

Adversarial Robustness of Transformer-Based Image Captioning Models

Authors: Vaida Danyil, Kohan Yurii

Mentor: Bryt Ori

Abstract

While it is well-known that transformer-based models are vulnerable to adversarial perturbations, the specific challenge of defending image captioning systems against such attacks has received limited attention.

In this project, we evaluate the robustness of several popular image captioning models (ExpansionNet, BLIP) against a range of perturbations.

We find that the models are largely unaffected by common distortions such as blur, noise, resizing, and rotation, maintaining high-quality captions. However, when exposed to adversarial perturbations, their performance degrades significantly, as measured by the METEOR score.

To address this vulnerability, we propose two defense strategies: (1) applying Gaussian blur to the input image before inference, and (2) filtering outputs based on a confidence score derived from the inverse perplexity of the model's output logits.

Our experiments show that while the confidence-based approach provides modest gains, input blurring is remarkably effective at neutralizing the impact of adversarial noise.

Introduction

Image captioning, the task of generating natural language descriptions for images, has seen remarkable advances with the advent of deep learning, particularly transformer-based architectures. Modern models such as ExpansionNet, BLIP demonstrate strong performance on standard datasets, producing semantically rich descriptions. However, as these models are increasingly deployed in applications, concerns about their reliability under malicious manipulation have grown.

While it is widely recognized that vision-language models are vulnerable to adversarial attacks-carefully crafted perturbations designed to fool a model without visibly altering the input-the specific problem of adversarial robustness in image captioning has received little attention. Unlike classification tasks, where the impact of an adversarial example is often easily quantifiable, assessing robustness in generative tasks such as captioning poses unique challenges.

In this work, we aim to bridge this gap by conducting a systematic evaluation of the robustness of popular image captioning models. We subject them to both basic perturbations (e.g. blur, noise, rotation) and adversarial pixel-level perturbations. Our results show that while models are generally resilient to basic transformations, they are highly susceptible to adversarial attacks, leading to significant degradation in caption quality.

To mitigate this vulnerability, we explore two simple yet effective defense strategies. First, we apply Gaussian blur to input images prior to inference, exploiting the tendency of adversarial noise to be high in frequency. Second, we introduce a confidence-based filtering mechanism that leverages the inverse perplexity of the output logits as a proxy for caption reliability. Our experiments reveal that while the confidence-based approach provides moderate improvements, input blurring offers substantial robustness gains with minimal impact on clean image captioning performance.

Through this work, we highlight the need for more attention to adversarial robustness in generative vision-language models and provide practical insights into lightweight defenses.

Related Work

Adversarial Attacks on Vision-Language Models.

The vulnerability of deep neural networks to adversarial examples has been extensively studied in the context of image classification. However, adversarial robustness in image captioning has only recently gained attention. Recent work by Hu et al. in [*“Non-Targeted Adversarial Attacks on Vision-Language Models via Maximizing Information Entropy”*](#) (2024) demonstrates that small, carefully designed perturbations can significantly distort the captions generated by image captioning models, often producing outputs that are semantically incorrect or completely unrelated to the image. Their approach highlights the difficulty of defending generative multimodal systems, especially under input-space attacks that exploit uncertainty in the vision-to-language mapping.

Out-of-Distribution Detection in Image Captioning.

While most robustness research in captioning focuses on the accuracy of generated descriptions, relatively little work has been done on confidence estimation. A notable exception is the work by Shalev et al. in [*“A Baseline for Detecting Out-of-Distribution Examples in Image Captioning”*](#), which introduces a method to detect OOD input images by setting a threshold on the likelihood scores of generated captions. Inspired by this approach, we adopt a similar idea to estimate caption confidence via the inverse perplexity of output logits, using it as a filtering mechanism to identify potentially unreliable captions.

Simple Defenses Against Adversarial Attacks.

Input transformation techniques, such as image resizing, JPEG compression, and Gaussian blur, have been explored in classification settings as lightweight defenses. These methods aim to remove high-frequency adversarial noise while preserving semantic content. We extend this line of defense to the captioning domain by applying Gaussian blur to input images and evaluating its effect on adversarial robustness, finding it to be highly effective.

Chosen models and datasets

To benchmark and analyze the robustness of image captioning systems, we selected two representative models based on their popularity and performance in the community.

Image Captioning Models

We selected two models that reflect the current state of the art and community adoption trends. The first is **BLIP** (Bootstrapped Language Image Pretraining), a widely used vision-language model available on HuggingFace, known for its strong generalization capabilities. The second is **ExpansionNet**, a transformer-based model specifically designed for image captioning, which achieves competitive results on standard benchmarks and is often cited as a state-of-the-art solution in this space. Both models are publicly available and were used without task-specific fine-tuning for fairness and reproducibility.

Datasets

For standard evaluation, we used the **COCO Captions** dataset, a large-scale benchmark containing richly annotated images with multiple captions per image. All images were resized to 256×256 to ensure compatibility with the models and to standardize evaluation. To simulate out-of-distribution (OOD) scenarios, we used two types of input shifts:

- **Perturbed COCO Images:** These include adversarial pixel-level perturbations produced by PGD attack.
- **Cartoon Dataset:** As an additional source of natural distribution shift, we included a dataset of cartoon images, which differ significantly in texture and semantics from COCO.

For one of our results, we did a sanity check on a subset of the caltech-101 dataset, to ensure the credibility of the results. Chosen classes can be found in [Appendix B](#).

Metrics

We used NLTK’s METEOR metric (METEOR v1.0) to understand how good our captions are:

Caption’s correspondence	Bad	Weak	Good	Very Good	SOTA
METEOR score	0-0.2	0.2-0.35	0.35-0.5	0.5-0.6	0.6+

Human study

To better understand how human perception compares to that of image captioning models in noisy environments, we conducted a small-scale user study. Our goal was to evaluate human ability to describe images affected by various levels of additive noise and to compare this with model-generated captions under similar conditions.

Study Setup

We selected 10 images from the COCO dataset and applied four levels of random noise: **clean (no noise)**, **slightly noisy**, **noisy**, **very noisy**. The full set of clean reference images used in the study is included in [Appendix A](#).

Each participant was shown 10 images, each with a different noise level. Importantly, no participant saw the same image at more than one noise level to avoid bias from repeated exposure. The time each image was displayed was limited to simulate realistic viewing conditions. Participants were instructed to write a single-sentence caption describing the content of each image as accurately as possible.

Annotation and Evaluation

Collected captions were evaluated by human annotators and classified into four quality categories: **very good**, **good**, **bad**, **very bad**.

The aggregate results are presented in the format **[very good – good – bad – very bad]** for each noise level. As expected, caption quality degraded as noise increased, indicating that even humans struggle to interpret heavily distorted images. The results, summarized below, validate the intuition that perceptual degradation is not unique to machine learning models:

Noise level	Caption Quality Distribution
Clean	[29 – 11 – 8 – 5]
Slight Noise	[18 – 19 – 14 -4]
Noise	[11 – 11 – 19 -6]
Heavy Noise	[11 – 9 – 16 – 15]

Perturbation benchmarks

To assess the robustness of image captioning models under more realistic distortions, we benchmarked performance using a variety of non-adversarial image perturbations.

We applied the following perturbations, grouped into two categories:

- **Noise and Blur Corruptions:**
 - Gaussian Blur
 - Normal Noise
 - Uniform Noise
 - Laplace Noise
 - Patch Occlusion
- **Geometric and Color Transformations:**
 - Horizontal Flip
 - Vertical Flip
 - Random Perspective
 - Random Rotation
 - Random Color Shift

All images were resized to 256×256 before processing. For each corruption type, we computed the METEOR score of the model-generated captions relative to ground-truth annotations in the COCO dataset.

Below is a summary of the averaged results for **ExpansionNet**:

Corruption Type	METEOR Score
Clean	0.547
Blur	0.527
Normal Noise	0.516
Uniform Noise	0.527
Laplace Noise	0.513
Patch Occlusion	0.541
Horizontal Flip	0.544
Vertical Flip	0.421
Random Perspective	0.536
Random Rotation	0.453
Random Color Shift	0.522

Adversarial benchmarks

To evaluate the vulnerability of image captioning models to adversarial examples, we conducted white-box attacks using the Projected Gradient Descent (PGD) method. We targeted two transformer-based captioning models: **BLIP** and **ExpansionNet v2**. The attacks were untargeted and optimized an adversarial loss defined as the **cross-entropy** between the model's caption's logits on clean and perturbed inputs. This approach encourages the model to diverge from its original caption when exposed to adversarial images.

We employed PGD under the ℓ^∞ norm constraint with a maximum perturbation bound of $\epsilon = 8/255$, using **25 to 50 steps** per attack and a fixed step size. All perturbations were imperceptible to the human eye but designed to maximally degrade the captioning performance. Importantly, these perturbations were generated directly on the image input without altering model parameters, preserving the architecture's original behavior.

We evaluated robustness using the **METEOR score**, comparing captions produced from adversarial images to those from clean inputs. Both models experienced significant performance degradation under attack. For instance, **BLIP's** METEOR score dropped from 0.46 to 0.30 while **ExpansionNet's** score decreased from 0.58 to 0.26. These results underscore the high sensitivity of captioning systems to even minimal pixel-level perturbations.

In addition to quantitative degradation, adversarial perturbations often led to semantically unrelated or nonsensical captions. For example, a clean image of a cat sitting on a car, originally captioned as "a cat sitting in the back of a car" was adversarially altered to produce "a dog wearing a bow tie and a face mask", despite the visual changes being imperceptible. Such qualitative shifts highlight the severity of the attack and the lack of robustness in current models.

Here are 2 examples of images attacked on **ExpansionNet**, that have the highest drop in quality of captions (based on METEOR score):



Original caption:

“a white plate with a half of an egg on it”

Adversarial caption:

“an anemone anemone anemone
anemone anemone anemone”



Original caption:

“a cat sitting in the back of a car”

Adversarial caption:

“a dog wearing a bow tie and a face
mask”

Confidence metric

To improve the robustness of image captioning models and detect potentially unreliable outputs, we explored two confidence scoring methods. These metrics were designed to quantify how much trust can be placed in a generated caption, particularly under adversarial perturbations.

The first method, applied to **ExpansionNet**, measured the **l2 of the output probabilities for each token, and then took the mean for all tokens generated**. The idea was that adversarial images will have a lot of uncertainty in them by having multiple options with relatively high probabilities. The higher that value is, the less uncertainty there is in the output. Unfortunately, we found that the adversarial examples did not have this kind of uncertainty in them.

To address this, we used a second, more practical confidence measure based on **inverse perplexity**. This was computed using the model's output probabilities for each predicted word in the caption. Specifically, we normalized the product of the highest-likelihood tokens' probabilities by the square root of the caption length:

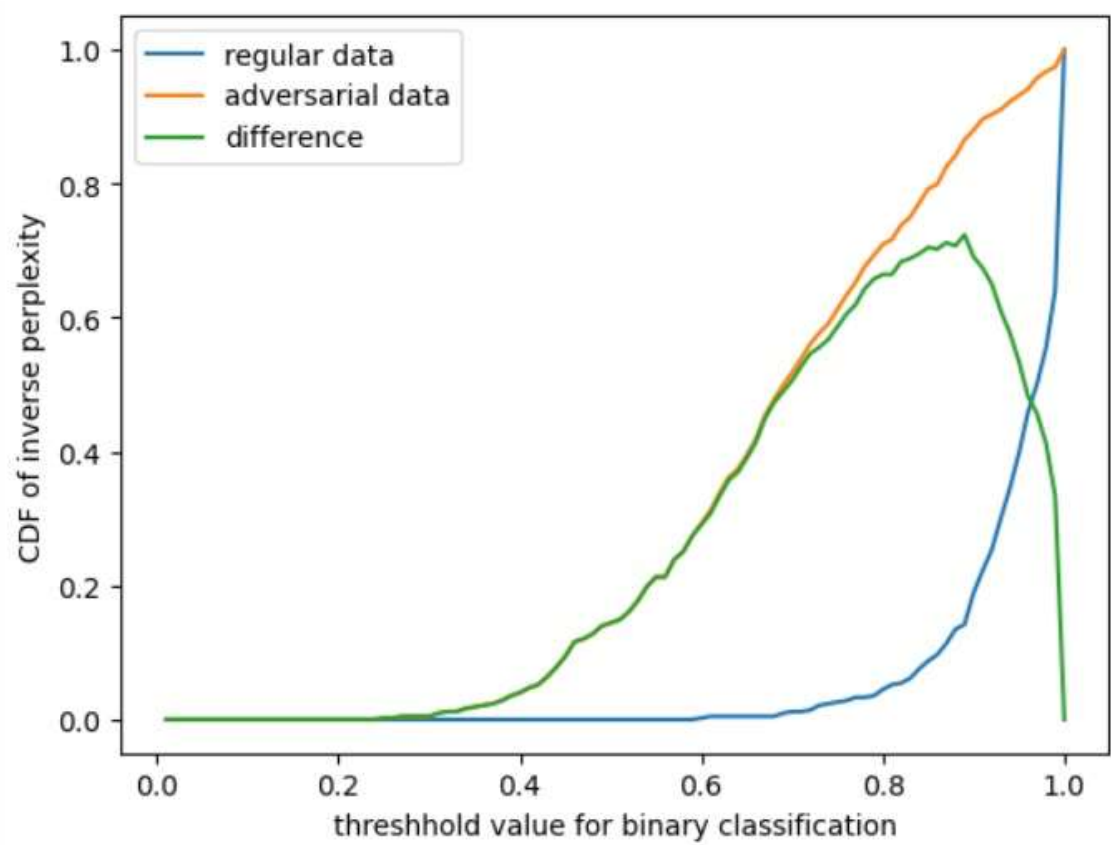
$$\frac{1}{\sqrt{outlen}} \sqrt{P(w_1) \cdot P(w_2) \dots P(w_{outlen})}$$

where *outlen* is the number of words in the caption, and $P(w_i)$ is the maximum predicted probability for word i at decoding time. This formulation penalizes uncertain or ambiguous predictions while remaining sensitive to the overall structure of the output.

We evaluated the effectiveness of this metric on both **ExpansionNet** and **BLIP**, using the **cumulative distribution function (CDF)** of the confidence scores across clean and adversarial examples. The resulting curves revealed a clear shift in the distribution under attack conditions, enabling the selection of meaningful thresholds to filter out low-confidence captions. For instance, by choosing a threshold of **0.89**, we were able to retain a large portion of clean captions while discarding most adversarial ones.

Here are some statistics that show the effectiveness of this method:

Optimal threshold	TP	FP	TN	FN	Precision	Recall	F_1
0.89	0.858	0.304	0.695	0.141	0.737	0.858	0.793



Blur

To defend against adversarial perturbations, we applied **Gaussian blur** to the input images prior to inference. This technique is motivated by the observation that adversarial noise tends to be high-frequency and spatially localized, while the core semantic content of natural images is typically low-frequency and spatially coherent. By applying a mild blur, we effectively suppress the high-frequency components of the image, diminishing the influence of adversarial perturbations without significantly affecting human-perceivable structure.

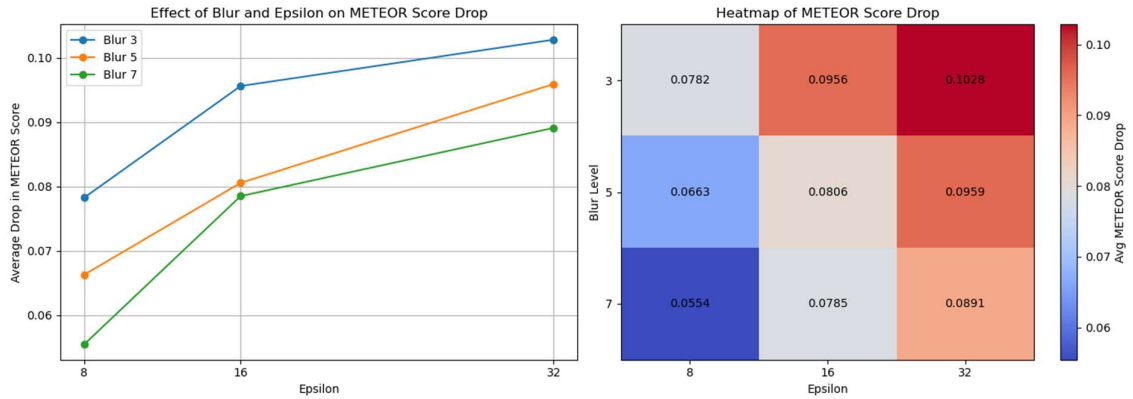
We found that this defense was surprisingly effective. When Gaussian blur was applied to adversarially perturbed inputs, the quality of generated captions improved noticeably across both **ExpansionNet** and **BLIP** models. This suggests that the models were no longer responding to the carefully crafted pixel-level noise introduced by the attacker. Importantly, **the same blurring operation had little to no negative impact on clean image performance**, preserving the METEOR and BLEU scores almost entirely. This balance—significant robustness gain with negligible accuracy loss—makes Gaussian blur an attractive lightweight defense.

The effectiveness of the blur can be attributed to the fact that most gradient-based attacks (like PGD) exploit small but precise pixel changes. These changes are often sensitive to input preprocessing steps; applying blur disrupts the gradient’s fine structure, making the adversarial signal less effective. Moreover, because transformer-based captioning models are relatively robust to mild spatial smoothing, the semantic integrity of clean images remains intact post-blur.

Overall, Gaussian blur serves as a simple yet powerful tool to neutralize adversarial perturbations, making it a practical first line of defense in real-world applications.

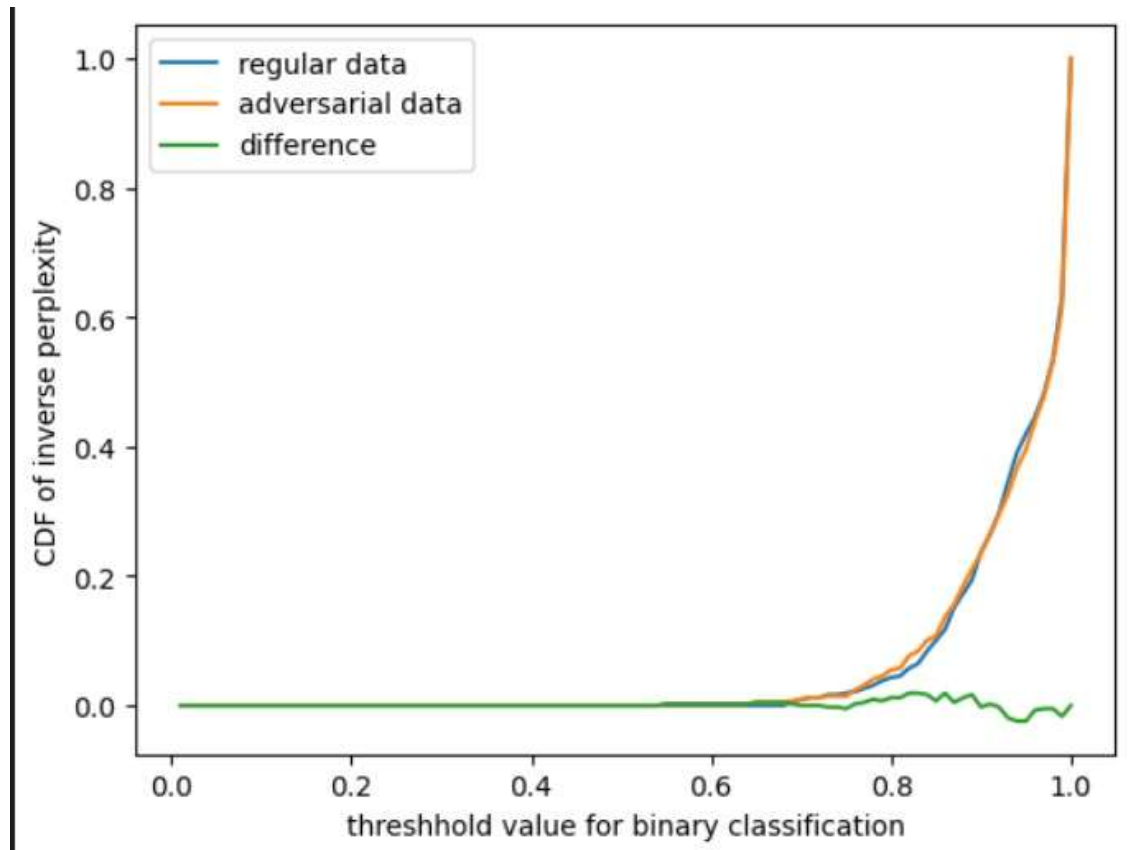
Here are the results of applying blur with different kernel sizes as a preprocessing step in our **ExpansionNet** and **BLIP** models:

Blur Level	Epsilon	Clean METEOR Score	Adversarial METEOR Score	Score Drop
3	8	0.5670	0.4888	0.0782
3	16	0.5648	0.4692	0.0956
3	32	0.5644	0.4616	0.1028
5	8	0.5591	0.4928	0.0663
5	16	0.5649	0.4844	0.0806
5	32	0.5613	0.4654	0.0959
7	8	0.5540	0.4986	0.0554
7	16	0.5578	0.4793	0.0785
7	32	0.5597	0.4706	0.0891



Confidence metric with Blur

We tried to apply confidence metric after blurring the images to further improve our results. Unfortunately, our confidence metric could not separate the images from clear dataset and adversarial examples.



This suggests that blur makes the distribution of adversarial images very close to the real ones, therefore not allowing it to impact output too much.

Sanity check

Although Gaussian blur proved highly effective in mitigating adversarial attacks in our image captioning experiments, its simplicity raised concerns about whether it truly neutralizes adversarial noise or merely obfuscates inputs. To validate whether this defense holds in a more standard classification context, we performed an additional sanity check using a transformer-based image classifier.

Experimental Setup.

We used the BLIP model's image encoder as a feature extractor and trained a simple linear classification head on top. The linear layer was fine-tuned on the **Caltech-101** dataset, a standard benchmark for object classification. The classifier achieved a clean test accuracy of around **90%** on unperturbed images.

Adversarial Evaluation.

To assess robustness, we used **AutoAttack**, a strong, standardized ensemble of adversarial attacks developed by Croce and Hein (2020). We applied AutoAttack under two adversarial perturbation budgets:

- For $\epsilon = 1/255$: the robust accuracy dropped to **0%** on adversarial inputs. However, when we applied **Gaussian blur** prior to classification, the model achieved **10–20% robust accuracy**, depending on the blur kernel size. This came at the cost of reduced clean accuracy, which ranged from **50–70%**, with stronger blur (larger kernels) yielding lower accuracy on clean inputs.
- For $\epsilon = 4/255$: the blur defense failed entirely. Regardless of the blur strength, the **robust accuracy dropped to 0%**, consistent with existing literature on the brittleness of transformer-based classifiers under moderate adversarial perturbations.

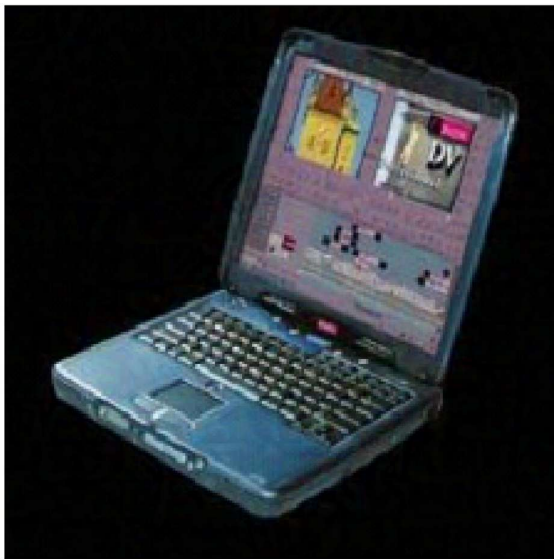
Decoder Experiment (Just for Fun).

To explore the semantic impact of these classifier-targeted adversarial perturbations, we passed the adversarial images—crafted to fool the Caltech-101 classifier—into the decoder of the original image captioning model.

Surprisingly, although these attacks were not designed to target the captioning model, the resulting captions frequently reflected **semantic shifts aligned with the misclassification**. In several cases, the caption not only became wrong, but also replaced the original object with one corresponding to the incorrect class label. For example, an image of a **laptop** adversarially perturbed to be classified as an **airplane** was captioned as *“a bunch of airplanes are parked in a parking”*.

This behavior resembles **targeted adversarial captioning**, despite the attack being applied solely to the classifier. Such examples suggest a transfer effect where classifier-focused perturbations indirectly manipulate the high-level linguistic content generated by captioning systems.

Predicted label: airplanes
Caption: a bunch of airplanes are parked in a parking



Predicted label: cup
Caption: a beer bottle with a logo on it



Findings.

As expected from the literature on transformer-based vision models, our classifier was highly vulnerable to adversarial perturbations, and applying blur to the adversarial inputs **did not** improve robustness. This result aligns with previous studies showing that input smoothing techniques like blur are not reliable defenses in classification settings. These findings suggest that while blur is effective in the specific context of image captioning, it should not be considered a general-purpose defense and must be interpreted with caution.

Appendix A: Images Used in the Human Study

The following images were selected from the COCO dataset and used in our human evaluation study. Each image was perturbed with one of four noise levels: *clean*, *slightly noisy*, *noisy*, or *very noisy*. Each participant saw only one version of each image. The original (clean) versions are shown below.



Example of perturbation levels is shown below:



Appendix B: Classes Used for the Sanity Check

In the classification-based sanity check, we constructed a lightweight classifier by fine-tuning a linear layer on top of the CLIP ViT-B/32 encoder using a subset of the **Caltech-101** dataset.

The following 9 object categories were selected to provide a diverse range of visual concepts: airplane, camera, car, cell phone, chair, cup, elephant, human, laptop

These categories were chosen to span both man-made and natural objects, with a mixture of rigid and deformable structures, to better evaluate the classifier's generalization and its robustness under adversarial perturbations and blur transformations.