

Benchmarking Robustness of Adaptation Methods on Pre-trained Vision-Language Models

Shuo Chen^{1,3*} Jindong Gu^{2*} Zhen Han¹ Yunpu Ma^{1,3} Philip Torr² Volker Tresp¹

¹LMU Munich ²University of Oxford ³Siemens AG

shuo.chen@campus.lmu.de, jindong.gu@eng.ox.ac.uk

Abstract

Various adaptation methods, e.g., LoRA and adapters, have been proposed to enhance the domain-specific performance of pre-trained vision-language models. The robustness of these adaptation methods against distribution shifts is underexplored. In this study, we assess the robustness of 11 widely-used adaptation methods across 4 vision-language datasets under multimodal corruptions. We introduce 7 benchmark datasets, including 96 visual and 87 textual corruptions, to investigate the robustness of different adaptation methods, the impact of available adaptation examples, and the influence of trainable parameter size during adaptation. Our analysis reveals that: 1) Adaptation methods are more sensitive to text corruptions than visual corruptions. 2) Full fine-tuning does not consistently provide the highest robustness; instead, adapters can achieve better robustness with comparable clean performance. 3) Contrary to expectations, increasing the number of adaptation data and parameters does not guarantee enhanced robustness; instead, it results in even lower robustness. The code and datasets used in this study are public¹.

1. Introduction

Large-scale pre-training of vision-language (VL) models has become the de facto framework for VL tasks [15, 13]. These models are typically trained in a self-supervised manner on unlabeled web-scale datasets in a general domain [15]. To improve performance on domain-specific downstream tasks, various model adaptation methods have been proposed [5, 10, 8, 9, 6].

Although adaptation methods can achieve promising results on various VL benchmark datasets, real-world applications introduce various distribution shifts [11], such as lighting conditions and text typos. Therefore, it is critical to ensure model robustness against distribution shifts, par-

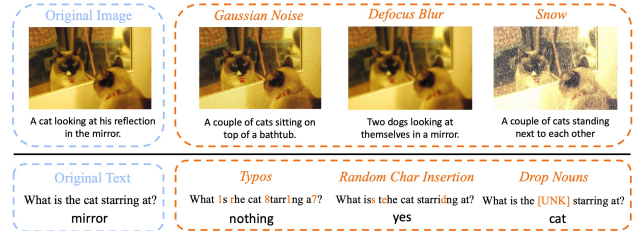


Figure 1: Multimodal adaptation methods are sensitive to image and text corruptions. The two rows show image captioning and visual question answering from Adapter [5].

ticularly in safety-critical applications. However, related research is still rare in multimodal models. To this end, this work investigates the robustness of various adaptation methods on VL models. We introduce a diverse set of 96 visual corruptions, including *impulse noise*, *snow* etc., and 87 textual corruptions encompassing *text addition*, *back translation*, etc. Moreover, extensive experiments have been conducted on 11 adaptation methods across 4 VL datasets.

Our analysis reveals several interesting findings: 1) Adaptation methods demonstrate a higher degree of sensitivity towards text corruptions compared to visual corruptions. 2) Full fine-tuning does not consistently yield the best relative robustness, whereas an adapter can achieve better robustness with comparable performance. 3) Surprisingly, the quantity of adaptation data and model parameters does not guarantee improved robustness. In fact, augmenting the amount of adaptation data might even lead to decreased robustness. To summarize, our contributions are as follows:

- We construct a suite of 7 large-scale robustness benchmark datasets including 96 different visual corruptions and 87 textual corruption methods.
- We evaluate the robustness of 11 adaptation methods on VL models with massive experiments.
- We release the benchmark, code as well as a leaderboard to the community.

*equal contribution

¹<https://adarobustness.github.io>

2. Benchmark and Evaluations

2.1. Corruption Methods

Image Corruptions. We use the corruption methods from ImageNet-C [4] and [1, 11]. A *blank* method is also added, which sets all pixel values to 255 and is used to examine the importance of visual information to VL models. All image corruptions can be categorized into five groups: **noise**, **blur**, **weather**, **digital**, and **extra**. In total, we have 96 types of visual corruption and we leave the details in the supplementary. By applying all image corruptions to 4 datasets used in this study, i.e., VQAv2 [3], GQA [7], NLVR² [12] and MSCOCO Caption [2], we construct 4 out-of-distribution (OOD) benchmark datasets.

Text Corruptions. We have incorporated a total of 35 text corruption methods, inspired by the approaches presented in [11, 1]. These methods can be categorized into three groups based on the level of corruption: **character-level**, **word-level**, and **sentence-level**. Furthermore, they are further subdivided into six sub-categories, namely *character modification*, *text style modification*, *text addition*, *dropping text based on POS*, *positional drop*, and *text swap*. Various severity levels for text corruptions are also introduced. In total, we have 87 different perturbations. After applying all text corruptions on VQAv2 [3], GQA [7], and NLVR² [12], we construct another 3 OOD benchmark datasets. Please refer to the supplementary for detailed information.

2.2. Robustness Evaluation Protocol

The performance decrease of a model F on the OOD test data D_O compared to the results on the in-distribution dataset D_I is measured to evaluate the robustness against distribution shifts. Concretely, the model performance P_I on D_I and P_O on D_O are first evaluated, where P is the corresponding evaluation metric for each task. Then, the **Relative Robustness** $RR = 1 - \Delta P / P_I$ [11, 1] is computed, where $\Delta P = (P_I - P_O)$. RR is a score ranging from 0 to 1, where $RR = 1$ indicates that F is totally robust and $R = 0$ means that F is not robust at all.

2.3. Experimental Settings

Tasks and Datasets. For VQA task, VQAv2 [3] and GQA [7] are adopted. We also incorporate NLVR² [12] for visual reasoning and MSCOCO [2] for image captioning.

Models. CLIP-BART (T5) [13] is as our base model. Because the model adaptation on VL models is mainly on language model component and the encoder-decoder architecture can tackle VL tasks via a unified text-generation task.

Model Adaptation Methods. The robustness of four mainstream adaptation methods is investigated, i.e., full fine-tuning, soft prompt [9], LoRA [6], and adapter-based methods, including Adapter [5], Hyperformer [10], and Compacter [8]. Shared adaptation methods are also investigated

Table 1: Clean performance and relative robustness (RR) of adaptation methods based on CLIP-BART against image (up) and text (down) corruptions. RR and the corresponding standard deviation are averaged and calculated over all image or text corruption methods. We strike out those high RR with quite low performance. The best RR is in bold.

Adaptation method	VQAv2		GQA		NLVR ²		MSCOCO Caption	
	Acc (%)	RR (%)	Acc (%)	RR (%)	Acc (%)	RR (%)	CIDEr	RR (%)
Full Fine-tuning	66.75	84.86 \pm 3.17	55.04	89.20 \pm 2.04	73.01	90.34 \pm 2.04	115.03	68.40 \pm 2.04
Multiple Adapters	65.30	85.33 \pm 4.06	53.39	86.16 \pm 2.04	69.41	92.02 \pm 2.04	114.47	68.72 \pm 2.04
Half-shared Adapters	65.20	85.18 \pm 5.01	52.96	89.37 \pm 2.04	70.03	91.72 \pm 2.04	114.50	68.45 \pm 2.04
Single Adapter	65.35	85.76 \pm 3.32	54.14	82.49 \pm 2.04	73.89	90.04 \pm 2.04	115.04	68.68 \pm 2.04
Hyperformer	65.38	85.38 \pm 4.44	52.52	90.05 \pm 2.04	72.21	90.13 \pm 2.04	114.89	68.74 \pm 2.04
Multiple Compacters	64.91	85.65 \pm 4.81	52.75	88.89 \pm 2.04	69.45	91.33 \pm 2.04	115.16	68.67 \pm 2.03
Single Compacter	64.47	85.47 \pm 4.06	52.90	82.62 \pm 2.04	69.94	92.04 \pm 2.04	113.06	69.92 \pm 2.03
Multiple LoRA	65.44	84.78 \pm 4.86	52.05	91.15 \pm 2.04	51.32	—	115.41	68.47 \pm 2.04
Single LoRA	65.34	84.78 \pm 4.81	53.19	82.58 \pm 2.04	73.58	90.05 \pm 2.04	114.54	69.26 \pm 2.03
Multiple Prompts	46.81	—	34.01	—	49.87	—	108.62	67.70 \pm 2.04
Single Prompt	44.00	—	37.54	—	51.95	—	103.70	68.56 \pm 2.03

Adaptation method	VQAv2		GQA		NLVR ²	
	Acc (%)	RR (%)	Acc (%)	RR (%)	Acc (%)	RR (%)
Full Fine-tuning	66.75	73.65 \pm 22.38	55.04	66.92 \pm 24.14	73.01	87.06 \pm 11.00
Multiple Adapters	65.30	76.62 \pm 20.66	53.39	66.93 \pm 22.43	69.41	90.14 \pm 10.19
Half-shared Adapters	65.20	76.78 \pm 20.79	52.96	68.20 \pm 24.78	70.03	89.16 \pm 10.12
Single Adapter	65.35	77.64 \pm 21.09	54.14	67.47 \pm 20.03	73.89	88.49 \pm 10.87
Hyperformer	65.38	75.06 \pm 21.29	52.52	70.30 \pm 23.13	72.21	87.27 \pm 11.27
Multiple Compacters	64.91	77.10 \pm 20.85	52.75	67.39 \pm 23.29	69.45	90.00 \pm 9.76
Single Compacter	64.47	77.17 \pm 20.40	52.90	67.90 \pm 20.33	69.94	90.10 \pm 9.81
Multiple LoRA	65.44	74.04 \pm 21.97	52.05	68.77 \pm 22.76	51.32	—
Single LoRA	65.34	74.50 \pm 21.42	53.19	63.94 \pm 20.99	73.58	87.64 \pm 11.04
Multiple Prompts	46.81	—	34.01	—	49.87	—
Single Prompt	44.00	—	37.54	—	51.95	—

(See supplementary materials). In total, there are 11 adaptation methods studied in this work.

Evaluation Metrics. Accuracy is used in VQAv2, GQA and NLVR² and CIDEr [14] is adopted in image captioning. The main paper reports the RR defined in Sec. 2.2 with severity 5; more detailed scores are in the supplementary.

3. Results and Analysis

3.1. Robustness of Multimodal Adaptation Methods

The relative robustness of adaptation methods against image and text corruptions are presented in Tab. 1. The reported relative robustness is the average value across all images or text corruption methods. **Although full fine-tuning generally achieves higher clean performance, our analysis reveals that its robustness is comparatively weaker than other adaptation methods.** In many cases, adapter and hyperformer achieve better robustness with much fewer parameters and comparable clean performance. For instance, full fine-tuning’s RR against text corruptions on the VQAv2 dataset is the smallest. Prompt tuning, despite exhibiting high robustness, fails to perform well on the clean test dataset. Please note that we have excluded robustness scores associated with very low task performance in Tab. 1. **Single Adapter vs Full Fine-tuning.** Previous research [13] has shown that a single adapter can achieve comparable performance on the four tasks with significantly fewer parameters. When it comes to robustness, *a single adapter is comparable to or slightly better than full fine-tuning on VQAv2, NLVR², and MSCOCO Caption given image cor-*

ruptions. The same goes for text corruption. For example, as shown in the 4th row and 1st row in Tab. 1 (lower panel), a single adapter on CLIP-BART achieves an average RR of 77.64% against text corruptions on VQAv2, while full fine-tuning’s RR is 73.65%. However, on GQA, a single adapter is less robust than full fine-tuning. Full fine-tuning achieves an average RR of 89.20% against image corruptions, while the RR of a single adapter is only 82.49%. In contrast, multiple and half-shared adapters have more parameters but achieve better robustness on the four tasks than a single adapter. **In conclusion, a single adapter can achieve similar or better robustness on VQAv2, NLVR², and MSCOCO Caption compared to full fine-tuning. On GQA, multiple and half-shared adapters are better.**

Adapter-based Methods. Although training multiple tasks with one set of adapter layers has the least parameters, *such single setting might hinder the robustness on certain tasks*. For instance, Single Adapter’s robustness on GQA against image corruptions (82.49%) is lower than that of the half-shared (89.37%) and multiple settings (86.16%). An explanation could be that the half-shared mechanism not only learns more general representation across tasks; it also maintains task-specific knowledge. On GQA, Single LoRA’s robustness against image corruptions is lower by 8.57% compared to Multiple LoRA’s. However, compared with multiple settings of LoRA and Adapter, the *Hyperformer has relatively fewer parameters but achieves comparable or better robustness*.

Vision-language Tasks. Among all datasets, MSCOCO Caption is the most vulnerable one against image corruptions. This is plausible as it only relies on visual information, whereas other tasks provide both visual and language information. Besides, GQA is the task with the lowest robustness performance against text corruptions. Moreover, *on GQA, the extreme single-module setting fails to achieve good robustness, such as Single Adapter and Single LoRA. This indicates that information sharing with other two datasets may hinder the robustness on GQA.*

3.2. Robustness Sensitivity to Image Corruptions and Text Corruptions

Our experimental findings suggest a potential vulnerability of adaptation methods on multimodal VL models to text corruptions, particularly those at the character level. Across all three tasks, the adaptation methods exhibit lower robustness indicators against text corruptions. For instance, Single Adapter based on CLIP-BART has the best robustness result 85.76% against image corruptions on VQAv2. However, although it is still the most robust adaptation method against text corruptions, the relative robustness is 77.64%. *Among image corruptions, zoom blur drops the robustness the most, and within text corruptions, char-level methods are the most challenging to these VL adaptation methods.*

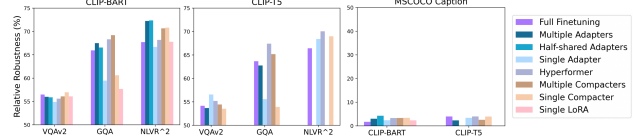


Figure 2: Relative robustness (%) of adaptation methods based on CLIP-BART (left) and CLIP-T5 (middle) against *blank* corruption. We group MSCOCO Caption results from CLIP-BART and CLIP-T5 together in the right sub-figure. We omit two bars in NLVR² from the middle figure as multiple adapters and multiple compacters did not perform well.

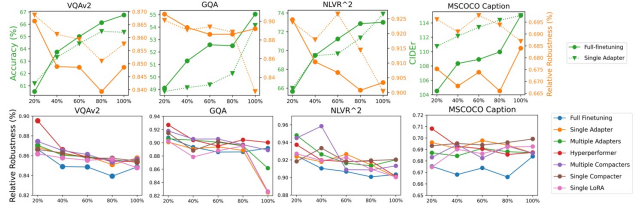


Figure 3: The first row represents the clean performance and RR of full fine-tuning and single adapter on CLIP-BART given different size of adaptation dataset. Green lines stand for performance in each task and the orange is robustness. The second row is RR given different size of adaptation dataset. X-axis shows the random subset ratio of training dataset during adaptation, ranging from 20% to 100%.

Blank Image Corruption. The results are shown in Fig. 2. The relative robustness on MSCOCO Caption is the lowest among all four datasets. This is plausible since image captioning relies only on visual information and is not supposed to perform well given a blank image. Apart from image captioning, adaptation methods on all three datasets could secure a relative robustness exceeding 50% without useful visual information. Several questions within the VL tasks can be accurately answered without relying on visual information, suggesting that **language information plays a more significant role than visual information**. This also explains the higher sensitivity to text corruptions compared to the sensitivity to image corruptions.

3.3. The Influence of Adaptation Data Size and Parameter Size on Robustness

Adaptation Data Size. Given more adaptation data, performance in all tasks has a steady increase (Fig. 3), and in most cases, the performance of full fine-tuning surpasses single adapter’s performance. Single adapter achieves better relative robustness compared to full fine-tuning against both image and text corruptions but has a robustness drop on GQA. We also investigate the robustness of other adaptation methods (Fig. 3). All lines present a steady declining tendency, which indicates that **increasing the size of the**

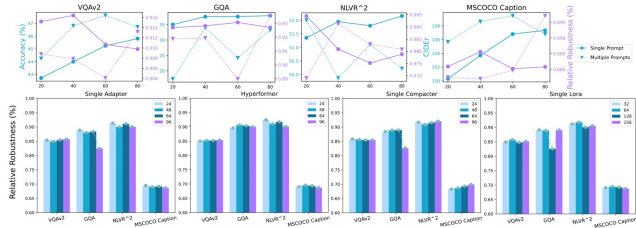


Figure 4: The top row shows the clean performance and RR with different prompt lengths on CLIP-BART. Blue and purple lines are performance and relative robustness respectively. The bottom row shows the RR given a different number of parameters in 4 adaptation methods. Different colors stand for different embedding sizes and larger numbers are with more parameters.

adaptation data does not consistently enhance relative robustness. In comparison to other adaptation methods, *full fine-tuning has relatively lower relative robustness.*

Adaptation Parameter Size. We inspect the robustness as well as the clean performance of prompt-tuning given different soft prompt lengths. The results conducted on CLIP-BART are shown in Fig. 4. We can see a steady increase in the performance on four tasks with longer prompt lengths which proves that prompt methods perform better given more parameters. Regarding relative robustness, such a steady increase does not apply to all tasks, and **longer soft prompts do not ensure better relative robustness.** We also conduct experiments on the other 4 adaptation methods given different sizes of trainable parameters. The results demonstrate that **more parameters do not ensure enhanced robustness and some even reduce it.**

4. Discussion and Conclusion

This study focuses on the robustness of adaptation methods on pre-trained VL models and provides 7 benchmark datasets containing 96 visual and 87 textual corruptions. We systematically inspect the robustness of 11 adaptation methods on 4 popular VL datasets. Potential future work includes investigation on more diverse pre-trained VL models, design of more robust adaptation methods, and integrating future model adaptation methods, etc.

References

- [1] Madeline Chantry, Shruti Vyas, Hamid Palangi, Yogesh Rawat, and Vibhav Vineet. Robustness analysis of video-language models against visual and language perturbations. *Advances in Neural Information Processing Systems*, 35:34405–34420, 2022.
- [2] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [3] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [4] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- [5] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.
- [6] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [7] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [8] Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–1035, 2021.
- [9] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- [10] Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*, 2021.
- [11] Jielin Qiu, Yi Zhu, Xingjian Shi, Florian Wenzel, Zhiqiang Tang, Ding Zhao, Bo Li, and Mu Li. Are multimodal models robust to image and text perturbations? *arXiv preprint arXiv:2212.08044*, 2022.
- [12] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Hua-jun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2018.
- [13] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. Vi-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5227–5237, 2022.
- [14] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015.
- [15] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *arXiv preprint arXiv:2202.03052*, 2022.

Benchmarking Robustness of Adaptation Methods on Pre-trained Vision-Language Models (Supplementary)

Shuo Chen^{1,3*} Jindong Gu^{2*} Zhen Han¹ Yunpu Ma^{1,3} Philip Torr² Volker Tresp¹

¹LMU Munich ²University of Oxford ³Siemens AG

shuo.chen@campus.lmu.de, jindong.gu@eng.ox.ac.uk

1. Related Work

Vision-language Models. Pre-trained VL models [36, 52, 6, 39, 35, 56, 12, 67] have shown outstanding performance on various VL tasks. Some use contrastive learning to align visual features with language representations and achieve surprising zero-shot performance [37, 47]. However, contrastive learning-based methods are limited to close-ended tasks and are not flexible enough. Another line of work follows BERT’s [10] pretrain-then-finetune paradigm [36, 52, 6, 56]. They treat object features extracted using pre-trained object detectors [17] as visual words sent to language models [10]. For example, VL-BART [6] uses BART [34] or T5 [48] as the text encoder and Faster-RCNN [17] as the visual backbone. Unlike other methods, VL-BART unifies VL tasks via a single text generation task. CLIP-BART [54] follows the same idea as VL-BART but adopts the CLIP [47] image encoder to extract pixel-level features. Recent approaches [58, 13, 1] follow such a unified view and freeze large language models (LLMs) to utilize the in-context learning ability of LLMs. However, as shown in [54], LM fine-tuning is still crucial to competitive performance on various downstream VL tasks. In this study, we follow the work [54] that benchmarks model adaptation methods. CLIP-BART [54] is selected as our VL model, given its generation flexibility and unified architecture.

Model Adaptation Methods. To enhance the performance of pre-trained VL models on downstream tasks and avoid infeasible computation, various adaptation methods have been proposed. Existing methods can be classified into three categories [54]: (1) adding a few trainable parameters while freezing other model parts [28, 40, 38, 32, 33]; (2) updating a few model parameters sparsely [24, 55, 65]; and (3) low-rank factorization of parameters to be updated, such as LoRA [29]. Adapters [28] belong to the first category and have been widely used in vision, language, and multimodal models [28, 66, 5]. Other representative methods in the

first category include Hyperformers [40], Compacters [32], and prompt-tuning [33, 2]. Although numerous adaptation methods have been proposed and widely adopted, their robustness against distribution shifts remains understudied.

Natural Robustness. The robustness of deep learning models against distribution shifts is critical for real-world applications [16, 25, 15]. Regarding vision robustness, researchers have investigated image classification models [25, 15, 27, 49, 20, 21], semantic segmentation [31, 23], object detection [41], video classification [64], and transformer-based architectures [11, 44, 45, 22, 62]. In the field of natural language processing (NLP), many robustness analysis toolboxes [50, 51, 61, 18], and various methods [14, 9, 43, 8] are available. The robustness investigation on multimodal models is gaining more attention but related studies are lacking. The literature includes the robustness of multimodal fusion models [63], audio-visual models [57], text-to-image generative models [7], text-to-video retrieval models [3] as well as image-text models [46].

In contrast to all the work above, our study focuses on *the robustness of adaptation methods integrated into large pre-trained vision-language models*. Understanding their robustness on different VL tasks will facilitate the design of more robust adaptation methods for multimodal models.

2. Corruption Details

2.1. Image Corruption Methods

We have incorporated 20 image corruption methods (Fig. 1) following [25, 3, 46] which can be categorized into 5 categories **noise**, **blur**, **weather**, **digital**, and **extra**. We also introduce 5 levels of severity (Fig. 2) to image corruption methods except for *blank* corruption.

Noise There are 4 different noises used in this study, namely *Impulse*, *Gaussian*, *Shot*, and *Speckle*. *Impulse noise* mimics the corruption caused by bit errors by introducing a mixture of salt and pepper noise, with varying in-

*equal contribution

tensities such as 0.03, 0.06, 0.09, 0.17, and 0.27. *Gaussian noise* simulates corruptions due to low-lighting conditions by pixel value normalization and normal noise addition. The intensity of this noise is scaled according to severity, with values of 0.08, 0.12, 0.18, 0.26, and 0.38. *Shot noise* simulates electronic noise caused by discrete light and is also called Poisson noise. *Speckle noise* is an additive noise where larger noise will be added if the original pixel value is larger.

Blur *Zoom blur* is observed when the camera swiftly moves toward an object. It blurs toward the center of the frame. *Defocus blur* occurs when the image is out of focus and its severity is defined by the disk radius convolved over the image ranging from (3, 0.1), (4, 0.5), (6, 0.5), (8, 0.5), (10, 0.5). *Motion blur* occurs when the camera is moving quickly. The blurring effect is created by a kernel with different radius and sigma ranging from (10, 3), (15, 5), (15, 8), (15, 12), (20, 15). *Frosted Glass blur* occurs when the glass of windows or panels frosts. *Gaussian blur* generates blurred pixels by a weighted average of its neighbors. The farther the neighbors are, the lower the weight in the average is.

Digital *JPEG* is a lossy image compression format that converts the original picture to JPEG format with quality ranging from 25, 18, 15, 10, 7 given different severity. *Contrast* simulates corruptions caused by the lighting conditions and object’s color. *Elastic* stretches small image regions for stretching effects. *Spatter* appears when the lens is occluded by rain or mud. *Pixelate* transforms the original images into a small number of large pixels.

Weather *Snow* mimics the visual obstruction caused by precipitation. *Frost* occurs when lenses are covered by ice crystals. *Fog* obscures the object and is rendered by the diamond-square algorithm. *Brightness* adds a bright effect to the image simulating daylight.

Extra we introduce a new corruption method called *blank* that sets all pixel values to 255, i.e. turning the original image into a blank picture.

2.2. Text Corruption Methods

We have included 35 text corruptions methods following [61, 46, 3] which corrupt on 3 levels (e.g. character, word, and sentence level) and can be categorized into 3 main categories, 6 sub-categories.

Character Modification. *Character modification* simulates common mistakes during typing and corrupts the text on a character level and contains 9 methods, namely *OCR*,

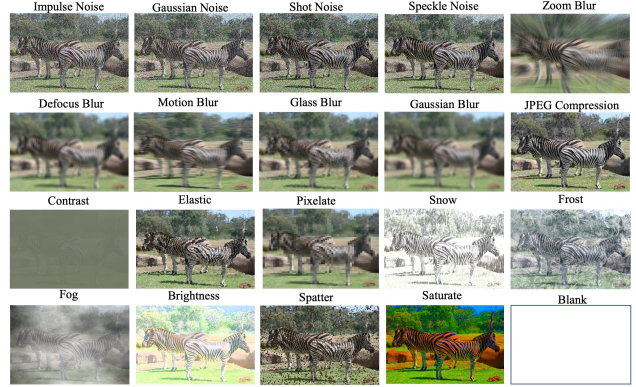


Figure 1: We introduce 20 corruption methods to image data following the methods in [25]. Except for the blank corruption, each type of corruption has five levels of severity. In total, there are 96 different corruptions, as shown in Appendix 1

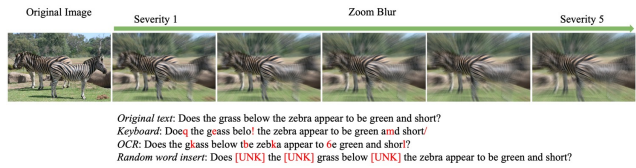


Figure 2: Examples of image and text corruptions. The top row shows an original image from GQA and images corrupted by *zoom blur* with 5 levels of severity. The second row presents text corruptions on the original texts where the red sign indicates the corrupted parts. More examples are shown in supplementary material.

Punct, *Typos*, *Keyboard*, *Spelling Error*, *char random insert*, *char random replace*, *char random swap*, *char random delete*. *OCR* replaces a character based on common Optical Character Recognition (OCR) errors. *Keyboard* substitutes original characters based on keyboard distance.

Text Style Transformation. Methods in the *text style* category modify the text on a sentence level and transform the sentence style to a target one, such as turning the original sentence to *passive*, *formal*, *causal*, to *double negative*, and changing the *tense*.

Text Addition. *Text addition* inserts additional words to the original text. *InsertAdv* adds an adverb before each verb. *AppendIrr* adds irrelevant phrases to the original texts. *Random Insert* randomly inserts token [UNK] to the original texts.

Text Drop based on POS. *Text drop* transforms the text on a word level based on POS tagging [61]. There are meth-

Table 1: Image Corruption Methods

Category	Corruption Method	Severity
Noise	Impulse	5
	Gaussian	5
	Shot	5
	Speckle	5
Blur	Zoom	5
	Defocus	5
	Motion	5
	Frosted Glass	5
	Gaussian Blur	5
Digital	JPEG	5
	Contrast	5
	Elastic	5
	Saturate	5
	Spatter	5
	Pixelate	5
Weather	Snow	5
	Frost	5
	Fog	5
	Brightness	5
Extra	Blank	1
5 Category	20 Methods	96 Severity

ods of dropping nouns (*DropNN*), verbs (*DropVB*), dropping both nouns and verbs (*Drop VB & NN*), dropping random nouns and verbs (*Drop Rand NN*, *Drop Rand VB*), keeping only nouns (*Only NN*), keeping only verbs (*Only VB*), keeping both (*Only NN & VB*).

Text Drop based on Position. According to the word position, there are methods of dropping the first (*Drop First*), dropping the last (*Drop Last*), dropping both the first and the last (*Drop First and Last*), shuffling word order (*Shuffle Order*), and randomly removing words (*Random Delete*).

Text Swap. Under the category of *text swap*, we can replace words randomly or with synonyms by word embedding (*SwapSyn Word Embd*) and WordNet [42] (*SwapSyn WordNet*). We also deploy back translation (*Back Trans*) from [61], which first translates English into French and then translates it back to English. *Random Swap* randomly swaps words position.

In addition to the above corruptions, Qiu et al. proposed in [46] that we should ensure the corrupted text has the same semantics as the original one to make sure the image-text pairs remain meaningful. We follow this setting and use the same fidelity guarantee mechanism as [46].

Table 2: Text Corruption Methods

Category	Sub Category	Level	Corruption Method	Severity	
Natural	Char Modification	character	OCR	5	
		character	Punct	1	
		character	Typos	5	
		character	Keyboard	5	
		character	Spelling Error	5	
		character	char random insert	5	
		character	char random replace	5	
		character	char random swap	5	
		character	char random delete	5	
	Text Style	sentence	Passive	1	
		sentence	Tense	1	
		sentence	Formal	1	
		sentence	Casual	1	
		sentence	Active	1	
		sentence	Double Neg	1	
	Text Addition	word	InsertAdv	1	
		word	Appendlrr	1	
		word	Random Insert	5	
	Synthetic	Drop Text based on POS tag	word	Drop NN	1
word			Drop Rand NN	1	
word			Drop VB	1	
word			Drop VB & NN	1	
word			Only NN	1	
word			Only VB	1	
word			Only NN & VB	1	
word			Drop Rand VB	1	
Positional Drop		word	Drop First	1	
		word	Drop Last	1	
		word	Drop First and Last	1	
		sent	Shuffle Order	1	
		word	Random Delete	5	
Machine		Text Swap	word	SwapSyn Word Embd	5
			word	SwapSyn WordNet	5
			sentence	Back Trans	1
			word	Random Swap	5
35 methods			87 levels of severity		

3. Model Implementations

4. Preliminary of Model Adaptation Methods

The pretrain-then-finetune paradigm on large models has shown dominant performance on multimodal tasks [36, 10, 47], yet the prohibitive costs of full fine-tuning have spurred intensive research efforts towards developing parameter-efficient adaptation methods [54, 29, 33, 28, 40, 32]. As the transformer architecture [59] is used for most state-of-the-art large pre-trained models, adaptation methods mainly focus on tweaking the input or the intermediate layers of the attention layers inside the large models. Formally, given a pre-trained large-scale model F parameterized by θ , we need to adapt F on a task-specific dataset \mathcal{D} , e.g., a VQA dataset. Then, we can obtain the output $\mathbf{y} = F'(\mathbf{x}; \theta)$ by providing an input $\mathbf{x} = \{x_1, \dots, x_n\}$ with n tokens from \mathcal{D} . Adaptation methods differ in how they interact with $F(\mathbf{x}; \theta)$ (Fig. 3). In general, full fine-tuning updates all θ . Prompt [33] concatenates the input \mathbf{x} with an extra prefix. LoRA [29] introduces modifications to the update mechanism of the model parameters θ and adapters [28, 40, 32] modify the intermediate output and input of θ .

Full fine-tuning directly updates the whole θ on \mathcal{D} and becomes prohibitive due to the rapidly growing model size. Therefore, the following adaptation methods are developed to achieve comparable performance while optimizing only

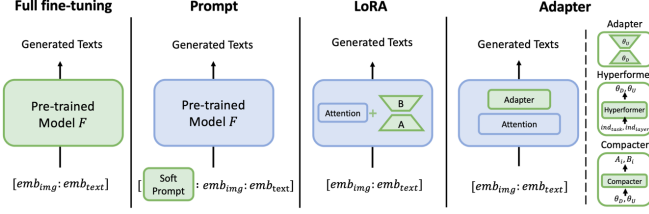


Figure 3: Illustration of adaptation methods used in our study. Green areas indicate trainable parameters whereas frozen parameters are in blue. The input is the concatenation of image and text embedding $[emb_{img} : emb_{text}]$ and the output is the generated texts.

a few parameters.

Prompt-based adaptation concatenates the input \mathbf{x} with either a trainable prefix (soft prompt) [33] or a manually designed prefix [2]. For the given input $\mathbf{x} = \{x_1, \dots, x_n\}$ with n tokens, the pre-trained model will first form an embedding matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ where d is the dimension of the embedding space. Soft-prompts [33] are then represented as a learnable parameter $\mathbf{P} \in \mathbb{R}^{p \times d}$, where p is the length of the prompt. Next, \mathbf{P} is concatenated with the original embedded input \mathbf{X} to form a new single matrix defined as $[\mathbf{P} : \mathbf{X}] \in \mathbb{R}^{(p+n) \times d}$. During adaptation, the model is trained to maximize the probability of the desired output while only updating \mathbf{P} .

LoRA [29] utilizes low-rank decomposition matrices to update parameters. For intermediate model parameters $\theta_0 \in \mathbb{R}^{d \times k}$, such as the parameters from a self-attention module in the transformer architecture, its update $\Delta\theta_0$ is represented by a low-rank decomposition $\Delta\theta_0 = \mathbf{B}\mathbf{A}$, $\mathbf{B} \in \mathbb{R}^{d \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times k}$, $r \ll \min(d, k)$. During adaptation, θ_0 is frozen while \mathbf{B} and \mathbf{A} are updated.

Adapter-based adaptation inserts sub-networks with a few learnable parameters into the large model. **Adapter** [28] consists of a pair of downsampling and upsampling layers as well as a residual connection. Suppose the original input to an intermediate layer θ_0 in model F is $\mathbf{x}_0 \in \mathbb{R}^{d_0}$, adapters insert a downsampling layer $\theta^D \in \mathbb{R}^{d_0 \times d_1}$ and an upsampling layer $\theta^U \in \mathbb{R}^{d_1 \times d_0}$, where d_0, d_1 are dimensions of the hidden embeddings, respectively. The output after injecting adapters is defined as $h = f_{\theta^U}(\sigma(f_{\theta^D}(\mathbf{x}_0))) + \mathbf{x}_0$, where f_{θ^U} denotes a function parameterized by θ^U and $\sigma(\cdot)$ is an activation function such as GELU [26]. To further reduce redundant parameters in adapters, **Compacter** [32] decomposes parameter matrices. It introduces *parameterized hypercomplex multiplication* (PHM) layers $\theta^D = \sum_{i=1}^k \mathbf{A}_i \otimes \mathbf{B}_i$, $\mathbf{A}_i \in \mathbb{R}^{k \times k}$, $\mathbf{B}_i \in \mathbb{R}^{\frac{d_0}{k} \times \frac{d_1}{k}}$, which decompose the layer in the adapter by Kronecker products. **Hyperformer** [40] relies on a hyper-network shared across tasks to generate the weights in adapters given a task index and a layer index. The hyper-

former maintains learnable embeddings for each task and each layer. For N_T tasks and N_L layers, the learnable embeddings in the hyperformer is denoted as $\mathbf{t}_1, \dots, \mathbf{t}_{N_T} \in \mathbb{R}^{d_e}$ and $\mathbf{l}_1, \dots, \mathbf{l}_{N_L} \in \mathbb{R}^{d_e}$, respectively. The hyperformer consists of a task projector $\theta^T \in \mathbb{R}^{(d_e + d_e) \times d_p}$ and a hyper-network $\theta^H \in \mathbb{R}^{d_p \times (d_0 \times d_1 + d_1 \times d_0)}$, and generates an adapter's weights in the i^{th} layer for the j^{th} task following $[\theta^D, \theta^U] = f_{\theta^H}(f_{\theta^T}([\mathbf{t}_j, \mathbf{l}_i]))$.

Adaptation shared over tasks [54] aims to further reduce redundant parameters by exploiting similar information shared across multiple tasks. In a multi-VL-task setting, an intuitive way is to train the adaptation modules per task using: **Multiple Adapters**, **Multiple Compacters**, **Multiple LoRA**, and **Multiple Prompts**. Additionally, We can train only one set of adapter layers for all tasks, and we have **Single Adapter**, **Single Compacter**, **Single LoRA**, and **Single Prompt**. Besides, **Half-shared adapter** [54] only shares the upsampling layers or downsampling layers across different tasks. Detailed information is presented in Supplementary Section 2.

4.1. CLIP-BART/T5

CLIP-BART/T5 [54], a combination of CLIP and BART/T5, follows VL-T5 [6] which unifies VL tasks to a text-generation problem. Specifically, given a pair of an image x^I and a sentence x^S , e.g. a picture and corresponding question texts in VQA, as the input to the model, CLIP-BART/T5 aims to maximize the agreement between the prediction and text label of M tokens $\mathbf{y} = (y_1, y_2, \dots, y_M)$. The primary generative model is an encoder-decoder language model, parameterized by θ^L . Visual representation from input images is extracted from a CLIP and a visual projection layer, parameterized by θ^V and $\theta^{V \rightarrow L}$ respectively. The concatenation of visual representation and sentence representation is fed into the encoder-decoder language model. The training goal is to minimize the cross-entropy loss [54]:

$$\begin{aligned} l(\mathbf{x}^I, \mathbf{x}^S, \mathbf{y}; \theta^L, \theta^V, \theta^{V \rightarrow L}) \\ = \text{CE}(f_{\theta^L}(\mathbf{x}^{V \rightarrow L}, \mathbf{x}^S), \mathbf{y}) \\ = - \sum_{i=1}^M y_i \log(f_{\theta^L}(\mathbf{x}^{V \rightarrow L}, \mathbf{x}^S)_i), \end{aligned} \quad (1)$$

where f_{θ} means a function parameterized by θ , $\mathbf{x}^{V \rightarrow L}$ is the projected visual representation.

The unified structure is beneficial to multi-task training where a universal dataset \mathcal{D} from N VL datasets is constructed and used to train the VL model. Under such a scenario, the parameters are optimized by minimizing the averaging loss on \mathcal{D} [54]:

Table 3: We deploy eleven distinct adaptation methods in total.

Type	Method
Full Fine-tuning	Full Fine-tuning
Adapter [28]	Single Adapter
	Half-shared Adapters
	Multiple Adapters
Compacter [32]	Single Compacter
	Multiple Compacter
Hyperformer [40]	Hyperformer
LoRA [29]	Single LoRA
	Multiple LoRA
Soft Prompt [33]	Single Prompt
	Multiple Prompts

$$\begin{aligned} \mathcal{L}(\mathcal{D}; \theta^L, \theta^V, \theta^{V \rightarrow L}) \\ = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}^I, \mathbf{x}^S, \mathbf{y}) \in \mathcal{D}} l(\mathbf{x}^I, \mathbf{x}^S, \mathbf{y}; \theta^L, \theta^V, \theta^{V \rightarrow L}) \end{aligned} \quad (2)$$

4.2. Adaptation Methods

Full fine-tuning directly updates the whole θ on \mathcal{D} and becomes prohibitive due to the rapidly growing model size. For instance, simply loading a GPT-3 language model with 175B parameters as the VL model backbone would require 700GB of memory¹. Therefore, the following more efficient adaptation methods are developed to achieve comparable performance while optimizing only a few parameters.

Prompt-based adaptation modifies the input \mathbf{x} to the model F by either concatenating a trainable prefix (Soft Prompt) [33] or a manually designed prefix [2]. For the given input $\mathbf{x} = \{x_1, \dots, x_n\}$ with n tokens, the pre-trained model will first form an embedding matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ where d is the dimension of the embedding space. Soft-prompts [33] are then represented as a learnable parameter $\mathbf{P} \in \mathbb{R}^{p \times d}$ where p is the length of the prompt. Next, \mathbf{P} is concatenated with the original embedded input \mathbf{X} to form a new single matrix defined as $[\mathbf{P}; \mathbf{X}] \in \mathbb{R}^{(p+n) \times d}$. During adaptation, the model is trained to maximize the probability of the desired output while only updating \mathbf{P} .

LoRA [29] also freezes the pre-trained model parameters θ , but it utilizes low-rank decomposition matrices to update gradients. For an intermediate model parameter $\theta_0 \in \mathbb{R}^{d \times k}$, which can be the parameters from one self-attention module in the transformer architecture, its update

is represented by a low-rank decomposition as shown in Formula 3. θ_0 is frozen, whereas B and A contain trainable parameters while adapting.

$$\theta_0 + \Delta\theta = \theta_0 + BA, B \in \mathbb{R}^{d \times r}, A \in \mathbb{R}^{r \times k}, r \ll \min(d, k). \quad (3)$$

Adapter-based adaptation inserts small modules between parameters in θ and modifies the intermediate learning process. Variants of adapter-based methods differ in the insertion manner. Adapters[28] consists of a pair of downsampling and upsampling layers and a residual connection. Suppose the original input to an intermediate layer θ_0 in model θ is $\mathbf{x}_0 \in \mathbb{R}^{d_0}$, adapters insert a downsampling layer $\theta^D \in \mathbb{R}^{d_0 \times d_1}$ and an upsampling layer $\theta^U \in \mathbb{R}^{d_1 \times d_0}$ where d_0, d_1 are dimensions of the hidden embeddings. The output after injecting adapters is defined in Formula 4 where $\sigma(\cdot)$ is an activation function such as GELU [26].

$$h = f_{\theta^U}(\sigma(f_{\theta^D}(\mathbf{x}_0))) + \mathbf{x}_0 \quad (4)$$

Compacter [32] is based on the mechanism of adapters [28], but it utilizes matrix decomposition and parameter sharing to further reduce redundant parameters in adapters. It introduces *parameterized hypercomplex multiplication* (PHM) layers (Formula 5), which decompose the layer in the adapter by Kronecker products. Compacter also shares the parameter of A_i across all layers and decomposes B_i even further with low-rank decomposition. However, as found in [54], such sharing and further decomposition severely decreases the VL performance. In our study, we only use PHM layers.

$$\theta^D = \sum_{i=1}^k A_i \otimes B_i, A_i \in \mathbb{R}^{k \times k}, B_i \in \mathbb{R}^{\frac{d_0}{k} \times \frac{d_1}{k}} \quad (5)$$

Hyperformer [40] also aims to reduce redundant parameters in adapters. It relies on a hyper-network that is shared across tasks to generate the weights in adapters given a task index and a layer index. The hyper-network maintains learnable embeddings for each task and each layer. For N_T tasks and N_L layers, the d_e -dimensional embeddings can be denoted as $\mathbf{t}_1, \dots, \mathbf{t}_{N_T} \in \mathbb{R}^{d_e}; \mathbf{l}_1, \dots, \mathbf{l}_{N_L} \in \mathbb{R}^{d_e}$. The Hyperformer consists of a task projector $\theta^T \in \mathbb{R}^{(d_e + d_e) \times d_p}$ and a hyper-network $\theta^H \in \mathbb{R}^{d_p \times (d_0 \times d_1 + d_1 \times d_0)}$, and generates an adapter's weights in the i^{th} layer for the j^{th} task following Formula 6.

$$[\theta^D, \theta^U] = f_{\theta^H}(f_{\theta^T}([\mathbf{t}_j, \mathbf{l}_i])) \quad (6)$$

Adaptation shared over tasks is inspired by Hyperformer and proposed in [54] which aims to exploit similar information shared across multiple tasks and to reduce redundant parameters. For vanilla adapters in a multi-task setting with N_T tasks, the collection of all inserted adapter

¹ $(175 \times 10^9) \times 4(\text{bytes}) \times 10^{-9} = 700\text{GB}$

modules can be denoted as $\Theta = \{\Theta^D, \Theta^U\}$ where Θ^D (Θ^U) stands for the subset of downsampling (upsampling) layers in adapters. The straightforward application is to train the adaptation modules per task, so we have independent $\{\Theta_i^D, \Theta_i^U\}$ for the i^{th} task, dubbed as **Multiple Adapters**. The same goes for other adaptation methods. By training one prompt, low-rank weights, and compacter layer for each task, we obtain **Multiple Prompts**, **Multiple LoRA**, and **Multiple Compacter** respectively. We can also train only one set of adapter layers for all tasks, and we have **Single Adapter** where $\Theta_i^D = \Theta_j^D$ and $\Theta_i^U = \Theta_j^U, i \neq j$. Also, if we use the same prompts, low-rank weights, and compacter layers, we have **Single Prompt**, **Single LoRA**, and **Single Compacter** respectively. For adapters, we can make parts of the weights shareable by making $\Theta_i^D = \Theta_j^U, i \neq j$ and. In this way, the task-specific information of i^{th} can still be learned by the rest Θ_i^U . Such sharing mechanism is called **half-shared adapter**. Lastly, Hyperformer already shares information from multiple tasks and therefore does not have such extensions.

4.3. Training and Evaluation

We follow the same experimental and hyperparameter settings as [54]. CLIP-ResNet101 is the vision encoder that takes the resized 224×224 images as input. The 7×7 gird features in the last convolutional layer are extracted as visual features and are downsampled to 6×6 by adaptive maximum pooling. $BART_{base}$ and $T5_{base}$ are both studied in this work as encoder-decoder language models. During training, AdamW is the optimizer along with a linear decay scheduler. Models are trained for 20 epochs, and the learning rate increases from 0 to the highest learning rate in the first 2 epochs. Training batch sizes are set as 500 and 250 for CLIP-BART and CLIP-T5 respectively. Models are trained in a multi-task setting where the training dataset includes all training split from 4 VL datasets. After training, we first evaluate the clean performance. Specifically, accuracy on the Karpathy-test split is evaluated for VQAv2. For GQA, accuracy on the test-dev split is evaluated, and accuracy on the test-P split is used for NLVR². In image captioning, we use CIDEr [60] on the Karpathy-test split. Then, we evaluate the corrupted performance on the corresponding corrupted test split from each dataset and calculate the relative robustness. All the training and evaluations are conducted on LRZ AI Systems at the Leibniz Supercomputing Centre of the Bavarian Academy of Sciences and Humanities.

5. Additional Analysis

5.1. BART vs. T5 language backbone

We compared the robustness of CLIP-BART and CLIP-T5 models in Table 6 to observe differences in performance.

Table 4: Dataset Statistics.

	VQAv2 [19]		GQA [30]		NLVR ² [53]		MSCOCO Caption [4]	
The Number of	Images	QA pairs	Images	QA pairs	Images	QA pairs	Images	Captions
Training set	113.2K	605.1K	72.1K	943.0K	103.2K	86.4K	113.2K	566.8K
Validation set	5.0K	26.7K	10.2K	132.1K	8.1K	7.0K	5.0K	5.0K
Test set	5.0K	26.3K	398	12.6K	8.1K	7.0K	5.0K	5.0K

Table 5: Relative robustness of adaptation methods based on CLIP-T5 against image (top) and text (bottom) corruptions.

Adaptation method	VQAv2		GQA		NLVR ²		COCO Caption	
	Acc (%)	RR (%)	Acc (%)	RR (%)	Acc (%)	RR (%)	CIDEr	RR (%)
Full Fine-tuning	66.29	85.11 ± 5.10	56.82	87.48 ± 0.04	74.06	89.36 ± 0.04	111.50	69.32 ± 0.14
Multiple Adapters	66.15	85.45 ± 4.84	55.66	87.70 ± 0.04	51.94	—	112.15	67.65 ± 0.15
Single Adapter	66.41	85.19 ± 5.16	55.90	78.57 ± 0.04	72.78	88.70 ± 0.08	111.70	68.52 ± 0.14
Hyperformer	65.18	86.02 ± 4.87	54.65	89.33 ± 0.04	70.56	91.07 ± 0.08	110.65	70.48 ± 0.13
Multiple Compacters	65.50	85.88 ± 4.86	54.68	88.09 ± 0.04	52.63	—	113.20	67.96 ± 0.14
Single Compacter	65.98	85.56 ± 5.07	55.33	80.57 ± 0.04	71.47	90.04 ± 0.04	111.61	69.15 ± 0.14

Adaptation Method	VQAv2		GQA		NLVR ²	
	Acc (%)	RR (%)	Acc (%)	RR (%)	Acc (%)	RR (%)
Full fine-tuning	66.29	72.52 ± 25.29	56.82	64.27 ± 25.75	74.06	87.67 ± 10.74
Multiple Adapters	66.15	76.68 ± 21.21	55.66	62.33 ± 25.92	51.94	—
Single Adapter	66.41	75.89 ± 21.20	55.90	63.30 ± 19.40	72.78	87.79 ± 10.68
Hyperformer	65.18	76.51 ± 20.80	54.65	66.96 ± 24.43	70.56	89.12 ± 10.06
Multiple Compacters	65.50	76.39 ± 21.14	54.68	67.66 ± 22.69	52.63	—
Single Compacter	65.98	76.16 ± 21.16	55.33	64.19 ± 20.22	71.47	88.70 ± 9.93

Generally, adaptations from CLIP-BART showed better robustness on GQA, NLVR², and COCO Caption against image and text corruptions. Specifically, all adaptation methods on CLIP-BART had higher robustness against text corruptions on the GQA dataset. For robustness against image corruptions, all adaptation methods with CLIP-BART, except for multiple adapters, achieved higher relative robustness scores. In contrast, CLIP-T5-based adaptations were more robust against image corruptions on VQAv2, while CLIP-BART showed more robustness against text corruptions. This may be due to the different language encoders used in BART and T5. Among all adaptation methods, Hyperformer seems to be a good choice for CLIP-T5, as it achieved better robustness on VQAv2, NLVR², and COCO Caption. A single adapter would likely be a more robust adaptation method when combined with CLIP-BART, as it showed better robustness on all datasets and against both types of corruption.

5.2. The Influence of Adaptation Hyperparameters on Robustness

To investigate the influence of parameter size on robustness, we conduct experiments on 6 adaptation methods, namely Single Prompt, Multiple Prompts, Single Adapter, Single Compacter, and Single LoRA, with various parameter sizes. For each setting, we adapt the model on multitask datasets given pre-trained CLIP and BART (T5) and test the relative robustness.

For prompt tuning, the prompt length p defined in Section 4.2 can be 20, 40, 60, and 80. The position of prompt embedding can be *front*, *middle*, and *back*. It specifies the

Table 6: RR(%) of adaptation methods based on CLIP-BART and CLIP-T5 against image (up) and text (down) corruptions with severity 5. The better relative robustness values for each comparison pair are in bold.

Adaptation Method	VQAv2		GQA		NLVR ²		COCO Caption	
	BART	T5	BART	T5	BART	T5	BART	T5
Full finetuning	84.86	85.11	89.20	87.48	90.34	89.36	68.40	69.32
Multiple adapters	85.33	85.45	86.16	87.70	92.02	—	68.72	67.65
Single adapter	85.76	85.19	82.49	78.57	90.04	88.70	68.68	68.52
Hyperformer	85.38	86.02	90.05	89.33	90.13	91.07	68.74	70.48
Multiple compacters	85.65	85.88	88.89	88.09	91.33	—	68.67	67.96
Single compacter	85.47	85.56	82.62	80.57	92.04	90.04	69.92	69.15

Adaptation method	VQAv2		GQA		NLVR ²	
	BART	T5	BART	T5	BART	T5
Full Fine-tuning	73.65	72.52	66.92	64.27	87.06	87.67
Multiple Adapters	76.62	76.68	66.93	62.33	90.14	—
Single Adapter	77.64	75.89	67.47	63.30	88.49	87.79
Hyperformer	75.06	76.51	70.30	66.96	87.27	89.12
Multiple Compacters	77.10	76.39	67.39	67.66	90.00	—
Single Compacter	77.17	76.16	67.90	64.19	90.10	88.7

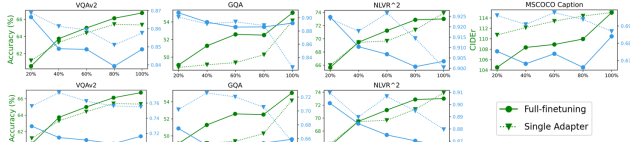


Figure 4: Performance and relative robustness of full-finetuning and single adapter on CLIP-BART given different size of adaptation dataset. The first row shows results given image corruptions and the second is from text corruptions. Green lines stand for performance in each task and the blue is robustness. X-axis shows the random subset ratio of training dataset during adaptation, ranging from 20% to 100%.

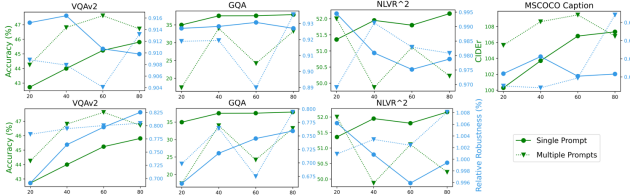


Figure 5: Performance and relative robustness from prompt adaptations with different prompt lengths on CLIP-BART. The top row shows the robustness against image corruptions and the bottom row is results against text corruptions. Blue lines stand for performance on each task and green lines represent relative robustness.

prompt position in the concatenated input to the encoder-decoder generative model. For example, *front* position means that the prompt embedding is at the beginning of the concatenation; *middle* position means the concatenation starts with visual embedding, then prompt embedding and the textual embedding is at the end. We also choose differ-

ent embedding dimensions for adapter-based methods. The dimension d_1 of adapter layers varies from $\{96, 64, 48, 24\}$. The dimension r of LoRA’s low-rank decomposition layer can be $\{32, 64, 128, 256\}$

Prompt embedding length Prompt-based adaptation methods have attracted more attention because they remove the burden of tuning any parameters in the pre-trained model. We inspect their robustness given different embedding lengths and sharing settings. We adjust the soft prompt length added to the concatenated embeddings and evaluate the performance along with the relative robustness. The results conducted on CLIP-BART are shown in Figure 5. We could see a steady increase in the performance on four tasks with longer prompt lengths which proves that prompt methods perform better given more parameters. Regarding relative robustness, such a steady increase does not apply to all tasks and longer soft prompts do not ensure better relative robustness. For instance, the relative robustness of single prompt drops with the increase of prompt length on VQAv2. Its relative robustness keeps relatively the same on GQA and MSCOCO Caption given image corruptions. We could also notice that for text corruptions, prompt methods gain better robustness on VQAv2 and GQA.

Task-specific vs. Universal Prompt. We conduct experiments using both single prompt and multiple prompts where the first shares a universal prompt embedding across all tasks and the second utilizes a specific prompt for each task. The single prompt has better performance on GQA and NLVR² whereas multiple prompts show better performance on VQAv2 and MSCOCO Caption. We observe that with less prompt embedding, a universal prompt is more robust against image corruptions but given longer prompt embedding, task-specific prompts regain the robustness and surpass the universal prompt.

5.3. Additional Results

The relative robustness of adaptation methods based on CLIP-T5 is presented in Table 5. Figure 5 presents the performance and relative robustness of full-finetuning and single adapter on CLIP-BART given different sizes of adaptation datasets. Table 9 and 8 present the corruption results from each image corruption method and each text corruption category. Table 11 and the left tables present detailed results on various severity levels.

Table 7: Relative robustness of adaptation methods based on CLIP-BART against image corruptions.

		Noise										Digital										Weather									
Adaptation method	Clean	ImageNet	Gaussian	Blur	Speckle	Zoom	Defocus	Motion	Gaussian	JPEG	Contrast	Defocus	Pixelate	Poster	Spatter	Snow	Frost	BBG	Brightness	RR	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	RR	RR	RR	RR
Full Finetuning	115.03	81.57	81.28	82.44	84.18	82.44	86.1	78.51	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Multiple Adapters	114.47	81.82	81.45	81.48	81.94	82.04	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Half-shared Adapters	114.5	79.94	79.94	84.56	84.56	84.56	85.72	73.87	64.18	82.16	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Single Adapter	115.08	80.28	81.57	84.01	84.01	84.01	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Hyperparameter	114.66	81.08	81.24	84.08	84.08	84.08	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Multiple Compacters	115.16	80.28	81.57	84.01	84.01	84.01	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Single Compacter	113.06	81.79	80.07	83.04	83.04	83.04	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Multiple LoRA	115.07	80.28	81.57	84.01	84.01	84.01	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Single LoRA	114.47	81.82	81.45	81.48	81.94	82.04	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	78.10	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04	66.04
Single Prompt	105.7	72.47	72.31	73.51	73.51	73.51	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	69.75	62.67	62.67	62.67	62.67	62.67	62.67	62.67	62.67	62.67	62.67
Multiple Prompts	110.49	78.88	78.48	78.48	78.48	78.48	86.29	78.44	65.06	83.06	100.57	64.22	71.91	62.57	72	68.52	72.81	71.96	74.86	100.19	70.49	63.41	63.41	63.41	63.41	63.41	63.41	63.41	63.41	63.41	63.41
Adaptation method	Clean <th colspan="10">Noise</th> <th colspan="10">Digital</th> <th colspan="10">Weather</th>	Noise										Digital										Weather									
ImageNet	Gaussian	Blur	Speckle	Zoom	Defocus	Motion	Gaussian	JPEG	Contrast	Defocus	Pixelate	Poster	Spatter	Snow	Frost	BBG	Brightness	RR	RR	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	RR	RR	RR	RR	RR
Full Finetuning	55.48	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Multiple Adapters	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Half-shared Adapters	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Single Adapter	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Hyperparameter	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Multiple Compacters	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Single Compacter	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Multiple LoRA	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Single LoRA	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Single Prompt	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	
Multiple Prompts	53.39	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	49.15	

Table 10: Relative robustness of adaptation methods based on CLIP-T5 against text corruptions grouped by corruption levels.

Adaptation method	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	RR	RR
Full finetuning	56.82	21.13	0.372	38.23	0.673	53.54	0.942	0.636	0.623
Multiple adapters	55.66	20.43	0.367	35.82	0.644	52.07	0.936	0.623	0.623
Single adapter	55.9	23.99	0.429	36.69	0.656	47.77	0.855	0.633	0.633
Hyperparameter	54.63	22.69	0.415	38.19	0.699	51.89	0.949	0.677	0.677
Single compacter	55.33	23.5	0.425	36.88	0.667	48.64	0.879	0.642	0.642
Multiple compacters	54.68	24.28	0.444	38.14	0.698	51.99	0.951	0.677	0.677
Adaptation method	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	RR	RR
Full finetuning	74.06	56.18	0.759	66.37	0.896	73.2	0.988	0.877	0.877
Multiple adapters	51.94	51.09	0.984	51.78	0.997	51.92	1	0.994	0.994
Single adapter	72.76	54.83	0.753	65.69	0.903	71.53	0.983	0.878	0.878
Hyperparameter	70.58	54.83	0.777	64.46	0.914	69.75	0.989	0.891	0.891
Single compacter	71.47	55.78	0.788	64.64	0.904	70.63	0.988	0.887	0.887
Multiple compacters	52.63	52.67	1.001	52.58	0.999	52.69	1.001	1	1
Adaptation method	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	RR	RR
Full finetuning	66.29	30.26	0.456	51.05	0.777	64.55	0.974	0.725	0.725
Multiple adapters	66.14	34.38	0.516	53.89	0.815	64.46	0.977	0.725	0.725
Single adapter	66.41	34.38	0.518	52.91	0.812	63.59	0.976	0.765	0.765
Hyperparameter	65.98	34.1	0.517	53.18	0.806	64.4	0.976	0.764	0.764
Single compacter	65.58	34.39	0.525	52.77	0.806	64.07	0.978	0.764	0.764

Table 8: Relative robustness of adaptation methods based on CLIP-BART against text corruptions grouped by corruption levels.

Adaptation method	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	Ave RR
Full Finetuning	55.04	23.47	0.426	37.9	0.689	52.91	0.961	0.669
Multiple Adapters	53.39	23.19	0.434	37.37	0.7	48.97	0.917	0.669
Half-shared Adapters	52.96	22.35	0.422	38.06	0.719	50.22	0.948	0.682
Single Adapter	54.14	24.11	0.445	38.54	0.712	48.18	0.89	0.675
Hyperformer	52.52	23.69	0.451	38.85	0.74	50.25	0.957	0.703
Multiple Compacters	52.75	22.8	0.432	37.21	0.705	49.13	0.931	0.674
Single Compacter	52.9	23.89	0.452	37.8	0.715	47.49	0.898	0.679
Multiple LoRA	52.05	22.79	0.438	37.74	0.725	48.72	0.936	0.688
Single LoRA	53.19	22.94	0.431	34.98	0.658	47.12	0.886	0.639
Single Prompt	37.54	21.68	0.586	26.94	0.718	36	0.959	0.725
Multiple Prompts	34.01	19.69	0.579	27.6	0.812	31.62	0.93	0.731
Adaptation method	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	Ave RR
Full Finetuning	73.01	55	0.753	64.88	0.889	71.99	0.986	0.871
Multiple Adapters	69.41	53.54	0.771	64.66	0.932	69.12	0.996	0.901
Half-shared Adapters	70.03	53.53	0.764	64.47	0.936	69.02	0.992	0.892
Single Adapter	73.89	56.24	0.761	67.19	0.909	73.09	0.989	0.885
Hyperformer	72.21	54.24	0.751	64.59	0.894	72.02	0.984	0.873
Multiple Compacters	69.45	53.84	0.775	64.54	0.929	68.67	0.989	0.9
Single Compacter	69.94	55.07	0.787	64.65	0.924	69.4	0.992	0.901
Multiple LoRA	51.32	51.41	1.002	51.52	1.004	51.25	0.999	1.002
Single LoRA	72.58	55.13	0.969	0.9	9	72.67	0.988	0.876
Single Prompt	51.91	51.87	1.001	52.01	1.001	52.1	1.003	1.003
Multiple Prompts	49.87	50.1	1.005	50.1	1.005	49.84	0.999	1.004
Adaptation method	Clean	Char-level	Char-level RR	Word-level	Word-level RR	Sentence-level	Sentence-level RR	Ave RR
Full Finetuning	66.75	32.23	0.484	51.99	0.779	64.68	0.969	0.736
Multiple Adapters	65.14	33.76	0.532	51.77	0.815	65.85	0.966	0.766
Half-shared Adapters	65.2	33.75	0.518	53.28	0.817	63.48	0.974	0.768
Single Adapter	65.35	33.65	0.515	54.42	0.833	63.83	0.977	0.776
Hyperformer	65.38	33.62	0.513	51.6	0.789	63.69	0.974	0.751
Multiple Compacters	64.91	33.63	0.518	53.38	0.822	63.27	0.975	0.771
Single Compacter	64.47	33.76	0.524	53.12	0.824	62.45	0.966	0.767
Multiple LoRA	65.34	33.51	0.512	53.51	0.772	63.61	0.974	0.74
Single LoRA	65.34	32.89	0.503	51.27	0.785	63.25	0.97	0.745
Single Prompt	44	28.89	0.657	33.49	0.761	42.75	0.972	0.77
Multiple Prompts	46.81	29.35	0.627	38.88	0.831	45.4	0.977	0.802

Table 12: Performance and decrease ratio of Half-shared Adapters on CLIP-BART against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	
impulse_noise	1	62.84	-2.36	-3.620%	51.654	-1.306	-2.460%	69.241	-0.789	-1.127%	109.345	-5.155	-4.502%				
gaussian_noise	1	63.97	-1.23	-1.887%	52.6	-0.60	-0.680%	70.217	-0.187	-0.267%	112.308	-2.192	-1.915%				
shot_noise	1	64.04	-1.16	-1.779%	51.701	-1.259	-2.377%	70.26	0.23	0.328%	111.183	-3.317	-2.897%				
speckle_noise	1	64.15	-1.05	-1.610%	52.139	-0.821	-1.551%	69.255	-1.107	-1.613%	111.887	-3.313	-2.893%				
zoom blur	1	58.33	-6.87	-10.577%	49.507	-3.453	-6.939%	57.227	-12.033	-18.282%	85.69	-28.88	-33.126%				
defocus blur	1	64.04	-1.16	-1.779%	52.155	-0.805	-1.521%	70.174	0.144	0.205%	110.785	-3.315	-3.244%				
motion blur	1	64.08	-1.12	-1.718%	52.457	-0.503	-0.950%	70.116	0.086	0.123%	110.903	-4.407	-3.849%				
glare blur	1	64.14	-1.06	-1.625%	52.385	-0.575	-1.086%	70.389	0.359	0.515%	111.411	-3.699	-3.249%				
gaussian blur	1	64.9	-0.3	-0.460%	52.95	-0.01	-0.020%	70.274	0.244	0.349%	113.665	-0.835	-0.729%				
jpeg_compression	1	64.63	-0.57	-0.874%	52.457	-0.503	-0.950%	70.36	0.33	0.472%	113.127	-1.373	-1.200%				
contrast	1	63.76	-1.44	-2.289%	52.027	-0.933	-1.761%	69.901	-0.128	-0.184%	111.485	-3.005	-2.643%				
elastic_transform	1	64.43	-0.77	-1.181%	52.274	-0.686	-1.296%	69.8	-0.23	-0.328%	112.4	-2.1	-1.834%				
pixelate	1	61.78	-3.42	-5.526%	51.137	-1.832	-3.424%	69.112	-0.918	-1.312%	102.902	-12.408	-10.836%				
snow	1	60.28	-4.92	-7.546%	50.477	-2.483	-4.888%	68.265	-1.765	-2.521%	99.576	-14.964	-13.069%				
frost	1	61.47	-3.73	-5.725%	50.469	-2.491	-4.793%	68.451	-1.579	-2.254%	103.326	-11.174	-9.759%				
fog	1	63.31	-1.89	-2.899%	51.964	-0.996	-1.881%	70.274	0.244	0.349%	111.251	-3.249	-2.875%				
brightness	1	64.79	-0.41	-0.629%	52.663	-0.297	-0.560%	70.748	0.718	1.025%	113.052	-4.148	-3.625%				
saturation	1	64.51	-0.69	-1.058%	52.552	-0.408	-0.770%	70.317	0.287	0.410%	112.966	-1.534	-1.349%				
saturation	1	62.37	-2.83	-4.349%	51.638	-1.322	-2.497%	69.269	-0.781	-1.086%	110.372	-4.128	-3.660%				
corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	
impulse_noise	2	62.14	-3.06	-4.891%	50.954	-2.006	-3.788%	67.978	-2.052	-2.931%	105.6	-8.9	-7.773%				
gaussian_noise	2	63.21	-1.99	-3.027%	51.467	-1.473	-2.782%	69.413	-0.617	-0.881%	109.248	-5.252	-4.586%				
shot_noise	2	63.13	-2.07	-3.175%	51.773	-1.187	-2.241%	68.939	-1.091	-1.577%	107.995	-6.505	-5.681%				
speckle_noise	2	63.53	-1.67	-2.561%	51.693	-1.267	-2.392%	69.399	-0.631	-0.902%	110.78	-3.72	-3.249%				
zoom blur	2	55.07	-8.53	-14.417%	47.604	-5.266	-9.843%	61.806	-4.224	-11.744%	69.462	-45.038	-39.334%				
defocus blur	2	63.37	-1.83	-2.807%	51.383	-1.577	-2.977%	70.001	-0.029	-0.044%	110.982	-4.184	-4.732%				
motion blur	2	62.96	-2.24	-3.436%	51.582	-1.728	-2.602%	69.872	-0.158	-0.225%	108.91	-7.609	-6.646%				
glare blur	2	63.17	-2.03	-3.113%	51.717	-1.243	-2.347%	69.126	-0.904	-1.291%	107.2	-7.3	-6.706%				
gaussian blur	2	63.88	-1.32	-2.025%	51.852	-1.108	-2.099%	70.231	0.201	0.278%	110.2	-4.3	-3.756%				
jpeg_compression	2	64.26	-0.94	-1.442%	52.083	-0.877	-1.556%	69.585	-0.445	-0.635%	112.207	-2.203	-1.924%				
contrast	2	63.14	-2.06	-3.160%	51.844	-1.116	-2.106%	69.327	-0.703	-1.004%	109.898	-5.402	-4.718%				
elastic_transform	2	61.55	-3.65	-5.598%	50.517	-2.443	-4.631%	67.547	-2.483	-3.546%	97.081	-47.119	-35.213%				
pixelate	2	61.27	-3.93	-6.020%	50.731	-2.29	-4.208%	66.284	-1.766	-2.540%	98.94	-44.966	-34.810%				
snow	2	57.4	-7.8	-11.863%	47.162	-5.798	-10.848%	66.815	-3.215	-4.914%	88.051	-49.449	-23.099%				
frost	2	58.16	-7.04	-10.798%	48.447	-4.463	-8.946%	67.002	-3.028	-4.324%	91.679	-22.621	-19.931%				
fog	2	62.61	-2.59	-3.972%	51.615	-1.545	-2.917%	65.638	-0.492	-0.727%	109.042	-5.458	-4.767%				
brightness	2	64.2	-1.0	-1.534%	52.003	-0.957	-1.806%	70.274	0.244	0.349%	112.985	-1.515	-1.323%				
saturation	2	63.27	-3.97	-6.099%	50.859	-2.101	-3.940%	69.212	-0.838	-1.168%	107.425	-11.575	-10.109%				
saturation	2	59.61	-5.59	-8.574%	50.708	-2.252	-4.253%	67.088	-2.942	-4.201%	104.939	-3.561	-3.350%				
corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	
impulse_noise	3	61.32	-3.88	-5.951%	50.175	-2.785	-5.259%	68.222	-1.868	-2.582%	103.077	-11.423	-9.977%				
gaussian_noise	3	61.62	-3.58	-5.691%	50.549	-2.411	-4.553%	68.58	-1.45	-2.076%	105.052	-9.448	-8.515%				
shot_noise	3	61.49	-3.71	-5.690%	51.264	-1.696	-3.202%	68.164	-1.906	-2.664%	103.781	-10.719	-9.362%				
speckle_noise	3	61.66	-3.24	-4.969%	50.867	-2.093	-3.950%	68.509	-1.521	-2.172%	104.433	-10.867	-8.726%				
zoom blur	3	52.68	-12.52	-19.202%	45.675	-7.285	-13.766%	59.165	-10.865	-15.515%	54.675	-39.825	-32.349%				
defocus blur	3	61.64	-3.56	-5.460%	50.875	-2.085	-3.938%	68.718	-1.292	-1.844%	102.239	-12.261	-10.708%				
motion blur	3	61.2	-4.0	-6.135%	51.153	-1.807	-3.412%	67.978	-2.052	-2.931%	108.62	-13.18	-12.123%				
glare blur	3	57.94	-2.86	-11.155%	47.742	-5.218	-9.853%	67.159	-0.871	-1.097%	97.072	-42.249	-23.656%				
gaussian blur	3	62.33	-2.87	-4.462%	51.423	-1.537	-2.898%	68.652	-1.378	-1.967%	106.626	-9.474	-8.624%				
jpeg_compression	3	63.99	-1.21	-1.856%	52.536	-0.824	-1.506%	69.643	-0.387	-0.553%	111.727	-2.425	-2.158%				
contrast	3	61.77	-3.43	-5.261%	50.755	-2.205	-4.163%	68.336	-1.694	-2.418%	105.203	-9.297	-8.120%				
elastic_transform	3	63.07	-2.2	-3.374%	51.574	-1.386	-2.417%	69.93	-0.3	-0.413%	107.179	-32.1	-24.975%				
pixelate	3	58.76	-6.44	-9.877%	49.372	-3.588	-6.775%	65.294	-3.766	-5.693%	90.913	-32.861	-20.600%				
snow	3	55.52	-9.68	-14.847%	47.289	-5.671	-10.708%	64.892	-5.138	-7.377%	79.752	-34.748	-30.347%				
frost	3	55.43	-9.77	-14.895%	46.772	-6.188	-11.044%	65.222	-4.804	-6.866%	82.156	-32.344	-28.248%				
fog	3	61.42	-3.78	-5.798%	51.328	-1.632	-3.082%	68.58	-1.45	-2.070%	105.091	-8.809	-7.604%				
brightness	3	63.4	-1.8	-2.714%	51.725	-1.235	-2.325%	69.274	-0.316	-0.453%	111.483	-3.017	-2.635%				
saturation	3	59.87	-5.33	-8.175%	49.96	-3.0	-5.646%	68.092	-1.958	-2.767%	96.967	-17.533	-15.319%				
saturation	3	64.25	-0.95	-1.457%	51.831	-1.139	-2.151%	70.518	0.488	0.691%	112.453	-2.047	-1.788%				
corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	
impulse_noise	4	58.84	-6.38	-9.755%	49.396	-3.364	-6.730%	67.303	-2.727	-3.894%	93.77	-20.73	-18.105%				
gaussian_noise	4	59.07	-6.13	-9.402%	49.833	-3.127	-5.904%	67.059	-2.971	-4.242%	95.573	-19.237	-16.530%				
shot_noise	4	58.48	-6.72	-10.307%	50.04	-2.92	-5.514%	65.896	-4.134	-5.903%	92.981	-21.519	-18.794%				
speckle_noise	4	60.43	-4.77	-7.101%	50.684	-2.276	-4.208%	67.317	-2.713	-3.874%	105.058	-13.902	-12.220%				
zoom blur	4	50.72	-14.48	-22.229%	44.538	-8.422	-15.020%	57.571	-12.459	-17.790%	44.394	-70.166	-61.228%				
defocus blur	4	59.77	-5.43	-8.328%	50.151	-2.809	-5.500%	67.246	-2.784	-3.976%	94.926	-19.574	-17.895%				
motion blur	4	58.71	-6.49	-9.954%	49.642	-3.318	-6.265%	66.413	-4.07	-5.165%	87.531	-26.969	-23.555%				
glare blur	4	56.08	-9.12	-13.988%	47.058	-5.901	-11.144%	66.485	-3.545	-5.062%	79.98	-34.52	-30.148%				
gaussian blur	4	60.1	-4.4	-6.748%	50.559	-2.451	-4.239%	67.978	-2.052	-2.931%	97.356	-17.144	-14.975%				
jpeg_compression	4	63.05	-1.25	-1.398%	51.169	-1.791	-3.382%	69.47	-0.56	-0.799%	107.974	-5.526	-5.099%				
contrast	4	57.65	-7.55	-11.580%	48.704	-4.256	-8.636%	65.423	-4.67	-6.579%	90.713	-23.787	-20.775%				
elastic_transform	4	60.82	-4.38	-6.718%	49.587	-3.73	-6.706%	69.47	-0.56	-0.799%	99.136	-15.364	-13.118%				
pixelate	4	54.89	-10.31	-15.813%	46.208	-6.752	-12.790%	61.678	-3.554	-5.192%	73.596	-40.004	-35.724%				
snow	4	52.27	-12.63	-19.202%	45.675	-7.285	-13.766%	59.165	-10.865	-15.515%	54.675	-39.825	-32.349%				
frost	4	52.02	-10.28	-15.767%	46.915	-6.045	-11.414%	65.431	-4.579	-6.388%	78.534	-35.976	-31.740%				
fog	4	59.85	-5.35	-8.266%	50.505	-2.817	-5.290%	67.317	-2.713	-3.874%	105.058	-13.902	-12.220%				
brightness	4	62.64	-1.8	-2.714%	51.725	-1.235	-2.325%	69.274	-0.316	-0.453%	111.483	-3.017	-2.635%				
saturation	4	57.78	-7.42	-11.389%	48.251	-4.79	-8.892%	66.808	-3.172	-4.529%	86.919	-27.841	-24.088%				
saturation	4	64.25	-0.95	-1.4													

Table 14: Performance and decrease ratio of Multiple Adapters on CLIP-BART against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	1	63.96	-1.34	-2.052%	50.525	-2.865	-5.307%	69.556	0.146	0.211%	112.085	-2.415	-2.110%	69.556	0.146	0.211%	112.085
gaussian_noise	1	64.1	-1.2	-1.838%	50.91	-2.46	-4.607%	69.628	0.218	0.314%	111.893	-2.377	-2.251%	69.628	0.218	0.314%	111.893
shot_noise	1	64.28	-1.02	-1.562%	51.018	-2.372	-4.443%	69.556	0.146	0.211%	112.372	-2.089	-1.833%	69.556	0.146	0.211%	112.372
speckle_noise	1	58.67	-6.63	-10.578%	47.822	-5.568	-10.430%	57.686	-11.724	-18.691%	86.317	-28.153	-24.549%	57.686	-11.724	-18.691%	86.317
defocus_blur	1	64.49	-0.81	-1.240%	50.602	-2.098	-5.044%	70.418	1.008	1.452%	110.992	-3.478	-3.030%	70.418	1.008	1.452%	110.992
motion_blur	1	64.31	-0.99	-1.516%	51.01	-2.38	-4.458%	69.901	0.491	0.707%	110.947	-3.523	-3.078%	69.901	0.491	0.707%	110.947
glass_blur	1	64.43	-0.87	-1.322%	51.057	-2.333	-4.599%	70.561	1.151	1.659%	112.562	-1.888	-1.649%	70.561	1.151	1.659%	112.562
gaussian_blur	1	65.1	-0.2	-0.306%	51.169	-2.221	-4.161%	70.217	0.807	1.162%	113.692	-0.778	-0.679%	70.217	0.807	1.162%	113.692
jpeg_compression	1	64.7	-0.6	-0.919%	51.312	-2.078	-3.892%	69.958	0.548	0.790%	113.934	-0.536	-0.469%	69.958	0.548	0.790%	113.934
contrast	1	63.91	-1.39	-2.129%	50.107	-3.223	-6.437%	69.657	0.247	0.356%	111.583	-2.992	-2.522%	69.657	0.247	0.356%	111.583
elastic_transform	1	64.3	-0.67	-1.026%	51.216	-2.174	-4.071%	69.427	0.071	0.025%	112.058	-2.412	-2.107%	69.427	0.071	0.025%	112.058
pixelate	1	61.72	-3.58	-5.826%	49.412	-3.978	-7.851%	68.466	0.944	1.361%	102.565	-11.905	-10.409%	68.466	0.944	1.361%	102.565
snow	1	60.6	-4.7	-7.188%	48.593	-4.797	-8.858%	68.193	-2.127	-3.115%	100.062	-14.408	-12.587%	68.193	-2.127	-3.115%	100.062
frost	1	61.74	-3.56	-5.452%	49.213	-4.177	-7.824%	68.394	-1.016	-1.464%	103.506	-10.964	-9.578%	68.394	-1.016	-1.464%	103.506
fog	1	63.34	-1.96	-3.082%	50.596	-2.794	-5.233%	69.241	-0.169	-0.244%	110.445	-4.025	-3.516%	69.241	-0.169	-0.244%	110.445
brightness	1	64.83	-0.47	-0.728%	51.177	-2.213	-4.146%	70.848	1.328	2.072%	113.991	-0.479	-0.419%	70.848	1.328	2.072%	113.991
saturation	1	64.72	-0.58	-0.889%	51.081	-2.309	-4.324%	70.461	1.051	1.514%	113.572	-0.808	-0.785%	70.461	1.051	1.514%	113.572
saturation	1	62.48	-2.82	-4.139%	50.588	-2.802	-5.248%	68.753	-0.405	-0.587%	110.884	-3.666	-3.203%	68.753	-0.405	-0.587%	110.884
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	2	62.36	-2.94	-4.502%	49.889	-3.501	-6.558%	68.466	-0.944	-1.361%	106.337	-7.933	-6.930%	68.466	-0.944	-1.361%	106.337
gaussian_noise	2	63.13	-1.97	-3.017%	50.183	-3.207	-6.007%	69.054	-0.256	-0.351%	109.084	-5.366	-4.750%	69.054	-0.256	-0.351%	109.084
shot_noise	2	63.37	-1.93	-2.956%	50.358	-3.032	-5.679%	68.724	-0.686	-0.988%	109.293	-5.177	-4.523%	68.724	-0.686	-0.988%	109.293
speckle_noise	2	63.79	-1.51	-2.312%	50.668	-2.722	-5.099%	68.695	-0.715	-1.030%	110.266	-4.204	-3.673%	68.695	-0.715	-1.030%	110.266
zoom_blur	2	55.81	-8.49	-14.335%	45.436	-7.954	-14.897%	61.633	-7.777	-11.204%	97.867	-41.603	-38.801%	61.633	-7.777	-11.204%	97.867
defocus_blur	2	63.91	-1.39	-2.129%	50.277	-3.12	-5.843%	69.485	0.075	0.088%	108.924	-5.486	-4.845%	69.485	0.075	0.088%	108.924
motion_blur	2	63.24	-2.06	-3.155%	50.533	-2.657	-5.252%	69.37	0.04	0.058%	108.19	-6.28	-5.487%	69.37	0.04	0.058%	108.19
glass_blur	2	62.78	-2.05	-3.129%	50.517	-2.873	-5.326%	69.7	0.29	0.418%	107.653	-6.817	-5.955%	69.7	0.29	0.418%	107.653
gaussian_blur	2	64.24	-1.06	-1.623%	50.429	-2.961	-5.545%	69.915	0.508	0.728%	109.886	-5.474	-3.996%	69.915	0.508	0.728%	109.886
jpeg_compression	2	64.64	-0.66	-1.011%	51.272	-2.118	-3.907%	69.7	0.29	0.418%	113.017	-1.453	-1.299%	69.7	0.29	0.418%	113.017
contrast	2	63.0	-2.3	-3.522%	50.024	-3.366	-6.305%	68.781	-0.620	-0.906%	109.711	-4.759	-4.149%	68.781	-0.620	-0.906%	109.711
elastic_transform	2	61.58	-3.72	-5.697%	49.229	-4.161	-7.794%	67.327	-1.083	-1.442%	99.907	-16.563	-14.886%	67.327	-1.083	-1.442%	99.907
pixelate	2	61.36	-3.94	-6.401%	49.084	-4.286	-8.076%	65.494	-3.916	-5.641%	98.54	-14.93	-13.042%	65.494	-3.916	-5.641%	98.54
snow	2	57.58	-7.72	-11.822%	46.152	-7.238	-13.557%	65.71	-3.7	-5.311%	88.051	-26.419	-23.097%	65.71	-3.7	-5.311%	88.051
frost	2	58.24	-6.76	-10.328%	47.384	-6.086	-11.249%	66.385	-3.825	-4.070%	99.346	-22.806	-19.923%	66.385	-3.825	-4.070%	99.346
fog	2	62.78	-2.52	-3.859%	49.968	-3.422	-6.409%	69.911	-1.499	-2.141%	109.097	-5.373	-4.693%	69.911	-1.499	-2.141%	109.097
brightness	2	64.17	-1.13	-1.730%	50.882	-2.508	-4.697%	70.159	0.749	1.080%	113.574	-0.806	-0.782%	70.159	0.749	1.080%	113.574
saturation	2	61.42	-3.88	-5.826%	49.595	-3.795	-7.109%	68.25	-1.16	-1.671%	104.738	-10.092	-8.816%	68.25	-1.16	-1.671%	104.738
saturation	2	59.52	-5.78	-8.831%	49.443	-3.947	-7.992%	67.03	-2.38	-3.428%	104.845	-9.825	-8.408%	67.03	-2.38	-3.428%	104.845
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	3	61.63	-3.67	-5.620%	49.952	-4.438	-6.439%	68.179	-1.231	-1.774%	102.223	-11.247	-9.825%	68.179	-1.231	-1.774%	102.223
gaussian_noise	3	61.83	-3.47	-5.314%	49.436	-3.954	-7.407%	67.848	-1.562	-2.250%	103.566	-10.314	-8.918%	67.848	-1.562	-2.250%	103.566
shot_noise	3	61.65	-3.65	-5.590%	49.913	-3.477	-6.513%	68.15	-1.26	-1.816%	105.464	-9.006	-7.888%	68.15	-1.26	-1.816%	105.464
speckle_noise	3	62.13	-3.17	-4.855%	49.634	-3.756	-7.245%	67.963	-1.447	-2.084%	104.688	-9.262	-8.145%	67.963	-1.447	-2.084%	104.688
zoom_blur	3	53.22	-12.08	-18.499%	43.671	-9.719	-18.203%	59.495	-9.915	-14.255%	87.833	-39.979	-49.479%	59.495	-9.915	-14.255%	87.833
defocus_blur	3	62.13	-3.17	-4.855%	49.609	-3.881	-6.174%	68.595	-0.815	-1.174%	102.018	-10.423	-9.078%	68.595	-0.815	-1.174%	102.018
motion_blur	3	61.22	-4.08	-6.248%	49.92	-3.47	-6.498%	68.494	-0.916	-1.319%	100.004	-14.466	-12.638%	68.494	-0.916	-1.319%	100.004
glass_blur	3	58.46	-6.84	-10.475%	49.923	-4.647	-12.112%	66.915	-2.495	-3.944%	86.333	-23.137	-24.580%	66.915	-2.495	-3.944%	86.333
gaussian_blur	3	62.75	-2.55	-3.605%	50.342	-3.048	-5.799%	68.399	-0.839	-1.046%	104.781	-9.689	-8.464%	68.399	-0.839	-1.046%	104.781
jpeg_compression	3	61.91	-1.39	-2.129%	51.002	-2.388	-4.173%	69.284	-0.126	-0.182%	109.071	-3.479	-3.057%	69.284	-0.126	-0.182%	109.071
contrast	3	61.82	-3.48	-5.329%	49.364	-4.026	-7.541%	67.92	-1.49	-2.146%	105.282	-9.188	-8.027%	67.92	-1.49	-2.146%	105.282
elastic_transform	3	63.33	-0.97	-1.301%	50.517	-2.873	-5.326%	69.399	0.017	0.025%	107.57	-6.026	-5.249%	69.399	0.017	0.025%	107.57
pixelate	3	58.99	-6.31	-9.663%	47.548	-5.982	-11.303%	64.518	-4.892	-7.047%	90.252	-21.141	-21.156%	64.518	-4.892	-7.047%	90.252
snow	3	55.69	-6.61	-14.177%	45.588	-7.802	-14.614%	64.203	-2.307	-3.502%	80.391	-33.879	-29.979%	64.203	-2.307	-3.502%	80.391
frost	3	56.14	-6.16	-14.028%	45.548	-7.842	-14.608%	64.877	-4.533	-6.506%	81.478	-32.902	-28.822%	64.877	-4.533	-6.506%	81.478
fog	3	61.34	-3.74	-5.772%	49.277	-4.113	-7.705%	68.695	-0.705	-1.030%	105.695	-8.773	-7.666%	68.695	-0.705	-1.030%	105.695
brightness	3	63.57	-1.73	-2.649%	50.557	-2.633	-5.203%	69.014	-0.204	-0.294%	111.594	-2.876	-2.513%	69.014	-0.204	-0.294%	111.594
saturation	3	60.94	-5.26	-8.055%	48.509	-4.821	-9.036%	67.062	-1.748	-2.519%	97.652	-16.818	-14.802%	67.062	-1.748	-2.519%	97.652
saturation	3	64.3	-1.0	-1.531%	50.572	-2.818	-5.277%	69.341	-0.069	-0.099%	112.884	-1.586	-1.385%	69.341	-0.069	-0.099%	112.884
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	4	58.4	-6.4	-8.901%	48.144	-4.646	-8.702%	67.322	-2.078	-2.994%	94.511	-19.459	-17.246%	67.322	-2.078	-2.994%	94.511
gaussian_noise	4	59.01	-6.29	-9.625%	48.768	-4.622	-8.658%	67.389	-2.021	-2.912%	95.8	-18.58	-16.231%	67.389	-2.021	-2.912%	95.8
shot_noise	4	58.67	-6.63	-10.153%	48.696	-4.694	-8.792%	65.222	-4.088	-6.044%	94.179	-20.2	-17.720%	65.222	-4.088	-6.044%	94.179
speckle_noise	4	58.78	-6.73	-10.269%	47.739	-5.099	-9.856%	64.885	-4.135	-5.912%	93.767	-20.729	-18.270%	64.885	-4.135	-5.912%	93.767
zoom_blur	4	51.14	-14.16	-21.625%	42.789	-10.601	-19.856%	58.043	-11.365	-16.374%	87.027	-41.418	-38.968%	58.043	-11.365	-16.374%	87.027
defocus_blur	4	61.7	-7.94	-11.747%	48.918	-4.906	-9.427%	68.494	-1.017	-1.459%	100.004	-14.466	-12.638%	68.494	-1.017	-1.459%	100.004
motion_blur	4	58.83	-6.47	-9.808%	48.346	-5.044	-9.447%	66.183	-3.277	-4.494%	87.833	-26.567	-23.200%	66.183	-3.277	-4.494%	87.833
glass_blur	4	58.46	-6.84	-10.335%	4												

Table 16: Performance and decrease ratio of Multiple LoRA on CLIP-BART against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	1	64.13	-1.31	-2.002%	52.011	-0.039	-0.074%	50.825	-0.495	-0.964%	112.679	-2.728	-2.363%				
gaussian_noise	1	64.17	-1.34	-2.048%	51.701	-0.149	-0.376%	50.466	-0.854	-1.663%	112.511	-2.896	-2.519%				
shot_noise	1	64.18	-1.26	-1.925%	51.808	-0.182	-0.349%	50.682	-0.638	-1.244%	112.6	-2.807	-2.432%				
zoom_blur	1	58.12	-7.32	-11.166%	48.958	-3.092	-5.939%	50.251	-2.083	-4.091%	86.91	-28.497	-24.693%				
defocus_blur	1	64.35	-1.09	-1.666%	51.98	-0.07	-0.135%	50.911	-0.409	-0.796%	112.609	-2.738	-2.371%				
motion_blur	1	64.09	-1.35	-2.063%	51.646	-0.04	-0.077%	51.816	0.496	0.966%	109.514	-5.893	-5.036%				
glass_blur	1	64.13	-1.31	-2.002%	51.596	-0.054	-0.105%	51.055	-0.265	-0.515%	111.24	-4.167	-3.619%				
jpeg_compression	1	65.0	-0.44	-0.672%	51.884	-0.166	-0.318%	51.73	0.1	0.798%	114.48	-0.927	-0.803%				
contrast	1	64.89	-0.55	-0.840%	51.829	-0.221	-0.425%	50.825	-0.495	-0.964%	113.78	-1.627	-1.410%				
elastic_transform	1	63.63	-1.81	-2.766%	52.147	0.097	0.186%	51.184	-0.138	-0.265%	112.217	-3.19	-2.764%				
pixelate	1	61.31	-4.13	-6.511%	50.62	-1.43	-2.747%	50.868	-0.452	-0.880%	105.256	-12.151	-10.529%				
snow	1	60.23	-5.21	-7.981%	49.467	-2.583	-4.962%	51.012	-0.308	-0.600%	101.293	-14.114	-12.114%				
frost	1	61.65	-3.79	-5.972%	50.453	-1.597	-3.068%	50.466	-0.454	-0.893%	103.413	-11.994	-10.393%				
fog	1	63.41	-2.03	-3.102%	51.932	-0.118	-0.227%	51.069	-0.251	-0.488%	111.682	-3.725	-3.228%				
brightness	1	65.16	-0.28	-0.428%	51.956	-0.094	-0.181%	51.471	0.151	0.295%	114.623	-0.734	-0.679%				
saturation	1	64.6	-0.84	-1.248%	51.087	-0.137	-0.267%	51.457	0.137	0.267%	113.683	-1.724	-1.494%				
severe	1	62.62	-2.82	-4.398%	50.843	-1.207	-2.319%	51.087	-0.367	-0.714%	111.041	-4.566	-3.783%				
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	2	62.37	-3.07	-4.691%	51.065	-0.983	-1.892%	50.825	-0.495	-0.964%	107.512	-7.895	-6.814%				
gaussian_noise	2	63.41	-2.0	-3.099%	50.986	-1.064	-2.044%	50.524	-0.796	-1.551%	110.583	-4.824	-4.180%				
shot_noise	2	63.01	-2.43	-3.713%	51.081	-0.969	-1.861%	50.28	-1.04	-2.027%	109.63	-5.777	-5.006%				
speckle_noise	2	63.76	-1.68	-2.567%	51.526	-0.524	-1.006%	50.825	-0.495	-0.964%	110.645	-4.762	-4.126%				
zoom_blur	2	55.11	-10.13	-15.785%	47.591	-4.459	-8.567%	50.309	-1.011	-1.971%	71.368	-44.019	-38.160%				
defocus_blur	2	63.52	-1.92	-2.934%	51.479	-0.571	-0.997%	51.27	-0.05	-0.097%	109.676	-5.731	-4.966%				
motion_blur	2	63.03	-2.41	-3.683%	51.685	-0.365	-0.709%	50.682	-0.638	-1.244%	107.584	-7.843	-6.796%				
glass_blur	2	62.88	-2.56	-3.912%	51.288	-0.762	-1.464%	50.696	-0.454	-0.893%	108.29	-7.117	-6.167%				
jpeg_compression	2	63.93	-1.51	-2.307%	51.924	-0.126	-0.242%	51.417	0.1	0.293%	111.553	-3.854	-3.339%				
contrast	2	64.62	-0.82	-1.251%	51.765	-0.285	-0.548%	50.825	-0.495	-0.964%	113.249	-2.688	-2.373%				
elastic_transform	2	63.1	-2.34	-3.576%	51.643	-0.587	-1.126%	51.27	-0.05	-0.097%	109.559	-5.848	-5.076%				
pixelate	2	61.44	-4.0	-6.121%	50.501	-1.549	-2.976%	50.438	-0.882	-1.719%	109.119	-7.288	-6.480%				
snow	2	57.53	-7.91	-12.087%	47.567	-4.483	-8.613%	50.122	-1.198	-2.344%	89.818	-25.589	-22.197%				
frost	2	57.92	-7.52	-11.491%	48.799	-3.251	-6.249%	49.864	-1.426	-2.784%	91.825	-23.582	-20.439%				
fog	2	62.57	-2.87	-4.386%	51.805	-0.345	-0.671%	50.438	-0.882	-1.719%	110.05	-5.87	-4.642%				
brightness	2	64.31	-1.13	-1.727%	51.733	-0.137	-0.266%	51.213	-0.107	-0.209%	114.086	-1.321	-1.144%				
saturation	2	61.08	-3.6	-5.663%	50.676	-1.374	-2.649%	51.227	-0.903	-1.811%	105.197	-10.21	-8.847%				
severe	2	59.48	-5.96	-9.108%	49.928	-2.122	-4.076%	50.797	-0.523	-1.020%	105.874	-8.333	-7.260%				
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	3	61.46	-3.98	-6.082%	50.622	-1.398	-2.686%	50.093	-1.227	-2.390%	104.261	-11.146	-9.658%				
gaussian_noise	3	61.81	-3.63	-5.547%	50.533	-1.517	-2.915%	50.083	-0.337	-0.656%	105.488	-9.919	-8.594%				
shot_noise	3	61.63	-3.81	-5.725%	50.286	-1.764	-3.389%	50.538	-0.722	-1.523%	105.624	-9.383	-8.783				
speckle_noise	3	61.7	-3.74	-5.715%	50.738	-1.342	-2.579%	50.194	-1.136	-2.195%	105.824	-9.473	-8.386%				
zoom_blur	3	52.57	-12.87	-19.667%	46.247	-5.803	-11.487%	49.993	-1.327	-2.586%	72.065	-38.342	-30.553%				
defocus_blur	3	61.82	-3.62	-5.532%	50.739	-1.311	-2.538%	50.538	-0.782	-1.523%	105.825	-9.562	-8.085%				
motion_blur	3	61.12	-4.32	-6.601%	50.662	-1.809	-3.409%	50.266	-1.054	-2.055%	101.771	-13.636	-11.141%				
glass_blur	3	58.06	-7.38	-11.278%	48.561	-3.489	-6.703%	49.997	-1.943%	-3.767%	87.648	-27.539	-24.615%				
jpeg_compression	3	62.43	-3.01	-4.469%	50.835	-1.215	-2.335%	50.725	-0.595	-1.160%	105.699	-7.798	-6.849%				
contrast	3	63.92	-1.52	-2.313%	51.638	-0.412	-0.792%	50.687	-0.513	-1.020%	112.413	-2.994	-2.594%				
elastic_transform	3	61.78	-3.66	-5.593%	50.517	-1.533	-2.946%	50.208	-1.122	-2.167%	105.746	-9.661	-8.371%				
pixelate	3	63.16	-2.26	-3.444%	51.701	-0.349	-0.670%	50.868	-0.452	-0.880%	107.495	-6.121	-5.351%				
snow	3	58.85	-4.59	-7.070%	49.784	-3.266	-6.276%	50.696	-0.624	-1.216%	99.566	-24.841	-21.524%				
frost	3	55.67	-9.77	-14.909%	47.233	-4.617	-9.254%	50.71	-0.61	-1.188%	81.666	-33.801	-29.299%				
fog	3	55.75	-9.69	-14.807%	47.106	-4.944	-9.480%	50.84	-0.38	-0.760%	82.465	-34.842	-28.544%				
brightness	3	61.33	-4.11	-6.211%	51.28	-0.77	-1.479%	50.969	-0.351	-0.684%	106.473	-8.777	-7.518%				
saturation	3	59.41	-6.03	-9.215%	49.555	-2.495	-4.784%	50.524	-0.796	-1.551%	98.194	-17.213	-14.915%				
severe	3	64.37	-1.07	-1.635%	51.566	-0.484	-0.929%	51.342	0.022	0.043%	113.511	-1.896	-1.643%				
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	
impulse_noise	4	58.89	-6.45	-9.856%	49.825	-2.225	-4.275%	49.519	-1.801	-3.509%	95.12	-20.267	-17.759%				
gaussian_noise	4	59.3	-6.14	-9.383%	49.269	-2.781	-5.344%	49.663	-1.657	-3.229%	96.911	-18.496	-16.027%				
shot_noise	4	58.57	-6.87	-9.808%	49.738	-2.312	-4.443%	49.677	-1.643	-3.201%	95.136	-20.271	-17.565%				
speckle_noise	4	60.4	-5.04	-7.702%	49.849	-2.201	-4.229%	49.702	-1.528	-2.978%	102.055	-13.352	-11.569%				
zoom_blur	4	50.6	-14.84	-22.677%	44.975	-7.073	-13.920%	50.38	-0.782	-1.523%	47.061	-36.346	-29.222%				
defocus_blur	4	60.02	-5.42	-8.252%	50.246	-1.404	-2.609%	49.534	-1.786	-3.481%	95.916	-19.491	-16.899%				
motion_blur	4	58.55	-6.89	-10.529%	49.602	-2.448	-4.702%	49.734	-1.566	-3.060%	90.346	-26.861	-22.582%				
glass_blur	4	56.15	-9.29	-14.196%	46.939	-5.111	-9.419%	49.677	-1.643	-3.201%	86.448	-34.959	-30.292%				
jpeg_compression	4	60.7	-4.74	-7.243%	50.398	-1.652	-3.175%	49.72	-1.6	-3.117%	98.408	-16.999	-14.730%				
contrast	4	57.84	-7.6	-11.614%	49.221	-2.829	-5.432%	49.872	-1.6	-3.113%	92.806	-22.001	-19.584%				
elastic_transform	4	60.89	-4.55	-6.935%	50.072	-1.796	-3.801%	50.323	-0.499	-0.949%	99.963	-16.344	-14.162%				
pixelate	4	56.18	-4.65	-7.699%	49.079	-3.979	-7.914%	50.84	-0.66	-1.290%	76.262	-34.993	-29.672%				
snow	4	52.46	-12.98	-19.835%	45.126	-6.924	-13.302%	49.132	-2.188	-4.266%	65.907	-49.3	-43.151%				
frost	4	51.28	-15.51	-23.909%	44.512	-7.595	-13.045%	49.182	-2.188	-4.266%	65.907	-49.3	-43.151%				
fog	4	60.03	-5.41	-8.257%	50.249	-1.621	-3.134%	50.38	-0.94	-1.831%	91.034	-18.013	-15.124%				
brightness	4	62.7	-2.74	-4.187%	51.713	-0.738	-1.418%	50.491	-0.491	-0.775%	109.223	-3.186	-2.568%				
saturation	4	57.84	-7.64	-11.614%	49.221	-2.829	-5.432%	49.872	-1.6	-3.113%	92.806	-22.001	-19.584%				
severe	4	61.37	-4.07	-6.219%	49.698	-2.352	-4.519%	50.594	-0.366	-0.712%	107.209	-6.188	-5.016%				

Table 18: Performance and decrease ratio of Single Adapter on CLIP-BART against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Caption
impulse_noise	1	64.03	-1.32	-2.020%	48.863	-5.277	-9.747%	73.26	-0.63	-0.853%	100.007	-6.031	-5.242%				
impulse_noise	2	64.88	-0.47	-0.719%	49.873	-4.267	-7.862%	73.446	-0.444	-0.601%	111.841	-3.197	-2.779%				
shot_noise	1	64.99	-0.36	-0.551%	49.785	-4.355	-8.843%	73.380	-0.501	-0.678%	112.405	-3.633	-2.288%				
speckle_noise	1	64.87	-0.48	-0.735%	49.587	-4.553	-8.140%	73.618	-0.272	-0.367%	112.232	-2.806	-2.489%				
zoom blur	1	59.13	-6.22	-9.518%	46.621	-7.519	-13.088%	57.155	-16.73	-22.648%	87.848	-27.19	-23.635%				
defocus blur	1	64.89	-0.46	-0.704%	49.833	-4.307	-7.955%	73.963	-0.073	-0.099%	111.714	-3.864	-3.358%				
motion blur	1	64.82	-0.53	-0.811%	49.928	-4.212	-7.779%	73.662	-0.228	-0.309%	110.564	-4.474	-3.889%				
glass blur	1	65.07	-0.28	-0.429%	50.016	-4.124	-7.717%	74.012	-0.13	-0.176%	112.163	-2.855	-2.482%				
gaussian blur	1	65.62	-0.27	-0.408%	50.358	-3.986	-7.436%	74.236	-0.346	-0.468%	113.593	-4.453	-3.256%				
jpeg_compression	1	65.26	-0.09	-0.138%	49.809	-4.331	-7.999%	73.719	-0.171	-0.231%	112.803	-2.235	-1.943%				
contrast	1	64.4	-0.75	-1.148%	49.777	-4.363	-8.808%	73.618	-0.272	-0.367%	111.844	-3.054	-2.655%				
elastic_transform	1	65.19	-0.16	-0.245%	49.817	-4.323	-7.985%	73.59	-0.3	-0.406%	112.36	-2.678	-2.328%				
pixelate	1	65.29	-2.86	-4.259%	48.474	-5.666	-10.466%	72.599	-2.291	-3.147%	102.729	-12.309	-10.709%				
snow	1	61.59	-3.76	-5.754%	41.802	-6.12	-11.303%	72.968	-1.822	-2.465%	100.139	-14.999	-12.951%				
frost	1	62.54	-2.81	-4.300%	48.434	-5.706	-10.540%	72.571	-1.319	-1.786%	102.648	-12.154	-10.565%				
fog	1	64.01	-1.34	-2.089%	49.936	-4.204	-7.764%	73.561	-0.329	-0.445%	110.89	-4.148	-3.669%				
brightness	1	65.53	0.18	0.275%	50.135	-4.005	-7.979%	74.48	0.59	0.798%	113.34	-1.698	-1.476%				
saturation	1	65.35	0.0	0.000%	49.889	-4.251	-7.852%	74.078	0.188	0.254%	112.884	-2.154	-1.872%				
saturation	1	65.25	-2.1	-3.213%	49.316	-4.824	-8.910%	72.838	-1.032	-1.397%	111.428	-3.64	-3.138%				
impulse_noise	2	63.07	-2.28	-3.489%	48.617	-5.523	-10.202%	72.025	-1.865	-2.524%	105.627	-9.411	-8.114%				
gaussian_noise	2	64.02	-1.33	-2.015%	49.428	-4.712	-8.704%	73.303	-0.587	-0.795%	109.739	-5.299	-4.696%				
shot_noise	2	63.9	-1.45	-2.219%	49.189	-4.951	-9.145%	72.858	-1.027	-1.397%	109.772	-5.266	-4.578%				
speckle_noise	2	64.33	-1.02	-1.561%	49.332	-4.808	-8.889%	73.044	-0.846	-1.144%	109.739	-5.319	-3.754%				
zoom blur	2	55.8	-8.55	-14.614%	44.451	-9.089	-17.807%	62.523	-11.367	-18.383%	73.635	-41.403	-35.900%				
defocus blur	2	64.4	-0.95	-1.454%	49.563	-4.577	-8.455%	73.432	-0.488	-0.620%	109.739	-4.702	-4.956%				
motion blur	2	63.69	-1.86	-2.540%	49.626	-4.514	-8.337%	73.69	-0.2	-0.270%	107.598	-7.74	-6.486%				
glass blur	2	63.67	-1.68	-2.571%	49.269	-4.871	-8.980%	72.8	-1.199	-1.475%	108.603	-6.975	-6.003%				
gaussian blur	2	64.58	-0.77	-1.178%	49.602	-4.538	-8.751%	74.121	0.231	0.312%	103.583	-4.685	-4.072%				
jpeg_compression	2	65.29	-0.08	-0.120%	49.65	-4.49	-8.203%	73.217	-0.057	-0.076%	111.506	-3.532	-3.071%				
contrast	2	64.14	-1.21	-1.852%	49.261	-4.879	-9.913%	73.561	-0.329	-0.445%	109.739	-5.423	-4.704%				
elastic_transform	2	62.25	-3.1	-4.744%	49.078	-5.062	-9.350%	70.289	-3.601	-4.876%	97.775	-17.263	-15.006%				
pixelate	2	61.77	-3.58	-5.478%	47.625	-6.215	-11.480%	69.169	-4.723	-6.389%	100.404	-14.634	-12.715%				
snow	2	58.49	-6.86	-10.497%	46.025	-8.115	-14.889%	69.413	-4.477	-6.059%	88.321	-26.717	-23.225%				
frost	2	59.16	-6.19	-9.472%	46.028	-8.06	-14.880%	70.676	-3.214	-3.590%	91.341	-23.607	-20.999%				
fog	2	63.35	-2.0	-3.069%	49.412	-4.72	-8.719%	69.729	-1.461	-1.652%	105.677	-5.461	-4.374%				
brightness	2	64.98	-0.37	-0.566%	50.024	-4.117	-7.603%	73.682	-0.028	-0.037%	113.192	-1.846	-1.604%				
saturation	2	61.39	-3.46	-5.250%	48.577	-5.363	-10.275%	71.824	-2.66	-2.766%	104.608	-10.192	-8.595%				
saturation	2	60.07	-5.28	-8.009%	48.243	-5.897	-10.902%	70.848	-3.042	-4.117%	105.558	-9.48	-8.241%				
impulse_noise	3	62.07	-3.28	-5.019%	48.481	-5.669	-10.452%	71.236	-2.654	-3.592%	102.932	-12.06	-10.245%				
gaussian_noise	3	62.18	-3.17	-4.851%	48.64	-5.5	-10.188%	71.308	-2.582	-3.495%	104.363	-10.652	-9.282%				
shot_noise	3	62.39	-2.96	-4.529%	48.672	-5.468	-10.099%	71.451	-2.439	-3.301%	104.749	-10.289	-8.944%				
speckle_noise	3	62.82	-2.53	-3.871%	48.716	-5.42	-10.011%	72.04	-1.85	-2.504%	104.894	-10.144	-8.183%				
zoom blur	3	53.05	-12.3	-18.822%	42.996	-11.44	-20.584%	59.538	-14.252	-19.424%	58.63	-36.4	-49.028%				
defocus blur	3	62.42	-2.93	-4.484%	48.394	-5.748	-10.613%	71.839	-2.051	-2.776%	102.174	-12.324	-10.713%				
motion blur	3	61.83	-3.52	-5.366%	48.792	-5.48	-9.799%	71.465	-2.425	-3.281%	101.501	-13.537	-11.767%				
glass blur	3	58.84	-6.51	-9.962%	45.937	-6.203	-15.515%	69.786	-4.104	-5.544%	86.792	-28.264	-24.535%				
gaussian blur	3	63.26	-2.09	-3.198%	48.704	-5.436	-10.400%	71.853	-2.077	-2.757%	105.781	-9.257	-8.047%				
jpeg_compression	3	64.83	-0.52	-0.766%	50.072	-4.068	-7.719%	73.145	-0.282	-0.382%	111.081	-3.957	-3.449%				
contrast	3	62.8	-2.55	-3.902%	48.537	-5.603	-10.349%	71.896	-1.994	-2.698%	105.345	-9.693	-8.420%				
elastic_transform	3	61.34	-4.01	-6.136%	48.444	-6.1	-11.777%	67.891	-5.999	-8.118%	81.86	-33.178	-28.841%				
pixelate	3	59.32	-6.03	-9.227%	46.128	-8.012	-14.790%	67.475	-4.415	-6.081%	90.111	-24.927	-21.669%				
snow	3	56.63	-8.72	-13.344%	44.069	-10.071	-18.602%	67.332	-6.558	-8.876%	81.487	-33.551	-29.658%				
frost	3	56.55	-8.1	-13.466%	44.84	-9.3	-17.177%	67.891	-5.999	-8.118%	81.86	-33.178	-28.841%				
fog	3	62.17	-3.18	-4.866%	48.855	-5.285	-9.719%	72.169	-1.721	-2.329%	106.133	-8.905	-7.741%				
brightness	3	64.41	-0.94	-1.438%	49.459	-4.681	-8.849%	72.872	-1.038	-1.378%	112.042	-2.996	-2.605%				
saturation	3	60.9	-4.65	-6.809%	48.1	-6.04	-11.157%	70.992	-2.898	-3.922%	97.078	-17.76	-15.412%				
saturation	3	65.18	-0.17	-0.260%	49.897	-4.243	-7.838%	74.178	-0.289	-0.392%	112.3	-2.78	-2.389%				
impulse_noise	4	59.67	-5.88	-8.692%	47.615	-6.225	-12.026%	69.887	-4.003	-5.418%	93.318	-21.72	-18.818%				
gaussian_noise	4	59.55	-5.8	-8.875%	47.702	-6.438	-11.891%	69.872	-4.037	-5.473%	95.42	-19.618	-17.654%				
shot_noise	4	59.34	-6.01	-9.191%	47.424	-6.716	-12.409%	69.155	-4.735	-6.409%	93.488	-21.55	-18.733%				
speckle_noise	4	61.44	-3.91	-5.580%	47.822	-6.318	-11.470%	70.877	-3.033	-4.078%	100.803	-14.235	-12.734%				
zoom blur	4	50.89	-14.46	-22.127%	41.207	-12.023	-23.888%	57.471	-16.419	-22.213%	61.471	-66.327	-57.657%				
defocus blur	4	60.54	-4.81	-7.369%	47.472	-6.668	-12.317%	69.915	-3.975	-5.379%	95.936	-20.372	-18.609%				
motion blur	4	59.18	-6.17	-9.441%	47.433	-6.31	-11.666%	68.681	-2.509	-3.448%	84.465	-36.753	-22.099%				
glass blur	4	56.77	-8.58	-13.129%	44.57	-9.57	-17.077%	68.624	-5.266	-7.127%	79.597	-35.441	-30.808%				
gaussian blur	4	61.34	-4.01	-6.136%	47.957	-6.183	-11.421%	70.504	-3.386	-4.583%	97.671	-17.027	-15.097%				
jpeg_compression	4	63.54	-1.81	-2.799%	49.118	-5.022	-9.777%	72.312	-1.578	-2.135%	107.37	-6.68	-6.066%				
contrast	4	58.46	-6.89	-10.543%	45.929	-8.211	-15.165%	68.724	-5.166	-6.991%	91.302	-23.736	-20.633%				
elastic_transform	4	51.62	-7.33	-11.308%	47.909	-6.231	-11.509%	72.011	-1.879	-2.543%	98.985	-16.083	-13.912%				
pixelate	4	55.01	-10.34	-15.822%	44.284	-9.856	-18.205%	64.274	-8.616	-13.013%	73.703	-41.43	-35.932%				
snow	4	53.43	-11.92	-18.249%	42.439	-11.701	-21.612%	64.877	-8.013	-12.197%	65.466	-49.576	-43.066%				
frost	4	55.81	-9.54	-14.280%	44.252	-9.888	-18.264%	68.222	-5.568	-7.755%	78.427	-36.611	-31.845%				
fog	4	60.53	-4.82	-7.376%	47.647	-6.493	-11.994%	71.451	-2.439	-3.301%	100.557	-14.481	-12.588%				
brightness	4	63.6	-1.75	-2.670%	49.213	-4.927	-9.101%	71.824	-2.66	-2.766%	104.608	-10.194	-8.595%				
saturation	4	58.17	-7.18	-10.897%	45.492	-8.648	-15.973%	68.408	-5.482	-7.419%	87.985	-27.053	-23.571%				
saturation	4	62.08	-3.27	-5.004%	48.092	-6.048	-11.171%	71.623	-2.607	-3.069%	107.11	-7.828					

Table 20: Performance and decrease ratio of Single LoRA on CLIP-BART against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CCider	Decrease	Decrease Ratio	CCider
impulse_noise	1	63.06	-2.28	-3.489%	48.036	-5.154	-9.689%	72.614	0.966	-1.313%	108.958	-5.585	-4.876%				
gaussian_noise	1	63.9	-1.44	-2.264%	48.179	-5.011	-9.420%	72.987	0.933	-0.806%	111.859	-2.684	-2.413%				
shot_noise	1	63.9	-1.44	-2.264%	48.172	-4.778	-8.199%	72.913	0.65	-0.846%	111.861	-3.111	-2.777%				
speckle_noise	1	63.97	-1.37	-2.097%	48.752	-4.438	-8.144%	72.671	0.901	-1.235%	111.538	-3.003	-2.623%				
zoom blur	1	58.19	-7.15	-10.405%	45.429	-7.761	-14.592%	58.677	-4.807	-8.255%	106.603	-2.794	-2.439%				
defocus blur	1	64.07	-1.27	-1.944%	48.879	-4.311	-8.105%	73.016	0.566	-0.767%	111.315	-3.228	-2.819%				
motion blur	1	64.08	-1.26	-1.928%	48.283	-4.907	-9.226%	73.403	0.177	-0.240%	109.757	-4.786	-4.178%				
glass blur	1	64.07	-1.27	-1.944%	48.736	-4.454	-8.734%	73.317	0.263	-0.357%	107.711	-3.832	-3.345%				
gaussian blur	1	64.78	-0.56	-0.837%	49.173	-4.017	-7.552%	73.791	0.214	-0.286%	113.37	-1.773	-1.024%				
jpeg_compression	1	64.52	-0.82	-1.255%	49.086	-4.104	-7.716%	72.872	0.708	-0.962%	112.734	-1.809	-1.579%				
contrast	1	63.72	-1.62	-2.476%	48.696	-4.494	-8.449%	73.001	0.579	-0.787%	111.418	-3.125	-2.726%				
elastic_transform	1	64.49	-0.85	-1.301%	49.181	-4.009	-7.537%	72.442	-1.138	-1.547%	111.973	-2.57	-2.244%				
pixelate	1	61.14	-4.2	-6.428%	47.265	-5.925	-11.299%	71.48	-2.1	-2.854%	102.811	-11.732	-10.242%				
snow	1	60.18	-5.16	-7.897%	46.701	-6.489	-12.200%	71.236	-2.344	-3.186%	99.386	-15.157	-13.233%				
frost	1	61.78	-3.56	-5.448%	47.392	-5.798	-10.900%	71.379	-2.201	-2.991%	102.884	-11.679	-10.196%				
fog	1	63.1	-2.24	-3.428%	48.402	-4.788	-9.902%	72.398	-1.182	-1.606%	110.773	-3.77	-3.291%				
brightness	1	64.85	-0.49	-0.739%	49.237	-3.953	-7.432%	73.504	-0.706	-0.104%	113.739	-3.404	-2.702%				
saturation	1	64.45	-0.89	-1.248%	49.284	-3.906	-7.343%	73.432	0.148	-0.201%	112.657	-1.888	-1.647%				
saturation	1	62.13	-3.21	-4.913%	47.809	-5.281	-9.828%	71.609	-1.971	-2.676%	110.161	-4.382	-3.826%				
impulse_noise	2	62.1	-3.24	-4.999%	47.957	-5.233	-9.393%	72.025	-1.555	-2.113%	105.736	-8.807	-7.689%				
gaussian_noise	2	63.05	-2.29	-3.565%	47.758	-5.432	-10.212%	72.169	-1.411	-1.918%	109.572	-4.971	-4.349%				
shot_noise	2	62.87	-2.47	-3.780%	47.845	-5.345	-10.404%	72.47	-1.11	-1.508%	108.551	-5.992	-5.231%				
speckle_noise	2	63.46	-1.88	-2.877%	48.307	-4.883	-9.118%	72.298	-1.282	-1.742%	109.917	-6.426	-5.039%				
zoom blur	2	55.27	-10.07	-15.426%	44.196	-8.994	-16.809%	63.212	-10.368	-14.690%	73.46	-41.083	-35.807%				
defocus blur	2	63.3	-2.04	-3.122%	48.728	-4.462	-8.389%	72.236	-1.354	-1.840%	109.109	-5.434	-4.744%				
motion blur	2	62.82	-2.52	-3.857%	48.012	-5.178	-9.749%	72.8	-0.78	-1.006%	106.62	-7.923	-6.917%				
glass blur	2	62.8	-2.54	-3.807%	48.124	-5.066	-9.529%	72.235	-1.325	-1.801%	106.694	-7.849	-6.835%				
gaussian blur	2	63.67	-1.67	-2.556%	48.696	-4.494	-8.449%	73.001	0.579	-0.787%	109.757	-4.786	-4.178%				
jpeg_compression	2	64.41	-0.93	-1.423%	48.315	-4.875	-9.166%	71.63	-0.5	-0.747%	112.232	-2.291	-2.000%				
contrast	2	62.82	-2.52	-3.857%	48.171	-5.019	-9.455%	72.269	-1.311	-1.781%	109.635	-6.809	-5.481%				
elastic_transform	2	61.37	-3.97	-6.076%	46.748	-6.442	-12.111%	69.513	-4.067	-5.922%	96.564	-17.889	-15.658%				
pixelate	2	60.63	-4.31	-6.788%	47.104	-5.996	-11.274%	68.624	-4.956	-6.756%	99.434	-15.109	-13.193%				
snow	2	57.21	-8.13	-12.449%	44.737	-8.453	-15.992%	68.667	-4.913	-6.678%	88.402	-26.141	-22.822%				
frost	2	58.26	-6.98	-10.618%	45.691	-7.499	-14.449%	69.513	-4.067	-5.922%	91.338	-23.205	-20.299%				
fog	2	62.28	-3.06	-4.603%	48.37	-4.42	-9.062%	68.207	-3.373	-4.902%	108.513	-6.03	-5.264%				
brightness	2	64.28	-1.06	-1.622%	49.197	-3.993	-7.975%	72.815	-0.765	-1.040%	113.05	-1.493	-1.303%				
saturation	2	61.17	-4.17	-6.362%	47.611	-5.519	-10.377%	71.552	-2.028	-2.757%	103.588	-10.955	-9.564%				
saturation	2	59.18	-6.16	-9.428%	47.337	-5.853	-11.005%	69.838	-3.722	-5.059%	85.23	-27.73	-74.15%				
impulse_noise	3	61.21	-4.13	-6.521%	47.384	-5.806	-10.915%	70.934	-2.446	-3.396%	103.522	-11.021	-9.621%				
gaussian_noise	3	61.36	-3.98	-6.091%	47.376	-5.814	-10.930%	71.422	-1.198	-1.632%	103.362	-11.181	-9.762%				
shot_noise	3	61.22	-4.12	-6.305%	47.71	-5.48	-10.302%	70.863	-2.177	-2.693%	103.026	-10.617	-9.269%				
speckle_noise	3	61.7	-3.64	-5.537%	47.615	-5.575	-10.482%	70.996	-2.474	-3.655%	104.299	-10.144	-8.569%				
zoom blur	3	52.63	-12.71	-19.423%	41.827	-11.363	-21.363%	60.729	-12.851	-17.465%	60.226	-31.371	-47.421%				
defocus blur	3	61.45	-3.89	-5.551%	47.36	-5.833	-10.960%	70.604	-2.976	-4.044%	104.541	-11.599	-10.126%				
motion blur	3	60.76	-4.58	-7.009%	47.297	-5.803	-11.079%	71.379	-2.0	-2.901%	99.834	-14.719	-12.849%				
glass blur	3	58.01	-7.33	-11.418%	44.959	-8.231	-15.474%	69.585	-5.959	-8.429%	87.965	-26.638	-23.256%				
gaussian blur	3	62.34	-3.1	-4.749%	48.004	-5.186	-9.699%	71.322	-2.258	-2.990%	104.888	-9.655	-8.429%				
jpeg_compression	3	63.99	-1.35	-2.066%	48.45	-4.74	-8.125%	72.663	-0.835	-1.114%	112.324	-2.902%	-2.531%				
contrast	3	61.6	-3.74	-5.724%	47.591	-5.599	-10.526%	70.834	-2.746	-3.732%	105.37	-8.973	-7.834%				
elastic_transform	3	59.7	-5.64	-8.632%	46.351	-6.939	-12.838%	70.877	-2.703	-3.674%	98.86	-17.657	-15.415%				
pixelate	3	58.51	-6.83	-10.453%	45.238	-7.079	-15.476%	68.844	-6.756	-9.155%	90.965	-24.478	-21.730%				
snow	3	55.74	-9.46	-14.026%	43.711	-9.479	-17.821%	68.844	-6.756	-9.155%	80.878	-33.965	-29.625%				
frost	3	55.94	-9.4	-13.860%	44.2	-9.47	-16.646%	67.705	-8.879	-7.985%	82.34	-32.219	-28.126%				
fog	3	61.13	-4.21	-6.443%	48.132	-5.058	-10.491%	71.552	-2.028	-2.757%	105.833	-6.49	-5.786%				
brightness	3	63.51	-1.83	-2.601%	48.338	-4.852	-9.212%	72.528	-1.052	-1.430%	110.715	-3.828	-3.342%				
saturation	3	59.7	-5.64	-8.632%	46.351	-6.939	-12.838%	70.877	-2.703	-3.674%	98.86	-17.657	-15.415%				
saturation	3	64.25	-1.09	-1.668%	48.259	-4.931	-9.721%	72.815	-0.765	-1.040%	111.852	-3.601	-2.232%				
impulse_noise	4	58.42	-6.92	-10.591%	46.327	-6.163	-12.903%	68.982	-4.599	-6.249%	93.878	-20.665	-18.041%				
gaussian_noise	4	59.01	-6.33	-9.688%	46.661	-6.529	-12.275%	69.054	-4.526	-6.151%	96.203	-18.34	-16.011%				
shot_noise	4	58.29	-7.05	-10.709%	46.311	-6.879	-12.933%	69.054	-4.526	-6.151%	96.592	-19.951	-17.418%				
speckle_noise	4	59.51	-5.41	-8.260%	47.133	-5.977	-11.050%	70.044	-3.536	-4.805%	101.619	-12.924	-11.283%				
zoom blur	4	50.88	-14.46	-22.110%	40.38	-12.81	-24.083%	59.15	-14.43	-19.611%	50.257	-64.286	-56.124%				
defocus blur	4	59.47	-5.87	-8.904%	46.637	-6.553	-12.200%	69.829	-3.751	-5.098%	95.621	-18.822	-16.230%				
motion blur	4	58.1	-7.04	-10.745%	46.343	-6.847	-12.873%	68.624	-4.596	-6.736%	97.717	-26.826	-23.242%				
glass blur	4	56.24	-9.1	-13.927%	43.743	-9.447	-17.761%	68.824	-4.756	-6.463%	79.445	-35.098	-30.642%				
gaussian blur	4	60.26	-5.08	-7.775%	46.78	-6.41	-12.951%	70.576	-3.008	-4.083%	97.958	-16.585	-14.480%				
jpeg_compression	4	62.55	-2.79	-4.270%	47.726	-5.464	-10.722%	71.738	-1.842	-2.503%	106.674	-7.869	-6.870%				
contrast	4	57.65	-7.69	-11.769%	45.611	-7.579	-14.248%	67.805	-5.775	-7.848%	92.032	-22.511	-19.635%				
elastic_transform	4	56.98	-8.39	-13.709%	46.828	-6.627	-11.961%	71.353	-1.827	-2.455%	99.181	-15.362	-13.112%				
pixelate	4	54.42	-10.92	-16.173%	42.19	-10.29	-19.345%	63.585	-9.978	-13.583%	73.761	-40.78	-35.604%				
snow	4	52.74	-12.6	-19.824%	42.137	-11.153	-20.789%	64.504	-12.336	-16.685%	65.68	-48.36	-42.675%				
frost	4	55.16	-10.13	-15.906%	43.799	-9.391	-17.609%	67.59	-8.599	-11.414%	79.57	-34.76	-30.730%				
fog	4	59.66	-5.68	-8.693%	46.732	-6.458	-12.141%	70.777	-2.803	-3.810%	101.097	-14.446	-11.739%				
brightness	4	62.69	-2.65	-4.056%	48.06	-5.13	-9.445%	71.566	-2.014	-2.757%	108.255	-6.308	-5.507%				
saturation	4	57.6	-7.74	-11.846%	44.884	-8.326	-15.653%	69.025	-4.555	-6.190%	87.767	-26.776	-23.730%				
saturation	4	61.19	-4.15	-6.362%	47.154	-6.063	-11.348%	71.035	-2.448	-3.459%	106.8						

Table 22: Performance and decrease ratio of Full Finetuning on CLIP-T5 against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		CIDder	Decrease	Decrease Ratio	
impulse_noise	1	64.29	-2.0	-3.017%	54.19	-2.63	-4.629%	72.413	-1.647	-2.224%	106.157	-5.343	-4.791%				
gaussian_noise	1	65.01	-1.28	-1.931%	54.699	-1.851	-3.258%	73.159	-0.901	-1.216%	105.984	-2.446	-2.194%				
shot_noise	1	65.07	-1.22	-1.840%	54.818	-2.002	-3.524%	73.188	-0.872	-1.176%	105.872	-2.522	-2.261%				
speckle_noise	1	65.05	-1.24	-1.871%	54.556	-2.264	-3.985%	73.044	-1.016	-1.371%	109.044	-2.586	-2.302%				
zoom blur	1	59.3	-6.99	-10.545%	51.375	-5.445	-9.952%	59.035	-1.025	-2.037%	84.714	-26.766	-24.023%				
defocus blur	1	65.01	-1.28	-1.931%	54.858	-1.962	-3.454%	73.475	-0.985	-1.360%	106.776	-2.724	-2.441%				
motion blur	1	64.56	-1.73	-2.610%	54.961	-1.859	-3.272%	73.016	-1.044	-1.410%	107.017	-4.483	-4.021%				
glass blur	1	65.2	-1.09	-1.644%	54.937	-1.883	-3.114%	73.26	-0.8	-1.081%	106.969	-2.511	-2.270%				
gaussian blur	1	65.83	-0.46	-0.694%	55.541	-2.259	-3.722%	74.107	-0.047	-0.063%	111.538	-0.038	-0.035%				
jpeg_compression	1	65.66	-0.63	-0.959%	55.406	-1.414	-2.486%	73.776	-0.284	-0.383%	111.035	-0.465	-0.417%				
contrast	1	65.1	-1.19	-1.795%	55.629	-1.191	-2.069%	73.288	-0.772	-1.042%	108.6	-5.9	-5.400%				
elastic_transform	1	65.31	-0.78	-1.177%	55.125	-1.628	-2.866%	73.489	-0.571	-0.771%	109.643	-2.037	-1.827%				
pixelate	1	62.4	-3.39	-5.114%	54.349	-2.471	-4.499%	71.609	-2.451	-3.399%	105.583	-9.017	-8.791%				
snow	1	61.47	-4.82	-7.717%	51.693	-5.127	-9.022%	72.068	-1.992	-2.689%	98.17	-13.33	-11.955%				
frost	1	62.45	-3.54	-5.340%	53.609	-3.751	-6.602%	72.269	-1.791	-2.418%	101.346	-10.154	-9.107%				
fog	1	64.43	-1.86	-2.880%	55.152	-1.668	-2.936%	73.245	-0.815	-1.100%	108.736	-3.124	-2.802%				
brightness	1	65.96	-0.33	-0.498%	55.184	-1.636	-2.880%	74.293	-0.233	-0.315%	111.749	-0.249	-0.223%				
saturation	1	65.54	-0.75	-1.131%	55.541	-1.279	-2.259%	74.566	-0.506	-0.683%	110.042	-1.458	-1.308%				
saturation	1	65.3	-2.99	-4.550%	54.412	-2.408	-4.373%	72.7	-1.36	-1.837%	108.141	-3.359	-3.042%				
impulse_noise	2	63.53	-2.76	-4.144%	53.602	-3.218	-5.664%	72.499	-1.561	-2.108%	103.546	-7.954	-7.113%				
gaussian_noise	2	64.59	-1.17	-1.764%	54.055	-2.765	-4.897%	72.786	-1.274	-1.720%	107.088	-3.802	-3.419%				
shot_noise	2	64.2	-2.09	-3.153%	53.665	-3.155	-5.522%	72.169	-1.891	-2.554%	106.268	-5.232	-4.692%				
speckle_noise	2	64.79	-1.5	-2.261%	54.468	-2.352	-4.139%	72.872	-1.188	-1.604%	107.593	-5.397	-4.940%				
zoom blur	2	56.2	-10.09	-15.221%	49.054	-7.766	-15.688%	63.944	-10.116	-15.659%	70.213	-41.287	-37.029%				
defocus blur	2	64.44	-1.85	-2.791%	54.556	-2.264	-3.985%	72.757	-1.303	-1.759%	106.958	-4.542	-4.074%				
motion blur	2	63.79	-2.5	-3.771%	54.667	-2.153	-3.799%	72.312	-1.748	-2.369%	104.51	-6.99	-6.299%				
glass blur	2	64.06	-2.23	-3.364%	54.238	-2.582	-4.549%	71.997	-2.063	-2.786%	105.223	-6.277	-5.836%				
gaussian blur	2	64.69	-1.6	-2.414%	54.754	-2.066	-3.635%	73.059	-1.101	-1.352%	108.336	-3.164	-2.837%				
jpeg_compression	2	65.47	-0.82	-1.227%	55.144	-1.676	-2.959%	73.618	-0.422	-0.569%	109.517	-1.983	-1.778%				
contrast	2	64.42	-1.87	-2.821%	54.969	-1.851	-3.258%	72.772	-1.208	-1.740%	106.565	-4.933	-4.426%				
elastic_transform	2	62.35	-3.94	-5.944%	53.037	-3.783	-6.585%	69.757	-4.303	-8.510%	96.449	-15.051	-13.498%				
pixelate	2	61.99	-4.3	-6.407%	53.586	-3.234	-5.092%	69.688	-4.902	-6.740%	98.122	-13.178	-11.980%				
snow	2	61.41	-7.88	-11.887%	49.881	-6.939	-12.213%	69.284	-4.776	-6.449%	86.106	-25.394	-22.775%				
frost	2	59.26	-7.03	-10.608%	50.811	-6.409	-10.578%	69.628	-4.432	-5.984%	88.499	-23.001	-20.820%				
fog	2	63.61	-2.68	-4.043%	54.762	-2.108	-3.621%	69.155	-1.905	-2.624%	106.338	-5.162	-4.629%				
brightness	2	65.01	-1.28	-1.931%	55.033	-1.787	-3.146%	73.475	-0.585	-0.790%	110.275	-1.225	-1.099%				
saturation	2	62.33	-3.96	-5.974%	53.021	-3.799	-6.686%	72.513	-1.547	-2.088%	107.702	-10.798	-9.644%				
saturation	2	60.63	-5.66	-8.538%	53.435	-3.385	-5.958%	69.987	-4.073	-5.499%	102.159	-31.341	-28.777%				
impulse_noise	3	62.53	-3.76	-5.672%	53.236	-3.584	-6.580%	71.623	-2.477	-3.200%	101.179	-9.321	-9.257%				
gaussian_noise	3	63.01	-3.28	-4.803%	52.965	-3.855	-6.745%	71.666	-2.784	-3.232%	101.83	-9.47	-8.473%				
shot_noise	3	62.72	-3.57	-5.385%	53.379	-3.441	-6.065%	70.791	-3.269	-4.414%	101.888	-10.612	-8.621%				
speckle_noise	3	62.82	-3.47	-5.235%	53.49	-3.323	-5.809%	71.422	-2.638	-3.561%	101.97	-8.53	-8.457%				
zoom blur	3	53.66	-12.63	-19.053%	47.567	-9.253	-16.267%	61.174	-12.886	-17.399%	56.507	-50.993	-49.211%				
defocus blur	3	62.72	-3.57	-5.385%	53.538	-3.282	-5.776%	70.848	-3.122	-4.337%	100.814	-10.668	-9.584%				
motion blur	3	62.1	-4.19	-6.216%	53.339	-3.481	-6.126%	70.705	-3.355	-4.536%	98.675	-12.425	-11.502%				
glass blur	3	59.02	-7.27	-10.967%	50.493	-6.237	-11.115%	68.939	-5.121	-6.914%	85.732	-26.128	-23.433%				
gaussian blur	3	63.15	-3.14	-4.777%	54.087	-2.733	-4.411%	71.753	-2.307	-3.116%	102.883	-8.517	-7.638%				
jpeg_compression	3	64.98	-1.31	-1.976%	54.921	-1.899	-3.242%	73.662	-0.798	-1.063%	108.661	-2.839	-2.547%				
contrast	3	62.82	-3.47	-5.235%	54.158	-2.662	-4.608%	71.509	-2.531	-3.445%	101.702	-9.798	-8.788%				
elastic_transform	3	62.38	-4.14	-6.294%	53.008	-3.012	-5.301%	72.47	-1.59	-2.147%	102.869	-13.515	-12.474%				
pixelate	3	59.64	-6.65	-10.023%	52.385	-4.435	-7.809%	67.906	-5.154	-7.130%	88.665	-22.835	-20.480%				
snow	3	56.65	-9.64	-14.525%	48.831	-7.989	-14.609%	66.838	-7.202	-9.724%	78.893	-32.067	-29.244%				
frost	3	56.82	-9.47	-14.268%	49.547	-7.273	-12.800%	67.447	-6.031	-8.906%	79.193	-32.307	-28.975%				
fog	3	62.46	-3.83	-5.778%	54.206	-2.614	-4.401%	72.054	-2.006	-2.703%	103.011	-8.469	-7.596%				
brightness	3	64.52	-1.77	-2.670%	54.961	-1.859	-3.272%	72.686	-1.74	-1.856%	108.419	-3.081	-2.763%				
saturation	3	60.84	-5.45	-8.221%	52.369	-4.451	-7.433%	71.882	-2.178	-2.941%	95.462	-15.858	-14.225%				
saturation	3	65.17	-1.12	-1.690%	54.953	-1.867	-3.286%	73.374	-0.686	-0.926%	110.019	-1.481	-1.328%				
impulse_noise	4	60.13	-4.18	-6.920%	51.813	-5.007	-8.131%	69.198	-4.862	-6.565%	91.621	-19.879	-17.299%				
gaussian_noise	4	60.37	-5.92	-8.930%	52.21	-4.61	-8.113%	69.887	-4.173	-5.635%	94.151	-17.349	-15.559%				
shot_noise	4	59.82	-6.47	-9.760%	51.701	-5.119	-9.008%	68.451	-5.609	-7.573%	97.1	-19.8	-17.758%				
speckle_noise	4	61.48	-4.81	-7.756%	53.109	-3.711	-6.325%	69.786	-4.274	-5.771%	97.774	-13.726	-12.311%				
zoom blur	4	51.58	-14.71	-22.100%	45.961	-10.859	-19.111%	59.409	-14.651	-19.783%	46.098	-45.802	-58.650%				
defocus blur	4	60.91	-5.38	-8.116%	52.759	-4.461	-7.148%	68.298	-4.762	-6.430%	93.846	-17.554	-15.749%				
motion blur	4	59.37	-6.92	-10.493%	51.529	-5.23	-9.240%	68.423	-5.037	-6.622%	97.289	-24.211	-21.713%				
glass blur	4	56.9	-9.39	-14.165%	49.733	-7.09	-12.479%	67.475	-5.885	-8.891%	78.169	-33.331	-29.893%				
gaussian blur	4	61.27	-5.02	-7.571%	52.822	-3.998	-7.636%	70.044	-4.016	-5.422%	96.796	-14.704	-13.188%				
jpeg_compression	4	63.91	-2.38	-3.590%	54.19	-2.63	-4.299%	72.729	-1.31	-1.798%	104.728	-6.772	-6.073%				
contrast	4	58.57	-7.72	-11.646%	51.598	-5.222	-9.140%	68.092	-5.968	-8.639%	84.528	-21.972	-19.760%				
elastic_transform	4	61.91	-4.38	-6.607%	52.361	-4.459	-7.847%	71.25	-2.81	-3.764%	96.029	-14.571	-13.008%				
pixelate	4	55.66	-10.63	-16.069%	49.873	-6.947	-12.227%	64.317	-6.743	-13.155%	74.119	-38.381	-33.525%				
snow	4	53.52	-12.77	-17.605%	47.217	-9.403	-16.909%	63.715	-10.345	-13.969%	64.344	-42.722	-42.280%				
frost	4	58.8	-4.99	-7.439%	49.165	-7.655	-13.472%	67.59	-6.47	-8.736%	76.626	-34.974	-31.277%				
fog	4	60.48	-5.49	-8.262%	52.783	-4.037	-7.106%	71.781	-2.379	-3.077%	97.573	-13.027	-12.491%				
brightness	4	63.48	-2.81	-4.229%	54.452	-2.368	-4.167%	71.939	-1.221	-1.664%	106.108	-5.362	-4.755%				
saturation	4	58.61	-7.68	-11.585%	50.747	-6.073	-10.088%	68.839	-5.271	-7.650%	84.966	-26.534	-23.798%				
saturation	4	62.25	-4.04	-6.094%	52.56	-4.27	-7.497%	71.681	-2.329	-3.212%	102.472	-9.028	-8.097%				

Table 24: Performance and decrease ratio of Multiple Adapters on CLIP-T5 against image corruptions given severity 1 to 5.

corruption	severity	VQA				GQA				NLVR				Caption			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIoU	Decrease	Decrease Ratio	
impulse_noise	1	64.5	-1.65	-2.49%	53.18	-2.48	-4.45%	51.801	0.139	0.267%	106.036	-6.113	-5.49%				
gaussian_noise	1	65.19	-0.96	-1.451%	53.15	-2.45	-4.213%	51.529	0.411	0.792%	108.756	-3.393	-3.026%				
shot_noise	1	65.05	-1.1	-1.663%	53.776	-1.884	-3.348%	51.414	0.528	1.013%	108.489	-3.3	-2.843%				
speckle_noise	1	65.12	-1.03	-1.557%	53.458	-2.202	-3.955%	51.988	0.048	0.092%	109.065	-3.084	-2.750%				
room_blur	1	59.62	-6.53	-9.872%	49.857	-5.803	-10.426%	49.792	-2.148	-4.136%	82.302	-2.987	-26.703%				
defocus_blur	1	65.27	-0.88	-1.339%	53.514	-2.146	-3.855%	51.342	0.498	0.951%	105.326	-3.823	-3.409%				
motion_blur	1	64.86	-1.29	-1.959%	53.22	-2.44	-4.384%	51.399	0.541	1.041%	107.751	-4.398	-3.922%				
glass_blur	1	65.28	-0.87	-1.315%	53.616	-1.844	-3.313%	51.873	0.067	0.129%	108.701	-3.448	-3.074%				
gaussian_blur	1	65.86	-0.29	-0.438%	54.055	-1.605	-2.884%	51.328	0.612	1.179%	110.904	-3.245	-3.111%				
jpeg_compression	1	65.69	-0.46	-0.695%	53.729	-1.931	-3.470%	52.131	0.191	0.369%	109.975	-2.174	-1.938%				
contrast	1	64.86	-1.29	-1.950%	53.657	-2.003	-3.589%	51.313	0.627	1.207%	109.128	-3.021	-2.684%				
elastic_transform	1	65.38	-0.77	-1.164%	53.721	-1.939	-3.484%	52.175	0.235	0.452%	109.828	-2.321	-2.070%				
pixelate	1	62.76	-3.39	-5.125%	52.663	-2.997	-5.584%	51.299	0.641	1.234%	99.662	-12.547	-11.876%				
snow	1	61.59	-4.56	-6.883%	50.564	-5.096	-9.155%	45.457	0.483	0.830%	96.677	-15.192	-13.547%				
fract	1	62.46	-3.69	-5.757%	51.781	-3.879	-6.969%	51.342	0.598	1.151%	100.831	-11.318	-10.092%				
fog	1	64.35	-1.8	-2.721%	53.506	-2.154	-3.707%	51.012	0.928	1.787%	107.446	-4.703	-4.194%				
brightness	1	65.85	-0.45	-0.684%	53.641	-2.019	-3.627%	52.203	0.263	0.507%	111.572	-3.977	-3.515%				
saturation	1	65.64	-0.51	-0.771%	53.57	-2.09	-3.755%	52.347	0.407	0.783%	110.657	-4.092	-3.131%				
saturation	1	63.49	-2.66	-4.021%	52.605	-2.965	-5.372%	51.701	0.429	0.866%	107.388	-3.761	-3.425%				
corruption	severity	VQA				GQA				NLVR				Caption			
impulse_noise	2	63.51	-2.64	-3.991%	52.759	-2.901	-5.212%	51.902	0.038	0.074%	103.427	-8.722	-7.777%				
gaussian_noise	2	64.45	-1.7	-2.579%	53.077	-2.883	-4.441%	51.773	0.167	0.322%	107.49	-4.459	-4.155%				
shot_noise	2	64.27	-1.88	-2.842%	53.244	-2.416	-3.811%	51.328	0.612	1.179%	106.21	-5.939	-5.290%				
speckle_noise	2	64.16	-1.55	-2.343%	53.506	-2.154	-3.707%	51.6	0.34	0.654%	108.274	-3.875	-3.455%				
room_blur	2	56.66	-8.49	-14.360%	48.06	-7.6	-13.654%	50.194	-1.746	-3.362%	66.462	-45.087	-40.718%				
defocus_blur	2	64.48	-1.67	-2.525%	52.942	-2.178	-3.884%	51.457	0.483	0.930%	106.236	-3.593	-3.281%				
motion_blur	2	63.68	-2.47	-3.734%	52.536	-3.124	-5.012%	52.16	0.22	0.424%	107.258	-7.891	-7.037%				
glass_blur	2	64.13	-2.02	-3.054%	53.212	-2.448	-4.389%	51.801	0.219	0.427%	104.629	-7.52	-6.706%				
gaussian_blur	2	65.49	-0.66	-0.999%	53.872	-1.788	-3.219%	53.021	1.081	2.082%	109.446	-2.703	-2.419%				
jpeg_compression	2	65.38	-0.77	-1.164%	53.721	-1.939	-3.484%	52.175	0.235	0.452%	109.828	-2.321	-2.070%				
contrast	2	62.82	-3.33	-5.034%	51.678	-3.982	-7.155%	51.084	-1.856	-3.649%	94.903	-17.246	-15.378%				
elastic_transform	2	62.26	-3.89	-5.881%	52.337	-3.323	-5.969%	51.5	0.44	0.847%	97.551	-14.988	-13.017%				
pixelate	2	58.42	-7.73	-11.660%	48.911	-6.749	-12.126%	50.171	-1.23	-2.367%	84.361	-27.788	-24.778%				
snow	2	59.3	-6.85	-10.355%	49.666	-5.994	-10.769%	51.213	-0.727	-1.409%	87.722	-24.429	-21.783%				
fract	2	63.61	-2.54	-3.849%	53.244	-2.416	-4.311%	50.423	-1.577	-3.020%	106.212	-5.937	-5.294%				
brightness	2	65.43	-0.72	-1.088%	53.882	-2.178	-3.913%	52.619	0.679	1.308%	103.536	-1.613	-1.483%				
saturation	2	62.4	-3.75	-5.669%	52.306	-3.354	-6.027%	52.634	0.694	1.336%	99.537	-12.617	-11.240%				
saturation	2	60.8	-5.35	-8.808%	52.043	-3.67	-6.490%	50.782	-1.158	-2.229%	102.179	-9.97	-8.809%				
corruption	severity	VQA				GQA				NLVR				Caption			
impulse_noise	3	62.57	-3.58	-5.412%	53.037	-2.623	-4.712%	51.342	0.598	1.151%	100.779	-11.37	-10.138%				
gaussian_noise	3	62.63	-3.52	-5.231%	52.767	-2.993	-5.189%	51.098	0.482	0.921%	101.76	-10.389	-9.264%				
shot_noise	3	62.71	-3.44	-5.209%	53.236	-2.424	-4.355%	50.667	-1.273	-2.450%	101.794	-10.55	-9.233%				
speckle_noise	3	62.7	-3.45	-5.215%	52.889	-2.679	-4.812%	50.596	-1.344	-2.588%	102.21	-9.859	-8.863%				
room_blur	3	53.85	-12.2	-18.434%	46.031	-8.729	-15.602%	50.179	-1.761	-3.390%	68.185	-60.264	-53.760%				
defocus_blur	3	62.69	-3.46	-5.231%	52.139	-3.521	-6.272%	51.414	-0.365	-0.703%	100.867	-11.462	-10.220%				
motion_blur	3	62.11	-4.04	-6.107%	52.226	-3.434	-6.169%	51.069	0.871	1.676%	98.529	-13.62	-12.145%				
glass_blur	3	58.93	-7.22	-10.915%	49.834	-6.026	-10.826%	50.682	-1.258	-2.422%	84.772	-27.377	-24.412%				
gaussian_blur	3	63.51	-2.64	-3.991%	52.274	-3.386	-6.045%	51.084	-0.856	-1.649%	101.916	-12.33	-9.125%				
jpeg_compression	3	65.16	-0.99	-1.497%	53.888	-1.812	-3.266%	52.663	0.721	1.391%	109.852	-2.997	-2.618%				
contrast	3	62.71	-3.44	-5.209%	52.568	-3.092	-5.555%	50.208	-1.732	-3.344%	101.41	-8.779	-7.792%				
elastic_transform	3	62.32	-3.3	-4.961%	52.346	-3.272	-5.089%	51.189	-2.284	-4.461%	96.23	-15.019	-13.493%				
pixelate	3	60.13	-6.02	-9.101%	51.526	-4.134	-7.426%	50.782	-1.158	-2.229%	87.638	-24.511	-21.855%				
snow	3	56.34	-9.51	-14.736%	47.941	-7.719	-13.808%	50.022	-1.918	-3.694%	76.488	-35.651	-31.789%				
fract	3	57.14	-9.01	-13.613%	48.625	-7.035	-12.640%	51.236	0.684	1.317%	76.995	-35.154	-31.340%				
fog	3	62.53	-3.62	-5.472%	52.862	-2.978	-5.072%	50.926	-1.014	-1.952%	102.556	-9.593	-8.544%				
brightness	3	64.65	-1.5	-2.284%	53.395	-2.263	-4.070%	52.505	0.565	1.076%	108.24	-3.969	-3.460%				
saturation	3	61.0	-5.15	-7.785%	51.256	-4.484	-7.912%	52.548	0.608	1.170%	93.794	-18.355	-16.860%				
saturation	3	65.26	-0.89	-1.345%	53.156	-2.504	-4.089%	52.878	0.938	1.806%	110.232	-2.017	-1.799%				
corruption	severity	VQA				GQA				NLVR				Caption			
impulse_noise	4	59.82	-6.35	-8.569%	51.757	-3.903	-7.012%	50.553	-1.367	-2.671%	91.794	-20.355	-18.159%				
gaussian_noise	4	60.12	-6.03	-9.161%	51.455	-4.205	-7.555%	51.177	-0.77	-1.483%	93.233	-18.916	-16.807%				
shot_noise	4	59.75	-6.4	-9.675%	51.193	-4.467	-8.026%	50.825	-1.115	-2.146%	90.816	-21.33	-19.022%				
speckle_noise	4	61.49	-4.66	-7.085%	52.369	-3.291	-5.912%	50.61	-1.33	-2.561%	98.103	-14.066	-12.324%				
room_blur	4	52.17	-13.98	-21.314%	45.254	-10.406	-18.696%	49.892	-2.048	-3.942%	62.244	-69.905	-62.333%				
defocus_blur	4	60.84	-5.31	-8.072%	51.034	-4.626	-8.125%	50.754	-1.186	-2.284%	93.41	-17.399	-16.209%				
motion_blur	4	59.58	-6.57	-9.972%	50.204	-5.566	-9.646%	50.452	-1.488	-2.865%	85.262	-26.867	-23.975%				
glass_blur	4	56.72	-9.43	-14.255%	48.434	-7.226	-12.983%	50.481	-1.499	-2.809%	75.714	-36.435	-32.488%				
gaussian_blur	4	61.6	-4.55	-6.879%	51.288	-4.372	-7.855%	50.754	-1.186	-2.284%	96.23	-15.019	-13.493%				
jpeg_compression	4	64.15	-2.0	-3.023%	53.132	-2.528	-4.514%	52.332	0.392	0.756%	104.967	-3.618	-3.404%				
contrast	4	58.4	-7.75	-11.716%	50.668	-4.992	-8.969%	50.719	-1.761	-3.390%	88.829	-23.32	-20.794%				
elastic_transform	4	62.01	-4.14	-6.299%	51.344	-3.116	-5.775%	51.285	-2.055	-3.961%	95.188	-16.961	-15.124%				
pixelate	4	55.87	-10.28	-15.549%	48.726	-6.94	-12.469%	50.39	-1.631	-3.144%	72.813	-39.336	-35.075%				
snow	4	51.58	-9.91	-16.972%	47.961	-7.061	-13.602%	50.179	-1.761	-3.390%	68.185	-60.264	-53.760%				
fract	4	56.21	-9.94	-15.540%	47.925	-7.753	-13.897%	51.543	-0.97	-0.764%	75.706	-37.063	-33.048%				
fog	4	60.89	-5.26	-7.982%	51.382	-4.08	-7.741%	51.41	-0.799	-1.536%	96.906	-15.153	-13.511%				
brightness	4	62.37	-3.68	-5.669%	52.306	-3.354	-6.027%	52.634	0.694	1.336%	99.537	-12.617	-11.240%				
saturation	4	58.63	-7.52	-11.368%	49.674	-5.986	-10.755%	52.261	-0.321	-0.617%	83.395	-28.845	-25.718%				
saturation	4	62.69	-3.48	-5.231%	51.932	-3.728	-6.698%	51.844	-0.098	-0.184%	103.63	-8.519	-7.596%				
corruption	severity	VQA				GQA				NLVR				Caption			
impulse_noise	5	56.02	-10.13	-15.314%	49.412	-6.248	-11.226%	50.007	-1.933	-3.721%	79.067	-33.082	-29.499%				
gaussian_noise	5	56.02	-10.13	-15.314%	49.507	-6.153	-11.054%	50.409	-1.521	-2.947%	78.336	-333.0					

Table 26: Performance and decrease ratio of Single Adapter on CLIP-T5 against image corruptions given severity 1 to 5.

corruption	severity	VQA			GQA			NLVR			Caption		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
impulse_noise	1	64.39	-2.02	-3.042%	48.179	-7.721	-13.812%	71.695	-1.491%	-0.001%	106.315	-5.385	-4.821%
gaussian_noise	1	65.09	-1.32	-1.988%	48.72	-7.18	-12.844%	72.04	-0.74	-1.017%	108.347	-3.353	-3.002%
shot_noise	1	65.2	-1.21	-1.822%	48.776	-7.124	-12.745%	71.81	-0.97	-1.333%	108.413	-3.287	-2.942%
speckle_noise	1	65.06	-1.35	-2.031%	48.943	-6.957	-12.446%	71.523	-1.257	-1.727%	108.585	-3.135	-2.807%
room_blur	1	59.35	-7.06	-10.611%	45.54	-10.36	-18.531%	56.064	-16.716	-22.967%	84.143	-27.557	-24.570%
defocus_blur	1	65.11	-1.3	-1.989%	48.863	-7.037	-12.588%	72.241	-0.59	-0.741%	108.771	-3.529	-3.196%
motion_blur	1	64.88	-1.53	-2.364%	49.022	-6.878	-12.304%	71.781	-0.999	-1.372%	107.027	-4.073	-4.184%
glass_blur	1	65.3	-1.11	-1.671%	48.895	-7.005	-12.531%	71.523	-1.257	-1.727%	108.294	-3.366	-2.969%
gaussian_blur	1	65.94	-0.47	-0.708%	49.38	-6.52	-11.664%	72.686	-0.094	-0.130%	110.921	-1.779	-0.608%
jpeg_compression	1	65.75	-0.66	-0.994%	49.046	-6.854	-12.772%	72.772	-0.008	-0.011%	110.994	-1.606	-1.437%
contrast	1	62.97	-1.34	-2.018%	49.308	-6.392	-11.792%	71.638	-1.142	-1.569%	108.489	-3.211	-2.814%
elastic_transform	1	66.04	-0.37	-0.557%	48.831	-7.069	-12.645%	72.04	-0.74	-1.017%	109.629	-2.071	-1.854%
pixelate	1	62.97	-1.34	-2.018%	48.831	-7.069	-12.645%	72.04	-0.74	-1.017%	109.629	-2.071	-1.854%
snow	1	61.89	-4.52	-6.806%	46.613	-8.267	-16.613%	69.614	-1.166	-1.550%	97.265	-14.435	-12.923%
fog	1	62.91	-1.35	-2.070%	47.631	-8.269	-14.793%	70.073	-2.707	-3.719%	109.419	-11.281	-10.100%
brightness	1	66.13	-0.28	-0.422%	49.865	-6.035	-10.796%	72.585	-0.195	-0.268%	111.78	-0.88	-0.072%
saturate	1	65.88	-0.53	-0.798%	49.348	-6.552	-11.721%	72.528	-0.252	-0.347%	110.277	-1.423	-1.274%
saturation	1	65.38	-2.83	-4.281%	48.752	-7.148	-12.787%	70.448	-1.932	-2.654%	106.852	-4.848	-4.340%
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
impulse_noise	2	63.59	-2.82	-4.246%	47.718	-8.182	-14.636%	71.035	-1.748	-2.398%	103.682	-8.048	-7.205%
gaussian_noise	2	64.125	-1.75	-2.573%	47.595	-7.185	-12.783%	71.185	-1.627	-2.141%	103.717	-7.418	-6.714%
shot_noise	2	64.18	-2.23	-3.358%	48.06	-7.84	-14.025%	70.906	-1.874	-2.575%	106.358	-5.342	-4.793%
speckle_noise	2	64.23	-1.99	-2.861%	48.096	-7.204	-12.887%	71.465	-1.315	-1.806%	107.147	-5.353	-3.811%
room_blur	2	56.34	-10.07	-15.103%	43.862	-12.018	-21.534%	60.499	-12.281	-16.874%	69.331	-42.369	-37.013%
defocus_blur	2	64.54	-1.87	-2.816%	48.402	-7.498	-13.413%	71.422	-1.358	-1.865%	106.266	-5.066	-4.599%
motion_blur	2	64.23	-2.23	-3.358%	47.595	-7.185	-12.783%	71.185	-1.627	-2.141%	103.717	-7.418	-6.714%
glass_blur	2	64.06	-2.35	-3.539%	48.346	-7.554	-13.519%	70.088	-2.052	-2.699%	104.12	-7.568	-6.776%
gaussian_blur	2	64.78	-1.63	-2.454%	48.672	-7.228	-12.930%	72.111	-0.609	-0.919%	107.162	-4.538	-4.063%
jpeg_compression	2	65.36	-0.85	-1.289%	49.144	-6.759	-12.091%	71.753	-1.027	-1.412%	108.73	-2.197	-2.049%
contrast	2	64.61	-1.8	-2.708%	48.206	-7.204	-12.930%	70.949	-1.831	-2.564%	106.802	-4.898	-4.355%
elastic_transform	2	62.7	-3.71	-5.587%	47.217	-8.083	-15.522%	68.25	-4.35	-6.224%	95.436	-16.264	-14.506%
pixelate	2	62.6	-4.41	-6.641%	47.559	-8.241	-14.921%	67.217	-5.863	-8.644%	97.363	-14.307	-12.809%
snow	2	58.66	-7.75	-11.670%	45.007	-10.893	-19.486%	66.829	-5.951	-8.766%	86.053	-25.647	-22.961%
fog	2	59.41	-7.0	-10.541%	45.182	-10.718	-19.178%	68.164	-6.068	-8.642%	88.041	-23.609	-21.183%
brightness	2	63.87	-2.54	-3.825%	48.545	-7.355	-13.157%	69.872	-5.908	-8.117%	105.731	-9.589	-8.344%
saturation	2	65.54	-0.87	-1.310%	49.03	-6.87	-12.720%	72.083	-0.697	-0.958%	110.209	-1.491	-1.335%
saturation	2	62.68	-3.73	-5.645%	47.214	-8.469	-15.490%	70.662	-2.118	-2.911%	100.122	-11.578	-10.340%
saturation	2	60.78	-5.63	-8.478%	47.599	-8.301	-14.850%	68.308	-4.472	-6.145%	102.909	-8.731	-7.817%
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
impulse_noise	3	62.77	-3.64	-5.481%	47.663	-8.237	-14.736%	69.643	-3.137	-4.311%	100.138	-10.542	-9.438%
gaussian_noise	3	62.75	-3.66	-5.511%	47.567	-8.333	-14.907%	69.872	-3.908	-5.955%	101.867	-9.833	-8.803%
shot_noise	3	62.77	-3.64	-5.481%	47.663	-8.237	-14.736%	69.643	-3.137	-4.311%	100.138	-10.542	-9.438%
speckle_noise	3	62.75	-3.66	-5.511%	47.567	-8.333	-14.907%	69.872	-3.908	-5.955%	101.867	-9.833	-8.803%
room_blur	3	55.54	-12.87	-19.300%	42.161	-13.739	-24.578%	57.93	-14.465	-25.044%	56.396	-30.304	-49.511%
defocus_blur	3	62.63	-3.78	-5.692%	47.241	-8.659	-15.490%	70.303	-2.477	-3.404%	100.458	-11.242	-10.064%
motion_blur	3	62.53	-3.88	-5.842%	47.138	-8.762	-15.675%	70.001	-2.779	-3.818%	99.777	-12.723	-12.285%
glass_blur	3	59.01	-7.4	-11.434%	45.023	-10.877	-19.458%	62.437	-5.987	-8.848%	87.217	-24.363	-20.930%
gaussian_blur	3	63.39	-3.02	-4.548%	47.639	-8.261	-14.739%	70.389	-3.291	-3.285%	109.319	-9.761	-8.739%
jpeg_compression	3	65.07	-1.34	-2.018%	49.205	-6.905	-11.977%	72.137	-0.41	-0.564%	108.003	-3.097	-3.106%
contrast	3	62.86	-3.55	-5.348%	47.798	-8.102	-14.494%	69.944	-3.836	-5.979%	102.8	-9.8	-7.988%
elastic_transform	3	64.31	-2.06	-3.102%	48.267	-7.633	-13.655%	70.963	-1.817	-2.496%	107.001	-6.963	-6.387%
pixelate	3	59.89	-6.52	-9.818%	46.438	-9.462	-16.926%	66.069	-6.711	-9.221%	88.354	-23.346	-20.901%
snow	3	56.58	-8.31	-14.802%	43.775	-12.225	-21.691%	64.834	-7.946	-11.983%	78.536	-33.164	-29.690%
fog	3	62.7	-3.64	-5.481%	47.359	-12.441	-21.719%	66.097	-6.083	-8.182%	94.485	-23.215	-20.841%
brightness	3	62.41	-4.0	-6.023%	47.631	-8.269	-14.793%	70.36	-2.42	-3.325%	103.165	-8.535	-7.641%
saturation	3	64.99	-1.42	-2.139%	49.066	-6.814	-12.710%	71.509	-1.271	-1.747%	108.955	-3.105	-2.789%
saturation	3	61.06	-5.35	-8.090%	46.494	-9.406	-16.827%	69.373	-3.023	-4.155%	94.485	-23.215	-20.841%
saturation	3	65.57	-0.84	-1.265%	48.823	-7.077	-12.599%	72.255	-0.525	-0.721%	109.806	-1.894	-1.669%
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
impulse_noise	4	60.08	-4.35	-6.552%	46.629	-9.271	-16.585%	67.92	-4.46	-6.077%	91.542	-20.158	-18.047%
gaussian_noise	4	60.52	-5.89	-8.899%	46.669	-9.231	-16.514%	67.92	-4.46	-6.077%	91.542	-20.158	-18.047%
shot_noise	4	59.86	-6.55	-9.596%	46.43	-9.47	-16.840%	66.838	-5.922	-8.137%	92.568	-19.432	-17.937%
speckle_noise	4	61.76	-4.65	-7.002%	47.368	-8.352	-13.262%	68.121	-4.609	-6.401%	97.848	-13.762	-12.311%
room_blur	4	51.82	-14.59	-21.970%	40.491	-15.409	-27.569%	56.284	-16.486	-22.625%	45.743	-34.943	-59.036%
defocus_blur	4	60.99	-5.42	-8.184%	46.343	-9.557	-17.097%	68.121	-4.609	-6.401%	97.848	-13.762	-12.311%
motion_blur	4	59.72	-6.69	-10.074%	45.818	-10.082	-18.060%	67.016	-5.764	-7.990%	96.809	-16.511	-15.229%
glass_blur	4	57.15	-9.26	-13.944%	43.401	-12.499	-22.359%	67.303	-4.477	-5.254%	77.581	-34.119	-30.545%
gaussian_blur	4	61.36	-4.85	-7.303%	46.661	-9.239	-16.529%	69.212	-3.568	-4.902%	95.388	-16.212	-14.603%
jpeg_compression	4	64.02	-2.39	-3.599%	48.227	-7.673	-13.760%	70.777	-0.203	-0.275%	104.095	-7.605	-6.808%
contrast	4	58.47	-7.94	-11.966%	45.5	-10.4	-17.629%	65.121	-7.659	-10.333%	90.667	-21.033	-18.830%
elastic_transform	4	62.94	-1.17	-1.629%	46.542	-8.958	-16.741%	71.107	-1.073	-1.296%	95.658	-16.062	-14.389%
pixelate	4	55.85	-10.56	-15.010%	44.141	-11.759	-21.073%	61.461	-11.319	-15.552%	72.151	-39.459	-35.400%
snow	4	53.48	-12.83	-19.093%	41.207	-14.093	-26.285%	62.71	-10.07	-13.836%	68.748	-48.132	-43.009%
fog	4	55.96	-10.45	-15.760%	43.212	-12.388	-21.600%	65.703	-7.013	-9.636%	76.794	-34.906	-31.250%
brightness	4	61.0	-5.41	-6.919%	46.677	-9.223	-16.500%	69.427	-3.353	-4.607%	97.95	-13.75	-12.103%
saturation	4	64.02	-2.38	-3.549%	48.489	-7.411	-13.257%	70.518	-2.262	-3.106%	106.423	-5.277	-4.724%
saturation	4	58.95	-7.46	-11.233%	44.514	-11.386	-20.368%	67.26	-5.27	-7.585%	85.064	-26.636	-23.840%
saturation	4	62.76	-3.65	-5.496%	47.122	-8.718	-15.703%	69.528	-2.552	-3.469%	103.763	-7.917	-7.088%
corruption	severity	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
impulse_noise	5	56.48	-9.93	-14.935%	44.657	-11.243	-20.112%	64.993	-7.845	-10.779%	78.788	-32.912	-29.465%
gaussian_noise	5	56.4	-10.01	-15.073%	43.91	-11.99	-21.449%	65.05	-7.73	-10.622%	79.632	-32.714	-28.709%
shot_noise	5	57.21	-9.2	-13.819%	44.586	-11.314	-20.249%	65.021	-7.799	-10.661%	81.986	-29.714	-26.023%
speckle_noise	5	59.94	-6.47	-9.743%	46.828	-9.072	-16.229%	64.885	-6.295	-8.650%	91.411	-20.289	-18.164%
room_blur	5	49.58	-16.83	-25.349%	39.066	-16.84	-32.555%	55.06	-17.249	-24.388%	37.009	-49.911	-66.807%
defocus_blur	5	59.16	-7.25	-10.917%	45.58	-10.32	-18.462%	66.284	-6.496	-8.926%	83.256	-28.444	-25.465%
motion_blur	5												

Table 28: Performance and decrease ratio of Full Finetuning on CLIP-BART against text corruptions given severity 1 to 4.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	1.0	55.75	-11.0	-16.479%	44.769	-10.271	-18.662%	56.581	-16.429	-22.502%			
typos	1.0	57.86	-8.89	-13.318%	46.128	-8.912	-16.192%	70.044	-2.966	-4.062%			
keyboard	1.0	56.36	-10.39	-15.566%	44.753	-10.287	-18.691%	69.671	-3.339	-4.573%			
spell_error	1.0	57.38	-9.37	-14.037%	46.375	-8.665	-15.744%	69.915	-3.095	-4.239%			
random_char_insert	1.0	52.55	-14.2	-21.273%	36.015	-19.025	-34.565%	64.605	-8.405	-11.513%			
random_char_replace	1.0	48.69	-18.06	-27.056%	32.35	-22.69	-41.224%	62.48	-10.53	-14.422%			
random_char_swap	1.0	42.09	-24.66	-36.944%	26.499	-28.541	-51.856%	58.088	-14.922	-20.438%			
random_char_delete	1.0	50.47	-16.28	-24.309%	34.982	-20.058	-36.443%	63.298	-9.712	-13.302%			
random_word_insert	1.0	66.71	-0.04	-0.060%	52.568	-2.472	-4.491%	71.308	-1.702	-2.323%			
random_word_delete	1.0	66.55	-0.2	-0.300%	53.896	-1.144	-2.079%	72.255	-0.755	-1.034%			
swap_syn_word_emb	1.0	55.57	-11.18	-16.749%	46.232	-8.808	-16.004%	67.604	-5.406	-7.404%			
swap_syn_word_net	1.0	58.71	-8.04	-12.456%	47.535	-7.505	-13.635%	70.533	-2.477	-3.393%			
random_word_swap	1.0	66.67	-0.08	-0.120%	54.071	-0.969	-1.761%	72.747	-0.54	-0.739%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	2.0	55.49	-11.26	-16.869%	39.998	-15.042	-27.328%	65.509	-7.501	-10.274%			
typos	2.0	56.43	-10.32	-15.461%	41.294	-13.746	-24.974%	66.241	-6.769	-9.272%			
keyboard	2.0	55.84	-10.91	-16.345%	40.022	-15.018	-27.285%	65.652	-7.338	-10.078%			
spell_error	2.0	57.18	-9.57	-14.337%	42.28	-12.76	-23.183%	66.671	-6.339	-8.682%			
random_char_insert	2.0	41.54	-25.21	-37.768%	26.101	-28.939	-52.578%	59.05	-13.96	-19.121%			
random_char_replace	2.0	36.15	-30.6	-45.843%	22.404	-32.636	-59.285%	56.796	-16.214	-22.207%			
random_char_swap	2.0	32.66	-34.09	-51.071%	20.782	-34.258	-62.241%	54.816	-18.194	-24.920%			
random_char_delete	2.0	38.71	-28.04	-42.007%	25.338	-29.702	-53.965%	57.887	-15.123	-20.713%			
random_word_insert	2.0	65.25	-1.5	-2.247%	47.253	-7.807	-14.184%	68.724	-4.286	-5.870%			
random_word_delete	2.0	60.35	-6.4	-9.588%	48.434	-6.606	-12.003%	69.686	-3.224	-4.553%			
swap_syn_word_emb	2.0	55.13	-11.62	-17.408%	42.113	-12.927	-23.486%	62.94	-10.07	-13.793%			
swap_syn_word_net	2.0	57.75	-9.0	-13.483%	44.443	-10.597	-19.254%	66.557	-6.453	-8.839%			
random_word_swap	2.0	65.49	-1.26	-1.888%	51.868	-3.172	-5.762%	72.126	-0.884	-1.211%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	3.0	50.29	-16.46	-24.659%	31.881	-23.159	-42.077%	59.968	-13.042	-17.863%			
typos	3.0	51.4	-15.35	-22.966%	34.306	-20.734	-37.671%	63.901	-9.109	-12.476%			
keyboard	3.0	50.14	-16.61	-24.884%	31.857	-23.183	-42.120%	62.981	-10.429	-14.285%			
spell_error	3.0	51.5	-15.43	-23.116%	34.616	-20.424	-37.108%	71.753	-10.857	-14.009%			
random_char_insert	3.0	33.4	-33.35	-49.963%	21.259	-33.781	-61.375%	55.82	-17.19	-23.544%			
random_char_replace	3.0	28.51	-38.28	-57.286%	18.771	-36.269	-65.896%	54.83	-18.18	-24.901%			
random_char_swap	3.0	28.26	-38.49	-57.663%	18.588	-36.452	-66.228%	51.859	-21.151	-28.970%			
random_char_delete	3.0	31.71	-35.024	-52.404%	21.021	-34.019	-61.808%	54.557	-18.453	-25.272%			
random_word_insert	3.0	63.51	-3.24	-4.854%	41.628	-13.412	-24.367%	66.341	-6.669	-9.134%			
random_word_delete	3.0	57.44	-9.31	-13.948%	43.862	-11.178	-20.308%	67.188	-5.822	-7.974%			
swap_syn_word_emb	3.0	50.7	-16.05	-24.045%	36.039	-19.001	-34.522%	59.61	-13.4	-18.354%			
swap_syn_word_net	3.0	54.09	-12.62	-18.966%	39.323	-13.717	-28.556%	64.375	-8.635	-11.370%			
random_word_swap	3.0	64.85	-1.9	-2.846%	50.024	-5.016	-9.114%	70.963	-2.047	-2.804%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	4.0	43.5	-23.25	-34.831%	26.88	-28.16	-51.162%	56.854	-16.156	-22.129%			
typos	4.0	42.21	-23.54	-35.366%	27.54	-27.5	-49.963%	59.925	-13.085	-17.922%			
keyboard	4.0	40.5	-26.25	-39.326%	25.298	-29.742	-54.037%	59.064	-13.946	-19.101%			
spell_error	4.0	43.0	-23.75	-35.581%	28.518	-26.522	-48.187%	58.72	-14.29	-19.573%			
random_char_insert	4.0	27.94	-38.81	-58.142%	19.145	-35.895	-65.217%	54.342	-18.668	-25.569%			
random_char_replace	4.0	23.72	-43.03	-64.646%	17.427	-37.613	-68.377%	53.108	-19.902	-27.600%			
random_char_swap	4.0	26.2	-40.55	-60.749%	18.047	-36.993	-67.210%	51.514	-21.496	-29.442%			
random_char_delete	4.0	26.09	-40.66	-60.914%	18.97	-36.07	-65.533%	53.222	-19.788	-27.103%			
random_word_insert	4.0	59.59	-7.16	-10.727%	37.963	-17.077	-31.026%	65.093	-7.917	-10.844%			
random_word_delete	4.0	50.03	-16.72	-25.409%	37.987	-17.053	-30.983%	66.04	-6.97	-9.547%			
swap_syn_word_emb	4.0	46.8	-19.05	-29.888%	32.994	-22.046	-40.054%	57.485	-15.525	-21.264%			
swap_syn_word_net	4.0	50.61	-16.14	-24.180%	36.476	-18.564	-33.727%	61.059	-11.951	-16.369%			
random_word_swap	4.0	64.18	-2.57	-3.850%	48.839	-6.201	-11.266%	71.537	-1.473	-2.017%			

Table 29: Performance and decrease ratio of Full Finetuning on CLIP-BART against text corruptions given severity 5.

corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	5.0	40.23	-26.52	-39.730%	23.66	-31.38	-57.012%	54.17	-18.84	-25.805%
pronunciation	5.0	65.94	-0.81	-1.213%	54.309	-0.731	-1.328%	72.355	-0.655	-0.897%
typos	5.0	30.35	-36.4	-54.532%	19.979	-35.061	-63.700%	53.136	-19.874	-27.221%
keyboard	5.0	27.34	-39.41	-59.041%	18.055	-36.985	-67.196%	50.854	-22.156	-30.346%
spell_error	5.0	36.24	-30.51	-45.708%	24.257	-30.783	-55.929%	55.691	-17.319	-23.721%
random_char_insert	5.0	23.51	-43.24	-64.779%	17.849	-37.191	-67.572%	52.591	-20.419	-27.968%
random_char_replace	5.0	19.82	-46.93	-70.307%	16.092	-38.948	-70.764%	51.974	-21.036	-28.813%
random_char_swap	5.0	25.33	-41.42	-62.052%	18.922	-36.118	-65.621%	52.591	-20.419	-27.968%
random_char_delete	5.0	22.22	-44.53	-66.712%	18.095	-36.945	-67.124%	51.615	-21.395	-29.305%
to-passive	5.0	62.82	-3.93	-5.888%	51.224	-3.816	-6.932%	70.131	-2.879	-3.944%
to-formal	5.0	65.81	-0.94	-1.408%	54.826	-0.214	-0.389%	73.044	0.034	0.047%
to-true	5.0	66.24	-0.51	-0.744%	53.935	-1.105	-2.007%	72.47	-0.54	-0.739%
to-casual	5.0	63.58	-1.17	-1.749%	51.566	-3.474	-6.311%	71.107	-1.903	-2.607%
to-active	5.0	64.88	-1.87	-2.801%	53.26	-1.78	-3.235%	72.384	-0.626	-0.857%
double-denial	5.0	66.63	-0.12	-0.180%	55.033	-0.007	-0.013%	73.016	0.006	0.008%
insert_adv	5.0	61.67	-5.08	-7.610%	49.61	-5.43	-9.865%	70.547	-2.463	-3.374%
append_ri	5.0	60.56	-6.19	-9.273%	47.194	-7.846	-14.256%	66.571	-6.439	-8.819%
random_word_insert	5.0	55.76	-10.99	-16.646%	32.0	-23.04	-41.860%	63.56	-9.654	-13.223%
drop_an	5.0	43.85	-22.9	-34.307%	30.919	-24.121	-43.824%	60.672	-12.338	-16.899%
drop_rand_an_jon	5.0	58.63	-16.25	-25.566%	40.356	-14.684	-26.678%	68.451	-4.559	-6.244%
drop_vb	5.0	45.83	-8.12	-12.165%	40.356	-14.684	-26.678%	68.451	-4.559	-6.244%
drop_vb_an	5.0	43.14	-23.61	-35.371%	27.373	-27.667	-50.267%	57.858	-15.152	-20.753%
only_an	5.0	32.64	-34.11	-51.101%	20.234	-34.806	-63.238%	52.16	-20.85	-28.557%
only_vb	5.0	31.97	-14.78	-22.124%	39.831	-15.209	-27.632%	65.265	-7.745	-10.608%
only_vb_an	5.0	37.61	-29.14	-43.655%	21.848	-31.192	-60.306%	54.084	-18.926	-27.600%
drop_rand_jon_vb	5.0	63.63	-3.12	-4.447%	47.257	-7.783	-14.410%	71.035	-1.975	-2.705%
drop_first	5.0	56.8	-9.95	-14.906%	39.744	-15.296	-27.791%	70.748	-2.262	-3.098%
drop_last	5.0	54.22	-12.53	-18.772%	45.548	-9.492	-17.246%	69.068	-3.942	-5.599%
drop_first_and_last	5.0	40.6	-26.15	-39.176%	31.436	-23.604	-42.885%	65.294	-7.716	-10.509%
shuffle_order	5.0	59.93	-6.82	-10.127%	41.724	-13.316	-24.194%	68.408	-4.602	-6.303%
random_word_delete	5.0	42.58	-24.17	-36.210%	31.849	-23.191	-42.134%	63.959	-9.051	-12.397%
swap_syn_word_emb	5.0	46.06	-20.69	-30.969%	32.0	-23.04	-41.860%	57.916	-15.094	-20.674%
swap_syn_word_net	5.0	46.31	-20.44	-30.622%	35.165	-19.875	-36.111%	58.619	-14.391	-19.711%
backtrans	5.0	62.77	-3.98	-5.963%	50.549	-4.491	-8.160%	71.781	-1.228	-1.683%
backtrans_net	5.0	63.58	-1.17	-1.749%	51.566	-3.474	-6.311%	71.107	-1.903	-2.607%
nonsense	5.0	60	-66.75	-100.00%	17.038	-38.002	-69.045%	51.227	-21.83	-29.835%

Table 32: Performance and decrease ratio of Hyperformer on CLIP-BART against text corruptions given severity 1 to 4.

corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	1.0	54.68	-10.7	-16.366%	42.821	-9.699	-18.468%	56.94	-15.27	-21.147%
typos	1.0	56.24	-9.14	-13.908%	45.047	-7.473	-14.229%	69.772	-2.438	-3.377%
keyboard	1.0	55.04	-10.34	-18.515%	43.95	-8.57	-16.318%	69.198	-3.012	-4.172%
spell_error	1.0	55.72	-9.66	-14.775%	44.753	-7.767	-14.789%	69.241	-2.969	-4.112%
random_char_insert	1.0	51.13	-14.25	-21.796%	35.983	-16.537	-31.486%	63.815	-8.395	-11.626%
random_char_replace	1.0	47.55	-17.83	-27.271%	31.825	-20.695	-39.403%	62.007	-10.203	-14.130%
random_char_swap	1.0	41.02	-24.36	-37.259%	26.244	-26.276	-50.030%	57.04	-15.17	-21.008%
random_char_delete	1.0	49.05	-16.33	-24.977%	34.449	-18.071	-34.408%	63.241	-8.969	-12.421%
random_word_insert	1.0	65.32	-0.60	-0.092%	50.636	-1.884	-3.587%	71.652	-0.558	-0.773%
random_word_delete	1.0	65.21	-0.17	-0.260%	51.36	-1.16	-2.210%	71.681	-0.529	-0.733%
swap_syn_word_emb	1.0	55.18	-10.2	-15.601%	45.5	-7.02	-13.366%	67.317	-4.893	-6.776%
swap_syn_word_net	1.0	57.82	-7.56	-11.563%	46.128	-6.392	-12.170%	69.427	-2.783	-3.854%
random_word_swap	1.0	65.28	-0.1	-0.153%	51.725	-0.795	-1.513%	71.781	-0.429	-0.594%
corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	2.0	53.7	-11.68	-17.865%	39.307	-13.213	-25.159%	65.308	-6.902	-9.558%
typos	2.0	54.86	-10.52	-16.091%	41.048	-11.472	-21.843%	66.169	-6.041	-8.366%
keyboard	2.0	54.58	-10.8	-16.519%	40.269	-12.251	-23.327%	66.857	-5.653	-7.829%
spell_error	2.0	55.08	-10.3	-15.754%	42.217	-10.303	-19.618%	65.624	-6.586	-9.121%
random_char_insert	2.0	41.12	-24.26	-37.106%	25.529	-26.991	-51.392%	57.744	-14.466	-20.034%
random_char_replace	2.0	36.76	-28.62	-43.775%	22.571	-29.949	-57.024%	55.935	-16.275	-22.538%
random_char_swap	2.0	33.93	-33.45	-51.162%	19.868	-32.652	-62.171%	52.82	-19.39	-26.852%
random_char_delete	2.0	39.61	-34.77	-53.181%	19.9	-32.62	-62.110%	51.687	-20.523	-28.422%
random_word_insert	2.0	63.6	-1.78	-2.723%	46.947	-5.573	-10.611%	70.26	-1.95	-2.701%
random_word_delete	2.0	59.43	-5.95	-9.101%	46.454	-6.066	-11.550%	68.911	-3.299	-4.569%
swap_syn_word_emb	2.0	54.64	-10.74	-16.427%	42.391	-10.129	-19.285%	62.466	-9.744	-13.494%
swap_syn_word_net	2.0	57.61	-7.77	-11.884%	44.292	-8.228	-15.667%	67.002	-5.208	-7.213%
random_word_swap	2.0	64.03	-1.35	-2.065%	49.428	-3.092	-5.888%	71.193	-1.017	-1.409%
corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	3.0	49.21	-16.17	-24.732%	32.509	-20.011	-38.101%	61.648	-10.562	-14.627%
typos	3.0	50.18	-15.2	-23.249%	34.958	-17.562	-33.439%	63.528	-8.682	-12.023%
keyboard	3.0	48.68	-16.7	-25.543%	33.09	-19.43	-36.966%	62.308	-9.902	-13.713%
spell_error	3.0	50.07	-15.31	-23.417%	35.165	-17.355	-33.045%	62.466	-9.744	-13.949%
random_char_insert	3.0	33.75	-31.63	-48.379%	21.506	-31.014	-59.052%	54.299	-17.911	-24.804%
random_char_replace	3.0	30.61	-34.77	-53.181%	19.9	-32.62	-62.110%	51.687	-20.523	-28.422%
random_char_swap	3.0	28.17	-37.21	-56.913%	18.031	-34.489	-65.667%	50.954	-21.256	-29.436%
random_char_delete	3.0	32.97	-32.41	-49.572%	20.822	-31.698	-60.354%	53.997	-18.213	-25.222%
random_word_insert	3.0	61.46	-3.92	-5.996%	42.439	-10.081	-19.194%	67.949	-4.261	-5.901%
random_word_delete	3.0	56.52	-8.86	-13.552%	42.924	-9.596	-18.271%	67.059	-5.151	-7.133%
swap_syn_word_emb	3.0	50.61	-14.77	-22.901%	37.001	-15.519	-29.549%	60.299	-11.911	-16.496%
swap_syn_word_net	3.0	53.92	-11.47	-17.528%	40.022	-12.498	-23.796%	63.37	-8.84	-12.242%
random_word_swap	3.0	63.19	-2.19	-3.350%	48.267	-4.253	-8.098%	70.418	-1.792	-2.482%
corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	3.0	49.21	-16.17	-24.732%	32.509	-20.011	-38.101%	61.648	-10.562	-14.627%
typos	3.0	50.18	-15.2	-23.249%	34.958	-17.562	-33.439%	63.528	-8.682	-12.023%
keyboard	3.0	48.68	-16.7	-25.543%	33.09	-19.43	-36.966%	62.308	-9.902	-13.713%
spell_error	3.0	50.07	-15.31	-23.417%	35.165	-17.355	-33.045%	62.466	-9.744	-13.949%
random_char_insert	3.0	33.75	-31.63	-48.379%	21.506	-31.014	-59.052%	54.299	-17.911	-24.804%
random_char_replace	3.0	30.61	-34.77	-53.181%	19.9	-32.62	-62.110%	51.687	-20.523	-28.422%
random_char_swap	3.0	28.17	-37.21	-56.913%	18.031	-34.489	-65.667%	50.954	-21.256	-29.436%
random_char_delete	3.0	32.97	-32.41	-49.572%	20.822	-31.698	-60.354%	53.997	-18.213	-25.222%
random_word_insert	3.0	61.46	-3.92	-5.996%	42.439	-10.081	-19.194%	67.949	-4.261	-5.901%
random_word_delete	3.0	56.52	-8.86	-13.552%	42.924	-9.596	-18.271%	67.059	-5.151	-7.133%
swap_syn_word_emb	3.0	50.61	-14.77	-22.901%	37.001	-15.519	-29.549%	60.299	-11.911	-16.496%
swap_syn_word_net	3.0	53.92	-11.47	-17.528%	40.022	-12.498	-23.796%	63.37	-8.84	-12.242%
random_word_swap	3.0	63.19	-2.19	-3.350%	48.267	-4.253	-8.098%	70.418	-1.792	-2.482%
corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	4.0	49.21	-16.17	-24.732%	32.509	-20.011	-38.101%	61.648	-10.562	-14.627%
typos	4.0	50.18	-15.2	-23.249%	34.958	-17.562	-33.439%	63.528	-8.682	-12.023%
keyboard	4.0	48.68	-16.7	-25.543%	33.09	-19.43	-36.966%	62.308	-9.902	-13.713%
spell_error	4.0	50.07	-15.31	-23.417%	35.165	-17.355	-33.045%	62.466	-9.744	-13.949%
random_char_insert	4.0	33.75	-31.63	-48.379%	21.506	-31.014	-59.052%	54.299	-17.911	-24.804%
random_char_replace	4.0	30.61	-34.77	-53.181%	19.9	-32.62	-62.110%	51.687	-20.523	-28.422%
random_char_swap	4.0	28.17	-37.21	-56.913%	18.031	-34.489	-65.667%	50.954	-21.256	-29.436%
random_char_delete	4.0	32.97	-32.41	-49.572%	20.822	-31.698	-60.354%	53.997	-18.213	-25.222%
random_word_insert	4.0	61.46	-3.92	-5.996%	42.439	-10.081	-19.194%	67.949	-4.261	-5.901%
random_word_delete	4.0	56.52	-8.86	-13.552%	42.924	-9.596	-18.271%	67.059	-5.151	-7.133%
swap_syn_word_emb	4.0	50.61	-14.77	-22.901%	37.001	-15.519	-29.549%	60.299	-11.911	-16.496%
swap_syn_word_net	4.0	53.92	-11.47	-17.528%	40.022	-12.498	-23.796%	63.37	-8.84	-12.242%
random_word_swap	4.0	63.19	-2.19	-3.350%	48.267	-4.253	-8.098%	70.418	-1.792	-2.482%

Table 33: Performance and decrease ratio of Hyperformer on CLIP-BART against text corruptions given severity 5.

corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	5.0	39.58	-25.8	-39.462%	24.948	-27.572	-52.497%	54.729	-17.481	-24.208%
punctuation	5.0	64.57	-0.81	-1.239%	51.304	-1.216	-2.316%	71.509	-0.701	-0.971%
typos	5.0	29.88	-35.5	-54.298%	19.924	-32.596	-62.065%	51.73	-20.48	-38.362%
keyboard	5.0	27.13	-38.25	-58.504%	18.898	-33.622	-64.017%	50.969	-21.241	-29.416%
spell_error	5.0	34.74	-30.64	-46.864%	24.67	-27.85	-53.027%	55.031	-17.179	-23.791%
random_char_insert	5.0	27.42	-37.96	-58.061%	18.747	-33.773	-64.305%	51.184	-21.026	-29.118%
random_char_replace	5.0	25.64	-39.74	-60.783%	18.286	-34.234	-65.183%	50.165	-22.045	-30.529%
random_char_swap	5.0	25.53	-39.85	-60.951%	18.58	-33.94	-64.623%	51.242	-20.968	-29.038%
random_char_delete	5.0	27.21	-38.17	-58.362%	17.825	-34.692	-66.061%	51.572	-20.638	-28.501%
to_passive	5.0	62.96	-2.42	-3.701%	47.209	-5.311	-10.112%	68.882	-3.328	-4.699%
to_active	5.0	64.51	-0.87	-1.331%	52.425	-0.095	-0.181%	72.14	-0.07	-0.097%
to_formal	5.0	64.84	-0.54	-0.826%	51.948	-0.572	-1.089%	71.681	-0.529	-0.733%
to_casual	5.0	62.6	-2.78	-4.252%	47.83	-4.69	-8.931%	70.174	-2.036	-2.820%
to_negative	5.0	63.74	-1.64	-2.508%	50.223	-2.297	-4.574%	70.854	-1.376	-1.906%
double_denial	5.0	65.27	-0.11	-0.168%	52.504	-0.016	-0.030%	72.198	-0.012	-0.017%
insert_adv	5.0	62.02	-3.36	-5.129%	45.203	-2.997	-5.706%	70.949	-1.261	-1.747%
append_ri	5.0	61.82	-3.56	-5.445%	46.446	-6.074	-11.565%	66.083	-6.127	-8.485%
random_word_insert	5.0	54.55	-10.43	-16.565%	34.576	-17.947	-34.166%	62.179	-10.031	-16.130%
drop_an	5.0	44.68	-20.71	-35.61%	33.567	-18.953	-36.088%	61.289	-10.921	-15.124%
drop_rand_one_an	5.0	56.79	-8.59	-13.139%	44.641	-7.879	-15.001%	70.533	-1.677	-2.323%
drop_vb	5.0	58.74	-6.64	-10.566%	41.779	-10.741	-20.451%	68.437	-3.773	-5.225%
drop_vbn	5.0	44.4	-20.98	-32.809%	30.601	-21.919	-41.734%	57.471	-14.739	-20.411%
only_an	5.0	31.7	-33.68	-51.34%	20.711	-31.809	-60.566%	50.911	-21.299	-29.495%
only_vb	5.0	50.58	-14.8	-22.637%	37.844	-14.676	-27.944%	65.164	-7.046	-9.757%
only_vbn	5.0	33.83	-31.55	-48.256%	22.794	-29.726	-56.600%	52.663	-19.547	-27.070%
drop_rand_one_vb	5.0	62.77	-2.61	-3.992%	47.822	-4.698	-8.946%	71.049	-1.161	-1.607%

Table 36: Performance and decrease ratio of Multiple Compacters on CLIP-BART against text corruptions given severity 1 to 4.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	1.0	53.92	-10.99	-16.931%	42.026	-10.724	-20.330%	55.964	-13.486	-19.419%			
typos	1.0	56.2	-8.71	-13.419%	44.125	-8.625	-16.351%	67.303	-21.477	-3.091%			
keyboard	1.0	54.99	-9.92	-15.283%	43.091	-9.659	-18.311%	67.375	-2.075	-2.988%			
spell_error	1.0	55.79	-9.12	-14.050%	44.244	-8.506	-16.125%	67.073	-2.377	-3.422%			
random_char_insert	1.0	50.54	-14.37	-22.138%	35.244	-17.506	-33.187%	62.236	-7.214	-10.387%			
random_char_replace	1.0	47.13	-17.78	-27.392%	32.104	-20.646	-39.140%	60.241	-9.209	-13.260%			
random_char_swap	1.0	40.85	-24.06	-37.067%	25.918	-26.832	-50.866%	55.792	-13.658	-19.667%			
random_char_delete	1.0	48.7	-16.21	-24.973%	33.439	-19.311	-36.606%	60.758	-8.692	-12.516%			
random_word_insert	1.0	64.94	-0.03	0.046%	50.692	-2.058	-3.902%	69.226	-0.224	-0.322%			
random_word_delete	1.0	64.79	-0.12	-0.185%	51.248	-1.502	-2.847%	68.982	-0.468	-0.673%			
swap_syn_word_emb	1.0	55.51	-9.4	-14.482%	44.912	-7.838	-14.859%	65.537	-4.113	-5.923%			
swap_syn_word_net	1.0	57.54	-7.37	-11.354%	45.103	-7.647	-14.498%	67.289	-2.161	-3.112%			
random_word_swap	1.0	64.82	-0.94	-0.139%	51.741	-1.009	-1.913%	69.327	-0.123	-0.177%			
corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	2.0	53.54	-11.37	-17.517%	38.361	-14.389	-27.278%	63.169	-6.281	-9.944%			
typos	2.0	54.82	-10.09	-15.545%	40.157	-12.593	-23.872%	64.361	-5.089	-7.328%			
keyboard	2.0	54.43	-10.48	-16.145%	39.736	-13.014	-24.671%	64.102	-5.348	-7.700%			
spell_error	2.0	54.95	-9.96	-15.344%	41.414	-11.336	-21.491%	64.274	-5.176	-7.452%			
random_char_insert	2.0	40.57	-24.34	-37.498%	25.791	-26.959	-51.107%	56.323	-13.127	-18.902%			
random_char_replace	2.0	36.4	-28.51	-43.922%	22.786	-29.964	-56.804%	54.787	-14.663	-21.113%			
random_char_swap	2.0	31.96	-32.95	-50.763%	19.423	-33.327	-63.180%	52.074	-17.376	-25.019%			
random_char_delete	2.0	38.32	-26.59	-40.964%	24.591	-28.159	-53.383%	56.538	-12.912	-18.592%			
random_word_insert	2.0	64.31	-0.6	-0.924%	47.233	-5.517	-10.458%	68.767	-0.683	-0.983%			
random_word_delete	2.0	59.42	-5.49	-8.458%	45.134	-7.616	-14.437%	67.375	-2.075	-2.988%			
swap_syn_word_emb	2.0	54.98	-9.93	-15.298%	42.391	-10.359	-19.637%	61.418	-8.032	-11.565%			
swap_syn_word_net	2.0	56.96	-7.95	-12.248%	42.884	-8.866	-18.703%	64.662	-4.788	-6.894%			
random_word_swap	2.0	63.68	-1.23	-1.895%	49.189	-3.561	-6.751%	68.968	-0.482	-0.694%			
corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	3.0	48.87	-16.04	-24.711%	31.595	-21.155	-40.105%	59.265	-10.185	-14.444%			
typos	3.0	50.18	-14.73	-22.693%	33.741	-19.009	-36.035%	61.016	-8.434	-12.144%			
keyboard	3.0	48.88	-16.03	-24.666%	32.446	-20.304	-38.492%	60.227	-9.223	-13.280%			
spell_error	3.0	49.82	-15.09	-22.348%	34.745	-18.015	-34.151%	60.657	-8.793	-12.640%			
random_char_insert	3.0	33.67	-31.24	-48.128%	21.482	-31.268	-59.276%	54.012	-15.438	-22.229%			
random_char_replace	3.0	30.79	-34.12	-52.565%	19.614	-33.136	-62.818%	53.064	-16.386	-23.593%			
random_char_swap	3.0	28.51	-36.4	-56.078%	16.903	-35.847	-67.957%	51.242	-18.208	-26.218%			
random_char_delete	3.0	32.61	-32.2	-49.761%	20.027	-32.723	-62.034%	52.892	-16.558	-23.841%			
random_word_insert	3.0	63.9	-1.01	-1.556%	44.945	-8.705	-16.502%	67.834	-1.616	-2.327%			
random_word_delete	3.0	56.79	-8.12	-12.510%	40.237	-12.513	-23.721%	65.164	-4.286	-6.171%			
swap_syn_word_emb	3.0	51.14	-13.77	-21.214%	36.834	-15.916	-30.172%	58.906	-10.544	-15.182%			
swap_syn_word_net	3.0	53.25	-11.66	-17.963%	37.772	-14.978	-28.394%	61.405	-7.845	-11.296%			
random_word_swap	3.0	63.16	-1.75	-2.696%	47.567	-5.183	-9.825%	68.494	-0.956	-1.376%			
corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	4.0	41.92	-22.69	-35.418%	25.934	-26.816	-50.836%	55.792	-13.658	-19.667%			
typos	4.0	41.66	-23.25	-35.190%	27.645	-25.146	-47.671%	57.873	-11.577	-18.096%			
keyboard	4.0	39.29	-25.99	-49.078%	25.688	-27.062	-51.303%	56.882	-12.568	-18.096%			
spell_error	4.0	41.32	-23.59	-36.343%	28.979	-23.771	-45.063%	57.012	-12.438	-17.910%			
random_char_insert	4.0	30.03	-34.88	-53.736%	18.795	-33.955	-64.370%	52.246	-17.204	-24.771%			
random_char_replace	4.0	28.09	-36.82	-56.725%	17.077	-36.673	-67.626%	51.328	-18.122	-26.094%			
random_char_swap	4.0	26.72	-38.19	-58.855%	15.599	-37.151	-70.429%	51.313	-18.137	-26.115%			
random_char_delete	4.0	29.35	-35.56	-54.784%	17.523	-35.227	-66.782%	52.017	-17.433	-25.102%			
random_word_insert	4.0	62.39	-2.52	-3.882%	42.073	-10.677	-20.240%	67.748	-1.702	-2.451%			
random_word_delete	4.0	50.93	-13.98	-21.538%	34.823	-17.927	-33.985%	63.404	-6.41	-9.239%			
swap_syn_word_emb	4.0	47.52	-17.39	-26.791%	34.568	-18.182	-34.468%	57.801	-11.649	-16.773%			
swap_syn_word_net	4.0	49.86	-15.05	-23.186%	34.878	-17.872	-33.880%	59.323	-10.127	-14.582%			
random_word_swap	4.0	62.66	-2.25	-3.466%	46.526	-6.224	-11.800%	68.265	-1.185	-1.707%			

Table 37: Performance and decrease ratio of Multiple Compacters on CLIP-BART against text corruptions given severity 5.

corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDr	Decrease	Decrease Ratio
ocr	5.0	38.61	-26.3	-40.518%	23.422	-29.328	-55.598%	53.538	-15.912	-22.911%
punctuation	5.0	64.59	-0.32	-0.493%	51.685	-1.065	-2.018%	68.997	-0.453	-0.653%
typos	5.0	29.7	-35.21	-54.244%	19.645	-33.105	-62.758%	51.486	-17.964	-25.867%
keyboard	5.0	27.11	-37.8	-58.234%	18.564	-34.186	-64.807%	50.768	-18.682	-26.900%
spell_error	5.0	34.81	-30.1	-46.322%	24.964	-27.786	-52.674%	54.385	-15.065	-21.692%
random_char_insert	5.0	27.63	-37.28	-57.433%	17.888	-34.862	-66.088%	51.17	-18.28	-26.321%
random_char_replace	5.0	26.65	-38.26	-58.943%	16.64	-36.11	-68.455%	50.782	-18.668	-26.879%
random_char_swap	5.0	26.26	-38.65	-59.544%	16.298	-36.452	-69.103%	52.045	-17.405	-25.061%
random_char_delete	5.0	27.34	-37.57	-57.880%	16.107	-36.643	-69.464%	51.356	-18.094	-26.053%
to-passive	5.0	62.46	-2.45	-3.774%	43.067	-9.683	-18.356%	67.317	-2.133	-3.071%
to-verb	5.0	64.12	-0.79	-1.217%	52.473	-0.277	-0.526%	69.241	-0.209	-0.301%
to-formal	5.0	64.44	-0.47	-0.724%	51.932	-0.818	-1.551%	69.628	0.178	0.257%
to-casual	5.0	62.14	-2.77	-4.267%	45.707	-7.043	-13.352%	67.877	-1.573	-2.265%
to-active	5.0	63.37	-1.54	-2.373%	48.887	-3.863	-7.323%	68.753	-0.697	-1.004%
double-denial	5.0	64.82	-0.09	-0.139%	52.719	-0.031	-0.059%	69.442	-0.008	-0.012%
insert_adv	5.0	61.69	-3.22	-4.961%	40.443	-3.307	-6.268%	68.078	-1.372	-1.975%
append_ri	5.0	61.1	-3.81	-5.870%	46.995	-5.755	-10.910%	63.585	-5.865	-8.444%
random_word_insert	5.0	60.19	-4.72	-7.722%	37.216	-15.534	-29.449%	67.059	-2.391	-3.443%
drop_an	5.0	46.48	-18.43	-29.993%	32.485	-20.265	-38.417%	63.643	-5.807	-8.362%
drop_main_one_an	5.0	61.59	-8.01	-12.340%	44.467	-8.283	-14.703%	68.638	-8.012	-11.669%
drop_vb	5.0	60.37	-4.54	-6.994%	34.911	-17.84	-33.820%	67.906	-1.544	-2.223%
drop_vb_an	5.0	47.84	-17.07	-26.298%	23.271	-29.479	-55.885%	61.964	-7.486	-10.780%
only_an	5.0	35.98	-28.93	-44.569%	16.314	-36.436	-69.073%	57.758	-11.692	-16.835%
only_vb	5.0	52.59	-12.32	-18.806%	38.890	-13.841	-26.238%	64.906	-4.544	-6.543%
only_vb_an	5.0	41.42	-23.49	-36.189%	21.515	-21.235	-44.025%	60.442	-8.008	-12.937%
drop_main_one_vb	5.0	63.62	-1.29	-1.987%	41.907	-10.843	-20.556%	68.236	-1.214	-1.748%
drop_first	5.0	58.62	-6.29	-9.690%	37.971	-14.779	-28.017%	68.322	-1.128	-1.624%
drop_last	5.0	54.76	-10.15	-15.637%	46.641	-6.709	-12.719%	67.777	-1.673	-2.409%
drop_first_and_last	5.0	45.91	-19.0	-29.271%	33.201	-19.549	-37.060%	66.298	-3.152	-4.538%
shuffle_word	5.0	57.91	-7.0	-10.784%	39.362	-13.388	-25.379%	64.992	-4.458	-6.419%
random_word_delete	5.0	45.22	-19.69	-30.334%	29.059	-23.691	-61.912%	61.72	-7.73	-11.31%
swap_syn_word_emb	5.0	46.64	-18.27	-28.147%	33.527	-19.223	-36.442%	57.17	-12.28	-17.682%
swap_syn_word_net	5.0	46.98	-17.93	-27.623%	33.487	-19.263	-36.517%	57.227	-12.223	-17.600%
back_trains	5.0	61.57	-3.34	-5.146%	49.149	-3.601	-6.826%	68.451	-0.999	-1.438%
random_word_swap	5.0	61.73	-3.35	-4.899%	45.047	-7.703	-14.603%	66.628	-2.822	-4.063%
token_shuffle	5.0	61.95	-1.96	-7.806%	40.823	-4.277	-9.234%	51.012	-1.848	-26.549%

Table 40: Performance and decrease ratio of Multiple Prompts on CLIP-BART against text corruptions given severity 1 to 4.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	1.0	41.44	-5.37	-11.472%	29.758	-4.252	-12.501%	50.251	0.381	0.764%			
typos	1.0	39.78	-7.07	-15.018%	29.599	-4.411	-12.969%	49.978	0.108	0.218%			
keyboard	1.0	39.47	-7.34	-15.806%	29.48	-4.53	-13.191%	49.978	0.108	0.218%			
spell_error	1.0	39.38	-7.43	-15.873%	30.259	-3.751	-11.029%	49.921	0.051	0.102%			
random_char_insert	1.0	38.69	-8.12	-17.347%	26.244	-7.766	-22.834%	50.452	0.582	1.167%			
random_char_replace	1.0	37.03	-9.78	-20.893%	24.575	-9.435	-27.743%	50.174	0.884	1.772%			
random_char_swap	1.0	33.51	-13.3	-28.413%	21.387	-12.623	-37.117%	50.194	0.324	0.649%			
random_char_delete	1.0	37.19	-9.62	-20.515%	24.36	-9.65	-28.374%	49.878	0.008	0.016%			
random_word_insert	1.0	46.83	-0.02	0.043%	33.455	-0.555	-1.631%	49.892	0.022	0.045%			
random_word_delete	1.0	46.69	-0.12	-0.256%	33.59	-0.42	-1.234%	49.964	0.094	0.189%			
swap_syn_word_zemb	1.0	43.73	-3.08	-6.580%	31.468	-2.542	-7.475%	49.921	0.051	0.102%			
swap_syn_word_net	1.0	44.91	-1.9	-4.059%	30.752	-3.258	-9.579%	49.964	0.094	0.189%			
random_word_swap	1.0	46.77	-0.04	-0.085%	33.821	-0.189	-0.556%	50.036	0.166	0.333%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	2.0	41.17	-5.64	-12.049%	28.462	-5.548	-16.312%	50.079	0.209	0.419%			
typos	2.0	39.39	-7.97	-15.851%	28.423	-5.587	-16.429%	49.95	0.08	0.160%			
keyboard	2.0	38.83	-7.98	-17.048%	28.089	-5.921	-17.410%	50.495	0.625	1.254%			
spell_error	2.0	38.69	-7.94	-16.962%	28.669	-5.341	-15.704%	49.433	-0.437	-0.876%			
random_char_insert	2.0	33.92	-12.89	-27.537%	21.657	-12.353	-36.322%	50.409	0.539	1.081%			
random_char_replace	2.0	31.53	-15.28	-32.643%	20.178	-13.832	-40.670%	50.725	0.855	1.714%			
random_char_swap	2.0	28.86	-17.95	-38.347%	18.357	-15.653	-46.023%	50.022	0.152	0.304%			
random_char_delete	2.0	31.78	-15.03	-32.109%	20.536	-13.474	-39.618%	50.194	0.324	0.649%			
random_word_insert	2.0	45.83	-0.98	-2.094%	31.953	-2.057	-6.049%	50.079	0.209	0.419%			
random_word_delete	2.0	42.35	-4.46	-9.528%	31.102	-2.908	-8.551%	49.734	-0.136	-0.272%			
swap_syn_word_zemb	2.0	43.35	-3.46	-7.392%	29.416	-4.594	-13.506%	49.935	0.065	0.131%			
swap_syn_word_net	2.0	44.62	-2.19	-4.678%	29.814	-4.196	-12.338%	49.763	-0.107	-0.214%			
random_word_swap	2.0	46.18	-0.63	-1.346%	33.479	-0.531	-1.561%	49.907	0.037	0.074%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	3.0	38.72	-8.09	-17.283%	24.781	-9.229	-27.135%	49.878	0.008	0.016%			
typos	3.0	36.34	-10.47	-22.367%	24.805	-9.205	-27.065%	50.05	0.18	0.361%			
keyboard	3.0	35.76	-11.05	-23.666%	24.495	-9.515	-27.977%	49.935	0.065	0.131%			
spell_error	3.0	35.55	-11.26	-24.055%	25.258	-8.752	-29.922%	49.902	0.022	0.045%			
random_char_insert	3.0	30.25	-16.56	-35.377%	19.645	-14.365	-42.236%	49.821	-0.049	-0.099%			
random_char_replace	3.0	28.76	-18.05	-38.560%	18.898	-15.112	-44.434%	50.266	0.396	0.793%			
random_char_swap	3.0	27.14	-19.67	-42.021%	17.602	-16.408	-48.244%	49.763	-0.107	-0.214%			
random_char_delete	3.0	28.88	-17.93	-38.304%	18.524	-15.486	-45.532%	50.366	0.096	0.095%			
random_word_insert	3.0	45.27	-1.54	-3.290%	30.649	-3.361	-9.883%	49.51	0.64	1.282%			
random_word_delete	3.0	40.34	-6.47	-13.822%	29.377	-4.633	-13.623%	49.993	0.123	0.246%			
swap_syn_word_zemb	3.0	41.79	-5.02	-10.724%	26.586	-7.424	-21.829%	49.978	0.108	0.218%			
swap_syn_word_net	3.0	43.28	-3.33	-7.541%	27.397	-6.613	-19.444%	50.108	0.238	0.477%			
random_word_swap	3.0	45.65	-1.16	-2.478%	33.018	-0.992	-2.917%	50.122	0.252	0.505%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	4.0	34.95	-11.86	-25.336%	21.919	-12.091	-35.551%	50.38	0.51	1.023%			
typos	4.0	32.05	-14.76	-32.523%	21.657	-12.353	-36.322%	50.782	0.912	1.829%			
keyboard	4.0	30.97	-15.84	-33.839%	21.076	-12.934	-38.029%	50.366	0.496	0.995%			
spell_error	4.0	31.32	-15.49	-33.091%	21.903	-12.107	-35.597%	49.964	0.094	0.189%			
random_char_insert	4.0	28.93	-17.88	-38.197%	18.787	-15.223	-44.761%	50.309	0.439	0.879%			
random_char_replace	4.0	27.32	-19.49	-41.636%	18.278	-15.732	-46.257%	49.964	0.094	0.189%			
random_char_swap	4.0	26.13	-20.68	-44.179%	17.332	-16.678	-49.039%	49.763	-0.107	-0.214%			
random_char_delete	4.0	27.23	-19.58	-41.829%	17.65	-16.36	-48.104%	50.136	0.266	0.534%			
random_word_insert	4.0	44.37	-2.44	-5.213%	29.138	-4.872	-14.325%	50.954	1.084	2.175%			
random_word_delete	4.0	36.63	-10.18	-21.747%	27.325	-6.085	-19.655%	49.835	-0.035	-0.070%			
swap_syn_word_zemb	4.0	40.46	-6.35	-13.565%	25.473	-8.537	-25.101%	49.821	-0.049	-0.099%			
swap_syn_word_net	4.0	42.49	-4.32	-9.229%	26.053	-7.957	-23.395%	50.395	0.525	1.052%			
random_word_swap	4.0	45.21	-1.6	-3.418%	32.819	-1.191	-3.501%	49.864	-0.006	-0.013%			

Table 41: Performance and decrease ratio of Multiple Prompts on CLIP-BART against text corruptions given severity 5.

corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDr	Decrease	Decrease Ratio
ocr	5.0	32.66	-14.15	-30.229%	20.496	-13.514	-39.735%	50.093	0.223	0.448%
punctuation	5.0	45.08	-1.73	-3.666%	33.344	-0.666	-1.958%	50.065	0.195	0.390%
typos	5.0	26.63	-20.18	-43.110%	16.855	-17.155	-50.442%	49.849	-0.021	-0.042%
keyboard	5.0	25.78	-21.03	-44.926%	15.789	-18.221	-53.574%	50.553	0.683	1.369%
spell_error	5.0	28.07	-18.74	-40.034%	19.28	-14.73	-43.312%	50.596	0.726	1.455%
random_char_insert	5.0	27.32	-19.49	-41.636%	18.19	-15.82	-46.514%	49.835	-0.035	-0.070%
random_char_replace	5.0	26.41	-20.4	-43.589%	17.721	-16.289	-47.894%	50.395	0.525	1.052%
random_char_swap	5.0	25.71	-21.1	-45.076%	18.016	-15.994	-47.029%	49.519	-0.351	-0.704%
random_char_delete	5.0	26.47	-20.34	-43.452%	17.515	-16.495	-48.501%	49.95	0.08	0.160%
to-passive	5.0	43.64	-3.17	-6.772%	27.508	-6.502	-19.117%	49.706	-0.164	-0.329%
to-tense	5.0	46.64	-0.17	-0.363%	33.837	-0.173	-0.509%	49.821	-0.049	-0.099%
to-formal	5.0	46.62	-0.19	-0.406%	33.718	-0.292	-0.860%	49.95	0.08	0.160%
to-casual	5.0	43.71	-3.17	-6.623%	28.073	-5.937	-17.457%	49.978	0.108	0.218%
to-active	5.0	44.44	-2.37	-5.003%	31.778	-2.232	-6.564%	49.734	-0.136	-0.272%
double-denial	5.0	46.79	-0.02	-0.043%	34.012	-0.002	-0.005%	49.878	0.008	0.016%
insert_adv	5.0	45.23	-1.58	-3.375%	32.931	-1.079	-3.174%	49.964	0.094	0.189%
insert_adj	5.0	35.62	-11.19	-23.905%	27.699	-6.311	-18.556%	50.165	0.295	0.592%
random_word_insert	5.0	43.28	-3.53	-7.541%	27.5	-6.51	-19.140%	50.452	0.582	1.167%
drop_an	5.0	40.11	-6.7	-14.133%	24.789	-9.221	-27.112%	50.28	0.41	0.822%
drop_rand_one_an	5.0	43.3	-3.33	-7.028%	24.366	-9.282	-36.046%	50.036	0.166	0.333%
drop_vb	5.0	44.87	-1.94	-4.144%	27.728	-6.772	-19.912%	50.022	0.152	0.304%
drop_vb_an	5.0	41.13	-5.68	-12.134%	22.333	-11.677	-34.335%	50.151	0.281	0.563%
only_an	5.0	27.59	-19.22	-41.060%	12.261	-11.749	-34.545%	50.251	0.381	0.764%
only_vb	5.0	38.75	-8.06	-17.129%	27.174	-6.836	-20.099%	50.28	0.41	0.822%
only_vb_an	5.0	26.41	-20.4	-43.589%	24.366	-9.282	-36.046%	50.036	0.166	0.333%
drop_rand_one_vb	5.0	46.18	-0.63	-1.346%	29.909	-4.101	-12.057%	49.935	0.065	0.131%
drop_first	5.0	27.59	-19.22	-41.060%	12.266	-11.749	-34.545%	50.288	0.338	0.678%
drop_last	5.0	45.33	-1.48	-3.162%	31.778	-2.232	-6.564%	50.165	0.295	0.592%
drop_first_and_last	5.0	27.17	-19.11	-40.825%	26.467	-7.543	-22.179%	49.495	0.625	1.254%
shuffle_order	5.0	40.81	-6.0	-12.818%	26.001	-3.409	-10.023%	49.993	0.123	0.246%
random_word_delete	5.0	33.65	-13.16	-28.114%	24.813	-9.197	-27.042%	49.548	-0.322	-0.646%
swap_syn_word_zemb	5.0	40.05	-6.76	-14.411%	24.924	-9.086	-26.714%	49.878	0.008	0.016%
swap_syn_word_net	5.0	39.41	-7.4	-15.809%	24.535	-9.475	-27.860%	50.366	0.496	0.995%
backtrans	5.0	45.95	-0.86	-1.837%	32.422	-1.588	-4.670%	49.935	-0.035	-0.070%
random_word_swap	5.0	44.71	-2.1	-4.486%	32.477	-1.533	-4.506%	49.706	-0.164	-0.329%
nonsense	5.0	24.93	-21.82	-46.742%	17.427	-16.583	-48.758%	49.404	-0.466	-0.934%

Table 44: Performance and decrease ratio of Single Com-
pacter on CLIP-BART against text corruptions given sever-
ity 1 to 4.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	1.0	54.28	-10.19	-15.806%	41.779	-11.121	-21.022%	56.122	-13.818	-19.757%			
typos	1.0	56.27	-8.2	-12.719%	43.171	-9.729	-18.392%	68.236	-1.704	-2.436%			
keyboard	1.0	55.05	-9.42	-14.611%	42.145	-10.755	-20.331%	67.447	-2.493	-3.565%			
spell_error	1.0	55.9	-8.57	-13.203%	42.892	-10.008	-18.918%	67.92	-2.02	-2.888%			
random_char_insert	1.0	51.37	-13.1	-20.320%	35.427	-17.473	-33.030%	64.633	-5.307	-7.588%			
random_char_replace	1.0	47.83	-16.64	-25.810%	33.002	-19.898	-37.614%	62.222	-7.718	-11.035%			
random_char_swap	1.0	41.69	-22.78	-35.344%	27.962	-24.938	-47.143%	57.815	-12.125	-17.336%			
random_char_delete	1.0	49.17	-15.3	-23.722%	34.743	-18.157	-34.323%	62.452	-7.488	-10.707%			
random_word_insert	1.0	64.38	-0.09	-0.140%	48.322	-4.578	-8.653%	70.418	0.478	0.683%			
random_word_delete	1.0	64.36	-0.11	-0.171%	48.593	-4.307	-8.142%	69.427	-0.513	-0.733%			
swap_syn_word_emb	1.0	54.63	-9.84	-15.263%	42.741	-10.159	-19.204%	65.767	-4.173	-5.966%			
swap_syn_word_net	1.0	56.92	-7.55	-11.711%	43.918	-8.982	-16.979%	67.682	-2.278	-3.257%			
random_word_swap	1.0	64.34	-0.13	-0.202%	49.014	-3.886	-7.346%	69.686	-0.254	-0.364%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	2.0	53.82	-10.65	-16.519%	39.251	-13.649	-25.801%	63.657	-6.283	-8.983%			
typos	2.0	55.05	-9.42	-14.611%	40.142	-12.758	-24.118%	64.92	-5.02	-7.177%			
keyboard	2.0	54.56	-9.91	-15.371%	38.893	-14.007	-26.478%	64.992	-4.948	-7.074%			
spell_error	2.0	55.03	-9.44	-14.642%	40.626	-12.274	-23.201%	64.432	-5.508	-7.875%			
random_char_insert	2.0	41.58	-22.89	-35.505%	27.238	-25.662	-48.510%	59.367	-10.373	-14.832%			
random_char_replace	2.0	36.69	-27.78	-43.009%	24.058	-28.842	-54.522%	56.997	-13.043	-18.609%			
random_char_swap	2.0	32.41	-30.47	-49.729%	22.229	-30.671	-57.979%	54.93	-15.01	-21.461%			
random_char_delete	2.0	38.9	-25.57	-39.662%	26.697	-26.203	-49.532%	57.126	-12.814	-18.321%			
random_word_insert	2.0	63.54	-0.03	-0.443%	45.516	-7.384	-13.958%	70.332	0.392	0.560%			
random_word_delete	2.0	59.56	-4.91	-7.616%	43.711	-9.189	-17.370%	67.82	-2.12	-3.032%			
swap_syn_word_emb	2.0	53.58	-10.89	-16.892%	39.513	-13.387	-25.305%	61.49	-8.45	-12.082%			
swap_syn_word_net	2.0	56.45	-8.02	-12.440%	41.803	-11.097	-20.977%	65.207	-4.733	-6.767%			
random_word_swap	2.0	63.41	-1.06	-1.644%	47.201	-5.699	-10.772%	69.04	-0.9	-1.287%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	3.0	49.39	-15.08	-23.391%	33.455	-19.445	-36.758%	59.739	-10.201	-14.586%			
typos	3.0	50.46	-14.01	-21.731%	34.576	-18.324	-34.638%	62.81	-7.13	-10.194%			
keyboard	3.0	49.29	-15.18	-23.546%	32.724	-20.176	-38.140%	62.078	-7.862	-11.241%			
spell_error	3.0	50.28	-14.19	-22.019%	34.195	-18.705	-35.360%	61.849	-6.991	-11.509%			
random_char_insert	3.0	34.67	-29.8	-46.223%	23.263	-29.637	-56.025%	56.524	-13.416	-19.183%			
random_char_replace	3.0	30.24	-34.23	-53.004%	20.687	-32.213	-60.894%	54.83	-15.11	-21.604%			
random_char_swap	3.0	28.78	-35.69	-55.359%	20.329	-32.571	-61.571%	52.232	-17.088	-25.319%			
random_char_delete	3.0	33.24	-31.23	-48.441%	22.18	-29.502	-55.920%	54.414	-15.526	-22.200%			
random_word_insert	3.0	62.35	-2.12	-3.288%	42.407	-10.493	-19.835%	69.255	-0.685	-0.979%			
random_word_delete	3.0	57.35	-7.12	-11.044%	44.817	-12.083	-22.841%	65.494	-4.446	-6.356%			
swap_syn_word_emb	3.0	50.21	-14.26	-22.119%	34.179	-18.721	-35.390%	59.265	-10.675	-15.263%			
swap_syn_word_net	3.0	53.02	-11.17	-17.760%	38.297	-14.603	-27.605%	62.193	-7.747	-11.076%			
random_word_swap	3.0	62.96	-1.51	-2.342%	46.224	-6.676	-12.621%	69.255	-0.685	-0.979%			

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	4.0	42.76	-21.71	-33.675%	29.067	-23.833	-45.054%	56.28	-13.66	-19.532%			
typos	4.0	42.99	-21.98	-34.093%	28.033	-24.867	-47.007%	58.047	-11.263	-16.104%			
keyboard	4.0	42.41	-22.06	-34.117%	28.073	-24.827	-46.932%	58.16	-11.78	-16.843%			
spell_error	4.0	40.63	-33.84	-52.490%	20.719	-32.181	-60.834%	54.988	-14.952	-21.379%			
random_char_insert	4.0	26.52	-37.95	-58.865%	18.079	-34.821	-65.824%	52.978	-16.962	-24.252%			
random_char_replace	4.0	26.33	-38.14	-59.159%	19.288	-33.612	-63.539%	53.104	-16.746	-23.944%			
random_char_swap	4.0	29.97	-34.5	-53.513%	20.981	-31.919	-60.338%	53.323	-16.617	-23.759%			
random_char_delete	4.0	60.56	-3.91	-6.065%	40.237	-12.663	-23.938%	68.365	-1.575	-2.252%			
random_word_insert	4.0	51.55	-12.92	-20.400%	35.077	-17.823	-33.692%	63.428	-6.512	-9.311%			
random_word_delete	4.0	46.5	-17.97	-27.873%	31.356	-21.544	-40.762%	58.102	-11.838	-16.925%			
swap_syn_word_emb	4.0	40.65	-14.82	-22.987%	35.952	-16.948	-32.036%	59.251	-10.689	-15.283%			
swap_syn_word_net	4.0	46.62	-1.85	-2.870%	45.277	-7.623	-14.409%	68.164	-1.776	-2.539%			
random_word_swap	4.0	62.96	-1.51	-2.342%	46.224	-6.676	-12.621%	69.255	-0.685	-0.979%			

Table 45: Performance and decrease ratio of Single Com-
pacter on CLIP-BART against text corruptions given sever-
ity 5.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	5.0	40.14	-24.33	-37.738%	26.475	-26.425	-49.953%	54.988	-14.952	-21.379%			
punctuation	5.0	64.18	-0.29	-0.449%	48.768	-4.132	-7.812%	69.858	-0.082	-0.117%			
typos	5.0	30.66	-33.81	-52.443%	20.83	-32.07	-60.624%	52.146	-17.794	-25.442%			
keyboard	5.0	27.91	-36.56	-56.709%	18.922	-33.978	-64.231%	51.931	-18.009	-25.750%			
spell_error	5.0	36.05	-28.42	-44.083%	24.09	-28.81	-54.462%	55.461	-14.479	-20.701%			
random_char_insert	5.0	27.56	-36.91	-57.251%	19.447	-33.453	-63.239%	53.625	-16.115	-23.041%			
random_char_replace	5.0	24.0	-40.47	-62.773%	16.529	-36.371	-68.755%	52.246	-17.694	-25.298%			
random_char_swap	5.0	25.68	-38.99	-60.168%	20.122	-32.778	-61.961%	52.863	-17.077	-24.416%			
random_char_delete	5.0	27.65	-36.82	-57.112%	19.86	-33.04	-62.457%	52.347	-17.593	-25.155%			
to-passive	5.0	59.84	-4.63	-7.182%	44.761	-8.139	-15.386%	68.638	-1.302	-1.862%			
tense	5.0	63.67	-0.8	-1.241%	49.459	-3.441	-6.504%	70.159	0.219	0.314%			
to-formal	5.0	63.94	-0.53	-0.822%	48.879	-4.021	-7.601%	69.973	0.033	0.047%			
to-casual	5.0	60.54	-3.93	-6.066%	45.572	-7.328	-13.853%	68.652	-1.288	-1.841%			
to-active	5.0	62.53	-1.94	-3.009%	47.798	-5.102	-9.645%	69.528	-0.412	-0.589%			
double-denial	5.0	64.43	-0.04	-0.062%	49.754	-3.146	-5.948%	69.944	0.004	0.006%			
insert_adv	5.0	60.74	-3.73	-5.786%	45.588	-7.312	-13.823%	68.997	-0.943	-1.349%			
append_ri	5.0	59.21	-5.26	-8.159%	43.171	-9.729	-18.392%	63.715	-6.225	-8.901%			
random_word_insert	5.0	56.42	-0.85	-1.2486%	36.484	-16.416	-31.031%	67.906	-2.034	-2.908%			
drop_an	5.0	47.12	-17.35	-26.912%	32.732	-20.168	-38.125%	63.815	-6.125	-8.757%			
drop_random_one_an	5.0	57.7	-10.019	-16.568%	43.568	-13.241	-17.641%	69.244	-3.324	-4.191%			
drop_vb	5.0	58.63	-5.84	-9.058%	42.447	-10.453	-19.760%	68.724	-1.216	-1.739%			
drop_vb_an	5.0	46.33	-18.14	-28.137%	30.72	-22.18	-41.928%	61.045	-8.895	-12.718%			
only_an	5.0	37.29	-27.18	-42.159%	23.072	-29.828	-56.386%	53.567	-16.373	-23.410%			
only_vb	5.0	51.43	-12.68	-19.668%	36.19	-16.71	-31.588%	63.673	-6.067	-8.675%			
only_vb_an	5.0	41.8	-22.76	-35.449%	25.783	-27.117	-51.261%	56.581	-13.519	-19.101%			
drop_random_one_vb	5.0	62.38	-2.09	-3.242%	46.812	-6.083	-11.509%	69.614	-0.326	-0.466%			
drop_first	5.0	59.78	-4.69	-7.725%	43.075	-9.825	-18.572%	69.772	-0.168	-0.241%			
drop_first_and_last	5.0	54.38	-10.09	-15.651%	43.306	-9.594	-18.137%	67.891	-2.049	-2.922%			
drop_first_and_last_shuffle_order	5.0	47.38	-17.09	-26.508%	37.224	-15.676	-29.634%	66.987	-2.953	-4.222%			
random_word_delete	5.0	59.98	-4.40	-6.964%	39.378	-13.522	-25.561%	66.14	-3.8	-5.433%			
swap_syn_word_emb	5.0	46.11	-18.36	-28.748%	31.332	-21.568	-40.770%	62.696	-7.244	-10.358%			
swap_syn_word_net	5.0	45.82	-18.65	-28.928%	30.529	-22.371	-42.288%	57.844	-12.096	-17.295%			
backtrans	5.0	46.94	-17.57	-27.191%	34.433	-18.467	-34.909%	57.643	-12.297	-17.582%			

Table 48: Performance and decrease ratio of Single Prompt on CLIP-BART against text corruptions given severity 1 to 4.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	1.0	39.57	-4.43	-10.68%	33.789	-3.751	-9.92%	51.629	-0.321	-0.618%			
types	1.0	38.99	-5.01	-11.36%	33.447	-4.093	-10.90%	52.339	0.339	0.653%			
keyboard	1.0	38.12	-5.88	-13.364%	32.787	-4.753	-12.660%	51.844	-0.106	-0.203%			
spell_error	1.0	38.42	-5.58	-12.682%	33.686	-3.854	-10.267%	52.189	0.239	0.460%			
random_char_insert	1.0	37.28	-6.72	-15.273%	29.925	-7.615	-20.284%	51.672	-0.278	-0.535%			
random_char_replace	1.0	35.6	-8.4	-19.091%	27.89	-9.65	-25.706%	52.361	0.411	0.791%			
random_char_swap	1.0	32.3	-11.7	-26.591%	24.694	-12.846	-34.220%	52.002	0.052	0.101%			
random_char_delete	1.0	35.77	-8.23	-18.705%	28.621	-8.919	-23.758%	51.758	-0.192	-0.369%			
random_word_insert	1.0	43.95	-0.05	-0.114%	36.953	-0.587	-1.563%	52.246	0.296	0.570%			
random_word_delete	1.0	43.86	-0.14	-0.318%	36.468	-1.072	-2.854%	51.787	-0.163	-0.314%			
swap_syn_word_emb	1.0	39.76	-4.24	-9.636%	34.998	-2.542	-6.722%	52.088	0.138	0.266%			
swap_syn_word_net	1.0	41.92	-2.08	-4.727%	35.085	-2.455	-6.540%	51.988	0.038	0.073%			
random_word_swap	1.0	43.96	-0.04	-0.091%	37.279	-0.261	-0.694%	51.916	-0.034	-0.065%			
corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	2.0	39.83	-4.17	-9.477%	32.644	-4.896	-13.041%	52.189	0.239	0.460%			
types	2.0	38.31	-5.69	-12.923%	31.857	-5.683	-15.138%	52.131	0.181	0.349%			
keyboard	2.0	37.82	-6.18	-14.045%	31.205	-6.335	-16.875%	51.514	-0.436	-0.839%			
spell_error	2.0	38.19	-5.81	-13.205%	32.43	-5.11	-13.613%	51.902	-0.048	-0.093%			
random_char_insert	2.0	33.14	-10.86	-24.682%	25.219	-12.321	-32.822%	51.945	-0.005	-0.010%			
random_char_replace	2.0	30.73	-13.27	-30.159%	23.231	-14.309	-38.117%	51.959	0.009	0.018%			
random_char_swap	2.0	28.42	-15.58	-35.409%	21.061	-16.479	-43.898%	52.131	0.181	0.349%			
random_char_delete	2.0	30.9	-13.1	-29.773%	23.875	-13.665	-36.401%	51.801	-0.149	-0.286%			
random_word_insert	2.0	42.38	-0.62	-1.409%	35.538	-2.002	-5.332%	52.175	0.225	0.432%			
random_word_delete	2.0	39.26	-4.74	-10.773%	32.724	-4.816	-12.830%	52.103	0.153	0.294%			
swap_syn_word_emb	2.0	39.47	-4.53	-10.205%	33.376	-4.164	-11.093%	51.787	-0.163	-0.314%			
swap_syn_word_net	2.0	41.73	-2.27	-5.159%	34.474	-3.266	-8.700%	52.49	0.54	1.040%			
random_word_swap	2.0	43.28	-0.72	-1.636%	36.445	-1.095	-2.918%	51.801	-0.149	-0.286%			
corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	3.0	38.04	-5.96	-13.545%	29.122	-8.418	-22.423%	52.361	0.411	0.791%			
types	3.0	36.97	-7.98	-18.136%	27.977	-9.563	-25.473%	52.045	0.395	0.754%			
keyboard	3.0	35.37	-8.63	-19.641%	27.373	-10.167	-27.083%	51.801	-0.149	-0.286%			
spell_error	3.0	35.73	-8.27	-18.795%	28.494	-9.046	-24.096%	52.017	0.067	0.128%			
random_char_insert	3.0	30.11	-13.89	-31.568%	23.462	-14.078	-37.502%	52.074	0.124	0.239%			
random_char_replace	3.0	28.27	-15.73	-35.750%	21.045	-16.495	-43.941%	51.902	-0.048	-0.093%			
random_char_swap	3.0	26.64	-17.26	-39.455%	19.661	-17.879	-47.626%	51.974	0.024	0.045%			
random_char_delete	3.0	28.16	-15.84	-36.000%	21.466	-16.074	-42.818%	52.734	0.784	1.510%			
random_word_insert	3.0	42.88	-1.12	-2.545%	34.409	-3.131	-8.340%	51.773	-0.177	-0.341%			
random_word_delete	3.0	36.96	-7.04	-16.000%	29.114	-8.426	-22.445%	51.816	-0.134	-0.259%			
swap_syn_word_emb	3.0	40.15	-3.46	-8.591%	36.074	-4.764	-13.016%	51.816	-0.134	-0.259%			
swap_syn_word_net	3.0	40.58	-3.42	-7.737%	31.897	-5.643	-15.032%	51.945	-0.005	-0.010%			
random_word_swap	3.0	43.21	-0.79	-1.795%	35.475	-2.065	-5.502%	51.773	-0.177	-0.341%			
corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio
ocr	4.0	35.04	-8.96	-20.364%	26.268	-11.272	-30.026%	51.974	0.024	0.045%			
keyboard	4.0	30.96	-13.04	-29.636%	23.048	-14.492	-38.604%	51.816	-0.134	-0.259%			
spell_error	4.0	31.55	-12.45	-28.295%	25.123	-12.417	-33.076%	51.801	-0.149	-0.286%			
random_char_insert	4.0	28.57	-15.43	-35.068%	21.053	-16.487	-43.919%	51.615	-0.355	-0.645%			
random_char_replace	4.0	27.02	-16.98	-38.919%	20.114	-17.426	-46.419%	52.361	0.411	0.791%			
random_char_swap	4.0	25.94	-18.06	-40.145%	19.534	-18.006	-47.965%	51.715	-0.235	-0.452%			
random_char_delete	4.0	26.91	-17.09	-38.841%	19.827	-17.72	-47.202%	52.189	0.239	0.460%			
random_word_insert	4.0	42.33	-1.67	-3.795%	32.946	-4.594	-12.237%	51.859	-0.091	-0.176%			
random_word_delete	4.0	35.12	-11.13	-31.682%	26.123	-11.407	-30.386%	51.931	-0.019	-0.037%			
swap_syn_word_emb	4.0	35.94	-8.06	-18.181%	29.766	-7.774	-20.708%	51.814	-0.34	-0.259%			
swap_syn_word_net	4.0	39.88	-4.12	-9.364%	30.037	-7.503	-19.988%	51.672	-0.278	-0.535%			
random_word_swap	4.0	42.56	-1.44	-3.273%	34.815	-2.725	-7.260%	52.017	0.067	0.128%			

Table 49: Performance and decrease ratio of Single Prompt on CLIP-BART against text corruptions given severity 5.

corruption	severity	VQA			GQA			NLVR		
		Acc	Decrease	Decrease Ratio	Acc	Decrease	Decrease Ratio	CIDder	Decrease	Decrease Ratio
ocr	5.0	34.13	-9.87	-22.432%	24.654	-12.886	-34.326%	51.873	-0.077	-0.148%
punctuation	5.0	42.47	-1.53	-3.477%	34.751	-2.789	-7.429%	51.959	0.009	0.018%
types	5.0	26.09	-17.91	-40.705%	19.367	-18.173	-48.409%	51.744	-0.206	-0.397%
keyboard	5.0	24.86	-19.14	-43.800%	18.501	-19.039	-50.718%	51.83	-0.12	-0.231%
spell_error	5.0	28.3	-15.7	-35.682%	22.516	-15.024	-40.023%	51.859	-0.091	-0.176%
random_char_insert	5.0	27.07	-16.93	-38.877%	20.17	-17.37	-46.270%	51.543	-0.407	-0.783%
random_char_replace	5.0	25.81	-18.19	-41.341%	19.463	-18.077	-48.155%	52.246	0.296	0.570%
random_char_swap	5.0	25.38	-18.62	-42.318%	18.898	-18.642	-49.659%	51.572	-0.378	-0.728%
random_char_delete	5.0	25.94	-18.06	-41.045%	19.502	-18.038	-48.049%	52.189	0.239	0.460%
to-passive	5.0	40.84	-3.16	-7.182%	33.439	-4.101	-10.923%	52.419	0.469	0.902%
tense	5.0	43.8	-0.2	-0.455%	37.462	-0.078	-0.207%	51.816	-0.134	-0.259%
to-formal	5.0	43.79	-0.21	-0.477%	37.359	-0.181	-0.482%	52.289	0.339	0.653%
to-causal	5.0	41.04	-2.96	-6.727%	33.749	-3.791	-10.097%	52.318	0.368	0.709%
to-active	5.0	42.52	-1.48	-3.364%	36.071	-1.469	-3.913%	51.945	-0.005	-0.010%
double_denial	5.0	43.93	-0.07	-0.159%	37.542	0.002	0.005%	51.974	0.024	0.045%
insert_adv	5.0	42.61	-1.39	-3.159%	35.912	-1.628	-4.337%	52.304	0.354	0.681%
append_1st	5.0	35.1	-8.9	-20.227%	26.833	-10.707	-38.525%	52.017	0.067	0.128%
random_word_insert	5.0	41.52	-2.48	-5.636%	32.096	-5.444	-14.503%	51.443	-0.507	-0.977%
drop_random_1st	5.0	35.18	-8.82	-20.045%	27.874	-9.666	-25.748%	52.074	0.124	0.239%
drop_random_2nd	5.0	39.3	-4.7	-10.682%	33.415	-4.125	-10.987%	51.974	0.024	0.045%
drop_vb	5.0	25.14	-18.6	-42.864%	20.385	-17.155	-45.698%	52.318	0.368	0.709%
drop_vb_1st	5.0	25.92	-18.08	-41.091%	17.539	-20.001	-53.280%	52.462	0.512	0.985%
only_vb	5.0	16.21	-27.79	-63.159%	12.824	-24.716	-65.839%	51.514	-0.436	-0.839%
only_vb_1st	5.0	37.9	-6.1	-13.864%	30.712	-6.828	-18.188%	51.959	0.009	0.018%
only_vb_2nd	5.0	35.09	-8.97	-20.205%	28.327	-9.213	-24.541%	51.543	-0.407	-0.783%
drop_random_vb	5.0	26.64	-14.36	-32.636%	25.54	-11.9	-37.629%	51.816	-0.134	-0.259%
drop_first	5.0	25.03	-23.97	-54.477%	15.249	-22.291	-59.380%	51.76	0.21	0.405%
drop_last	5.0	42.14	-1.86	-4.227%	36.683	-0.857	-2.283%	52.203	0.253	0.487%
drop_first_and_last	5.0	17.78	-26.22	-59.591%	14.756	-22.784	-60.693%	52.218	0.268	0.515%
shuffle_order	5.0	40.09	-3.91	-8.886%	28.955	-5.585	-22.868%	52.304	0.354	0.681%
random_word_delete	5.0	29.11	-14.89	-33.841%	23.414	-14.126	-37.629%	51.959	0.009	0.018%
swap_syn_word_emb	5.0	35.43	-8.57	-18.723%	29.146	-8.394	-22.360%	52.002	0.052	0.101%
swap_syn_word_net	5.0	35.96	-8.04	-18.273%	28.6	-8.394	-22.360%	52.002	0.052	0.101%
swap_syn_word_syn	5.0	43.43	-0.67	-1.537%	37.957	-1.183	-3.115%	51.951	-0.011	-0.023%
random_word_swap	5.0	42.27	-1.73	-3.932%	33.956	-3.584	-9.547%	51.945	-0.005	-0.010%
nonsense	5.0	24.61	-19.39	-44.068%	18.0	-19.54	-52.052%	51.959	0.009	0.018%

Table 52: Performance and decrease ratio of Hyperformer on CLIP-T5 against text corruptions given severity 1 to 4.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		CIDder	Decrease	Decrease Ratio	
ocr	1.0	54.25	-10.93	-16.769%		44.705	-9.945	-18.198%		58.691	-11.869	-16.821%	
punctuation	1.0	64.64	-0.54	-0.828%		52.457	-2.193	-4.013%		70.647	0.087	0.124%	
typos	1.0	56.88	-8.3	-12.734%		47.05	-7.6	-13.906%		68.982	-1.578	-2.236%	
spell_Error	1.0	56.94	-8.24	-12.642%		46.573	-8.077	-14.779%		68.064	-2.496	-3.538%	
random_char_insert	1.0	50.88	-14.3	-21.939%		37.001	-17.649	-32.294%		64.762	-5.798	-8.216%	
random_char_replace	1.0	46.86	-18.32	-28.107%		33.765	-20.885	-38.215%		61.705	-8.855	-12.549%	
random_char_swap	1.0	40.92	-24.26	-37.220%		27.699	-26.951	-49.315%		58.174	-12.386	-17.554%	
random_char_delete	1.0	49.57	-15.61	-23.949%		36.723	-17.927	-32.804%		63.557	-7.003	-9.925%	
random_word_insert	1.0	65.13	-0.05	-0.077%		50.97	-3.68	-6.734%		70.188	-0.372	-0.527%	
random_word_delete	1.0	65.01	-0.17	-0.261%		52.409	-2.241	-4.101%		69.528	-1.032	-1.463%	
swap_syn_word_emb	1.0	54.91	-10.27	-15.756%		46.836	-7.814	-14.299%		65.868	-4.692	-6.650%	
swap_syn_word_net	1.0	58.65	-6.53	-10.018%		47.456	-7.194	-13.164%		68.58	-1.98	-2.805%	
random_word_swap	1.0	65.09	-0.09	-0.138%		52.767	-1.883	-3.446%		69.844	-0.716	-1.015%	

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		CIDder	Decrease	Decrease Ratio	
ocr	2.0	53.56	-11.62	-17.828%		40.444	-14.206	-25.995%		64.978	-5.582	-7.911%	
punctuation	2.0	64.71	-0.47	-0.721%		52.139	-2.511	-4.595%		70.274	-0.286	-0.405%	
typos	2.0	55.96	-9.22	-14.145%		42.36	-12.29	-22.489%		66.399	-4.161	-5.897%	
spell_Error	2.0	56.7	-8.48	-13.010%		44.093	-10.557	-19.318%		65.652	-4.908	-6.955%	
random_char_insert	2.0	40.34	-24.84	-38.110%		27.588	-27.062	-49.519%		58.49	-12.07	-17.106%	
random_char_replace	2.0	35.49	-29.69	-45.531%		24.249	-30.401	-55.629%		55.519	-15.041	-21.317%	
random_char_swap	2.0	31.74	-33.44	-51.304%		18.85	-35.8	-65.507%		52.605	-17.955	-25.446%	
random_char_delete	2.0	38.59	-26.59	-40.795%		26.578	-28.072	-51.367%		57.643	-12.917	-18.306%	
random_word_insert	2.0	63.76	-1.42	-2.179%		47.973	-6.677	-12.218%		69.14	-1.42	-2.012%	
random_word_delete	2.0	60.14	-5.04	-7.722%		47.265	-7.385	-13.513%		67.159	-3.401	-4.819%	
swap_syn_word_emb	2.0	54.15	-11.03	-16.922%		43.608	-11.042	-20.205%		62.007	-8.553	-12.122%	
swap_syn_word_net	2.0	58.08	-7.1	-10.893%		45.158	-9.492	-17.368%		66.212	-4.348	-6.162%	
random_word_swap	2.0	63.54	-1.64	-2.516%		50.7	-3.95	-7.228%		69.528	-1.032	-1.463%	

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		CIDder	Decrease	Decrease Ratio	
ocr	3.0	48.9	-16.28	-24.977%		33.908	-20.742	-37.954%		62.322	-8.238	-11.675%	
punctuation	3.0	64.59	-0.59	-0.905%		52.361	-2.289	-4.188%		70.159	-0.401	-0.568%	
typos	3.0	51.02	-14.16	-21.724%		35.864	-18.786	-34.375%		63.772	-6.788	-9.620%	
spell_Error	3.0	51.79	-13.39	-20.543%		37.995	-16.655	-30.476%		63.227	-7.333	-10.393%	
random_char_insert	3.0	33.48	-31.7	-48.635%		22.571	-32.079	-58.699%		54.816	-15.744	-22.314%	
random_char_replace	3.0	30.02	-35.16	-53.943%		19.097	-35.553	-65.056%		52.232	-18.328	-25.975%	
random_char_swap	3.0	27.97	-37.21	-57.088%		14.653	-39.997	-73.188%		51.328	-19.232	-27.257%	
random_char_delete	3.0	32.67	-32.51	-49.877%		20.798	-33.852	-61.943%		54.213	-16.347	-23.168%	
random_word_insert	3.0	63.1	-2.08	-3.191%		43.64	-11.01	-20.147%		69.112	-1.448	-2.053%	
random_word_delete	3.0	57.38	-7.8	-11.967%		43.314	-11.336	-20.743%		66.356	-4.204	-5.959%	
swap_syn_word_emb	3.0	50.18	-15.0	-23.013%		38.384	-16.266	-29.763%		58.375	-12.185	-17.269%	
swap_syn_word_net	3.0	54.48	-10.7	-16.416%		40.571	-14.079	-25.762%		63.341	-7.219	-10.230%	
random_word_swap	3.0	63.53	-2.76	-4.164%		49.801	-7.019	-12.353%		71.178	-2.882	-3.891%	

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		CIDder	Decrease	Decrease Ratio	
ocr	4.0	42.03	-23.15	-35.178%		28.717	-25.933	-47.453%		58.49	-12.07	-17.106%	
punctuation	4.0	61.62	-0.56	-0.859%		52.298	-2.352	-4.304%		70.518	-0.402	-0.559%	
typos	4.0	42.24	-22.94	-35.195%		29.997	-24.653	-45.111%		60.456	-10.104	-14.319%	
spell_Error	4.0	44.06	-21.12	-32.403%		32.318	-22.332	-40.863%		60.729	-9.831	-13.933%	
random_char_insert	4.0	29.46	-35.72	-54.802%		19.145	-35.505	-64.969%		52.361	-18.199	-25.792%	
random_char_replace	4.0	27.43	-37.75	-57.917%		16.243	-38.407	-70.279%		50.754	-19.806	-28.070%	
random_char_swap	4.0	26.3	-38.88	-59.650%		12.514	-42.136	-77.102%		50.897	-19.663	-27.867%	
random_char_delete	4.0	29.24	-35.94	-55.140%		15.837	-38.813	-71.021%		52.002	-18.558	-26.301%	
random_word_insert	4.0	61.76	-3.47	-5.247%		40.642	-14.008	-25.631%		68.997	-1.563	-2.162%	
random_word_delete	4.0	51.41	-13.77	-17.126%		37.311	-17.339	-31.727%		63.93	-6.63	-9.396%	
swap_syn_word_emb	4.0	56.37	-18.61	-28.859%		35.602	-19.048	-34.855%		56.61	-13.95	-19.711%	
swap_syn_word_net	4.0	51.49	-13.69	-21.003%		38.122	-16.528	-30.243%		60.643	-9.917	-14.055%	
random_word_swap	4.0	61.82	-3.36	-5.155%		47.71	-6.94	-12.698%		68.265	-2.295	-3.253%	

Table 53: Performance and decrease ratio of Hyperformer on CLIP-T5 against text corruptions given severity 5.

corruption	severity	VQA				GQA				NLVR			
		Acc	Decrease	Decrease Ratio		Acc	Decrease	Decrease Ratio		CIDr	Decrease	Decrease Ratio	
ocr	5.0	38.87	-26.31	-40.365%		26.109	-28.541	-52.225%		55.189	-15.371	-21.785%	
punctuation	5.0	64.59	-0.59	-0.905%		52.488	-2.162	-3.955%		70.375	-0.185	-0.263%	
typos	5.0	29.7	-35.48	-54.434%		20.639	-34.011	-62.234%		52.806	-17.754	-25.161%	
keybaord	5.0	26.47	-38.71	-59.389%		18.532	-36.118	-66.089%		51.758	-18.802	-26.646%	
spell_Error	5.0	37.8	-27.38	-42.007%		27.914	-26.736	-48.923%		57.844	-12.716	-18.021%	
random_char_insert	5.0	27.43	-37.75	-57.917%		17.189	-37.461	-68.548%		51.17	-19.39	-27.480%	
random_char_replace	5.0	26.18	-39.0	-59.834%		15.169	-39.481	-72.243%		50.84	-19.72	-27.948%	
random_char_swap	5.0	25.84	-39.34	-60.356%		12.403	-42.247	-77.305%		51.572	-18.988	-26.911%	
random_char_delete	5.0	27.25	-37.93	-58.193%		13.786	-40.864	-74.774%		51.945	-18.615	-26.382%	
to-passive	5.0	62.68	-2.5	-3.836%		50.318	-4.332	-7.927%		68.107	-2.453	-3.477%	
to-tense	5.0	64.34	-0.84	-1.289%		53.458	-1.192	-2.180%		70.304	-0.056	-0.080%	
to-formal	5.0	64.71	-0.47	-0.721%		52.687	-1.963	-3.592%		69.987	-0.573	-0.812%	
to-casual	5.0	62.48	-2.7	-4.142%		50.111	-4.539	-8.305%		69.485	-1.075	-1.524%	
to-active	5.0	63.79	-1.39	-2.133%		52.52	-2.13	-3.897%		69.614	-0.946	-1.341%	
double-denial	5.0	65.13	-0.05	-0.077%		53.602	-1.048	-1.919%		70.561	0.001	0.002%	
insert_adv	5.0	61.04	-4.14	-6.352%		50.62	-4.03	-7.374%		68.537	-2.023	-2.867%	
append_inr	5.0	62.95	-2.23	-3.421%		45.262	-9.388	-17.179%		64.016	-6.544	-9.274%	
random_word_insert	5.0	59.59	-5.59	-8.576%		37.534	-17.116	-31.320%		68.25	-2.31	-3.273%	
drop_an	5.0	45.81	-19.37	-29.718%		33.376	-21.274	-38.928%		63.054	-7.506	-10.637%	
drop_rand_one_an	5.0	56.58	-8.6	-13.104%		46.009	-8.641	-15.812%		68.193	-2.367	-3.355%	
drop_vb	5.0	57.72	-7.46	-11.445%		37.2	-17.45	-31.931%		68.781	-1.779	-2.521%	
drop_vb_an	5.0	46.33	-18.85	-28.920%		24.312	-30.338	-55.513%		61.547	-9.013	-12.773%	
only_an	5.0	33.76	-31.42	-48.205%		13.85	-40.8	-74.658%		53.165	-17.395	-24.653%	
only_vb	5.0	52.2	-12.98	-19.914%		39.521	-15.129	-27.683%		63.715	-6.845	-9.701%	
only_vb_an	5.0	41.64	-23.54	-36.115%		29.369	-25.281	-46.369%		55.993	-14.567	-20.645%	
drop_rand_one_vb	5.0	62.16	-3.02	-4.633%		45.421	-9.229	-16.888%		69.413	-1.147	-1.626%	
drop_first	5.0	58.75	-6.43	-9.865%		41.334	-13.316	-24.366%		69.384	-1.176	-1.666%	
drop_end_last	5.0	52.99	-12.19	-18.702%		47.042	-7.608	-13.920%		67.719	-2.841	-4.026%	
drop_first_and_last	5.0	46.58	-18.6	-28.366%		33.646	-21.004	-38.434%		66.959	-3.601	-5.104%	
shuffle_order	5.0	57.7	-7.48	-11.476%		39.426	-15.224	-27.857%		65.839	-4.721	-6.691%	
random_word_delete	5.0	46.23	-18.95	-29.073%		31.762	-22.888	-41.881%		62.265	-8.295	-11.756%	
swap_syn_word_emb	5.0	45.5	-19.68	-30.193%		34.568	-20.082	-36.746%		56.926	-13.634	-19.323%	
swap_syn_word_net	5.0	47.61	-17.57	-26.566%		36.572	-18.078	-33.080%		58.088	-12.472	-17.676%	
back_trans	5.0	61.99	-3.19	-4.894%		50.541	-4.109	-7.519%		69.987	-0.573	-0.812%	
random_word_swap	5.0	61.02	-4.16	-6.382%		46.414	-8.236	-15.070%		67.36	-3.2	-4.535%	
nonsense	5.0	15.93	-48.95	-76.481%		3.999	-56.051	-92.682%		50.768	-19.792	-28.050%	

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Madeline Chantry, Shruti Vyas, Hamid Palangi, Yogesh Rawat, and Vibhav Vineet. Robustness analysis of video-language models against visual and language perturbations. *Advances in Neural Information Processing Systems*, 35:34405–34420, 2022.
- [4] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015.
- [5] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022.
- [6] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. In *International Conference on Machine Learning*, pages 1931–1942. PMLR, 2021.
- [7] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers. *arXiv preprint arXiv:2202.04053*, 2022.
- [8] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 120–128, 2019.
- [9] Dorottya Demszky, Devyani Sharma, Jonathan H Clark, Vinodkumar Prabhakaran, and Jacob Eisenstein. Learning to recognize dialect features. *arXiv preprint arXiv:2010.12707*, 2020.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [12] Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. A survey of vision-language pre-trained models. *arXiv preprint arXiv:2202.10936*, 2022.
- [13] Constantin Eichenberg, Sidney Black, Samuel Weinbach, Letitia Parcalabescu, and Anette Frank. Magma—multimodal augmentation of generative models through adapter-based finetuning. *arXiv preprint arXiv:2112.05253*, 2021.
- [14] Matt Gardner, Yoav Artzi, Victoria Basmova, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, et al. Evaluating models’ local decision boundaries via contrast sets. *arXiv preprint arXiv:2004.02709*, 2020.
- [15] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.
- [16] Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A Wichmann. Generalisation in humans and deep neural networks. *Advances in neural information processing systems*, 31, 2018.
- [17] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.
- [18] Karan Goel, Nazneen Rajani, Jesse Vig, Samson Tan, Jason Wu, Stephan Zheng, Caiming Xiong, Mohit Bansal, and Christopher Ré. Robustness gym: Unifying the nlp evaluation landscape. *arXiv preprint arXiv:2101.04840*, 2021.
- [19] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [20] Jindong Gu and Volker Tresp. Improving the robustness of capsule networks to image affine transformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7285–7293, 2020.
- [21] Jindong Gu, Volker Tresp, and Han Hu. Capsule network is not more robust than convolutional network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14309–14317, 2021.
- [22] Jindong Gu, Volker Tresp, and Yao Qin. Are vision transformers robust to patch perturbations? In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XII*, pages 404–421. Springer, 2022.
- [23] Jindong Gu, Hengshuang Zhao, Volker Tresp, and Philip HS Torr. Segpgd: An effective and efficient adversarial attack for evaluating and boosting segmentation robustness. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXIX*, pages 308–325. Springer, 2022.
- [24] Demi Guo, Alexander M Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning. *arXiv preprint arXiv:2012.07463*, 2020.
- [25] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- [26] Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian error linear units. *CoRR, abs/1606.08415*, 3, 2016.

- [27] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021.
- [28] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.
- [29] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [30] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [31] Christoph Kamann and Carsten Rother. Benchmarking the robustness of semantic segmentation models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8828–8838, 2020.
- [32] Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–1035, 2021.
- [33] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- [34] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [35] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- [36] Liunan Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [37] Wei Li, Can Gao, Guocheng Niu, Xinyan Xiao, Hao Liu, Jiachen Liu, Hua Wu, and Haifeng Wang. Unimo: Towards unified-modal understanding and generation via cross-modal contrastive learning. *arXiv preprint arXiv:2012.15409*, 2020.
- [38] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- [39] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.
- [40] Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint arXiv:2106.04489*, 2021.
- [41] Claudio Michaelis, Benjamin Mitzkus, Robert Geirhos, Evgenia Rusak, Oliver Bringmann, Alexander S Ecker, Matthias Bethge, and Wieland Brendel. Benchmarking robustness in object detection: Autonomous driving when winter is coming. *arXiv preprint arXiv:1907.07484*, 2019.
- [42] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [43] John Miller, Karl Krauth, Benjamin Recht, and Ludwig Schmidt. The effect of natural distribution shift on question answering models. In *International Conference on Machine Learning*, pages 6905–6916. PMLR, 2020.
- [44] Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2071–2081, 2022.
- [45] Francesco Pinto, Philip HS Torr, and Puneet K. Dokania. An impartial take to the cnn vs transformer robustness contest. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, pages 466–480. Springer, 2022.
- [46] Jieli Qiu, Yi Zhu, Xingjian Shi, Florian Wenzel, Zhiqiang Tang, Ding Zhao, Bo Li, and Mu Li. Are multimodal models robust to image and text perturbations? *arXiv preprint arXiv:2212.08044*, 2022.
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [48] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [49] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pages 5389–5400. PMLR, 2019.
- [50] Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. Beyond accuracy: Behavioral testing of nlp models with checklist. *arXiv preprint arXiv:2005.04118*, 2020.
- [51] Barbara Rychalska, Dominika Basaj, Alicja Gosiewska, and Przemysław Biecek. Models in the wild: On corruption robustness of neural nlp systems. In *Neural Information Processing: 26th International Conference, ICONIP 2019, Sydney, NSW, Australia, December 12–15, 2019, Proceedings, Part III* 26, pages 235–247. Springer, 2019.
- [52] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.
- [53] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Hua-jun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2018.

- [54] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. VI-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5227–5237, 2022.
- [55] Yi-Lin Sung, Varun Nair, and Colin A Raffel. Training neural networks with fixed sparse masks. *Advances in Neural Information Processing Systems*, 34:24193–24205, 2021.
- [56] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.
- [57] Yapeng Tian and Chenliang Xu. Can audio-visual integration strengthen robustness under multimodal attacks? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5601–5611, 2021.
- [58] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021.
- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [60] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015.
- [61] Xiao Wang, Qin Liu, Tao Gui, Qi Zhang, 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*, pages 347–355, Online, aug 2021. Association for Computational Linguistics.
- [62] Boxi Wu, Jindong Gu, Zhifeng Li, Deng Cai, Xiaofei He, and Wei Liu. Towards efficient adversarial training on vision transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, pages 307–325. Springer, 2022.
- [63] Karren Yang, Wan-Yi Lin, Manash Barman, Filipe Condessa, and Zico Kolter. Defending multimodal fusion models against single-source adversaries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3340–3349, 2021.
- [64] Chenyu Yi, Siyuan Yang, Haoliang Li, Yap-peng Tan, and Alex Kot. Benchmarking the robustness of spatial-temporal models against corruptions. *arXiv preprint arXiv:2110.06513*, 2021.
- [65] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *arXiv preprint arXiv:2106.10199*, 2021.
- [66] Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *arXiv preprint arXiv:2111.03930*, 2021.
- [67] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419*, 2023.