Supplementary Materials of Grounding Language Models for Visual Entity Recognition

Zilin Xiao¹, Ming Gong², Paola Cascante-Bonilla¹, Xingyao Zhang², Jie Wu², and Vicente Ordonez¹

¹ Department of Computer Science, Rice University, USA {zilin, pc51, vicenteor}@rice.edu ² Microsoft STCA, China {migon, xingyaozhang, jiewu1}@microsoft.com

A Experimental Details

Table 1: Hyperparameter settings for training of AUTOVER	Table 1:	Hyperparameter	settings for	r training of	Autover
---	----------	----------------	--------------	---------------	---------

Hyperparameter	Value
learning rate	2e-5
$weight_decay$	0
batch size per device	8
effective batch size	256
learning rate strategy	Cosine
optimizer	AdamW
gradient clipping	Disabled

We present the hyperparameter choice in Table 1. As the effectiveness of contrastive learning only reveals in a large batch size setting, we conduct all pieces of training in a batch size of 256 on 32 V100-SXM2-32GB GPUs. Utilizing DeepSpeed [2] for distributed training management, we train AUTOVER-7B for 2 days and AUTOVER-13B for 6.5 days with 2.5 million training samples for one epoch.

B Hard Negative Groups Construction Details

We first illustrate the process of constructing KB-HARD group construction. Specifically, we consider three types of relations in Wikidata for category hierarchy construction, "instance of (P31)", "parent taxon (P171)" and "subclass of (P279)". To obtain category hierarchy labels, we execute the SPARQL statement shown on the left of Figure 2 on Wikidata knowledge graph. To exclude generic categories, we limit the depth of traversal on the knowledge graph to 3, and manually filter generic categories³. On the right of Figure 2, we visualize the 3-hop

³ ["Wikidata metaclass", "second-order class", "taxon", "classification scheme", "metaclass", "third-order class", "human", "direct anatomical metaclass", "organisms known

2 Z. Xiao et al.



Fig. 1: Error analysis of two query image-question pairs. **Left:** two query imagequestion pairs. **Right:** some of retrieved candidates. While the correct entity answer is retrieved indicated in green font, the MLLM decides to produce the wrong entity identifier depicted in red font.

category hierarchy map for the entity ATR 42. Finally, we collect entities that share any parent nodes in the category hierarchy as knowledge-similar entities.

We use ViTL/16-224px⁴ fine-tuned on ImageNet-1k as the visual classifier in VISION-HARD. While we collect the entities that share the same prediction label as vision-similar entities, we filter those entities with its image top-1 confidence score lower than 0.4, those entities without an infobox image, and those prediction labels with less than 10 entities. In Figure 3, we pick some entity images that share the image prediction "Alps" and "Goose".

C Error Analysis

We present two typical types of errors in Figure 1. Although the model is able to retrieve the correct entity, the MLLM predicts a too specific (but actually correct as it is indeed a great blue heron) entity in the upper sample. This is due to the misunderstood of intent by the MLLM, as the query does not ask for species but category. The lower sample has the MLLM make a wrong prediction on "Wetsuit" instead of "Diving suit".

by a particular common name", "Wikimedia disambiguation page", "variable-order class", "concept", "food", "artificial physical object", "artificial object", "physical object", "physical substance", "product", "object", "artificial physical structure", "architectural structure", "equipment", "tool", "physical location", "Wikimedia list article"]

⁴ https://huggingface.co/google/vit-large-patch16-224



Fig. 2: KB-HARD group construction illustration. Left: The SPARQL query template that used for constructing the category hierarchy of entities. We only consider category labels within 3 hops to exclude general categories like "object", "food", etc. Right: Visualized 3-hop category hierarchy map for the entity ATR 42. Code highlight and category visualization are from Wikidata Graph Builder.



Fig. 3: Example VISION-HARD groups constructed with predictions from ViTL/16-224px.

D More Case Study

We present 7 cases for side-by-side comparison between GPT-4V and AUTOVER-7B in Figure 5. They include situations where AUTOVER-7B has an advantage (first 4 rows), GPT-4V outperforms (row 5), both make accurate predictions (row 6) and both commit mistakes (row 7). We also conclude the reasons behind GPT-4V's failures, such as its inability to ground accurately to entities, refusal to predict with insufficient clues, hallucinatory responses, struggles with disambiguating challenging entities and lack of specific domain knowledge.

E Qualitative Analysis on A-OKVQA-ENT Dataset

We present a qualitative analysis in Fig. 4 for zero-shot baselines reported in the main paper. For all baselines under **multi-choice** settings except AUTOVER, we append an explicit instruction⁵ that is omitted in the figure. Still, LLaVA-v1-7B

 $^{^5}$ "Please strictly follow the options you are given. Only output the option number or the option itself."

4 Z. Xiao et al.

and OpenFlamingo-9B fail to follow the instructions, leading to their extremely low evaluation results. Our AUTOVER does not need any prompt engineering as it is guided by a prefix tree which never allows it to generate answers out of scope. We highlight the keywords where the model exhibits strong hallucinations in the analysis.

F Current Limitations

Although AUTOVER demonstrates superior performance using fewer trainable parameters compared to PaLI-17B, the accuracy on the UNSEEN subset remains suboptimal, particularly when compared with the almost identical human performance on the SEEN and UNSEEN subsets. We expect future work that focuses on exploring visual entity recognition for emerging entities that never show in the training set, and even closing the gap between model and human performance. In addition, the autoregressive nature makes it challenging to predict through a single cheap vector retrieval like a bi-encoder. Instead, it requires multiple forward passes on a billion-size language model. We believe that efficient inference such as in [1] could be an interesting direction toward real-time autoregressive visual recognition systems.

Image, Query and Choices	Input Prompt	OpenFlamingo-9B	InstructBLIP-7B	LLaVA-v1-7B	AutoVER-7B	LLaVA-v1.5-7B
	C: What type of animal is on a leash on the sidewalk? A. dog B. cat C. tiger D. lion A:	Answer: horse	A	A horse is on a leash on the sidewalk.	dog (Q144)	Dog
What type of animal is on a leash on the sidev dog cat	Q: What type of animal is on a leash on the sidewalk? Answer it using a single word or phrase. A:	A horse.	A horse.	A horse is on a leash on the sidewalk.	dog (Q144)	Horse
tiger lion						
	C: What does the man have on his head? A. shower cap B. hat C. earnuffs D. headphones A:	A shower cap.	D	The man has headphones on his head.	headphones (Q186819)	Headphones
25.00						
What does the man have on his head?	Q: What does the man have on his head? Answer it using a single word or phrase.	A hat.	2 earphones.	The man has a hat on his head.	headphones	Headphones
shower cap hat	A:				(Q100019)	
earmuffs bootsboros						
lieaupiones						
	Q: What style of hat is the man on the left wearing? A. top hat B. derby C. baseball cap D. fedora A:	The man on the left is wearing a fedora.	в	The man on the left is wearing a derby hat.	baseball cap (Q639686)	Baseball cap
What style of hat is the man on the left wearin top hat derby	Q: What style of hat is the man on the left using a single word or phrase. A:	A rainbow umbrella.	ibow umbrella. A bucket hat.	The man on the left is wearing a rainbow- colored hat.	baseball cap (Q639686)	Baseball cap
hand half and						
baseball cap fedora						
	C: What type of boat is on the beach? A. barge B. row C. sail D. cance A:	A rowboat.	A	There is a saliboat on the beach.	canoe (Q171529)	Canoe
the second s						
	Q: What type of boat is on the beach? Answer		A small boat.	There is a sallboat on the beach.	canoe (0171520)	Row boat
what type or boat is on the beach?	it using a single word or phrase.	A Seagull.				
barge row	A:				[Q1/1529]	

Fig. 4: Qualitative analysis on A-OKVQA-ENT dataset. For each triplet of image, query and choices, the upper row illustrates the input prompt and response from different baselines under **multi-choice** evaluation strategy, while the lower row is under **entity match**. Background color in **red** indicates an incorrect answer under the specific evaluation setting and that in green denotes a correct response. **Bold** highlights the model is heavily hallucinating.



Fig. 5: Case study on comparison between the GPT-4V response and AUTOVER-7B decision. Bold indicates the entity keyword. In the GPT-4V response, red indicates a false entity keyword and green indicates a correct one. In retrieved candidates and model decisions, red indicates a false entity prediction and green denotes a correct prediction. We also include the top-k position of the GPT-4V predicted entity for reference. The answer for the last query is *Matteuccia*.

References

- He, Z., Zhong, Z., Cai, T., Lee, J.D., He, D.: Rest: Retrieval-based speculative decoding. arXiv preprint arXiv: 2311.08252 (2023)
- Rasley, J., Rajbhandari, S., Ruwase, O., He, Y.: Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. p. 3505–3506. KDD '20, Association for Computing Machinery, New York, NY, USA (2020). https://doi.org/10.1145/3394486.3406703, https://doi.org/10.1145/3394486.3406703