

A. Appendix / Supplemental Material

A.1. Pseudocode

In this section, we provide pseudocodes to illustrate the workflow of PSA. Algorithm 1 presents the initialization process, detailing how discrete prototypes are identified from paired image–text data, while Algorithm 2 demonstrates the query–response mechanism generating approximate textual semantics purely from image embeddings.

Algorithm 1 PSA Initialization

```

1: Define:
2:   - A set of  $K$  paired samples  $\{(I_i, T_i)\}_{i=1}^K$ ;
3:   - image-only samples  $\{I_j\}$ ;
4:   - pretrained encoders  $f_{\text{enc}}^I$  and  $f_{\text{enc}}^T$ ;
5:   - cross-attention module (Language-Guided U-Net)
6:   - token selection threshold  $\tau$ ;
7:   - number of semantic clusters  $N$  (HDBSCAN);
8:   - number of sub-clusters  $M$  (K-means).
9: Return: Prototype space  $\mathcal{S} = (\mathcal{S}^Q, \mathcal{S}^R)$ 
10: Step 1: Encode Paired Samples
11: for  $i = 1$  to  $K$  do
12:    $e_i^I \leftarrow f_{\text{enc}}^I(I_i)$ 
13:    $e_i^T \leftarrow f_{\text{enc}}^T(T_i)$ 
14: end for
15: Step 2: Extract Segmentation-Relevant Tokens
16: for  $i = 1$  to  $K$  do
17:   Compute cross-attention scores  $\alpha_j$  for each token  $t_j$ 
   in  $T_i$ 
18:    $T_i^{\text{selected}} \leftarrow \{t_j \mid \alpha_j > \tau\}$ 
19:    $e_i^{\text{sem}} \leftarrow f_{\text{enc}}^T(T_i^{\text{selected}})$ 
20: end for
21: Step 3: Cluster Textual Semantics (HDBSCAN)
22: Perform HDBSCAN on  $\{e_i^{\text{sem}}\}$  to form  $N$  clusters
 $\{\mathcal{C}_1, \dots, \mathcal{C}_N\}$ 
23: Step 4: Form Image Sub-Clusters (K-means)
24:  $\mathcal{S}^Q \leftarrow \emptyset, \mathcal{S}^R \leftarrow \emptyset$ 
25: for  $i = 1$  to  $N$  do
26:   Extract embeddings  $\{(e_j^I, e_j^T) \mid j \in \mathcal{C}_i\}$ 
27:   Run K-means with  $M$  sub-clusters:  $\mathcal{C}_{i1}, \dots, \mathcal{C}_{iM}$ 
28:   for  $j = 1$  to  $M$  do
29:     Identify representative  $c_{ij} = (e_k^I, e_k^T)$  closest to
     sub-cluster centroid
30:      $q_{ij} \leftarrow e_k^I$  (query prototype)
31:      $r_{ij} \leftarrow e_k^T$  (response prototype)
32:      $\mathcal{S}^Q \leftarrow \mathcal{S}^Q \cup \{q_{ij}\}, \mathcal{S}^R \leftarrow \mathcal{S}^R \cup \{r_{ij}\}$ 
33:   end for
34: end for
35: Step 5: Output Prototype Space
36:  $\mathcal{S} \leftarrow (\mathcal{S}^Q, \mathcal{S}^R)$ 
37: return  $\mathcal{S}$ 

```

Algorithm 2 PSA Query and Response

Require: Prototype space $\mathcal{S} = (\mathcal{S}^Q, \mathcal{S}^R)$; pretrained image encoder f_{enc}^I ; Language-Guided U-Net f_{seg} ; query image I^* ; top- k integer k .

Ensure: Approximated textual feature r^* for guiding segmentation

```

1: Step 1: Encode the Query Image
2:  $q^* \leftarrow f_{\text{enc}}^I(I^*)$ 
3: Step 2: Compute Similarity Scores
4: for all  $q_{ij}$  in  $\mathcal{S}^Q$  do
5:    $s_{ij} \leftarrow \text{cosine\_similarity}(q^*, q_{ij})$ 
6: end for
7: Step 3: Select Top- $k$  Queries
8:  $Q^* \leftarrow \arg \text{top}_k(\{s_{ij}\})$ 
9: Step 4: Retrieve Corresponding Responses
10:  $R^* \leftarrow \{r_{ij} \mid q_{ij} \in Q^*\}$ 
11: Step 5: Aggregate Responses (Weighted Sum)
12:  $r^* \leftarrow \sum_{(q_{ij}, r_{ij}) \in Q^* \times R^*} w_{ij} r_{ij}$ 
13: where  $w_{ij} = \frac{\exp(s_{ij})}{\sum_{q_{i'j'} \in Q^*} \exp(s_{i'j'})}$ 
14: return  $r^*$ 

```

A.2. Implementation Details

Following the previous design of language-guided segmentation networks [? ?], we adopt a U-Net backbone with feature fusion at the decoder stage. The image is resized into 224×224 , and textual reports are tokenized, truncated, and padded to a fixed length of 256 tokens. To construct the prototype space we set the number of surrogate labels to 6, with each label containing 64 prototypes. During inference, the PSA module retrieves the top 10 prototype candidates per query for semantic approximation. We use the AdamW optimizer with an initial learning rate of 10^{-4} , which is scheduled to decay using cosine annealing.

A.3. Limitations and Future Works

In this work, we focused on demonstrating the core idea of ProLearn in single-label 2D segmentation. Future directions involve exploring multi-label and volumetric data, broader imaging modalities, and extending to more language-guided vision tasks.

A.4. Visualization

To further demonstrate the effectiveness of our proposed ProLearn, we provide additional visual comparisons of segmentation results in the next page. Specifically, we show the performance of LViT, GuideSeg, SGSeg, and our ProLearn on the QaTa-COV19 and MosMedData+ dataset under different text availability (1%, 5%, 10%, 25%, and 50%).

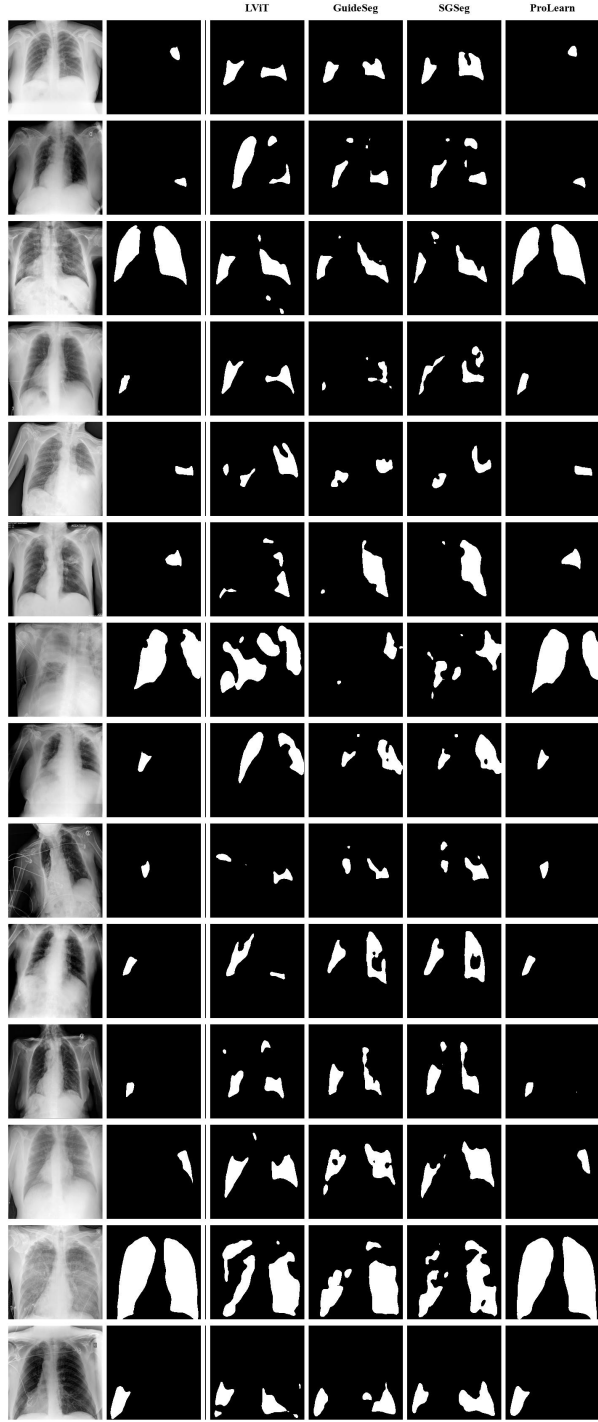


Figure A1. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on QaTa-COV19 under 1% text availability.

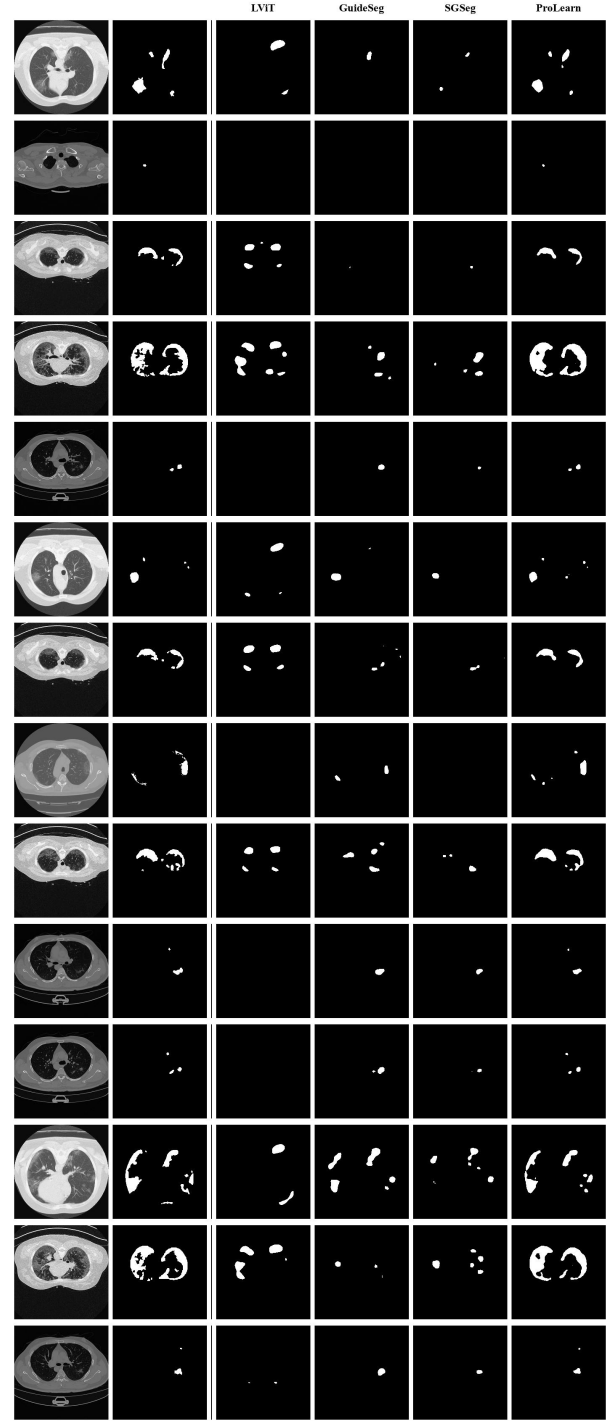


Figure A2. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on MosMedData+ dataset under 1% text availability.

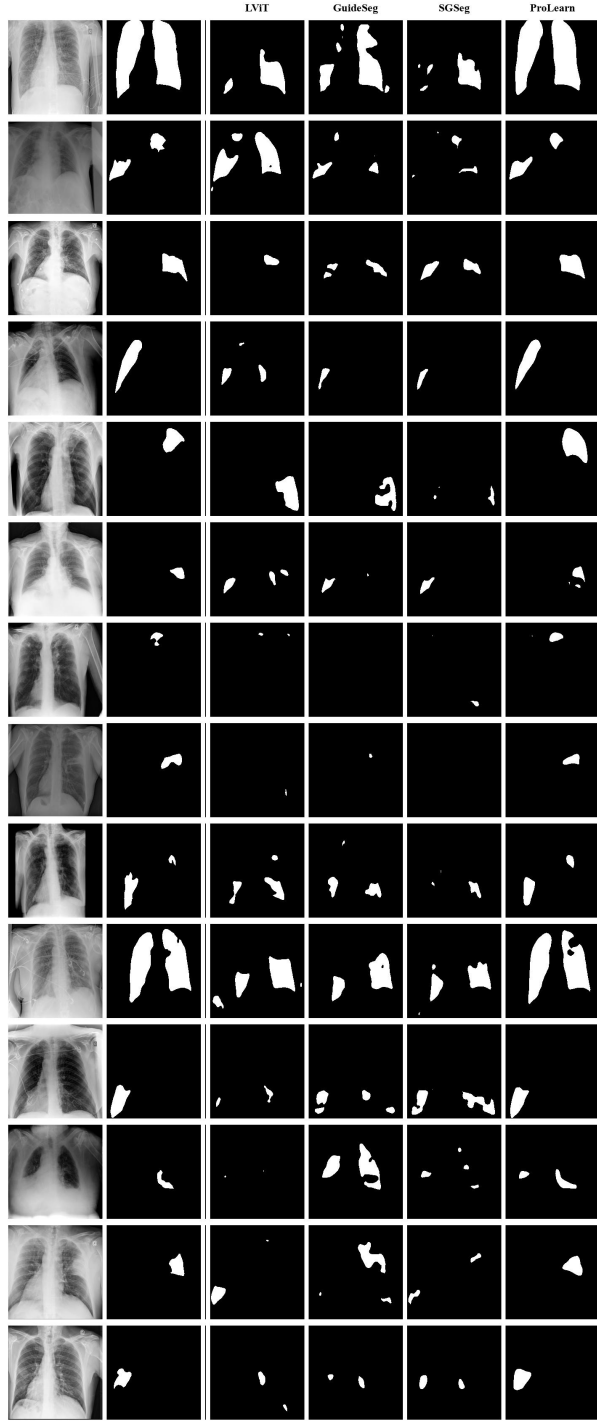


Figure A3. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on QaTa-COV19 dataset under 5% text availability.

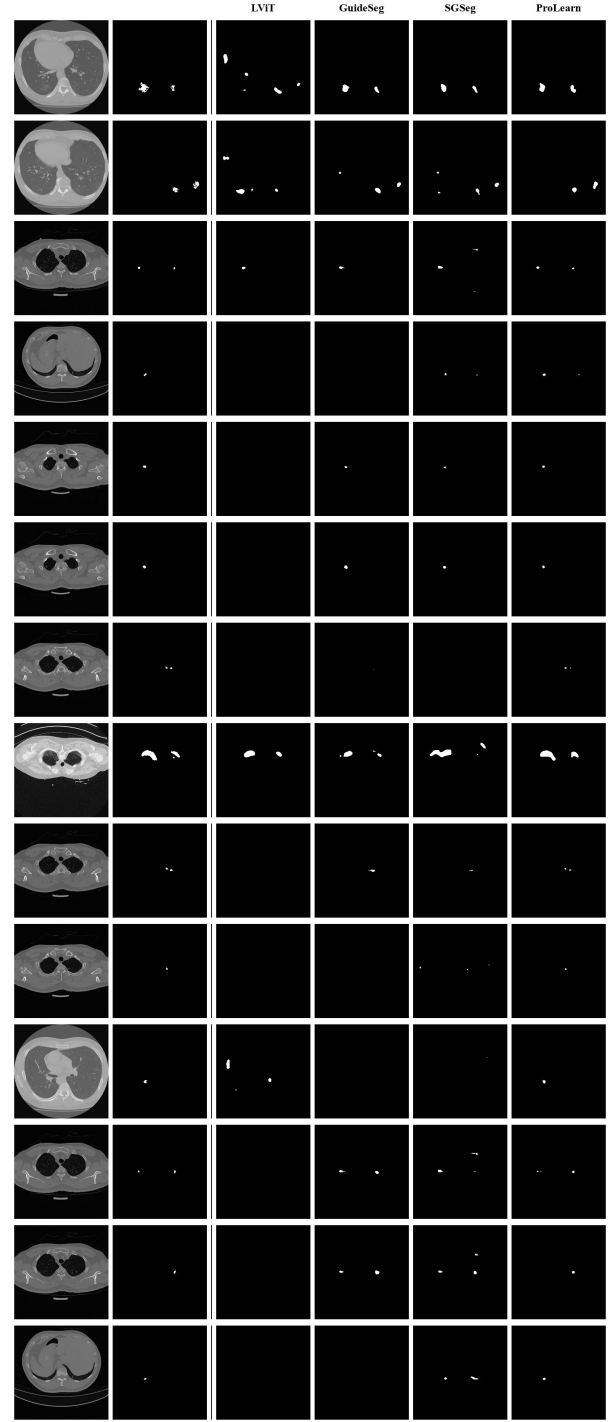


Figure A4. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on MosMedData+ dataset under 5% text availability.

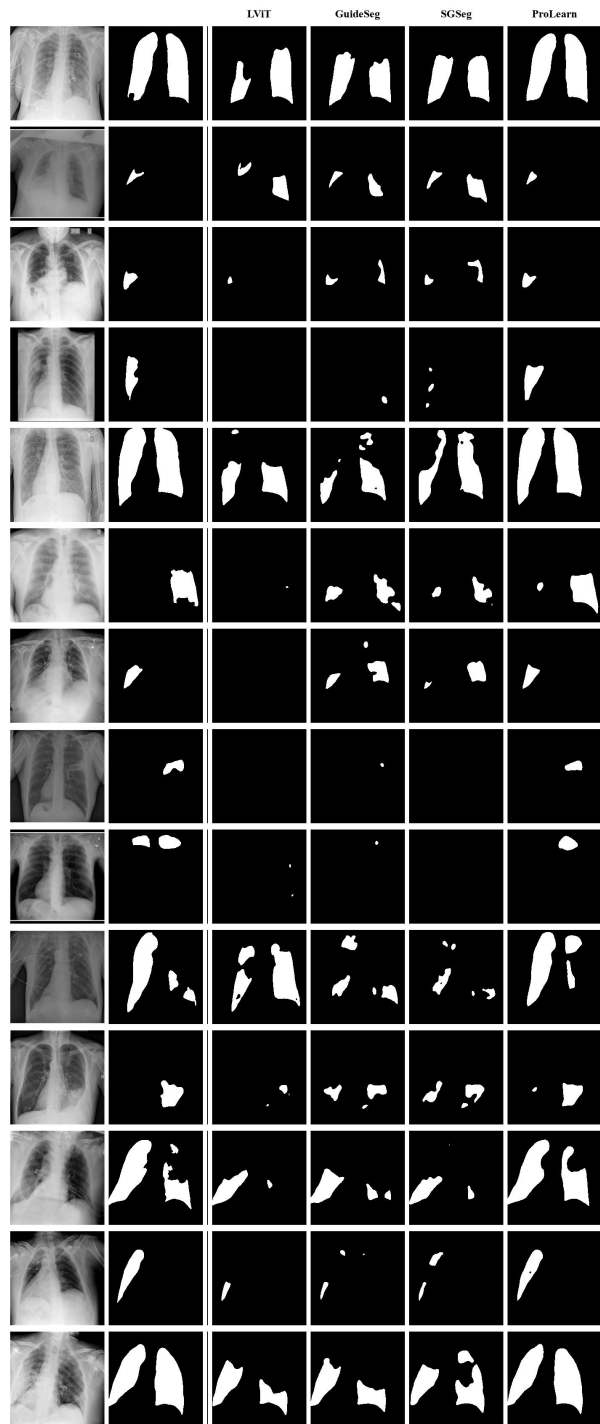


Figure A5. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on QaTa-COV19 dataset under 10% text availability.

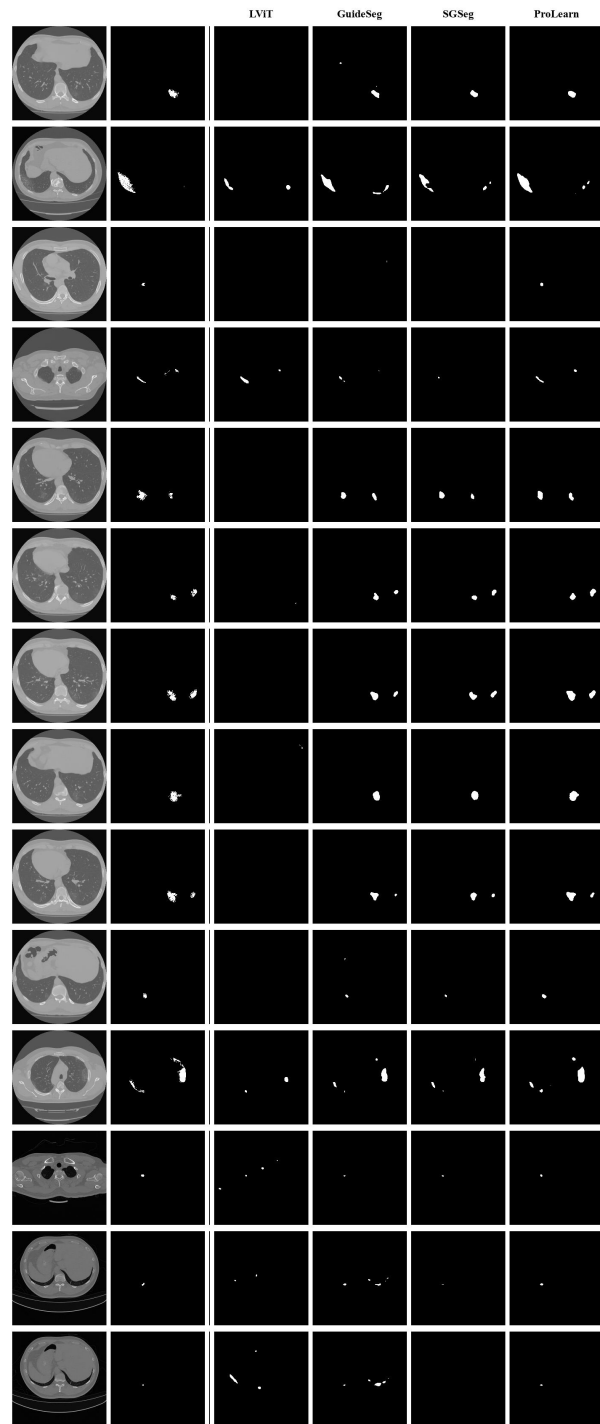


Figure A6. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on MosMedData+ dataset under 10% text availability.

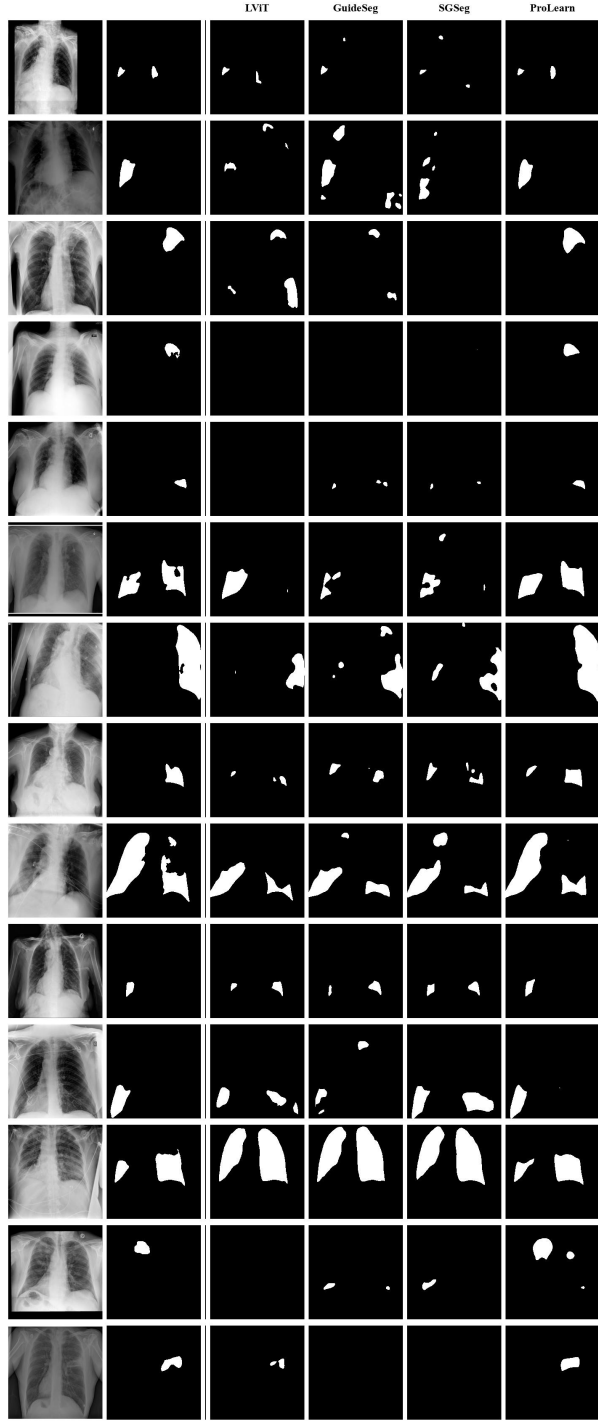


Figure A7. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on QaTa-COV19 dataset under 25% text availability.

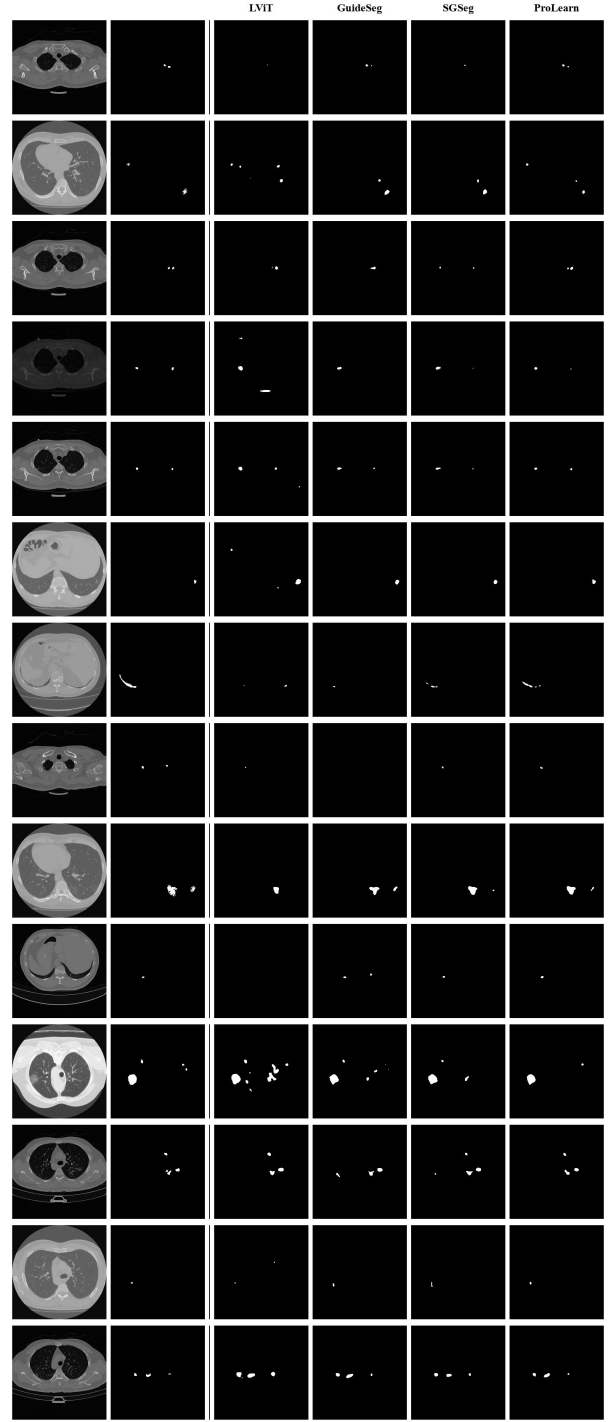


Figure A8. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on MosMedData+ dataset under 25% text availability.

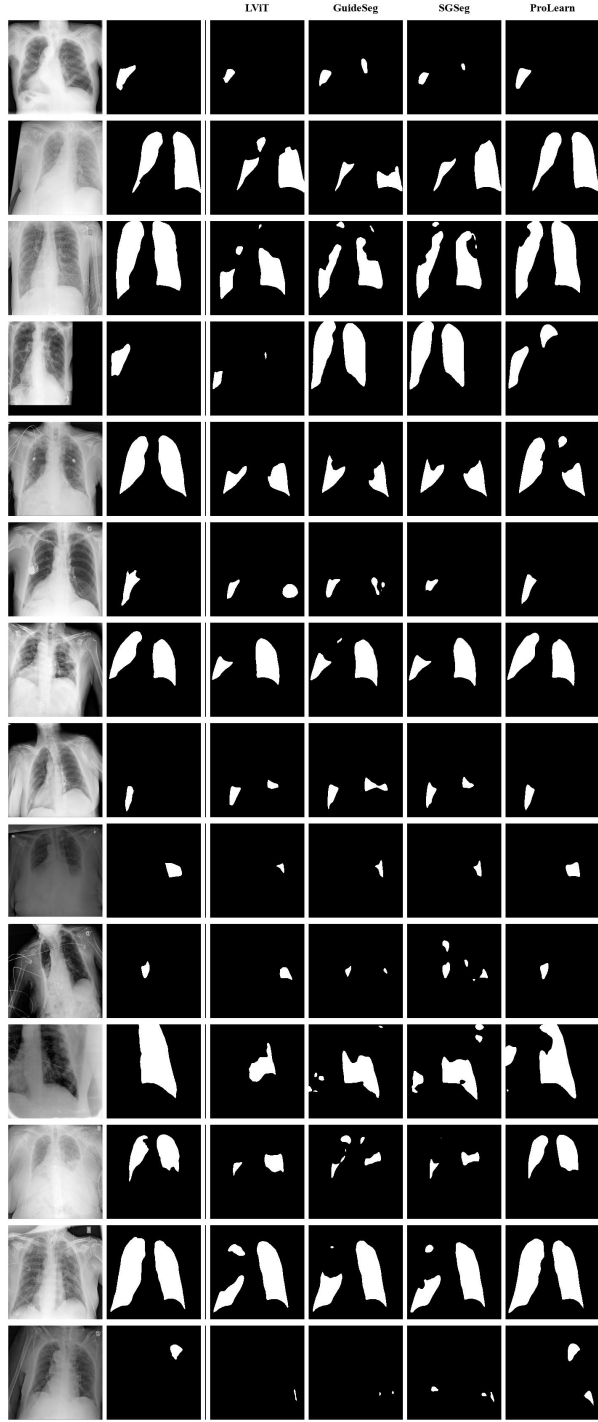


Figure A9. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on QaTa-COV19 dataset under 5% text availability.

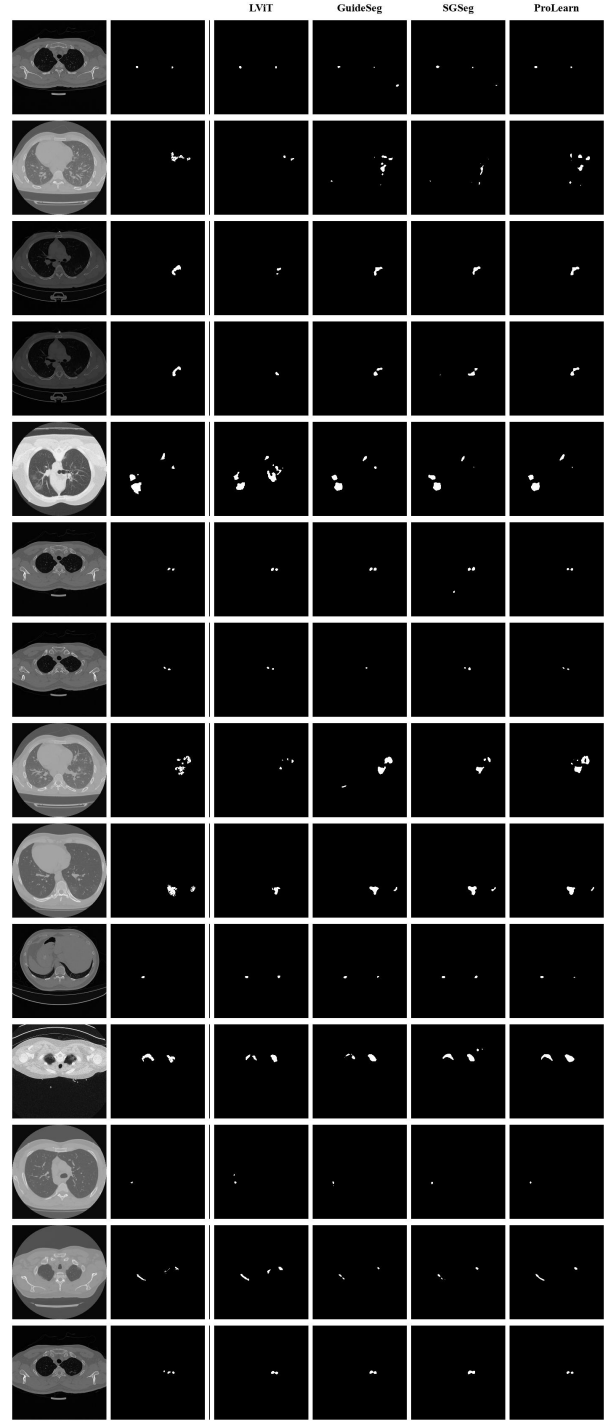


Figure A10. Comparison of segmentation results among LViT, GuideSeg, SGSeg, and our ProLearn on MosMedData+ dataset under 5% text availability.