*All textual metadata extracted from DICOM headers — no manual annotation required*

# MGLL Framework: Three Key Components

## Soft CLIP Loss

Extends contrastive loss with soft labels from label co-occurrence. At optimality, image feature  $\rightarrow$  weighted centroid of text features.



## Point-wise Constraint

Binary CE ensuring correct pairwise ranking between positive and negative image-text pairs at each granularity.



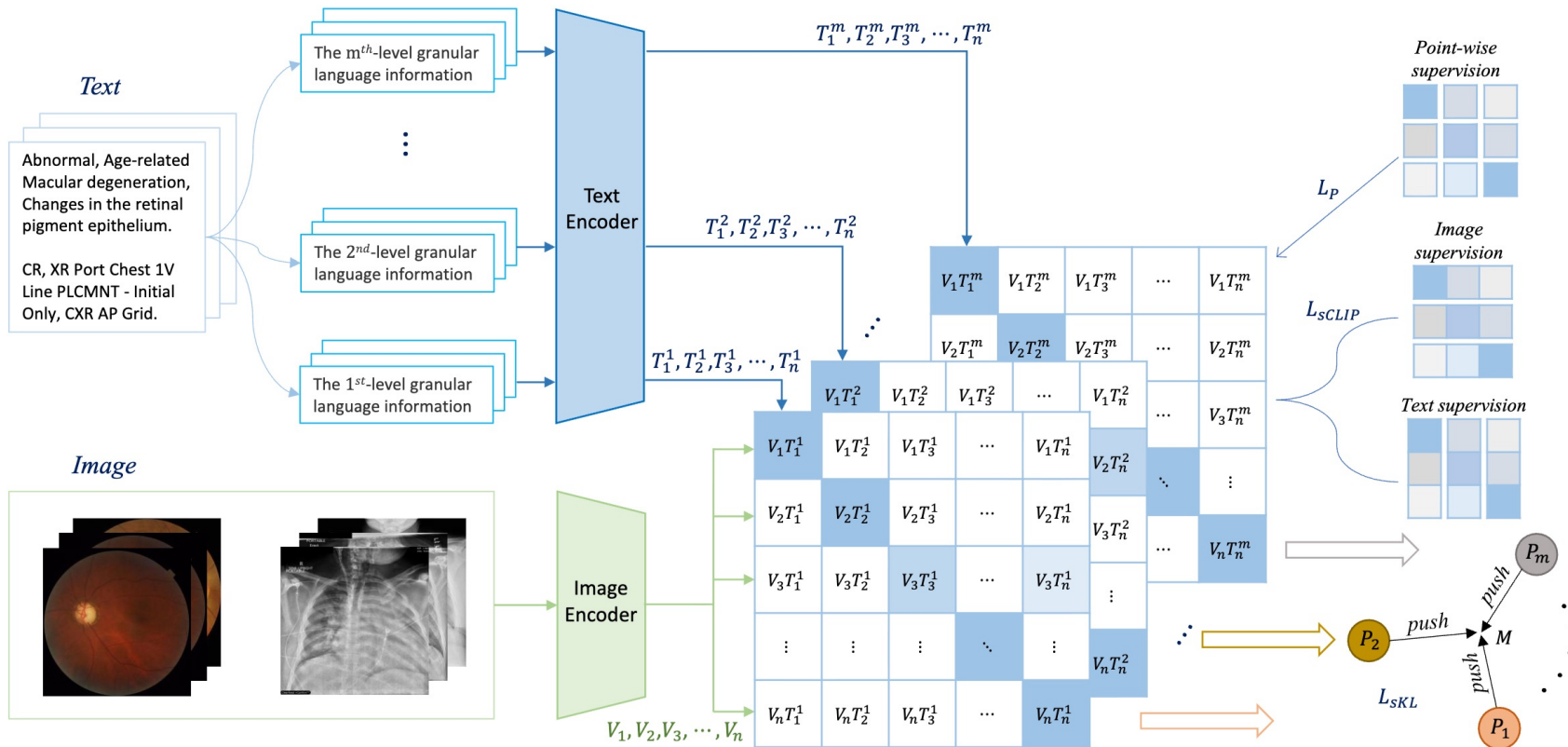
## Smooth KL Divergence

Cross-granularity consistency: fine predictions must be compatible with coarse ones. Smooth KL avoids numerical instability.

Paper: <https://arxiv.org/pdf/2511.15943>

Code: <https://github.com/HUANGLIZI/MGLL>

# MGLL Framework: Three Key Components



# Key Theoretical Insight

At optimality under Soft CLIP Loss, the image feature converges to the weighted centroid of its associated text features:

$$V_i = \frac{\sum_k w_{ik} \cdot T_{ik}}{|\sum_k w_{ik} \cdot T_{ik}|}$$

## Why this matters:

- CLIP forces Visual representation to match a single text T — biased representation
- MGLL allows Visual embedding to sit at the semantic center of ALL its labels — balanced representation
- Weights  $w_{ik}$  ensure each label contributes proportionally to its clinical relevance
- This principled formulation explains why soft alignment is fundamentally more appropriate for multi-label medical domains

# Large-Scale Multi-Granular Datasets

## MGLL-Fundus

**246,389**

image-text pairs

- 49 public datasets aggregated
- 50+ disease categories
- 2 granularity levels
- Clinical explanations from EyeWiki

## MGLL-Xray

**190,882**

X-ray images from MIDRC

- DICOM → PNG with metadata
- 3 granularity levels from headers
- CR + DX modality types
- Zero manual annotation needed

# Results: State-of-the-Art Across 10+ Benchmarks

**+11%**

AUC on ADAM  
(LP: 90% vs 79%)

**+16%**

AUC on RFMiD  
(LP: 80% vs 63%)

**99.1%**

AUC on MIDRC-XR  
(FT, vs 94% UniMed)

**+34%**

LLaVA-Med gain  
with MGLL encoder

Dataset	CLIP	Best Baseline	MGLL (Ours)	Improvement
ADAM (LP AUC)	61.43%	79.33%	<b>90.02%</b>	<b>+10.69%</b>
REFUGE (LP AUC)	56.82%	84.59%	<b>92.42%</b>	<b>+7.83%</b>
MIDRC-XR (FT AUC)	88.52%	94.15%	<b>99.08%</b>	<b>+4.93%</b>
MIDRC-Portable (FT AUC)	91.83%	97.26%	<b>99.75%</b>	<b>+2.49%</b>
RFMiD Multi-label (LP AUC)	44.66%	63.38%	<b>79.62%</b>	<b>+16.24%</b>

# Results: State-of-the-Art Across 10+ Benchmarks

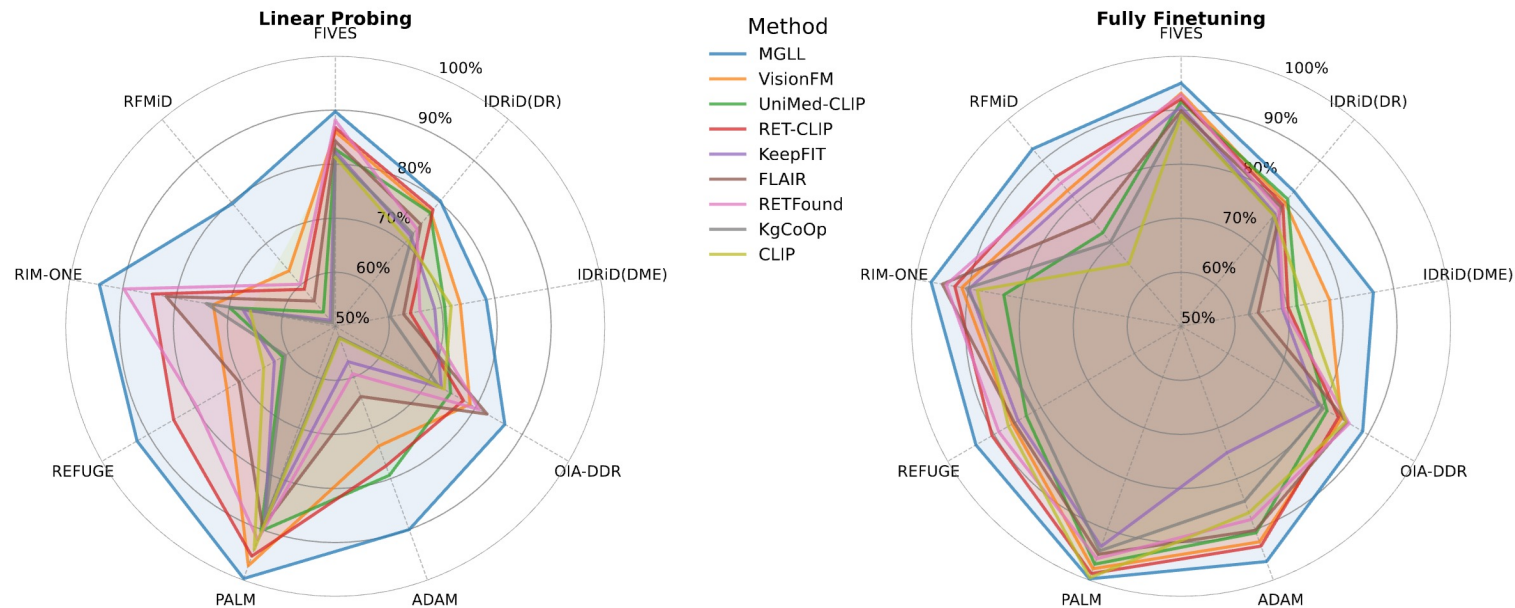


Figure 3: The quantitative comparison (AUC) between baseline methods and proposed MGLL on nine fundus downstream datasets.

# Integration with Multimodal LLMs

Replacing CLIP encoders with MGLL in 7 different MLLMs consistently improved performance — even for medical-specialized models.

Up to **34.1% improvement** for LLaVA-Med with MGLL encoder swap

## Key takeaway:

- MGLL's plug-and-play design means better visual features transfer directly to downstream models
- Works across architectures: LLaVA, Med-Flamingo, LLaVA-Med, and 4 others
- Even medical-specialized LLMs (trained on medical data) benefit from MGLL's multi-granular features
- No architectural changes needed — just swap the vision encoder

# Integration with Multimodal LLMs

Table 2: Comparison of multiple-choice accuracy with MGLL in multimodal large language models on selected ten representative diseases.

Method	AMD	Cataract	CSR	DR	Glaucoma	Media Haze	Myopia	Retinitis	DME	Tessellation	Average $\uparrow$
<a href="#">InstructBLIP Dai et al. (2023)</a>	80.17%	80.00%	0.00%	76.51%	59.30%	16.13%	44.25%	11.11%	63.79%	41.67%	47.29%
+ MGLL	<b>83.63%</b>	<b>85.00%</b>	<b>28.57%</b>	<b>82.55%</b>	<b>65.43%</b>	<b>45.16%</b>	<b>52.65%</b>	<b>44.44%</b>	<b>74.14%</b>	<b>58.33%</b>	<b>61.99%</b> (14.7% $\uparrow$ )
<a href="#">Mini-Gemini Li et al. (2024)</a>	76.61%	85.00%	14.29%	79.87%	67.90%	38.71%	58.41%	33.33%	60.34%	33.33%	54.78%
+ MGLL	<b>82.46%</b>	<b>85.00%</b>	<b>42.86%</b>	<b>84.56%</b>	<b>72.22%</b>	<b>58.06%</b>	<b>64.16%</b>	<b>55.56%</b>	<b>65.52%</b>	<b>41.67%</b>	<b>65.21%</b> (10.4% $\uparrow$ )
<a href="#">Qwen-VL Bai et al. (2023)</a>	81.87%	75.00%	28.57%	80.54%	78.40%	54.84%	76.55%	22.22%	84.48%	25.00%	60.75%
+ MGLL	<b>85.96%</b>	<b>80.00%</b>	<b>42.86%</b>	<b>89.93%</b>	<b>87.04%</b>	<b>70.97%</b>	<b>80.97%</b>	<b>33.33%</b>	<b>89.66%</b>	<b>41.67%</b>	<b>70.24%</b> (9.5% $\uparrow$ )
<a href="#">InternVL Chen et al. (2024a)</a>	81.29%	85.00%	71.43%	94.63%	89.51%	64.52%	88.05%	44.44%	87.93%	66.67%	77.35%
+ MGLL	<b>86.55%</b>	<b>90.00%</b>	<b>71.43%</b>	<b>96.64%</b>	<b>90.74%</b>	<b>67.74%</b>	<b>91.15%</b>	<b>55.56%</b>	<b>94.83%</b>	<b>75.00%</b>	<b>81.96%</b> (4.6% $\uparrow$ )
<a href="#">LLaVA Liu et al. (2023)</a>	83.04%	90.00%	42.86%	87.25%	91.36%	48.39%	88.50%	44.44%	93.10%	58.33%	72.73%
+ MGLL	<b>84.80%</b>	<b>90.00%</b>	<b>57.14%</b>	<b>93.96%</b>	<b>91.98%</b>	<b>61.29%</b>	<b>90.71%</b>	<b>66.67%</b>	<b>96.55%</b>	<b>66.67%</b>	<b>79.98%</b> (7.3% $\uparrow$ )
<a href="#">LLaVA-Med Li et al. (2023)</a>	16.37%	15.00%	42.86%	26.85%	25.31%	25.81%	23.89%	33.33%	16.67%	16.67%	24.28%
+ MGLL	<b>58.48%</b>	<b>65.00%</b>	<b>57.14%</b>	<b>77.18%</b>	<b>59.26%</b>	<b>51.61%</b>	<b>57.08%</b>	<b>44.44%</b>	<b>55.17%</b>	<b>58.33%</b>	<b>58.37%</b> (34.1% $\uparrow$ )
<a href="#">Med-Flamingo Moor et al. (2023)</a>	25.73%	30.00%	57.14%	36.91%	24.07%	22.58%	18.58%	22.22%	24.14%	8.33%	26.97%
+ MGLL	<b>69.01%</b>	<b>75.00%</b>	<b>71.43%</b>	<b>80.54%</b>	<b>61.11%</b>	<b>54.84%</b>	<b>45.58%</b>	<b>44.44%</b>	<b>51.72%</b>	<b>33.33%</b>	<b>58.70%</b> (31.7% $\uparrow$ )
<a href="#">Janus-Pro Chen et al. (2025)</a>	88.30%	75.00%	42.86%	93.29%	90.74%	58.06%	87.17%	33.33%	62.07%	58.33%	68.92%
+ MGLL	<b>90.64%</b>	<b>85.00%</b>	<b>71.43%</b>	<b>96.64%</b>	<b>95.06%</b>	<b>67.74%</b>	<b>90.27%</b>	<b>55.56%</b>	<b>70.69%</b>	<b>75.00%</b>	<b>79.80%</b> (10.88% $\uparrow$ )

# Connecting MGLL to Clinical Research

## Clinical Foundation Models

Multi-granular alignment for clinical notes

*focus: Foundation models for clinical text using private unique data*

## Multimodal Disease Prediction

MGLL handles multi-label conditions at different stages of progression naturally

*focus: Predicting disease onset from multimodal data with uncertainty*

## Data Harmonization at Scale

DICOM/LOINC harmonization as multi-granular alignment — directly transferable to EHR unification

*focus: Unifying data across millions of patients and providers*

**MGLL is plug-and-play, scalable, and designed for real clinical workflows**

# Summary

**Problem:**

CLIP's single-granularity alignment fails for complex medical images

**Solution:**

MGLL — multi-label + cross-granularity alignment via 3 complementary losses

**Impact:**

SOTA on 10+ benchmarks, 34% MLLM improvement, clinically grounded features

**Key Properties:**

Architecture-agnostic, plug-and-play, minimal computational overhead