# Benchmarking Large Multimodal Models against Common Corruptions

**Jiawei Zhang** [* 1]  **Tianyu Pang** [2]  **Chao Du** [2]  **Yi Ren** [† 3]  **Bo Li** [1 4]  **Min Lin** [2]

## Abstract

This technical report aims to fill a deficiency in the assessment of large multimodal models (LMMs) by specifically examining the self-consistency of their outputs when subjected to common corruptions. We investigate the cross-modal interactions between text, image, and speech, encompassing four essential generation tasks: *text-to-image*, *image-to-text*, *text-to-speech*, and *speech-to-text*. We create a comprehensive benchmark, named **MMCBench**, that covers **more than 100** popular LMMs (totally **over 150** model checkpoints). A thorough evaluation under common corruptions is critical for practical deployment and facilitates a better understanding of the reliability of cutting-edge LMMs. The benchmarking code is available at https://github.com/sail-sg/MMCBench.

## 1. Introduction

The introduction of the CLIP framework [89] has significantly accelerated recent advances in multimodal learning. The effective integration of visual and textual data by CLIP has not only marked a paradigm shift in machine learning's approach to diverse data modalities, but has also introduced unprecedented zero-shot prediction capabilities. Following it, models such as BLIP [59] and Stable Diffusion [97] expanded on these ideas, improving the interaction between visual and textual elements and setting new standards for high-quality generative modeling.

The field has advanced significantly in recent years, with notable advances in large language models (LLMs) [14, 80]. Consequently, large multimodal models (LMMs) such as InstructBLIP [25], LLaVA [64, 65], MiniGPT-4 [127], CogVLM [115], and ShareGPT4V [18], building on the foundations of LLMs like LLaMA [107, 108] or Vicuna [19, 125], have made remarkable progress in managing

cross-modality interactions, particularly between image and text. This progress has resulted in the ability to perform more complex tasks, such as contextual understanding, multimodal conversation, and visual question answering, among others. With these advancements and a greater emphasis on practical applications, the robustness and real-world efficacy of these models are becoming increasingly important.

For early-stage multimodal models that do not involve LLMs, MMRobustness [87] evaluates the performance impact of common corruptions on the input modality across five downstream tasks. Following that, MME [35] offers a comprehensive benchmark for evaluating LMMs, encompassing 14 subtasks in perception and cognition. Several efforts have also been made to investigate the possibility of manipulating LMMs via adversarial input images [10, 16, 30, 124]. Existing benchmarks, however, lack a comprehensive evaluation of the self-consistency of LMMs' outputs under common corruptions. Our goal is to bridge this gap by providing a comprehensive benchmark covering a wide range of popular LMMs.

To ensure efficient evaluation, we intend to exploit the diversity of large datasets such as LAION [99, 100] and Common Voice [3] while controlling the number of testing samples. To that end, we propose utilizing text as a semantic anchor. Because humans typically grasp complex concepts through text alone, this strategy is based on the belief that text adequately conveys semantic information across various modalities. We hope to select examples with significant textual changes for evaluation by mapping all modalities into text and calculating similarity as shown in the first row on Figure 1.

After selecting representative examples from the input modality, we subject them to various common corruptions. These corrupted inputs are then fed into LMMs to generate outputs. A critical aspect of our evaluation is the quantitative measurement of self-consistency in these models. Primarily, we focus on cross-modality similarity—if such a model exists between the input and output modalities, it is utilized to measure self-consistency. However, in scenarios lacking a cross-modality model or the performance of such model is poor, we assess consistency within the output modality itself. This method, while straightforward, has its limitations; for instance, if a model consistently outputs

---

*Work done during an associate membership at Sea AI Lab.
†Work done while at Sea AI Lab. [1]UIUC [2]Sea AI Lab [3]ByteDance [4]University of Chicago. Correspondence to: Jiawei Zhang <jiaweiz7@illinois.edu>, Tianyu Pang <tianyupang@sea.com>, Chao Du <duchao@sea.com>.
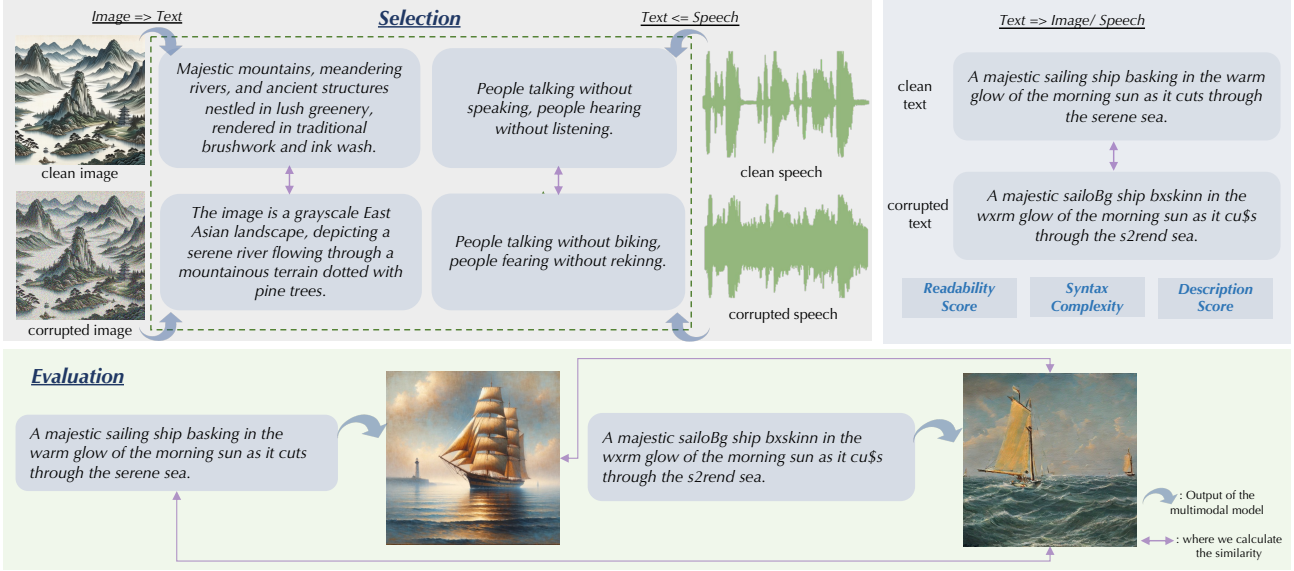
Technical report, work in progress.

*Figure 1.* Overview of the selection and evaluation process for cross-modality consistency. The selection process (*top row*) involves determining similarity based on text modality, using model-generated captions or transcriptions for non-text inputs, while directly comparing text inputs before and after corruption. The evaluation process (*bottom row*) measures self-consistency by comparing clean inputs with outputs from corrupted inputs and by comparing outputs from clean and corrupted inputs against each other.

poor-quality captions, it may still achieve a high consistency score on the output modality, which can be misleading. In our study, we present both cross-modality and output-only modality consistency scores for a comprehensive reference as shown in the second on Figure 1, allowing for a nuanced interpretation of model performance.

We mainly focus on cross-modality interactions involving text, image, and speech in this paper, encompassing four key generative tasks: *text-to-image*, *image-to-text*, *text-to-speech*, and *speech-to-text*. The released benchmark is named as **MMCBench** (**M**ulti**M**odal **C**orruption **Bench**mark). Our MMCBench performs extensive analyses of popular LMMs as well as traditional non-LLM-based multimodal models, providing insights into their robustness and self-consistency under various corruption scenarios. Specifically, we test against **23** text corruptions, **29** image corruptions, and **16** speech corruptions. For the text-to-image task, we assess **27** models across **37** different checkpoints; in the image-to-text category, **39** models are evaluated over **58** checkpoints; for text-to-speech, **14** models are scrutinized through **15** checkpoints; and for speech-to-text, we examine **41** models across **47** checkpoints. This extensive analysis offers deep insights into the resilience of these models in diverse and challenging conditions.

## 2. Related Work

**Robustness evaluation on unimodal models.** Evaluating the robustness of unimodal models, such as vision and lan-

guage models, is critical for their practical deployment.

Robustness is typically measured in the context of vision models along three key dimensions: *common corruption*, *adversarial robustness*, and *distribution shift robustness*. Datasets such as ImageNet-C [44] and 3DCC [51] are used to test whether classification models can maintain accurate predictions in the presence of common corruptions such as Gaussian noise, motion blur, and brightness changes. In terms of adversarial robustness, the emphasis shifts to how models respond to intentional perturbations. This can be evaluated using datasets like ImageNet-Patch [83], which involves attaching malicious patches to images, or by subjecting models to direct $\ell_p$-norm attacks, as assessed by ARES [29] and RobustBench [23]. The evaluation of distribution shift robustness tests the generalisation ability of vision models. Datasets like OOD-CV [123] and ImageNet-O [45] enable the evaluation of model performance on images that differ significantly from their training distribution. This evaluation helps assess the models' ability to adapt to new and unseen environments.

Similarly, in the context of language models, common corruption robustness can be evaluated by introducing, e.g., typos, word deletions, or whitespace insertions. On the other hand, adversarial robustness can be evaluated using model-based perturbations, as demonstrated by works such as BERT-Attack [62] and TextFooler [49]. Platforms including TextAttack [75], TextFlint [116], Robustness Gym [38], and NL-Augmenter [27] all provide convenient tools to fa-
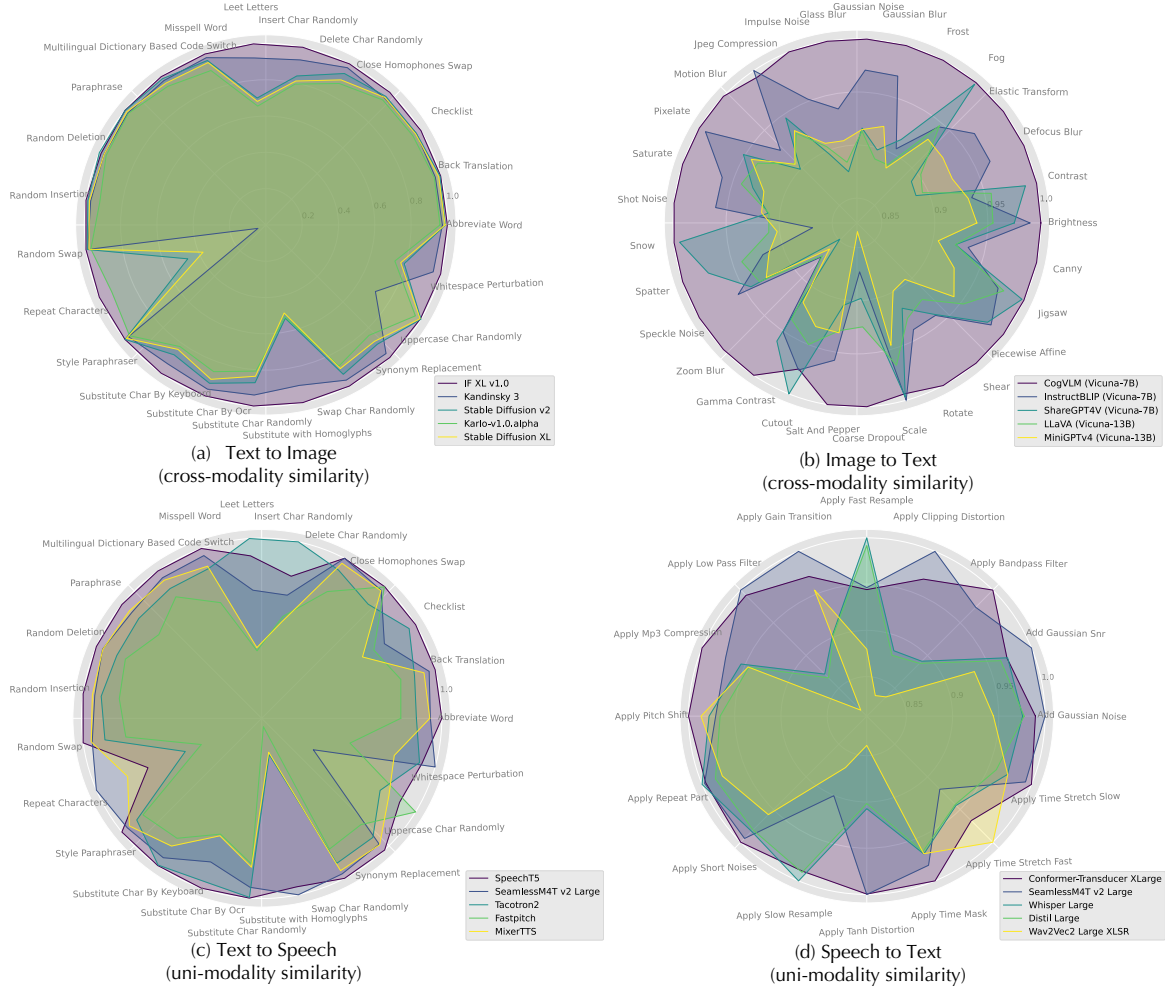
*Figure 2.* Radar charts depicting the relative consistency scores of five selected models for various corruptions across four cross-modality tasks: (a) text-to-image, (b) image-to-text, (c) text-to-speech, and (d) speech-to-text. The scores are normalized with the highest scoring model set as the baseline (ratio of 1) for each type of corruption, allowing for a comparative analysis of each model's resilience to specific corruptions relative to the best performer in that category.

cilitate these evaluations. The emergence of large language models (LLMs) has broadened the scope of robustness in language models. This expansion can be seen in platforms ranging from AdvGlue [111], which offers a comprehensive set of 14 textual adversarial attacks at the word or sentence level, challenging language models in tasks like sentiment classification or natural language inference, to PromptBench [128]. PromptBench focuses on adversarial prompts that are intended mimic potential user errors, thereby testing the consistency of LLM outputs. There are also platforms like HADES [66] and HaluEval [60], which are aimed at detecting hallucinations in the contents generated by LLMs. Recently, DecodingTrust [112] broadens its evaluation criteria to include toxicity, stereotype bias, privacy, machine ethics, and fairness.

**Robustness evaluation of multimodal models.** A number of benchmarks for multimodal model robustness have

been established, each focusing on a different aspect of model resilience [16, 124]. MMRobustness [87] primarily focuses on the relative performance drop observed when common corruptions are introduced to the input image or text. BenchLMM [15] instead investigates the robustness of large multimodal models (LMMs) such as GPT-4V [80] and LLaVA [64, 65]. It investigates how these models respond to stylistic shifts such as artistic, sensor, and application styles, thereby expanding our understanding of LMM resilience. The Bingo benchmark [24], on the other hand, focuses on hallucination phenomena in visual language models. It delves into the biases and interference issues that are prevalent in LMMs, providing a specialized perspective on these unique challenges. The introduction of the VLAA framework [109], which shifts the emphasis to out-of-distribution generalization and adversarial robustness, shedding light on critical vulnerabilities in LMMs. VLAA

*Table 1.* Performance comparison of various **text-to-image** models evaluated by self-consistency scores (cross-modality) across different corruption intensities and data selection levels. Scores represent the average multiplied cosine similarities (max 2300) between original captions and the generated images for the captions under different corruption conditions.

| Models | Hard | | Random | | Average |
|---|---|---|---|---|---|
| | Heavy | Light | Heavy | Light | |
| 🏆 IF XL v1.0 [26] | 630 | 722 | 577 | 679 | 652.00 |
| 🥈 IF L v1.0 [26] | 610 | 700 | 571 | 670 | 637.75 |
| 🥉 IF M v1.0 [26] | 584 | 666 | 564 | 661 | 618.75 |
| Kandinsky 3 [4] | 563 | 664 | 555 | 670 | 613.00 |
| Stable Diffusion v2 [97] | 549 | 666 | 532 | 657 | 601.00 |
| SDXL Base [84] | 529 | 649 | 527 | 651 | 589.00 |
| SDXL Turbo [98] | 521 | 637 | 533 | 658 | 587.25 |
| SDXL Refiner [84] | 526 | 645 | 522 | 645 | 584.50 |
| Karlo-v1.0.alpha [58] | 520 | 627 | 522 | 637 | 576.50 |
| Openjourney v4 [86] | 518 | 613 | 528 | 638 | 574.25 |
| LCM (SDXL) [69, 84] | 515 | 626 | 515 | 638 | 573.50 |
| LCM LoRA (SDXL) [70, 84] | 503 | 613 | 517 | 640 | 568.25 |
| Stable Diffusion Turbo [98] | 493 | 602 | 519 | 641 | 563.75 |
| Anything Midjourney v4.1 [86] | 491 | 595 | 510 | 623 | 554.75 |
| Dreamlike Photoreal 2.0 [31] | 491 | 590 | 513 | 623 | 554.25 |
| Stable Diffusion v1 [97] | 493 | 591 | 511 | 620 | 553.75 |
| Kandinsky 2.2 [101] | 488 | 588 | 510 | 627 | 553.25 |
| Small Stable Diffusion [52] | 473 | 565 | 518 | 625 | 545.25 |
| Unidiffuser [11] | 485 | 583 | 495 | 602 | 541.25 |
| LCM (SSD-1B) [69, 42] | 468 | 564 | 490 | 599 | 530.25 |
| LCM LoRA (SSD-1B) [70, 42] | 462 | 560 | 485 | 600 | 526.75 |
| SSD 1B [42] | 462 | 560 | 485 | 599 | 526.50 |
| LCM LoRA (SD v1) [70, 97] | 426 | 522 | 472 | 581 | 500.25 |
| Dreamshaper v7 [71] | 425 | 522 | 472 | 581 | 500.00 |
| LCM (Dreamshaper v7) [69, 71] | 417 | 507 | 472 | 577 | 493.25 |
| Lafite [126] | 415 | 480 | 474 | 558 | 481.75 |
| Glide [78] | 334 | 390 | 391 | 465 | 395.00 |

focuses on LMMs' responses to visually unrelated or linguistically perturbed inputs, as well as their resilience to adversarial attacks that threaten model safety and reliability.

Despite these comprehensive benchmarks, a significant gap in the field remains: the lack of tools specifically designed to measure the consistency of multimodal model output when confronted with commonly corrupted inputs. Our research aims to close this gap, making contribution to the larger landscape of multimodal model evaluation.

## 3. Experiments

This section describes our MMCBench in details, including the common corruption methods used for each modality, the strategies of data selection, and the evaluation metrics (i.e., quantification of consistency scores) for each multimodal

model. In general, two primary modalities are considered in multimodal model evaluation: the input modality and the output modality. Our approach is twofold: first, we identify a representative subset from large base datasets such as LAION [56, 100] for consistency evaluation; and second, we assess the consistency scores of multimodal models on the output modality when the input modality is corrupted.

**Data selection.** We use consistency scores in the text modality for ranking data samples for data selection, which is applicable to all cross-modalities, including image-to-text and speech-to-text. We calculate text similarities by mapping all data into the text modality and then rank and select samples with the lowest text similarity. This design choice is due to the critical role of text in representing semantic information across multiple modalities, as well as its utility in reducing computational complexity for selection. Fur-

*Table 2.* Performance comparison of various **image-to-text** models evaluated by self-consistency scores (cross-modality) across different corruption intensities and data selection levels. Scores represent the average multiplied cosine similarities (max 2900) between original images and the captions generated for the images under different corruption conditions.

| Models | LLM | Hard | | Random | | Average |
|---|---|---|---|---|---|---|
| | | Heavy | Light | Heavy | Light | |
| 🏆 CogVLM [115, 46] | Vicuna 7B [125] | 815 | 887 | 841 | 886 | 857.25 |
| 🥈 InstructBLIP [25] | Vicuna 7B [125] | 760 | 856 | 799 | 848 | 815.75 |
| 🥉 ShareGPT4V [18] | Vicuna 7B [125] | 724 | 843 | 798 | 853 | 804.50 |
| InstructBLIP [25] | Flan T5 XXL [20] | 751 | 841 | 774 | 820 | 796.50 |
| InstructBLIP [25] | Vicuna 13B [125] | 744 | 839 | 777 | 825 | 796.25 |
| LLaVA [65] | Vicuna 13B [125] | 733 | 836 | 783 | 827 | 794.75 |
| LLaVA-v1.5 [64] | Vicuna 13B [125] | 736 | 827 | 788 | 827 | 794.50 |
| LLaVA-v1.5 [64] | Vicuna 7B [125] | 728 | 821 | 781 | 820 | 787.50 |
| MiniGPT-4 [127] | Vicuna 13B [125] | 713 | 808 | 790 | 835 | 786.50 |
| BLIP2 [61] | OPT-2.7b [122] | 704 | 790 | 791 | 840 | 781.25 |
| MiniGPT-4 [127] | Vicuna 7B [125] | 707 | 801 | 780 | 827 | 778.75 |
| BLIP2 [61] | OPT-6.7b [122] | 703 | 792 | 784 | 833 | 778.00 |
| BLIP2 [61] | Flan T5 XXL [20] | 694 | 784 | 788 | 834 | 775.00 |
| BLIP2 [61] | Flan T5 XL [20] | 689 | 778 | 783 | 830 | 770.00 |
| LLaVA [64] | LLaMA2 13B [108] | 700 | 803 | 763 | 810 | 769.00 |
| VisualGLM [32, 28] | ChatGLM-6B [32, 121] | 717 | 793 | 753 | 786 | 762.25 |
| MiniGPT-4 [127] | LLaMA2 7B [108] | 685 | 767 | 773 | 810 | 758.75 |
| Qwen-VL-Chat [9] | Qwen-7B [8] | 631 | 803 | 738 | 844 | 754.00 |
| LLaVA [64] | LLaMA2 7B [108] | 687 | 786 | 747 | 796 | 754.00 |
| InstructBLIP [25] | Flan T5 XL [20] | 608 | 834 | 758 | 804 | 751.00 |
| LLaVA [64] | MPT 7B [105] | 681 | 778 | 744 | 791 | 748.50 |
| mPLUG-Owl2 [120] | LLaMA2 7B [108] | 608 | 746 | 722 | 806 | 720.50 |
| mPLUG-Owl [119] | LLaMA 7B [107] | 642 | 742 | 723 | 770 | 719.25 |
| LLaMA-Adapter v2 [37] | LLaMA 7B [107] | 646 | 744 | 712 | 759 | 715.25 |
| mPLUG-Owl (multilingual) [119] | LLaMA 7B [107] | 607 | 700 | 693 | 741 | 685.25 |
| GIT Large [114] | - | 581 | 683 | 696 | 753 | 678.25 |
| BLIP Large [59] | - | 558 | 686 | 692 | 774 | 677.50 |
| OpenFlamingo [6, 1] | RedPajama 3B [106] | 600 | 685 | 681 | 731 | 674.25 |
| LaVIN [67, 68] | LLaMA 13B [107] | 584 | 680 | 648 | 702 | 653.50 |
| Unidiffuser [11] | - | 539 | 624 | 682 | 758 | 650.75 |
| OpenFlamingo [6, 1] | MPT 1B [76] | 607 | 672 | 633 | 688 | 650.00 |
| BLIP Base [59] | - | 468 | 563 | 645 | 735 | 602.75 |
| GIT Base [114] | - | 468 | 546 | 630 | 698 | 585.50 |
| ImageBind-LLM [43] | Open Chinese LLaMA 7B [81, 107] | 489 | 569 | 613 | 654 | 581.25 |
| Multimodal-GPT [39] | LLaMA 7B [107] | 513 | 573 | 547 | 575 | 552.00 |
| OpenFlamingo[1] [6, 1] | MPT 7B [105] | 477 | 494 | 477 | 505 | 488.25 |
| ViT-GPT2 [79] | GPT2 [88] | 357 | 399 | 503 | 551 | 452.50 |
| MiniGPT v2[2] [17] | LLaMA2 7B [108] | 375 | 426 | 468 | 498 | 441.75 |

[1] OpenFlamingo's inferior performance in instance tasks is attributed to the fact that their utilized language model, MPT 7B, is not tuned for instructions. Therefor, sometimes the model fails to follow our instructions accurately.

[2] MiniGPTv2's inferior performance in this task is due to the model's tendency to avoid providing captions for corrupted images.

thermore, we include a comparative random sample from the base dataset. As a result, our data selection includes two levels: *hard*: selecting the subset with the lowest text similarity post-corruption; and *random*: randomly selecting a subset from the larger base dataset.

**Model evaluation.** We apply various common corruptions to the input modality at different intensities (*heavy* and *light*) for the evaluation process. The consistency between the output modality's generated samples under input corruption and the original uncorrupted data is then measured. We generate four distinct result categories based on the dual levels of dataset and corruption intensity. Models are ranked according to their overall performance in these categories. Additionally, we will report the best performance of each model, particularly in cases where multiple checkpoints are available, such as those trained on different datasets.

**Disclaimer**: In this study, models were evaluated using greedy decoding wherever feasible, primarily to ensure reproducibility and to reduce computational costs. However, it is important to note that greedy decoding may not always yield the optimal output for each model, potentially leading to an underestimation of the true performance capabilities. Consequently, the results presented in our tables should be regarded as *a conservative estimate or a lower bound* of the actual potential of the models.

Below, we introduce the common corruptions, data selection methods, and model evaluation criteria for each type of multimodal model in details.

### 3.1. Text-to-Image Generation

**Text corruptions.** We incorporate a total of 23 distinct types of text corruption, drawn from various platforms including NlpAug [72], TextAugment [73], and NL-Augmenter [27]. These corruptions are systematically categorized into three complexity levels, as follows:

- **Char Level**: *Substitute Char by OCR*, *Substitute Char by Keyboard*, *Insert Char Randomly*, *Substitute Char Randomly*, *Swap Char Randomly*, *Delete Char Randomly*, *Uppercase Char Randomly*, *Repeat Characters*, *Leet Letters* [33], *Whitespace Perturbation*, *Substitute with Homoglyphs*

- **Word Level**: *Synonym Replacement*, *Random Deletion*, *Random Swap*, *Random Insertion*, *Misspell Word*, *Abbreviate Word*, *Multilingual Dictionary Based Code Switch*, *Close Homophones Swap*

- **Sentence Level**: *CheckList* [96], *Back Translation* [77, 103], *Style Paraphraser* [54], *Paraphrase* [110]

Examples illustrating the corrupted text of each corruption type are provided in Appendix A.1.

**Data selection.** We start with a randomly selected subset of 10 million caption-image pairs from the LAION-COCO dataset [56], which originally contains 600 million images. This dataset is chosen for its large volume and high-quality COCO-style captions [63]. From these 10 million pairs, we then select $1,000$ captions that present more challenges in maintaining output consistency in *text-to-image* generation under various text corruptions. These captions are distinguished by their inherent complexity in language or syntax, as well as their richness in description. Specifically, we select them based on four scores:

1) *The inconsistency score* is calculated by one minus the average cosine text similarity between the original captions and their heavily corrupted versions across the 23 text corruption types. To compute text similarity, we utilize the `sentence-t5-large` model from sentence-transformer [93]. This metric quantifies the degree of semantic alteration in the captions due to corruption.

2) *The readability score* is calculated using `textstat`[1] to evaluate the inherent semantic complexity of the captions.

3) *The syntax complexity score* is measured by simply calculating the depth of the parse trees for each caption.

4) *The description score* assesses the content's richness by counting # of nouns, verbs, and adjectives in each caption.

The latter three scores are normalized to a range of $0$ to $1$, ensuring a balanced evaluation of caption quality and complexity. The overall selection score for each caption is computed as a weighted sum of the four components, with the inconsistency score carrying a weight of $6$, and the other scores each having a weight of $1$. During the initial selection process, however, we discovered that captions with the highest selection scores typically describe diagrams, charts, or posters. To avoid selecting captions that are not semantically rich, we add an additional filtering step based on aesthetic scores. We first choose the top $10,000$ captions based on their selection scores, and then choose the $1,000$ captions with the highest aesthetics scores from this subset using the improved aesthetic predictor[2]. This procedure ensures that the final caption selection is not only linguistically challenging, but also visually relevant and representative of the images. The image size is maintained at $256 \times 256$ pixels across all experiments.

**Evaluation methodology.** Self-consistency of text-to-image models in the face of specified corruptions is quantitatively assessed using a defined metric: the average cosine similarity between the original captions and the images generated after applying the 23 types of corruptions to the $1,000$ selected captions. This similarity score is then

---

[1] https://github.com/textstat/textstat
[2] https://github.com/christophschuhmann/improved-aesthetic-predictor

multiplied by 100, resulting in a maximum possible value of $2,300$. To calculate this cross-modality similarity, we utilize the CLIP model[3] trained on the LAION-2B dataset. A higher score in this metric indicates that a model is more capable of dealing with common corruptions. The final evaluation results on various text-to-image models are shown in Table 1, while we also provide the evaluation results on uni-modality similarity in Appendix B.1.

As observed, the self-consistency scores for models processing our carefully selected captions (hard) tend to be lower than those for randomly sampled captions (random) across identical corruption levels. This trend validates the effectiveness of our selection criteria, which is based solely on text degradation without any interference from text-to-image models during the selection phase. Notably, high-performing models such as IF [26], Kandinsky 3 [4], and Stable Diffusion v2 [97] exhibit a narrower performance gap between hard and random conditions, underscoring their robustness and superior handling of input corruptions.

### 3.2. Image-to-Text Generation

**Image corruptions.** In our robustness evaluation, we include a wide range of corruptions from the imagecorruptions library [74] which is based on ImageNet-C [44], and the imgaug library [50]. The selected types of corruptions are categorized as follows:

- **Noise-Related:** *Gaussian Noise*, *Shot Noise*, *Impulse Noise*, *Speckle Noise*

- **Blur-Related:** *Defocus Blur*, *Glass Blur*, *Motion Blur*, *Zoom Blur*, *Gaussian Blur*

- **Weather Conditions:** *Snow*, *Frost*, *Fog*

- **Digital:** *Brightness*, *Contrast*, *Pixelate*, *JPEG Compression*, *Spatter*, *Saturate*, *Gamma Contrast*

- **Arithmetic:** *Cutout*, *Salt and Pepper*, *Coarse Dropout*

- **Geometric:** *Scale*, *Rotate*, *Shear*, *Piecewise Affine*, *Jigsaw*

- **Edge:** *Canny*

This integration results in a total of 29 distinct and varied corruption methods for our comprehensive study. Examples illustrating the corrupted images of each corruption type are provided in Appendix A.2.

**Data selection.** We select a subset of 3 million images from the LAION-Aesthetics dataset[4] for our study. We intend to select $1,000$ images from this subset, focusing on those with high visual quality and inherently pose challenges for

multimodal models. The inconsistency score is determined by one minus the average cosine similarity between the text generated from the original uncorrupted images and the text generated from images subjected to the 15 common corruptions of ImageNet-C with severity level 3. The $1,000$ images with the highest inconsistency score are chosen. This assessment utilizes outputs from three baseline models: `vitgpt2` [79], `blip-base` [59], and `git-base` [114]. The image size is maintained at $384 \times 384$ pixels across all experiments.

**Evaluation methodology.** Similarly, the self-consistency of image-to-text models in response to specified corruptions is quantified using a metric based on the average sum of the cosine similarities between the original images and the captions generated for the chosen $1,000$ images under the 29 types of corruptions. This similarity score is then multiplied by 100, resulting in a maximum possible value of $2,900$. For calculating this cross-modality similarity, we employ the same CLIP model as used in Section 3.1. A higher score indicates a model's superior ability to deal with common corruptions. Note that for MLLMs, we consistently employed the instruction *"Describe this image as detailed as possible."* Detailed evaluation results for various image-to-text models are presented in Table 7, while the evaluation results on uni-modality similarity are deferred to Appendix B.2.

The comparative analysis of model performances reveals that the self-consistency scores for models evaluated on our carefully selected images are significantly lower than those for models tested with randomly selected images, under the same levels of corruption, which attests to the effectiveness of our selection methodology. Besides, as we can see, large multimodal models (LMMs) consistently outperform their counterparts that do not utilize LLMs, highlighting the added resilience that LLM integration confers in the face of corrupted inputs. However, a larger LLM size does not necessarily equate to higher consistency; instead, some models with smaller LLMs (e.g., 7B model) actually achieve comparable or even superior scores than the LMMs equipped with larger LLMs (e.g., 13B model).

### 3.3. Text-to-Speech Generation

**Text corruptions.** We use the same 23 text corruptions as previously described in Section 3.1, which is consistent with our text-to-image experiments.

**Data selection.** For this study, we utilize 1.75 million validated text-speech pairs from the Common Voice 15.0 dataset [3]. From this extensive collection, we extract a subset of $1,000$ pairs, each selected for its complexity and potential to challenge text-to-speech (TTS) models. Similar to text-to-image generation as described in Section 3.1, the selection process is still based on the four scores: inconsistency, readability, syntax complexity, and description. This

---

[3] https://huggingface.co/laion/CLIP-ViT-L-14-laion2B-s32B-b82K

[4] https://laion.ai/blog/laion-aesthetics/

*Table 3.* Performance of various **speech-to-text** models evaluated by self-consistency scores (uni-modality) across different corruption intensities and data selection levels. Scores represent the average multiplied cosine similarities (max 1600) between the transcriptions generated based on the original speech and the transcriptions generated for the speeches under different corruption conditions.

| Models | Hard | | Random | | Average |
|---|---|---|---|---|---|
| | Heavy | Light | Heavy | Light | |
| 🏆 Conformer-Transducer XLarge [41] | 1299 | 1488 | 1429 | 1551 | 1441.75 |
| 🥈 SeamlessM4T v2-Large [12] | 1281 | 1476 | 1420 | 1543 | 1430.00 |
| 🥉 FastConformer-Transducer XXLarge [94] | 1271 | 1477 | 1411 | 1549 | 1427.00 |
| Conformer-Transducer XXLarge [41] | 1274 | 1471 | 1414 | 1543 | 1425.50 |
| Conformer-CTC XLarge [41] | 1272 | 1464 | 1407 | 1538 | 1420.25 |
| Conformer-Transducer Large [41] | 1259 | 1457 | 1399 | 1537 | 1413.00 |
| SeamlessM4T Medium [13] | 1260 | 1450 | 1405 | 1531 | 1411.50 |
| FastConformer-CTC Large [94] | 1250 | 1458 | 1394 | 1537 | 1409.75 |
| Conformer-CTC Large [41] | 1251 | 1450 | 1394 | 1532 | 1406.75 |
| FastConformer-CTC XXLarge [94] | 1245 | 1452 | 1391 | 1534 | 1405.50 |
| FastConformer-Transducer Large [94] | 1245 | 1450 | 1390 | 1534 | 1404.75 |
| FastConformer-CTC XLarge [94] | 1237 | 1454 | 1384 | 1533 | 1402.00 |
| SeamlessM4T Large [13] | 1248 | 1442 | 1392 | 1525 | 1401.75 |
| FastConformer-Transducer XLarge [94] | 1228 | 1442 | 1376 | 1529 | 1393.75 |
| Whisper Large [91] | 1211 | 1414 | 1395 | 1530 | 1387.50 |
| Distil Large [36] | 1206 | 1402 | 1383 | 1521 | 1378.00 |
| Conformer-Transducer Medium [41] | 1203 | 1399 | 1355 | 1504 | 1365.25 |
| Conformer-CTC Medium [41] | 1203 | 1402 | 1353 | 1503 | 1365.25 |
| Whisper Medium [91] | 1175 | 1384 | 1361 | 1519 | 1359.75 |
| Distil Medium [36] | 1169 | 1376 | 1347 | 1509 | 1350.25 |
| Wav2Vec2 Large XLSR [21, 40] | 1191 | 1383 | 1304 | 1489 | 1341.75 |
| Conformer-CTC Small [41] | 1174 | 1361 | 1321 | 1477 | 1333.25 |
| Conformer-Transducer Small [41] | 1168 | 1360 | 1316 | 1477 | 1330.25 |
| Whisper Small [91] | 1141 | 1346 | 1323 | 1499 | 1327.25 |
| Distil Small [36] | 1141 | 1347 | 1319 | 1494 | 1325.25 |
| MMS 1B All [85] | 1179 | 1355 | 1296 | 1467 | 1324.25 |
| Hubert Large [47] | 1154 | 1335 | 1277 | 1465 | 1307.75 |
| Wav2Vec2 Conformer Rel Pos Large [113] | 1142 | 1321 | 1265 | 1459 | 1296.75 |
| Wav2Vec2 Conformer Rope Large [113] | 1142 | 1311 | 1261 | 1445 | 1289.75 |
| Robust Wav2Vec2 Large [48] | 1148 | 1311 | 1259 | 1440 | 1289.50 |
| Wav2Vec2 Large Self-Training [7, 118] | 1151 | 1298 | 1262 | 1432 | 1285.75 |
| S2T Small [113] | 1182 | 1273 | 1250 | 1364 | 1267.25 |
| Whisper Base [91] | 1084 | 1260 | 1263 | 1450 | 1264.25 |
| Wav2Vec2 Large [7] | 1124 | 1244 | 1228 | 1390 | 1246.50 |
| Whisper Tiny [91] | 1084 | 1219 | 1220 | 1404 | 1231.75 |
| Wav2Vec2 Base [7] | 1118 | 1207 | 1212 | 1346 | 1220.75 |
| ESPnet [117, 5, 90] | 1086 | 1232 | 1190 | 1364 | 1218.00 |
| S2T Medium [113] | 1107 | 1211 | 1205 | 1328 | 1212.75 |
| SpeechT5 [2] | 1077 | 1195 | 1198 | 1359 | 1207.25 |
| S2T Large [113] | 1087 | 1200 | 1188 | 1326 | 1200.25 |

[1] The Conformer CTC/Transducers and FastConformer CTC/Transducers models here are sourced from NeMo [55].

*Table 4.* Performance comparison of various **text-to-speech** models evaluated by self-consistency scores (uni-modality) across different corruption intensities and data selection levels. Scores represent the average multiplied cosine similarities (max 2300) between the transcriptions for the speeches generated for the clean transcriptions and the transcriptions for the speeches generated for the transcriptions under different corruption conditions. The name in parentheses indicates the source from which the checkpoint for this model was obtained.

| Models | Hard | | Random | | Average |
|---|---|---|---|---|---|
| | Heavy | Light | Heavy | Light | |
| 🏆 SpeechT5 [2] | 2023 | 2130 | 2007 | 2115 | 2068.75 |
| 🥈 SeamlessM4T v2 Large [12] | 2007 | 2118 | 1989 | 2101 | 2053.75 |
| 🥉 Tacotron2 (SpeechBrain) [102, 92] | 2008 | 2111 | 1981 | 2083 | 2045.75 |
| Fastpitch (NeMo) [57, 55] | 1996 | 2103 | 1975 | 2083 | 2039.25 |
| MixerTTS (NeMo) [104, 55] | 1994 | 2101 | 1973 | 2082 | 2037.50 |
| Tacotron2 (NeMo) [102, 55] | 1989 | 2098 | 1972 | 2080 | 2034.75 |
| MixerTTS-X (NeMo)[104, 55] | 1985 | 2091 | 1965 | 2072 | 2028.25 |
| FastSpeech2 (Facebook) [95, 113] | 1976 | 2076 | 1952 | 2053 | 2014.25 |
| XTTS v2 [22] | 1968 | 2068 | 1931 | 2027 | 1998.50 |
| SeamlessM4T Large [13] | 1951 | 2056 | 1932 | 2043 | 1995.50 |
| ESPnet2 [117] | 1963 | 2053 | 1928 | 2016 | 1990.00 |
| VITS (NeMo) [53, 55] | 1947 | 2043 | 1921 | 2014 | 1981.25 |
| MMS [85] | 1949 | 2039 | 1918 | 2007 | 1978.25 |
| SeamlessM4T Medium [13] | 1875 | 1974 | 1880 | 1987 | 1929.00 |

multifaceted approach ensures a comprehensive evaluation by focusing on pairs that present significant linguistic and acoustic challenges to TTS models.

**Evaluation methodology.** When assessing the self-consistency of speech-to-text (STT) models, we initially considered cross-modality similarity using the CLAP model [34] from LAION. However, it was determined to be less representative of true performance, given that CLAP's training predominantly involved classification-focused audio data, not speech. Consequently, we instead mainly focus to unimodality similarity here. This is measured as the average cosine similarity between the clean original transcriptions and the transcriptions of speech generated from texts with 16 types of text corruptions. Transcription of the generated speech relies on the `wav2vec2-base-960h` [7], while the text similarity is still computed using the `sentence-t5-large` model from the sentence-transformer library. The resulting similarity scores are multiplied by 100, yielding a maximum potential score of 2, 300. This method accounts for potential variations in speaker voices across STT models and places the emphasis on transcription consistency. The final evaluation results are shown in Table 4; while for self-consistency scores on cross modality, we will include enhanced models for more accurate cross-similarity calculations between speech and text in the later version.

The results shows that the level of corruption has a more

pronounced impact on model performance compared to the selection method used for evaluation data. One potential explanation is that varying corruption levels in the input transcriptions may lead to differences in only certain segments of the generated speech. Consequently, the transcriptions of these speeches could differ by only a few characters or words. While the current models used for calculating text similarity are likely more attuned to semantic discrepancies rather than minor character or word changes, which could account for the observed differences in performance. We plan to propose more sensitive measures of similarity that can better capture these subtle variations in future work.

### 3.4. Speech-to-Text Generation

**Speech corruptions.** Our study incorporates 16 distinct audio corruptions from audiomentations[5], which are categorized as follows:

- **Noise Additions and Interference:** *Gaussian Noise*, *Short Noises* [6], *Gaussian SNR*

- **Filtering and Frequency Adjustments:** *Bandpass Filter*, *Low Pass Filter*

- **Distortion and Audio Quality Effects:** *Clipping Distortion*, *MP3 Compression*, *Tanh Distortion*

---

[5]https://github.com/iver56/audiomentations
[6]The short noise sounds are sampled from ESC-50 dataset [82].

- **Temporal and Speed Modifications:** *Fast Resample*, *Slow Resample*, *Time Stretch (Fast)*, *Time Stretch (Slow)*

- **Pitch and Dynamic Range Adjustments:** *Pitch Shift*, *Gain Transition*

- **Repetitive and Temporal Effects:** *Repeat Part*, *Time Mask*

**Data selection.** We use the Common Voice 15.0 dataset [3], which includes around 1.75 million validated text-speech pairs and is known for its diversity. As a result, it is an ideal choice for speech-to-text (STT) analysis. We chose $1,000$ speeches from this large collection to demonstrate the difficulties in STT processing. The selection process relies on the average cosine similarity between texts generated from corrupted and original audio, using baseline models such as `speecht5_asr` [2], `wav2vec2-base-960h` [7], and `whisper-base.en` [91]. We prioritize speeches with the greatest drop in text similarity after corruption. The original sampling rate for the audio is maintained at 16,000 Hz (16 kHz), which is a standard configuration for Automatic Speech Recognition (ASR) systems.

**Evaluation methodology.** For the speech-to-text (STT) models, our evaluation is still focused in uni-modality similarity, focusing on the consistency of the generated transcriptions amid speech corruptions. Specifically, we calculate the average cosine similarity between transcriptions generated from clean speech and those produced after applying 16 distinct types of speech corruptions. These similarity scores are multiplied by 100, allowing for a maximum score of $1,600$. A higher score suggests a model's increased robustness against speech corruption. The text similarity here is determined using the `sentence-t5-large` model from the sentence-transformer library. The corresponding results are presented in Table 3; while for self-consistency scores on cross modality, we will include enhanced models for more accurate cross-similarity calculations between speech and text in the later version.

From the results, we can see that the scores for our selectively chosen speech samples are consistently lower than those for randomly sampled speech under the same levels of corruption, confirming the challenging nature of our selected samples across various TTS models.

## 4. Discussion

Returning self-consistent outputs in the face of input corruption is an essential requirement for the practical use of LMMs. We position MMCBench as a lightweight and comprehensive benchmark that concentrates on common corruptions in the multimodal learning literature. We intend to constantly update new models and incorporate more

modalities, such as video, into MMCBench. Furthermore, future iterations of our research will aim to develop a more effective method for evaluating cross-modality similarity between speech and text. Regular updates and progress will be documented in our GitHub repository to facilitate community engagement and contributions.

## References

[1] Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., Ring, R., Rutherford, E., Cabi, S., Han, T., Gong, Z., Samangooei, S., Monteiro, M., Menick, J., Borgeaud, S., Brock, A., Nematzadeh, A., Sharifzadeh, S., Binkowski, M., Barreira, R., Vinyals, O., Zisserman, A., and Simonyan, K. Flamingo: a visual language model for few-shot learning. *ArXiv*, abs/2204.14198, 2022.

[2] Ao, J., Wang, R., Zhou, L., Wang, C., Ren, S., Wu, Y., Liu, S., Ko, T., Li, Q., Zhang, Y., et al. Speecht5: Unified-modal encoder-decoder pre-training for spoken language processing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5723–5738, 2022.

[3] Ardila, R., Branson, M., Davis, K., Kohler, M., Meyer, J., Henretty, M., Morais, R., Saunders, L., Tyers, F., and Weber, G. Common voice: A massively-multilingual speech corpus. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 4218–4222, 2020.

[4] Arkhipin, V., Filatov, A., Vasilev, V., Maltseva, A., Azizov, S., Pavlov, I., Agafonova, J., Kuznetsov, A., and Dimitrov, D. Kandinsky 3.0 technical report, 2023.

[5] Arora, S., Futami, H., weon Jung, J., Peng, Y., Sharma, R., Kashiwagi, Y., Tsunoo, E., and Watanabe, S. Universlu: Universal spoken language understanding for diverse classification and sequence generation tasks with a single network, 2023.

[6] Awadalla, A., Gao, I., Gardner, J., Hessel, J., Hanafy, Y., Zhu, W., Marathe, K., Bitton, Y., Gadre, S., Sagawa, S., Jitsev, J., Kornblith, S., Koh, P. W., Ilharco, G., Wortsman, M., and Schmidt, L. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023.

[7] Baevski, A., Zhou, Y., Mohamed, A., and Auli, M. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural*

*information processing systems*, 33:12449–12460, 2020.

[8] Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

[9] Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.

[10] Bailey, L., Ong, E., Russell, S., and Emmons, S. Image hijacks: Adversarial images can control generative models at runtime. *arXiv preprint arXiv:2309.00236*, 2023.

[11] Bao, F., Nie, S., Xue, K., Li, C., Pu, S., Wang, Y., Yue, G., Cao, Y., Su, H., and Zhu, J. One transformer fits all distributions in multi-modal diffusion at scale. *arXiv preprint arXiv:2303.06555*, 2023.

[12] Barrault, L., Chung, Y.-A., Meglioli, M. C., Dale, D., Dong, N., Duppenthaler, M., Duquenne, P.-A., Ellis, B., Elsahar, H., Haaheim, J., et al. Seamless: Multilingual expressive and streaming speech translation. *arXiv preprint arXiv:2312.05187*, 2023.

[13] Barrault, L., Chung, Y.-A., Meglioli, M. C., Dale, D., Dong, N., Duquenne, P.-A., Elsahar, H., Gong, H., Heffernan, K., Hoffman, J., et al. Seamlessm4t-massively multilingual & multimodal machine translation. *arXiv preprint arXiv:2308.11596*, 2023.

[14] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[15] Cai, R., Song, Z., Guan, D., Chen, Z., Luo, X., Yi, C., and Kot, A. Benchlmm: Benchmarking cross-style visual capability of large multimodal models. *arXiv preprint arXiv:2312.02896*, 2023.

[16] Carlini, N., Nasr, M., Choquette-Choo, C. A., Jagielski, M., Gao, I., Awadalla, A., Koh, P. W., Ippolito, D., Lee, K., Tramer, F., et al. Are aligned neural networks adversarially aligned? *arXiv preprint arXiv:2306.15447*, 2023.

[17] Chen, J., Zhu, D., Shen, X., Li, X., Liu, Z., Zhang, P., Krishnamoorthi, R., Chandra, V., Xiong, Y., and Elhoseiny, M. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023.

[18] Chen, L., Li, J., Dong, X., Zhang, P., He, C., Wang, J., Zhao, F., and Lin, D. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.

[19] Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https://lmsys.org/blog/2023-03-30-vicuna/.

[20] Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.

[21] Conneau, A., Baevski, A., Collobert, R., Mohamed, A., and Auli, M. Unsupervised cross-lingual representation learning for speech recognition. *arXiv preprint arXiv:2006.13979*, 2020.

[22] coqui. Xtts-v2. https://huggingface.co/coqui/XTTS-v2/tree/main, 2023.

[23] Croce, F., Andriushchenko, M., Sehwag, V., Debenedetti, E., Flammarion, N., Chiang, M., Mittal, P., and Hein, M. Robustbench: a standardized adversarial robustness benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

[24] Cui, C., Zhou, Y., Yang, X., Wu, S., Zhang, L., Zou, J., and Yao, H. Holistic analysis of hallucination in gpt-4v (ision): Bias and interference challenges. *arXiv preprint arXiv:2311.03287*, 2023.

[25] Dai, W., Li, J., Li, D., Tiong, A. M. H., Zhao, J., Wang, W., Li, B., Fung, P., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.

[26] deep floyd. If. https://github.com/deep-floyd/IF, 2023.

[27] Dhole, K. D., Gangal, V., Gehrmann, S., Gupta, A., Li, Z., Mahamood, S., Mahendiran, A., Mille, S., Shrivastava, A., Tan, S., et al. Nl-augmenter: A framework for task-sensitive natural language augmentation. *arXiv preprint arXiv:2112.02721*, 2021.

[28] Ding, M., Yang, Z., Hong, W., Zheng, W., Zhou, C., Yin, D., Lin, J., Zou, X., Shao, Z., Yang, H., et al. Cogview: Mastering text-to-image generation via transformers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.

[29] Dong, Y., Fu, Q.-A., Yang, X., Pang, T., Su, H., Xiao, Z., and Zhu, J. Benchmarking adversarial robustness on image classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[30] Dong, Y., Chen, H., Chen, J., Fang, Z., Yang, X., Zhang, Y., Tian, Y., Su, H., and Zhu, J. How robust is google's bard to adversarial image attacks? *arXiv preprint arXiv:2309.11751*, 2023.

[31] Dreamlike Art. Dreamlike photoreal 2.0. `https://huggingface.co/dreamlike-art/dreamlike-photoreal-2.0`, 2023.

[32] Du, Z., Qian, Y., Liu, X., Ding, M., Qiu, J., Yang, Z., and Tang, J. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 320–335, 2022.

[33] Eger, S., Sahin, G. G., Rücklé, A., Lee, J., Schulz, C., Mesgar, M., Swarnkar, K., Simpson, E., and Gurevych, I. Text processing like humans do: Visually attacking and shielding NLP systems. *CoRR*, abs/1903.11508, 2019. URL `http://arxiv.org/abs/1903.11508`.

[34] Elizalde, B., Deshmukh, S., Al Ismail, M., and Wang, H. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.

[35] Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Yang, J., Zheng, X., Li, K., Sun, X., et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.

[36] Gandhi, S., von Platen, P., and Rush, A. M. Distilwhisper: Robust knowledge distillation via large-scale pseudo labelling, 2023.

[37] Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., Li, H., and Qiao, Y. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023.

[38] Goel, K., Rajani, N. F., Vig, J., Taschdjian, Z., Bansal, M., and Ré, C. Robustness gym: Unifying the nlp evaluation landscape. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Demonstrations*, pp. 42–55, 2021.

[39] Gong, T., Lyu, C., Zhang, S., Wang, Y., Zheng, M., Zhao, Q., Liu, K., Zhang, W., Luo, P., and Chen, K. Multimodal-gpt: A vision and language model for dialogue with humans, 2023.

[40] Grosman, J. Fine-tuned XLSR-53 large model for speech recognition in English. `https://huggingface.co/jonatasgrosman/wav2vec2-large-xlsr-53-english`, 2021.

[41] Gulati, A., Qin, J., Chiu, C.-C., Parmar, N., Zhang, Y., Yu, J., Han, W., Wang, S., Zhang, Z., Wu, Y., et al. Conformer: Convolution-augmented transformer for speech recognition. *Interspeech 2020*, 2020.

[42] Gupta, Y., Jaddipal, V. V., Prabhala, H., Paul, S., and Platen, P. V. Progressive knowledge distillation of stable diffusion xl using layer level loss, 2024.

[43] Han, J., Zhang, R., Shao, W., Gao, P., Xu, P., Xiao, H., Zhang, K., Liu, C., Wen, S., Guo, Z., et al. Imagebind-llm: Multi-modality instruction tuning. *arXiv preprint arXiv:2309.03905*, 2023.

[44] Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations (ICLR)*, 2018.

[45] Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15262–15271, 2021.

[46] Hong, W., Wang, W., Lv, Q., Xu, J., Yu, W., Ji, J., Wang, Y., Wang, Z., Dong, Y., Ding, M., and Tang, J. Cogagent: A visual language model for gui agents, 2023.

[47] Hsu, W.-N., Bolte, B., Tsai, Y.-H. H., Lakhotia, K., Salakhutdinov, R., and Mohamed, A. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29: 3451–3460, 2021.

[48] Hsu, W.-N., Sriram, A., Baevski, A., Likhomanenko, T., Xu, Q., Pratap, V., Kahn, J., Lee, A., Collobert, R., Synnaeve, G., et al. Robust wav2vec 2.0: Analyzing domain shift in self-supervised pre-training. *arXiv preprint arXiv:2104.01027*, 2021.

[49] Jin, D., Jin, Z., Zhou, J. T., and Szolovits, P. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, 2020.

[50] Jung, A. B., Wada, K., Crall, J., Tanaka, S., Graving, J., Reinders, C., Yadav, S., Banerjee, J., Vecsei, G., Kraft, A., Rui, Z., Borovec, J., Vallentin, C., Zhydenko, S., Pfeiffer, K., Cook, B., Fernández, I., De Rainville, F.-M., Weng, C.-H., Ayala-Acevedo, A., Meudec, R., Laporte, M., et al. imgaug. https://github.com/aleju/imgaug, 2020. Online; accessed 01-Feb-2020.

[51] Kar, O. F., Yeo, T., Atanov, A., and Zamir, A. 3d common corruptions and data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18963–18974, 2022.

[52] Kim, B.-K., Song, H.-K., Castells, T., and Choi, S. On architectural compression of text-to-image diffusion models. *arXiv preprint arXiv:2305.15798*, 2023.

[53] Kim, J., Kong, J., and Son, J. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *International Conference on Machine Learning*, pp. 5530–5540. PMLR, 2021.

[54] Krishna, K., Wieting, J., and Iyyer, M. Reformulating unsupervised style transfer as paraphrase generation. In *Empirical Methods in Natural Language Processing*, 2020.

[55] Kuchaiev, O., Li, J., Nguyen, H., Hrinchuk, O., Leary, R., Ginsburg, B., Kriman, S., Beliaev, S., Lavrukhin, V., Cook, J., et al. Nemo: a toolkit for building ai applications using neural modules. *arXiv preprint arXiv:1909.09577*, 2019.

[56] Laion. Laion-COCO: A large-scale custom object dataset for computer vision, 2022. URL https://laion.ai/blog/laion-coco/.

[57] Łańcucki, A. Fastpitch: Parallel text-to-speech with pitch prediction. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6588–6592. IEEE, 2021.

[58] Lee, D., Kim, J., Choi, J., Kim, J., Byeon, M., Baek, W., and Kim, S. Karlo-v1.0.alpha on coyo-100m and cc15m. https://github.com/kakaobrain/karlo, 2022.

[59] Li, J., Li, D., Xiong, C., and Hoi, S. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pp. 12888–12900. PMLR, 2022.

[60] Li, J., Cheng, X., Zhao, W. X., Nie, J.-Y., and Wen, J.-R. Halueval: A large-scale hallucination evaluation benchmark for large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 6449–6464, 2023.

[61] Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.

[62] Li, L., Ma, R., Guo, Q., Xue, X., and Qiu, X. Bert-attack: Adversarial attack against bert using bert. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6193–6202, 2020.

[63] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.

[64] Liu, H., Li, C., Li, Y., and Lee, Y. J. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023.

[65] Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.

[66] Liu, T., Zhang, Y., Brockett, C., Mao, Y., Sui, Z., Chen, W., and Dolan, B. A token-level reference-free hallucination detection benchmark for free-form text generation. *arXiv preprint arXiv:2104.08704*, 2021.

[67] Luo, G., Huang, M., Zhou, Y., Sun, X., Jiang, G., Wang, Z., and Ji, R. Towards efficient visual adaption via structural re-parameterization. *arXiv preprint arXiv:2302.08106*, 2023.

[68] Luo, G., Zhou, Y., Ren, T., Chen, S., Sun, X., and Ji, R. Cheap and quick: Efficient vision-language instruction tuning for large language models. *arXiv preprint arXiv:2305.15023*, 2023.

[69] Luo, S., Tan, Y., Huang, L., Li, J., and Zhao, H. Latent consistency models: Synthesizing high-resolution images with few-step inference, 2023.

[70] Luo, S., Tan, Y., Patil, S., Gu, D., von Platen, P., Passos, A., Huang, L., Li, J., and Zhao, H. Lcm-lora: A universal stable-diffusion acceleration module. *arXiv preprint arXiv:2311.05556*, 2023.

[71] Lykon. dreamshaper-7. https://huggingface.co/Lykon/dreamshaper-7, 2023.

[72] Ma, E. Nlp augmentation. https://github.com/makcedward/nlpaug, 2019.

[73] Marivate, V. and Sefara, T. Improving short text classification through global augmentation methods. In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, pp. 385–399. Springer, 2020.

[74] Michaelis, C., Mitzkus, B., Geirhos, R., Rusak, E., Bringmann, O., Ecker, A. S., Bethge, M., and Brendel, W. Benchmarking robustness in object detection: Autonomous driving when winter is coming. *arXiv preprint arXiv:1907.07484*, 2019.

[75] Morris, J. X., Lifland, E., Yoo, J. Y., Grigsby, J., Jin, D., and Qi, Y. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. *arXiv preprint arXiv:2005.05909*, 2020.

[76] mosaicml. mpt-1b-redpajama-200b. https://huggingface.co/mosaicml/mpt-1b-redpajama-200b, 2023.

[77] Ng, N., Yee, K., Baevski, A., Ott, M., Auli, M., and Edunov, S. Facebook fair's wmt19 news translation task submission. *arXiv preprint arXiv:1907.06616*, 2019.

[78] Nichol, A. Q., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., Mcgrew, B., Sutskever, I., and Chen, M. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pp. 16784–16804. PMLR, 2022.

[79] NLP Connect. vit-gpt2-image-captioning (revision 0e334c7), 2022. URL https://huggingface.co/nlpconnect/vit-gpt2-image-captioning.

[80] OpenAI. Gpt-4 technical report, 2023.

[81] openlmlab. open-chinese-llama-7b-patch. https://huggingface.co/openlmlab/open-chinese-llama-7b-patch, 2023.

[82] Piczak, K. J. Esc: Dataset for environmental sound classification. In *Proceedings of the 23rd ACM international conference on Multimedia*, pp. 1015–1018, 2015.

[83] Pintor, M., Angioni, D., Sotgiu, A., Demetrio, L., Demontis, A., Biggio, B., and Roli, F. Imagenet-patch: A dataset for benchmarking machine learning robustness against adversarial patches. *Pattern Recognition*, 134:109064, 2023.

[84] Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

[85] Pratap, V., Tjandra, A., Shi, B., Tomasello, P., Babu, A., Kundu, S., Elkahky, A., Ni, Z., Vyas, A., Fazel-Zarandi, M., Baevski, A., Adi, Y., Zhang, X., Hsu, W.-N., Conneau, A., and Auli, M. Scaling speech technology to 1,000+ languages. *arXiv*, 2023.

[86] promputhero. openjourney-v4. https://huggingface.co/promputhero/openjourney-v4, 2023.

[87] Qiu, J., Zhu, Y., Shi, X., Wenzel, F., Tang, Z., Zhao, D., Li, B., and Li, M. Are multimodal models robust to image and text perturbations? *arXiv preprint arXiv:2212.08044*, 2022.

[88] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

[89] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

[90] Radford, A., Kim, J. W., Xu, T., Brockman, G., Mcleavey, C., and Sutskever, I. Robust speech recognition via large-scale weak supervision. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 28492–28518. PMLR, 23–29 Jul 2023.

[91] Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., and Sutskever, I. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pp. 28492–28518. PMLR, 2023.

[92] Ravanelli, M., Parcollet, T., Plantinga, P., Rouhe, A., Cornell, S., Lugosch, L., Subakan, C., Dawalatabad, N., Heba, A., Zhong, J., Chou, J.-C., Yeh, S.-L., Fu, S.-W., Liao, C.-F., Rastorgueva, E., Grondin, F., Aris, W., Na, H., Gao, Y., Mori, R. D., and Bengio, Y. SpeechBrain: A general-purpose speech toolkit, 2021. arXiv:2106.04624.

[93] Reimers, N. and Gurevych, I. 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, 11 2019. URL https://arxiv.org/abs/1908.10084.

[94] Rekesh, D., Kriman, S., Majumdar, S., Noroozi, V., Juang, H., Hrinchuk, O., Kumar, A., and Ginsburg, B. Fast conformer with linearly scalable attention for efficient speech recognition. *arXiv preprint arXiv:2305.05084*, 2023.

[95] Ren, Y., Hu, C., Tan, X., Qin, T., Zhao, S., Zhao, Z., and Liu, T.-Y. Fastspeech 2: Fast and high-quality end-to-end text to speech. In *International Conference on Learning Representations*, 2020.

[96] Ribeiro, M. T., Wu, T., Guestrin, C., and Singh, S. Beyond accuracy: Behavioral testing of NLP models with CheckList. In Jurafsky, D., Chai, J., Schluter, N., and Tetreault, J. (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4902–4912, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.442. URL https://aclanthology.org/2020.acl-main.442.

[97] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.

[98] Sauer, A., Lorenz, D., Blattmann, A., and Rombach, R. Adversarial diffusion distillation. *arXiv preprint arXiv:2311.17042*, 2023.

[99] Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., and Komatsuzaki, A. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.

[100] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al. Laion-5b: An open large-scale dataset for training next generation image-text models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.

[101] Shakhmatov, A., Razzhigaev, A., Nikolich, A., Arkhipin, V., Pavlov, I., Kuznetsov, A., and Dimitrov, D. kandinsky 2.2. https://huggingface.co/kandinsky-community/kandinsky-2-2-decoder, 2023.

[102] Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., Chen, Z., Zhang, Y., Wang, Y., Skerrv-Ryan, R., et al. Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 4779–4783. IEEE, 2018.

[103] Sugiyama, A. and Yoshinaga, N. Data augmentation using back-translation for context-aware neural machine translation. In Popescu-Belis, A., Loáiciga, S., Hardmeier, C., and Xiong, D. (eds.), *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pp. 35–44, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-6504. URL https://aclanthology.org/D19-6504.

[104] Tatanov, O., Beliaev, S., and Ginsburg, B. Mixer-tts: non-autoregressive, fast and compact text-to-speech model conditioned on language model embeddings. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7482–7486. IEEE, 2022.

[105] Team, M. et al. Introducing mpt-7b: a new standard for open-source, commercially usable llms, 2023.

[106] togethercomputer. Redpajama-incite-instruct-3b-v1. https://huggingface.co/togethercomputer/RedPajama-INCITE-Instruct-3B-v1, 2023.

[107] Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[108] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[109] Tu, H., Cui, C., Wang, Z., Zhou, Y., Zhao, B., Han, J., Zhou, W., Yao, H., and Xie, C. How many unicorns are in this image? a safety evaluation benchmark for vision llms. *arXiv preprint arXiv:2311.16101*, 2023.

[110] Vorobev, M. K. V. and Kuznetsov, M. A paraphrasing model based on chatgpt paraphrases. *A paraphrasing model based on ChatGPT paraphrases*, 2023.

[111] Wang, B., Xu, C., Wang, S., Gan, Z., Cheng, Y., Gao, J., Awadallah, A. H., and Li, B. Adversarial glue:

A multi-task benchmark for robustness evaluation of language models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

[112] Wang, B., Chen, W., Pei, H., Xie, C., Kang, M., Zhang, C., Xu, C., Xiong, Z., Dutta, R., Schaeffer, R., et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *arXiv preprint arXiv:2306.11698*, 2023.

[113] Wang, C., Tang, Y., Ma, X., Wu, A., Popuri, S., Okhonko, D., and Pino, J. Fairseq s2t: Fast speech-to-text modeling with fairseq. *arXiv preprint arXiv:2010.05171*, 2020.

[114] Wang, J., Yang, Z., Hu, X., Li, L., Lin, K., Gan, Z., Liu, Z., Liu, C., and Wang, L. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*, 2022.

[115] Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., Xu, J., Xu, B., Li, J., Dong, Y., Ding, M., and Tang, J. Cogvlm: Visual expert for pretrained language models, 2023.

[116] Wang, X., Liu, Q., Gui, T., Zhang, Q., Zou, Y., Zhou, X., Ye, J., Zhang, Y., Zheng, R., Pang, Z., et al. Textflint: Unified multilingual robustness evaluation toolkit for natural language processing. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pp. 347–355, 2021.

[117] Watanabe, S., Hori, T., Karita, S., Hayashi, T., Nishitoba, J., Unno, Y., Enrique Yalta Soplin, N., Heymann, J., Wiesner, M., Chen, N., Renduchintala, A., and Ochiai, T. ESPnet: End-to-end speech processing toolkit. In *Proceedings of Interspeech*, pp. 2207–2211, 2018. doi: 10.21437/Interspeech.2018-1456. URL http://dx.doi.org/10.21437/Interspeech.2018-1456.

[118] Xu, Q., Baevski, A., Likhomanenko, T., Tomasello, P., Conneau, A., Collobert, R., Synnaeve, G., and Auli, M. Self-training and pre-training are complementary for speech recognition. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3030–3034. IEEE, 2021.

[119] Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.

[120] Ye, Q., Xu, H., Ye, J., Yan, M., Liu, H., Qian, Q., Zhang, J., Huang, F., and Zhou, J. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. *arXiv preprint arXiv:2311.04257*, 2023.

[121] Zeng, A., Liu, X., Du, Z., Wang, Z., Lai, H., Ding, M., Yang, Z., Xu, Y., Zheng, W., Xia, X., Tam, W. L., Ma, Z., Xue, Y., Zhai, J., Chen, W., Liu, Z., Zhang, P., Dong, Y., and Tang, J. GLM-130b: An open bilingual pre-trained model. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. URL https://openreview.net/forum?id=-Aw0rrrPUF.

[122] Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., Mihaylov, T., Ott, M., Shleifer, S., Shuster, K., Simig, D., Koura, P. S., Sridhar, A., Wang, T., and Zettlemoyer, L. Opt: Open pre-trained transformer language models, 2022.

[123] Zhao, B., Yu, S., Ma, W., Yu, M., Mei, S., Wang, A., He, J., Yuille, A., and Kortylewski, A. Oodcv: a benchmark for robustness to out-of-distribution shifts of individual nuisances in natural images. In *European Conference on Computer Vision*, pp. 163–180. Springer, 2022.

[124] Zhao, Y., Pang, T., Du, C., Yang, X., Li, C., Cheung, N.-M., and Lin, M. On evaluating adversarial robustness of large vision-language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

[125] Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*, 2023.

[126] Zhou, Y., Zhang, R., Chen, C., Li, C., Tensmeyer, C., Yu, T., Gu, J., Xu, J., and Sun, T. Lafite: Towards language-free training for text-to-image generation. *arXiv preprint arXiv:2111.13792*, 2021.

[127] Zhu, D., Chen, J., Shen, X., Li, X., and Elhoseiny, M. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

[128] Zhu, K., Wang, J., Zhou, J., Wang, Z., Chen, H., Wang, Y., Yang, L., Ye, W., Gong, N. Z., Zhang, Y., et al. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts. *arXiv preprint arXiv:2306.04528*, 2023.

# A. Corrupted Examples

## A.1. Examples of Text Corruptions

Examples of the 23 distinct text corruptions are presented in Table 5.

## A.2. Examples of Image Corruptions

Examples of the 29 distinct image corruptions are illustrated in Figure 3.

# B. Additional Experimental Results

## B.1. Text to Image

The inconsistency score, calculated based on unimodality, is presented in Table 6. This score is derived from the comparison between the image generated from the clean input text and the images generated from the corrupted texts. The image similarity calculations are performed using the image encoder from the CLIP model [7]. We observe that the score for various models on our selected challenging dataset is significantly lower than that on randomly sampled data, which further demonstrate the effectiveness of our selection strategy. However, one limitation of this metric is that a model consistently yielding subpar results may still attain a high consistency score. For instance, despite its inferior quality in correlating the original caption with its generated images, the model ranked first, Lafite, manages to secure a high similarity score between its generated images. Additionally, it is noteworthy that the high-quality results of the IF model are contingent upon extensive computation, necessitating prolonged image generation times. In contrast, the LCM model, which operates with greater speed, manages to achieve a comparable level of consistency performance.

## B.2. Image to Text

The inconsistency score, based on unimodality, is presented in Table 7. This metric is derived from the semantic similarity between captions generated for clean input images and those generated for corrupted images. The similarity computations utilize the text encoder from the CLIP model referenced in Appendix B.1. The results show that models perform considerably better on our meticulously chosen challenging data than on randomly selected data, thereby confirming the efficacy of our selection strategy. Furthermore, it is evident that MLLMs still significantly outperform non-LLM-based models when dealing with corrupted image inputs, showcasing the robustness of LLMs in resisting the impact of such distortions.

---

[7] https://huggingface.co/laion/CLIP-ViT-L-14-laion2B-s32B-b82K

*Table 5.* Examples of our 23 distinct text corruptions categorized into three levels: character level, word level, and sentence level.

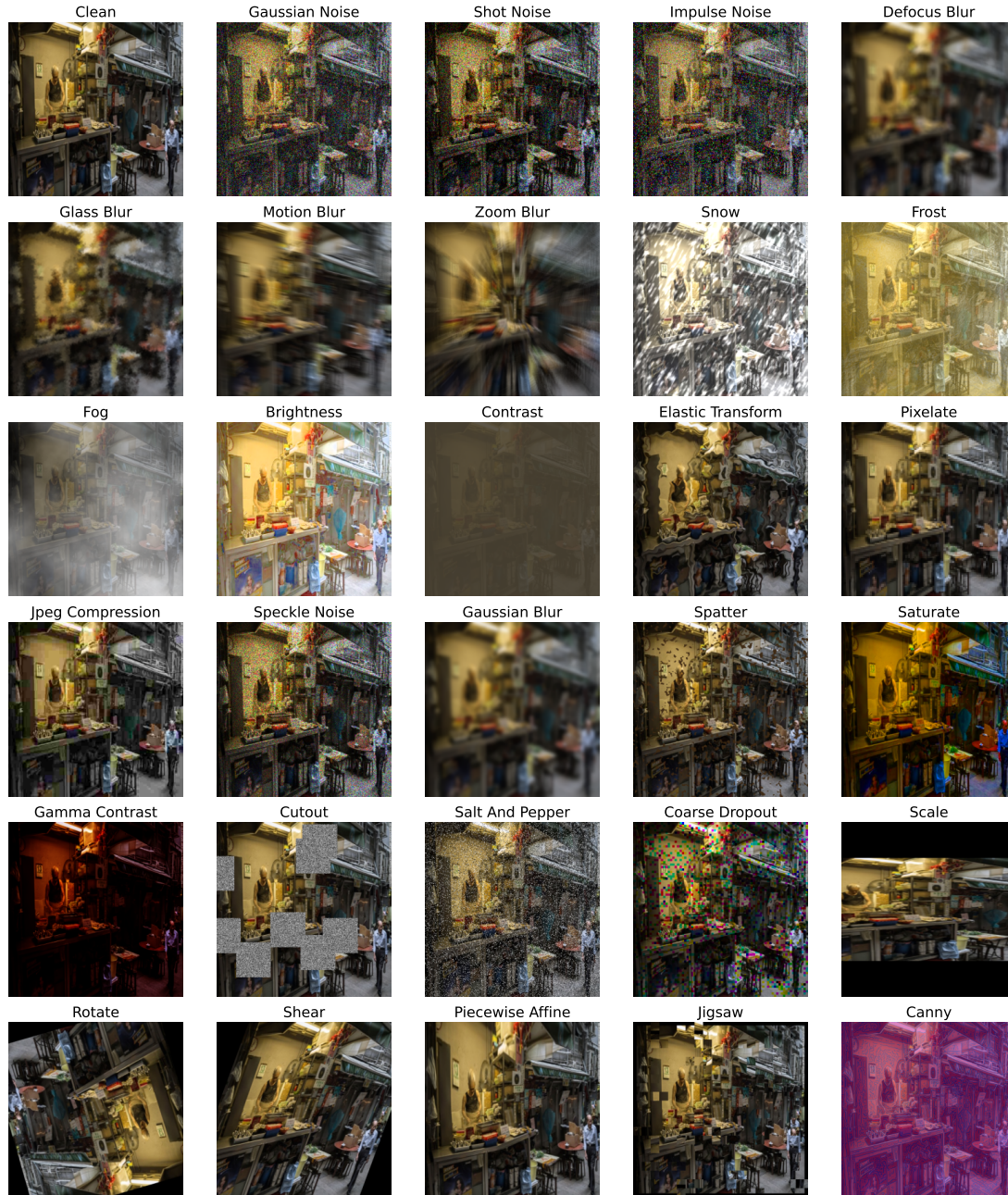| Level | Perturbation Type | Caption |
|---|---|---|
| - | clean | The band, including one man in black jacket and two women standing next to blue corrugated walls. |
| Char | Substitute Char by Ocr | The band, inclodin9 one man in b1ack jacket and two women standing next tu blue corrugated walls. |
| | Substitute Char by Keyboard | The band, including one man in black jacket and two Qomen zranding n3xt to blue cotrugatDd walls. |
| | Insert Char Randomly | The barnd, including one man in black ljackdet and two wom3en standing next to blue corru9gatEed walls. |
| | Substitute Char Randomly | The band, imclu(ing one man in black jycOet and two women stanxinb next to blue corrugated w&lls. |
| | Swap Char Randomly | The band, including one man in balck jacket and two women standing enxt to blue corruagtde walsl. |
| | Delete Char Randomly | The bad, including one man in black jacket and two woem stadig next to blu corrugated walls. |
| | Uppercase Char Randomly | ThE banD, incLudIng onE MaN in BlAck jacket and two wOMEn standing nExt to bluE cOrrUgatEd walls. |
| | Repeat Characters | TThhee bbaanndd,, iinncclluuddiinngg oonnee mmaann iinn bbllaacckk jjaacckkeett ... |
| | Leet Letters | The band, including on3 3an in black jack3t and tw0 3o33n standing next t0 6lu3 corrugated 3alls. |
| | Whitespace Perturbation | The band, including one ma n in black  jack et and two w omen standing next to b lue corrugated walls. |
| | Substitute with Homoglyphs | The band, including one man in black jacket and two women standing next to blue corrugated walls. |
| Word | Synonym Replacement | The band, admit unmatchable homo in pitch dark jacket crown and deuce char brook adjacent to blue sky corrugate walls. |
| | Random Deletion | [The] band, including one man in black jacket [and two] women standing next [to] blue corrugated walls. |
| | Random Swap | The band, man and including in black jacket one two women standing next to blue corrugated walls. |
| | Random Insertion | The band, let in including one succeeding man in black jacket and two women standing next to blue corrugated walls. |
| | Misspell Word | The banda, including onw man In black jacket and two women standing nexto to blue corrugated walls. |
| | Abbreviate Word | The band , including one man in blk jacket and two women standing next tuh blu corrugated walls . |
| | Multilingual Dictionary Based Code Switch | The band, compresa one man in black jaqueta and two women standing next to blue corrugated walls. |
| | Close Homophones Swap | Thee Banned, including one man Inn black Jackett and Thuy women standing Nex to Blew corrugated Wall'S. |
| Sentence | CheckList | @YKIW The band, including one man in black jacket and two women standing next to blue corrugated https://t.co/DdIQ walls. |
| | Back Translation | The band, which includes a man in a black jacket and two women, stands next to blue corrugated iron walls. |
| | Style Paraphraser | The band, with one man in black and two women in blue corrugated walls, stands in the street. |
| | Paraphrase | One man was in a black jacket, and two women were standing next to blue corrugated walls as members of the band. |

18

*Figure 3.* Examples of 29 distinct image corruptions applied to a single image.

*Table 6.* Performance comparison of various **text-to-image** models evaluated by self-consistency scores (uni-modality) across different corruption intensities and data selection levels. Scores represent the average multiplied cosine similarities (max 2300) between the images generated from original captions and the generated images for the captions under different corruption conditions.

| Models | Hard | | Random | | Average |
|---|---|---|---|---|---|
| | Heavy | Light | Heavy | Light | |
| Lafite [126] | 1697 | 1790 | 1679 | 1783 | 1737.25 |
| IF XL v1.0 [26] | 1444 | 1615 | 1533 | 1732 | 1581.00 |
| IF L v1.0 [26] | 1414 | 1581 | 1529 | 1732 | 1564.00 |
| SDXL Turbo [98] | 1412 | 1614 | 1474 | 1697 | 1549.25 |
| IF M v1.0 [26] | 1388 | 1549 | 1522 | 1725 | 1546.00 |
| Kandinsky 3 [4] | 1377 | 1573 | 1478 | 1706 | 1533.50 |
| LCM (Dreamshaper v7) [69, 71] | 1421 | 1560 | 1483 | 1646 | 1527.50 |
| LCM (SSD-1B) [69, 42] | 1388 | 1547 | 1463 | 1649 | 1511.75 |
| Stable Diffusion Turbo [98] | 1384 | 1564 | 1418 | 1630 | 1499.00 |
| SSD 1B [42] | 1368 | 1518 | 1441 | 1621 | 1487.00 |
| LCM LoRA (SSD-1B) [70, 42] | 1363 | 1522 | 1438 | 1621 | 1486.00 |
| LCM LoRA (SDXL) [70, 84] | 1328 | 1503 | 1394 | 1586 | 1452.75 |
| Small Stable Diffusion [52] | 1344 | 1488 | 1393 | 1572 | 1449.25 |
| LCM (SDXL) [69, 84] | 1307 | 1475 | 1398 | 1595 | 1443.75 |
| SDXL Refiner [84] | 1305 | 1482 | 1378 | 1572 | 1434.25 |
| Kandinsky 2.2 [101] | 1290 | 1454 | 1389 | 1598 | 1432.75 |
| SDXL Base [84] | 1307 | 1476 | 1375 | 1568 | 1431.50 |
| Dreamlike Photoreal 2.0 [31] | 1279 | 1445 | 1370 | 1560 | 1413.50 |
| LCM LoRA (SD v1) [70, 97] | 1276 | 1416 | 1395 | 1566 | 1413.25 |
| Dreamshaper v7 [71] | 1272 | 1419 | 1397 | 1562 | 1412.50 |
| Openjourney v4 [86] | 1258 | 1413 | 1340 | 1517 | 1382.00 |
| Stable Diffusion v2 [97] | 1211 | 1380 | 1316 | 1516 | 1355.75 |
| Anything Midjourney v4.1 [86] | 1225 | 1385 | 1313 | 1492 | 1353.75 |
| Glide [78] | 1278 | 1342 | 1326 | 1420 | 1341.50 |
| Stable Diffusion v1 [97] | 1217 | 1364 | 1300 | 1479 | 1340.00 |
| Unidiffuser [11] | 1168 | 1298 | 1293 | 1461 | 1305.00 |
| Karlo-v1.0.alpha [58] | 1123 | 1333 | 1220 | 1476 | 1288.00 |

*Table 7.* Performance comparison of various **image-to-text** models evaluated by self-consistency scores (uni-modality) across different corruption intensities and data selection levels. Scores represent the average multiplied cosine similarities (max 2900) between the captions generated for the original images and the captions generated for the images under different corruption conditions.

| Models | LLM | Hard | | Random | | Average |
|---|---|---|---|---|---|---|
| | | Heavy | Light | Heavy | Light | |
| 🏆 InstructBLIP [25] | Vicuna 13B [125] | 2017 | 2254 | 2238 | 2389 | 2224.50 |
| 🥈LLaVA-v1.5 [64] | Vicuna 13B [125] | 1956 | 2192 | 2275 | 2404 | 2206.75 |
| 🥉InstructBLIP [25] | Vicuna 7B [125] | 1980 | 2223 | 2229 | 2386 | 2204.50 |
| InstructBLIP [25] | Flan T5 XL [20] | 1945 | 2241 | 2200 | 2344 | 2182.50 |
| LLaVA-v1.5 [64] | Vicuna 7B [125] | 1918 | 2160 | 2249 | 2388 | 2178.75 |
| BLIP2 [61] | Flan T5 XXL [20] | 1819 | 2123 | 2215 | 2399 | 2139.00 |
| LLaVA [65] | Vicuna 13B [125] | 1885 | 2122 | 2203 | 2337 | 2136.75 |
| BLIP2 [61] | Flan T5 XL [20] | 1777 | 2131 | 2201 | 2427 | 2134.00 |
| InstructBLIP [25] | Flan T5 XXL [20] | 1966 | 2172 | 2125 | 2251 | 2128.50 |
| LLaVA [64] | LLaMA2 13B [108] | 1838 | 2080 | 2182 | 2331 | 2107.75 |
| BLIP2 [61] | OPT-6.7b [122] | 1776 | 2088 | 2141 | 2332 | 2084.25 |
| mPLUG-Owl [119] | LLaMA 7B [107] | 1842 | 2107 | 2107 | 2274 | 2082.50 |
| LLaVA [64] | MPT 7B [105] | 1827 | 2056 | 2141 | 2288 | 2078.00 |
| LLaMA-Adapter v2 [37] | LLaMA 7B [107] | 1840 | 2059 | 2133 | 2264 | 2074.00 |
| mPLUG-Owl2 [120] | LLaMA2 7B [108] | 1748 | 2054 | 2104 | 2324 | 2057.50 |
| BLIP2 [61] | OPT-2.7b [122] | 1714 | 2030 | 2084 | 2284 | 2028.00 |
| LLaVA [64] | LLaMA2 7B [108] | 1769 | 1997 | 2102 | 2243 | 2027.75 |
| LaVIN [67, 68] | LLaMA 13B [107] | 1759 | 2036 | 2047 | 2249 | 2022.75 |
| mPLUG-Owl (multilingual) [119] | LLaMA 7B [107] | 1765 | 2008 | 2071 | 2229 | 2018.25 |
| ShareGPT4V [18] | Vicuna 7B [125] | 1699 | 1956 | 2105 | 2277 | 2009.25 |
| CogVLM [115, 46] | Vicuna 7B [125] | 1819 | 2093 | 1937 | 2119 | 1992.00 |
| MiniGPT-4 [127] | Vicuna 13B [125] | 1765 | 1938 | 2052 | 2150 | 1976.25 |
| MiniGPT-4 [127] | Vicuna 7B [125] | 1759 | 1936 | 1999 | 2089 | 1945.75 |
| MiniGPT-4 [127] | LLaMA2 7B [108] | 1722 | 1875 | 2015 | 2084 | 1924.00 |
| ImageBind-LLM [43] | Open Chinese LLaMA 7B [81, 107] | 1585 | 1833 | 2005 | 2165 | 1897.00 |
| BLIP Large [59] | - | 1487 | 1805 | 1983 | 2238 | 1878.25 |
| VisualGLM [32, 28] | ChatGLM-6B [32, 121] | 1598 | 1811 | 1849 | 1953 | 1802.75 |
| Unidiffuser [11] | - | 1484 | 1669 | 1894 | 2085 | 1783.00 |
| ViT-GPT2 [79] | GPT2 [88] | 1262 | 1468 | 1918 | 2169 | 1704.25 |
| GIT Large [114] | - | 1336 | 1556 | 1878 | 2046 | 1704.00 |
| BLIP Base [59] | - | 1192 | 1410 | 1887 | 2153 | 1660.50 |
| Qwen-VL-Chat [9] | Qwen-7B [8] | 1317 | 1611 | 1696 | 1901 | 1631.25 |
| OpenFlamingo [6, 1] | MPT 7B [105] | 1506 | 1681 | 1530 | 1708 | 1606.25 |
| Multimodal-GPT [39] | LLaMA 7B [107] | 1400 | 1591 | 1639 | 1771 | 1600.25 |
| GIT Base [114] | - | 1232 | 1394 | 1769 | 1965 | 1590.00 |
| OpenFlamingo [6, 1] | RedPajama 3B [106] | 1367 | 1537 | 1488 | 1660 | 1513.00 |
| OpenFlamingo [6, 1] | MPT 1B [76] | 1219 | 1404 | 1312 | 1519 | 1363.50 |
| MiniGPT v2 [17] | LLaMA2 7B [108] | 1131 | 1204 | 1318 | 1381 | 1258.50 |