# Open-domain Visual Entity Recognition:
# Towards Recognizing Millions of Wikipedia Entities

Hexiang Hu♠   Yi Luan♠   Yang Chen♠♡†   Urvashi Khandelwal♠

Mandar Joshi♠   Kenton Lee♠   Kristina Toutanova♠   Ming-Wei Chang♠

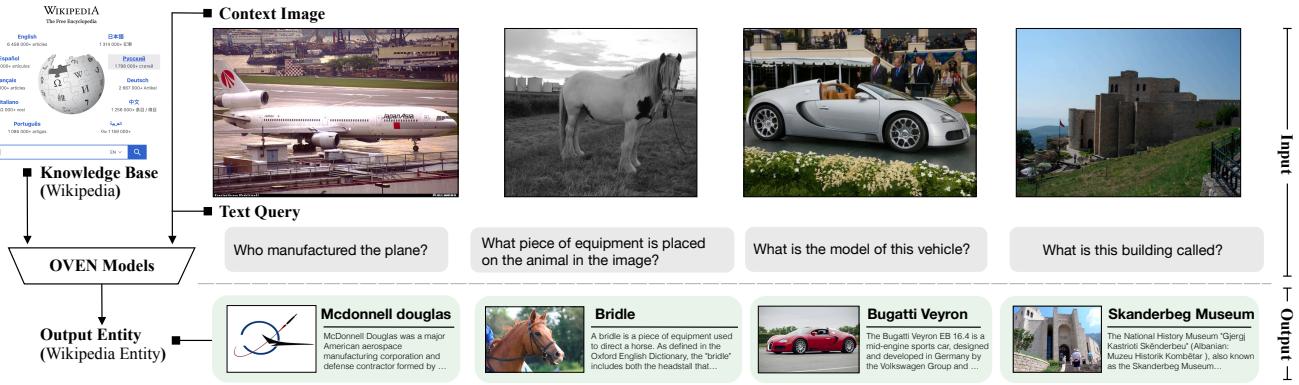♠ **Google Research**   ♡ **Georgia Institute of Technology**

Figure 1: An illustration of the proposed OVEN task. Examples on the right are sampled from the constructed OVEN-Wiki dataset. OVEN aims at recognizing entities *physically presented* in the image or can be *directly inferred* from the image.

## Abstract

*Large-scale multi-modal pre-training models such as CLIP [35] and PaLI [9] exhibit strong generalization on various visual domains and tasks. However, existing image classification benchmarks often evaluate recognition on a specific domain (e.g., outdoor images) or a specific task (e.g., classifying plant species), which falls short of evaluating whether pre-trained foundational models are universal visual recognizers. To address this, we formally present the task of Open-domain Visual Entity recognitioN (OVEN), where a model need to link an image onto a Wikipedia entity with respect to a text query. We construct OVEN-Wiki‡ by re-purposing 14 existing datasets with all labels grounded onto one single label space: Wikipedia entities. OVEN-Wiki challenges models to select among six million possible Wikipedia entities, making it a general visual recognition benchmark with the largest number of labels. Our study on state-of-the-art pre-trained models reveals large headroom in generalizing to the massive-scale label space. We show that a PaLI-based auto-regressive visual recognition model performs surprisingly well, even on Wikipedia entities that have never been seen during fine-tuning. We also find existing pre-trained models yield different strengths: while PaLI-based models obtain higher overall performance, CLIP-based models are better at recognizing tail entities.*

## 1. Introduction

Pre-trained large language models [4, 12], *inter alia*, have shown strong transferable text processing and generation skills in tackling a wide variety of natural language tasks [43, 50, 54] across languages and task formats, while requiring very few manually labeled per-task examples. At the same time, while there has been equally impressive progress in multi-modal pre-training [9, 35], it remains unclear whether similarly universal visual skills, *i.e.*, recognizing millions of coarse-grained and fine-grained visual concepts, have emerged. *Are pre-trained multi-modal models capable of recognizing open-domain visual concepts?*

Answering this question requires a visual recognition dataset with broad coverage of visual domains and tasks, under a universally defined semantic space. Existing recognition benchmarks such as ImageNet [39, 41], Stanford Cars [24], or SUN database [58] represent a large number of visual concepts, but make specific assumptions about the granularity of the target concepts (e.g. building type such as "castle" in ImageNet but not a specific building in the world such as "Windsor Castle"), or limit attention to concepts of the same type such as car models/years. Visual question answering (VQA) datasets test models' abilities to

---

recognize concepts which can be of more flexible granularities and object types, but in practice existing VQA datasets tend to focus on higher-level categories. We aim to assess models' abilities to recognize visual concepts from a close to universal, unified space of labels that covers nearly all visual concepts known to humankind, and at a flexible level of granularity, specified by a user or a downstream application. Given a short specification of each element in the target space of visual concepts (such as a textual description), multimodal pre-trained models could in principle recognize concepts without seeing labeled instances covering each of them.

Towards evaluating models on such universal visual recognition abilities, we introduce the task of **O**pen-domain **V**isual **E**ntity recognitio**N** (OVEN), targeting a wide range of entities and entity granularities, including animals, plants, buildings, locations and much more. Particularly, we construct OVEN-Wiki by building on existing image recognition and visual QA datasets and unifying their label spaces/granularities and task formulations. For our unified label space, we use English Wikipedia which covers millions of visual entities of various levels of granularity and also includes a specification of each entity via its Wikipedia page (containing entity name, text description, images, etc.). Wikipedia also evolves as new entities appear or become known in the world, and can be used as a first approximation of a universal visual concept space.

We re-purpose 14 existing image classification, image retrieval, and visual QA datasets, and ground all labels to Wikipedia. In addition to unifying labels, we unify input recognition intent specifications, which is necessary when combining specialized datasets with the goal of evaluating universal recognition. Given an image showing a car and a tree behind it, OVEN makes the recognition intent explicit via a natural language query such as "What is the model of the car?" or "What is the species of the tree?". Therefore, the OVEN task takes as input an image and a text query[1] that expresses visual recognition intent with respect to the image. The goal is to provide an answer by linking to the correct entity (e.g. BUGATTI VEYRON or BACTRIS GASIPAES) out of the millions of possible Wikipedia entities, each coming with descriptions and a relevant set of images from its Wikipedia page (see Figure 1). Importantly, OVEN requires recognition of entities that were UNSEEN in the training data. Models can still take advantage of the text description and/or images on the Wikipedia page of the UNSEEN entities, as well as knowledge acquired through pre-training.

Human annotators were hired to help create OVEN-Wiki for two reasons. First, grounding labels from the component datasets into Wikipedia entities is non-trivial due to language ambiguity. For example, 'Tornado' can be a weather phe-

nomenon or a type of airplane (PANAVIA TORNADO). To reduce such ambiguity in the grounding, we take multiple steps to refine the labels, including the use of human annotators, a state-of-the-art textual entity linking system [13], and heavy filtering. Second, creating unambiguous textual query intents is also challenging. In many cases, a text query can lead to multiple plausible answers (e.g. of various granularities), and a human often needs to make revisions to make sure no other objects could be correct answers. For our training and development/test sets we rely on semi-automatic processing, but additionally introduce a gold evaluation set, for which annotators thoroughly corrected entity linking errors and rewrote ambiguous input query intents.

Based on OVEN-Wiki, we examine two representative multi-modal pre-trained models, PaLI [9] and CLIP [35], to establish an empirical understanding of the state-of-the-art in universal entity recognition. Particularly, these two models are used for creating an auto-regressive visual entity recognition model (similar to [13]) and a visual entity retrieval model, respectively. Our study suggests that there is a large room for improvement in generalizing to the massive label space. We show that the PaLI-based auto-regressive visual recognition model performs surprisingly well, even on Wikipedia entities that have never been seen during fine-tuning. Digging deeper, we discover that CLIP variants and PaLI-based models make very different kinds of errors. Particularly, PaLI dominates in recognizing popular Wikipedia entities, whereas CLIP models can win consistently on recognizing tail entities.

## 2. Open Domain Visual Entity Recognition

To drive progress in universal entity recognition, we propose the task of Open-domain Visual Entity recognitioN (OVEN). There are two desiderata that we would like to meet for the OVEN task. First, there should exist a universal label space. In OVEN, we make use of a multi-modal knowledge base, such as Wikipedia, to serve as the universal label space, covering millions of entities. Second, the answer label for each OVEN input should be unambiguous. This is particularly challenging when the label space is very large and multi-granular. To accomplish this, OVEN makes use of input text queries to define the recognition intent (*e.g.*, identifying car types or car models), allowing visual concepts from different granularities to be unambiguously specified.

**Task Definition** The input to an OVEN model is an image-text pair $x = (x^p, x^t)$, with the text query $x^t$ expressing intent with respect to the corresponding image $x^p$. Given a unified label space $\mathcal{E}$ which defines the set of all possible entities, the knowledge base $\mathcal{K} = \{(e, p(e), t(e)) \mid e \in \mathcal{E}\}$ is a set of triples, each containing an entity $e$, its corresponding text description $t(e)$ (*i.e.*, name of the entity, description, etc.) and a (possibly empty) set of relevant
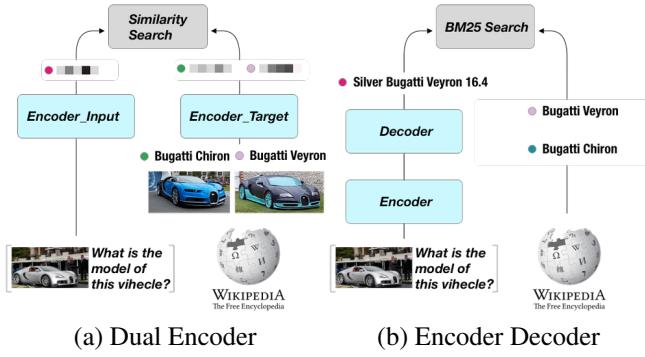
---

(a) Dual Encoder      (b) Encoder Decoder

Figure 2: **Illustration on two OVEN Models.**

images $p(e)$. For instance, an entity $e = $ Q7395937 would have a corresponding textual description $t(e) = $ 'Name: Sabatia campestris; Description:...'[2] and a set $p(e)$ containing one or more images from the corresponding Wikipedia page[3] of SABATIA CAMPESTRIS. We consider the combination of $t(e)$ and $p(e)$ the *multi-modal knowledge* for the entity $e$. As OVEN is a recognition task, we focus on recognizing and linking entities that are *physically* present in the image.[4]

The goal of learning for OVEN is to optimize a function $f_\Theta$ that predicts the entity $e$ from a given test example $x = (x^p, x^t)$ and the associated knowledge base of triples $\mathcal{K}$. There are different ways to utilize the information available in $\mathcal{K}$, and models may choose to use only a subset of this information. Figure 2 presents two typical ways of modeling OVEN. For encoder-decoder models [9, 55], the most straight-forward utilization is to memorize the entities of the database $\mathcal{K}$ into model parameters $\Theta$ via pre-training and fine-tuning, and then *generate* entity names directly during inference. Given that the generated name might not appear in the database, BM25 is used to map the prediction to the entity with the closet name in the available database For dual-encoder models [8, 17, 22, 35], an alternative is to explicitly compare a given test example $x$ to representations of entities $e \in \mathcal{E}$, making the prediction an *entity retrieval* problem. We refer to Section 4 for concrete examples of how to implement OVEN models.

**Data Split and Evaluation** Due to OVEN's goal of evaluating pre-trained multi-modal models, we only provide a partial set of visual concepts (*i.e*., SEEN categories) for model training or fine-tuning. For evaluation, an OVEN model is tested on generalization to entities not present in the fine-

---

[2]In this paper, we only consider using the name of the entity as its textual representation, despite the fact that more textual descriptions are available.

[3]https://en.wikipedia.org/wiki/File:Sabatia_campestris_Arkansas.jpg

[4]Extending this framework to entities that are not physically present in the image (e.g. the inventor of the airplane) is also valid and useful. See a follow-up works [10] for more details.

---

tuning data (thus UNSEEN), without forgetting the SEEN concepts. The models need to either acquire information from the knowledge base, or make a prediction using knowledge obtained during pretraining. We evaluate OVEN with a metric aiming to balance performance between SEEN and UNSEEN entities using a harmonic mean, as shown below:

$$\text{HM}(\text{ACC}_\text{SEEN}, \text{ACC}_\text{UNSEEN}) = 2 \ / \ \left(\frac{1}{\text{ACC}_\text{SEEN}} + \frac{1}{\text{ACC}_\text{UNSEEN}}\right) \ (1)$$

Harmonic mean equally weighs the importance of the SEEN and UNSEEN subsets, and penalizes models with a short barrel. Further details are provided in §3.

**OVEN versus recognition benchmarks** Given that an OVEN model need to generalize to UNSEEN entities, it is required to predict over all KB entities, which can exceed 6 million in our experiments (*e.g*., the size of English Wikipedia). This is orders of magnitude larger than existing benchmarks. Second, the large label space has made the generalization to UNSEEN entities the most critical criterion for a successful OVEN model, which also allows future open-domain evaluation[5]. Third, OVEN requires models to do multi-modal reasoning, *i.e*., comprehending the text query within its visual context, to predict the answer entity.
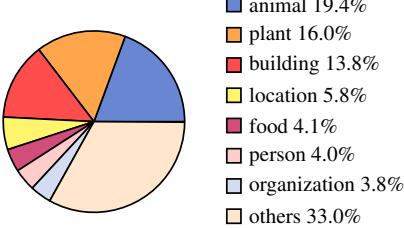
**OVEN versus Visual QA tasks** OVEN can be considered as a VQA task because its input format is the same as that of standard VQA models (*e.g*., text query + image). However, OVEN is specialized and focuses solely on recognition, with the text input serving mainly for intent disambiguation. Moreover, OVEN models are required to generate the name of an entity that exists in a given KB (like models for text entity linking tasks), while VQA models output free-form answers (such as yes/no for verification questions and numbers for counting questions).

**From OVEN to Knowledge-Intensive VQA** Although this paper aims to evaluate pre-trained multi-modal models on universal visual entity recognition, we highlight that models that excel at OVEN can serve as foundational components for systems that can answer knowledge-intensive questions. For example, given an image and a question "When was the church built?", one could apply an OVEN model to link the image to a concrete church's Wikipedia page and then extract the answer from that document. A follow-up work has conducted a thorough study on the value of Wikipedia grounding for answering knowledge-intensive visual questions [10].

## 3. The OVEN-Wiki dataset

Based on the task formulation of OVEN, we create the OVEN-Wiki dataset by combining 14 existing datasets, grounding their labels to Wikipedia, resolving label ambiguities, and providing unambiguous textual query intents for all

---

[5]One can collect and label a new set of entities from Wikipedia, to serve as a new evaluation data for OVEN models

| | Train Set | Val Set | Test Set | Human Set |
|---|---|---|---|---|
| # unique queries | 19,129 | 3,124 | 18,341 | 17,669 |
| # SEEN entities | 7,943 | 1,942 | 10,137 | 2,487 |
| # SEEN examples | 4,958,569 | 63,366 | 366,350 | 14,016 |
| # UNSEEN entities | 0 | 2,004 | 10,156 | 2,174 |
| # UNSEEN examples | 0 | 66,124 | 362,909 | 10,851 |
| # Total examples | 4,958,569 | 129,490 | 729,259 | 24,867 |

| | Wiki$_{\text{EN}}$ |
|---|---|
| # entities | 6,063,945 |
| # images | 2,032,340 |
| # title | 6,063,945 |
| `AvgLen`(title) | 2.93 |

Figure 3: Dataset Statistics of the OVEN-Wiki. **Left:** Distribution of super-categories of entities that have positive examples (See Appendix for more details). **Mid:** Statistics of different splits of the OVEN-Wiki. **Right:** Properties of the Wikipedia dump-`2022/10/01`.

examples. The 14 datasets were originally created for image recognition/retrieval, and visual question answering. Below is the complete list:

- **Image Recognition Datasets**: ImageNet21k-P [39, 41], iNaturalist2017 [51], Cars196 [24], SUN397 [58], Food101 [2], Sports100 [19], Aircraft [30], Oxford Flower [34], Google Landmarks v2 [56].
- **Visual QA Datasets**: VQA v2 [20], Visual7W [66], Visual Genome [25], OK-VQA [32], Text-VQA [48].

These datasets belong to two groups: image recognition (or retrieval) which provides *diverse visual entities*, defined as the **Entity Split** (ES); and VQA which provides *visually-situated natural language queries*, defined as the **Query split** (QS). For examples that originate from VQA datasets, we employ human annotators to write templated rules and filter out questions that do not lead to visual entity answers that are present in the image. For examples from recognition datasets, we first extract the super-category of their label (using the Wikipedia database), and then apply a templated query generation engine to generate a query with unambiguous intent that leads to the label (details in the Appendix).

**Label Disambiguation and Human Annotation** Grounding the labels of 14 datasets to Wikipedia entities is challenging, and we perform the following steps to accomplish this. We first apply a state-of-the-art textual entity linking system [13] to recognize text labels and map them into Wikipedia. Human annotators are used to write rules to detect bad linking results or unlinkable labels (e.g. numbers), and correct entity linking errors. The union of original dataset labels were linked to 20,549 unique Wikipedia entities, each with a number of examples for the purpose of training and evaluation. Meanwhile, we construct the candidate label space using the English Wikipedia snapshot from *Oct. 1 2022*, by removing all disambiguation, redirect, and media file pages. As shown in Figure 1 (right), this left us with 6,063,945 Wikipedia entities in total. Note that we only consider using the first Infobox images [57] from each page to serve as the visual support for each Wikipedia entity; these are available for 2,032,340 entities.

We further perform human annotation to create a high-quality evaluation dataset. Specifically, we hired over 30 dedicated annotators to validate the entity links in <image, query, answer> triplets sampled from the test split. They were asked to re-annotate the triplets with access to the visual context, ensuring that the query leads to the correct Wikipedia entity answer. Through this process, we collected 24,867 natural language queries, equally distributed over triplets originally sampled from the Entity and Query splits (*i.e.*, test splits). We asked the annotators to rewrite the queries so that no other object in the image could be a valid answer. As a result, the percentage of unique queries in the total examples (17,669 out of 24,867) as shown in Table 3 (mid) is significantly higher in the human set than in the other sets. This brings higher query generalization challenges for the human eval set. We report results using the same evaluation metrics on the human data, with respect to SEEN and UNSEEN entities. Figure 1 provides a glance at the human annotated data.

**Dataset Statistics** Figure 3 (left) presents the general distribution of the super-categories for our final collection of Wikipedia entities that have positive examples. Figure 3 (right) shows detailed statistics for queries and entities for each of the fine-tuning (train), validation, test, and human splits. Note that the models do not know which entities are present in the val/test/human set, and must scan through the whole KB to make predictions. The # of SEEN/UNSEEN examples indicates the # of examples of which the positive entity labels are in the SEEN/UNSEEN split.

**Evaluation Details** As aforementioned, we evaluate models by asking them to predict one out of over 6 million English Wikipedia entries. While our data does not cover all 6 million labels as positive examples, models still need to consider all possible outputs due to the presence of UNSEEN entities. We measure the models' performance using both the Entity Split (ES) and Query Split (QS). Specifically, we first compute the harmonic mean of accuracy over examples from the SEEN and UNSEEN classes, as $\text{Acc}_{\text{ES}} = \text{HM}(\text{Acc}_{\text{ESSEEN}}, \text{Acc}_{\text{ESUNSEEN}})$ and $\text{Acc}_{\text{QS}} = \text{HM}(\text{Acc}_{\text{QSSEEN}}, \text{Acc}_{\text{QSUNSEEN}})$ as the Equation 1. Then we further calculate the harmonic mean between splits $\text{HM}(\text{Acc}_{\text{ES}}, \text{Acc}_{\text{QS}})$ to reward models that do well on both splits. We use the validation data, which contains examples from subsets of both SEEN and UNSEEN entities, for model

selection, and we measure performance on the test split and the human evaluation set.

## 4. Fine-tuning Pre-trained Models for OVEN

We evaluate two prominent pre-trained multi-modal models: CLIP [35], a widely-used dual encoder model for image and text, and PaLI [9], a state-of-the-art pre-trained encoder-decoder model. Figure 2 has illustrated high-levelly on how encoder-decoder and dual encoder models can model the task of OVEN. In the following, we demonstrate with more details about how these two models can be fine-tuned for OVEN.

### 4.1. Dual encoders: CLIP and its variants for OVEN

One can naturally apply CLIP on OVEN by treating it as an image-to-text retrieval task. For an input image $x^p$, the image encoder is used to form an image embedding. Then the predicted entity could be retrieved by finding the entity that has the maximum dot product value between the entity text embeddings and entity image embeddings among the entire entity database. However, this naive implementation ignores the input intent $x^t$ and the entity images $p(e)$.

In the following, we present two variants of CLIPs: CLIP Fusion and CLIP2CLIP. The goal of these two variants is to use all of the information provided in the OVEN task. Both variants learn a function $f_\Theta$ that maximizes the score of the target entity for the given input image-query pair, using multimodal knowledge from the knowledge base. Given a test example $x = (x^p, x^t)$ and the knowledge base of triples $\mathcal{K}$, the function is used to make a prediction,

$$e' = \arg\max_{e \in \mathcal{E}} f_\Theta(x^p, x^t, p(e), t(e)) \quad (2)$$

**CLIP Fusion** adopts the pre-trained CLIP model as the featurizer to develop this system, via adding a 2-layer Multi-Modal Transformer on top of the CLIP image and text features as a mixed-modality encoder. The left encoder (for an input image-query pair) and the right encoder (for multi-modal knowledge information) use the same architecture, but do not share parameters. We fine-tune all of their parameters on the OVEN-Wiki, which includes both the pre-trained CLIP weights and randomly initialized Transformer weights.

**CLIP2CLIP** relies more heavily on the pre-trained CLIP model and introduces only a minimal set of new parameters (*i.e.*, four) to re-weigh and combine CLIP similarity scores. Particularly, it computes the cosine similarity between $<x^p, t(e)>$, $<x^t, p(e)>$, $<x^p, p(e)>$, and $<x^t, t(e)>$, using the image and text encoders of CLIP, respectively. Then it aggregates these similarities by multiplying them with a learnable vector that reflects importance weights.

**Scaling to 6 million candidates.** It is expensive to perform dot product scoring with respect to 6 million webpages on-the-fly. Fortunately, there exist approximate algorithms for maximum inner product search whose running time and storage space scale sub-linearly with the number of documents [38, 46, 47]. In all our experiments, we use ScaNN [21] as our library for entity retrieval.

### 4.2. Encoder-Decoder: PaLI for OVEN

PaLI [9] is a sequence-to-sequence model pre-trained on web text, image-text pairs (*i.e.*, WebLI) and other sources. PaLI can accept both an image and text as input and generates text as output. In order to map the PaLI predictions to the knowledge base, we run a BM25 [40] model to retrieve the most similar Wikipedia entity name for every generated text output. We found that this can slightly but consistently improve the entity recognition results. Note that we directly fine-tune PaLI on the OVEN training data, which does not cover all entities and questions appearing in our Dev and Test splits. However, we found that PaLI is still able to handle entities that are unseen during fine-tuning due to the knowledge acquired during pre-training. To make the comparison with CLIP more comprehensive, we report results on both PaLI-3B and PaLI-17B. The former PaLI variant is at the same magnitude (in its number of parameters) as the largest CLIP model, and the latter PaLI variant is one magnitude larger, and much stronger based on other evaluation [9].

## 5. Experiments

We first describe the essential experimental setups in §5.1, and then present the main benchmark results in §5.2.

### 5.1. Experimental Setups

**Pre-trained Model Details.** For all the CLIP variants, we employ the largest CLIP checkpoints, *i.e.*, ViT-L14, which leverages Vision Transformer [16, 52] as its visual backbone. For the PaLI model [9], we make use of the 3B and 17B parameter pre-trained models provided by the original authors, for fine-tuning on OVEN.

**Data Processing Details.** We process all images in our dataset by resizing them to 224×224, linearize them into a sequence of 14×14 patches, and apply the normalization technique consistent with each model's pretraining to pre-process the images. For natural language text, we perform tokenization based on the adopted pre-trained model's original vocabulary. For CLIP variants that encode Wikipedia images for entity retrieval, we apply the same image processing pipeline whenever the image is available. When the Wikipedia entity does not have an infobox image, we use a black image to represent the visual support.

### 5.2. Benchmark Results

---

[6]The human study is done on a random sampling of 100 examples.

|  | # Params | Entity Split(Dev) | | | Query Split(Dev) | | | Overall(Dev) |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
|  |  | SEEN | UNSEEN | HM | SEEN | UNSEEN | HM | HM |
| **Dual Encoders:** | | | | | | | | |
| ● CLIP$_{ViTL14}$ | 0.42B | 5.4 | 5.3 | 5.4 | 0.8 | 1.4 | 1.0 | 1.7 |
| ● CLIP Fusion$_{ViTL14}$ | 0.88B | 32.7 | 4.3 | 7.7 | 33.4 | 2.2 | 4.2 | 5.4 |
| ● CLIP2CLIP$_{ViTL14}$ | 0.86B | 12.6 | 10.1 | 11.2 | 4.1 | 2.1 | 2.8 | 4.4 |
| **Encoder Decoder:** | | | | | | | | |
| ◆ PaLI-3B | 3B | 21.6 | 6.6 | 10.1 | 33.2 | 14.7 | 20.4 | 13.5 |
| ◆ PaLI-17B | 17B | 30.6 | 12.4 | 17.6 | 44.2 | 22.4 | 29.8 | 22.1 |

Table 1: Comparison between the fine-tuned models on the OVEN-Wiki **validation** set.

|  | # Params | Entity Split(Test) | | Query Split(Test) | | Overall(Test) | Human Eval | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
|  |  | SEEN | UNSEEN | SEEN | UNSEEN | HM | SEEN | UNSEEN | HM |
| **Dual Encoders:** | | | | | | | | | |
| ● CLIP$_{ViTL14}$ | 0.42B | 5.6 | 4.9 | 1.3 | 2.0 | 2.4 | 4.6 | 6.0 | 5.2 |
| ● CLIP Fusion$_{ViTL14}$ | 0.88B | 33.6 | 4.8 | 25.8 | 1.4 | 4.1 | 18.0 | 2.9 | 5.0 |
| ● CLIP2CLIP$_{ViTL14}$ | 0.86B | 12.6 | 10.5 | 3.8 | 3.2 | 5.3 | 14.0 | 11.1 | 12.4 |
| **Encoder Decoder:** | | | | | | | | | |
| ◆ PaLI-3B | 3B | 19.1 | 6.0 | 27.4 | 12.0 | 11.8 | 30.5 | 15.8 | 20.8 |
| ◆ PaLI-17B | 17B | 28.3 | 11.2 | 36.2 | 21.7 | 20.2 | 40.3 | 26.0 | 31.6 |
| **Human+Search** [6] | - | - | - | - | - | - | 76.1 | 79.3 | 77.7 |

Table 2: Results of methods on the OVEN-Wiki **test** set and **human evaluation** set. Human+Search represents human performances with information retrieval tools such as search engines and others, on a random subset of OVEN-Wiki$_{Human\_Eval}$.

**Main Results** Results on the validation set are presented in Table 1, and include performance on the Entity and Query splits, as well as the overall combined scores.

There are several interesting (perhaps surprising) observations from Table 1. First, while CLIP variants such as CLIP Fusion and CLIP2CLIP are utilizing more information from Wikipedia (*i.e.*, entity names and entity images), they are weaker than the auto-regressive PaLI-3B and PaLI-17B model, across most evaluation data splits. This suggests that high-capacity generative multi-modal pre-trained models are capable of recognizing visual entities. Second, this performance gap is more apparent on the query split than the entity split, potentially due to the VQ2A pre-training objectives [7] and the underlying powerful language models [36] employed by the PaLI model.

Comparing all CLIP-based models, we observe that CLIP Fusion and CLIP2CLIP, which uses all Wikipedia information are generally performing better than the vanilla CLIP model, showcasing the benefits of multimodal information from Wikipedia. Meanwhile, we also observe that CLIP Fusion, where two new layers have been added on top of pretrained CLIP, shows very strong results on SEEN entities for both the Entity and the Query splits, but weak results on UNSEEN entities, thus leading to lower overall performance.

The CLIP2CLIP model, on the other hand, is capable of retaining the cross-entity generalization performance while improving its prediction accuracy on SEEN entities.

Comparing the PaLI models, we observe a drastic improvement as the number of parameters in the models increased. Particularly, PaLI-17B has a double-digit performance gain in the overall performances, against the PaLI-3B model. This suggests that scaling the capacity of the model is one of the most important factors, and should be considered as a top priority in future multi-modal dual encoder research.

**Results on Human Set and Human Performance.** Table 2 shows that the results on the test set and human set are generally aligned with observations on the validation set. We conduct a study to estimate the human performance on OVEN-Wiki, via requesting 3 dedicated human annotators to answer 100 examples (sampled from human evaluation set, answers are non-overlapping). We allow the annotators to use search engines (*e.g.*, Google Image Search, Wikipedia Search, etc.)[7], as long as the annotators can provide a valid Wikipedia entity name as the answer. As a result of this study, human achieves 77.7% harmonic mean accuracy, which is

---

[7]Even with search engines, each annotator has used 254 seconds to complete one example.

significantly higher than the best comparison systems shown in Table 2.

# 6. Analysis

In this section, we perform empirical studies to analyze the pre-trained CLIP2CLIP and PaLI models, and conduct a detailed analysis of these two models' common errors.

**Does fine-tuning always help generalization?** Figure 4 presents the validation scores of the PaLI model (left) and the CLIP2CLIP model (right), during fine-tuning on OVEN-Wiki's training split. It shows that a longer training schedule does not lead to better generalization performance, particularly when evaluated on the UNSEEN entities. Because of this, we employ the early stopping strategy for model selection, and pick the model with the best harmonic mean combined score on the validation set. However, due to this early stopping strategy, both fine-tuned models are not utilizing 100% of the examples in OVEN's training data because their UNSEEN performance starts to degenerate within one epoch. This has indicated that more advanced fine-tuning strategies that use better regularization techniques to encourage generalization across Wikipedia entities, could be a promising research to explore in the future.

**How would the number of entities in KB influence the model's prediction?** Figure 5 presents the accuracy of CLIP2CLIP, as a function of the # of total candidates to retrieve from. Here, we compute the accuracy by sub-sampling the negative candidates from KB to different sizes. We observe that when the retrieval candidate entities are only the positive entities (with the # of candidates being 20K), the performance of the CLIP2CLIP model is significantly higher than the open-domain setting (with 6M entities in total). Beyond this, as the KB size increases, model accuracy decreases. Concretely, it shows an approximately linear decline along the log-scale x-axis in Figure 5. This indicates that as the KB size increases, the models' accuracy first drops significantly and then follows with a gradual decline. On the other hand, PaLI's performance is generally more steady as the size of KB grows, potentially because its prediction has already matched up entity names inside KB, so narrowing down the set of candidates does not help the BM25 post-processing. One potential direction is to employ constrained decoding for the PaLI-based model, which we leave for future works.

**How would models perform on head vs. tail entities?** We evaluate the visual entity recognition performances of CLIP2CLIP and PaLI, on entities of different popularity. Particularly, Figure 6 presents a histogram according to models' performance on the entity that has different average monthly Wikipedia page views in 2022 [31]. From the comparison, we can see that PaLI is significantly more accurate compared to CLIP2CLIP, on the head entities (that have more than 5K

| | PALI-17B | CLIP2CLIP |
|---|---|---|
| CORRECT | 29% | 15% |
| IN-CORRECT | 71% | 85% |
| → (A) WRONG BUT RELEVANT | 23% | 27% |
| → (B) TOO GENERIC | 15% | 1% |
| → (C) MISUNDERSTAND QUERY | 7% | 37% |
| → (D) MISCELLANEOUS | 24% | 20% |

Table 3: Error type distribution for difference models. PaLI predicts more answers with less granularity (less granularity), while most of the CLIP errors are due to not understanding the questions.

monthly page views). However, we observe that CLIP2CLIP can perform on par or even outperform PaLI on tail-ish entities (that have less than 2.5K monthly views). This suggests that the retrieval-based visual entity recognition model has its own advantages, in recognizing the difficult and tail entities. Meanwhile, this result also provides a hint that potentially a frequency calibrated evaluation should be developed to reward models more with strong recognition capability on the tail entities.

**Error analysis** To better understand the errors that CLIP2CLIP and PaLI models are making, we sampled a random 100 examples on the human evaluation set, and manually categorize and analyze the errors that PaLI and CLIP2CLIP are making. Particularly, we categorize the errors of the pre-trained models into four categories: (a) erroneous but relevant prediction, on concepts of the same granularity; (b) errors due to predicting very generic concepts; (c) errors due to misunderstanding the intent behind the query. (d) other miscellaneous errors. Note that errors type (d) are mostly mistakes that are unrelated and not easily interpretable. The results are shown in Table 3. Figure 4 has provided some concrete examples of the above types of mistakes made by CLIP2CLIP and PaLI. Interestingly, it shows that the two models, *i.e.*, CLIP2CLIP and PaLI, are making very different types of errors in their predictions. Particularly, CLIP based model is good at capturing the right granularity of the entity, but often fails to understand the true intent of the text query. For instance, Figure 4 (c) shows that CLIP2CLIP ignores the text query and starts to predict the name of the barrel racer. In contrast, PaLI is good at following the text query, but can usually predict generic concepts when it does not know the answer confidently (see Figure 4 (b)).

# 7. Related Works

**Learning to Recognize UNSEEN Categories** There has been a significant amount of prior work [26, 28, 53] focusing on the generalization situation where information of
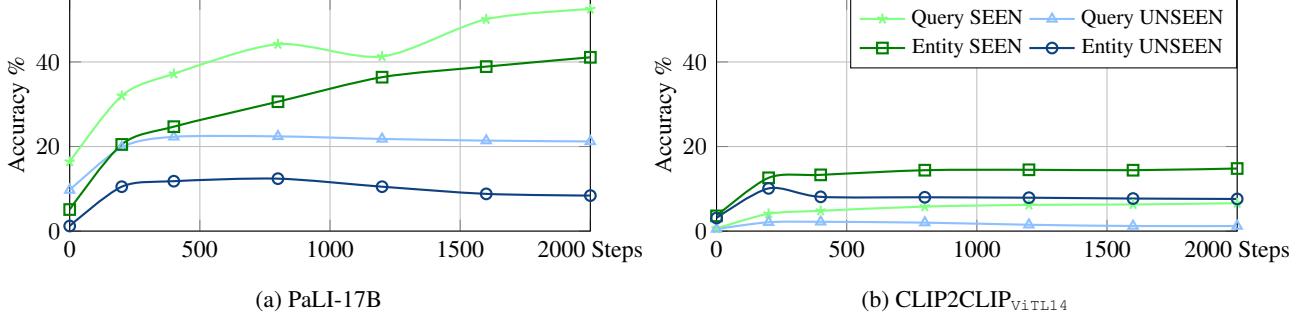
**(a) PaLI-17B**



**(b) CLIP2CLIP$_{\text{ViTL14}}$**

Figure 4: **Fine-tuning PaLI or CLIP2CLIP for large # of steps** increases the SEEN entity accuracy but hurts the UNSEEN entity accuracy.



**(a) Query Split**
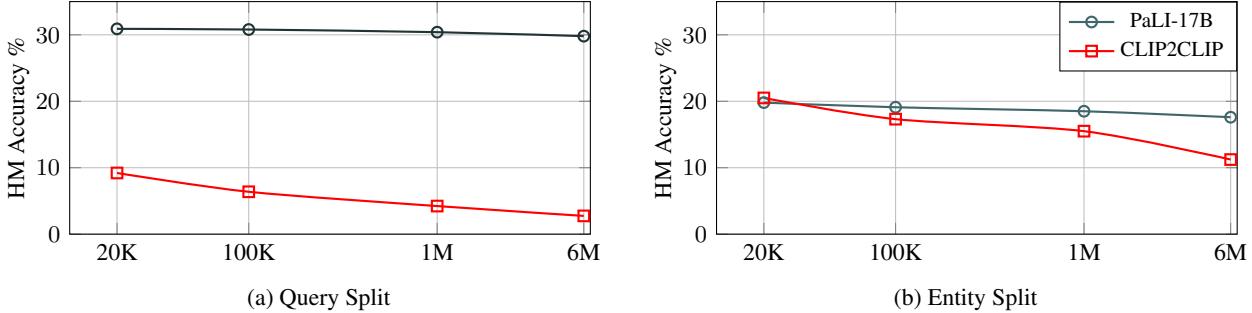


**(b) Entity Split**

Figure 5: **Impact of # Wikipedia Candidates on PaLI and CLIP2CLIP.** Increasing the size of Wikipedia makes the tasks difficult.



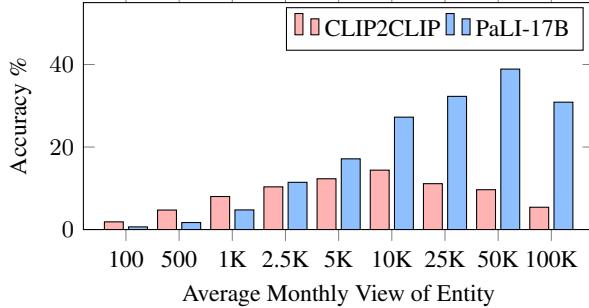Figure 6: **Comparison of Performances on Head vs. Tail Entities (results on Validation set).** PaLI wins over CLIP2CLIP on popular (*i.e.*, high monthly page view) Wikipedia entities, but loses on rare (*i.e.*, low monthly page view) Wikipedia entities.

novel categories are presented at test time. Zero-shot learning (ZSL) is one of such attempts that tackles learning new categories with zero images for training. To achieve such transfer, ZSL methods typically rely generating UNSEEN image classifiers based on corresponding semantic representations, in the format of manually labeled attributes [26], unsupervised learned word vectors [6], or pre-trained sentence embeddings [23, 35]. Few-shot learning (FSL) [53] proposes a more realistic setup, where learners have access to a limited number of visual exemplars during the model deployment. With this goal, FSL methods aim to extract the inductive bias of learning from the SEEN classes, such that the model can leverage it in learning the UNSEEN classes, to avoid severe over-fitting. Particularly, prior works either use adapted non-parametric classifiers [42, 49, 61], or meta-optimized linear classifiers [18, 37] to incorporate the few-shot UNSEEN support examples. Comparing to them, our proposed task exposes different challenges as we ask the model to make the best use of open-world Web knowledge (*i.e.*, Wikipedia pages with images & figures), which contains textual semantic information and visual appearance of the entities in the open world.

**Vision and Language + Knowledge** There have been efforts in combining knowledge into vision and language tasks, such as Visual QA [5, 11, 32, 44] and entity-focused image captioning [1, 27]. Among them, knowledge-based VQA is most related to OVEN, but also differs in many aspects. Specifically, [5] presents a text QA dataset that requires understanding multi-modal knowledge in a KB. [44] propose to perform knowledge-based question answer tasks, centered around questions that resolve relational query over public Figures. Meanwhile, [32] propose to answer questions where the answer is outside of the image context, to assess model's capability in understanding real-world knowledge More recently, [11] studies the zero-shot visual QA setting where
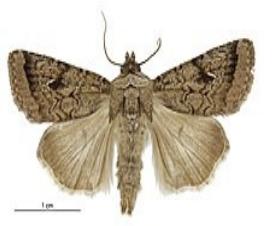
| Error Type | (a) Wrong but Relevant | (b) Too Generic | (c) Misunderstand Query |
|---|---|---|---|
| Input Query | *What is the name of the model of this aircraft?* | *What is the species of this animal?* | *What sports event is displayed in the picture?* |
| Input Image | | | |

| | | | |
|---|---|---|---|
| PaLI-17B: | WikiID: Q589498<br>Name: *BAe 146* | WikiID: Q255496<br>Name: *Butterfly* | WikiID: Q2529836<br>Name: *Barrel racing* |
| CLIP2CLIP: | WikiID: Q937949<br>Name: *Dornier 328* | WikiID: Q13510645<br>Name: *Proteuxoa comma* | WikiID: Q****4678<br>Name: *E. W. (barrel racer)†* |
| Ground-Truth: | WikiID: Q218637<br>Name: *ATR 42* | WikiID: Q592001<br>Name: *Hoary comma* | WikiID: Q2529836<br>Name: *Barrel racing* |

Table 4: **Visualization of mistakes made by the CLIP2CLIP and PaLI-17B Model.** We visualize the Wikipedia infobox images for each of model's predictions, to provide more context about the visual similarity between the prediction/ground-truth and the input image. Correct predictions are marked as green, whereas incorrect predictions are marked as red. (†: Since no infobox image is available for this Wikipedia entity, a face-anonymized Web image of the entity is visualized for reference.)

some answers (out of a total of 500 frequent answers of general concepts) are unseen during the training, where a KB is supplied to assist the model in answering unseen answers. Comparing to them, OVEN steps back to the more fundamental problem of establishing the link between visual content and entity in the KB, but at a larger scale and broader coverage. We believe that stronger models developed on OVEN would benefit such knowledge-intensive visual QA tasks.

**Entity Linking** Entity linking (EL) is the task of grounding entity mentions in the text by linking them to entries in a given knowledge base. Supervised EL [33] has demonstrated

its strong performance when all entities are in-distribution during the evaluation. Because KB is updating all the time, recent works [3, 13, 15, 29, 64] focus on a more realistic setting where entity linking needs to be achieved in the zero-shot, with a large portion of entities (to be evaluated) completely unseen during the training. OVEN is a visual analog of zero-shot EL, and targets at developing generalizable models that recognize entities unseen in the training. Among all EL literature, visually assisted EL [65] is most relevant to this work, whose goal is to use the associated image of text to improve the precision of text EL. OVEN is different as its text queries do not mention the name of the entities, which put visual understanding and reasoning into the central position.

## 8. Discussion

In this paper, we have introduced OVEN, a task that aims to unambiguously link visual content to the corresponding entities in a web-scale knowledge base (*i.e.*, Wikipedia), covering a total of more than 6 millions of entities. To facilitate the evaluation of OVEN, we created the OVEN-Wiki dataset, via combining and re-annotating 14 existing visual recognition, retrieval, and visual QA datasets, and linked over 20K labels to the Wikipedia entities. With OVEN-Wiki, we evaluate state-of-the-art multi-modal pre-trained models, *i.e.*, the CLIP [35]-based entity retrieval models and the PaLI [9]-based entity generation model, via fine-tuning them for the OVEN task, to examine their capability on recognizing open-domain visual concepts. As a result, PaLI models have presented significantly stronger performances than the CLIP variants, even on unseen visual entities during the fine-tuning. Meanwhile, although the CLIP-based entity retrieval model is overall weaker, it shows advantages in recognizing the tail visual entities.

One additional nice property of OVEN-Wiki is its strong extensibility. As a result of grounding of all recognition labels to Wikipedia entities, we as a community can keep growing the member recognition datasets of OVEN-Wiki, by adding positive instances to Wikipedia entities that do not have examples by far. Moreover, successful OVEN models can generalize to recognize emerging entities (*e.g.*, `iPhone 14 Pro`), as long as the corresponding Wikipedia page is created. In summary, we hope OVEN will drive future research on knowledge-infused multimodal representation learning via visual entity recognition.

## Ethics Statement

As our dataset, *i.e.*, OVEN-Wiki, is composed of existing image recognition, image retrieval, and visual question answering datasets, we have introduced minimum risk of exposing additional social bias in our data. However, OVEN-Wiki is still at the risk of inheriting existing dataset biases. As a re-

sult, we employed existing data curation strategies [60] to reduce such potential risks. Besides such risk, OVEN-Wiki also opens up new possibilities that can alleviate ethical concerns in AI systems. Specifically, OVEN-Wiki is a dataset that targets advancing research for establishing stronger grounding between the visual content and knowledge base, which can potentially contribute to building more attributed visual systems, such as a visual question answering model that produces answers based on the linked Wikipedia page, with improved interpretability and controllability.

## Acknowledgement

## References

[1] Ali Furkan Biten, Lluis Gomez, Marçal Rusinol, and Dimosthenis Karatzas. Good news, everyone! context driven entity-aware captioning for news images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12466–12475, 2019. 8

[2] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *ECCV*, 2014. 4

[3] Jan A Botha, Zifei Shan, and Dan Gillick. Entity linking in 100 languages. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7833–7845, 2020. 10

[4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 1

[5] Yingshan Chang, Mridu Narang, Hisami Suzuki, Guihong Cao, Jianfeng Gao, and Yonatan Bisk. Webqa: Multihop and multimodal qa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16495–16504, 2022. 8

[6] Soravit Changpinyo, Wei-Lun Chao, Boqing Gong, and Fei Sha. Synthesized classifiers for zero-shot learning. In *CVPR*, pages 5327–5336, 2016. 8

[7] Soravit Changpinyo, Doron Kukliansky, Idan Szpektor, Xi Chen, Nan Ding, and Radu Soricut. All you may need for vqa are image captions. *arXiv preprint arXiv:2205.01883*, 2022. 6

[8] Jiacheng Chen, Hexiang Hu, Hao Wu, Yuning Jiang, and Changhu Wang. Learning the best pooling strategy for visual semantic embedding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15789–15798, 2021. 3

[9] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni,

Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022. 1, 2, 3, 5, 10, 15

[10] Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang. Can pre-trained vision and language models answer visual information-seeking questions? *Technical report*, 2023. 3

[11] Zhuo Chen, Jiaoyan Chen, Yuxia Geng, Jeff Z Pan, Zonggang Yuan, and Huajun Chen. Zero-shot visual question answering using knowledge graph. In *International Semantic Web Conference*, pages 146–162. Springer, 2021. 8

[12] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 1

[13] Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. Autoregressive entity retrieval. *arXiv preprint arXiv:2010.00904*, 2020. 2, 4, 10

[14] Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. Autoregressive entity retrieval. In *ICLR*, 2021. 13

[15] Nicola De Cao, Ledell Wu, Kashyap Popat, Mikel Artetxe, Naman Goyal, Mikhail Plekhanov, Luke Zettlemoyer, Nicola Cancedda, Sebastian Riedel, and Fabio Petroni. Multilingual autoregressive entity linking. *Transactions of the Association for Computational Linguistics*, 10:274–290, 2022. 10

[16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 5

[17] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. Vse++: Improved visual-semantic embeddings. In *BMVC*, 2017. 3

[18] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017. 8

[19] Gerry. Sports100: 100 sports image classification. https://www.kaggle.com/datasets/gpiosenka/sports-classification/metadata, 2021. Accessed: 2022-09-26. 4

[20] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 4

[21] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. Accelerating large-scale inference with anisotropic vector quantization. In *International Conference on Machine Learning*, 2020. 5

[22] 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. 3

[23] Jihyung Kil and Wei-Lun Chao. Revisiting document representations for large-scale zero-shot learning. In *EMNLP*, 2021. 8

[24] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *3dRR-13*, 2013. 1, 4

[25] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017. 4

[26] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot visual object categorization. *TPAMI*, 36(3):453–465, 2014. 7, 8

[27] Fuxiao Liu, Yinghan Wang, Tianlu Wang, and Vicente Ordonez. Visual news: Benchmark and challenges in news image captioning. *arXiv preprint arXiv:2010.03743*, 2020. 8

[28] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *CVPR*, 2019. 7

[29] Lajanugen Logeswaran, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, Jacob Devlin, and Honglak Lee. Zero-shot entity linking by reading entity descriptions. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3449–3460, 2019. 10

[30] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013. 4

[31] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Hannaneh Hajishirzi, and Daniel Khashabi. When not to trust language models: Investigating effectiveness and limitations of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511*, 2022. 7

[32] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019. 4, 8

[33] David Milne and Ian H Witten. Learning to link with wikipedia. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pages 509–518, 2008. 9

[34] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*. IEEE, 2008. 4

[35] 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. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 1, 2, 3, 5, 8, 10, 15

[36] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 6

[37] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of maml. 2020. 8

[38] Parikshit Ram and Alexander G Gray. Maximum inner-

product search using cone trees. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 931–939, 2012. 5

[39] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. In *NeurIPS*, 2021. 1, 4

[40] Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009. 5

[41] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 1, 4

[42] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *ICLR*, 2019. 8

[43] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020. 1

[44] Sanket Shah, Anand Mishra, Naganand Yadati, and Partha Pratim Talukdar. Kvqa: Knowledge-aware visual question answering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 8876–8884, 2019. 8

[45] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4596–4604. PMLR, 2018. 15

[46] Fumin Shen, Wei Liu, Shaoting Zhang, Yang Yang, and Heng Tao Shen. Learning binary codes for maximum inner product search. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4148–4156, 2015. 5

[47] Anshumali Shrivastava and Ping Li. Asymmetric lsh (alsh) for sublinear time maximum inner product search (mips). In *Advances in Neural Information Processing Systems*, pages 2321–2329, 2014. 5

[48] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *CVPR*, 2019. 4

[49] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, 2017. 8

[50] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*, 2022. 1

[51] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2018. 4

[52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5

[53] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.

7, 8

[54] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018. 1

[55] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. SimVLM: Simple visual language model pretraining with weak supervision. In *ICLR*, 2021. 3

[56] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval. In *CVPR*, 2020. 4

[57] Inc. Wikipedia Foundation. Help:infobox picture. https://en.wikipedia.org/wiki/Help:Infobox_picture. 4

[58] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010. 1, 4

[59] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*, 2020. 15

[60] Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, and Olga Russakovsky. Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the imagenet hierarchy. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pages 547–558, 2020. 10, 13

[61] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *CVPR*, 2020. 8

[62] Shuhei Yokoo, Kohei Ozaki, Edgar Simo-Serra, and Satoshi Iizuka. Two-stage discriminative re-ranking for large-scale landmark retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 1012–1013, 2020. 13

[63] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12104–12113, 2022. 15, 16

[64] Wenzheng Zhang, Wenyue Hua, and Karl Stratos. Entqa: Entity linking as question answering. In *International Conference on Learning Representations*, 2021. 10

[65] Qiushuo Zheng, Hao Wen, Meng Wang, and Guilin Qi. Visual entity linking via multi-modal learning. *Data Intelligence*, 4(1):1–19, 2022. 10

[66] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016. 4

# 9. Appendix

# 10. Dataset Construction, Annotation, and Additional Statistics

In this section, we describes the complete details on data collection, curation, entity linking, and show additional statistics of the processed dataset (§10.1). Then we also discuss how we train annotators to annotate our task, and provide the concrete annotation interface(§10.2).

## 10.1. Data Collection & Pre-processing

**Data Filtering** Some of our member datasets have been reported to include non-imageable classes, classes with undesired social bias [60], or non-entity classes (*e.g.*, numbers). Therefore, we apply a filtering process to compose our dataset, based on the individual condition of each source dataset. Overall, to create the Entity split, we first apply a general safety filter [60] to remove non-imageable labels, non-entity labels, and labels with social bias. To create the Query split, we employed three expert annotators to write heuristic policies to filter each VQA dataset, and ensure our task is focusing on entity related questions. Concretely, questions related to counting, verification, or querying non-entity attributes (*e.g.*, dates), are removed. Then we apply the same safety filter.

**Linking labels to Wikipedia Entities** Based on the filtered data, we developed a two-staged entity linking strategy to connect the label text to Wikipedia entities, on both Entity and Query splits. First, we obtain exact match based entity candidates by querying the Wikipedia search API (with the auto-suggestion disabled) with the raw label text. We reject candidates whose landing pages are identified as disambiguation pages. The Wikipedia API[8] automatically redirects queries (in our case, labels) matching entity aliases to their canonical form. For the labels which do not have an exact match in Wikipedia, we use a state-of-the-art text-based entity linker (*i.e.*, GENRE [14]) to obtain top candidate Wikipedia entity names. Finally, we link the label to the top ranked entity whose landing page is not a disambiguation page.

**Preparing Multi-Modal Knowledge** Using the entity linking process described earlier, we successfully connect a total of 24,895 class labels in OVEN-Wiki to corresponding Wikipedia entities. Overall, our dataset contains 20,801 unique entities. For the Entity split data, we generate a synthetic text query based on the super-category information of the label (either provided by source dataset or mined from Wikidata[9]), using templated language. For example, iNaturalist has provided detailed supercategory annotation on each class, such as `Plantae`, `Reptilia`, *etc*. For dataset that

---

do not provide this information, we use the super-category mined from Wikidata, which is publiclly crowd sourced and maintained. As a result, our templated query generator produces the query ``what is the species of the plant in this image?'' for the entity ``Eryngium alpinum'', whose super-category is `Plantae`. Due to space limit, we provide more explanation in Appendix. For all Wikipedia entities, we use the corresponding Wikipedia page and its associating multi-media content (*e.g.*, information box images, *etc*.) as the source of *multi-modal knowledge* about entities.

**Statistics on Entities** Specifically, Figure 7 shows the number of unique entities in both the Entity and Query splits, where we compare the total number of entities in each source dataset against its original population (after applied safety filter). Note that for the Google Landmarks v2 (Gldv2) dataset, we employed the cleaned data split from [62], where the total number of unique entities is significantly reduced. Because Gldv2 is automatically generated and has reported to contain noises particularly with tail entities [62], we removed entities with less than 50 instances for a improved precision (further reduces the # of entities in Gldv2 to ∼6k).

**Entity Super-Categories** To give more details for the Figure 3 in the main text, we further present full super-category grouping information in Figure 8. As aforementioned, we have combined entities that belongs to general groups (*e.g.*, "object", "item" groups) or unpopular groups (*e.g.*, groups with less than 5 entities) into the "others" group. We also merged some sub-categories into super-categories, *e.g.*, "location"+"park"+"lake"+"river"+"mountain"→"location", "building"+"bridge"→"building".

## 10.2. Human Annotation Procedure & Interface

In order to verify the quality of OVEN-Wiki and to provide a human verified test set to evaluate on, we conduct human annotation on a subset of test set. The annotators are asked to correct the errors in the ¡image, query, answer¿ triplets. The details are as follows.
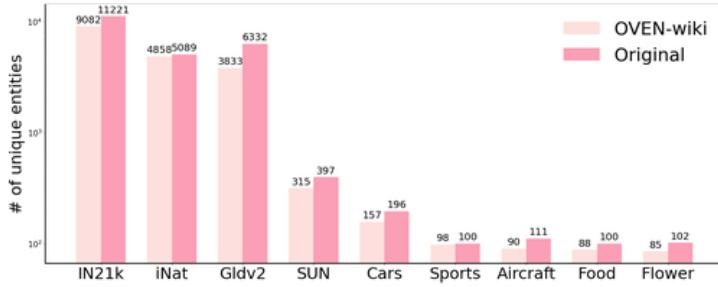
**Annotation interface** Figure 9 illustrates the annotation interface. The left side of Figure 9 are the input to the annotators which includes the original question, image and the answer (together with the wikipedia hyperlink). The annotators are asked to complete the following questions:

1. *Does the Wikipedia represent the correct meaning of the answer? Provide the Wikipedia link if not.*

   This question requires the annotators to correct the entity linking errors. The annotators use Google search to find the most suitable Wikipedia link if the provided one is not adequate. In our dataset, 8.4% of the entity links are reported wrong by more than 2 annotators, which are manually corrected later.

---

| | # Original Answers | # Entity Answers |
|---|---|---|
| VG | 50,130 | 3,460 |
| OK-VQA | 4,214 | 1,600 |
| Text-VQA | 19,500 | 3,562 |
| VQA v2 | 26,748 | 4,337 |
| Visual7W | 7,588 | 1,945 |

Figure 7: Number of unique entities on Entity split (left) and Query split (right). We compare it against the # of entities before applying pre-processing. Note that VQA datasets contain massive non-entity answers, or collapsed answers, which leads to a large reduction in numbers after pre-processing.
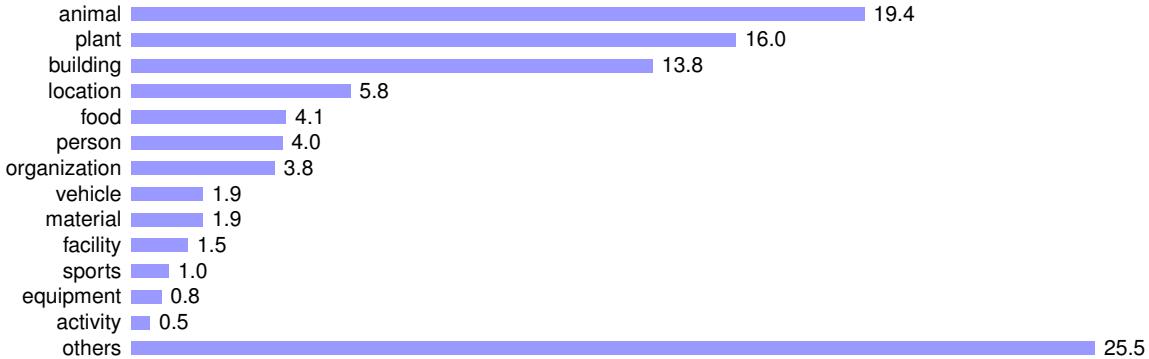


Figure 8: Distribution of the entities in our datasets (Grouped by their super category).

2. *Is the Wikipedia answer physically present in the image.*

This question is mainly aimed at filtering out the OCR examples which are out of our scope. One example is that the image about a wall painted with the word "love" and the linked entity is the "love" Wikipedia. In our dataset, 10.3% of the answers are reported not physically present in the image by more than 2 annotators, which are discarded from the human evaluation set.

3. *Rewrite the question so that no other object can be the answer.*

The annotators will rewrite the question is the answer is wrong or ambiguous. Annotators will make sure that the question can not be answered without the image and that the answers can not be included in the rewritten questions. In our annotation, 99.9% of the questions are being rewritten.

**Instruction and Training**   We carefully design the training procedure to improve the annotation quality. We first conduct a "Self-study session" where the annotators will read the instructions and annotate a few toy examples. Then we conduct a "In-person tutorial" where we have an online video session in which we walk annotators through the full

version of the instructions and discuss mistakes made in the self-study annotations. Finally we conduct a "Test exam" and the qualified annotators are accepted. In total, 30 annotators went through our training procedure and all of them were eventually accepted to work full-time on the main task.

**Quality control**   We have a three way annotations where each examples are annotated by three annotators. We were giving regular feedback on the questions the annotators may have during the annotation and pointed out mistakes identified in annotators' past answers.

On average, it took annotators 4.6 minutes to answer each question with the time consumption slightly decreasing as annotators get familiar with the task. The compensation rate for the task was set to be $17.8/hour which is higher than the minimum hourly wage in the US.

We filtered out all the examples where the wikipedia links are marked as wrong or the Wikipedia answers are marked as "Not physically present in the image".

## 11. Implementation Details of the baseline systems

In this section, we provide implementation details on the CLIP variants and PaLI model for the OVEN task.
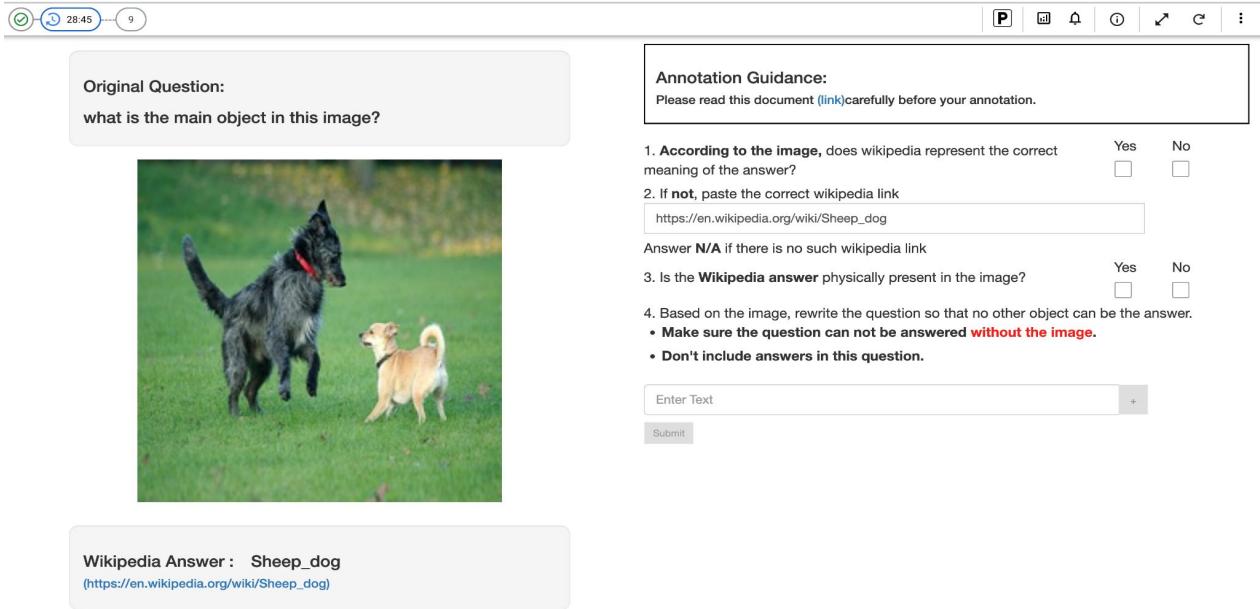
Figure 9: Annotation inferface

## 11.1. CLIP Fusion Model

As aforementioned, we implemented this multi-modal dual encoder via taking pre-trained CLIP image and text encoders as featurizers. The CLIP model is based on a ViT-Large, with a total of over 400ᴍ parameters, pre-trained on a 400ᴍ prviate image-text dataset collected by OpenAI. Based on this model, we build two 2-layer Transformer models, on top of two CLIP models as the left and right encoder, for encoding the query representation and the entity representation, respectively. The 2-layer Transformers follows the same architecture as T5 Transformer [35], but with 2 layers, 12 attention heads, with each attention head of 64 dimensions, and the embedding size of 768. We then fine-tune this composed model on the Oᴠᴇɴ-Wiki's training data, using a in-batch contrastive learning objective [35], with a batch size of 4,096. We optimize the model for 10K steps in the fine-tuning stage, with Adafactor optimizer [45] and a initial learning rate of 0.001. There are 1k steps for the warmup, followed by a square root LR decay schedule with final learning rate of 1e-6.

## 11.2. CLIP2CLIP Model

Different from CLIP Fusion, CLIP2CLIP is a model that adds minimum new parameters to the pre-trained CLIP encoders. Same as other models, we initialize both the query encoder and the target encoder separately with the pre-trained CLIP model. Specifically, we use the pre-trained CLIP encoders for both left and right encoders, to encode the image and text modality for both the query representation

and the entity representation. We then compute the four dot product similarity scores on the <input image, target text>, <input text, target image>, <input image, target image>, and <input text, target text> pairs, which is then combined via a learnable similarity weights into one logit score. The make sure that the learnable similarity weights is initialized properly, we perform a grid search to find a roughly good similarity weights for the CLIP2CLIP model (using Oᴠᴇɴ-Wiki's training data). Then we took this similarity weights to initialize the CLIP2CLIP model and fine-tune all parameters on Oᴠᴇɴ-Wiki's training set, under the same contrastive learning objective. Different from other models, given that this model has most of its parameters pre-trained, we realized that it works the best to early stop the model. As a result, we only fine-tune this model for 2k steps, with an initial learning rate of 1e-4, and a square root LR decay schedule with final learning rate of 1e-6.

## 11.3. PaLI Model

As aforementioned, we have evaluated two variants of PaLI models, the model with 3B total parameters (*i.e.*, PaLI-3B) and the model with 17B parameters (*i.e.*, PaLI-17B). The PaLI-17B model reuses 13B parameter from the mT5-XXL [59] and 4B parameters from the ViT-e [63], which were pre-trained Web Text and JFT-3B datasets, and then jointly trained on the WebLI [9] dataset with 10ʙ image and text pairs, under a variety of pre-training objectives, including object recognition, split captioning, visual question answering, etc. Similarly, the PaLI-3B model reuses 1B parameters from mT5-Large [59], and 1.8B pa-

rameters from the ViT-G [63], under the same pre-training recipe. To fine-tune PaLI on our dataset, we finetue the pre-trained PaLI model using its Visual QA interface, and inject the OVEN text queries into the PaLI's VQA prompt. As a concrete example, we convert a original query of `what species is the animal in the image?` into the format of `Answer in en:  what species is the animal in the image?`, as input to the PaLI model. The objective of fine-tuning process is then to maximize the likelihood of answer generation, same as its standard VQA fine-tuning practices. Similarly, we employ the Adafactor optimizer to optimize the fine-tuning, with a total of 2K fine-tuning steps, with a warmup of 1K steps and linear LR decay schedule.