

# Learning Dynamics of Multitask Training Data in Vision Language Models

Anonymous Author(s)

Affiliation

Address

email



Figure 1: We label each question in the LLaVA dataset as corresponding to one of five skills, allowing us to understand the learning dynamics of individual skills and their effects on each other. (Left) Example LLaVa questions corresponding to our five target skills. (Right) Model accuracy on training (solid lines) and validation (dashed lines) examples throughout the training process.

## Abstract

Vision language models (VLMs) are trained on massive amounts of data to perform many visual tasks simultaneously. Accordingly, many VLM benchmarks have been recently created to properly evaluate the models' capabilities. However, relatively little has been done to understand how and when the model acquires particular skills during training. We evaluate checkpoints throughout a one-epoch VLM training on recently seen and unseen datapoints to capture the generalization dynamics during model learning. We categorize the training data into five broad visual reasoning groups (Bounding, Complex, Object, OCR, and Semantic questions) and observe when these skills are learned. We note for example that despite not being explicitly trained to do OCR, VLMs can quickly learn to perform OCR tasks better than object recognition tasks. Digging deeper, we perform a case study on how VLMs use visual cues to solve OCR questions, indicating a form of shortcut that is not captured by standard VLM benchmarks. In contrast to OCR questions which are quickly learned, bounding capabilities are inefficiently learned due to the complexity of the bounding box format – despite the fact that bounding box questions comprise the majority of the training data. Our work provides a glimpse into the underlying learning process of VMs on the LLaVA dataset.

## 1 Introduction

Vision Language Models (VLMs), such as Flamingo and LLaVA [1, 9], achieve impressive performance across diverse visual tasks through massive multi-task pretraining on image-text pairs. These models can seemingly learn any visual reasoning task, from object detection to OCR [2, 7, 15], given

Category	Ratio	Description
Bounding	34.5%	Provide a $[x_1, y_1, x_2, y_2]$ bounding box or describe such a region
Complex	13.8%	Multi-step reasoning and logical deduction, often image-free
Object	19.8%	Relating to specific object(s), e.g., recognizing, counting
OCR	14.2%	Reading printed text and semantic queries about them
Spatial	17.7%	Spatial relationships of object(s) and between them

Table 1: **Category Details.** Ratios averaged over 7 sampled sets and 8 random seeds (gpt-4o-mini).

sufficient data. However, the source of this generalization remains poorly understood, and current evaluation methods provide limited insight into how VLMs actually learn during training.

The standard approach of evaluating on held-out validation sets [16] cannot distinguish between what knowledge the training data provides and how well the model learns from it. This is particularly problematic for VLMs trained under large language model loss objectives, where cross entropy differs significantly from model accuracy. We lack metrics to track accurate learning dynamics on the training data itself. Some works are beginning to identify how metrics such as perplexity and math multiple choice accuracy correlate to study the learning dynamics of LLMs for reasoning [4], but similar approaches for VLMs remain unexplored and elusive. We address this gap by developing an evaluation protocol that captures VLM training dynamics in the one-epoch setting. Our approach uses an LLM to judge the accuracy of the model on the open-ended training data, and categorizes visual reasoning skills into five distinct categories (Bounding, Complex, Object, OCR, and Spatial, summarized in Table 1) to identify how different visual skills perform. We track this performance on both recently seen and unseen examples throughout training, demonstrating how much information has been successfully captured from the recently-seen training data and to what extent has this information been generalized to new unseen examples.

Our key findings challenge common assumptions about VLM learning. Despite constituting only 14.2% of training data, OCR questions achieve the highest accuracy by the end of training, while bounding questions, despite being the majority (34.5%) of training data, struggle due to the inefficiency of coordinate-based representations. We demonstrate that VLMs learn OCR through visual shortcuts rather than traditional text recognition, revealing that high accuracy scores can be misleading about the underlying learning mechanisms, and that some skills like Object and Spatial reasoning generalize well while others like Complex reasoning plateau early. These insights provide a foundation for more data-effective VLM training and highlight the importance of evaluating training dynamics directly.

## 2 Creating categories and an evaluation protocol

We quickly establish the standard procedure for training VLMs on the LLaVA dataset. Our categorization scheme factors the multitask LLaVA data into corresponding visual capabilities to analyze their individual learning dynamics. This is made possible by separating examples into “seen” and “unseen” sets throughout training to measure generalization.

For the base of our studies, we adopt the LLaVA 1.5 dataset [10] using the Prismatic VLM design [5] of a one epoch, Stage 2 only training recipe with a DINOv2-SigLIP vision backbone [13, 17].

**LLaVA Dataset.** We adopt only the visual instruction tuning portion of the dataset, as these directly correspond to specific visual reasoning skills. This portion consists of image-text pairs derived from five vision datasets: COCO [8], GQA [3], OCRVQA [12], TextVQA [14], and VisualGenome [6]. Each example consists of 0-1 images and an average of five turns of questions and answers which are generally independent. In total, the dataset has 665k examples and 3.4 million question answer pairs.

**Visual Skill Categories.** While the LLaVA dataset is comprised of multiturn conversations, there are rarely dependencies between turns. We filter out those examples and split the remaining examples by turn. Analyzing the resulting questions, we found 5 overarching task categories: Bounding, Complex, Object, OCR, and Spatial. Examples of each category are shown in Figure 1 and detailed in Table 1.

**Evaluation Protocol.** Evaluating how well the model answers open-ended questions is a non-trivial task. We adopt the field standard of using an LLM as a judge of the model’s correctness [11].

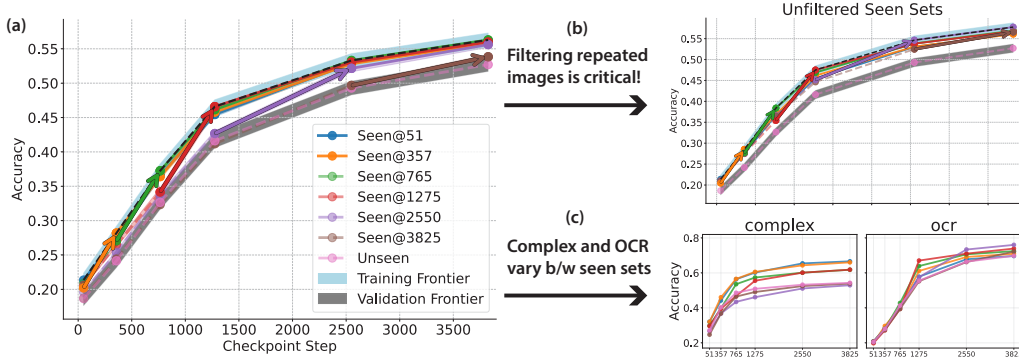


Figure 2: **Our LLaVA training dynamics exhibit a “train” frontier and a “validation” frontier.** As sets become seen for each checkpoint, each line becomes solid and jumps from the unseen level to the seen level. There is still a clear seen-unseen generalization gap, hinting that more learning is possible. Results are averaged over 8 random seeds with  $1\sigma$  widths.

Each ground truth question, answer and model response was provided to text-only gpt-4o-mini, and judgements were calibrated using a rubric which aligned well with human ratings (Appendix A.1).

To model the training dynamics, we fix a set of checkpoint steps as snapshots of different phases of the model’s training to evaluate on, and to prevent evaluating a prohibitively large number of questions. These are steps 51, 357, 765, 1275, 2550, and 3825 (out of a total of 5120 steps).

However, designing training and validation sets in a one-epoch regime requires care. A standard multi-epoch setup has a static training and validation set, the former being trained on in its entirety right before validation while the latter is never. In the one-epoch regime, we never see the same training example twice within a training run. Therefore, to mirror the standard setup as closely as possible, we create one seen set for each checkpoint consisting of the 5,000 most recently seen samples to ensure the freshest evaluation of capabilities, denoted as  $\text{Seen}@k$  for step  $k$ . Our validation set, or unseen set, is sampled from beyond the last selected checkpoint to ensure the questions are unseen. In this way,  $\text{Seen}@k$  is a validation set for checkpoints before  $k$  and a training set afterwards.

### 3 Results

**Emerging Training and Validation Frontiers.** Our LLaVA training dynamics reveal two distinct curves (Figure 2a): a “training” frontier tracking examples seen during training, and a “validation” frontier tracking unseen examples. These are highlighted as bands based on the highest seen and unseen accuracies for each checkpoint. Notably, there is an observable train-val gap, offering insight into how much the model generalizes from the data.

As validation of our methodology, we observe the following. For each  $\text{Seen}@k$  set, prior to checkpoint  $k$ , performance matches the validation frontier as expected for an untrained set. Once the model trains on the set, performance jumps to the training frontier, remaining consistent across all seen sets. We mark this transition with arrows in Figure 2a. The last checkpoint is an exception, which we attribute to the low learning rate near training end due to cosine annealing.

Determining what constitutes an “unseen” example proves critical to observing this behavior. A natural criteria would be to consider each  $(\mathcal{I}, \mathcal{T})$  image-text pair jointly, making every example different. However, images are used in an average of 3 different conversations, and seeing an image once is enough to alter our validation frontier (Figure 2b). The same is not true of text, likely due to the vast pre-training the LLM has already received.

**OCR learns quickly while Bounding struggles and Complex plateaus early.** Using our categorization scheme, we analyze each category’s performance throughout training. A clear pattern emerges (Figure 1, right). Object, Spatial, and Complex questions all start with similar accuracy at step 51, while OCR is slightly lower and Bounding is nearly 0%. Over time, Object, Spatial and OCR questions improve to around 65% accuracy, with OCR surprisingly becoming the best performer.

Complex questions plateau quickly, suggesting the model cannot learn these effectively. Bounding improves slowly, which is disappointing given that 34.5% of questions are Bounding questions, though this is not unexpected due to the difficulty of the bbox format and pseudo-regression objective.

Examining generalization gaps, Object and Spatial questions show acceptable gaps, while OCR has a large gap between seen and unseen. Complex demonstrates odd behavior, sometimes performing better on unseen than seen sets, suggesting the model fails to learn the data or that patterns are weak in this category. Bounding performs similarly between seen and unseen, which makes sense given the specificity of object detection tasks leaves little to overfit on.

**Complex and OCR dynamics differ throughout training.** One concern with our protocol is that seen sets may drift from the overall data distribution, as we are resampling and filtering them for each checkpoint. To investigate this, we plot each category’s learning dynamics across every seen set (Figure 2c). Only Complex and OCR show strong variance with seen sets, while Bounding, Object, and Spatial questions remain tightly grouped without temporal artifacts from the one-epoch setting. Complex questions perform worse for later sets than earlier ones without much difference between checkpoints, which we attribute to a lack of strong language-centric data in LLaVA causing the model to relinquish complex text understanding abilities and plateau early. OCR exhibits different dynamics where the best performing set for each checkpoint is the corresponding seen set, followed by previous ones in order of recency, pointing to strong memorization and weak generalization.

**VLMs take surprising paths to learn Bounding and OCR.** We conclude with an enlightening case study of Bounding and OCR questions. Bounding questions divide into two forms: “Describe” questions where the model describes a provided bounding box region, and “Bbox” questions where the model returns bounding box coordinates for a specified description.

Intuitively, one might expect describe questions to perform much better due to stronger text supervision, but both types perform poorly. As shown in Figure 3, the model performs slightly better when asked to return a region description as expected. However, the performance gap is small because describe questions also involve floating coordinates.

OCR questions in the LLaVA dataset are nearly entirely about books and divide into two subcategories: recognition and semantic questions. Recognition questions require direct transcription of the title or author. Semantic questions require external knowledge or logical reasoning, such as determining a book’s genre or whether it’s for children.

Surprisingly, VLMs learn semantic questions much faster and earlier than recognition questions, circumventing the expected learning order entirely! This shows that recognition questions are difficult, especially given the large generalization gap, and that the VLM likely relies on other visual cues to solve semantic questions rather than first learning to read text.

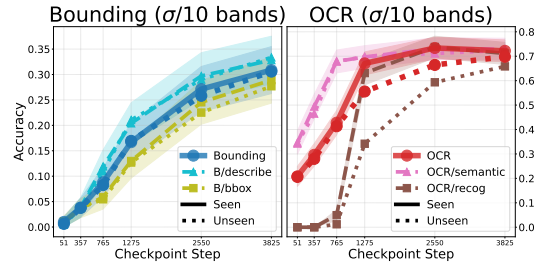


Figure 3: **Bounding and OCR have counter-intuitive learning dynamics** when investigating their subcategories.

## 4 Conclusion

We have presented a new methodology for evaluating the training dynamics of VLMs in a one-epoch setting. Under this setting, we have shown that two distinct frontiers emerge which can capture generalization: a “training” frontier, which tracks examples seen during training, and a “validation” frontier, which tracks unseen examples. We take this one step further by categorizing the questions to understand the learning dynamics of each category and how some categories learn faster than others. Limitations of our method are that we only perform analysis on a single type of VLM architecture and on one dataset, albeit one of the most popular ones. Additionally, LLM-as-a-judge is expensive and not feasible for development, so an equally informative signal or automatic metric could make such understanding more widespread. This is a promising direction for future work.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *NeurIPS*, 2022. 1
- [2] Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. *arXiv preprint arXiv:2109.10852*, 2021. 1
- [3] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019. 2
- [4] Katie Kang, Amrith Setlur, Dibya Ghosh, Jacob Steinhardt, Claire Tomlin, Sergey Levine, and Aviral Kumar. What do learning dynamics reveal about generalization in llm reasoning? *arXiv preprint arXiv:2411.07681*, 2024. 2
- [5] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic vlms: Investigating the design space of visually-conditioned language models. In *ICML*, 2024. 2
- [6] 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*, 123:32–73, 2017. 2
- [7] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022. 1
- [8] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft COCO: Common Objects in Context, February 2015. 2
- [9] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 1
- [10] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>. 2
- [11] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *ACL*, 2024. 2, 6
- [12] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *ICDAR*, 2019. 2
- [13] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. In *TMLR*, 2023. 2
- [14] 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. 2
- [15] Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, Ziteng Wang, Rob Fergus, Yann LeCun, and Saining Xie. Cambrian-1: A fully open, vision-centric exploration of multimodal llms, 2024. URL <https://arxiv.org/abs/2406.16860>. 1
- [16] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 2
- [17] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023. 2



## A Technical Appendices and Supplementary Material

### A.1 LLM-as-Judge

To evaluate the performance during training, we utilized text-only GPT-4o-mini as a judge. For each question, we pass in a tuple consisting of the question, ground truth answer, and predicted answer, and ask the LLM to decide if the question was answered correctly or not. Following Video-ChatGPT [11], we create a rubric that guides the model to first provide an overall score before returning a yes/no prediction for better calibration. See Figure 5 for the rubric.

	Acc	Prec	Recall	$F_1$ Score
LLM (w/ rubric)	<b>89.4%</b>	93.0%	<b>89.8%</b>	<b>0.91</b>
LLM (w/o rubric)	85.1%	<b>97.9%</b>	78.0%	0.87

Table 2: **An LLM (gpt-4o-mini) serves as a suitable judging replacement for humans.** Our rubric helps improve the overall accuracy and consistency as shown by the better  $F_1$  score.

The rubric consists of three scores ranging from 1-5 which help the LLM better judge under the text-only setting: Missed details, Hallucinations, and Major Subjects. Missed details describes how many ground truth answer details were left out in the predicted answer. Hallucinations describes how plausible the added details (hallucinations) are in the predicted answer. Major subjects compares the major subjects in ground truth and predicted answer. A score of 5 is the most accurate, and score of 1 represents a bad and weak answer. Each score is averaged to get a general prediction score for each question-answer pair.

We conducted a small human study to validate both the use of GPT-4o-mini as a judge, as well as our rubric for improved calibration. The human study was conducted on 4 individuals, each receiving the same 100 question-answer-prediction tuples as the judge. The subjects were asked to rate if each prediction was plausibly correct given the ground truth, and a threshold of 75% agreement was used for the final human ratings. As shown in Table 2, we found that the LLM on its own was a suitable replacement for humans, achieving 85.1% accuracy. However, including the rubric boosts its accuracy by 4.3% to 89.4%, and improves its overall precision-recall tradeoff as shown by the 4 point increase in  $F_1$  score. We adopt this rubric evaluation scheme by default for our methodology.

For completeness, we also include the prompt used to categorize the training data examples into categories in Figure 6.

### A.2 Full Figures

We present the full plots of Figure 2(c) here in Figure 4, showing that the other three categories are consistent across seen sets while Complex and OCR vary.

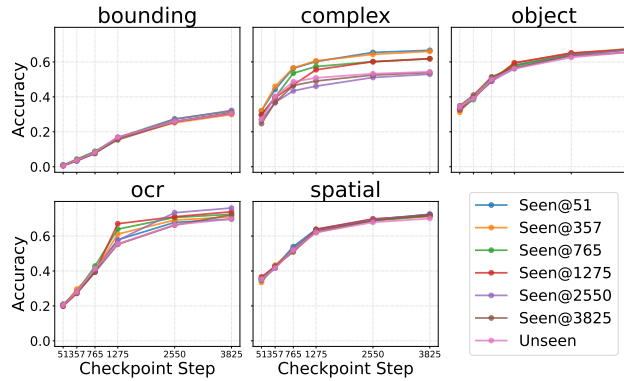


Figure 4: **Complex and OCR questions vary with seen set, while other categories are consistent.** Complex questions show degraded performance in later sets, suggesting interference from newly acquired visual knowledge. OCR performance indicates memorization, as seen accuracy never drops below unseen performance.

System: You are an intelligent chatbot designed for evaluating the correctness of generative outputs for question-answer-prediction tuples. Your task is to grade the pred answer given a correct answer and a rubric, then provide a final score.

-----

##Rubric:

Do not grade too harshly

In example sentence, 'a woman in a brown shirt, holding a white purse, is hugging a bear', major subjects are 'a woman' and 'a bear'. The details would include everything related to major subjects. Such as 'in a brown shirt', 'hugging', 'holding a white purse'.

Also subjects can be the same/similar but have different names. For example, a bag and suitcase would be the same subject.

When comparing, focus on the meaning rather than exact language.

Major Subjects: Graded 1 - 5, with 5 being the highest score. Give a score of:

- 1 if no major subjects remotely similar.
- 2 if a few (>2) subjects are not the same.
- 3 if some (2) subjects are not the same.
- 4 if subjects that are different are also semi-plausible given context of answer
- 5 if subjects that are different are also plausible given context of answer

Missed Details: Graded 1 - 5, with 5 being the highest score. Give a score of:

- 1 if many (>5) key details missed that change the meaning of the answer
- 2 if some (~5) key details missed
- 3 if some (2 to 4) non-important details missed
- 4 if few (2) non-important details are missed
- 5 if very few to none (<2) non-important details are missed

Hallucinations: Graded 1 - 5, with 5 being the highest score. Give a score of:

- 1 if all details that are added are not plausible
- 2 if some (>3) details added that are not plausible
- 3 if few (~3) details added or if added details are semi-plausible
- 4 if very few (1 to 3) details are added or that details added are plausible
- 5 if no details added or details added are plausible

Please evaluate the following image-based question-answer-prediction tuples with the given rubric.

Question: ..., Correct Answer: ..., Predicted Answer: ...

Provide your evaluation only as a yes/no, score for Major Subjects, score for Hallucinations, and a score for Missed Details where both scores are an integer value between 1 and 5, with 5 indicating the highest meaningful match. If the pred answer could be true given the context of the correct answer, then evaluation should be yes, otherwise it should be no.

Please generate the response in the form of a Python dictionary string with keys 'p', 'MS', 'MD', and 'H', where value of 'p' is a string of 'yes' or 'no' and values of 'MS', 'MD', and 'H' are in INTEGER, not STRING.

p stands for prediction, MS for Major Subjects, H for Hallucinations and MD for Missed Details

DO NOT PROVIDE ANY OTHER OUTPUT TEXT OR EXPLANATION. Only provide the Python dictionary string.

For example, your response should look like this:

{ 'p': 'yes', 'MS': 5, 'MD': 4, 'H': 3 }.

**Figure 5: The rubric used to judge the VLM responses.** We use the scores as a way to calibrate the LLM's responses before asking it to give a final score. The score is generally calibrated with the accuracies, so that the accuracies within each score band are roughly proportional to the score.

System: You are an intelligent chatbot designed for categorizing different types of questions. Your task is to group questions/tasks into 4 categories: Object Analysis, Spatial Analysis, Bounding Box, and Complex Reasoning.

-----

##INSTRUCTIONS:

- Focus on the goal of the questions.
- Look at the examples of each category to aid in categorization. Examples of each category are the following:

Object Analysis:

Which kind of appliance is it?  
Are there either any doors or windows that are made of metal?  
Is this a transportation engineering book?  
What color are the umbrellas in the picture?

Spatial Analysis:

Are there any players to the left of the helmet on the right?  
What kind of device is to the left of the computer monitor?  
What is near the bottle of alcohol?  
A. bunny  
B. toilet  
C. whistle  
D. man  
Answer with the option's letter from the given choices directly.  
Who is in the water on the beach?

Bounding Box:

Please provide the bounding box coordinate of the region this sentence describes: man with royal blue and white toothbrush in mouth.  
Please provide a short description for this region: [0.06, 0.78, 0.93, 0.83].

Complex Reasoning:

Why are the cats resting?  
What can we infer about the elephants' social behavior from this scene?  
What potential reasons might explain the unattended devices in this scene?  
What is the genre of this book?  
What activities might someone enjoy in this well-lit room?

Please categorize the following question into the 4 categories: Object Analysis, Spatial Analysis, Bounding Box, and Complex Reasoning. DO NOT PROVIDE ANY OTHER OUTPUT TEXT OR EXPLANATION. Return O if Object Analysis, S if Spatial Analysis, B if Bounding Box, or C if Complex Reasoning.

Figure 6: The prompt used to categorize the training data examples into categories. We use the scores as a way to calibrate the LLM's responses before asking it to give a final score. The score is generally calibrated with the accuracies, so that the accuracies within each score band are roughly proportional to the score.

### 225 A.3 Societal Impacts

226 Like any work involving generative models, specifically LLMs, there are potential harms without  
227 significant safeguards. The models we use should have been tuned to be safe, but there are always  
228 risks. The datasets we use are among some of the most widely used, so they have been through much  
229 scrutiny. In general, we hope that our work helps to shed a light on how exactly these models gain  
230 their capabilities, which previously has been a relatively opaque process.

### 231 A.4 Compute Resources

232 All experiments were run on a cluster of 48GB L40s. The storage of the checkpoints was the limiting  
233 factor, which required some smart management of the checkpoints.



## NeurIPS Paper Checklist

### 1. Claims

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

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

Justification: The paper's abstract and introduction are written with the goal of reflecting the content of the paper.

Guidelines:

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

### 2. Limitations

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

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

Justification: This is included in the conclusion as suggestions for future work.

Guidelines:

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

### 3. Theory assumptions and proofs

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

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

Justification: No theoretical results.

Guidelines:

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

#### 4. Experimental result reproducibility

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

Answer: [Yes]

Justification: Yes, plus code is planned to be released. Seeds were carefully chosen to ensure reproducibility, and the labeled data from GPT is also planned to be released (how the models are updated with their knowledge is outside of our control).

Guidelines:

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

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

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Yes, our code is based on publically available code and we plan to also release the labeled data from GPT.

Guidelines:

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

## 6. Experimental setting/details

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

Answer: [Yes]

Justification: Most details are present in the paper, in Section 2, however the specific seeds should not be a concern with enough runs (the final seeds will still be present in the code, as well as the specific hyperparameters).

Guidelines:

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

## 7. Experiment statistical significance

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

Answer: [Yes]

Justification: Yes,  $1\sigma$  error bars are reported over 8 random seeds as mentioned in the paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

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

Answer: [Yes]

Justification: Yes, the compute resources are reported in the appendix.

Guidelines:

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

## 9. Code of ethics

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

Answer: [Yes]

Justification: Yes, all of our data and methods are publicly available, and we hope to shed better light on how these models are trained on what their capabilities are.

Guidelines:

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

## 10. Broader impacts

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

Answer: [Yes]

Justification: Yes, we discuss the potential societal impacts of our work in the paper in the appendix.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.

- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

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

Answer: [NA]

Justification: NA, we simply follow prior work in this area as the models we use should already have safeguards in place.

Guidelines:

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

## 12. Licenses for existing assets

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

Answer: [Yes]

Justification: Yes, we cite the original paper that produced the code package or dataset.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New assets

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

Answer: [NA]

Justification: No new assets are introduced in the paper.

Guidelines:

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

### 14. Crowdsourcing and research with human subjects

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

Answer: [Yes]

Justification: Yes, we perform a brief small-scale human study as described in the appendix.

Guidelines:

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

### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

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

Answer: [NA]

Justification: No IRB needed for our small-scale (<4 participants, <1 hour) human study.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

544           • We recognize that the procedures for this may vary significantly between institutions  
545           and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
546           guidelines for their institution.  
547           • For initial submissions, do not include any information that would break anonymity (if  
548           applicable), such as the institution conducting the review.

549   **16. Declaration of LLM usage**

550   Question: Does the paper describe the usage of LLMs if it is an important, original, or  
551   non-standard component of the core methods in this research? Note that if the LLM is used  
552   only for writing, editing, or formatting purposes and does not impact the core methodology,  
553   scientific rigorousness, or originality of the research, declaration is not required.

554   Answer: [\[Yes\]](#)

555   Justification: Yes our work seeks to analyse how large multimodal LLMs, which are trained  
556   on top of LLMs, learn during their training.

557   Guidelines:

558           • The answer NA means that the core method development in this research does not  
559           involve LLMs as any important, original, or non-standard components.  
560           • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
561           for what should or should not be described.