Analyze, Generate, Improve: Failure-Based Data Generation for Large Multimodal Models

Gabriela Ben-Melech Stan1Estelle Aflalo1Avinash Madasu1Vasudev Lal1Phillip Howard2†¹Intel Labs²Thoughtworks

{gabriela.ben.melech.stan,estelle.aflalo}@intel.com phillip.howard@thoughtworks.com

Abstract

Training models on synthetic data is an effective strategy for improving large multimodal models (LMMs) due to the scarcity of high-quality paired image-text data. Existing methods generate multimodal datasets but do not address specific reasoning deficiencies in LMMs. In contrast, humans learn efficiently by focusing on past failures. Inspired by this, we propose a synthetic data generation approach that analyzes an LMM's reasoning failures using frontier models to generate and filter high-quality examples. Our method produces a 553k-example multimodal instruction tuning dataset, leading to improved LMM performance, even surpassing models trained on equivalent real data demonstrating the high value of generating synthetic data targeted to specific reasoning failure modes in LMMs.

1. Introduction

Recent advancements in large language models (LLMs) can be attributed largely to scaling models and training data. However, high-quality data availability limits further scaling [38, 41]. As a result, synthetic data generation is gaining traction for augmenting datasets. High-quality synthetic data generation has grown more feasible with increasing LLM capabilities, enabling synthetic data production at scale [7, 13, 14, 20, 34, 44]. Synthetic data is particularly valuable for training large multimodal models (LMMs), which combine an LLM with a vision encoder to enable text generation conditioned on multimodal inputs. Real data for training LMMs is relatively scarce due to the lack of naturally-occurring images paired with high-quality text.

Existing approaches for generating multimodal synthetic data face two key limitations. First, they depend on real images paired with synthetic text generated by another LMM, restricting their use when image data is scarce. Second, they generate data arbitrarily without prioritizing useful examples, leading to inefficiencies in both data generation and

training. For a broader discussion of previous methods, see Section 6 in the supplementary material. Unlike standard synthetic data generation, humans learn efficiently by focusing on examples tied to past reasoning failures. Ericsson et al. [8] suggest expert learning benefits from explicit instruction, error diagnosis, and targeted feedback. Since humans learn from failures [3, 6], seeking examples related to past errors aids mastery, while problems involving familiar reasoning are often ignored.

Inspired by this observation, we propose a synthetic data generation approach based on an LMM's reasoning failures. We first evaluate the LMM on a benchmark dataset and use a strong frontier model to analyze its errors. The frontier model then generates related question-answer pairs and image descriptions, which can link to existing images or guide synthetic image generation. Finally, we ensure quality using an LMM-as-a-judge filtering process. Our approach generates a 553k-example multimodal dataset from LLaVA-1.5-7B's [22] reasoning failures. Training experiments show our synthetic data improves LLaVA's performance across various downstream tasks, even surpassing training on real datasets, unlike prior work requiring significantly more examples to match real data performance [12]. Our method achieves greater gains compared to previously-proposed synthetic datasets, demonstrating the value of grounding synthetic data in model failure analysis.

2. Dataset construction

Diagnosing model failures To generate our tailored synthetic data, we first analyze reasoning failures in a baseline LMM using a more advanced frontier LMM. The frontier model is selected for its superior multimodal reasoning capabilities and high accuracy on diverse vision-language benchmarks. Reasoning failures are identified by evaluating both models on the training sets of vision-language benchmarks, and selecting samples where the baseline LMM produces incorrect responses, while the frontier model succeeds. This process generates a subset of failure cases per benchmark, providing a focused challenging dataset, which

[†]Work completed while at Intel Labs.



Figure 1. Illustration of our approach. Given a sample from an existing dataset which LLaVA answers incorrectly, we prompt a frontier model to analyze LLaVA's reasoning failures and propose new synthetic samples which require similar types of reasoning.

we denote as Model Failure Sets (MFS).

Synthetic data generation Using the resulting MFS, we guide the frontier model through a structured multi-turn process to diagnose the failure reasons of the baseline LMM and generate new training samples that address these failure modes, as illustrated in Figure 2 in Section 7.3 of the supplementary material (SM). In the first steps, the frontier model is prompted to describe the image and analyze reasoning errors by examining the question, ground truth answer, and the incorrect response generated by the baseline LMM. Next, the frontier model is instructed to propose new challenging samples designed to target the identified failure modes, which consist of a detailed image description, a clear question, and a deterministic answer. We explore two different approaches for sourcing images: utilizing existing real images (Method 1), or using synthetically generated images (Method 2).

Specifically, in Method 1, we leverage the original image from a failed sample and prompt the frontier model to generate 10 new question-answer pairs per sample, following the prompt in Figure 2, omitting the final step. This is especially effective for benchmarks like InfoVQA and ScienceQA, since text-to-image models often struggle with precise text rendering and spatial accuracy. In Method 2, the frontier model generates a question-answer pair and a detailed image prompt using the prompt in Figure 2. Each image prompt is then fed into a text-to-image diffusion model to generate 10 synthetic images at varying guidance scales, producing 100 fully synthetic samples per failed sample.

To enhance data diversity, we use a variation of our

prompt instructing the frontier LMM to "provide examples that challenge Model A's weaknesses using scenarios from entirely different domains or situations". This instruction is added to Step 4 in Figure 2 to encourage samples generation in different domains. Domain-similar samples preserve the original theme, while non-similar samples offer broader contextual diversity for improved generalization, see Figure 6 in Section 7.3 (SM). We also enforce constraints on question format and instructions, as detailed in Section 7.3.

Filtering We apply a filtering process using the same frontier LMM which produced the samples. The LMM is instructed to evaluate each synthetic sample given an image, question, and answer on a scale of 1 to 3, where 1 indicates an incorrect sample, 2 being partially correct, and 3 fully correct. The prompt used for filtering is provided in Figure 3 in Section 7.4 (SM). Only samples rated 3 were included in the final dataset to ensure quality and reliability.

Dataset overview Our synthetic dataset consists of 553,992 samples incorporating both real and generated images derived from the MFS of LLaVA-1.5-7B, with Vicuna-1.5-7B [46] base LLM, on four benchmark training sets: VizWiz [11], InfoVQA [27], ScienceQA [25], and OK-VQA [26]. These benchmarks were chosen to ensure a wide range of visual and reasoning challenges, as detailed in Section 7.2 (SM). We use GPT-4o [30] as the frontier LMM for analyzing the reasoning failures of LLaVA-7B due to its strong multimodal reasoning capabilities, and FLUX.1-schnell Labs [18] text-to-image model to generate the synthetic images. Table 3 in Section 7.2 (SM) provides an

| Dataset | Model | N | N_{syn} | EM Score | |
|-----------|-----------------------|---------|-----------|----------|--|
| | LLaVA | 624,610 | 0 | 26.7 | |
| | LLaVA _{real} | 634,684 | 0 | 31.6 | |
| InfoVQA | LLaVA _{syn} | 634,684 | 10,074 | 30.8 | |
| | LLaVA _{syn} | 687,071 | 62,461 | 33.0 | |
| | LLaVA _{syn} | 710,610 | 86,000 | 34.3 | |
| | LLaVA | 624,610 | 0 | 70.7 | |
| SaianaaOA | LLaVA _{real} | 630,195 | 0 | 70.0 | |
| ScienceQA | LLaVA _{syn} | 630,195 | 5,585 | 71.8 | |
| | LLaVA _{syn} | 646,594 | 21,984 | 73.0 | |
| | LLaVA | 624,610 | 0 | 57.0 | |
| OK-VQA | LLaVA _{real} | 633,619 | 0 | 54.3 | |
| | LLaVA _{syn} | 633,619 | 9,009 | 61.3 | |
| | LLaVA _{syn} | 687,071 | 62,461 | 61.3 | |
| | LLaVA _{syn} | 749,532 | 124,922 | 61.5 | |

Table 1. In-domain evaluation of baseline LLaVA, LLaVA models trained using synthetically augmented data (LLaVA_{syn}) and using training data augmented with real in-domain data (LLaVA_{real}).

overview of each benchmark, including the original size of the training split, the MFS size of LLaVA-7b, and the total number of synthetic samples retained after filtering.

As previously described, we utilize two methods for synthetic data generation, depending on the dataset: for InfoVQA and ScienceQA, we apply Method 1, which generates new question-answer pairs based on the original real images. For OK-VQA, we employ Method 2, generating synthetic images along with corresponding question-answer pairs, and for VizWiz we apply both methods. Overall, our final dataset consists of 42% real-image-based samples and 58% fully synthetic samples. All samples in the final dataset contain a single-turn conversation, featuring a question and a short answer. Our filtering approach effectively removes lower-quality synthetic samples, with synthetic-image samples exhibiting higher removal rates compared to real-image samples, further details in Section 7.2 (SM). Examples of our dataset can be seen in Section 11 (SM).

Additionally, the large volume of filtered data ensures broad coverage of LLaVA-7B's reasoning failures across benchmarks, as listed in Section 9.5 (SM).

3. Experiments

Details of experimental setting We fine-tune LLaVA-1.5-7B using a subset of the LLaVA-Instruct-1.5-mix-665K dataset [23], which contains 665K user-GPT conversations focused on visual prompts. Since our approach targets image-grounded reasoning, we use only the 624K samples with images. We then augment this visually oriented subset with our synthetic MFS data (Section 2). The final training set combines real image-text conversations with targeted synthetic examples to address reasoning failures.

For fine-tuning, we used Vicuna-1.5-7B weights as the LLM backbone and leveraged the pretrained multimodal projector from LLaVA-1.5-7b. We followed the original

LLaVA training procedure, ensuring a fair comparison to existing methods, with training details provided in Section 8 (SM). We refer to models trained with a mixture of the original LLaVA-Instruct dataset and our synthetic data as LLaVA_{syn}. As a baseline, we report the performance of LLaVA trained under the same setting but without any additional synthetic data added to the training dataset (i.e., $N_{syn} = 0$). Additionally, for each dataset from which reasoning failures were derived for synthetic data generation, we report the performance of a LLaVA model trained on an equivalent amount of real data sourced from the corresponding training dataset (denoted as LLaVA_{real}). This provides a measure of the efficiency of our synthetic data relative to training on real in-domain data.

We used a variety of multimodal reasoning benchmarks for evaluating models. For in-domain evaluations (i.e., evaluating on a withheld validation set corresponding to the training set from which reasoning failures were derived), we used the InfoVQA, OK-VQA and ScienceQA validation sets. Since InfoVQA and ScienceQA questions rely heavily on reading text contained in the images, we further evaluated models on TextVQA [33] and OCR-Bench [9]. Finally, because our synthetic dataset was designed to enhance the model's reasoning capabilities, we chose MMBench [24] and MMMU [45] as additional OOD benchmarks.

Synthetic data augmentation results Table 1 provides in-domain evaluation results utilizing synthetic data derived from InfoVQA, ScienceQA, and OK-VQA. Notably, augmenting the LLaVA-Instruct dataset with our synthetic data achieves performance comparable to or better than using an equivalent amount of real domain-specific data in most cases. This result is particularly significant given that the synthetic samples were generated using only a small subset of the original training data: specifically, only those examples where LLaVA scored 0.0 while GPT scored 1.0. For instance, with OK-VQA, the original training set consists of 9,009 samples, but we utilized only 607 training samples which LLaVA failed on in order to generate 9009 synthetic samples, resulting in a performance boost of 13% on the OK-VQA test set. Similarly, our approach utilized only 28% of the ScienceQA training dataset to generate full synthetic replacements, yet still resulted in better performance than training directly on the real dataset. We also observe that performance improves as the amount of synthetic data used for data augmentation increases. In practice, it may be desirable to combine synthetically generated data which was derived from reasoning failures across different datasets. We therefore provide results for two different sized mixtures of our synthetic data derived from InfoVQA, ScienceQA, and OK-VQA reasoning failures in Table 2. In addition to in-domain evaluations for these three datasets, we provide results for the four other datasets mentioned

| Base LLM | N | N_{syn} | Augmentation Data | TextVQA | OCR-Bench | InfoVQA | OK-VQA | ScienceQA | MMBench | MMMU |
|-----------|--------------------|--------------|--------------------------------|-----------------------------|-----------------------------|----------------------|----------------------|----------------------|-----------------------------|------------------------------|
| Vicuna-7B | 624,610 | 0 | N/A (baseline) | 47.0 | 31.9 | 26.7 | 57.0 | 70.7 | 52.3 | 36.4 |
| Vicuna-7B | 687,071 | 62,461 | ALLaVA CoSyn-400K SimVQA | 47.9 47.1 46.8 | 34.0 31.8 31.6 | 28.4 28.5 26.6 | 50.4 55.8 54.4 | 71.2 71.5 71.1 | 50.2 52.4 53.5 | 36.7 34.6 34.7 |
| Vicuna-7B | 749,532 | 124,922 | ALLaVA CoSyn-400K | 47.4 47.2 46.8 | 33.2 34.1 32.7 | 28.8 29.6 | 49.4 57.7 | 66.5 70.9 | 52.5 43.5 51.9 | 36.2 34.2 36.9 27.4 |
| Gemma-2B | 624,610 749,532 | 0 124,922 | Ours | 39.9 | 28.3 | 21.8 | 51.7 | 62.3 | 29.2 | 37.4 |
| | | | | 40.9 | 29.9 | 29.8 | 54.8 | 65.2 | 30.7 | 31.7 |
| Qwen2-7B | 624,610 749,532 | 0 124,922 | Ours | 45 46.2 | 31.4 32.5 | 26.7 27 | 59.2 60.6 | 77.9 80.7 | 63 63.5 | 42.3 42.2 |

Table 2. Training data augmentation experimental results. N denotes the total number of training examples, N_{syn} denotes the number of synthetic examples in the training dataset generated using our approach. The first section of the table compared the same base LLM (Vicuna) trained on various datasets with our dataset, while the second section compares different LLM backbones trained on our dataset.

previously to measure the impact on OOD generalization. We also provide results for LLaVA models trained on three alternative synthetically generated datasets: ALLaVA, CoSyn-400k, and SimVQA. From Table 2, we observe that our synthetic data outperforms all baselines in the maximum data augmentation setting ($N_{syn} = 124, 922$). When augmenting training sets with half as much data, we observe that our synthetic data produces the greatest improvements in-domain, whereas the ALLaVA dataset, which contains only real images paired with synthetic text, performs slightly better in other OOD settings. Finally, training LLaVA on the mixed synthetic dataset yields similar indomain performance to training it separately on failures from each individual dataset (Table 1).

To evaluate the generalization ability of our synthetic dataset, we trained models using the same datasets but with different backbone LLMs. In the following experiments, we used Gemma-2B [35] and Qwen2-7B [42] as base LLMs. We adopted the LLaVA two-phase training procedure: pretraining on the LLaVA 558k dataset [23] followed by instruction fine-tuning. Table 2 indicates that our synthetic data in its maximum augmentation setting, $N_{syn} = 124,922$, outperform the baseline for both LLaVA-Gemma-2B and LLaVA-Qwen-7B almost on all benchmarks. Notably, despite being generated based on LLaVA-Vicuna-7B failure mechanisms, our synthetic data enhance the performance of other models whether they are of same size (Qwen2-7B) or smaller (Gemma-2B).

4. Analysis and Ablations

We analyzed the quality of our dataset, as detailed in Section 9 (SM). We conducted a human evaluation of our dataset (Section 9.1) and found that fully synthetic samples match or slightly exceed the fidelity of real-imagebased sample. The evaluation also demonstrates excellent alignment between generated questions, answers, and image prompts. In Section 9.2, we tested our data in low resource setting by replacing portions of the original dataset with fully synthetic samples. Even with up to 25% substitution, the model finetuned on our dataset matches or outperforms the baseline LLaVA, highlighting the efficiency of our targeted dataset. Section 9.3 quantitatively describes the impact of filtering on data quality and shows that our filtering approach improves model performance. Section 9.4 compares two frontier LMMs for synthetic data generation, showing that Qwen2-VL [39] produces samples with lower downstream performance than GPT-40 (See example in Figure 5). Finally, Section 9.5 identifies reasoning failures and shows that targeted synthetic data improves LLaVA-Instruct's performance on related benchmarks.

5. Conclusion

We introduced a new approach for generating multimodal synthetic data based on analyzing a model's reasoning failures. This led to a multimodal instruction tuning dataset with over 553k synthetic examples derived from LLaVA's failures. Experiments show that our data significantly improves LLaVA's performance on InfoVQA, ScienceQA, and OK-VQA, even surpassing training on an equivalent amount of real data. Furthermore, models trained on our synthetic dataset exhibit improvements in OOD evaluations and outperform training on other existing synthetic datasets when training data augmentation is scaled, showing consistent improvements across different base models. We also showed that training LLaVA only on examples derived from specific failure modes improves its performance on tasks which require corresponding forms of reasoning. Ablations and human evaluations confirm the effectiveness and quality of our approach, highlighting the potential of targeted synthetic data generation to address model deficiencies.

References

- [1] Deepseed AI. Deepspeed. https://github.com/ deepspeedai/DeepSpeed, 2025. Accessed: 2025-03-07. 2
- [2] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101-mining discriminative components with random forests. In *Computer vision–ECCV 2014: 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part VI 13*, pages 446–461. Springer, 2014. 4
- [3] John Seely Brown, Allan Collins, and Paul Duguid. Situated cognition and the culture of learning. *1989*, 18(1):32–42, 1989.
- [4] Atoosa Chegini and Soheil Feizi. Identifying and mitigating model failures through few-shot clip-aided diffusion generation. arXiv preprint arXiv:2312.05464, 2023. 1
- [5] Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for lite vision-language models. arXiv preprint arXiv:2402.11684, 2024. 1
- [6] Aubteen Darabi, Thomas Logan Arrington, and Erkan Sayilir. Learning from failure: A meta-analysis of the empirical studies. *Educational Technology Research and Development*, 66:1101–1118, 2018. 1
- [7] Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Junjie Hu, Luu Anh Tuan, and Shafiq Joty. Data augmentation using llms: Data perspectives, learning paradigms and challenges. In *Findings* of the Association for Computational Linguistics ACL 2024, pages 1679–1705, 2024. 1
- [8] K Anders Ericsson, Ralf T Krampe, and Clemens Tesch-Römer. The role of deliberate practice in the acquisition of expert performance. *Psychological review*, 100(3):363, 1993. 1
- [9] Ling Fu, Biao Yang, Zhebin Kuang, Jiajun Song, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu Wang, Hao Lu, Mingxin Huang, Zhang Li, Guozhi Tang, Bin Shan, Chunhui Lin, Qi Liu, Binghong Wu, Hao Feng, Hao Liu, Can Huang, Jingqun Tang, Wei Chen, Lianwen Jin, Yuliang Liu, and Xiang Bai. Ocrbench v2: An improved benchmark for evaluating large multimodal models on visual text localization and reasoning, 2024. 3
- [10] Himanshu Gupta, Kevin Scaria, Ujjwala Anantheswaran, Shreyas Verma, Mihir Parmar, Saurabh Arjun Sawant, Chitta Baral, and Swaroop Mishra. Targen: Targeted data generation with large language models. arXiv preprint arXiv:2310.17876, 2023. 1
- [11] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617, 2018. 2, 1
- [12] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and Xiaojuan Qi. Is synthetic data from generative models ready for image recognition? *arXiv preprint arXiv:2210.07574*, 2022. 1, 3

- [13] Phillip Howard, Gadi Singer, Vasudev Lal, Yejin Choi, and Swabha Swayamdipta. Neurocounterfactuals: Beyond minimal-edit counterfactuals for richer data augmentation. arXiv preprint arXiv:2210.12365, 2022. 1
- [14] Phillip Howard, Junlin Wang, Vasudev Lal, Gadi Singer, Yejin Choi, and Swabha Swayamdipta. Neurocomparatives: Neuro-symbolic distillation of comparative knowledge. arXiv preprint arXiv:2305.04978, 2023. 1
- [15] Saachi Jain, Hannah Lawrence, Ankur Moitra, and Aleksander Madry. Distilling model failures as directions in latent space. arXiv preprint arXiv:2206.14754, 2022. 1
- [16] Gyeongman Kim, Doohyuk Jang, and Eunho Yang. Promptkd: Distilling student-friendly knowledge for generative language models via prompt tuning. arXiv preprint arXiv:2402.12842, 2024. 1
- [17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 4
- [18] Black Forest Labs. Flux. https://github.com/ black-forest-labs/flux, 2024. 2, 1
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
 4
- [20] Nicholas Lee, Thanakul Wattanawong, Sehoon Kim, Karttikeya Mangalam, Sheng Shen, Gopala Anumanchipalli, Michael W Mahoney, Kurt Keutzer, and Amir Gholami. Llm2llm: Boosting llms with novel iterative data enhancement. arXiv preprint arXiv:2403.15042, 2024. 1
- [21] Zhuowan Li, Bhavan Jasani, Peng Tang, and Shabnam Ghadar. Synthesize step-by-step: Tools templates and llms as data generators for reasoning-based chart vqa. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13613–13623, 2024. 1
- [22] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in neural information processing systems, 36:34892–34916, 2023. 1
- [23] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2024. 3, 4
- [24] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player?, 2024. 3
- [25] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The* 36th Conference on Neural Information Processing Systems (NeurIPS), 2022. 2, 1
- [26] 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. 2, 1
- [27] Minesh Mathew, Viraj Bagal, Rubèn Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. Infographicvqa.

In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 1697–1706, 2022. 2,

- [28] Rakesh R Menon and Shashank Srivastava. Discern: Decoding systematic errors in natural language for text classifiers. arXiv preprint arXiv:2410.22239, 2024. 1
- [29] Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Codas, Yadong Lu, Weige Chen, Olga Vrousgos, Corby Rosset, et al. Agentinstruct: Toward generative teaching with agentic flows. arXiv preprint arXiv:2407.03502, 2024. 1
- [30] OpenAI. Gpt-4o system card, 2024. 2
- [31] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In 2012 IEEE conference on computer vision and pattern recognition, pages 3498–3505. IEEE, 2012. 4
- [32] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing.* Association for Computational Linguistics, 2019. 3
- [33] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings* of the IEEE/CVF conference on computer vision and pattern recognition, pages 8317–8326, 2019. 3
- [34] Xin Su, Man Luo, Kris W Pan, Tien Pei Chou, Vasudev Lal, and Phillip Howard. Sk-vqa: Synthetic knowledge generation at scale for training context-augmented multimodal llms. *arXiv preprint arXiv:2406.19593*, 2024. 1
- [35] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikuła, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle

Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. 4

- [36] Solomon Ubani, Suleyman Olcay Polat, and Rodney Nielsen. Zeroshotdataaug: Generating and augmenting training data with chatgpt. arXiv preprint arXiv:2304.14334, 2023. 1
- [37] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018. 4
- [38] Pablo Villalobos, Jaime Sevilla, Lennart Heim, Tamay Besiroglu, Marius Hobbhahn, and Anson Ho. Will we run out of data? an analysis of the limits of scaling datasets in machine learning. *arXiv preprint arXiv:2211.04325*, 1, 2022.
- [39] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution, 2024. 4, 3
- [40] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashionmnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
 4
- [41] Fuzhao Xue, Yao Fu, Wangchunshu Zhou, Zangwei Zheng, and Yang You. To repeat or not to repeat: Insights from scaling llm under token-crisis. *Advances in Neural Information Processing Systems*, 36:59304–59322, 2023. 1
- [42] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report, 2024.
- [43] Yue Yang, Ajay Patel, Matt Deitke, Tanmay Gupta, Luca Weihs, Andrew Head, Mark Yatskar, Chris Callison-Burch, Ranjay Krishna, Aniruddha Kembhavi, et al. Scaling textrich image understanding via code-guided synthetic multimodal data generation. arXiv preprint arXiv:2502.14846, 2025. 1
- [44] Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander J Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. Large language model as attributed training data generator:

A tale of diversity and bias. *Advances in Neural Information Processing Systems*, 36:55734–55784, 2023. 1

- [45] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2024. 3
- [46] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. 2
- [47] Ingo Ziegler, Abdullatif Köksal, Desmond Elliott, and Hinrich Schütze. Craft your dataset: Task-specific synthetic dataset generation through corpus retrieval and augmentation. *arXiv preprint arXiv:2409.02098*, 2024. 1