

- 810 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
811 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer
812 Vision*, pp. 11975–11986, 2023.
- 813
- 814 Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. Investigating
815 the catastrophic forgetting in multimodal large language model fine-tuning. In *Conference on
816 Parsimony and Learning*, pp. 202–227. PMLR, 2024.
- 817 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models
818 are not robust multiple choice selectors. In *The Twelfth International Conference on Learning
819 Representations*, 2023a.
- 820
- 821 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
822 Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
823 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023b.
- 824
- 825 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
826 Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information
827 Processing Systems*, 36, 2024.
- 828
- 829 Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe
830 Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, et al. Sotopia: Interactive evaluation for
831 social intelligence in language agents. *arXiv preprint arXiv:2310.11667*, 2023a.
- 832
- 833 Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit
834 Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language
835 models. *arXiv preprint arXiv:2310.00754*, 2023b.

836 A DESCRIPTION OF EVALUATION BENCHMARKS

- 837
- 838 • **MM-Vet** (Yu et al., 2023) dataset is a benchmark designed to evaluate large vision-language
839 models (LVLMs) across six core vision-language (VL) capabilities: recognition, knowl-
840 edge, optical character recognition (OCR), spatial awareness, language generation, and
841 mathematical reasoning. The dataset includes open-ended, real-world questions based on
842 image-text pairs, requiring models to integrate multiple capabilities to solve complex tasks.
843 MM-Vet benchmark consists of 200 images paired with 218 open-ended questions.
 - 844 • **Q-Bench** (Wu et al., 2023) evaluates the capabilities of large vision-language models in
845 three main areas related to low-level vision tasks. These tasks focus on evaluating how
846 well LVLMs can perform basic low-level perception tasks that are traditionally associated
847 with human visual perception. In the Q-Bench dataset, the questions are of three types:
848 Yes-or-No, What, and How.
 - 849 – **Low-Level Visual Perception:** Assesses how accurately LVLMs can answer ques-
850 tions about low-level image attributes (e.g., clarity, color, distortion). LLVisionQA
851 dataset includes 2,990 images, each with a corresponding question about low-level
852 features.
 - 853 – **Low-Level Visual Description:** Evaluates the ability of LVLMs to describe images.
854 LLDescribe dataset has 499 images with expert-labeled descriptions averaging 58
855 words each. LVLMs are compared against these to assess completeness, preciseness,
856 and relevance.
 - 857 – **Visual Quality Assessment:** Evaluates LVLMs’ ability to predict quantifiable quality
858 scores for images by assessing how well they align with human-rated mean opinion
859 scores (MOS) on low-level visual appearances, using 81,284 samples.
 - 860 • **SQA-IMG** (Lu et al., 2022a) is a portion of the Science Question Answering (SQA) dataset
861 that contains questions from a wide range of scientific domains, each paired with corre-
862 sponding image contexts. The dataset includes 10,332 examples of multimodal multiple-
863 choice questions, along with lectures and explanations that detail the reasoning behind the
correct answers.

- 864 • **ChartQA** (Masry et al., 2022) dataset is a benchmark designed to test AI models on their
865 ability to perform question-answering tasks over various types of charts. It focuses specif-
866 ically on questions requiring complex reasoning, such as visual and logical interpretation,
867 going beyond simpler template-based datasets. ChartQA includes 9,608 human-authored
868 open-ended questions as well as 23,111 questions that are automatically generated from
869 chart summaries.
- 870 • **SEED-IMG** (Li et al., 2023), a subset of SEED-Bench, focuses on evaluating spatial com-
871 prehension of images by testing models on various dimensions like scene understanding,
872 object identification, and spatial relationships. In terms of scale, the dataset includes 19,000
873 multiple-choice questions that evaluate both image and video comprehension, covering 12
874 evaluation dimensions such as scene understanding, instance identity, spatial relations, and
875 action recognition.
- 876 • **MME** (Fu et al., 2023) evaluates both perception and cognition abilities of LVLMS. It
877 features 14 subtasks, including recognition tasks (such as object existence, count, position,
878 color) and reasoning tasks (such as commonsense reasoning, numerical calculation, text
879 translation, and code reasoning). MME uses manually created instruction-answer pairs,
880 ensuring no overlap with public datasets. MME uses "yes/no" responses for quantitative
881 evaluations.
- 882 • **MathVista** (Lu et al., 2023) is a benchmark designed to evaluate the mathematical rea-
883 soning capabilities of foundation models in visual contexts. It integrates challenges from
884 diverse mathematical and visual tasks, with a focus on fine-grained, deep visual under-
885 standing and compositional reasoning. MathVista consists of 6,141 examples including
886 3,392 multiple-choice questions and 2,749 free-form questions derived from 28 existing
887 multimodal datasets and 3 newly created datasets: IQTest, FunctionQA, and PaperQA.
- 888 • **LLaVA-W** (Liu et al., 2024c) is a challenging evaluation benchmark created to assess the
889 generalization and instruction-following capabilities LVLMS in complex, real-world sit-
890 uations. It consists of 24 images and 60 questions, including diverse scenes like indoor
891 environments, outdoor settings, memes, paintings, and sketches. Each image is associated
892 with a highly detailed and manually curated description, and the questions focus on extract-
893 ing intricate details and reasoning about the visual content. LLaVA-W involves a variety of
894 tasks, including detailed descriptions, conversational answers, and complex reasoning.
- 895 • **MMStar** (Chen et al., 2024a) is a vision-dependent multimodal benchmark designed to
896 evaluate the multimodal capabilities of LVLMS. It addresses two main issues identified
897 in previous benchmarks: the reliance on textual information without visual input and data
898 leakage during training. MMStar is composed of 1,500 samples carefully selected to en-
899 sure that visual content is necessary for solving each problem. MMStar evaluates six core
900 capabilities across 18 detailed axes, which include tasks like image perception and logical
901 reasoning. MMStar uses multiple-choice as the primary answer type.
- 902 • **MMVP** (Tong et al., 2024) evaluates the visual grounding capabilities of large vision-
903 language models by identifying scenarios where they fail to distinguish simple visual pat-
904 terns in images. These patterns include aspects like orientation, counting, viewpoint, and
905 relational context. The benchmark is constructed using 150 pairs of images, resulting in
906 300 multiple-choice questions.

908 B DESCRIPTION OF EVALUATION LVLMS

- 910 • **LLaVA-1.5** (Liu et al., 2024a) incorporates academic task-oriented datasets to enhance
911 performance in VQA tasks and features an MLP vision-language connector, which im-
912 proves upon the original linear layer utilized in LLaVA (Liu et al., 2024c). It uses CLIP
913 ViT-L/14 (Radford et al., 2021) with a 336px resolution as its vision encoder, resulting in
914 a total of $(336/14)^2 = 576$ visual tokens. LLaVA-1.5 is built on Vicuna with either 7B or
915 13B parameters. The training dataset includes 558K samples for pre-training and 665K for
916 fine-tuning, totaling 1.2M image-text pairs from publicly available datasets
- 917 • **LLaVA-NeXT** (Liu et al., 2024b) (also known as LLaVA-1.6) enhances visual reasoning,
OCR, and world knowledge, offering four times higher image resolution (up to 1344x336)

918 and improved performance in visual conversations. Its architecture includes a CLIP ViT-
919 L/14 as a vision encoder, paired with Vicuna models ranging from 7B to 34B as a back-
920 bone language model. It utilizes 1.3M visual instruction tuning data samples for training,
921 maintaining efficiency with approximately one day of training on 32 A100 GPUs. The arch-
922 itecture’s high resolution and dynamic grid scheme improve detailed image processing
923 capabilities.

- 924 • **LLaVA-OneVision** (Li et al., 2024b) is a LVLM designed for task transfer across single-
925 image, multi-image, and video scenarios, with strong capabilities in video understand-
926 ing through image-to-video task transfer. Its architecture consists of a Qwen2 language
927 model (Yang et al., 2024) with 8B to 72B parameters, and the SigLIP vision encoder (Zhai
928 et al., 2023), which processes images at a base resolution of 384x384, producing 729 visual
929 tokens. The model employs a 2-layer MLP projector. The training utilized 3.2M single-
930 image data samples and 1.6M multi-modal data samples, focusing on high-quality visual
931 instruction tuning data to enhance its multimodal capabilities.
- 932 • **Meteor** (Lee et al., 2024c) is a large vision-language model that uniquely embeds multi-
933 faceted rationales using a Mamba-based architecture (Gu & Dao, 2023), enabling efficient
934 processing of lengthy rationales to enhance its vision-language understanding. This ap-
935 proach allows Meteor to achieve superior performance without scaling up model size or
936 employing additional vision encoders. Its architecture includes a CLIP-L/14 vision en-
937 coder with an image resolution of 490x490, comprising 428M parameters, and InternLM2-
938 7B (Cai et al., 2024) as a foundational LLM. Meteor was trained on 2.1M question-answer
939 pairs, with 1.1M curated triples.
- 940 • **TroL** (Lee et al., 2024b) uses a unique characteristic called layer traversing, which reuses
941 layers in a token-wise manner, allowing it to simulate retracing the answering process with-
942 out physically adding more layers, making it efficient despite smaller model sizes. TroL
943 uses CLIP-L and InternViT as vision encoders, containing 428M and 300M parameters,
944 respectively, and supports 24 layers. The image resolution is adjusted using MLPs in the
945 vision projector. For its foundational LLM, TroL utilizes Phi-3-mini with 3.8B parameters
946 and InternLM2 with 1.8B and 7B parameters. The training dataset comprises 2.3M visual
947 instruction tuning samples.

948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971