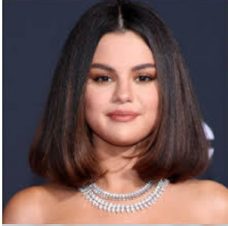


## A PROMPTING THE IMAGE-TO-TEXT MODEL

To ensure that captions are descriptive and composed of short sentences, we prompt our image-to-text model with the following for all experiments:



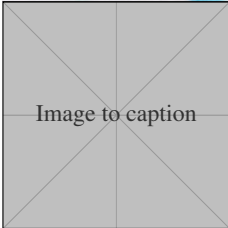
is a realistic picture of two penguins. They are holding hands. They are standing in front of the sea. The picture is mostly grey. The penguins are facing away from the camera. They take up most of the image.



is a portrait photograph of a famous person. She is wearing two necklaces. She has dark hair and is wearing makeup. She is facing the camera and the background is black.



is a cute photograph of three kittens. They are under a blanket. The background is blurred but it seems white and orange. The blanket is purple. The two cats on the right are orange and the one on the left is grey. The orange cats have open eyes and the grey cat has closed eyes. They are all super cute.



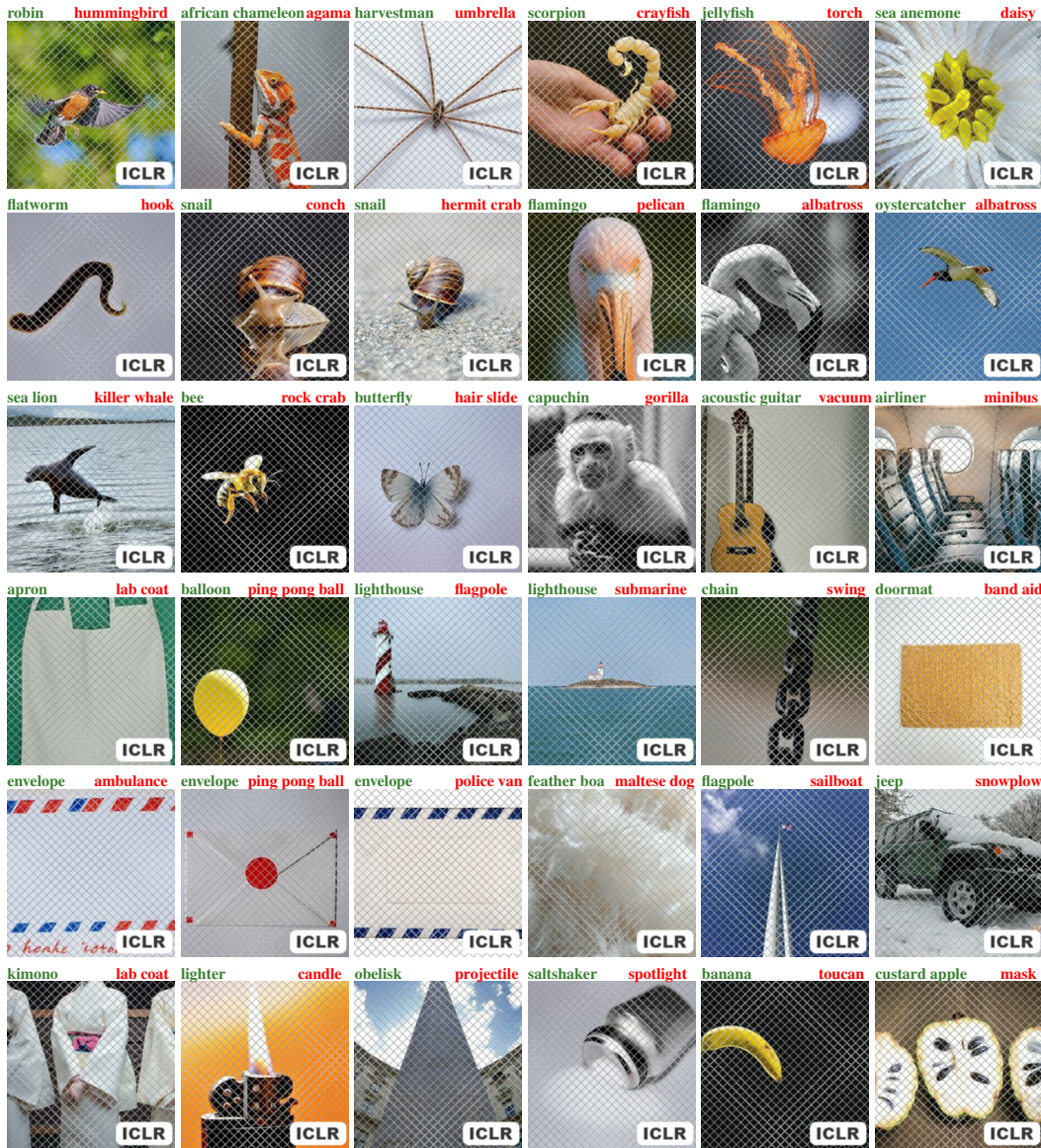
is a realistic photograph of a [label name]. [...]

These images are not from IMAGENET but downloaded from the web. Choosing these captions is only done once and then fixed for all experiments in the paper. We also set the decoding strategy to be *greedy* (as we did not observe significant improvements from using beam search). We highlight that any expert knowledge, if needed, is only required to annotate these three images. Compared to annotating a full test set, the cost is negligible.

## B ADDITIONAL RESULTS

### B.1 OPEN-ENDED FAILURE SEARCH

**Large scale experiments.** Similarly to Fig. 2 and Table 1, Fig. 6 and Table 3 show failure cases automatically found by our pipeline for a RESNET-50 available on TF-HUB in a fully open-ended manner (i.e., without leveraging an external dataset). The labels considered are a subset of the 200 labels present in IMAGENET-A. We let the reader interpret these failure cases themselves. The failures are diverse and are due to different factors, such as: (i) misleading color patterns (e.g., sea anemone → daisy), (ii) spurious context (e.g., jeep → snowplow), (iii) missing knowledge (e.g., custard apple → mask), or (iv) hallucinations (e.g., feather boa → maltese dog).



**Figure 6: Illustration of failure cases listed in Table 3.** The correct label is to the left in green. The incorrect prediction is to the right in red. The model used is a RESNET-50 found on TF-HUB.

True label	Target label	Caption	Failure rate (target)
robin	hummingbird	a realistic photograph of a robin (oscine). — " — It is flying.	0.0032% 1× 7.35% 2264.3×
african chameleon	agama	a realistic photograph of an african chameleon (lizard). — " — He is holding a stick. The chameleon is orange and white.	0.15% 1× 1.01% 6.7×
harvestman	umbrella	a realistic photograph of a harvestman (arthropod). — " — It is shot from above. The harvestman is on a white background.	0.45% 1× 1.32% 2.9×
scorpion	crayfish	a realistic photograph of a scorpion (arthropod). — " — It is on a person's hand.	0.0042% 1× 0.13% 29.5×
jellyfish	torch	a realistic photograph of a jellyfish (invertebrate). — " — The background is black. The jellyfish is orange.	0.14% 1× 0.45% 3.2×
sea anemone	daisy	a realistic photograph of a sea anemone (coelenterate). — " — It is yellow and white. The background is blurred.	0.32% 1× 1.33% 4.2×
flatworm	hook	a realistic photograph of a flatworm (invertebrate). — " — It is on a white background.	0.61% 1× 1.58% 2.6×
snail	conch	a realistic photograph of a snail (mollusk). — " — It is on a black background. The snail is reflected on the floor.	0.039% 1× 0.88% 22.5×
snail	hermit crab	a realistic photograph of a snail (mollusk). — " — It is on a grey road.	0.0082% 1× 0.10% 12.2×
flamingo	pelican	a realistic photograph of a flamingo (aquatic bird). — " — It is a close up of the head. The flamingo is facing the camera.	0.023% 1× 0.87% 37.6×
flamingo	albatross	a realistic photograph of a flamingo (aquatic bird). — " — It is black and white. The flamingo is looking to the right.	0.081% 1× 1.30% 16.1×
oystercatcher	albatross	a realistic photograph of an oystercatcher (wading bird). — " — It is flying.	0.0025% 1× 0.52% 208.2×
sea lion	killer whale	a realistic photograph of a sea lion (seal). — " — It is jumping out of the water.	0.14% 1× 4.31% 31.3×
bee	rock crab	a realistic photograph of a bee (insect). — " — It is flying. The background is black.	0.086% 1× 0.18% 2.1×
cabbage butterfly	hair slide	a realistic photograph of a cabbage butterfly (butterfly). — " — It is on a white background. It is in the middle of the image.	0.099% 1× 1.98% 20.0×
capuchin	gorilla	a realistic photograph of a capuchin (monkey). — " — It is a black and white photograph.	0.011% 1× 0.84% 73.9×
acoustic guitar	vacuum	a realistic photograph of an acoustic guitar (stringed instrument). — " — It is leaning against a wall.	0.20% 1× 1.02% 5.0×
airliner	minibus	a realistic photograph of an airliner (heavier-than-air craft). — " — There are seats in the foreground.	0.16% 1× 3.05% 19.2×
apron	lab coat	a realistic photograph of an apron (clothing). — " — It is white.	0.44% 1× 3.57% 8.1×
balloon	ping-pong ball	a realistic photograph of a balloon (aircraft). — " — It is yellow. The background is blurred.	0.53% 1× 17.86% 33.8×
lighthouse	flagpole	a realistic photograph of a beacon (structure). — " — The lighthouse has red and white stripes.	0.095% 1× 2.12% 22.4×
lighthouse	submarine	a realistic photograph of a beacon (structure). — " — It is on a small island at the horizon.	0.041% 1× 12.5% 308.0×
chain	swing	a realistic photograph of a chain (attachment). — " — The chain is vertical. The chain is in focus.	1.22% 1× 12.5% 10.2×
doormat	band aid	a realistic photograph of a doormat (floor cover). — " — The doormat is rectangular and is on a white background.	0.94% 1× 5.68% 6.1×
envelope	ambulance	a realistic photograph of an envelope (instrumentality). — " — It has white and has red and blue stripes at the top and bottom.	0.210% 1× 17.86% 60.0×
envelope	ping-pong ball	a realistic photograph of an envelope (instrumentality). — " — It is white and has a red dot on it.	1.04% 1× 75.5% 72.0×
envelope	police van	a realistic photograph of an envelope (instrumentality). — " — It has white and has white and blue diagonal stripes at the top and bottom.	0.510% 1× 6.94% 11.7×
feather boa	maltese dog	a realistic photograph of a feather boa (garment). — " — It is white and fluffy.	4.59% 1× 41.67% 9.1×
flagpole	sailboat	a realistic photograph of a flagpole (stick). — " — It is white and the sky is blue.	0.19% 1× 2.19% 11.5×
jeep	snowplow	a realistic photograph of a jeep (motor vehicle). — " — It is parked in the snow.	0.30% 1× 17.86% 59.3×
kimono	lab coat	a realistic photograph of a kimono (garment). — " — It is white.	0.69% 1× 2.84% 4.1×
lighter	candle	a realistic photograph of a lighter (instrumentality). — " — It has a flame coming out of it.	5.37% 1× 41.67% 7.8×
obelisk	projectile	a realistic photograph of an obelisk (structure). — " — It is pointing up. The sky is blue.	0.14% 1× 1.09% 7.7×
saltshaker	spotlight	a realistic photograph of a saltshaker (container). — " — It has a silver lid. The salt shaker is on a white background. The salt is spilling out of the jar.	1.02% 1× 13.89% 13.7×
banana	toucan	a realistic photograph of a banana (produce). — " — It is yellow and is floating in the air. The background is black.	0.0058% 1× 0.047% 8.2×
custard apple	mask	a realistic photograph of a custard apple (produce). — " — The fruit is cut in half.	0.32% 1× 4.46% 14.0×

**Table 3: Absolute failure rates of a RESNET-50 for 36 additional true and target label pairs.** We show the target failure rate (i.e., the model prediction is the target label). Captions are automatically discovered using the method detailed in [Sec. 3](#). Note that to the contrary of [Table 1](#), we consider an image to be misclassified when the top-1 prediction is wrong (and not from the same WORDNET parent) rather than when the true label is not part of the top-3 predictions.



**Generalization to other vision architectures.** We investigate whether our approach extends to other, more challenging backbones, such as a ViT-B/32 and ViT-B/8 which obtain much better performance than a RESNET50 on IMAGENET. We use the `fly` and `bee` failure case as our case study and run the same experiment as the one presented in Table 1. Results are reported in Table 4. First, we observe that both ViT models exhibit the same bias than the one found for the RESNET-50 model (that a fly on a flower is more often confused as a bee). Second, while the failure rates increase significantly compared to the baseline, the bias seems to be less pronounced (with only a  $114\times$  and  $13.5\times$  increase in failure rates compared to  $497\times$  for the RESNET) which highlights the qualities of both models.

Architecture	True label	Target label	Caption	Failure rate (target)	
ViT-B/32	fly	bee	a realistic photograph of a fly (insect).	0.0002%	$1\times$
			— " — it is on a flower.	0.02278%	$113.9\times$
ViT-B/8	fly	bee	a realistic photograph of a fly (insect).	0.0002%	$1\times$
			— " — it is on a flower.	0.0027%	$13.5\times$

**Table 4: Absolute failure rates of a ViT-B/8, ViT-B/32 on the `fly/bee` failure case.** We show the target failure rate (i.e., the model prediction is the target label). Captions are automatically discovered using the method detailed in Sec. 3.

**Significance of results.** We further investigate the significance of the results on two of the open-ended failure cases in the main paper (the ones exhibiting larger failure rates). Here, we evaluate our RESNET-50 and generate samples until we either find 10 images that cause the classifier to mispredict the class towards the target class (using top-1 accuracy instead of the typical top-3 to allow us to run many experiments efficiently) or find no misclassification towards the target class within 20K samples. We report the failure rate for the original and discovered captions and compute p-values using the Mann-Whitney U test (Mann & Whitney, 1947) at a significance level of 0.005 to determine if the differences in failure rates are statistically significant. Results are in Table 5. As the p-values ( $0.00015$ ,  $6.34 \cdot 10^{-5}$ ) are lower than 0.005, we find a *significant* result that images of the discovered caption (e.g., "... crayfish (crustacean). it is in a net.") are more often misclassified for the target class (e.g., `chainlink fence`) than the original caption (e.g., "... crayfish (crustacean).").

True label	Target label	Caption	Failure rate (target)	p-value
fly	bee	a realistic photograph of a fly (insect).	$0.00\% \pm 0.00$	0.00015
		— " — it is on a flower.	$0.58\% \pm 0.15$	
crayfish	chainlink fence	a realistic photograph of a crayfish (crustacean).	$0.00\% \pm 0.00$	$6.34 \cdot 10^{-5}$
		— " — it is in a net.	$2.09\% \pm 0.56$	

**Table 5: Significance of failure rates.** We report the mean and standard deviation of the target failure rate for the original and discovered captions over ten runs. We then compute the p-value to determine if the difference in failure rates between the original and discovered caption is statistically significant. We do this for two failure cases and find that our results are statistically significant.

## B.2 ADDITIONAL COMPARISONS ON IN-G-RN, IN-G-ViT

**Models considered.** We collate a large set of models trained on IMAGENET with differing size, pretraining, augmentation, and architectures in addition to the two RESNETs and ViTs we trained:

- ViT-B\*, ViT-L\*, ViT-S\* (Dosovitskiy et al., 2020): ViTs pretrained on IMAGENET21K.
- ViT-R\* (Steiner et al., 2022): a hybrid ViT and RESNET model pretrained on IMAGENET21K.
- BIT-\* (Kolesnikov et al., 2020): BIT models pretrained either on IMAGENET21K (BIT-M \*) or not pretrained (BIT-S \*).
- INCEPTION\_RESNET V2 (Szegedy et al., 2017): a hybrid INCEPTION, RESNET model with no pretraining.
- INCEPTION\* (Szegedy et al., 2015): INCEPTION models with no pretraining.
- RESNET\* (He et al., 2016): RESNET models with no pretraining.

**Transferability of errors on IN-G-RN and IN-G-VIT.** We evaluate how often failures in IN-G-RN and IN-G-VIT transfer to these models in Fig. 7. We can see that failures from both IN-G-RN and IN-G-VIT transfer across model architectures. However, the VITs and BITS which are pretrained on IMAGENET21K and achieve lower error on IMAGENET are fooled the least often. Within a model class, larger versions of the model seem more robust. For example, VIT-B/16 is more robust than VIT-B/32 and similarly the larger BITS (those of size 101x1) are more robust than their smaller counterparts (those of size 50x1). Thus, stronger pretraining and larger models seem to lead to improved (but not complete) robustness against these generated datasets.

**Error consistency on IN-G-RN and IN-G-VIT.** Finally, we measure the error consistency of the models in Fig. 8. We combine IN-G-RN and IN-G-VIT into one dataset and evaluate how often models make similar errors while accounting for the accuracy of each model (see Eq. 3 in Geirhos et al., 2021). A value of 100% indicates that the errors two models make are perfectly correlated and -100% that they are perfectly anti-correlated. It is striking that errors are most consistent within a model class: RESNETs make similar errors to other RESNETs trained in a similar manner and similarly BITS make similar errors to other BITS, especially BITS trained in the same manner.

IN-G-RN1	39%	53%	34%	35%	63%	56%	49%	50%	46%	39%	51%	41%	63%	60%	62%	66%	60%	61%	77%	73%	69%	69%	69%	75%	70%	71%	75%	100%	78%	59%	54%
IN-G-VIT1	50%	61%	44%	45%	67%	62%	56%	61%	55%	49%	59%	51%	69%	70%	71%	70%	68%	68%	77%	76%	74%	73%	74%	77%	74%	76%	77%	79%	77%	100%	77%
ImageNet	5%	8%	4%	4%	7%	7%	5%	7%	9%	8%	11%	8%	13%	11%	12%	13%	11%	10%	14%	11%	14%	12%	12%	14%	13%	12%	14%	11%	10%	6%	6%
	VIT-B/16	VIT-B/32	VIT-B/8	VIT-L/16	VIT-R26-S/32 (light aug)	VIT-R26-S/32 (med. aug)	VIT-R50-L/32	VIT-S/16	BIT-M (101x1)	BIT-M (101x3)	BIT-M (50x1)	BIT-M (50x3)	BIT-S (101x1)	BIT-S (101x3)	BIT-S-r152x4	BIT-S (50x1)	BIT-S (50x3)	Inception_ResNet V2	Inception V1	Inception V2	Inception V3	ResNet101 V1	ResNet152 V1	ResNet50 V1	ResNet101 V2	ResNet152 V2	ResNet50 V2	RN1	RN2	VIT1	VIT2

**Figure 7: Failure rates (top-3) for different models on two generated datasets and IMAGENET.** We report the failure rates of different models trained on IMAGENET on both IN-G-RN and IN-G-VIT as well as IMAGENET.

Error consistency for IN-G-RN, IN-G-ViT																																	
VIT-B/16	100%	52%	63%	62%	38%	48%	53%	54%	47%	50%	44%	50%	31%	32%	30%	33%	36%	20%	24%	27%	29%	27%	23%	28%	26%	24%	7%	21%	27%	38%			
VIT-B/32	52%	100%	47%	48%	45%	53%	53%	55%	43%	41%	43%	43%	34%	33%	34%	35%	39%	26%	29%	31%	33%	32%	28%	34%	30%	30%	10%	25%	26%	36%			
VIT-B/8	63%	47%	100%	65%	33%	43%	51%	49%	47%	51%	42%	52%	29%	30%	27%	32%	34%	18%	22%	25%	26%	26%	21%	25%	24%	21%	7%	20%	25%	36%			
VIT-L/16	62%	48%	65%	100%	33%	42%	52%	48%	45%	49%	40%	50%	28%	29%	26%	30%	33%	17%	21%	25%	26%	25%	20%	25%	24%	21%	6%	19%	25%	35%			
VIT-R26-S/32 (light aug)	38%	45%	33%	33%	100%	55%	50%	44%	35%	32%	37%	33%	29%	29%	30%	30%	32%	22%	23%	25%	29%	27%	25%	27%	26%	28%	11%	22%	18%	25%			
VIT-R26-S/32 (med. aug)	48%	53%	43%	42%	55%	100%	58%	52%	43%	40%	44%	41%	34%	35%	35%	35%	37%	25%	27%	30%	32%	31%	27%	31%	29%	29%	11%	24%	24%	32%			
VIT-R50-L/32	53%	53%	51%	52%	50%	58%	100%	49%	40%	42%	39%	42%	29%	29%	29%	30%	34%	21%	24%	27%	29%	28%	24%	28%	26%	26%	8%	21%	22%	30%			
VIT-S/16	54%	55%	49%	48%	44%	52%	49%	100%	46%	43%	45%	44%	35%	36%	35%	36%	38%	26%	29%	30%	33%	31%	27%	33%	30%	28%	9%	26%	30%	39%			
BIT-M (101x1)	47%	43%	47%	45%	35%	43%	40%	46%	100%	53%	59%	58%	43%	44%	42%	45%	42%	26%	30%	32%	33%	33%	27%	33%	32%	28%	10%	26%	27%	37%			
BIT-M (101x3)	50%	41%	51%	49%	32%	40%	42%	43%	53%	100%	48%	57%	34%	36%	33%	38%	35%	20%	23%	28%	28%	28%	22%	28%	26%	23%	6%	21%	25%	35%			
BIT-M (50x1)	44%	43%	42%	40%	37%	44%	39%	45%	59%	48%	100%	53%	49%	47%	47%	49%	43%	30%	33%	36%	37%	36%	31%	37%	35%	32%	11%	30%	27%	36%			
BIT-M (50x3)	50%	43%	52%	50%	33%	41%	42%	44%	58%	57%	53%	100%	39%	40%	36%	42%	39%	22%	26%	29%	31%	30%	25%	30%	29%	26%	8%	24%	27%	37%			
BIT-S (101x1)	31%	34%	29%	28%	29%	34%	29%	35%	43%	34%	49%	39%	100%	51%	54%	53%	42%	35%	36%	36%	39%	38%	34%	39%	36%	36%	14%	33%	27%	34%			
BIT-S (101x3)	32%	33%	30%	29%	29%	35%	29%	36%	44%	36%	47%	40%	51%	100%	48%	53%	42%	31%	34%	35%	37%	38%	33%	38%	37%	34%	11%	31%	32%	37%			
BIT-S (50x1)	30%	34%	27%	26%	30%	35%	29%	35%	42%	33%	47%	36%	54%	48%	100%	52%	42%	38%	39%	37%	39%	39%	37%	40%	38%	38%	16%	34%	26%	33%			
BIT-S (50x3)	33%	35%	32%	30%	30%	35%	30%	36%	45%	38%	49%	42%	53%	53%	52%	100%	44%	33%	36%	38%	40%	40%	35%	40%	38%	36%	13%	33%	30%	37%			
Inception_ResNet V2	36%	39%	34%	33%	32%	37%	34%	38%	42%	35%	43%	39%	42%	42%	42%	44%	100%	36%	41%	45%	47%	46%	40%	46%	44%	41%	15%	37%	31%	39%			
Inception V1	20%	26%	18%	17%	22%	25%	21%	26%	26%	20%	30%	22%	35%	31%	38%	33%	36%	100%	44%	36%	40%	40%	42%	40%	38%	42%	25%	41%	19%	26%			
Inception V2	24%	29%	22%	21%	23%	27%	24%	29%	30%	23%	33%	26%	36%	34%	39%	36%	41%	44%	100%	38%	41%	41%	42%	42%	41%	39%	22%	41%	26%	32%			
Inception V3	27%	31%	25%	25%	25%	30%	27%	30%	32%	28%	36%	29%	36%	35%	37%	38%	45%	36%	38%	100%	43%	43%	39%	42%	41%	40%	16%	33%	27%	33%			
ResNet101 V1	29%	33%	26%	26%	29%	32%	29%	33%	33%	28%	37%	31%	39%	37%	39%	40%	47%	40%	41%	43%	100%	53%	51%	53%	50%	50%	19%	39%	25%	34%			
ResNet152 V1	27%	32%	26%	25%	27%	31%	28%	31%	33%	28%	36%	30%	38%	38%	39%	40%	46%	40%	41%	43%	53%	100%	50%	51%	49%	49%	19%	38%	28%	34%			
ResNet50 V1	23%	28%	21%	20%	25%	27%	24%	27%	27%	22%	31%	25%	34%	33%	37%	35%	40%	42%	42%	39%	51%	50%	100%	50%	47%	50%	22%	39%	22%	29%			
ResNet101 V2	28%	34%	25%	25%	27%	31%	28%	33%	33%	28%	37%	30%	39%	38%	40%	40%	46%	40%	42%	42%	53%	51%	50%	100%	51%	51%	19%	41%	26%	33%			
ResNet152 V2	26%	30%	24%	24%	26%	29%	26%	30%	32%	26%	35%	29%	36%	37%	38%	38%	44%	38%	41%	41%	50%	49%	47%	51%	100%	46%	18%	39%	27%	34%			
ResNet50 V2	24%	30%	21%	21%	28%	29%	26%	28%	28%	23%	32%	26%	36%	34%	38%	36%	41%	42%	39%	40%	50%	49%	50%	51%	46%	100%	22%	39%	21%	28%			
RN1	7%	10%	7%	6%	11%	11%	8%	9%	10%	6%	11%	8%	14%	11%	16%	13%	15%	25%	22%	16%	19%	19%	22%	19%	18%	22%	100%	31%	-15%	5%			
RN2	21%	25%	20%	19%	22%	24%	21%	26%	26%	21%	30%	24%	33%	31%	34%	33%	37%	41%	41%	33%	39%	38%	39%	41%	39%	39%	31%	100%	21%	28%			
VIT1	27%	26%	25%	25%	18%	24%	22%	30%	27%	25%	27%	27%	27%	32%	26%	30%	31%	19%	26%	27%	25%	28%	22%	26%	27%	21%	-15%	21%	100%	56%			
VIT2	38%	36%	36%	35%	25%	32%	30%	39%	37%	35%	36%	37%	34%	37%	33%	37%	39%	26%	32%	33%	34%	34%	29%	33%	34%	28%	5%	28%	56%	100%			

Figure 8: Error consistency for all models on the combined IN-G-RN and IN-G-ViT dataset.

## C INTERPRETING FAILURE CASES

In this section, we aim at further characterizing failure cases by investigating why the models considered in our work yield wrong predictions on IN-G-RN and IN-G-VIT instances. For that, we compared a few failure cases with their respective nearest neighbors within the training set of IMAGENET in order to find patterns that shed light on the reasons behind wrong predictions.

We find the ten nearest neighbors of an IN-G-VIT instance in the embedding space induced by the second-to-last layer of a VIT trained on IMAGENET, using cosine similarity as the distance measure. This particular model achieves 82.7% top-1 accuracy on the IMAGENET validation set and has a failure rate of 100% on IN-G-VIT.

Fig. 9-11 show IN-G-VIT failure cases, along with their respective ten nearest neighbors within the *full* IMAGENET training set and the ten nearest neighbors with the same label. Results suggest that failure cases found by our approach induce errors by generating images that have elements in their background which are more often found in other classes within IMAGENET. We further observe that all failure cases are closer to examples containing objects semantically related with cues present in images that are not commonly found in the training set of IMAGENET for these classes. In Fig. 9(a), for example, we show a failure case labeled as `mushroom` for which the VIT predicts the label `snail`. All nearest neighbors shown in Fig. 9(b) are labeled as `snail` and contain elements such as human skin and grass in the background, which do not appear in the nearest neighbors from the label `mushroom`, as shown in Fig. 9(c). The VIT appears to be capturing spurious features in its representations (presence of human skin and grassy background) and relying on them to make predictions, which lead it to yield the wrong label for the IN-G-VIT instance presented in Fig. 9(a). In Fig. 10 and 11, we observe a similar pattern, where the VIT focuses on spurious cues such as the presence of a net in the background in Fig. 10(a) and snow in Fig. 11(a). Exploiting such correlations made the model mistake the particular instances of `cabbage butterfly` and `flagpole` as `barn spider` and `ski`, respectively.



(a) IN-G-ViT failure case.



(b) 10 nearest neighbors of (a) in the IMAGENET train set.



(c) 10 nearest neighbors of (a) in the IMAGENET train set and in its respective label (mushroom).

**Figure 9: Interpreting failure cases by inspecting nearest neighbors in the train set of IMAGENET.** We analyze the failure case in IN-G-ViT shown in panel (a). The example is labeled as `mushroom` and classified as `snail` by a ViT trained on IMAGENET. In panel (b), we show the 10 nearest neighbors of (a) in the train set of IMAGENET. All 10 neighbors are from the class `snail` and have similar features to the failure case, such as the background (e.g., the human hand), while the 10 nearest neighbors with the label `mushroom` showed in panel (c) do not have those features. This suggests that the ViT correlates such features with the label `snail`, and these spurious correlations likely induced it to misclassify the image in (a).





(a) IN-G-ViT failure case.



(b) 10 nearest neighbors of (a) in the IMAGENET train set.

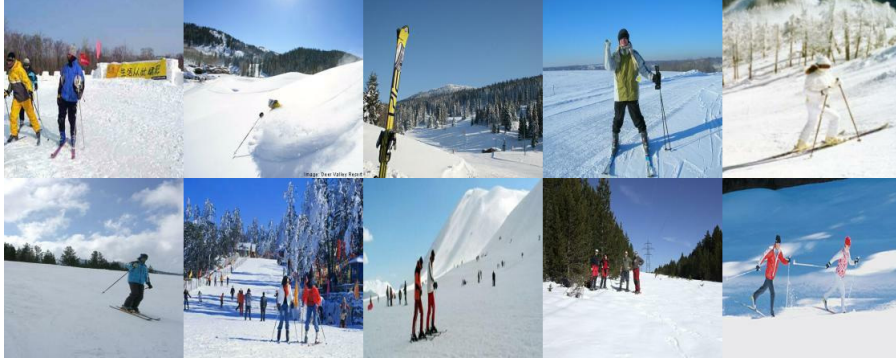


(c) 10 nearest neighbors of (a) in the IMAGENET train set and in its respective class (cabbage butterfly).

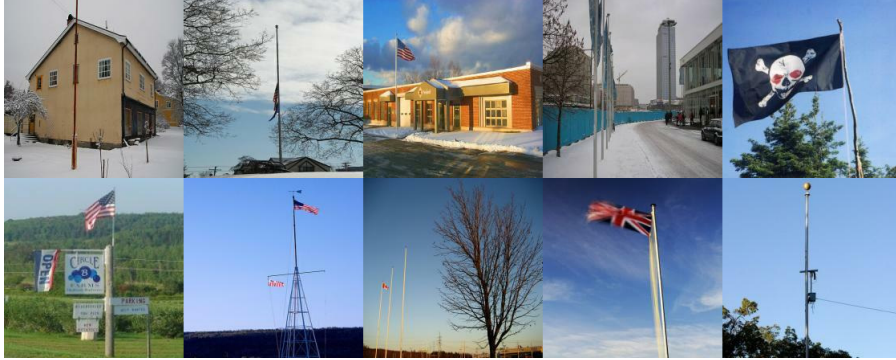
**Figure 10: Interpreting failure cases by inspecting their nearest neighbors in the train set of IMAGENET.** We analyze the failure case in IN-G-ViT shown in panel (a). The example is labeled as `cabbage butterfly` and classified as `barn spider` by a ViT trained on IMAGENET. In panel (b) we show the 10 nearest neighbors of (a) in the train set of IMAGENET. All 10 neighbors are from the classes `barn spider` or `black and gold garden spider` and have similarities to the failure case such as the white net in the background, while the 10 nearest neighbors in the class `cabbage butterfly` shown in panel (c) do not present those common features. This suggests that the ViT correlates such features with instances from the labels `barn spider` and `black and gold garden spider`, and exploiting these spurious correlations likely induced the model to misclassify the image in (a).



(a) IN-G-VIT failure case.



(b) 10 nearest neighbors of (a) in the IMAGENET train set..



(c) 10 nearest neighbors of (a) in the IMAGENET train set and in its respective class (flagpole)

**Figure 11: Interpreting failure cases by inspecting their nearest neighbors in the train set of IMAGENET.** We analyze the failure case in IN-G-VIT shown in panel (a). The example is labeled as `flagpole` and classified as `ski` by a VIT trained on IMAGENET. In panel (b) we show the 10 nearest neighbors of (a) in the train set of IMAGENET. All 10 neighbors have the label `ski` and have similarities to the failure case, such as simultaneously having snow on the ground and blue sky in the background, while the 10 nearest neighbors with the label `flagpole` shown in (c) do not present those common features. This suggests that the VIT correlates such features with instances from the label `ski`, and exploiting such spurious correlations likely induced it to misclassify the image in (a).

## D MALICIOUS USAGE AND MITIGATION STRATEGIES

This work demonstrates how to find failure cases in vision classifiers with the help of large-scale generative models. Much like adversarial examples (Biggio et al., 2013; Szegedy et al., 2013), malicious actors could leverage the proposed approach to build adversarial images that bypass automated online filtering mechanism. In this section, we discuss how to make classifiers robust to these failure cases.

First, classifiers can be trained with discovered failure cases to make them more robust to generated inputs. As a demonstration, we split the IN-G-VIT dataset into a train and test set (80% train, 20% test). We train the original VIT model in the exact same manner as before, except that batches are now made of 95% IMAGENET data and 5% IN-G-VIT data. We report results with and without additional synthetic data in Table 6. Training with additional generated data leads to a minimal loss of performance on IMAGENET while achieving nearly 90% top-1 accuracy on the IN-G-VIT test set. This demonstrates that adding the generated failure cases into the training set is an effective mitigation strategy.

Second, we note that our approach is computationally expensive. It requires hundreds to thousands of calls to the generative model and vision classifier to find a *single* failure case. Hiding the underlying classifier behind a rate-limited API can act as a first line of defense.

Training Set	top-1 on IMAGENET $\uparrow$	top-1 on IN-G-VIT $\uparrow$
IMAGENET (train)	<b>82.57 <math>\pm</math> 0.09</b>	5.60 $\pm$ 2.80
IMAGENET (train) + IN-G-VIT (train)	82.11 $\pm$ 0.05	<b>88.11 <math>\pm</math> 0.44</b>

**Table 6: top-1 accuracy on IMAGENET and IN-G-VIT.** We train a VIT model on either just IMAGENET or IMAGENET and IN-G-VIT. By training on IN-G-VIT, we achieve nearly 88% top-1 accuracy on IN-G-VIT (test) while minimally hurting performance on IMAGENET. To obtain standard deviations, we run the experiment with 5 random seeds.

## E ADDITIONAL VISUALIZATIONS



**Figure 12: Images from the text-to-image model used in this manuscript.** Images are generated with captions identical to those used in Fig. 2(b) and Fig. 2(c). A comparison with DALL·E 2 is shown in Fig. 13.

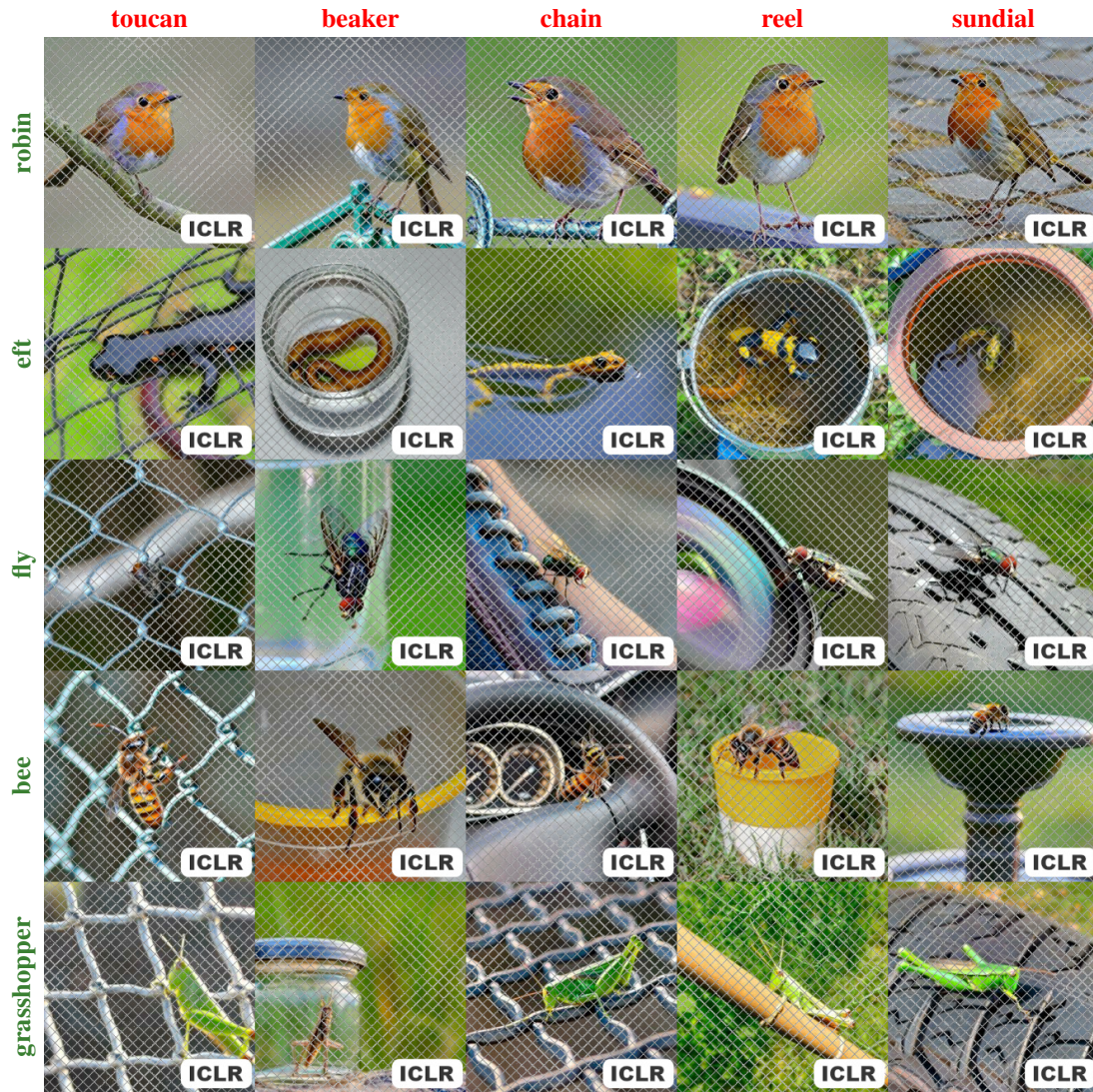


**Figure 13: DALL·E 2 images.** Images are generated with captions identical to those used in Fig. 2(b) and Fig. 2(c). A comparison with the text-to-image model used in the paper is shown in Fig. 12.



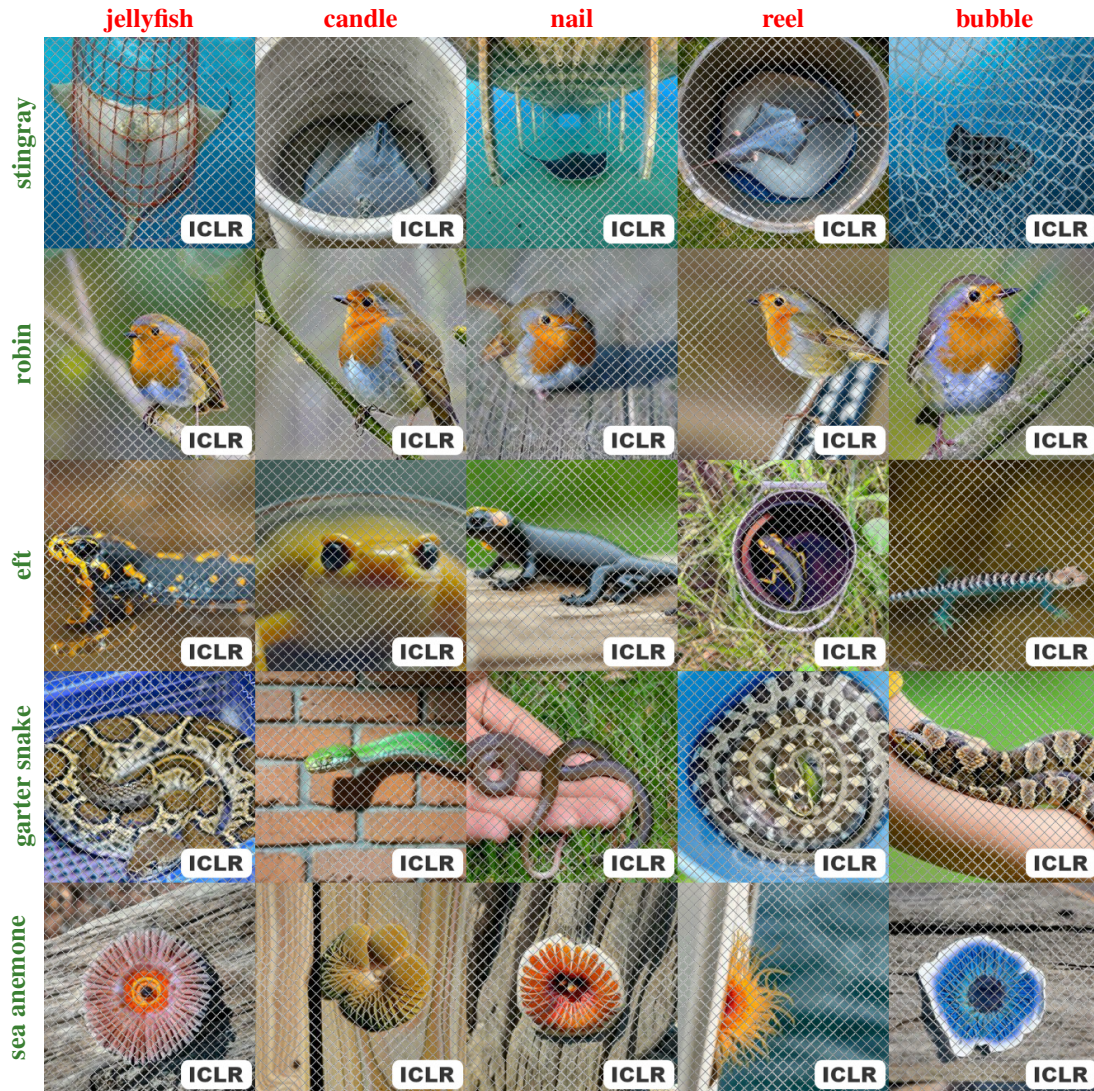
**Figure 14: STABLE-DIFFUSION images.** Images are generated with captions identical to those used in Fig. 2(b) and Fig. 2(c). A comparison with the text-to-image model used in the paper is shown in Fig. 12.





**Figure 15: Further examples from IN-G-RN.** The label at the top of the column is one of the incorrectly predicted top-3 labels and the label on the left is the true label.





**Figure 16: Examples from IN-G-ViT.** The label at the top of the column is one of the incorrectly predicted top-3 labels and the label on the left is the true label.