# L³Seg: Lean Linear Layers for Language-Guided Vision Transformer in Medical Image Segmentation

Rahul Bhardwaj[1], Utkarsh Yashwant Tambe[2], Debanga Raj Neog[1]

[1]Mehta Family School of Data Science & Artificial Intelligence, Indian Institute of Technology Guwahati, India

[2]Department of Data Science & Business Systems, SRM Institute of Science & Technology, Kattankulathur, India

## Introduction

**Problem Statement:**
- Vision-Language models are heavy (params, FLOPs).
- Need cross-modality generalization.

**Motivation:**
- Compute is concentrated in dense linear projections.
- Fine-tuning/PEFT change few weights while base matrix multiplications still run, keeping compute high and adaptation limited.

**Key Contributions:**
- Replace all dense projections with L³
- Trainables: $O(d_{in} d_{out}) \rightarrow O(r(d_{in} + d_{out}))$

---

**Algorithm 1** Lean Linear Layer (L³)

```
1:  Input: x ∈ ℝ^{B×N×d_in}
2:  Frozen base: W₀ ∈ ℝ^{d_out×d_in}, b₀ ∈ ℝ^{d_out}
3:  B: batch size,  N: number of tokens,  d: feature dimension
4:  class LeanLinearLayer(Module):
5:    def __init__(self, d_in, d_out, r, W₀, b₀):
6:      super().__init__()
7:      self.W₀ ← W₀                          [d_out, d_in], frozen
8:      self.b₀ ← b₀                          [d_out], frozen
9:      # trainable low-rank factors
10:     Notation: 𝒩(0, 10⁻³)_{m×n} = randn(m, n) × 10⁻³
11:     Notation: 0_{m×n} = zeros(m, n)
12:     self.A_g ← 𝒩(0, 10⁻³)^{d_in×r}
13:     self.B_g ← 0_{r×d_out}
14:     self.A_b ← 𝒩(0, 10⁻³)^{d_in×r}
15:     self.B_b ← 0_{r×d_out}
16:   def forward(self, x):
17:     # 1. frozen baseline
18:     y₀ ← xW₀ᵀ + b₀                        [B, N, d_out]
19:     # 2. low-rank scale and shift
20:     γ ← (xA_g)B_g                         [B, N, d_out]
21:     β ← (xA_b)B_b                         [B, N, d_out]
22:     # 3. scaled-offset fusion
23:     return (1 + γ) ⊙ y₀ + β
```

## Method


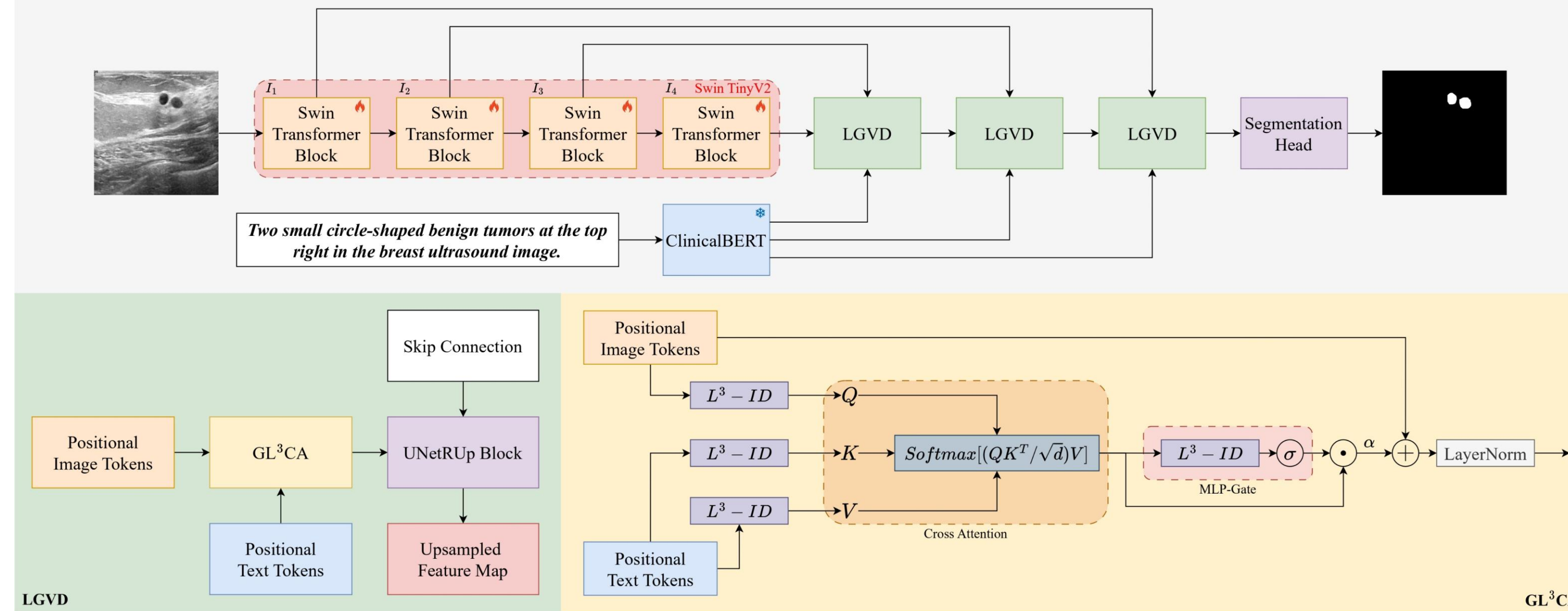
Fig. 1. Overview of L³Seg: Language-Guided Vision Decoder fuses image and text using Gated L³ Cross-Attention

## Comparative Analysis

**Table 1. Quantitative Comparison on QaTa-COV19 (X-ray), Kvasir-SEG (endoscopy) and BUSI (ultrasound) dataset. CNN-based (◇), SAM-based (¶), and hybrid CNN-Transformer (†).**

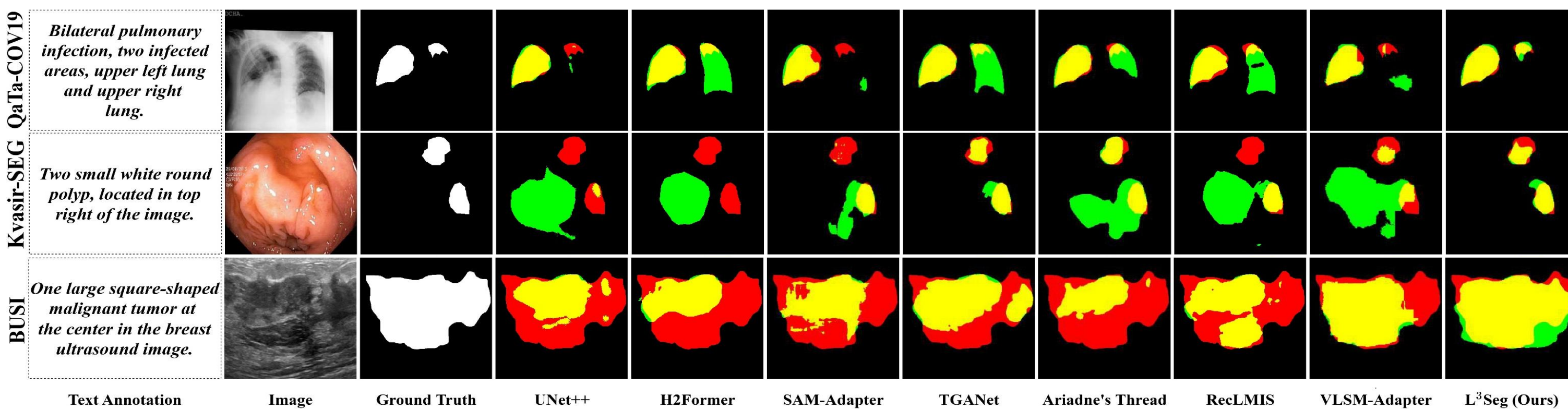| Method | Venue | Text | Params (M) ↓ | FLOPs (G) ↓ | QaTa-COV19 (XRay) Dice(%) ↑ | QaTa-COV19 (XRay) mIoU(%) ↑ | Kvasir-SEG (Endoscopy) Dice(%) ↑ | Kvasir-SEG (Endoscopy) mIoU(%) ↑ | BUSI (Ultrasound) Dice(%) ↑ | BUSI (Ultrasound) mIoU(%) ↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| UNet◇ [21] | MICCAI'15 | ✗ | 14.8 | 50.3 | 79.02 | 69.46 | 81.83 | 74.60 | 57.28 | 49.19 |
| UNet++◇ [27] | IEEE TMI'19 | ✗ | 74.5 | 94.6 | 79.62 | 70.25 | 82.10 | 74.43 | 63.46 | 56.59 |
| Swin-UNet† [3] | ECCV'22 | ✗ | 82.3 | 67.3 | 78.07 | 68.34 | 85.90 | 77.56 | 63.67 | 55.54 |
| H2Former† [9] | IEEE TMI'23 | ✗ | 33.7 | 24.6 | 77.86 | 68.35 | 80.03 | 72.23 | 63.72 | 56.71 |
| SAM¶ [13] | ICCV'23 | ✗ | 93.6 | 50.9 | 71.85 | 56.06 | 77.83 | 70.72 | 49.93 | 33.27 |
| SAM-Adapter¶ [4] | ICCV'23 | ✗ | 104.3 | 55.2 | 84.76 | 73.55 | 83.42 | 71.55 | 77.47 | 63.22 |
| CLIPSeg† [17] | CVPR'22 | ✓ | 150.0 | 23.0 | 78.92 | 71.55 | 83.71 | 76.02 | 62.06 | 57.91 |
| TGANet◇ [22] | MICCAI'22 | ✓ | 19.8 | 41.9 | 79.87 | 70.75 | 89.51 | 82.49 | 69.33 | 62.32 |
| Ariadne's Thread† [26] | MICCAI'23 | ✓ | 44.0 | 22.4 | 89.78 | 81.45 | 87.61 | 77.95 | 79.36 | 65.78 |
| LViT† [14] | IEEE TMI'23 | ✓ | 29.7 | 54.1 | 83.66 | 75.11 | 88.62 | 81.90 | 65.51 | 58.73 |
| RecLMIS† [11] | IEEE TMI'24 | ✓ | 23.7 | 24.1 | 85.22 | 77.00 | 85.78 | 78.76 | 63.66 | 55.96 |
| SGSeg◇ [24] | MICCAI'24 | ✓ | 76.9 | 19.3 | 87.41 | 77.85 | 86.99 | 77.27 | 68.39 | 63.68 |
| VLSM-Adapter† [7] | MICCAI'24 | ✓ | 136.9 | 38.3 | 79.98 | 76.69 | 82.34 | 74.91 | 65.02 | 57.20 |
| **L³Seg (Ours)†** | ICCVW'25 | ✓ | **8.2** | **5.1** | **90.98** | **83.46** | **90.10** | **82.67** | **85.53** | **74.72** |



Fig. 2. Qualitative Comparison on QaTa-COV19, Kvasir-SEG and BUSI dataset. (TP, FN, FP)
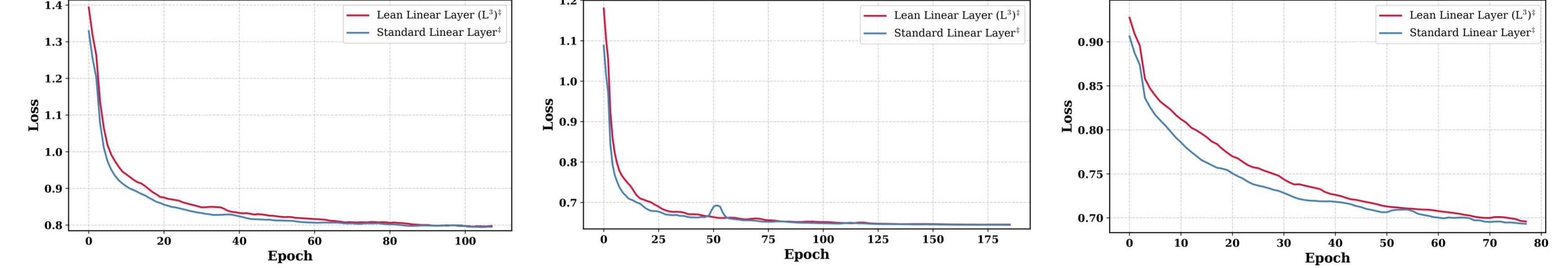
## Experimental Results



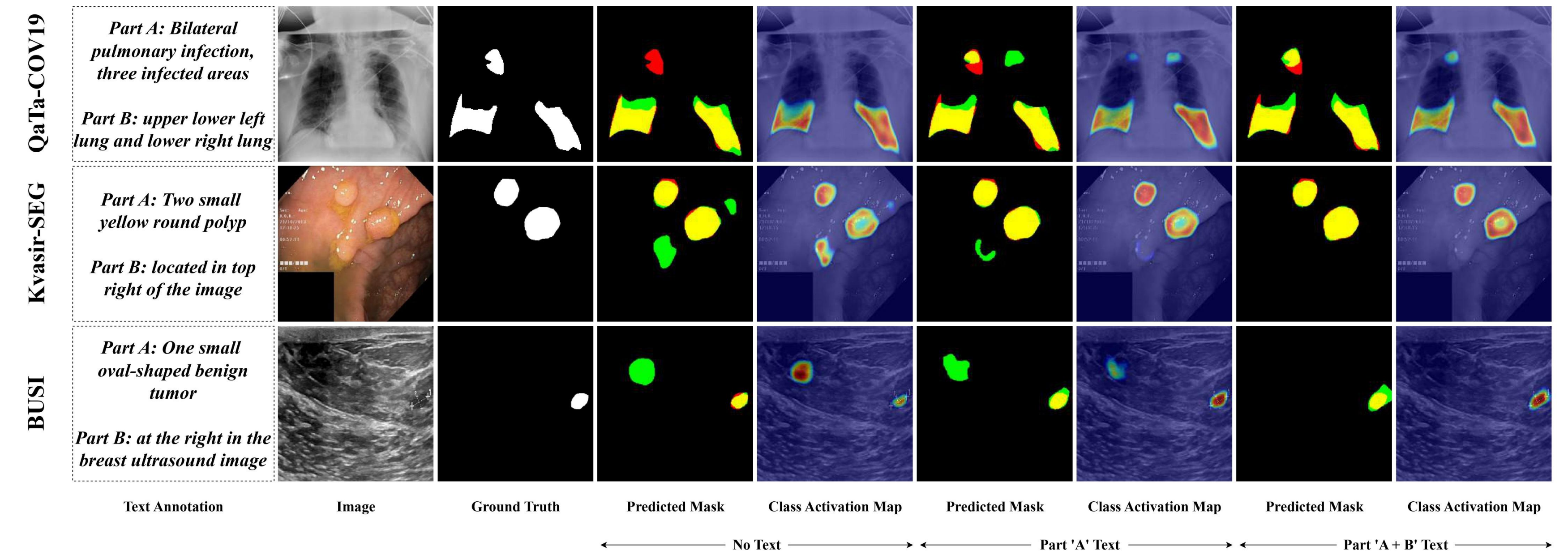Fig. 3. Training loss curves for both the Standard Linear Layer and the Lean Linear Layer (L³).



Fig. 4. Segmentation Visualizations with Varying Text Inputs. (TP, FN, FP)

**Table 2. Impact of Training Data Size.**

| Data Usage | QaTa-COV19 Dice(%) ↑ | QaTa-COV19 mIoU(%) ↑ | Kvasir-SEG Dice(%) ↑ | Kvasir-SEG mIoU(%) ↑ | BUSI Dice(%) ↑ | BUSI mIoU(%) ↑ |
|---|---|---|---|---|---|---|
| SAM-Adapter [4] (100% Training) | 84.76 | 73.55 | 83.42 | 71.55 | 77.47 | 63.22 |
| VLSM-Adapter [7] (100% Training) | 79.98 | 76.69 | 82.34 | 74.91 | 65.02 | 57.20 |
| L³Seg (25% Training) | 86.15 | 77.43 | 83.06 | 72.50 | 77.29 | 62.98 |
| L³Seg (50% Training) | 87.10 | 80.98 | 84.99 | 73.90 | 82.05 | 69.57 |
| L³Seg (75% Training) | 89.59 | 81.80 | 87.96 | 78.50 | 83.61 | 71.83 |
| **L³Seg (100% Training)** | **90.98** | **83.46** | **90.10** | **82.67** | **85.53** | **74.72** |

## At a Glance

- Also accepted at **ICCV 2025 CVAMD Workshop.**

- For more information, please visit project webpage.