

Is CLIP the main roadblock for fine-grained open-world perception?

Lorenzo Bianchi*[†] , Fabio Carrara* , Nicola Messina* , Fabrizio Falchi* 
*CNR-ISTI, Pisa, Italy [†] University of Pisa, Italy
Email: <name>.<surname>@isti.cnr.it

Abstract—Modern applications increasingly demand flexible computer vision models that adapt to novel concepts not encountered during training. This necessity is pivotal in emerging domains like extended reality, robotics, and autonomous driving, which require the ability to respond to open-world stimuli. A key ingredient is the ability to identify objects based on free-form textual queries defined at inference time – a task known as *open-vocabulary object detection*. Multimodal backbones like CLIP are the main enabling technology for current open-world perception solutions. Despite performing well on generic queries, recent studies highlighted limitations on the *fine-grained* recognition capabilities in open-vocabulary settings – i.e., for distinguishing subtle object features like color, shape, and material. In this paper, we perform a detailed examination of these open-vocabulary object recognition limitations to find the root cause. We evaluate the performance of CLIP, the most commonly used vision-language backbone, against a fine-grained object-matching benchmark, revealing interesting analogies between the limitations of open-vocabulary object detectors and their backbones. Experiments suggest that the lack of fine-grained understanding is caused by the poor separability of object characteristics in the CLIP latent space. Therefore, we try to understand whether fine-grained knowledge is present in CLIP embeddings but not exploited at inference time due, for example, to the unsuitability of the cosine similarity matching function, which may discard important object characteristics. Our preliminary experiments show that simple CLIP latent-space re-projections help separate fine-grained concepts, paving the way towards the development of backbones inherently able to process fine-grained details. The code for reproducing these experiments is available at <https://github.com/lorebianchi98/FG-CLIP>.

Index Terms—fine-grained understanding, open-vocabulary object detection, image-text matching, evaluation study

I. INTRODUCTION

Nowadays, pivotal technologies such as extended reality, autonomous driving, and robotics necessitate more than adherence to closed-set assumptions; they demand the ability to adapt to novel concepts not encountered during the training phase. Open-vocabulary object detection (OVD) stands out as a critical task for achieving this adaptability, as underlined by several open challenges on egocentric data [5], [9]. It involves recognizing objects not included in the training dataset, thus

overcoming the limitations inherent in traditional detectors confined to a predefined set of objects.

This flexibility is typically achieved by common-space multi-modal models. These models embed images and texts in a shared latent space, thanks to the contrastive pre-training performed on large datasets of image-text pairs scraped from the web. This simplifies similarity calculations, which can be achieved through an efficient dot product between the representations. CLIP [26] stands out as the most widely utilized model in this category. Its capabilities become crucial in open-vocabulary object detection, where models typically perform the task by i) detecting regions of the image that are likely to contain objects, ii) computing the similarity between the embedding of the detected image region and that of a set of free-form texts defined at test time, called *vocabulary*.

While open-vocabulary detectors excel in generalizing to new concepts not encountered during training, recent studies indicate limitations in capturing fine-grained properties of objects [2]. For instance, they may encounter difficulties in distinguishing between a *light brown* dog and a *dark brown* one (Figure 1). One potential explanation for these shortcomings is that CLIP representations may exhibit bias towards category-level concepts while overlooking attribute-level nuances [4]. To the best of our knowledge, it is not well studied in the literature whether fine-grained properties are absent in the latent space or if these characteristics exist, but trivial matching methods (e.g., dot product, cosine similarity) are insufficient to extract this information.

In this work, we assess the performance of CLIP, the most used backbone in open-vocabulary object detectors, on a fine-grained benchmark to scrutinize its capacity to accurately discern intricate properties of objects. The findings reveal that the performance of CLIP in fine-grained understanding mirrors that of an open-vocabulary object detector based on CLIP, indicating shared challenges. This suggests that the fine-grained issues observed in the open-vocabulary object detector may be attributed to the image-text alignment carried out by CLIP rather than to detector localization failures.

Subsequently, by incorporating additional layers on top of frozen visual and textual CLIP encoders and training them on a fine-grained dataset, we illustrate the model’s proficiency in accurately assigning the relevant attribute to the object. This shows that the original CLIP embeddings contain fine-grained information that is ignored during the matching phase.

In summary, our paper contributes to the field in the

This work was partially supported by the following projects: SUN – Social and hUman ceNtered XR (Horizon Europe R&I GA:101092612), FAIR – Future Artificial Intelligence Research - Spoke 1 (EU NextGenerationEU PNRR M4C2 PE00000013), ITSERR – ITalian Strengthening of the Esfri Ri Resilience (EU NextGenerationEU CUP:B53C22001770006). MUCES – a MUltimedia platform for Content Enrichment and Search in audiovisual archives (EU NextGenerationEU - PRIN 2022 PNRR P2022BW7CW - CUP: B53D23026090001)

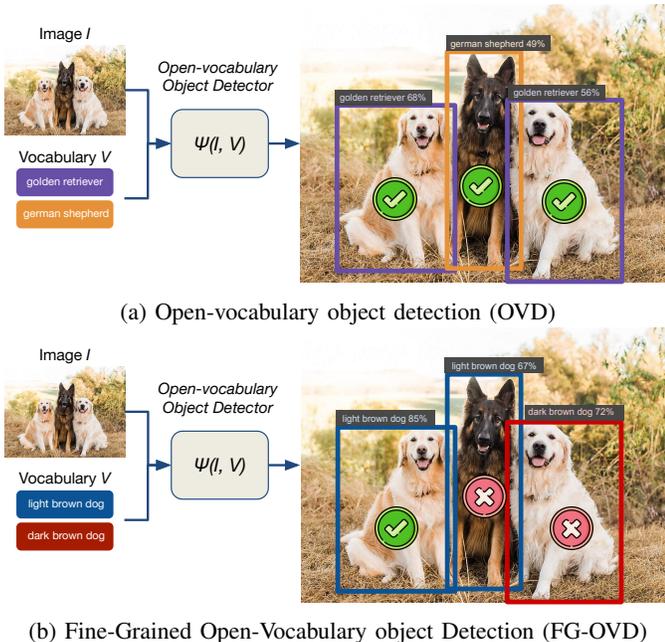


FIG. 1: **OVD (a) and FG-OVD (b)**: in the latter, fine-grained details about the categories to detect are given as free-form text in the input vocabulary.

following ways:

- We comprehensively evaluate and analyze CLIP’s performance on a fine-grained open-vocabulary object detection benchmark. This investigation sheds light on the possibility that challenges faced by open-vocabulary object detectors may be attributed to issues within the CLIP latent space.
- We showcase the existence of fine-grained information within CLIP’s latent space through the implementation of lightweight architectures trained on frozen CLIP embeddings. This demonstration is substantiated by the model’s ability to perform fine-grained matching successfully.

II. RELATED WORK

A. Image-Text matching

In recent years, the focus of researchers on image-text matching increased. The foundation of the shared space approach for cross-modal matching begins with the exploration of hinge-based triplet ranking loss with hard-negative mining. This was first attempted with GRU as text extractor [8], and later with Transformer Encoder [15], [18]–[20], [25], [29], [32].

With the growing strength of Transformers also in vision tasks [6], many works exploited an early-fusion approach, leveraging Transformers encoder to jointly process images and texts from the very beginning of the proposed architectures [14], [16], [17], [30], [31], [35]. These methods treat image-text matching as a binary classification problem, where, given as input an image-text pair, they train the Transformer architecture to predict the probability that the text correctly

describes the image. While achieving good performance, these models cannot be used in many real-case scenarios since they are computationally expensive at inference time, as they require processing every image-text pair to obtain the score on the whole test set. Consequently, many methods preferred to exploit a late-fusion approach, keeping the visual and textual pipeline separated [11], [18]–[20], [26], [28], [32]. This allows separate pipelines for producing the images and text representations, and the similarity score can be computed in a second stage with a simple dot product.

Among these models, CLIP [26] stands out as one of the most widely used for performing image-text matching. The majority of open-vocabulary object detectors rely on the knowledge acquired by CLIP to embed object regions and vocabulary entries in the same feature space [1], [7], [10], [21], [22], [33], [36], [37]. This approach circumvents the need to conduct inference for each pair of detected image regions and vocabulary entries.

B. Fine-grained understanding

Although CLIP shows strong performance in tasks such as classification and coarse-grained retrieval, it has shortcomings in associating nuanced properties between sentences and images [23], [24], [27], [36]. Yuksekgonul et al. [34] suggest that contrastive pretraining does not optimize the model’s understanding of the relationships between objects and their attributes. Furthermore, standard retrieval benchmarks are considered inadequate for assessing the compositional understanding of such models.

Krojer et al. [13] highlight that vision-language models tend to overlook fine-grained visual information. Similarly, Chen et al. [4] show that representations learned from contrastive pretraining in common-space multimodal models are biased toward category-level concepts rather than attributes, making attribute recognition difficult.

These limitations are even more pronounced in tasks where models usually rely on CLIP latent space to work in an open-world environment, such as open-vocabulary object detection. Some recent advanced benchmarks reveal the weaknesses of these models [2], [3]. These benchmarks not only test the models’ ability to generalize to unseen objects during training but also assess their ability to recognize object attributes. These benchmarks show that current methods are far from achieving satisfactory results. Our work follows up on those recent findings and investigates potential causes for the performance gap in fine-grained settings.

III. METHOD

In this work, we aim to address two pivotal questions:

Q1. *Can the limitations observed in open-vocabulary object detection regarding fine-grained understanding be traced back to deficiencies within the CLIP latent space?*

Q2. *If yes, do these limitations arise from an absence of such information in the latent space itself, or is it a consequence of the inadequacy of trivial matching methods (e.g., dot product,*

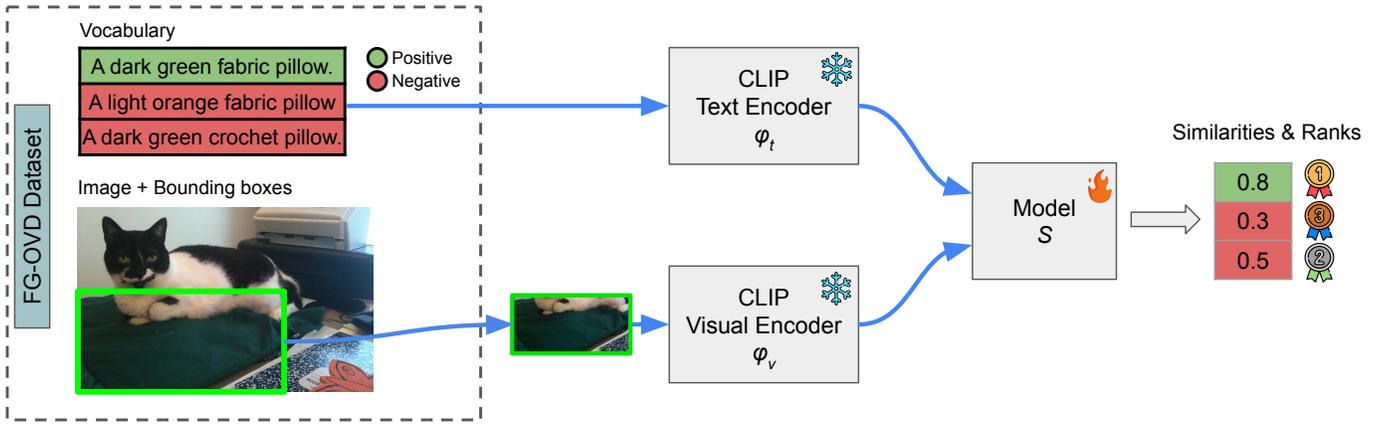


Fig. 2: **FG-OVD Dataset for CLIP Matching.** We leverage the Fine-Grained Open-Vocabulary object Detection (**FG-OVD**) benchmark suite and training set to investigate our two research questions **Q1** and **Q2**. For each object, we extract the corresponding bounding box crop and compute its visual encoding using CLIP. Text embeddings are then generated for the assigned vocabulary entries (composed by positive + negative captions) associated with the object. Finally, we calculate the similarity and rank of the positive caption between the image crop and the vocabulary entries. To address **Q1** we use a cosine similarity as model S , and the entire pipeline is used only during inference. To address **Q2**, we choose model S from the solutions described in subsection III-B and train it on the FG-OVD training set.

cosine similarity) to extract this nuanced information?

In this section, we outline the methodologies employed to address these inquiries.

A. CLIP fine-grained evaluation

Our objective in this part is to address question **Q1**, aiming to discern whether the failure of open-vocabulary object detectors can be attributed to the image-text matching in the embedding space or to the localization phase.

To assess the fine-grained knowledge contained in the CLIP latent space and to facilitate a comprehensive analysis in comparison to open-vocabulary object detectors, we used the benchmark suite tailored for Fine-Grained Open-Vocabulary Detection (FG-OVD) proposed in [2]. These benchmarks provide object bounding boxes, each associated with a detailed natural language caption that provides attributes about the object, referred to as a *positive caption*. In addition, each positive caption is associated with a set of semantically similar but subtly different *negative captions*.

For the purposes of matching rather than detection, we evaluate each object individually. This adaptation involved retaining a vocabulary specific to each object and cropping the associated bounding box, sparing the model from performing region proposal and detection. Then, we compute the similarity between the cropped image embedding and the captions embedding using the cosine similarity. This evaluation pipeline is outlined in Figure 2.

To benchmark these results, we compare CLIP with OWL [21], [22], an open-vocabulary object detector that relies on CLIP. The performance gap between the two is only explained by the errors introduced by the region proposal and object localization phases, thus providing an estimate of their contribution to open-vocabulary detection performance.

B. Latent Space Characteristics and Matching Approaches

This section provides the strategy to answer question **Q2**. Assuming that the fine-grained knowledge is present within the CLIP latent space, we hypothesize that the matching scheme used to compare the representations, i.e., the typical cosine similarity, is insufficient to extract this specific information. To explore this possibility, our strategy involves learning a customized similarity function $S(\mathbf{v}, \mathbf{t})$, which takes as input the two embeddings \mathbf{v} and \mathbf{t} obtained from the frozen visual and textual encoders, ϕ_v and ϕ_t respectively. By forcing S to recognize nuanced object properties based only on the embedded information, we can state that successful results in this regard mean that the embeddings inherently encode fine-grained knowledge. To this aim, we will use two distinct datasets having similar image distribution but different annotations. The first one comprises general image-text pairs, which can be used to train and validate our model on standard coarse-grained category-centric classification. Differently, the second dataset is dedicated to training and evaluating the learned function S for fine-grained understanding and is organized in positive and negative captions for each object as explained in Section III-A.

Therefore, the overall training strategy is composed of two steps. First, we perform a warm-up phase in which we train S on the coarse-grained image-text pairing dataset. Then, we fine-tune S using the fine-grained matching dataset. The first warm-up phase works as a fine-tuning of the original CLIP model, which repurposes the CLIP features to work well with our image distribution and consequently initializes the S function. This creates a strong and reliable baseline that we can employ as a reference to track the subsequent decline in coarse-grained performance after the fine-tuning step.

We perform the warm-up on the coarse-grained dataset

using the *hinge-based triplet loss* as a loss function. Namely, we optimize

$$\mathcal{L} = \sum_{\substack{i,j \in \mathcal{B} \\ i \neq j}} [\alpha + S(\mathbf{v}_i, \mathbf{t}_j) - S(\mathbf{v}_i, \mathbf{t}_i)]_+ + [\alpha + S(\mathbf{v}_j, \mathbf{t}_i) - S(\mathbf{v}_i, \mathbf{t}_i)]_+ \quad (1)$$

where S is the similarity function, $\mathcal{B} = \{1 \dots B\}$ is a batch, i is the index of an image I_i described by the text T_i , while $j \neq i$ is the index of a negative image I_j and a negative text T_j taken from the batch. ϕ_v and ϕ_t are respectively the visual and textual CLIP encoders, and $\mathbf{v}_i = \phi_v(I_i)$ and $\mathbf{t}_i = \phi_t(T_i)$ are the frozen encoded text and image, respectively. α is the margin of separation.

We then perform training on the fine-grained dataset, where each image I_i is associated with a vocabulary composed of a positive caption and a set of N similar but slightly different captions, called negative captions. Again, we rely on the *hinge-based triplet loss*, but this time, we adapt it for the fine-grained discrimination task: for each tuple (image, vocabulary), we take as anchor the image, as positive the correct caption, and as negative the incorrect ones. Formally,

$$\mathcal{L}_{\text{FG}} = \sum_{i=1}^B \sum_{j=1}^N [\alpha + S(\mathbf{v}_i, \mathbf{t}_{i,j}^{\text{neg}}) - S(\mathbf{v}_i, \mathbf{t}_i^{\text{pos}})]_+, \quad (2)$$

where B is the batch size, \mathbf{v}_i is the encoding of the i -th image, $\mathbf{t}_i^{\text{pos}}$ is the encoding of the positive caption associated with the i -th image, $\mathbf{t}_{i,j}^{\text{neg}}$ is the encoding of the j -th negative caption associated with image i , α is the minimum separation margin that should hold between positive and negative captions, and $[x]_+ \equiv \max(x, 0)$.

We evaluate different implementations of the matching function S :

1) *Baseline (CLIP matching function):*

$$S(\mathbf{v}, \mathbf{t}) = \cos(\mathbf{v}, \mathbf{t}) \quad (3)$$

The vanilla cosine similarity represents the commonly used matching function used in CLIP and open-vocabulary object detectors to match visual and textual representations.

2) *Linear projection layer:*

$$S(\mathbf{v}, \mathbf{t}) = \cos(W_v \mathbf{v} + b_v, W_t \mathbf{t} + b_t) \quad (4)$$

We propose two linear projection layers that operate on top of the frozen visual (\mathbf{v}) and textual (\mathbf{t}) features, aiming to project the embeddings into a space of the same dimensionality. The final matching function is kept as the cosine similarity between the transformed feature vectors. We aim to explore the feasibility of linearly separating fine-grained concepts (if they exist) embedded in the CLIP embeddings.

3) *Linear projection layer only above text encoder:*

$$S(\mathbf{v}, \mathbf{t}) = \cos(\mathbf{v}, W_t \mathbf{t} + b_t) \quad (5)$$

We introduce a linear projection layer solely above the text encoder while keeping the visual embedding fixed. This setup forces the image latent space to maintain coherence with the original CLIP space while repurposing only textual representations.

4) *Linear projection layer only above visual encoder:*

$$S(\mathbf{v}, \mathbf{t}) = \cos(W_v \mathbf{v} + b_v, \mathbf{t}) \quad (6)$$

Similarly, we apply a linear projection layer solely above the visual encoder, following the same principle as the previous scheme. By projecting only one embedding modality at a time, we can explore any potential asymmetry in the results and gain insights into the role of each modality in capturing fine-grained information.

5) *MLPs layer:*

$$S(\mathbf{v}, \mathbf{t}) = \cos(\text{MLP}_v(\mathbf{v}), \text{MLP}_t(\mathbf{t})) \quad (7)$$

We incorporate two Multi-Layer Perceptrons (MLPs) with a non-linear activation function to repurpose the embeddings before performing the cosine similarity. This scheme aims to explore the effects of non-linearity on both coarse-grained and fine-grained performance and to compare these results with the effects observed with linear projection approaches.

6) *Attention layer:*

$$S(\mathbf{v}, \mathbf{t}) = \sigma \left(\text{MHA}([\text{CLS}, \mathbf{v}, \mathbf{t}])^{(0,0)} \right) \quad (8)$$

We build a multi-head attention layer above the two encoders to compute self-attention utilizing the two embeddings along with a to-be-learned CLS token of the same dimensionality. Then, after the attention layer, we apply the sigmoid function to the first element of the CLS token output. The score of the sigmoid represents the similarity between the two embeddings. This approach explores a more complex and expressive non-linear alternative to the MLP that exploits attention to automatically weight the contribution of image and text features.

IV. EXPERIMENTS

A. Dataset and Metrics

For the warm-up phase and to evaluate the coarse-grained capabilities of our model, we selected the MS-COCO dataset as our image-text pair dataset. We followed the partitioning introduced by Karpathy et al. [12], reserving 113,287 images for training, 5,000 for validation, and 5,000 for testing, each with five captions. We evaluated retrieval performance on COCO using the Recall@ k metric to measure the ability of our model to retrieve relevant text or images accurately. Specifically, Recall@ k assesses the proportion of queries that successfully retrieve the correct item within the first k results.

We used the Fine-Grained Open-Vocabulary Object Detection (FG-OVD) [2] suite for fine-grained training and evaluation, adapting it for classification using crops of object regions instead of detection. This suite associates each object with a customized vocabulary consisting of a detailed sentence describing the object and its attributes, called the positive caption, and a set of negative captions subtly altering certain attributes (see Figure 2). The benchmarks within the suite are categorized by difficulty level (Trivial, Easy, Medium, and Hard), with the degree of change in the captions decreasing as the selected benchmark becomes more difficult. For example, the Hard benchmark indicates that only one

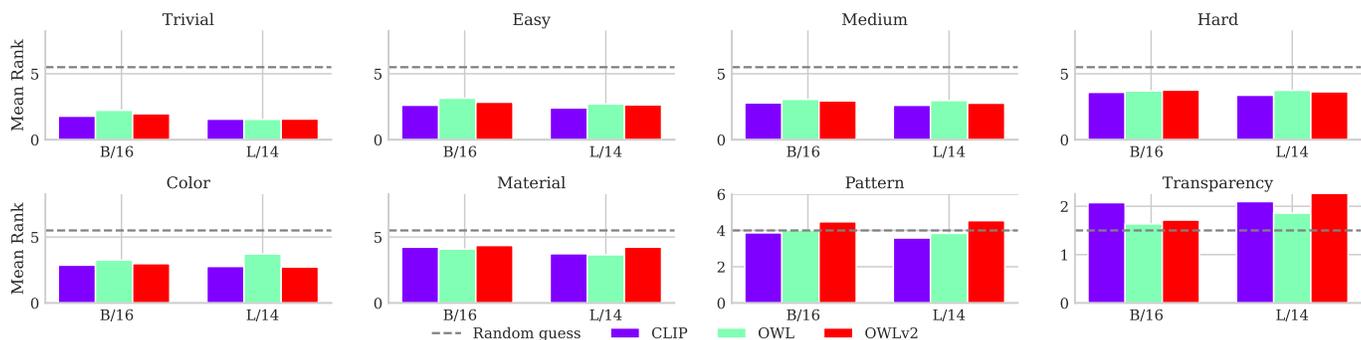


Fig. 3: **CLIP vs. OWL in fine-grained understanding.** We evaluate CLIP and OWL, configured as B/16 and L/14, against the Difficulty-based (first row) and Attribute-based (second row) FG-OVD benchmarks. The bar graph shows the Mean Rank of the positive label (**lower is better**), which represents the average position assigned by the model to the correct label within the overall vocabulary. Vocabulary lengths vary, with 3 for transparency, 8 for pattern, and 11 for other attributes.

TABLE I: **Coarse-grained and Fine-grained performance.** We analyze the performance of the investigated similarity functions S after a warm-up train on COCO and a subsequent fine-tuning on the fine-grained dataset FG-OVD. In the fine-tuned configuration (+FG-OVD rows), we denote the delta between these results and those obtained during the warm-up in parentheses.

	COCO Retrieval						
	FG-OVD	I→T			T→I		
		Mean Rank ↓	R@1 ↑	R@5 ↑	R@10 ↑	R@1 ↑	R@5 ↑
CLIP B/16	2.98	41.5	65.9	76.2	22.6	44.1	54.9
Linear (both)	3.78	48.0	76.6	85.4	37.2	65.6	76.6
+FG-OVD	1.46 (-2.32)	37.1 (-10.9)	66.8 (-9.8)	78.4 (-7.0)	35.6 (-1.6)	63.9 (-1.7)	75.0 (-1.6)
Linear (visual only)	3.53	45.8	74.2	83.7	35.4	64.2	75.3
+FG-OVD	1.54 (-1.99)	39.4 (-6.4)	69.5 (-4.7)	79.8 (-3.9)	34.3 (-1.1)	62.9 (-1.3)	74.1 (-1.2)
Linear (text only)	3.48	47.3	75.1	84.9	36.0	64.3	75.6
+FG-OVD	1.57 (-1.91)	41.1 (-6.2)	70.4 (-4.7)	80.6 (-4.3)	34.7 (-1.3)	63.2 (-1.1)	74.6 (-1.0)
MLP	3.49	45.9	75.5	84.6	36.5	64.9	76.2
+FG-OVD	1.43 (-2.06)	31.9 (-14.0)	60.2 (-15.3)	72.7 (-11.9)	33.6 (-2.9)	62.0 (-2.9)	73.9 (-2.3)
MHA	4.08	36.3	66.1	78.1	29.1	57.6	70.3
+FG-OVD	1.54 (-2.54)	22.3 (-14.0)	48.3 (-17.8)	61.2 (-16.9)	22.6 (-6.5)	49.2 (-8.4)	62.0 (-8.3)

attribute is changed in the negative labels. In addition, the suite includes attribute-based benchmarks (Color, Material, Pattern, and Transparency), where only attributes of a particular type are changed, making it easier to evaluate model performance in each attribute category. To evaluate on this benchmark, we measured the Mean Rank of the positive label within the vocabulary when ranked by descending matching score. As for the number of negative labels in the vocabulary N , we follow [2] and choose 10 for Trivial, Easy, Medium, Hard, Color and Material, 7 for Pattern, and 2 for Transparency. In subsection IV-C, we plot the Mean Rank for each benchmark, while in subsection IV-D, we report the mean of the values obtained from the eight benchmarks.

For training, we used the Hard training set with $N = 10$ negatives.

B. Implementation details

For the warm-up phase on COCO, we train the matching function S with the Adam optimizer, using a learning rate of $5e^{-4}$ for 10 epochs and a triplet loss margin $\alpha = 0.2$. For fine-

tuning on the FG-OVD, we use Adam with a reduced learning rate ($1e^{-5}$) and set the triplet loss separation margin $\alpha = 0.05$. Maintaining a deliberately low margin is critical, as higher values were observed to disrupt the alignment established during the warm-up phase, significantly reducing retrieval performance on COCO. The fine-tuning process includes 10 epochs. Regarding the newly added layers, we set up the attention layer with 64 heads, while the MLPs feature 2 layers with 512 neurons each and the \tanh activation function. We train all the proposed architectures using embeddings extracted from CLIP B/16.

C. Impact of object localization is marginal in FG-OVD

Figure 3 compares a CLIP-based open-vocabulary object detector with vanilla CLIP applied to pre-segmented image regions on the FG-OVD benchmark suite. It is important to note that the task of the detector is more challenging than that of CLIP. Indeed, the detector must also localize the object rather than solely classifying it. We make the comparison with OWL, an open vocabulary object detector based on CLIP,

presented in both its original [22] and second [21] versions, with the same backbones (B/16 and L/14).

Looking at the results, the performance of CLIP mirrors the pattern shown by the detectors. However, the overall performance of CLIP remains relatively low. For example, in the Hard benchmark, the model ranks the correct caption on average around 4th out of 11 possible captions. This highlights the significant limitations of classification using the CLIP latent space and points to the need for significant improvements in open-world fine-grained classification.

Despite facing a more challenging task than CLIP, the detectors' performance is not significantly lower. Their results are consistent with CLIP's patterns, showing similar trends in both the Difficulty-based and Attribute-based benchmarks, with better performance observed in the Color benchmark (an attribute more commonly found in web-scraped images), while the Material and Pattern benchmarks show lower performance. This suggests that the recent challenges open vocabulary object detectors face in fine-grained understanding are more related to classification within the shared image-text embedding space than to the localization phase.

D. A linear projection is enough for fine-grained matching

The results presented in Table I illustrate the performance of the proposed architecture after the warm-up on COCO and the subsequent fine-tuning on FG-OVD.

Observing the results, the warm-up on COCO improves the performance of COCO retrieval compared to the original CLIP, indicating increased specialization in the characteristics of COCO captions and images. Conversely, performance on the fine-grained benchmark is poor because COCO captions generally lack detailed attributes about objects.

After fine-tuning on the fine-grained dataset, there's a slight decrease in retrieval performance on COCO. However, the key observation is the significant decrease in the FG-OVD Mean Rank. This suggests that the newly added layers, relying solely on information from the embeddings, effectively discriminate the correct attributes for object mapping. This addresses our question **Q2**, where we wondered if fine-grained knowledge was missing in the CLIP latent space or if we were not using the right tools to extract it. These results show that we can learn a more complex similarity matching between the representations and that nuanced information is indeed present in CLIP embeddings. In addition, the fact that these results can be achieved only with a linear projection, demonstrate that this type of information can be linearly separated in the embedding space.

Comparing the MLPs with the linear projection, it is evident that nonlinearity does not provide any advantage in maintaining a favorable trade-off between fine- and coarse-grained performance. Although the results on the fine-grained benchmarks are slightly better, the retrieval performance on COCO worsens.

The results obtained with a linear projection applied to only one embedding modality reveal interesting observations.

During the warm-up phase, these projections demonstrate superior fine-grained performance compared to models with both embeddings repurposed. This phenomenon can be attributed to the fact that only one projection needs to be learned, which is forced to maintain coherence with the original CLIP latent space. Consequently, this leads to a smaller fine-grained performance drop during the warm-up phase, which is usually due to the limited fine-grained attribute information within the COCO dataset. In addition, this adherence to the original embedding space results in a minor coarse-grained performance degradation after fine-tuning compared to other tested configurations. What is most compelling is that despite the fine-grained results being slightly lower compared to the other matching functions, there is a great improvement in the fine-grained results despite the need to adhere to the CLIP latent space. This suggests that fine-grained properties are not only present and linearly separable within the CLIP embedding space but can even emerge with a simple linear adjustment of the embedding of only one modality.

Experiments with the multi-head attention layer suggest that more complex and expressive architecture can better adapt to the novel fine-grained task but with a higher risk of overfitting the training data and disrupting the original space.

V. CONCLUSION

We studied the challenges confronting current state-of-the-art cross-modal image-text models in achieving open-world understanding. We began our analysis by examining the relationship between open-vocabulary object detectors and their vision-language backbone, specifically focusing on CLIP. Our results suggest that localization is marginal in the limitations observed in fine-grained open-vocabulary object detection. This demonstrates that the primary problem lies in the interaction between vision and language within the shared latent space.

Furthermore, while fine-grained information exists within the CLIP latent space, the representation is heavily biased towards coarse-grained concepts. This bias causes similar concepts to be positioned too closely within the latent space, making it difficult to detect nuanced differences using traditional cosine similarity. Importantly, we demonstrated that fine-grained information is still linearly separable within latent space despite this bias.

In the future, we aim to explore better pre-training strategies to construct more balanced image-text representations that effectively incorporate fine- and coarse-grained features. In addition, we will investigate alternative matching functions capable of extracting fine-grained features within the CLIP latent space without the need for task-specific datasets to learn this function.

REFERENCES

- [1] Relja Arandjelović, Alex Andonian, Arthur Mensch, Olivier J Hénaff, Jean-Baptiste Alayrac, and Andrew Zisserman. Three ways to improve feature alignment for open vocabulary detection. *arXiv preprint arXiv:2303.13518*, 2023.

- [2] Lorenzo Bianchi, Fabio Carrara, Nicola Messina, Claudio Gennaro, and Fabrizio Falchi. The devil is in the fine-grained details: Evaluating open-vocabulary object detectors for fine-grained understanding. In *[Accepted for publication in CVPR 2024]*, 2024.
- [3] Maria A Bravo, Sudhanshu Mittal, Simon Ging, and Thomas Brox. Open-vocabulary attribute detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7041–7050, 2023.
- [4] Keyan Chen, Xiaolong Jiang, Yao Hu, Xu Tang, Yan Gao, Jianqi Chen, and Weidi Xie. Ovarnet: Towards open-vocabulary object attribute recognition. In *CVPR*, 2023.
- [5] Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, Sanja Fidler, David Fouhey, and Dima Damen. Epic-kitchens visor benchmark: Video segmentations and object relations. In *Proceedings of the Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks*, 2022.
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.
- [7] Yu Du, Fangyun Wei, Zihe Zhang, Miaoqing Shi, Yue Gao, and Guoqi Li. Learning to prompt for open-vocabulary object detection with vision-language model. pages 14084–14093, 2022.
- [8] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. Vse++: Improving visual-semantic embeddings with hard negatives. 2018.
- [9] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022.
- [10] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. 2022.
- [11] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021.
- [12] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137, 2015.
- [13] Benno Krojer, Vaibhav Adlakha, Vibhav Vineet, Yash Goyal, Edoardo Ponti, and Siva Reddy. Image retrieval from contextual descriptions. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3426–3440, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [14] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- [15] Kungpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. Visual semantic reasoning for image-text matching. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4654–4662, 2019.
- [16] Xiujuan Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 121–137. Springer, 2020.
- [17] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.
- [18] Nicola Messina, Giuseppe Amato, Andrea Esuli, Fabrizio Falchi, Claudio Gennaro, and Stéphane Marchand-Maillet. Fine-grained visual textual alignment for cross-modal retrieval using transformer encoders. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 17(4):1–23, 2021.
- [19] Nicola Messina, Fabrizio Falchi, Andrea Esuli, and Giuseppe Amato. Transformer reasoning network for image-text matching and retrieval. In *2020 25th International conference on pattern recognition (ICPR)*, pages 5222–5229. IEEE, 2021.
- [20] Nicola Messina, Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, Fabrizio Falchi, Giuseppe Amato, and Rita Cucchiara. Aladin: Distilling fine-grained alignment scores for efficient image-text matching and retrieval. In *International Conference on Content-based Multimedia Indexing*, pages 64–70, 2022.
- [21] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. Scaling open-vocabulary object detection. 2023.
- [22] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. pages 728–755. Springer, 2022.
- [23] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3170–3180, 2023.
- [24] Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. Valse: A task-independent benchmark for vision and language models centered on linguistic phenomena. *arXiv preprint arXiv:2112.07566*, 2021.
- [25] Leigang Qu, Meng Liu, Da Cao, Liqiang Nie, and Qi Tian. Context-aware multi-view summarization network for image-text matching. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1047–1055, 2020.
- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. pages 8748–8763. PMLR, 2021.
- [27] Kanchana Ranasinghe, Brandon McKinzie, Sachin Ravi, Yinfei Yang, Alexander Toshev, and Jonathon Shlens. Perceptual grouping in contrastive vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5571–5584, 2023.
- [28] Nikolaos Sarafianos, Xiang Xu, and Ioannis A Kakadiaris. Adversarial representation learning for text-to-image matching. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5814–5824, 2019.
- [29] Matteo Stefanini, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. A novel attention-based aggregation function to combine vision and language. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 1212–1219. IEEE, 2021.
- [30] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations*, 2020.
- [31] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. Image as a foreign language: BEiT pretraining for vision and vision-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [32] Keyu Wen, Xiaodong Gu, and Qingrong Cheng. Learning dual semantic relations with graph attention for image-text matching. *IEEE transactions on circuits and systems for video technology*, 31(7):2866–2879, 2020.
- [33] Xiaoshi Wu, Feng Zhu, Rui Zhao, and Hongsheng Li. CORA: Adapting clip for open-vocabulary detection with region prompting and anchor pre-matching. pages 7031–7040, 2023.
- [34] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International Conference on Learning Representations*, 2022.
- [35] Pengchuan Zhang, Xiujuan Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5579–5588, 2021.
- [36] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. RegionCLIP: Region-based language-image pretraining. pages 16793–16803, 2022.
- [37] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. pages 350–368. Springer, 2022.