

---

# Parameter-efficient Fine-tuning in Hyperspherical Space for Open-vocabulary Semantic Segmentation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Open-vocabulary semantic segmentation seeks to label each pixel in an image  
2 with arbitrary text descriptions. Vision-language foundation models, especially  
3 CLIP, have recently emerged as powerful tools for acquiring open-vocabulary  
4 capabilities. However, fine-tuning CLIP to equip it with pixel-level prediction  
5 ability often suffers three issues: 1) high computational cost, 2) misalignment  
6 between the two inherent modalities of CLIP, and 3) degraded generalization ability  
7 on unseen categories. To address these issues, we propose H-CLIP, a symmetrical  
8 parameter-efficient fine-tuning (PEFT) strategy conducted in hyperspherical space  
9 for both of the two CLIP modalities. Specifically, the PEFT strategy is achieved  
10 by a series of efficient block-diagonal learnable transformation matrices and a  
11 dual cross-relation communication module among all learnable matrices. Since  
12 the PEFT strategy is conducted symmetrically to the two CLIP modalities, the  
13 misalignment between them is mitigated. Furthermore, we apply an additional  
14 constraint to PEFT on the CLIP text encoder according to the hyperspherical energy  
15 principle, i.e., minimizing hyperspherical energy during fine-tuning preserves the  
16 intrinsic structure of the original parameter space, to prevent the destruction of  
17 the generalization ability offered by the CLIP text encoder. Extensive evaluations  
18 across various benchmarks show that H-CLIP achieves new SOTA open-vocabulary  
19 semantic segmentation results while only requiring updating approximately 4% of  
20 the total parameters of CLIP.

## 21 1 Introduction

22 The aim of open-vocabulary semantic segmentation is to create a segmentation model capable of  
23 labeling each pixel in an image with categories that are not limited to a specific closed set according to  
24 text descriptions. Vision-language foundation models [43, 5, 34, 39, 11, 17, 26, 21, 27, 13, 18, 29, 10,  
25 28, 45], especially CLIP [39], are often utilized to endow open-vocabulary recognition capabilities.  
26 Consequently, open-vocabulary semantic segmentation essentially boil down to transferring these  
27 vision-language foundation models, originally trained with image-level supervision, to perform  
28 pixel-level predictions.

29 To this end, current methods [52, 48, 7, 50] typically fine-tune CLIP on a benchmark dataset with  
30 segmentation annotations, i.e., COCO [2], to equip it with the segmentation ability. However, this  
31 often leads to three main issues. First, fine-tuning CLIP on limited categories would affect its  
32 generalization ability, resulting in significant performance degradation on unseen categories. Second,  
33 current fine-tuning strategies are usually asymmetrical, which inevitably causes a misalignment  
34 between the two inherent modalities of CLIP, i.e., image and text [52], which may lead to sub-  
35 optimal performance. Third, although remarkable performance gains, these approaches often rely on  
36 computationally extensive full fine-tuning, which raises concerns about scalability and affordability.

37 To address these issues, we propose a symmetric parameter-efficient fine-tuning (PEFT) strategy  
38 for CLIP, dubbed H-CLIP. Specifically, we implement this PEFT through a partial orthogonal fine-  
39 tuning (POF) strategy, which introduces a series of efficient block-diagonal learnable transformation  
40 matrices into the hyperspherical space. Then, to preserve CLIP’s generalization ability, we leverage  
41 the hyperspherical energy principle [32, 38], which suggests that maintaining the same hyperspherical  
42 energy during fine-tuning preserves the intrinsic structure, i.e., generalization ability. In light of this,  
43 we upgrade our POF by incorporating orthogonal constraints in the learnable matrices for updating  
44 CLIP’s text encoder, as orthogonal transformations keep the hyperspherical energy unchanged during  
45 fine-tuning. Subsequently, we introduce a dual cross-relation communication (DCRC) module to  
46 explicitly encourage cross-modal and cross-layer communications within all learnable matrices.  
47 This communication not only preserves the hyperspherical energy but also further mitigating the  
48 misalignment problem.

49 Extensive results demonstrate that H-CLIP achieves new state-of-the-art open-vocabulary semantic  
50 segmentation results across three benchmarks by fine-tuning CLIP with approximately 4% of the  
51 total parameters of CLIP.

## 52 **2 Related Work**

### 53 **2.1 Open-vocabulary Semantic Segmentation**

54 Prior open-vocabulary semantic segmentation works typically perform this task through leveraging  
55 CLIP [39]. initial efforts like [56] directly fine-tune CLIP on mainstream segmentation datasets, e.g.,  
56 COCO [2]. However, they claim that fine-tuning CLIP’s encoder significantly reduces its ability  
57 to generalize to unseen classes. To address this issue, some methods [15, 8, 51, 49] swing to the  
58 opposite extreme, fine-tuning an additional mask generator [6] for segmentation while keeping CLIP  
59 frozen to maintain generalization-oriented recognition. However, this frozen parameter space lacks  
60 segmentation awareness, resulting in a misalignment between regions and text descriptions [30].  
61 Other studies [52, 50, 7] propose an advanced solution that fine-tunes only selected parameters, e.g.,  
62 certain layers of CLIP, to enable pixel-level predictions while keeping most of CLIP’s parameters  
63 fixed, thus minimizing losing of generalization. Although the advantages are remarkable, these  
64 methods often work with a very small learning rate, implicitly encouraging a small deviation from  
65 the pre-trained CLIP, limiting the segmentation performance. In a nutshell, the trade-off between  
66 preserving CLIP’s generalization and learning segmentation knowledge persists, hindering the final  
67 performance. Based on the paradigm of existing fine-tuning-based methods, our method explores a  
68 better trade-off from a fresh viewpoint: hyperspherical space.

### 69 **2.2 Large-scale Model Fine-tuning**

70 Along with the improvement of large-scale foundation models [26, 34, 28, 23, 53, 42, 41, 40, 60],  
71 e.g., segment anything model [23], numerous fine-tuning works [37, 36, 4, 57, 58, 14, 54, 47, 31, 61]  
72 are proposed to adapt these models to various downstream scenarios. The core of these approaches  
73 lies in updating only limited parameters to capture the specific characteristics of different scenarios,  
74 while keeping most parameters fixed to maintain generalization. In contrast, fine-tuning CLIP  
75 for open-vocabulary semantic segmentation often meets a dilemma. On the one hand, limited  
76 parameters typically fall short in facilitating the transition from a classification model, i.e., CLIP, to a  
77 segmentation task. On the other hand, directly increasing the number of trainable parameters risks  
78 undermining CLIP’s ability to generalize to unseen classes, as experimented in CAT-Seg [7]. Most  
79 methods [52, 48] solve this issue by simply freezing CLIP’s text encoder and fine-tuning its image  
80 encoder, inevitably causing misalignment between the two modalities of CLIP. In this paper, we shed  
81 light on how to preserve generalization in a symmetric parameter-efficient fine-tuning manner and  
82 strive to explore an appropriate fine-tuning method for open-vocabulary semantic segmentation.

## 83 **3 Preliminaries**

### 84 **3.1 Hyperspherical Energy**

85 Existing fine-tuning methods implicitly assume that a smaller Euclidean distance between the fine-  
86 tuned model and the pre-trained model indicates better preservation of the pre-trained ability. However,

87 the Euclidean difference is unable to fully capture the degree of semantic preservation. According to  
 88 the inspiration from Thomson problem[44] which is to determine the minimum electrostatic potential  
 89 energy configuration of  $N$  mutually-repelling electrons on the surface of a unit sphere, we adopt the  
 90 *Hyperspherical Energy* to characterize the diversity of the model. The hyperspherical energy function  
 91 of a fully connected layer  $\mathbf{W}$  is defined as  $\text{HE}(\mathbf{W}) := \sum_{i \neq j} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\|^{-1}$ , where  $\hat{\mathbf{w}}_i := \mathbf{w}_i / \|\mathbf{w}_i\|$   
 92 denotes the  $i$ -th normalized neuron. The power of the model representation can be characterized by the  
 93 hyperspherical energy of its neurons. Higher energy implies higher redundancy, while lower energy  
 94 indicates that these neurons of the model are more diverse. For the original semantic information not  
 95 to be destroyed in the case of fine-tuning, we hypothesize that a good fine-tuning model should have  
 96 a minimal difference in hyperspherical energy compared to the pre-trained model:

$$\min_{\mathbf{W}} \|\text{HE}(\mathbf{W}) - \text{HE}(\mathbf{W}^0)\| \Leftrightarrow \min_{\mathbf{W}} \left\| \sum_{i \neq j} \|\hat{\mathbf{w}}_i - \hat{\mathbf{w}}_j\|^{-1} - \sum_{i \neq j} \|\hat{\mathbf{w}}_i^0 - \hat{\mathbf{w}}_j^0\|^{-1} \right\|. \quad (1)$$

97 One can easily observe that the attainable minimum is zero for Eq. (1). In this case, the hyperspherical  
 98 energy should satisfy an invariance property (the application of the same orthogonal transformation  
 99 for all neurons demonstrates the pairwise hyperspherical similarity). Based on the hyperspherical  
 100 energy invariance property, the minimum of zero can be achieved as long as  $\mathbf{W}$  and  $\mathbf{W}^0$  differ  
 101 only up to a rotation or reflection, i.e.,  $\mathbf{W} = \mathbf{R}\mathbf{W}^0$  in which  $\mathbf{R} \in \mathbb{R}^{d \times d}$  is an orthogonal matrix (The  
 102 determinant 1 or  $-1$  means rotation or reflection, respectively).

### 103 3.2 Notation of Tensor Product

104 In this section, we introduce the fundamental concept underlying our DCRC (Sec. 4.3): tensor product.  
 105 A  $p$ -order tensor is indexed by  $p$  indices and can be represented as a multidimensional array of data.  
 106 Formally, a  $p$ -order tensor  $\mathcal{A}$  can be written as  $\mathcal{A} = (a_{i_1, i_2, \dots, i_p}) \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_p}$ . Slices of a tensor  
 107 are matrices defined from the tensor by holding all but two indices constant. For a 3-order tensor,  
 108  $\mathcal{A}(:, :, k)$  corresponds the  $k^{\text{th}}$  frontal slice. For  $p$ -order tensors, matrix slices of  $p$ -order tensors can be  
 109 referenced using linear indexing by reshaping the tensor into an  $n_1 \times n_2 \times n_3 n_4 \dots n_p$  3-order tensor  
 110 and referring to the  $k^{\text{th}}$  frontal slice as  $\mathcal{A}(:, :, k)$ .  $\mathcal{A}_i \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_{p-1}}$  for  $i = 1, \dots, n_p$  denotes the  
 111  $(p-1)$ -order tensor created by holding the  $p$ th index of  $\mathcal{A}$  fixed at  $i$ . It is possible to create a tensor  
 112 in a block circulant pattern, where each block is a tensor of  $(p-1)$ -order:

$$\text{circ}(\mathcal{A}) = \begin{bmatrix} \mathcal{A}_1 & \mathcal{A}_{n_p} & \mathcal{A}_{n_p-1} & \dots & \mathcal{A}_2 \\ \mathcal{A}_2 & \mathcal{A}_1 & \mathcal{A}_{n_p} & \dots & \mathcal{A}_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathcal{A}_{n_p} & \mathcal{A}_{n_p-1} & \mathcal{A}_{n_p-2} & \dots & \mathcal{A}_1 \end{bmatrix},$$

113 where  $\text{circ}(\cdot)$  creates a block circulant tensor and the size of  $\text{circ}(\mathcal{A})$  is  $(n_1 n_p \times n_2 n_p \times \dots \times n_{p-2} n_p \times$   
 114  $n_{p-1})$ . define  $\text{unfold}(\cdot)$  to take an  $n_1 \times \dots \times n_p$  tensor  $\mathcal{A}$  and return an  $n_1 n_p \times n_2 \times \dots \times n_{p-1}$  block  
 115 tensor in the following way:

$$\text{unfold}(\mathcal{A}) = [\mathcal{A}_1 \quad \mathcal{A}_2 \quad \dots \quad \mathcal{A}_{n_p}]^T.$$

116 The operation that takes  $\text{unfold}$  back to tensor form is the ‘‘fold’’ command. Specially,  $\text{fold}(\cdot, n_p)$   
 117 takes an  $n_1 n_p \times n_2 \times \dots \times n_{p-1}$  block tensor and returns an  $n_1 \times \dots \times n_p$  tensor, defined as:

$$\text{fold}(\text{unfold}(\mathcal{A}), n_p) = \mathcal{A}.$$

## 118 4 Methodology

### 119 4.1 Overview of H-CLIP

120 Fig. 2 illustrates the proposed H-CLIP framework, which is based on two core components: (1) POF  
 121 updates the pre-trained parameter space of CLIP using a series of block-diagonal transformation  
 122 matrices. According to analysis in Sec. 1, each parameter matrix in CLIP’s text encoder is orthogonal  
 123 to preserve generalization. (2) DCRC incorporates cross-modal and cross-layer communication  
 124 within all tunable matrices, facilitating alignment between different modalities.

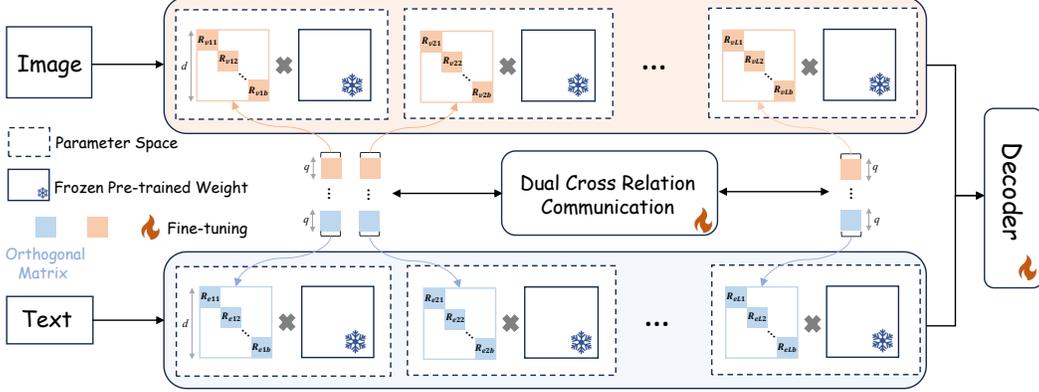


Figure 1: A schematic representation of H-CLIP. In the H-CLIP framework, we propose a partial orthogonal fine-tuning strategy, where each pre-trained weight matrix is paired with a tuned block-diagonal transformation matrix, some of which are orthogonal to preserve generalization. Then, we introduce a dual cross-relation communication mechanism to facilitate communication among all matrices, enabling alignment between different modalities.

## 125 4.2 Partial Orthogonal Fine-tuning

126 The core idea of partial orthogonal fine-tuning (POF) is to introduce the concept of hyperspherical  
 127 space for fine-tuning CLIP. In this hyperspherical space, we fine-tune CLIP’s text encoder under an  
 128 orthogonality design principle from OFT [38] to preserve the hyperspherical energy of the pre-trained  
 129 parameter space. Similarly, we use Cayley parameterization [3] to ensure a tunable matrix  $\mathbf{R}$  is  
 130 strictly orthogonal, formally as:

$$\mathbf{R} = (\mathbf{I} + \mathbf{Q})(\mathbf{I} - \mathbf{Q})^{-1}, \quad (2)$$

131 where  $\mathbf{Q}$  is skew-symmetric. Here, for  $\mathbf{R}$  in CLIP’s image encoder, we remove the orthogonality  
 132 constraint, defined as:

$$\mathbf{R}^\top \mathbf{R} = \mathbf{R} \mathbf{R}^\top = \mathbf{I}, \quad (3)$$

133 where  $\mathbf{I}$  is an identity matrix. Considering the relatively large dimension  $d$  of the pre-trained matrix,  
 134 for better efficiency, we introduce a block-diagonal structure by parameterizing  $\mathbf{R}$  with  $b$  blocks,  
 135 formally as:

$$\mathbf{R} = \text{diag}(\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_i, \dots, \mathbf{R}_b) = \begin{bmatrix} \mathbf{R}_1 & & \\ & \ddots & \\ & & \mathbf{R}_b \end{bmatrix}, \quad (4)$$

136 where  $\mathbf{R}_i \in \mathbb{R}^{d/b \times d/b}$ . Specifically, denote  $\mathcal{R}^V = \{\mathbf{R}_{v1}, \dots, \mathbf{R}_{v\ell}, \dots, \mathbf{R}_{vL}\}$  and  $\mathcal{R}^E =$   
 137  $\{\mathbf{R}_{e1}, \dots, \mathbf{R}_{e\ell}, \dots, \mathbf{R}_{eL}\}$  as the sets of block-diagonal matrices in CLIP’s image encoder and  
 138 text encoder, respectively, where  $L$  is its number of Transformer layers,  $\mathbf{R}_{v\ell} \in \mathbb{R}^{d_v \times d_v}$ , and  
 139  $\mathbf{R}_{e\ell} \in \mathbb{R}^{d_e \times d_e}$ . For simplicity, we set  $d_v = d_e = d$ . Overall, we develop a H-CLIP framework,  
 140 and for an input feature map  $\mathbf{M}_\ell$  in the  $\ell^{\text{th}}$  Transformer layer of CLIP, the right branch produces the  
 141 adjusted feature map via H-CLIP,  $\tilde{\mathbf{M}}_\ell$ , formally via:

$$\tilde{\mathbf{M}}_\ell = \begin{cases} \mathcal{F}_\ell(\mathbf{M}_\ell; \mathbf{R}_\ell \mathbf{W}_\ell), & \text{if } \mathbf{R}_\ell \in \mathcal{R}^V \\ \mathcal{F}_\ell(\mathbf{M}_\ell; \mathbf{R}_\ell \mathbf{W}_\ell), & \text{s.t. } \mathbf{R}_\ell^\top \mathbf{R}_\ell = \mathbf{R}_\ell \mathbf{R}_\ell^\top = \mathbf{I} \text{ otherwise,} \end{cases} \quad (5)$$

142 where  $\mathbf{W}_\ell$  is a pre-trained weight matrix in  $\ell^{\text{th}}$  layer of CLIP’s encoder, and  $\mathcal{F}_\ell$  represents  $\ell^{\text{th}}$  layer of  
 143 CLIP’s encoder. During the fine-tuning phase, H-CLIP is fine-tuned in conjunction with the original  
 144 parameter space of CLIP, which is loaded from the pre-trained checkpoint and remains frozen.

## 145 4.3 Dual Cross Relation Communication

146 Although in POF, we relax the orthogonal constraint for CLIP’s image encoder to learn segmentation  
 147 knowledge, each layer of the image encoder still incorporates a limited number of parameters, which  
 148 largely restricts the flexibility of the projection adjustment due to the limitation of Hidden Markov

149 Chain along layers [24, 46, 36]. To address this limitation, one might consider fully fine-tuning  
 150 instead of using a small number of parameters. However, this approach can cause a misalignment  
 151 between image and text features in CLIP, resulting in sub-optimal performance [52]. Based on  
 152 the above analysis, we introduce Dual Cross-Relation Communication (DCRC), which facilitates  
 153 interaction among different layers and modalities (i.e., text and image). DCRC explicitly enhances  
 154 the flexibility of fine-tuned projection adjustments and prevents misalignment issues.

155 DCRC introduces cross-layer and cross-modality communication among different block-diagonal  
 156 matrices, achieved through two relation projections. To do this, we first treat all blocks in  $\ell^{\text{th}}$  layer as  
 157 an individual slice in this 3-order tensor  $\mathcal{T}_\ell$ , which is derived as follows:

$$\mathcal{T}_\ell = [\mathbf{R}_{v\ell 1}, \mathbf{R}_{e\ell 1}, \dots, \mathbf{R}_{v\ell i}, \mathbf{R}_{e\ell i}, \dots, \mathbf{R}_{v\ell b}, \mathbf{R}_{e\ell b}] \in \mathbb{R}^{q \times q \times (b+b)}, \quad (6)$$

158 Where  $q = d/b$ . Then, we treat the tensor  $\mathcal{T}_\ell$  as an individual slice within a 4-order tensor  $\mathcal{T}$ , defined  
 159 as follows:

$$\mathcal{T} = [\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_\ell, \dots, \mathcal{T}_L] \in \mathbb{R}^{q \times q \times (b+b) \times L}. \quad (7)$$

160 Initially, according to the characteristics of gradient propagation in deep learning theory, i.e., chain  
 161 rule, each frontal slice  $\mathbf{R}_{\cdot\ell i} \in \{\mathbb{R}^{q \times q}\}^{(b+b) \times L}$  is updated sequentially in CLIP’s encoder. As a result,  
 162 updating the  $\mathcal{T}$  lacks cross-frontal-slice communication, limiting the flexibility of adjusting fine-tuned  
 163 projection. To address this, we introduce two special tensor products, i.e., **3-order T-product** and  
 164 **Higher-order T-product**.

165 **Definition 4.1(3-order T-product)** For  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$  and  $\mathcal{B} \in \mathbb{R}^{n_2 \times l \times n_3}$ , the 3-order T-product  
 166  $\mathcal{C} \in \mathbb{R}^{n_1 \times l \times n_3} = \mathcal{A} * \mathcal{B}$  is defined as:

$$\mathcal{C} = \mathcal{A} * \mathcal{B} = \text{fold}(\text{circ}(\mathcal{A}) \cdot \text{unfold}(\mathcal{B})), \quad (8)$$

167 where “ $\cdot$ ” represents standard matrix product.

168 **Definition 4.2(Higher-order T-product)** For  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3 \cdots \times n_p}$  and  $\mathcal{B} \in \mathbb{R}^{n_2 \times l \times n_3 \cdots \times n_p}$ , the  
 169 High-order T-product  $\mathcal{C} \in \mathbb{R}^{n_1 \times l \times n_3 \cdots \times n_p} = \mathcal{A} * \mathcal{B}$  is defined as:

$$\mathcal{C} = \mathcal{A} * \mathcal{B} = \text{fold}(\text{circ}(\mathcal{A}) * \text{unfold}(\mathcal{B})). \quad (9)$$

170 If  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$ , according to the **3-order T-product**, there is an invertible transform  $S_3(\cdot) :$   
 171  $\mathbb{R}^{n_1 \times n_2 \times n_3} \rightarrow \mathbb{R}^{n_1 \times n_2 \times n_3}$  in third dimension and it transform the Eq. (8) as:

$$\mathcal{C} = S_3^{-1}(S_3(\mathcal{A}) \odot S_3(\mathcal{B})) = S_3^{-1}(\bar{\mathcal{A}} \odot \bar{\mathcal{B}}) = S_3^{-1}(\bar{\mathcal{C}}), \quad (10)$$

172 where  $\bar{\mathcal{C}} = \bar{\mathcal{A}} \odot \bar{\mathcal{B}}$  denotes the frontal-slice-wise product (Definition 2.1 refers to [19])  $\bar{\mathcal{C}}(;, ;, i) =$   
 173  $\bar{\mathcal{A}}(;, ;, i) \cdot \bar{\mathcal{B}}(;, ;, i), i = 1, 2, \dots, n_3$  and  $S_3^{-1}(\cdot)$  is the inverse transform of  $S_3(\cdot)$ . According to the  
 174 definition of the frontal-slice-wise product, the invertible transform  $S_3(\cdot)$  is formulated as:

$$\bar{\mathcal{A}} = S_3(\mathcal{A}) = \mathcal{A} \times_3 \mathbf{S}_3, \quad (11)$$

175 where “ $\times_3$ ” denotes the mode-3 product and  $\mathbf{S}_3 \in \mathbb{R}^{n_3 \times n_3}$  is an arbitrary invertible matrix. Similarly,  
 176 the inverse transform of Eq. (11) is derived as:

$$\mathcal{A} = S_3^{-1}(\bar{\mathcal{A}}) = \bar{\mathcal{A}} \times_3 \mathbf{S}_3^{-1}. \quad (12)$$

177 Similarly, if  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_p}$ , according to the **Higer-order T-product**, there are invertible  
 178 transform  $S_i(\cdot) : \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_p} \rightarrow \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_p}, i = 3, 4, \dots, p$  in  $i^{\text{th}}$  dimension and they  
 179 transform the Eq. (9) as:

$$\mathcal{C} = \tilde{S}^{-1}(\tilde{S}(\mathcal{A}) \odot \tilde{S}(\mathcal{B})) = \tilde{S}^{-1}(\bar{\mathcal{A}} \odot \bar{\mathcal{B}}) = \tilde{S}^{-1}(\bar{\mathcal{C}}), \quad (13)$$

180 where  $\tilde{S}(\mathcal{A}) = S_p(S_{p-1}(\cdots S_3(\mathcal{A}) \cdots))$ ,  $\bar{\mathcal{C}} = \bar{\mathcal{A}} \odot \bar{\mathcal{B}}$  denotes the frontal-slice-wise product  
 181  $\bar{\mathcal{C}}(;, ;, i) = \bar{\mathcal{A}}(;, ;, i) \cdot \bar{\mathcal{B}}(;, ;, i), i = 1, 2, \dots, n_3 n_4 \cdots n_p$  and  $\tilde{S}^{-1}(\cdot)$  is the inverse transform  
 182 of  $\tilde{S}(\cdot)$ . Similarly, the inverse transform  $\tilde{S}(\cdot)$  is formulated as:

$$\bar{\mathcal{A}} = \tilde{S}(\mathcal{A}) = \mathcal{A} \times_3 \mathbf{S}_3 \times_4 \mathbf{S}_4 \cdots \times_p \mathbf{S}_p, \quad (14)$$

183 and its inverse transform is derived as:

$$\mathcal{A} = \tilde{S}^{-1}(\bar{\mathcal{A}}) = \bar{\mathcal{A}} \times_3 \mathbf{S}_3^{-1} \times_4 \mathbf{S}_4^{-1} \cdots \times_p \mathbf{S}_p^{-1}. \quad (15)$$

| Model                                  | VLM           | Additional Backbone | A-847       | PC-459      | A-150       | PC-59       | PAS-20      | PAS-20 <sup>b</sup> |
|--|---------------|---------------------|-------------|-------------|-------------|-------------|-------------|---------------------|
| <b>Traditional Fine-Tuning</b>         |               |                     |             |             |             |             |             |                     |
| ZS3Net [1]                             | -             | ResNet-101          | -           | -           | -           | 19.4        | 38.3        | -                   |
| LSeg [25]                              | CLIP ViT-B/32 | ResNet-101          | -           | -           | -           | -           | 47.4        | -                   |
| ZegFormer [8]                          | CLIP ViT-B/16 | ResNet-101          | 4.9         | 9.1         | 16.9        | 42.8        | 86.2        | 62.7                |
| ZSseg [51]                             | CLIP ViT-B/16 | ResNet-101          | 7.0         | -           | 20.5        | 47.7        | 88.4        | -                   |
| OpenSeg [15]                           | ALIGN         | ResNet-101          | 4.4         | 7.9         | 17.5        | 40.1        | -           | 63.8                |
| OVSeg [30]                             | CLIP ViT-B/16 | ResNet-101c         | 7.1         | 11.0        | 24.8        | 53.3        | 92.6        | -                   |
| ZegCLIP [59]                           | CLIP ViT-B/16 | -                   | -           | -           | -           | 41.2        | 93.6        | -                   |
| CAT-Seg [7]                            | CLIP ViT-B/16 | -                   | <u>12.0</u> | <u>19.0</u> | <u>31.8</u> | <u>57.5</u> | <u>94.6</u> | <u>77.3</u>         |
| <b>Parameter-efficient Fine-Tuning</b> |               |                     |             |             |             |             |             |                     |
| SAN [50]                               | CLIP ViT-B/16 | -                   | 10.1        | 12.6        | 27.5        | 53.8        | 94.0        | -                   |
| Ours                                   | CLIP ViT-B/16 | -                   | <b>12.4</b> | <b>19.3</b> | <b>32.4</b> | <b>57.9</b> | <b>95.2</b> | <b>78.2</b>         |
| <b>Traditional Fine-Tuning</b>         |               |                     |             |             |             |             |             |                     |
| LSeg [25]                              | CLIP ViT-B/32 | ViT-L/16            | -           | -           | -           | -           | 52.3        | -                   |
| OpenSeg [15]                           | ALIGN         | Eff-B7              | 8.1         | 11.5        | 26.4        | 44.8        | -           | 70.2                |
| OVSeg [30]                             | CLIP ViT-L/14 | Swin-B              | 9.0         | 12.4        | 29.6        | 55.7        | 94.5        | -                   |
| SAN [50]                               | CLIP ViT-L/14 | -                   | 12.4        | 15.7        | 32.1        | 57.7        | 94.6        | -                   |
| ODISE [49]                             | CLIP ViT-L/14 | Stable Diffusion    | 11.1        | 14.5        | 29.9        | 57.3        | -           | -                   |
| CAT-Seg [7]                            | CLIP ViT-L/14 | -                   | <u>16.0</u> | <u>23.8</u> | <u>37.9</u> | <u>63.3</u> | <u>97.0</u> | <u>82.5</u>         |
| <b>Parameter-efficient Fine-Tuning</b> |               |                     |             |             |             |             |             |                     |
| SAN [50]                               | CLIP ViT-L/14 | -                   | 12.4        | 15.7        | 32.1        | 57.7        | 94.6        | -                   |
| Ours                                   | CLIP ViT-L/14 | -                   | <b>16.5</b> | <b>24.2</b> | <b>38.4</b> | <b>64.1</b> | <b>97.7</b> | <b>83.2</b>         |

Table 1: **Comparison with state-of-the-art methods on standard benchmarks.** The best-performing results are presented in bold, while the second-best results are underlined. ‘‘VLM’’: visual language model.

184 *Derivation.* please refer to supplementary material. ■

185 According to Eqs. (29), (30) and (31), we adopt its idea and design arbitrary invertible relation matrix  
186  $\mathbf{S}_3 \in \mathbb{R}^{(b+b) \times (b+b)}$  and  $\mathbf{S}_4 \in \mathbb{R}^{L \times L}$  to capture the cross-modality and cross-layer information in  $\mathcal{T}$ .  
187 Then the updated tensor  $\mathcal{T}_w$  is formulated as:

$$\mathcal{T}_w = \mathcal{T} \times_3 \mathbf{S}_3 \times_4 \mathbf{S}_4 \in \mathbb{R}^{q \times q \times (b+b) \times L}, \quad (16)$$

188 where the relation matrix  $\mathbf{S}_3$  and  $\mathbf{S}_4$  are learnable. To better capture the nonlinear interactions inside  
189 the whole parameter space, we further adopt  $k$  layers deep neural network (DNN)  $f_3(\cdot)$  and  $f_4(\cdot)$  to  
190 replace the transform  $\times_3 \mathbf{S}_3$  and  $\times_4 \mathbf{S}_4$ , respectively, and the DNN  $f_3(\cdot)$  is formulated as:

$$f_3(\mathcal{T}) = \sigma(\sigma(\cdots \sigma(\sigma(\mathcal{A} \times_3 \mathbf{W}_1) \times_3 \mathbf{W}_2) \cdots) \times \mathbf{W}_{k-1}) \times \mathbf{W}_k, \quad (17)$$

191 where  $\sigma(\cdot)$  is a nonlinear scalar function and matrices  $\{\mathbf{W}_j \in \mathbb{R}^{(b+b)}\}_{j=1}^k$ . The DNN  $f_4(\cdot)$  is similar.  
192 Finally, the  $\mathcal{T}$  is updated by  $\mathcal{T} = \mathcal{T} + \alpha \mathcal{T}_w$ , where  $\alpha \in \mathbb{R}^{(b+b) \times L}$  is a learnable parameter.

## 193 5 Experiments

### 194 5.1 Experimental Setup

195 **Datasets.** Following previous studies [7, 48], we utilizes the COCO-Stuff dataset [2] as our training  
196 set. This dataset comprises approximately 118,000 densely annotated images across 171 distinct  
197 semantic categories. During inference, we carry out comparisons with state-of-the-art methods across  
198 several semantic segmentation datasets, including ADE20K [55], PASCAL VOC [12], and PASCAL-  
199 Context [35]. **ADE20K [55]** is a classical semantic segmentation dataset comprising around 20,000  
200 training images and 2,000 validation images. Besides, it includes two different test sets: A-150 and  
201 A-847. The test set A-150 has 150 common categories, while the test set A-847 has 847 categories.  
202 **PASCAL VOC [12]** is a small dataset for semantic segmentation, which includes 1464 training  
203 images and 1449 validation images. The dataset contains 20 different foreground categories. We  
204 name it as PAS-20. In line with [7], we also report a score on PAS-20<sup>b</sup>, which involves ‘‘background’’  
205 as the 21st category. **PASCAL-Context [35]** is upgraded from the original PASCAL VOC dataset.  
206 It includes two different test sets: PC-59 and PC-459 for evaluation. The test set PC-59 has 59  
207 categories, while the test set PC-459 has 459 categories.

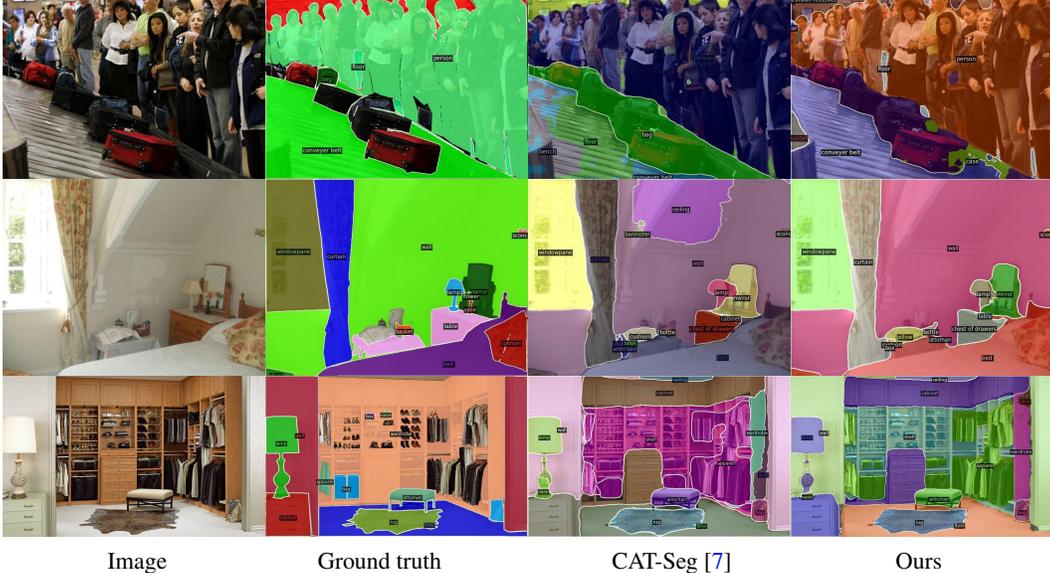


Figure 2: Comparison of qualitative results on ADE20K [55] with 150 categories. we compare Our method with CAT-Seg [7].

208 **Evaluation metric.** Following prior works [7, 48], we adopt mean Intersection over Union (mIoU)  
 209 to evaluate the semantic segmentation performance on the three benchmarks.

210 **Implementation Details.** We implement our method using the Transformer-based CLIP model.  
 211 Following the protocol established in [7], we evaluate our results on two versions of the CLIP model:  
 212 ViT-B/16 and ViT-L/14. For training, we use the Adam optimizer [22] with an initial learning rate  
 213 of  $5 \times 10^{-6}$  for CLIP, and a weight decay of  $10^{-4}$ . Training is conducted with one image per  
 214 mini-batch. We set  $q = 128$  for balancing efficiency and performance. The function  $f_3(\cdot)$  and  $f_4(\cdot)$   
 215 are implemented using two 2-layer MLPs. We act the cost-based approach provided in [7] as our  
 216 decoder. All models are trained over 80,000 iterations on 4 NVIDIA RTX 3090 GPUs.

## 217 5.2 Main Results

218 **Comparing to SOTAs.** Here, we compare  
 219 our proposed H-CLIP with several state-of-  
 220 the-art methods, as shown in Table 1, using  
 221 six test sets across three benchmarks. Over-  
 222 all, we achieve the best results. Most exist-  
 223 ing open-vocabulary semantic segmentation  
 224 methods employ traditional fine-tuning approaches, i.e., full or partial fine-tuning (tuning certain lay-  
 225 ers of CLIP). While these methods offer sufficient flexibility for learning new knowledge, they often  
 226 result in a significant performance drop on unseen classes, as observed with OVSeg [30]. Among  
 227 these methods, CAT-Seg [7] achieves performance comparable to ours. However, its fine-tuning  
 228 scheme is manually controlled through different layer combinations, necessitating a careful design to  
 229 balance generalization and flexibility, while ours does not suffer from such an issue. Then, compared  
 230 to SAN [50], another parameter-efficient fine-tuning method that introduces only a limited number of  
 231 tunable parameters, our approach significantly outperforms it, achieving improvements of 6.6% on  
 232 the PC-459 dataset and 3.9% on the PC-59 dataset with ViT-B/16 as the base model. These results  
 233 demonstrate the effectiveness of our method in preserving generalization while learning segmentation  
 234 knowledge.

235 **Qualitative results.** Here, we visualize our method’s representative example segmentation results  
 236 against prevailing methods, e.g., CAT-Seg [7] in the PC-459 dataset. As shown in Figs. 2 - 4, we  
 237 observe that our approach is able to generalize on diverse scenarios and produce more accurate  
 238 results.

| Methods    | OVSeg [30] | CAT-Seg [7] | SAN [50] | Ours       |
|------------|------------|-------------|----------|------------|
| Param. (M) | 147.2      | 25.6        | 8.4      | <b>5.6</b> |

Table 2: **Efficiency comparison** in terms of learnable parameters.

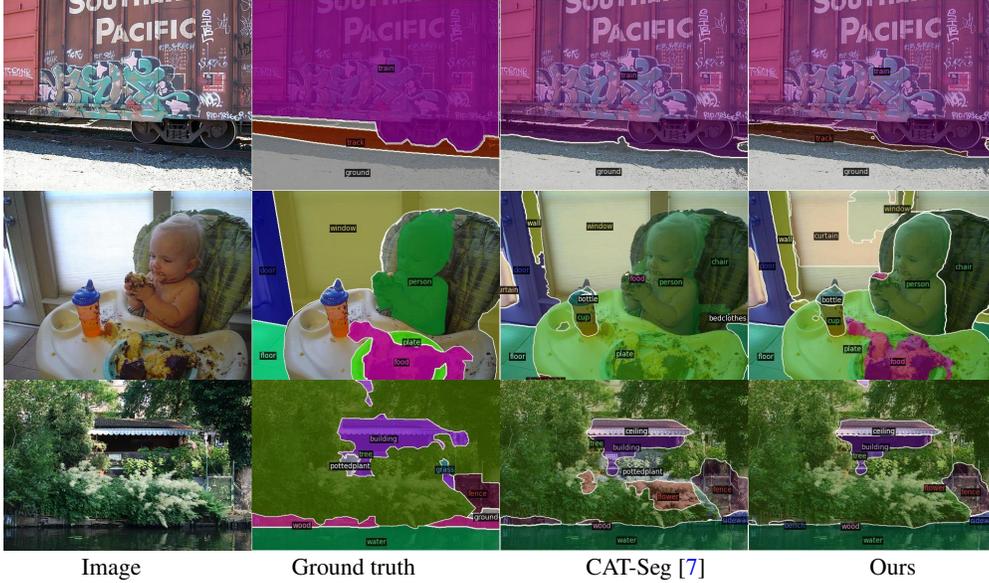


Figure 3: Comparison of qualitative results on VOC2010 [12] with 59 categories.

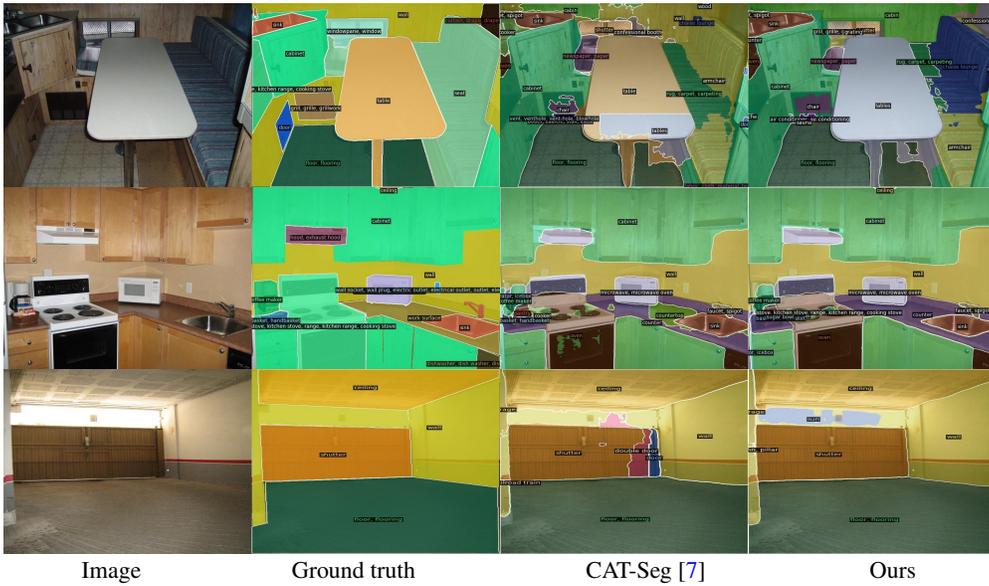


Figure 4: Comparison of qualitative results on ADE20K [55] with 847 categories.

| Method    | POF | DCRC | Param. (M) | A-847       | PC-459      | A-150       | PC-59       | PAS-20      | PAS-20 <sup>b</sup> |
|-----------|-----|------|------------|-------------|-------------|-------------|-------------|-------------|---------------------|
| Freeze    | ✗   | ✗    | 0          | 4.4         | 6.6         | 24.8        | 49.4        | 92.5        | 71.9                |
| LoRA [16] | ✗   | ✗    | 7.5        | 11.4        | 17.6        | 28.6        | 55.1        | 94.2        | 76.7                |
| H-CLIP    | ✓   | ✗    | 5.62       | 12.3        | 19.0        | 31.6        | 56.4        | 94.6        | 76.3                |
|           | ✗   | ✓    | 0.01       | 7.6         | 10.9        | 26.8        | 53.6        | 92.7        | 74.5                |
|           | ✓   | ✓    | 5.63       | <b>12.4</b> | <b>19.3</b> | <b>32.4</b> | <b>57.9</b> | <b>95.2</b> | <b>78.2</b>         |

Table 3: **Ablation study on the components of H-CLIP.** “LoRA”: a mainstream parameter-efficient tuning method with a comparable number of parameters for comparison. “POF”: Partial Orthogonal Fine-tuning. “DCRC”: Dual Cross Relation Communication. The base model is ViT-B/16.

239 **Efficiency comparison.** We compare the efficiency of our method with other approaches, including  
 240 OVSeg [30], CAT-Seg [7], and SAN [50], all of which utilize CLIP ViT models. The comparison,

|     | Block dimension $q$   | Param. (M) | A-847       | PC-459      | A-150       | PC-59       | PAS-20      | PAS-20 <sup>b</sup> |
|-----|-----------------------|------------|-------------|-------------|-------------|-------------|-------------|---------------------|
| (a) | $256 \times 256$      | 22.52      | 12.4        | 19.2        | <b>32.7</b> | 57.6        | <b>95.4</b> | 77.9                |
|     | $128 \times 128$      | 5.63       | <b>12.4</b> | <b>19.3</b> | 32.4        | <b>57.9</b> | 95.2        | <b>78.2</b>         |
|     | $64 \times 64$        | 1.41       | 11.7        | 18.4        | 31.7        | 56.9        | 95.0        | 76.4                |
|     | Orthogonal Constraint | Param. (M) | A-847       | PC-459      | A-150       | PC-59       | PAS-20      | PAS-20 <sup>b</sup> |
| (b) | w/o                   | 7.51       | 11.9        | 18.5        | 32.2        | 57.5        | <b>95.3</b> | 76.9                |
|     | with                  | 3.76       | 12.2        | 19.1        | 31.4        | 57.1        | 94.3        | 76.8                |
|     | POF                   | 5.63       | <b>12.4</b> | <b>19.3</b> | <b>32.4</b> | <b>57.9</b> | 95.2        | <b>78.2</b>         |

Table 4: **Ablation study on different designs in POF.** We show the impact of (a) different block dimensions  $q$  and (b) orthogonal constraints. The base model is ViT-B/16.

241 summarized in Table 2, shows that our method employs the fewest trainable parameters while  
242 balancing the generalization of the pre-trained model and the flexibility for learning new knowledge.  
243 Additionally, since we introduce a lightweight architecture for calculating relations, specifically two  
244 relation matrices, the inference overhead is negligible during the inference phase.

### 245 5.3 Ablative Studies

246 **Ablation of Main Components.** Here, we conduct an ablation study to demonstrate the benefits  
247 of each component of our proposed H-CLIP: partial orthogonal fine-tuning (POF) and dual cross-  
248 relation communication (DCRC). We use the ViT-B/16 [9] version of CLIP as the baseline, shown in  
249 row 1 of Table 3. Additionally, we implement a mainstream parameter-efficient fine-tuning (PEFT)  
250 method, LoRA [16], for comparison with a similar number of learnable parameters, as shown in row  
251 2. Note that LoRA can improve performance compared to the baseline, demonstrating that PEFT is a  
252 viable approach for this task. Then, comparing row 5 to row 2, we observe significant performance  
253 gains, indicating that our results are driven by our targeted solution rather than merely the number of  
254 parameters. Moreover, row 3 shows that using only POF preserves generalization on unseen classes,  
255 particularly in the A-847 dataset. Meanwhile, solely adapting DCRC shows limited improvement, as  
256 it only enhances communication among frozen weight matrices. Finally, integrating DCRC with POF  
257 yields clear performance gains, e.g., a 12.6% improvement on the PC-459 dataset.

258 **Different Design of POF.** Table 4 presents experiments introducing different designs into POF. The  
259 design of POF is related to (1) block dimension, i.e.,  $q$ , and (2) how orthogonality constraints are  
260 applied. In (a), the results show that larger  $q$  generally performs better than smaller  $q$ . However, we find a good trade-off between performance and parameter  
261 efficiency, with  $q = 128$  working well across datasets and tasks. Therefore, we maintain this setting  
262 in other experiments. In (b), we show that both blindly applying orthogonality constraints to the  
263 learnable matrices of all layers and not using any constraints at all can degrade performance on most  
264 test sets, demonstrating the value of our analysis with the hyperspherical energy principle.  
265

## 266 6 Conclusion

267 In this paper, we propose a H-CLIP framework to address three issues: 1) high computational cost, 2)  
268 misalignment between the two inherent modalities of CLIP, and 3) degraded generalization ability  
269 on unseen categories when equipping CLIP with pixel-level prediction ability for open-vocabulary  
270 semantic segmentation. Specifically, we propose a symmetrical parameter-efficient fine-tuning (PEFT)  
271 strategy conducted in hyperspherical space for both of the two CLIP modalities. Specifically, the  
272 PEFT strategy is achieved by a series of efficient block-diagonal learnable transformation matrices and  
273 a dual cross-relation communication module among all learnable matrices to mitigate misalignment  
274 between different modalities. Furthermore, we apply an additional constraint to PEFT on the CLIP  
275 text encoder according to the hyperspherical energy principle, i.e., minimizing hyperspherical energy  
276 during fine-tuning preserves the intrinsic structure of the original parameter space, to prevent the  
277 destruction of the generalization ability offered by the CLIP text encoder. Extensive experiments  
278 demonstrate that the proposed H-CLIP framework generalized improves segmentation performance  
279 across several benchmarks while introducing approximately 4% of CLIP’s total parameters. We hope  
280 our approach will provide a new direction and inspire future research in this field.

## References

- [1] Maxime Bucher, Tuan-Hung Vu, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. *Advances in Neural Information Processing Systems*, 32, 2019.
- [2] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018.
- [3] Arthur Cayley. Sur quelques propriétés des déterminants gauches. 1846.
- [4] Tianrun Chen, Lanyun Zhu, Chaotao Ding, Runlong Cao, Shangzhan Zhang, Yan Wang, Zejian Li, Lingyun Sun, Papa Mao, and Ying Zang. Sam fails to segment anything? – sam-adapter: Adapting sam in underperformed scenes: Camouflage, shadow, and more, 2023.
- [5] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *European conference on computer vision*, pages 104–120. Springer, 2020.
- [6] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1290–1299, 2022.
- [7] Seokju Cho, Heeseong Shin, Sunghwan Hong, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation, 2024.
- [8] Jian Ding, Nan Xue, Gui-Song Xia, and Dengxin Dai. Decoupling zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11583–11592, 2022.
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [10] Zi-Yi Dou, Aishwarya Kamath, Zhe Gan, Pengchuan Zhang, Jianfeng Wang, Linjie Li, Zicheng Liu, Ce Liu, Yann LeCun, Nanyun Peng, et al. Coarse-to-fine vision-language pre-training with fusion in the backbone. *Advances in neural information processing systems*, 35:32942–32956, 2022.
- [11] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, et al. An empirical study of training end-to-end vision-and-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18166–18176, 2022.
- [12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010.
- [13] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. *Advances in Neural Information Processing Systems*, 33:6616–6628, 2020.
- [14] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *International Journal of Computer Vision*, 132(2):581–595, 2024.
- [15] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *European Conference on Computer Vision*, pages 540–557. Springer, 2022.
- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.

- 329 [17] Zhicheng Huang, Zhaoyang Zeng, Yupan Huang, Bei Liu, Dongmei Fu, and Jianlong Fu.  
330 Seeing out of the box: End-to-end pre-training for vision-language representation learning. In  
331 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
332 12976–12985, 2021.
- 333 [18] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas  
334 Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings*  
335 *of the IEEE/CVF International Conference on Computer Vision*, pages 1780–1790, 2021.
- 336 [19] Eric Kernfeld, Misha Kilmer, and Shuchin Aeron. Tensor–tensor products with invertible linear  
337 transforms. *Linear Algebra and its Applications*, 485:545–570, 2015.
- 338 [20] Misha E Kilmer and Carla D Martin. Factorization strategies for third-order tensors. *Linear*  
339 *Algebra and its Applications*, 435(3):641–658, 2011.
- 340 [21] Wonjae Kim, Bokyoung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without  
341 convolution or region supervision. In *International conference on machine learning*, pages  
342 5583–5594. PMLR, 2021.
- 343 [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
344 *arXiv:1412.6980*, 2014.
- 345 [23] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,  
346 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In  
347 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026,  
348 2023.
- 349 [24] Jungbeom Lee, Jooyoung Choi, Jisoo Mok, and Sungroh Yoon. Reducing information bottleneck  
350 for weakly supervised semantic segmentation. *Advances in Neural Information Processing*  
351 *Systems*, 34:27408–27421, 2021.
- 352 [25] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-  
353 driven semantic segmentation. In *International Conference on Learning Representations*,  
354 2022.
- 355 [26] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-  
356 image pre-training for unified vision-language understanding and generation. In *International*  
357 *conference on machine learning*, pages 12888–12900. PMLR, 2022.
- 358 [27] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven  
359 Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum  
360 distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- 361 [28] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu  
362 Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image  
363 pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
364 *Recognition*, pages 10965–10975, 2022.
- 365 [29] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang,  
366 Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for  
367 vision-language tasks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow,*  
368 *UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 121–137. Springer, 2020.
- 369 [30] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yanan Zhao, Hang Zhang, Peizhao Zhang,  
370 Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted  
371 clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
372 pages 7061–7070, 2023.
- 373 [31] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*  
374 *in neural information processing systems*, 36, 2024.
- 375 [32] Weiyang Liu, Rongmei Lin, Zhen Liu, Lixin Liu, Zhiding Yu, Bo Dai, and Le Song. Learning  
376 towards minimum hyperspherical energy. *Advances in neural information processing systems*,  
377 31, 2018.

- 378 [33] Canyi Lu, Xi Peng, and Yunchao Wei. Low-rank tensor completion with a new tensor nuclear  
379 norm induced by invertible linear transforms. In *Proceedings of the IEEE/CVF conference on*  
380 *computer vision and pattern recognition*, pages 5996–6004, 2019.
- 381 [34] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic  
382 visiolinguistic representations for vision-and-language tasks. *Advances in neural information*  
383 *processing systems*, 32, 2019.
- 384 [35] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Wan Lee, Sanja Fidler,  
385 Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic  
386 segmentation in the wild. In *Proceedings of the IEEE conference on computer vision and*  
387 *pattern recognition*, pages 891–898, 2014.
- 388 [36] Zelin Peng, Zhengqin Xu, Zhilin Zeng, Lingxi Xie, Qi Tian, and Wei Shen. Parameter efficient  
389 fine-tuning via cross block orchestration for segment anything model. In *IEEE/CVF Conference*  
390 *on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- 391 [37] Zelin Peng, Zhengqin Xu, Zhilin Zeng, Xiaokang Yang, and Wei Shen. Sam-parser: Fine-tuning  
392 sam efficiently by parameter space reconstruction. In *Proceedings of the AAAI Conference on*  
393 *Artificial Intelligence*, volume 38, pages 4515–4523, 2024.
- 394 [38] Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian  
395 Weller, and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning.  
396 *Advances in Neural Information Processing Systems*, 36:79320–79362, 2023.
- 397 [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
398 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
399 models from natural language supervision. In *International conference on machine learning*,  
400 pages 8748–8763. PMLR, 2021.
- 401 [40] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan,  
402 Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models  
403 for code. *arXiv preprint arXiv:2308.12950*, 2023.
- 404 [41] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba,  
405 Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment  
406 model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-*  
407 *tion*, pages 15638–15650, 2022.
- 408 [42] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. VI-bert:  
409 Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*,  
410 2019.
- 411 [43] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from  
412 transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- 413 [44] JJ Thomson. On the structure of the atom: an investigation of the stability and periods of  
414 oscilation of a number of corpuscles arranged at equal intervals around the circumference of  
415 a circle; with application of the results to the theory atomic structure. *Philos. Mag. Series 6*,  
416 7(39):237, 1904.
- 417 [45] Jianfeng Wang, Xiaowei Hu, Zhe Gan, Zhengyuan Yang, Xiyang Dai, Zicheng Liu, Yumao  
418 Lu, and Lijuan Wang. Ufo: A unified transformer for vision-language representation learning.  
419 *arXiv preprint arXiv:2111.10023*, 2021.
- 420 [46] Zifeng Wang, Xi Chen, Rui Wen, Shao-Lun Huang, Ercan Kuruoglu, and Yefeng Zheng.  
421 Information theoretic counterfactual learning from missing-not-at-random feedback. *Advances*  
422 *in Neural Information Processing Systems*, 33:1854–1864, 2020.
- 423 [47] Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Pmc-llama: Further  
424 finetuning llama on medical papers. *arXiv preprint arXiv:2304.14454*, 2023.
- 425 [48] Bin Xie, Jiale Cao, Jin Xie, Fahad Shahbaz Khan, and Yanwei Pang. Sed: A simple encoder-  
426 decoder for open-vocabulary semantic segmentation, 2023.

- 427 [49] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello.  
428 Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *Proceedings*  
429 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2955–2966,  
430 2023.
- 431 [50] Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xiang Bai. Side adapter network  
432 for open-vocabulary semantic segmentation. In *Proceedings of the IEEE/CVF Conference on*  
433 *Computer Vision and Pattern Recognition*, pages 2945–2954, 2023.
- 434 [51] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A  
435 simple baseline for open-vocabulary semantic segmentation with pre-trained vision-language  
436 model. In *European Conference on Computer Vision*, pages 736–753. Springer, 2022.
- 437 [52] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die  
438 hard: Open-vocabulary segmentation with single frozen convolutional clip. In *NeurIPS*, 2023.
- 439 [53] Haotian Zhang, Pengchuan Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Li, Xiyang Dai,  
440 Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. Glipv2: Unifying localization  
441 and vision-language understanding. *Advances in Neural Information Processing Systems*,  
442 35:36067–36080, 2022.
- 443 [54] Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan  
444 Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with  
445 zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023.
- 446 [55] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio  
447 Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal*  
448 *of Computer Vision*, 127:302–321, 2019.
- 449 [56] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *European*  
450 *Conference on Computer Vision*, pages 696–712. Springer, 2022.
- 451 [57] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learn-  
452 ing for vision-language models. In *IEEE/CVF Conference on Computer Vision and Pattern*  
453 *Recognition (CVPR)*, 2022.
- 454 [58] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for  
455 vision-language models. *International Journal of Computer Vision (IJCV)*, 2022.
- 456 [59] Ziqin Zhou, Yinjie Lei, Bowen Zhang, Lingqiao Liu, and Yifan Liu. Zegclip: Towards  
457 adapting clip for zero-shot semantic segmentation. *Proceedings of the IEEE/CVF Conference*  
458 *on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 459 [60] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: En-  
460 hancing vision-language understanding with advanced large language models. *arXiv preprint*  
461 *arXiv:2304.10592*, 2023.
- 462 [61] Xiangyang Zhu, Renrui Zhang, Bowei He, Aojun Zhou, Dong Wang, Bin Zhao, and Peng Gao.  
463 Not all features matter: Enhancing few-shot clip with adaptive prior refinement. In *Proceedings*  
464 *of the IEEE/CVF International Conference on Computer Vision*, pages 2605–2615, 2023.

## Appendix of H-CLIP

### 466 A Derivation of the Definition

467 In this section, we provide derivations of definitions in the main paper.

468 **Definition 4.1(3-order T-product)** For  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$  and  $\mathcal{B} \in \mathbb{R}^{n_2 \times l \times n_3}$ , the 3-order T-product  
469  $\mathcal{C} \in \mathbb{R}^{n_1 \times l \times n_3} = \mathcal{A} * \mathcal{B}$  is defined as:

$$\mathcal{C} = \mathcal{A} * \mathcal{B} = \text{fold}(\text{circ}(\mathcal{A}) \cdot \text{unfold}(\mathcal{B})), \quad (18)$$

470 where “ $\cdot$ ” represents standard matrix product.

471 **Definition 4.2(Higher-order T-product)** For  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3 \cdots \times n_p}$  and  $\mathcal{B} \in \mathbb{R}^{n_2 \times l \times n_3 \cdots \times n_p}$ , the  
472 High-order T-product  $\mathcal{C} \in \mathbb{R}^{n_1 \times l \times n_3 \cdots \times n_p} = \mathcal{A} * \mathcal{B}$  is defined as:

$$\mathcal{C} = \mathcal{A} * \mathcal{B} = \text{fold}(\text{circ}(\mathcal{A}) * \text{unfold}(\mathcal{B})). \quad (19)$$

473 *Derivation.* According to [20], if  $\mathcal{A}$  is  $n_1 \times n_2 \times n_3$ ,  $\mathcal{A}$  can be block diagonalized by using Discrete  
474 Fourier Transformer (DFT) matrix  $\mathbf{F}_{n_3} \in \mathbb{R}^{n_3 \times n_3}$  as:

$$(\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1}) \cdot \text{circ}(\text{unfold}(\mathcal{A})) \cdot (\mathbf{F}_{n_3}^* \otimes \mathbf{I}_{n_2}) = \mathbf{D} = \begin{bmatrix} \mathbf{D}_1 & & \\ & \ddots & \\ & & \mathbf{D}_{n_3} \end{bmatrix} \in \mathbb{R}^{n_1 n_3 \times n_2 n_3}, \quad (20)$$

475 where “ $\otimes$ ” denotes the Kernecker product, “ $\mathbf{F}_{n_3}^*$ ” denotes the conjugate transpose of  $\mathbf{F}_{n_3}$ , “ $\cdot$ ” means  
476 standard matrix product and  $\mathbf{D}$  is a block diagonal matrix. In fact, the  $i$ -th block matrix  $\mathbf{D}_i$  of  $\mathbf{D}$  can  
477 be computed by applying DFT of  $\mathcal{A}$  along 3-rd dimension. The **3-order T-product** in Eq. (18) can  
478 be computed as:

$$(\mathbf{F}_{n_3}^* \otimes \mathbf{I}_{n_1}) \cdot ((\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1}) \cdot \text{circ}(\text{unfold}(\mathcal{A})) \cdot (\mathbf{F}_{n_3}^* \otimes \mathbf{I}_{n_2})) \cdot (\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_2}) \cdot \text{unfold}(\mathcal{B}). \quad (21)$$

479 It is readily shown that  $(\mathbf{F}_{n_3} \otimes \mathbf{I}_{n_2}) \text{unfold}$  can be computed by applying DFT of  $\mathcal{B}$  along 3-rd  
480 dimension: the result called  $\bar{\mathcal{B}}$ . Thus, Eq. (21) remains to multiply each block matrix  $\mathbf{D}_i$  of  $\mathbf{D}$  with  
481 each block matrix  $\mathbf{B}_i$  of  $\bar{\mathcal{B}}$ , then take an inverse DFT along the 3-rd dimension of the result. Hence,  
482 the **3-order T-product** in Eq. (18) can be re-formulated as:

$$\mathcal{C} = \text{DFT}_3^{-1}(\text{DFT}_3(\mathcal{A}) \odot \text{DFT}_3(\mathcal{B})) = \text{DFT}_3^{-1}(\bar{\mathcal{A}} \odot \bar{\mathcal{B}}) = \text{DFT}_3^{-1}(\bar{\mathcal{C}}), \quad (22)$$

483 where  $\text{DFT}_3(\cdot)$  is DFT along 3-rd dimension and  $\text{DFT}_3^{-1}(\cdot)$  is the inverse DFT along 3-rd dimension.  
484 In mathematics, the DFT of  $\mathcal{A}$  along 3-rd dimension is formulated as:

$$\bar{\mathcal{A}} = \text{DFT}_3(\mathcal{A}) = \mathcal{A} \times_3 \mathbf{F}_{n_3}. \quad (23)$$

485 Similarly, the inverse DFT of  $\bar{\mathcal{A}}$  along 3-rd dimension is derived as:

$$\mathcal{A} = \text{DFT}_3^{-1}(\bar{\mathcal{A}}) = \bar{\mathcal{A}} \times_3 \mathbf{F}_{n_3}^{-1}. \quad (24)$$

486 By the detailed theoretical analysis in [33], the DFT has been extended to a general invertible  
487 transform  $S$  with an invertible transform matrix  $\mathbf{S}$ . In mathematics, the invertible transform of  $\mathcal{A}$   
488 along 3-rd dimension is formulated as:

$$\bar{\mathcal{A}} = S_3(\mathcal{A}) = \mathcal{A} \times_3 \mathbf{S}_{n_3}. \quad (25)$$

489 Similarly, the inverse transform of  $\bar{\mathcal{A}}$  along 3-rd dimension is derived as:

$$\mathcal{A} = S_3^{-1}(\bar{\mathcal{A}}) = \bar{\mathcal{A}} \times_3 \mathbf{S}_{n_3}^{-1}. \quad (26)$$

490 Similarly, if  $\mathcal{A} \in \mathbb{R}^{n_1 \times n_2 \times \cdots \times n_p}$ ,  $\mathcal{A}$  can be block diagonalized by using a sequence of DFT matrices  
491  $\mathbf{F}_{n_i} \in \mathbb{R}^{n_i \times n_i}$ ,  $i = 3, 4, \dots, p$  as:

$$(\mathbf{F}_{n_p} \otimes \mathbf{F}_{n_{p-1}} \otimes \cdots \otimes \mathbf{F}_{n_3} \otimes \mathbf{I}_{n_1}) \cdot \tilde{\mathcal{A}} \cdot (\mathbf{F}_{n_p}^* \otimes \mathbf{F}_{n_{p-1}}^* \otimes \cdots \otimes \mathbf{F}_{n_3}^* \otimes \mathbf{I}_{n_2}) = \mathbf{D}, \quad (27)$$

492 where  $\tilde{\mathcal{A}} = \text{circ}(\text{unfold}(\mathcal{A})) \in \mathbb{R}^{n_1 n_3 n_4 \cdots n_p \times n_2 n_3 \cdots n_p}$ . Since the matrix  $\mathbf{D}$  is block diagonal with  
 493  $n_3 n_4 \cdots n_p$  blocks each of size  $n_1 \times n_2$ , the **Higher-order T-product** in Eq. (19) can be computed  
 494 as:

$$(\tilde{\mathbf{F}}^* \otimes \mathbf{I}_{n_1}) \cdot ((\tilde{\mathbf{F}} \otimes \mathbf{I}_{n_1}) \cdot \tilde{\mathcal{A}} \cdot (\tilde{\mathbf{F}}^* \otimes \mathbf{I}_{n_2})) \cdot (\tilde{\mathbf{F}}_{n_3} \otimes \mathbf{I}_{n_2}) \cdot \tilde{\mathcal{B}}, \quad (28)$$

495 where  $\tilde{\mathbf{F}} = \mathbf{F}_{n_p} \otimes \mathbf{F}_{n_{p-1}} \otimes \cdots \otimes \mathbf{F}_{n_3}$ . Using the DEF, it is straightforward to show that the block  
 496 diagonal matrix  $\mathbf{D}$  in Eq. (27) can be obtained by repeated DFTs of  $\mathcal{A}$  along each dimension expect  
 497 for 1-st and 2-nd dimension. Similarly, by using a sequence invertible transform  $S_j(\cdot), i = 3, 4, \dots, p$   
 498 with invertible transform matrix  $\mathbf{S}_i$ , the **Higher-order T-product** in Eq. (19) can be re-formulated as:

$$\mathcal{C} = \tilde{S}^{-1}(\tilde{S}(\mathcal{A}) \odot \tilde{S}(\mathcal{B})) = \tilde{S}^{-1}(\bar{\mathcal{A}} \odot \bar{\mathcal{B}}) = \tilde{S}^{-1}(\bar{\mathcal{C}}), \quad (29)$$

499 where  $\tilde{S}(\mathcal{A}) = S_p(S_{p-1}(\cdots S_3(\mathcal{A}) \cdots))$ ,  $\bar{\mathcal{C}} = \bar{\mathcal{A}} \odot \bar{\mathcal{B}}$  denotes the frontal-slice-wise product  
 500  $\bar{\mathcal{C}}(;, ;, i) = \bar{\mathcal{A}}(;, ;, i) \cdot \bar{\mathcal{B}}(;, ;, i), i = 1, 2, \dots, n_3 n_4 \cdots n_p$  and  $\tilde{S}^{-1}(\cdot)$  is the inverse transform  
 501 of  $\tilde{S}(\cdot)$ . The inverse transform  $\tilde{S}(\cdot)$  is formulated as:

$$\bar{\mathcal{A}} = \tilde{S}(\mathcal{A}) = \mathcal{A} \times_3 \mathbf{S}_3 \times_4 \mathbf{S}_4 \cdots \times_p \mathbf{S}_p, \quad (30)$$

502 and its inverse transform is derived as:

$$\mathcal{A} = \tilde{S}^{-1}(\bar{\mathcal{A}}) = \bar{\mathcal{A}} \times_3 \mathbf{S}_3^{-1} \times_4 \mathbf{S}_4^{-1} \cdots \times_p \mathbf{S}_p^{-1}. \quad (31)$$

503

■

## 504 **NeurIPS Paper Checklist**

### 505 **1. Claims**

506 Question: Do the main claims made in the abstract and introduction accurately reflect the  
507 paper’s contributions and scope?

508 Answer: [\[Yes\]](#)

509 Justification: The abstract clearly states the following claims about the paper’s contribu-  
510 tions and scope. We aim to address a significant challenge in open-vocabulary semantic  
511 segmentation: fine-tuning CLIP to achieve per-pixel predictions without compromising  
512 its generalization capabilities. To this end, we propose a novel H-CLIP, which introduces  
513 a partial orthogonal fine-tuning strategy that prevents a drop in hyperspherical energy,  
514 thereby preserving generalization. Subsequently, H-CLIP employs dual cross-relation com-  
515 munication to increase projection flexibility, facilitating the acquisition of segmentation  
516 knowledge.

### 517 **2. Limitations**

518 Question: Does the paper discuss the limitations of the work performed by the authors?

519 Answer: [\[NA\]](#)

520 Justification: Those are not discussed in the paper.

521 Guidelines:

- 522 • The answer NA means that the paper has no limitation while the answer No means that  
523 the paper has limitations, but those are not discussed in the paper.
- 524 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 525 • The paper should point out any strong assumptions and how robust the results are to  
526 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
527 model well-specification, asymptotic approximations only holding locally). The authors  
528 should reflect on how these assumptions might be violated in practice and what the  
529 implications would be.
- 530 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
531 only tested on a few datasets or with a few runs. In general, empirical results often  
532 depend on implicit assumptions, which should be articulated.
- 533 • The authors should reflect on the factors that influence the performance of the approach.  
534 For example, a facial recognition algorithm may perform poorly when image resolution  
535 is low or images are taken in low lighting. Or a speech-to-text system might not be  
536 used reliably to provide closed captions for online lectures because it fails to handle  
537 technical jargon.
- 538 • The authors should discuss the computational efficiency of the proposed algorithms  
539 and how they scale with dataset size.
- 540 • If applicable, the authors should discuss possible limitations of their approach to  
541 address problems of privacy and fairness.
- 542 • While the authors might fear that complete honesty about limitations might be used by  
543 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
544 limitations that aren’t acknowledged in the paper. The authors should use their best  
545 judgment and recognize that individual actions in favor of transparency play an impor-  
546 tant role in developing norms that preserve the integrity of the community. Reviewers  
547 will be specifically instructed to not penalize honesty concerning limitations.

### 548 **3. Theory Assumptions and Proofs**

549 Question: For each theoretical result, does the paper provide the full set of assumptions and  
550 a complete (and correct) proof?

551 Answer: [\[Yes\]](#)

552 Justification: All assumptions are proved clearly in the main paper or the supplemental  
553 material.

554 Guidelines:

- 555 • The answer NA means that the paper does not include theoretical results.

- 556 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
557 referenced.
- 558 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 559 • The proofs can either appear in the main paper or the supplemental material, but if  
560 they appear in the supplemental material, the authors are encouraged to provide a short  
561 proof sketch to provide intuition.
- 562 • Inversely, any informal proof provided in the core of the paper should be complemented  
563 by formal proofs provided in appendix or supplemental material.
- 564 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 565 4. Experimental Result Reproducibility

566 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
567 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
568 of the paper (regardless of whether the code and data are provided or not)?

569 Answer: [Yes]

570 Justification: The detailed experimental settings and information are provided in Sec. 5, and  
571 this information is sufficient to reproduce the main experimental results.

572 Guidelines:

- 573 • The answer NA means that the paper does not include experiments.
- 574 • If the paper includes experiments, a No answer to this question will not be perceived  
575 well by the reviewers: Making the paper reproducible is important, regardless of  
576 whether the code and data are provided or not.
- 577 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
578 to make their results reproducible or verifiable.
- 579 • Depending on the contribution, reproducibility can be accomplished in various ways.  
580 For example, if the contribution is a novel architecture, describing the architecture fully  
581 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
582 be necessary to either make it possible for others to replicate the model with the same  
583 dataset, or provide access to the model. In general, releasing code and data is often  
584 one good way to accomplish this, but reproducibility can also be provided via detailed  
585 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
586 of a large language model), releasing of a model checkpoint, or other means that are  
587 appropriate to the research performed.
- 588 • While NeurIPS does not require releasing code, the conference does require all submis-  
589 sions to provide some reasonable avenue for reproducibility, which may depend on the  
590 nature of the contribution. For example
  - 591 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
592 to reproduce that algorithm.
  - 593 (b) If the contribution is primarily a new model architecture, the paper should describe  
594 the architecture clearly and fully.
  - 595 (c) If the contribution is a new model (e.g., a large language model), then there should  
596 either be a way to access this model for reproducing the results or a way to reproduce  
597 the model (e.g., with an open-source dataset or instructions for how to construct  
598 the dataset).
  - 599 (d) We recognize that reproducibility may be tricky in some cases, in which case  
600 authors are welcome to describe the particular way they provide for reproducibility.  
601 In the case of closed-source models, it may be that access to the model is limited in  
602 some way (e.g., to registered users), but it should be possible for other researchers  
603 to have some path to reproducing or verifying the results.

#### 604 5. Open access to data and code

605 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
606 tions to faithfully reproduce the main experimental results, as described in supplemental  
607 material?

608 Answer: [No]

609 Justification: The code will be made public after acceptance.

610 Guidelines:

- 611 • The answer NA means that paper does not include experiments requiring code.
- 612 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 613
- 614 • While we encourage the release of code and data, we understand that this might not be
- 615 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 616 including code, unless this is central to the contribution (e.g., for a new open-source
- 617 benchmark).
- 618 • The instructions should contain the exact command and environment needed to run to
- 619 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 620
- 621 • The authors should provide instructions on data access and preparation, including how
- 622 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 623 • The authors should provide scripts to reproduce all experimental results for the new
- 624 proposed method and baselines. If only a subset of experiments are reproducible, they
- 625 should state which ones are omitted from the script and why.
- 626 • At submission time, to preserve anonymity, the authors should release anonymized
- 627 versions (if applicable).
- 628 • Providing as much information as possible in supplemental material (appended to the
- 629 paper) is recommended, but including URLs to data and code is permitted.

630 **6. Experimental Setting/Details**

631 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-

632 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the

633 results?

634 Answer: [Yes]

635 Justification: All the experimental settings are clarified in Sec. 5. (and Appendix).

636 Guidelines:

- 637 • The answer NA means that the paper does not include experiments.
- 638 • The experimental setting should be presented in the core of the paper to a level of detail
- 639 that is necessary to appreciate the results and make sense of them.
- 640 • The full details can be provided either with the code, in appendix, or as supplemental
- 641 material.

642 **7. Experiment Statistical Significance**

643 Question: Does the paper report error bars suitably and correctly defined or other appropriate

644 information about the statistical significance of the experiments?

645 Answer: [No]

646 Justification: error bars are not reported because it would be too computationally expensive.

647 Guidelines:

- 648 • The answer NA means that the paper does not include experiments.
- 649 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 650 dence intervals, or statistical significance tests, at least for the experiments that support
- 651 the main claims of the paper.
- 652 • The factors of variability that the error bars are capturing should be clearly stated (for
- 653 example, train/test split, initialization, random drawing of some parameter, or overall
- 654 run with given experimental conditions).
- 655 • The method for calculating the error bars should be explained (closed form formula,
- 656 call to a library function, bootstrap, etc.)
- 657 • The assumptions made should be given (e.g., Normally distributed errors).
- 658 • It should be clear whether the error bar is the standard deviation or the standard error
- 659 of the mean.
- 660 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 661 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 662 of Normality of errors is not verified.

- 663
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
  - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
- 664
- 665
- 666
- 667

## 668 8. Experiments Compute Resources

669 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

670

671

672 Answer: [Yes]

673 Justification: The computing requirements are provided in Sec. 5.

- The answer NA means that the paper does not include experiments.
  - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
  - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
  - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681

## 682 9. Code Of Ethics

683 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

684

685 Answer: [Yes]

686 Justification: We ensure our research adheres to the guidelines.

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
  - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
  - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 687
- 688
- 689
- 690
- 691

## 692 10. Broader Impacts

693 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

694

695 Answer: [NA]

696 Justification: there is no societal impact of the work performed.

697 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
  - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715

- 716
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 717  
718  
719

## 720 11. Safeguards

721 Question: Does the paper describe safeguards that have been put in place for responsible  
722 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
723 image generators, or scraped datasets)?

724 Answer: [NA]

725 Justification: the paper poses no such risks.

726 Guidelines:

- 727 • The answer NA means that the paper poses no such risks.
- 728 • Released models that have a high risk for misuse or dual-use should be released with  
729 necessary safeguards to allow for controlled use of the model, for example by requiring  
730 that users adhere to usage guidelines or restrictions to access the model or implementing  
731 safety filters.
- 732 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
733 should describe how they avoided releasing unsafe images.
- 734 • We recognize that providing effective safeguards is challenging, and many papers do  
735 not require this, but we encourage authors to take this into account and make a best  
736 faith effort.

## 737 12. Licenses for existing assets

738 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
739 the paper, properly credited and are the license and terms of use explicitly mentioned and  
740 properly respected?

741 Answer: [NA]

742 Justification: All the assets used are properly credited and are the license and terms of use  
743 explicitly mentioned and properly respected.

- 744 • The answer NA means that the paper does not use existing assets.
- 745 • The authors should cite the original paper that produced the code package or dataset.
- 746 • The authors should state which version of the asset is used and, if possible, include a  
747 URL.
- 748 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 749 • For scraped data from a particular source (e.g., website), the copyright and terms of  
750 service of that source should be provided.
- 751 • If assets are released, the license, copyright information, and terms of use in the  
752 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
753 has curated licenses for some datasets. Their licensing guide can help determine the  
754 license of a dataset.
- 755 • For existing datasets that are re-packaged, both the original license and the license of  
756 the derived asset (if it has changed) should be provided.
- 757 • If this information is not available online, the authors are encouraged to reach out to  
758 the asset's creators.

## 759 13. New Assets

760 Question: Are new assets introduced in the paper well documented and is the documentation  
761 provided alongside the assets?

762 Answer: [NA]

763 Justification: The paper does not release new assets.

- 764 • The answer NA means that the paper does not release new assets.
- 765 • Researchers should communicate the details of the dataset/code/model as part of their  
766 submissions via structured templates. This includes details about training, license,  
767 limitations, etc.

- 768
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- 769
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 770
- 771

772 **14. Crowdsourcing and Research with Human Subjects**

773 Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

774

775

776 Answer: [NA]

777 Justification: This paper does not involve crowdsourcing nor research with human subjects.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 778
- 779
- 780
- 781
- 782
- 783
- 784
- 785

786 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

787

788 Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

789

790

791

792 Answer: [NA]

793 Justification: This paper does not involve crowdsourcing nor research with human subjects.

794 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
  - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
  - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 795
- 796
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804