
Class Concept Representation from Contextual Texts for Training-Free Multi-Label Recognition

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The power of large vision-language models (VLMs) has been demonstrated for
2 diverse vision tasks including multi-label recognition with training-free approach or
3 prompt tuning by measuring the cosine similarity between the text features related
4 to class names and the visual features of images. Prior works usually formed the
5 class-related text features by averaging simple hand-crafted text prompts with *class*
6 *names* (e.g., “a photo of {class name}”). However, they may not fully exploit the
7 capability of VLMs considering how humans form the concepts on words using rich
8 contexts with the patterns of co-occurrence with other words. Inspired by that, we
9 propose *class concept* representation for zero-shot multi-label recognition to better
10 exploit rich contexts in the massive descriptions on images (e.g., captions from MS-
11 COCO) using large VLMs. Then, for better aligning visual features of VLMs to our
12 class concept representation, we propose context-guided visual representation that
13 is in the same linear space as class concept representation. Experimental results
14 on diverse benchmarks show that our proposed methods substantially improved
15 the performance of zero-shot methods like Zero-Shot CLIP and yielded better
16 performance than zero-shot prompt tunings that require additional training like
17 TaI-DPT. In addition, our proposed methods can *synergetically* work with existing
18 prompt tuning methods, consistently improving the performance of DualCoOp and
19 TaI-DPT in a training-free manner with negligible increase in inference time.

20 1 Introduction

21 The goal of multi-label image recognition is to assign all semantic labels (or class names) within an
22 image [10, 44, 48, 11, 27, 33, 31]. Differing from single-label recognition, multi-label recognition
23 addresses a broader range of practical applications such as image retrieval [36, 39], recommendation
24 systems [52, 8], medical diagnosis recognition [43] and retail checkout recognition [17, 45]. However,
25 one of the challenges in multi-label recognition is the difficulty of collecting full label annotations,
26 which is laborious and prone to missing. To alleviate it, recent works have investigated training with
27 incomplete labels such as partial labels [37, 6, 31, 15, 9] or a single positive label [13, 46].

28 Recent advances of large vision-language models (VLMs) [32, 2, 22, 25, 47, 49] has demon-
29 strated their strong transferability on various downstream tasks with great performance. Contrastive
30 Language-Image Pretraining (CLIP) achieved impressive performance in zero-shot classification by
31 measuring the cosine similarity between images and class-related hand-crafted text prompts [32].
32 Fine-tuning VLMs for adapting desired downstream datasets [32] can further improve performance
33 for targeted tasks, but tuning millions of parameters is usually undesirable due to computation burden
34 and possible forgetting. Prompt tuning has been investigated as an efficient and low-cost training
35 paradigm [54, 53], learning only a few context tokens of VLMs for a given task. In multi-label
36 recognition, prompt tuning with CLIP has been investigated for distinguishing multiple objects in an

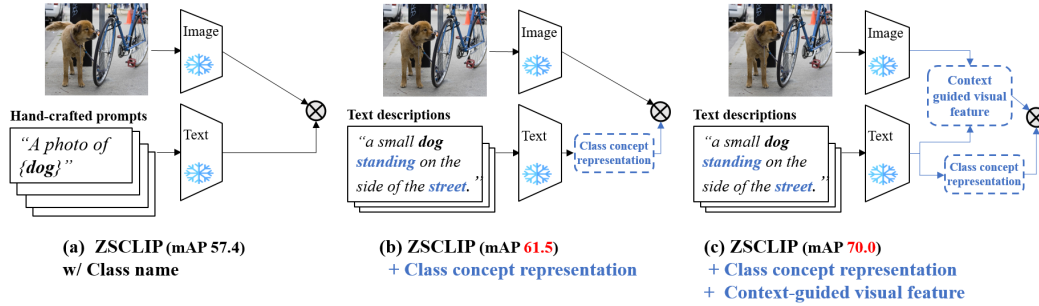


Figure 1: Illustration of our methods applied to zero-shot CLIP (ZSCLIP) [32]. (a→b) Class concept is formed from the text descriptions that contain rich contextual information with relevant class names and other related words, yielding substantially improved performance without aligning with visual features yet. (b→c) Context-guided visual feature is transformed from visual feature so that it is in the same linear space as class concept representation, yielding significantly improved performance.

37 image [37, 18, 41], mitigating the difficulty of acquiring annotated samples. However, prompt tuning
 38 inherently requires labeled data with additional training and may be susceptible to overfitting for
 39 context tokens, hindering generalization. The capability of VLMs for label-free and/or training-free
 40 classification has been exploited using prompt engineering [32, 34, 50, 4]. However, prompt ensem-
 41 bles by averaging text features from simple hand-crafted prompts (e.g., “a sketch of {class name}”)
 42 yielded marginal improvements and struggled with multi-label recognition. Thus, the approach of
 43 prior works on zero-shot or prompt-tuning based multi-label recognition using *class names* to obtain
 44 class-related text features from VLMs may not use the full capacity of VLMs properly.

45 Humans form concepts on words from past experience, especially using their patterns of co-occurrence
 46 with *other words* [5, 29, 20]. Inspired by this perspective in cognitive neuroscience, we propose a
 47 novel approach of exploiting VLMs for multi-label recognition by replacing single *class name*-related
 48 hand-crafted prompts with our proposed *class concept* representation using text descriptions such
 49 as “A **person** holding a large pair of **scissors**,” capturing rich contextual information with target
 50 class names (e.g., person) as well as related words (e.g., holding, scissors). Our class concept will
 51 be constructed from rich contextual descriptions on classes that may contain diverse and realistic
 52 patterns of co-occurrence with target class name and other related class names. Then, this novel text
 53 features with class concept representation requires aligned visual features with them for multi-label
 54 recognition to properly match them with our class concepts. Thus, we propose context-guided visual
 55 features to bring VLM’s visual features to the same representation domain as our class concept
 56 representation by using our sequential attention. See Fig. 1 for the differences of performing multi-
 57 label recognition using (a) prior zero-shot approach (ZS-CLIP), (b) our proposed class concepts from
 58 text descriptions and (c) our proposed context-guided visual features on the same space as the class
 59 concepts. We demonstrated that our proposed methods achieved improved performance on multiple
 60 benchmark datasets without additional training (tuning), without additional labels (text-image pairs)
 61 and with negligible increase in inference time. Here is the summary of the contributions:

- 62 • Proposing a novel class concept representation for training-free multi-label recognition tasks
 63 using VLMs from massive text descriptions inspired by how human forms concept on words.
- 64 • Proposing a context-guided visual feature, transformed onto the same text feature space as
 65 class concepts using sequential attention for better aligning multi-modal features.
- 66 • Demonstrating that our methods synergistically improve the performance of ZSCLIP and
 67 other state-of-the-art prompt tuning methods with a negligible increase in inference time.

68 2 Related Works

69 **Multi-label image recognition with CLIP.** Multi-Label Recognition (MLR) aims to identify all
 70 semantic labels within an image. However, it is difficult to collect the annotation of multi-label images
 71 which involve complex scenes and diverse objects. Recently, prompt tuning with the pre-trained vision-

72 language model CLIP has been developed to address the high labeling costs of multi-label images in
73 incomplete label setting. Among them, DualCoOp [37] proposed a novel prompt tuning approach
74 that trains positive and negative learnable contexts with class names in the partially labeled setting.
75 For mitigating data-limited or label-limited issues, TaI-DPT [18] proposed effective dual-grained
76 prompt tuning method using easily accessible text descriptions. It is worth noting that TaI-DPT
77 used the same text descriptions as ours not for performing training-free multi-label recognition
78 itself, but for label-free prompt tuning by replacing the image features with the contextual text
79 features (text as image) under the conventional framework of multi-label recognition with class
80 name. SCPNet [14] is designed to leverage the structured semantic prior from CLIP to complement
81 deficiency of label supervision for MLR with incomplete labels. CDUL [1] proposed unsupervised
82 multi-label recognition through pseudo-labeling using CLIP, alleviating the annotation burden. Even
83 though recent works has demonstrated outstanding performance of multi-label recognition task, they
84 still require tuning costs or labeled dataset to adapt pre-trained CLIP to various downstream tasks. In
85 this work, our method enables training-free and label-free adaptation of CLIP into downstream tasks,
86 utilizing the text descriptions.

87 **Training-free enhancement with CLIP.** For single-label recognition, recent works has developed
88 the training-free enhancement of CLIP. ZPE [4] weighted-averaged many prompts by automatically
89 scoring the importance of each prompt in zero-shot manner for improving prompt ensemble technique.
90 CALIP [19] designed a simple parameter-free attention module for zero-shot enhancement over CLIP
91 without any tuning of model parameter. With few-shot samples, Tip-Adapter [51] proposed training-
92 free approach for fast adaptation to target task, obtaining the weights of adapter using few-shot
93 samples during inference. Since these methods were originally developed for single-label recognition,
94 it is difficult to be directly applied to multi-label recognition. In multi-label recognition, our method
95 enables training-free enhancement and demonstrated its effectiveness on the benchmark dataset.

96 3 Method

97 First of all, we propose *class concept* representation as a training-free approach for multi-label
98 recognition instead of *class name* by exploiting pre-trained VLM and rich contextual text descriptions.
99 Secondly, we also propose context-guided visual feature that can enhance the alignment of the
100 visual feature of VLM with our novel class concept. Our proposed methods are label-free as well as
101 training-free so that they can be applicable *synergetically* for most existing VLM-based multi-label
102 recognition methods. The overall pipeline of our method is illustrated in Figure 2.

103 3.1 Class Concept Representation

104 Humans form concepts on words from past experience, often using their patterns of co-occurrence
105 with *other words* [5, 29, 20]. For example, the word “apple” does not exist alone, but often comes
106 with the verb “eat” or the noun “basket.” However, it may not well associate with other words such
107 as “fly” or “space.” Fortunately, we can easily obtain rich contextual text descriptions from various
108 public sources, including captions from benchmark datasets [26, 23, 24, 30], web crawling and large
109 language models [38, 7, 40, 28]. These text descriptions do not only contain *class names*, but also
110 include *other words* like class-related verbs and nouns in real-world contexts.

111 Assume that rich contextual text descriptions were gathered from the public sources that include one
112 or multiple class names. We denote the set of text descriptions as $Z^{all} = \{z_1, z_2, \dots, z_M\}$ where z_i
113 refers to an individual text description. M denotes the total number of text descriptions across all
114 classes. Note that M can be dynamically changed at inference since our proposed method does not
115 require additional training, thus can be seen as test-time adaptation. Assuming that the target task
116 uses the class names of person, scissors, clock, building and cake, the examples of the contextual text
117 descriptions from Z^{all} are as follows:

118 “A **person** holding a large pair of **scissors**.”
119 “A **clock** mounted on top of a **building** in the city.”
120 “Half of a white **cake** with coconuts on top.”

121 TaI-DPT [18] used these descriptions with rich contextual information as a surrogate for images
122 to propose a label-free prompt tuning. In this work, we propose to use these descriptions to form
123 concepts on class names to compare with images, so that ways of using them are completely different.

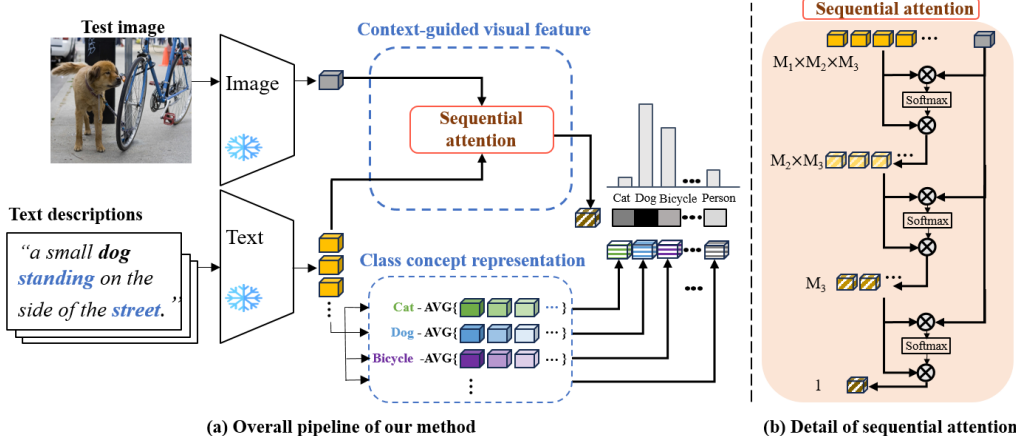


Figure 2: (a) Overall pipeline of our method. 1) Class concept representation: VLM’s text features from the rich contextual descriptions associated with each class name are used to construct the class concept. 2) Context-guided visual features: VLM’s visual features are sequentially transformed onto the class concept representation space using (b) sequential attention mechanism.

124 We define the class concept as a vector in the space constructed by the text descriptions as fol-
 125 lows. Firstly, the linear space \mathcal{Z} can be constructed by spanning the VLM’s text features from
 126 all text descriptions z_i in \mathcal{Z}^{all} using the VLM’s text encoder $\mathcal{E}_{\text{txt}}(z_i) \in R^{1 \times D}$, leading to
 127 $\mathcal{Z} = \text{span}\{\mathcal{E}_{\text{txt}}(z_1), \mathcal{E}_{\text{txt}}(z_2), \dots, \mathcal{E}_{\text{txt}}(z_M)\}$. Secondly, we propose the class concept for a target
 128 class name c as a vector t_c^{concept} in the space \mathcal{Z} by defining it as follows:

$$t_c^{\text{concept}} = \sum_{i=1}^M w_{c,i} \mathbb{1}_c(z_i) \mathcal{E}_{\text{txt}}(z_i) \in R^{1 \times D} \quad (1)$$

129 where $\mathbb{1}_c(z_i)$ an indicator function such that $\mathbb{1}_c(z_i) = 1$ if the text description z_i contains the class
 130 name c and $\mathbb{1}_c(z_i) = 0$ otherwise. The weight $w_{c,i}$ is assigned to the text feature of each text
 131 description within a class c and it is assumed to be normalized within the class. In this work, we set
 132 $w_{c,i} = 1 / \sum_j^M \mathbb{1}_c(z_j)$ for $\forall i$, thus will be the same for all i for each class, which was guided by the
 133 prior work on prompt ensembling [4], demonstrated that the prompt ensembling with equal weights
 134 achieved significant performance gains that were comparable to weighted ensembling for single-label
 135 recognition. Each class concept can be stored individually or together as a matrix.

136 Our class concept representation thus consists of various text features including diverse contextual
 137 information related to the target class name. For instance, the descriptions for the class name “dog”
 138 should contain the target class name as the following examples of the text descriptions:

139 “A **dog** greets a sheep that is in a sheep pen.”
 140 “A woman walks her **dog** on a city sidewalk.”
 141 “A **dog** with goggles is in a motorcycle side car.”

142 Note that the descriptions include the target class name (bold) as well as other related words in class-
 143 related contexts (underline) as intended. We expect that our novel class concepts will be beneficial for
 144 multi-label recognition due to other nouns (other class names) as well as other verbs to better explain
 145 the context where the target class name is used. In this work, we obtain the texts from two sources to
 146 collect the sufficient contextual text descriptions. The first source is the MS-COCO dataset [26] that is
 147 publicly available and the second source is large language model (*i.e.*, GPT-3.5[28]) that can generate
 148 the several sentences quickly if the set of class names related to the target task were provided.

149 3.2 Context-Guided Visual Feature

150 Our novel class concept representation forms new vectors for diverse class names in the linear space
 151 \mathcal{Z} instead of the embedding space of the VLM where the text and image encoders were relatively
 152 well-aligned. Thus, it is expected that the class concept representation and the VLM’s visual feature

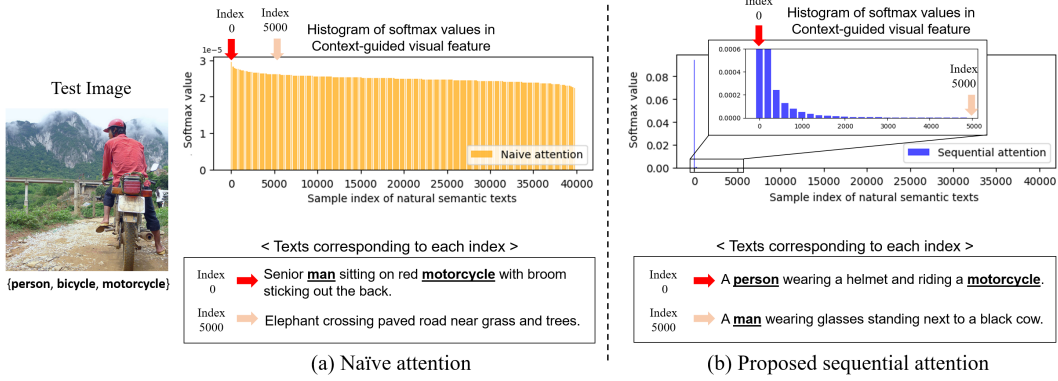


Figure 3: Softmax values can be used to weigh the relevance with the given image. However, (a) naive attention mechanisms yielded almost equal softmax values, thus may include texts with low relevance. The proposed sequential attention method focuses on a subset of texts most relevant to the test image, thus can transform visual features to context-guided visual features for multi-label recognition by assigning very high softmax value to the relevant text at index 0 while very low softmax value to the irrelevant text at index 5000.

153 may not be aligned well. Here, we propose context-guided visual feature by transforming the visual
 154 features of the VLM onto the same space as the class concept representation \mathcal{Z} by using our sequential
 155 attention with the text descriptions \mathcal{Z}^{all} that were used for class concept construction.

156 For the target image q and the VLM’s visual encoder $\mathcal{E}_{img}(q)$, the L2-normalized global visual feature
 157 f is obtained by using $\mathcal{E}_{img}(q) \in R^{1 \times D}$ and the flatten local visual feature $F \in R^{HW \times D}$ is also
 158 constructed by using $\mathcal{E}_{img}(P_{i,j}(q))$ where $P_{i,j}(\cdot)$ is an extractor of the (i, j) th patch of the input image.
 159 Then, we aim to transform both the global visual feature vector f and the local visual feature matrix F
 160 onto the same linear space \mathcal{Z} as our class concept representation. One easy way is to “project” these
 161 visual features f and F onto the space \mathcal{Z} by computing the cosine similarity between visual features
 162 (f and the column vectors of F) and all the text features $t_i = \mathcal{E}_{txt}(z_i) \in R^{1 \times D}$, $i = 1, \dots, M$
 163 that constructed \mathcal{Z} . Unfortunately, when the softmax function is applied to the cosine similarity
 164 values, they tend to become similar, thus weigh both relevant and irrelevant texts almost equally
 165 as illustrated in Figure 3 (a). To address this challenge, we propose sequential attention, applying
 166 the softmax function to part of the cosine similarity values by dividing them into G groups. For the
 167 text feature matrix $T = [t_1 \ t_2 \ \dots \ t_M] \in R^{M \times D}$, let us determine M_i for $i = 1, \dots, G$ such that
 168 $M = \prod_{i=1}^G M_i$ and reshape the text feature matrix to be $T \in R^{M_1 \times \dots \times M_G \times D}$. Then, propose to
 169 sequentially apply the following attention process for G iterations for estimating both global and
 170 local context-guided visual features $v^{(k)}$ and $V^{(k)}$, respectively, at the k th iteration:

$$v^{(k)} = \begin{cases} T & \text{if } k = 0, \\ \text{Softmax}_{\dim_k} \left(\frac{f(v^{(k-1)})^t}{\alpha_f} \right) v^{(k-1)} & \text{if } k > 0, \end{cases} \quad (2)$$

171

$$V^{(k)} = \begin{cases} T & \text{if } k = 0, \\ \text{Softmax}_{\dim_k} \left(\frac{F(V^{(k-1)})^t}{\alpha_F} \right) V^{(k-1)} & \text{if } k > 0, \end{cases} \quad (3)$$

172 where α_f and α_F denote the modulation parameters, Softmax_{M_k} refers to the softmax operation
 173 applied along the dimension corresponding to M_k . In this work, we utilize $v^{(3)}$ and $V^{(3)}$ to compute
 174 classification score. The sequential attention process is illustrated in Figure 2 (b). Figure 3 further
 175 demonstrates that our sequential attention is particularly effective in handling massive text descriptions.
 176 Without sequential attention, weighted averaging essentially becomes equal averaging.

177 3.3 Multi-Label Recognition with Class Concepts

178 **Architecture of model.** Two encoders of CLIP are denoted as \mathcal{E}_{img} and \mathcal{E}_{txt} for the visual encoder
 179 and text encoder, respectively. Following Tai-DPT [18], we adopt the structure of double-grained
 180 prompts (DPT), which has been shown effective for enhancing zero-shot multi-label recognition
 181 performance. To obtain visual representations at both coarse-grained and fine-grained levels, we

182 extract the local visual feature map $F = \mathcal{E}_{\text{img}}(x) \in R^{HW \times D}$ is extracted before attention pooling
 183 layer, where H and W are spatial dimension of visual feature. After attention pooling layer, we
 184 obtain the global visual feature $f \in R^{1 \times D}$. Similarly, text features $t = \mathcal{E}_{\text{txt}}(z) \in R^{1 \times D}$ are obtained
 185 by projecting the End-of-Sentence (EOS) token of the text prompt. Thus, we leverage both global
 186 and local visual features for multi-label recognition.

187 **Inference.** Through our sequential attention, we obtain the context-guided visual features $v^{(G)}$ and
 188 $V^{(G)}$ at both global and local levels, respectively. The similarity score S^{glo} and S^{loc} are calculated
 189 between the transformed context-guided visual features $v^{(G)}$, $V^{(G)}$ and the class concepts $t_c^{concept}$
 190 using the cosine similarity $\Psi(\cdot, \cdot)$ as follows:

$$S_c^{tot} = S_c^{glo} + S_c^{loc} = \Psi(v^{(G)}, t_c^{concept}) + \sum_{j=1}^{HW} \text{Softmax}(s_{c,j}^{loc}) \cdot s_{c,j}^{loc} \quad (4)$$

191 where S_c^{tot} is the classification score for the class c and $s_{c,j}^{loc} = \Psi([V^{(G)}]_j, t_c^{concept})$ for the class c .
 192 For obtaining S_c^{loc} , we employ the spatial aggregation over HW [37].

193 Finally, we combined ZSCLIP[32] and other prompt tuning methods with our training-free approach
 194 through simple logit ensemble. In our experiments, we demonstrate the effectiveness of integrating of
 195 our method with existing methods, thereby boosting the performance of multi-label recognition.

196 4 Experiments

197 4.1 Implementation Details

198 **Architecture.** We employ CLIP ResNet-50 in the Table. 2 and Table. 3 and ResNet-101 in other
 199 experiments as the visual encoders and the CLIP transformer as the text encoder for ZSCLIP[32],
 200 TaI-DPT [18], DualCoOP [37] and our method in the paper. In addition, ZSCLIP[32], TaI-DPT [18]
 201 and our method are based on the double-grained prompt [18] for both global and local features¹.

202 **Datasets.** For evaluation, we performed multi-label recognition experiments on 3 benchmark datasets.
 203 MS-COCO [26] consists of 80 classes with 82,081 images for training and 40,504 images for test.
 204 VOC2007[16] consists of 20 object classes with 5,011 image for training and 4,952 images for test.
 205 NUS-WIDE[12] consists of 81 concepts with 161,789 image for training and 107,859 image for
 206 test. For MS-COCO [26] and VOC2007 [26], text description source is from MS-COCO [26]. For
 207 NUS-WIDE[12], we gathered the text descriptions from GPT-3.5. Note that there is example of text
 208 template for extracting sentence from GPT-3.5 in supplementary.

209 **Inference Details.** In the paper, we set the total number of text descriptions, denoted as M , for
 210 the MSCOCO[26], VOC2007[16], and NUS-WIDE[12] at 40,000, 64,000, and 57,600, respectively.
 211 Note that we prepared the text embeddings of every text descriptions from CLIP text encoder in
 212 advance. We set values of modulation parameter α via validation.

213 4.2 Evaluation on Limited Data Setting

214 To evaluate our method, we conducted the experiments in limited data scenarios, including zero-shot
 215 and few-shot settings for data-limited cases and partially labeled setting for label-limited cases. Note
 216 that only our method provides training-free enhancement of CLIP without tuning cost for multi-label
 217 recognition. Therefore, our method can be easily combined with existing methods to improve their
 218 performance.

219 **Evaluation on Zero-Shot Setting.** We performed comparison studies for different zero-shot and fully
 220 supervised methods in multi-label image recognition. To evaluate the effectiveness of our method
 221 which, we combined our method with existing zero-shot methods, ZSCLIP[32] and TaI-DPT [18],
 222 for zero-shot setting, as shown in Table 1. Additionally, we utilized the fully supervised method,
 223 DualCoOp[37] with our method, for zero-shot learning setting (ZSL) as presented in Table 2.

224 Table 1 summarizes the results of the zero-shot experiment on benchmark datasets. In MS-COCO [26]
 225 and VOC2007 [16], TaI-DPT [18] and our method utilized the public language data from MS-
 226 COCO [26]. By applying our method to ZSCLIP[32] and TaI-DPT [18] during inference, we yield
 227 performance improvements without tuning costs. Especially, the performance of ZSCLIP[32] with

¹<https://github.com/guozix/TaI-DPT>

Table 1: Multi-label recognition with zero-shot methods on MS-COCO [26], VOC2007 [16] and NUS-WIDE [12]. Without training, our method significantly enhances the performance of existing zero-shot methods. The evaluation is based on mAP.

Training-free	Methods	MS-COCO[26]	VOC2007[16]	NUS-WIDE[12]
✓	ZSCLIP[32]	57.4	82.8	37.3
✓	+Ours	70.0 (+12.6)	89.2 (+6.4)	46.6 (+9.3)
✗	TaI-DPT[18]	68.0	88.9	46.5
✓	+Ours	70.9 (+2.9)	90.1 (+1.2)	49.1 (+2.6)

Table 2: Multi-label recognition with 17 unseen classes on MS-COCO [26]. In zero-shot learning (ZSL, recognizing only unseen classes) and generalized ZSL (GZSL, recognizing both seen and unseen classes), our method effectively supplements the complementary information of unseen classes to the supervised DualCoOp[37] on 48 seen classes. The evaluation is based on mAP.

Methods	ResNet-50		ResNet-101	
	ZSL	GZSL	ZSL	GZSL
DualCoOp[37]	78.2	70.2	82.9	74.9
+Ours	82.9 (+4.7)	73.2 (+3.0)	87.6 (+4.7)	78.0 (+3.1)

228 our method is notably enhanced, achieving better and comparable performance to TaI-DPT [18],
 229 which requires mild tuning. In NUSWIDE [12], we incorporate contextual text descriptions from
 230 a large language model (GPT-3.5) to validate the potential of utilizing generated texts instead of
 231 well-curated caption data. With provided class name of NUSWIDE [12], we readily gathered the
 232 massive set of text descriptions within a short amount of time. TaI-DPT [18] is trained with the
 233 public caption data from OpenImages[23]. Our method exceeds the performance of ZSCLIP[32] and
 234 TaI-DPT [18] by a large margin, with improvements of 9.3 mAP and 2.6 mAP, respectively.

235 Table 2 shows the results of the zero-shot learning setting for unseen classes. In MS-COCO [26],
 236 we follow the DualCoOp[37] and split the dataset into 48 seen classes and 17 unseen classes.
 237 The evaluation is conducted in both zero-shot setting (ZSL, recognizing only unseen classes) and
 238 generalized zero-shot setting (GZSL, recognizing both seen and unseen classes). Based on prompt
 239 tuning, DualCoOp[37] trains learnable context tokens on 48 seen classes and achieves the state-of-the-
 240 art performance on both ZSL and GZSL. Our method was originally designed to handle novel classes
 241 (unseen classes) by leveraging text descriptions. As a result, our method significantly improved
 242 the ZSL and GZSL performance of the supervised DualCoOp[37] by providing complementary
 243 information. Table 1 and Table 2 demonstrate the effectiveness of our method performing training-
 244 free enhancement of CLIP with only text descriptions that are easily obtained.

245 **Evaluation on Few-Shot Setting.** We performed comparison study with few-shot methods in multi-
 246 label recognition. In TaI-DPT [18], they have investigate to confirm the effectiveness of their zero-shot
 247 method. Here, we further validate our method, which is zero-shot test-time task adaption without
 248 tuning costs.

249 Table 3 summarizes the results of the few-shot methods on MS-COCO dataset [26], especially using
 250 1 and 5 shot samples for all classes. While existing few-shot methods [3, 35, 54, 51] demonstrated
 251 the performance enhancements with an increase of labeled samples, TaI-DPT [18] and our method
 252 are performed within the zero-shot setting. By applying our method with existing zero-shot methods
 253 (ZSCLIP[32] and TaI-DPT [18]), we consistently enhance performance, as already demonstrated in a
 254 zero-shot setting. In the absence of labeled samples and tuning, we achieved comparable performance
 255 with ML-FSL[35] and better results than other few-shot methods utilizing 5-shot samples.

256 **Evaluation on Partially Labeled Setting.** Due to high costs of annotation in multi-label image
 257 recognition, training with partially labeled samples [37, 21, 31, 6] has been studied. Following
 258 DualCoOp [37], we performed the evaluation of partially labeled setting. As shown in Table 4,
 259 our method supplements the decreased performance of DualCoOp [37] caused by partially labeled
 260 samples by providing complementary information during inference. Through zero-shot test time task
 261 adaptation without tuning costs, we consistently enhance the the performance of DualCoOp [37] on

Table 3: Comparison with few-shot methods on MS-COCO [26]. The evaluation is based on mAP with 16 novel classes. For each shot, we highlighted the best performance in bold.

Training-free	Methods	0-shot	1-shot	5-shot
✗	LaSO[3]	-	45.3	58.1
✗	ML-FSL[35]	-	54.4	63.6
✗	CoOp[54]	-	46.9	55.6
✓	Tip-Adapter[51]	-	53.8	59.7
✓	ZSCLIP[32]	49.7	-	-
✓	+Ours	58.5 (+8.8)	-	-
✗	TaI-DPT[18]	59.2	-	-
✓	+Ours	61.4 (+2.2)	-	-

Table 4: Performance of multi-label recognition based on the partially labeled dataset [26, 16, 12]. Without training and labeled samples, our method consistently enhanced the performance of supervised DualCoOp [37] over all partial label ratio. DualCoOp [37] is reproduced and the evaluation is based on mAP.

Datasets	Method	Partial label									Avg.
		10%	20%	30%	40%	50%	60%	70%	80%	90%	
MS-COCO	SARB[31]	71.2	75.0	77.1	78.3	79.6	79.6	80.5	80.5	80.5	77.9
	DualCoOp[37]	80.8	82.2	82.8	83.0	83.5	83.8	83.9	84.1	84.2	82.7
	DualCoOp[37]+Ours	81.5	82.8	83.3	83.5	84.0	84.2	84.4	84.5	84.6	83.6
VOC2007	SARB[31]	83.5	88.6	90.7	91.4	91.9	92.2	92.6	92.8	92.9	90.7
	DualCoOp[37]	91.6	93.3	93.7	94.3	94.5	94.7	94.8	94.9	94.8	94.0
	DualCoOp[37]+Ours	92.5	93.9	94.3	94.7	94.9	95.0	95.1	95.2	95.1	94.5
NUS-WIDE	DualCoOp[37]	54.0	56.1	56.9	57.4	57.9	57.8	58.0	58.4	58.8	57.3
	DualCoOp[37]+Ours	55.0	56.9	57.7	58.2	58.6	58.6	58.8	59.2	59.5	58.1

262 all benchmark dataset. Furthermore, we achieved the performance of DualCoOp [37] trained with
 263 90% labels by applying our method with DualCoOp trained with 60% labels from MS-COCO [26],
 264 50% labels from VOC2007 [16], and 70% labels from NUSWIDE [12].

265 4.3 Ablation Study and Analysis

266 4.3.1 Effectiveness of our method

267 To verify the effectiveness of components of our method, we conducted an ablation study for analyzing
 268 our method. As shown in Table 5, we first proposed a novel class concept representation with text
 269 descriptions by class to ZSCLIP[32]. Since the text descriptions contain the semantic meaning among
 270 multiple class names and contextual information for multi-label recognition, the alignment between
 271 visual features of test image and text features are improved compared to the hand-crafted prompts as
 272 shown in the Fig.1. Thus, the performance is increased by 4.1 mAP and 1.1 mAP on MS-COCO [26]
 273 and VOC2007 [16], respectively. Then, we performed the context-guided visual feature using a large
 274 set of text descriptions, Z^{all} . Transforming the visual features into same text feature space as our class
 275 concept representation is essential to minimize the gap between visual feature from task-agnostic
 276 visual encoder and text features for each class. Constructing context-guided visual feature, our method
 277 yield remarkable performance gain by 8.5 mAP and 5.3 mAP on MS-COCO [26] and VOC2007 [16],
 278 respectively. Thus, we effectively designed our method that improves the alignment between visual
 279 and text features.

280 4.3.2 The Number of Text Descriptions

281 We investigate the effect of the number of text descriptions for our method. As shown in Table 6,
 282 we evaluated performance by increasing the number of randomly selected text descriptions from 1K
 283 to 32K texts. With only 1K text descriptions, our method enhances performance by approximately

Table 5: Effectiveness of our method on MS-COCO [26] and VOC2007 [16]. Each component of our method consistently improves performance, with significant enhancements achieved particularly in context-guided visual feature through narrowing the gap between visual and text features. The evaluation is based on mAP.

Method	MS-COCO [26]	VOC2007 [16]
Baseline (ZSCLIP[32])	57.4	82.8
+Class concept representation	61.5(+4.1)	83.9(+1.1)
+Context-guided visual feature	70.0(+8.5)	89.2(+5.3)

Table 6: Ablation studies in terms of the number of the text descriptions. As increasing the number of texts, we measured the performance of ZSCLIP[32] with our method in mAP on MS-COCO [26] and VOC2007 [16]. Note that ZSCLIP[32] achieves 57.4 mAP and 82.8 mAP for MS-COCO [26] and VOC2007 [16], respectively.

Dataset	Number of text descriptions					
	1K	2K	4K	8K	16K	32K
MS-COCO [26]	65.8	68.4	68.5	69.1	69.6	69.9
VOC2007 [16]	88.1	88.5	88.8	88.9	89.0	89.1

284 8 mAP on MS-COCO [26] and 5 mAP on VOC2007 [16], respectively. As the number of text
 285 descriptions ranges from 1K to 32K, the text embeddings of Z^{all} can cover the wider range of text
 286 dataset, resulting in increased performance gains. For adapting to novel classes during inference, our
 287 method not only achieves a significant performance improvement with only 1K texts but also further
 288 enhances performance as the quantity of texts increases.

289 4.3.3 Analysis of Inference Time

290 We analyzed the inference time of our method depending on the number of text descriptions. When
 291 extracting text embeddings from the text descriptions in advance, we measure the inference time
 292 as the number of text descriptions increases. ZSCLIP[32], as the baseline model, processes each
 293 sample for classification in 7.2ms. When the number of texts increases from 1K to 32K, integrating
 294 ZSCLIP[32] with our method only increases the inference time by 0.4-0.5ms, with tests conducted on
 295 the RTX3090. In addition, Our method (6.8GB) requires slightly more memory than ZSCLIP (6.5GB)
 296 on VOC2007 [16]. Therefore, our method presents a simple and efficient approach for training-free
 297 enhancement approach at inference.

298 5 Conclusion

299 In this paper, we propose a novel class concept representation from massive text descriptions for
 300 training-free multi-label recognition tasks. Inspired by how humans form concepts based on words,
 301 as studied in cognitive neuroscience, we replace single class name prompts with the class concept
 302 representation that capture various patterns of co-occurrence with other words. To further enhance
 303 alignment between multi-modal features of VLMs, we propose a context-guided visual representation
 304 that is transformed onto the same linear space as the class concept representation. Remarkably,
 305 our proposed method outperforms zero-shot prompt tuning methods such as TaI-DPT and achieves
 306 significant enhancements over ZSCLIP and other state-of-the-art prompt tuning methods without
 307 requiring parameter tuning or labeled samples, and with minimal inference time overhead.

308 **Limitations.** While our method achieved impressive results with training-free enhancement of CLIP,
 309 it exhibits limitations. First, a significant performance gap exists compared to prompt tuning methods
 310 with full samples, like DualCoOp [37]. Second, the computational memory demands of our method
 311 grow at a faster rate than ZSCLIP[32] as the batch size increases.

References

- 312
- 313 [1] Rabab Abdelfattah, Qing Guo, Xiaoguang Li, Xiaofeng Wang, and Song Wang. Cdul: Clip-driven
314 unsupervised learning for multi-label image classification. In *ICCV*, 2023.
- 315 [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc,
316 Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for
317 few-shot learning. *NeurIPS*, 2022.
- 318 [3] Amit Alfassy, Leonid Karlinsky, Amit Aides, Joseph Shtok, Sivan Harary, Rogerio Feris, Raja Giryes, and
319 Alex M Bronstein. Laso: Label-set operations networks for multi-label few-shot learning. In *CVPR*, 2019.
- 320 [4] James Urquhart Allingham, Jie Ren, Michael W Dusenberry, Xiuye Gu, Yin Cui, Dustin Tran, Jeremiah Zhe
321 Liu, and Balaji Lakshminarayanan. A simple zero-shot prompt weighting technique to improve prompt
322 ensembling in text-image models. In *ICML*, 2023.
- 323 [5] Lawrence W Barsalou. Perceptual symbol systems. *Behavioral and brain sciences*, 22(4):577–660, 1999.
- 324 [6] Emanuel Ben-Baruch, Tal Ridnik, Itamar Friedman, Avi Ben-Cohen, Nadav Zamir, Asaf Noy, and Lihi
325 Zelnik-Manor. Multi-label classification with partial annotations using class-aware selective loss. In *CVPR*,
326 2022.
- 327 [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
328 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
329 In *NeurIPS*, 2020.
- 330 [8] Dolly Carrillo, Vivian F López, and María N Moreno. Multi-label classification for recommender systems.
331 In *Trends in Practical Applications of Agents and Multiagent Systems: 11th International Conference on*
332 *Practical Applications of Agents and Multi-Agent Systems*, pages 181–188. Springer, 2013.
- 333 [9] Tianshui Chen, Tao Pu, Hefeng Wu, Yuan Xie, and Liang Lin. Structured semantic transfer for multi-label
334 recognition with partial labels. In *AAAI*, 2022.
- 335 [10] Tianshui Chen, Muxin Xu, Xiaolu Hui, Hefeng Wu, and Liang Lin. Learning semantic-specific graph
336 representation for multi-label image recognition. In *CVPR*, pages 522–531, 2019.
- 337 [11] Zhao-Min Chen, Xiu-Shen Wei, Peng Wang, and Yanwen Guo. Multi-label image recognition with graph
338 convolutional networks. In *CVPR*, 2019.
- 339 [12] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. Nus-wide: a
340 real-world web image database from national university of singapore. In *CIVR*, pages 1–9, 2009.
- 341 [13] Elijah Cole, Oisín Mac Aodha, Titouan Lorieul, Pietro Perona, Dan Morris, and Nebojsa Jojic. Multi-label
342 learning from single positive labels. In *CVPR*, 2021.
- 343 [14] Zixuan Ding, Ao Wang, Hui Chen, Qiang Zhang, Pengzhang Liu, Yongjun Bao, Weipeng Yan, and Jungong
344 Han. Exploring structured semantic prior for multi label recognition with incomplete labels. In *CVPR*,
345 2023.
- 346 [15] Thibaut Durand, Nazanin Mehrasa, and Greg Mori. Learning a deep convnet for multi-label classification
347 with partial labels. In *CVPR*, 2019.
- 348 [16] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The
349 pascal visual object classes (voc) challenge. *IJCV*, 2010.
- 350 [17] Marian George and Christian Floerkemeier. Recognizing products: A per-exemplar multi-label image
351 classification approach. In *ECCV*, 2014.
- 352 [18] Zixian Guo, Bowen Dong, Zhilong Ji, Jinfeng Bai, Yiwen Guo, and Wangmeng Zuo. Texts as images in
353 prompt tuning for multi-label image recognition. In *CVPR*, 2023.
- 354 [19] Ziyu Guo, Renrui Zhang, Longtian Qiu, Xianzheng Ma, Xupeng Miao, Xuming He, and Bin Cui. Calip:
355 Zero-shot enhancement of clip with parameter-free attention. In *AAAI*, 2023.
- 356 [20] Paul Hoffman, James L. McClelland, and Matthew A. Lambon Ralph. Concepts, control, and context: A
357 connectionist account of normal and disordered semantic cognition. *Psychological Review*, 125(3):293–328,
358 Apr.
- 359 [21] Ping Hu, Ximeng Sun, Stan Sclaroff, and Kate Saenko. Dualcoop++: Fast and effective adaptation to
360 multi-label recognition with limited annotations. *arXiv preprint arXiv:2308.01890*, 2023.

- 361 [22] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung,
362 Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text
363 supervision. In *ICML*. PMLR, 2021.
- 364 [23] Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Haija, Alina Kuznetsova, Hassan
365 Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, et al. Openimages: A public dataset for large-scale
366 multi-label and multi-class image classification. *Dataset available from <https://github.com/openimages>*,
367 2017.
- 368 [24] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen,
369 Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision
370 using crowdsourced dense image annotations. *IJCV*, 2017.
- 371 [25] Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie
372 Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm.
373 *arXiv preprint arXiv:2110.05208*, 2021.
- 374 [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
375 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- 376 [27] Yongcheng Liu, Lu Sheng, Jing Shao, Junjie Yan, Shiming Xiang, and Chunhong Pan. Multi-label image
377 classification via knowledge distillation from weakly-supervised detection. In *ACMMM*, 2018.
- 378 [28] OpenAI. Gpt-4 technical report, 2023.
- 379 [29] Karalyn Patterson, Peter J Nestor, and Timothy T Rogers. Where do you know what you know? the
380 representation of semantic knowledge in the human brain. *Nature reviews neuroscience*, 8(12):976–987,
381 2007.
- 382 [30] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana
383 Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence
384 models. In *ICCV*, 2015.
- 385 [31] Tao Pu, Tianshui Chen, Hefeng Wu, and Liang Lin. Semantic-aware representation blending for multi-label
386 image recognition with partial labels. In *AAAI*, 2022.
- 387 [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
388 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
389 natural language supervision. In *ICML*, 2021.
- 390 [33] Nikolaos Sarafianos, Xiang Xu, and Ioannis A Kakadiaris. Deep imbalanced attribute classification using
391 visual attention aggregation. In *ECCV*, 2018.
- 392 [34] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei
393 Xiao. Test-time prompt tuning for zero-shot generalization in vision-language models. *Advances in Neural
394 Information Processing Systems*, 35:14274–14289, 2022.
- 395 [35] Christian Simon, Piotr Koniusz, and Mehrtash Harandi. Meta-learning for multi-label few-shot classifica-
396 tion. In *WACV*, 2022.
- 397 [36] Josef Sivic and Andrew Zisserman. Video google: Efficient visual search of videos. *Toward category-level
398 object recognition*, pages 127–144, 2006.
- 399 [37] Ximeng Sun, Ping Hu, and Kate Saenko. Dualcoop: Fast adaptation to multi-label recognition with limited
400 annotations. *NeurIPS*, 2022.
- 401 [38] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,
402 and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model. *Stanford Center for
403 Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 2023.
- 404 [39] Ivona Tautkute, Tomasz Trzcinski, Aleksander P Skorupa, Łukasz Brocki, and Krzysztof Marasek. Deep-
405 style: Multimodal search engine for fashion and interior design. *IEEE Access*, 2019.
- 406 [40] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
407 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and
408 fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 409 [41] Ao Wang, Hui Chen, Zijia Lin, Zixuan Ding, Pengzhang Liu, Yongjun Bao, Weipeng Yan, and Guiguang
410 Ding. Hierarchical prompt learning using clip for multi-label classification with single positive labels. In
411 *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5594–5604, 2023.

- 412 [42] Shuai Wang, Daoan Zhang, Zipei Yan, Jianguo Zhang, and Rui Li. Feature alignment and uniformity for
413 test time adaptation. In *CVPR*, 2023.
- 414 [43] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, and Ronald M Summers. Tienet: Text-image embedding
415 network for common thorax disease classification and reporting in chest x-rays. In *CVPR*, 2018.
- 416 [44] Ya Wang, Dongliang He, Fu Li, Xiang Long, Zhichao Zhou, Jinwen Ma, and Shilei Wen. Multi-label
417 classification with label graph superimposing. In *AAAI*, 2020.
- 418 [45] Xiu-Shen Wei, Quan Cui, Lei Yang, Peng Wang, and Lingqiao Liu. Rpc: A large-scale retail product
419 checkout dataset. *arXiv preprint arXiv:1901.07249*, 2019.
- 420 [46] Ning Xu, Congyu Qiao, Jiaqi Lv, Xin Geng, and Min-Ling Zhang. One positive label is sufficient:
421 Single-positive multi-label learning with label enhancement. *NeurIPS*, 2022.
- 422 [47] Lewei Yao, Runhui Huang, Lu Hou, Guansong Lu, Minzhe Niu, Hang Xu, Xiaodan Liang, Zhenguo Li,
423 Xin Jiang, and Chunjing Xu. Filip: Fine-grained interactive language-image pre-training. *arXiv preprint*
424 *arXiv:2111.07783*, 2021.
- 425 [48] Vacit Oguz Yazici, Abel Gonzalez-Garcia, Arnau Ramisa, Bartlomiej Twardowski, and Joost van de Weijer.
426 Orderless recurrent models for multi-label classification. In *CVPR*, 2020.
- 427 [49] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong
428 Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv*
429 *preprint arXiv:2111.11432*, 2021.
- 430 [50] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and
431 Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. In *Proceedings of the IEEE/CVF*
432 *Conference on Computer Vision and Pattern Recognition*, pages 18123–18133, 2022.
- 433 [51] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hong-
434 sheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *arXiv preprint*
435 *arXiv:2111.03930*, 2021.
- 436 [52] Yong Zheng, Bamshad Mobasher, and Robin Burke. Context recommendation using multi-label classi-
437 fication. In *IEEE/WIC/ACM International Joint Conferences on WI and IAT*, volume 2, pages 288–295,
438 2014.
- 439 [53] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for
440 vision-language models. In *CVPR*, 2022.
- 441 [54] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language
442 models. *IJCV*, 2022.

443 A Generation of Text Descriptions using LLMs

444 Our proposed method leverages the text descriptions for enhancing the alignment between the
445 visual and text features. In practice, gathering the proper text descriptions is an essential process for
446 replacing the hand-crafted prompts. As mentioned in the main paper, the text descriptions can be
447 readily gathered from benchmark dataset, web crawling, or large language models. Recent advances
448 in large language models (LLMs) enable to rapidly generate text descriptions that are similar to
449 image captions in MS-COCO [26]. Therefore, we utilized the generated text descriptions from large
450 language model. With provided class name of NUSWIDE [12], Fig. 6 illustrates the example of input
451 prompt template and corresponding generated text descriptions using GPT3.5. We carefully designed
452 the instruction of input prompt including main description, constraints, examples of bad and good
453 cases, class names of target task and output format.

454 B Implementing Other Zero-Shot Training-Free Method

455 In single-label recognition, CALIP [19] proposed zero-shot alignment enhancement of CLIP for adapt-
456 ing target task without few-shot samples or additional training. The parameter-free attention module of
457 cross-modal interaction effectively enhances the alignment of visual and text features. CALIP utilized
458 the visual feature $F = \text{Enc}_v(x_k) \in R^{HW \times D}$ via reshaping and the text feature $T = \text{Enc}_t(P^h) \in R^{C \times D}$

Table 7: Ablation study of hyperparameter searching on validation set. We varied the modulation parameters $\alpha_{f,t}$ and $\alpha_{F,t}$ and searched the proper values for context-guided visual feature.

$1/\alpha_{f,t}, 1/\alpha_{F,t}$	MS-COCO	VOC2007	NUS-WIDE
100, 50	69.45	87.62	47.32
80, 40	69.51	87.94	48.31
60, 30	69.25	88.06	49.05
40, 20	67.47	87.37	49.82
20, 10	64.13	85.04	47.33

459 where P^h is a hand-crafted description and C denotes the number of classes. The parameter-free
 460 attention module is formulated as follows:

$$F^a = \text{Softmax}(A/\alpha_t)T, \quad (5)$$

$$T^a = \text{Softmax}(A^T/\alpha_v)F \quad (6)$$

461 where the attention matrix is $A=FV^T \in R^{HW \times C}$, α_t and α_v are the modulation parameters of textual
 462 and visual features, respectively, and T^a and F^a are bidirectionally updated textual and visual features.
 463 After pooling the updated visual feature $F_v^a \in R^{1 \times D}$ and the global visual feature $F_v \in R^{1 \times D}$, the
 464 classification logit S is obtained as below:

$$S = \beta_1 \cdot F_v T^T + \beta_2 \cdot F_v T^{aT} + \beta_3 \cdot F_v^a T^T, \quad (7)$$

465 where $\beta_1, \beta_2, \beta_3$ are the weights for the three logits.

466 CALIP [19] tuned the hyperparameters β_2, β_3 for each dataset while fixed β_1 to be 1 for simplicity.
 467 As shown in Fig. 4, we have explored the value of β_2, β_3 for multi-label recognition setting on
 468 MS-COCO [26] and have observed that the parameter-free attention module consistently decreases
 469 the mAP performance since multi-label recognition covers the identification of multiple objects
 470 within an image, involving complex scene and diverse objects.

471 C Exploring Modulation Parameters

472 For hyperparameter searching, following existing methods for classification tasks, such as zero-
 473 shot [18, 19], training-free [51, 19], and test-time adaptation [42], we explore the modulation pa-
 474 rameters α_t by conducting ablation studies on validation set. For simplicity, we set the value of
 475 $\alpha_{f,t}$ to be half of $\alpha_{F,t}$. As shown in Table 7, the value of $(1/\alpha_{f,t}, 1/\alpha_{F,t})$ is suitable in the range
 476 of (40~80, 20~40). In the experiments of main paper, we set the $(1/\alpha_{f,t}, 1/\alpha_{F,t})$ as (80,40) for
 477 MS-COCO [26], (60,30) for VOC2007 [16] and (40,20) for NUSWIDE [12].

478 D Examples of Local Alignment Enhancement

479 In Fig. 5, we visualized the examples of local alignment enhancement by applying our method.
 480 Enhancing local alignment is important to recognize multiple objects in a test image [37]. Our
 481 proposed method enhances the local alignment between the visual features of test image and the
 482 text features of each class name, thereby suppressing the false-positive prediction. Therefore, Fig. 5
 483 demonstrates the effectiveness of our method.

484 E Positive and Negative Societal Impacts

485 As a positive societal impact, our method can allow people with limited computing resources to
 486 achieve better performance in multi-label classification using existing vision-language models. This
 487 is because it does not require extensive training or labeled data. However, as a negative societal
 488 impact, the failure of classification could produce the negative side effects. For example, in security
 489 applications, incorrect classification of objects could lead to false alarms or missed detections,
 490 potentially compromising safety and security.

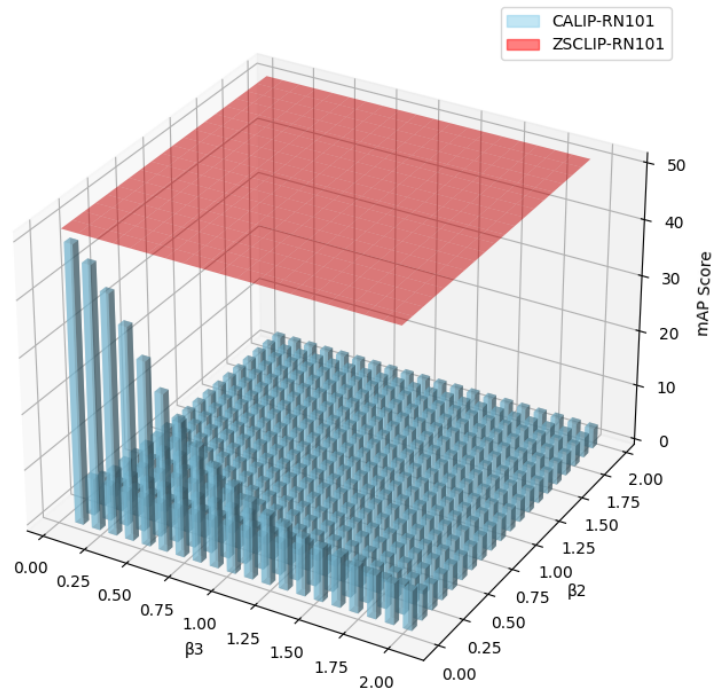


Figure 4: Results of hyperparameter searching of CALIP [19] on MS-COCO [26] on β_2 and β_3 . Applying the parametric-free attention module of CALIP consistently decreases performance as compared to the zero-shot CLIP (ZSCLIP) [32].

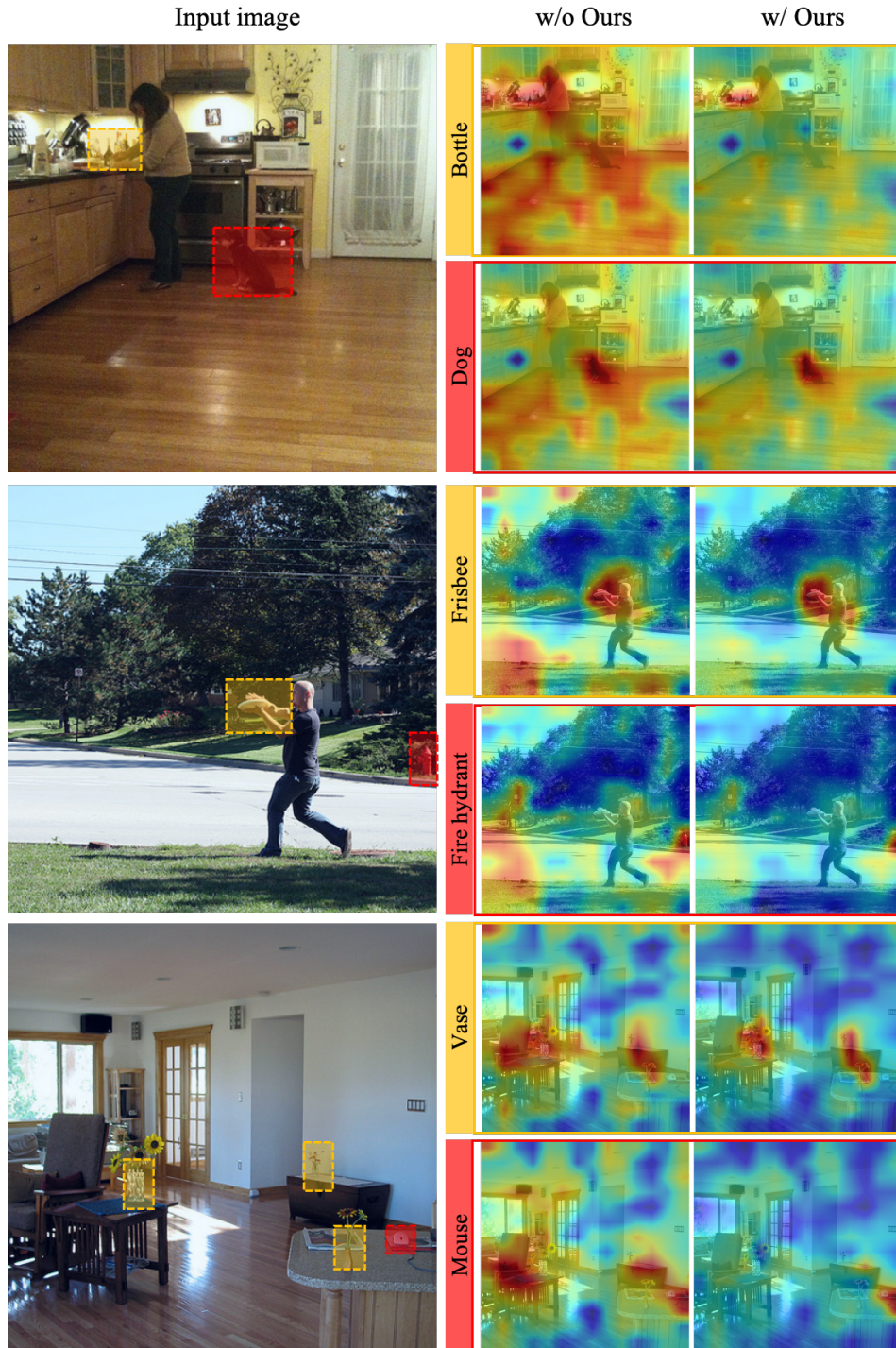



Figure 5: Additional examples of local alignment enhancement via our method. We visualized the test image in the left column and its corresponding spatial similarity map of each class name in the right column. The yellow and red boxes refer to the bounding boxes for different labels in a multi-label setting. By applying our method, the local alignment is enhanced across multiple objects in a test image, thereby suppressing false-positive predictions.



Such as MSCOCO caption, generates the diverse 100 sentences including following class names for multi-label classification caption. The sentence should include proper context and multiple classes in a sentence.


Constraints:

1. Must be in JSON file format.
2. No code is required.
3. Sentence length should be between 40-50 characters.
4. Sentences must start with a lowercase letter.
5. Each sentence must contain at least one class name or more.
6. Must describe a scene that is likely to occur in reality.
7. Do not start the subject with 'I'.
8. Sentences must be visually descriptive.
9. Avoid abstract descriptions and ensure sentences depict concrete scenes.
10. A total of 100 sentences are required.
11. must include text of "{class name}" for all sentences but it's not need to start with this text and should be natural sentence.

Examples:
 bad example : boats dot the harbor near the sleepy town
 good example : a cat lounges in the sunny apartment window

Class names:
 airport animal beach ... whales window zebra

Output format :
 [{"sentence": "an airport buzzes with planes and travelers at dawn"}, {"sentence": "surf crashes onto the beach as swimmers look on"}, {"sentence": "a person reads a book in the quiet public library"}, {"sentence": "tigers roam the rocky terrain under the hot sun"},...]



Here are 100 sentences that meet your requirements, each containing one or more of the provided class names. These sentences describe various scenes likely to occur in reality:

```
[ {"sentence": "a dog runs through the park chasing a frisbee"},
  {"sentence": "birds soar high in the clear blue sky"},
  {"sentence": "surfers catch waves at the crowded beach"},...]
```

Figure 6: Example of text description generation using GPT3.5 for contextual text descriptions of NUSWIDE [12]. We carefully designed the input prompt to ensure that the generated sentences include the class name of the target task. The elements considered in designing the input prompt include the main description, constraints, examples, class names, and the desired output format.

491 **NeurIPS Paper Checklist**

492 **1. Claims**

493 Question: Do the main claims made in the abstract and introduction accurately reflect the
494 paper's contributions and scope?

495 Answer: [Yes]

496 Justification: We clearly state our contributions in both the abstract and introduction. Espe-
497 cially, we summarize our contributions in the last part of the introduction.

498 Guidelines:

- 499 • The answer NA means that the abstract and introduction do not include the claims
500 made in the paper.
- 501 • The abstract and/or introduction should clearly state the claims made, including the
502 contributions made in the paper and important assumptions and limitations. A No or
503 NA answer to this question will not be perceived well by the reviewers.
- 504 • The claims made should match theoretical and experimental results, and reflect how
505 much the results can be expected to generalize to other settings.
- 506 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
507 are not attained by the paper.

508 **2. Limitations**

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

510 Answer: [Yes]

511 Justification: We discuss the limitations of our method in conclusion.

512 Guidelines:

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

539 **3. Theory Assumptions and Proofs**

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

542 Answer: [NA]

543 Justification: Our paper does not include theoretical results, assumptions and proof.

544 Guidelines:

- 545 • The answer NA means that the paper does not include theoretical results.
- 546 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 547 referenced.
- 548 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 549 • The proofs can either appear in the main paper or the supplemental material, but if
- 550 they appear in the supplemental material, the authors are encouraged to provide a short
- 551 proof sketch to provide intuition.
- 552 • Inversely, any informal proof provided in the core of the paper should be complemented
- 553 by formal proofs provided in appendix or supplemental material.
- 554 • Theorems and Lemmas that the proof relies upon should be properly referenced.

555 4. Experimental Result Reproducibility

556 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

557 perimental results of the paper to the extent that it affects the main claims and/or conclusions

558 of the paper (regardless of whether the code and data are provided or not)?

559 Answer: [Yes]

560 Justification: We provide the details of the used model, hyperparameters, source of datasets

561 and proposed algorithm for reproducing main experimental results.

562 Guidelines:

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

594 5. Open access to data and code

595 Question: Does the paper provide open access to the data and code, with sufficient instruc-

596 tions to faithfully reproduce the main experimental results, as described in supplemental

597 material?

598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648

Answer: [Yes]

Justification: We provide open access to the code for our proposed method. In our experiments, we utilize a publicly accessible benchmark dataset described in Section ??.

Guidelines:

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

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide the details of the experimental setting, including hyperparameters, in Sections 4.1 and C. The generation details of the texts used in our method are also provided in Section A.

Guidelines:

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

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We do not report the statistical significance of the experimental results, as our method does not rely on statistical variables for inference.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- 649 • The method for calculating the error bars should be explained (closed form formula,
650 call to a library function, bootstrap, etc.)
- 651 • The assumptions made should be given (e.g., Normally distributed errors).
- 652 • It should be clear whether the error bar is the standard deviation or the standard error
653 of the mean.
- 654 • It is OK to report 1-sigma error bars, but one should state it. The authors should
655 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
656 of Normality of errors is not verified.
- 657 • For asymmetric distributions, the authors should be careful not to show in tables or
658 figures symmetric error bars that would yield results that are out of range (e.g. negative
659 error rates).
- 660 • If error bars are reported in tables or plots, The authors should explain in the text how
661 they were calculated and reference the corresponding figures or tables in the text.

662 8. Experiments Compute Resources

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

666 Answer: [Yes]

667 Justification: We provide the information of types of compute worker (GPU model), memory
668 usage and inference time in Section. 4.3.3.

669 Guidelines:

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

678 9. Code Of Ethics

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

681 Answer: [Yes]

682 Justification: We abide by the NeurIPS Code of Ethics.

683 Guidelines:

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

689 10. Broader Impacts

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

692 Answer: [Yes]

693 Justification: We discuss the positive and negative societal impacts of our paper in the
694 supplementary material.

695 Guidelines:

- 696 • The answer NA means that there is no societal impact of the work performed.
- 697 • If the authors answer NA or No, they should explain why their work has no societal
698 impact or why the paper does not address societal impact.

- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 714
- 715
- 716
- 717
- 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.
 - 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).

718 11. Safeguards

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

722 Answer: [NA]

723 Justification: Our paper does not pose a high risk for misuse in terms of model and dataset.

724 Guidelines:

- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

735 12. Licenses for existing assets

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

739 Answer: [Yes]

740 Justification: We cite the papers that provide datasets, code and models in the Section. 4.

741 Guidelines:

- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 753
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 754
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 755
- 756

757 **13. New Assets**

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

760 Answer: [Yes]

761 Justification: We provide the detail documentation for the code in our submission.

762 Guidelines:

- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 763
- 764
- 765
- 766
- 767
- 768
- 769
- 770

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

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

775 Answer: [NA]

776 Justification: Our paper does not involve neither crowdsourcing nor research with human
777 subjects.

778 Guidelines:

- 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.
- 779
- 780
- 781
- 782
- 783
- 784
- 785
- 786

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

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

793 Answer: [NA]

794 Justification: Our paper does not involve neither crowdsourcing nor research with human
795 subjects.

796 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.
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804

805
806

- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.