# Supplementary – ClearCLIP: Decomposing CLIP Representations for Dense Vision-Language Inference

Mengcheng Lan<sup>1</sup>, Chaofeng Chen<sup>1</sup>, Yiping Ke<sup>2</sup>, Xinjiang Wang<sup>3</sup>, Litong Feng<sup>3</sup>\*, and Wayne Zhang<sup>3</sup>

 <sup>1</sup> S-Lab, Nanyang Technological University
<sup>2</sup> CCDS, Nanyang Technological University <sup>3</sup> SenseTime Research lanm0002@e.ntu.edu.sg {chaofeng.chen, ypke}@ntu.edu.sg {wangxinjiang, fenglitong, wayne.zhang}@sensetime.com

## Appendix

#### A Ablation study with different backbones and datasets

We showcase the results of the ablation study for each dataset across different CLIP models in Fig. 1. It's clear that our method, which involves removing the residual connection and FFN, markedly enhances the open-vocabulary semantic segmentation capability of CLIP throughout all datasets. This enhancement is especially pronounced within the ViT-L/14 architecture, characterized by a larger norm of residual connection. These findings conclusively affirm the efficacy of our proposed methodology.

#### **B** Impact of channel-wise residual features

In this part, we investigate the effect of residual features with low intensity. Specifically, we conduct experiments by selectively reintroducing channels from residual features that have lower average values. We report the results of eliminating the top  $\beta$  high-value channels and the effect of normalizing  $X_{\rm res}$  in Tab. 1. The best performance is achieved when  $\beta \geq 70\%$ . Additionally, normalizing  $X_{\rm res}$  significantly reduces its scale, resulting in performance comparable to  $\beta \geq 70\%$ . These findings support our hypothesis that high-level supervision in CLIP emphasizes global feature direction in the residual latent space, which introduces noise into the residual features. For simplicity, we eliminate all channels in  $X_{\rm res}$ .

#### C Integration across models

Our solution serves as a free lunch applicable to various architectures and segmentation models with *just 2-3 lines of code modification*. Specifically, for

<sup>\*</sup> Corresponding author.



**Fig. 1:** Ablation study on each dataset under different architectures and attention mechanisms.  $\bigcirc$ : original CLIP;  $\triangle$ : CLIP w/o residual connection;  $\Leftrightarrow$ : CLIP w/o residual connection and FFN.

MaskCLIP and SCLIP, we achieve this by eliminating the residual connection and Feed-Forward Network (FFN) of the last self-attention layer. For GEM, we utilize the attention output from the final layer as the final representation. Importantly, we preserve the original attention mechanisms of these methods. For baseline models, *i.e.*, CLIP, BLIP, OpenCLIP, and MetaCLIP, we enhance them by incorporating our complete solution. The performance of different models on five datasets is summarized in Tab. 2. The results demonstrate that our solution consistently enhances the performance of existing models in open-vocabulary semantic segmentation tasks, showcasing its exceptional generalizability.

## D Visualization of feature maps

To intuitively demonstrate how the residual connections affect the performance, we visualize the feature maps of  $X_{\rm res}$ ,  $X_{\rm attn}$ , and  $X_{\rm sum}$  for two randomly selected samples in Fig. 2. It is obvious that the  $X_{\rm res}$  feature maps associated with the residual connections are characterized by peak values in one channel (highlighted in a red box), significantly surpassing the other channels. And  $X_{\rm sum}$  is similar

$\beta$ (%)	0 5	10	30	50	70	100	Norm
Avg. 2	22.1 30.2	2 33.5	37.4	38.0	38.1	38.1	38.1

Table 1: Average performance (mIoU) over all 8 datasets.

to  $X_{\rm res}$ , indicating the big influence of  $X_{\rm res}$  to the final feature. Conversely, the feature maps in  $X_{\rm attn}$  demonstrate a more uniform distribution across channels. Given that the segmentation map is derived from the cosine similarity of feature vectors at each spatial location, such a disparity implies that the features in  $X_{\rm sum}$ and  $X_{\rm res}$  are less discernible compared to those in  $X_{\rm attn}$ , thereby introducing noise into the segmentation results. This observation supports our proposal that the high-level supervision in CLIP emphasizes the global feature direction in the residual latent space, making local feature vectors less distinguishable and leading to noise in residual features.

#### E Additional qualitative examples

In this part, we present more qualitative results comparison between ClearCLIP and state-of-the-art methods. Figs. 3 and 4 show the results from COCOStuff, ADE20K and Pascal Context59 datasets respectively. Similar to the findings in the main text, the results of ClearCLIP exhibit much less noise than other methods, further underscoring the superiority of our method.

#### References

- Bousselham, W., Petersen, F., Ferrari, V., Kuehne, H.: Grounding everything: Emerging localization properties in vision-language transformers. arXiv preprint arXiv:2312.00878 (2023)
- Cherti, M., Beaumont, R., Wightman, R., Wortsman, M., Ilharco, G., Gordon, C., Schuhmann, C., Schmidt, L., Jitsev, J.: Reproducible scaling laws for contrastive language-image learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2818–2829 (2023)
- Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In: International Conference on Machine Learning. pp. 12888–12900. PMLR (2022)
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
- Wang, F., Mei, J., Yuille, A.: Sclip: Rethinking self-attention for dense visionlanguage inference. arXiv preprint arXiv:2312.01597 (2023)
- Xu, H., Xie, S., Tan, X.E., Huang, P.Y., Howes, R., Sharma, V., Li, S.W., Ghosh, G., Zettlemoyer, L., Feichtenhofer, C.: Demystifying clip data. arXiv preprint arXiv:2309.16671 (2023)
- Zhou, C., Loy, C.C., Dai, B.: Extract free dense labels from clip. In: European Conference on Computer Vision. pp. 696–712. Springer (2022)

### 4 M.Lan et al.

	VOC20	Context59	Stuff	Cityscape	ADE20K	Avg.
CLIP [4] +ClearCLIP	41.8 80.9	$9.2 \\ 35.9$	4.4 23.9	$5.5 \\ 30.0$	$2.1 \\ 16.7$	12.6 37.5 + 24.9
BLIP [3] +ClearCLIP	$37.3 \\ 73.5$	$7.8 \\ 31.4$	$5.4 \\ 21.3$	4.3 23.8	$2.0 \\ 13.5$	11.4 32.7 +21.3
OpenCLIP [2] +ClearCLIP	47.2 81.4	$9.0 \\ 34.1$	$5.0 \\ 23.1$	$5.1 \\ 31.8$	$2.9 \\ 18.9$	13.8 37.9 +24.1
MetaCLIP [6] +ClearCLIP	35.4 78.3	8.1 34.8	$4.3 \\ 23.5$	$5.0 \\ 27.9$	$2.2 \\ 17.4$	11.0 36.4 + 25.4
MaskCLIP [7] +ClearCLIP	$74.9 \\ 61.4$	$26.4 \\ 28.3$	$\begin{array}{c} 16.4 \\ 18.4 \end{array}$	$12.6 \\ 24.7$	9.8 13.6	28.0 29.5 + 1.8
SCLIP [5] +ClearCLIP	78.2 77.9	$33.0 \\ 35.6$	21.1 23.6	29.1 31.0	$\begin{array}{c} 14.6 \\ 17.0 \end{array}$	35.2 37.9 + 1.6
GEM [1] +ClearCLIP	79.9 80.2	$\begin{array}{c} 35.9\\ 36.5\end{array}$	$23.7 \\ 24.4$	$\begin{array}{c} 30.8\\ 30.5\end{array}$	$\begin{array}{c} 15.7 \\ 17.4 \end{array}$	$37.2 \\ 37.8 + 0.6$
CLIP [4] +ClearCLIP	$\begin{array}{c} 15.8\\ 80.0\end{array}$	$4.5 \\ 29.6$	2.4 19.9	$2.9 \\ 27.9$	$1.2 \\ 15.0$	5.4 34.5 + 29.1
BLIP [3] +ClearCLIP	$22.5 \\ 67.5$	$5.8 \\ 16.8$	$2.4 \\ 11.5$	$\begin{array}{c} 3.8\\ 9.3\end{array}$	$1.5 \\ 7.1$	7.2 22.4 +15.2
OpenCLIP [2] +ClearCLIP	$39.7 \\ 65.3$	$7.0 \\ 27.9$	4.1 19.5	$3.9 \\ 26.4$	$2.3 \\ 16.0$	$\frac{11.4}{31.0 + 19.6}$
MetaCLIP [6] +ClearCLIP	22.7 78.2	$6.2 \\ 30.3$	$3.6 \\ 20.5$	$5.1 \\ 25.6$	$\begin{array}{c} 2.2\\ 16.4 \end{array}$	8.0 34.2 + 26.2
MaskCLIP [7] +ClearCLIP	$30.1 \\ 65.1$	$12.6 \\ 26.5$	8.9 17.6	$10.1 \\ 21.2$	$6.9 \\ 15.1$	13.7 29.1 +11.1
SCLIP [5] +ClearCLIP	60.3 79.2	$20.5 \\ 30.6$	$13.1 \\ 20.5$	$17.0 \\ 27.8$	$7.1 \\ 15.6$	$\begin{array}{c} 23.6\\ 34.7 + 15.4\end{array}$
GEM [1] +ClearCLIP	80.3 79.7	26.4 29.9	$17.6 \\ 19.4$	$22.6 \\ 25.9$	$11.6 \\ 14.2$	31.7 33.8 +2.1

**Table 2:** Average performance (mIoU) over 5 datasets without background class basedon ViT-BaseandLargearchitectures.



Fig. 2: Visualization of feature maps with CLIP for two randomly selected examples from the COCOStuff dataset. The first row shows the first 64 feature maps of each type, while the second row displays all 768 feature maps of each type.



Fig. 3: Qualitative comparison between different open-vocabulary segmentation methods on (a) COCOStuff and (b) ADE20K datasets.



 ${\bf Fig. 4: Qualitative \ comparison \ between \ different \ open-vocabulary \ segmentation \ methods \ on \ the \ Pascal \ Context59 \ dataset.}$