
Ambient Diffusion nni: Training Good Models with Bad Data

Anonymous Author(s)

Affiliation

Address

email

Abstract

We show how to use low-quality, synthetic, and out-of-distribution images to improve the quality of a diffusion model. Typically, diffusion models are trained on curated datasets that emerge from highly filtered data pools from the Web and other sources. We show that there is immense value in the lower-quality images that are often discarded. We present Ambient Diffusion Omni, a simple, principled framework to train diffusion models that can extract signal from all available images during training. Our framework exploits two properties of natural images – spectral power law decay and locality. We first validate our framework by successfully training diffusion models with images synthetically corrupted by Gaussian blur, JPEG compression, and motion blur. We then use our framework to achieve state-of-the-art ImageNet FID and we show significant improvements in both image quality and diversity for text-to-image generative modeling. The core insight is that noise dampens the initial skew between the desired high-quality distribution and the mixed distribution we actually observe. We provide rigorous theoretical justification for our approach by analyzing the trade-off between learning from biased data versus limited unbiased data across diffusion times.

1 Introduction

Large-scale, high-quality training datasets have been a primary driver of recent progress in generative modeling. These datasets are typically assembled by filtering massive collections of images sourced from the web or proprietary databases [25, 43, 53, 58, 59]. The filtering process is crucial to the quality of the resulting models [13, 27, 25, 32, 27]. However, filtering strategies are often heuristic and inefficient, discarding large amounts of data [51, 43, 25, 13]. We demonstrate that the data typically rejected as low-quality holds significant, underutilized value.

Extracting meaningful information from degraded data requires algorithms that explicitly model the degradation process. In generative modeling, there is growing interest in approaches that learn to generate directly from degraded inputs [18, 17, 14, 15, 7, 47, 39, 52, 5, 1, 2, 55, 71, 46, 64, 45, 11, 48]. A key limitation of existing methods is their reliance on knowing the exact form of the degradation. In real-world scenarios, image degradations—such as motion blur, sensor artifacts, poor lighting, and low resolution—are often complex and lack a well-defined analytical description, making this assumption unrealistic. Even within the same dataset, from ImageNet to internet scale text-to-image datasets, there are samples of varying qualities [28], as shown in Figures 3, 24, 27, 25. Given access to this mixed-bag of datapoints, we would like to sample from a tilted continuous measure of high-quality images, without sacrificing the diversity present in the training points.

The training objective of diffusion models naturally decomposes sampling from a target distribution into a sequence of supervised learning tasks [30, 61, 62, 16, 19, 9, 10]. Due to the power-law structure

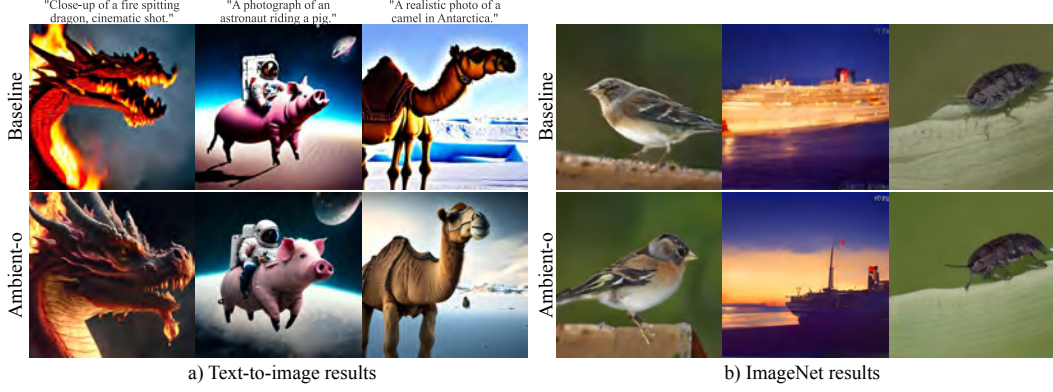


Figure 1: Effect of using Ambient-o for (a) training a text-to-image model (Micro-Diffusion [54]) and (b) a class-conditional model for ImageNet (EDM-2 [35]). All generations are initialized with the same noise. The baseline models are trained using all the data equally. Ambient-o changes the way the data is used during the diffusion process based on its quality. This leads to significant visual improvements without sacrificing diversity, as would happen with a filtering approach (see Fig. 29).

of natural image spectra [65], high diffusion times focus on generating globally coherent, semantically meaningful content [22], while low diffusion times emphasize learning high-frequency details.

Our first key theoretical insight is that low-quality samples can still be valuable for training in the high-noise regime. As noise increases, the diffusion process contracts distributional differences (see Theorem B.2), reducing the mismatch between the high-quality target distribution and the available mixed-quality data. At the same time, incorporating low-quality data increases the sample size, reducing the variance of the learned estimator. Our analysis formalizes this bias–variance trade-off and motivates a principled algorithm for training denoisers at high diffusion times using noisy, heterogeneous data.

For low diffusion times, our algorithm leverages a second key property of natural images: locality. We show a direct relationship between diffusion time and the optimal receptive field size for denoising. Specifically, small image crops suffice at lower noise levels. This allows us to borrow high-frequency details from out-of-distribution or synthetic images, as long as the marginal distributions of the crops match those of the target data.

We introduce Ambient Diffusion Omni (Ambient-o), a simple and principled framework for training diffusion models using arbitrarily corrupted and out-of-distribution data. Rather than filtering samples based on binary ‘good’ or ‘bad’ labels, Ambient-o retains all data and modulates the training process according to each sample’s utility. This enables the model to generate diverse outputs without compromising image quality. Empirically, Ambient-o advances the state of the art in unconditional generation on ImageNet and enhances diversity in text-conditional generation without sacrificing fidelity. Theoretically, it achieves improved bounds for distribution learning by optimally balancing the bias–variance trade-off: low-quality samples introduce bias, but their inclusion reduces variance through increased sample size.

2 Background and Related Work

Diffusion Modeling. Diffusion models transform the problem of sampling from p_0 into the problem of learning *denoisers* for smoothed versions of p_0 defined as $p_t = p_0 \otimes \mathcal{N}(0, \sigma^2(t)\mathbf{I})$. We typically denote with $X_0 \sim p_0$ the R.V. distributed according to the distribution of interest and $X_t = X_0 + \sigma(t)Z$, the R.V. distributed according to p_t . The target is to estimate the set of optimal l_2 denoisers, i.e., the set of the conditional expectations: $\{\mathbb{E}[X_0|X_t = \cdot]\}_{t=1}^T$. Typically, this can be achieved through supervised learning by minimizing the following loss (or a re-parametrization of it):

$$J(\theta) = \mathbb{E}_{t \in \mathcal{U}[0, T]} \mathbb{E}_{x_0, x_t | t} \left[\|h_\theta(x_t, t) - x_0\|^2 \right], \quad (2.1)$$

that is optimized over a function family $\mathcal{H} = \{h_\theta : \theta \in \Theta\}$ parametrized by network parameters θ . For sufficiently expressive families, the minimizer is indeed: $h_{\theta^*}(x, t) = \mathbb{E}[X_0|X_t = x]$.

Learning from noisy data. The diffusion modeling framework described above assumes access to samples from the distribution of interest p_0 . An interesting variation of this problem is to learn to sample from p_0 given access to samples from a tilted measure \tilde{p}_0 and a known degradation model. In Ambient Diffusion [18], the goal is to sample from p_0 given pairs (Ax_0, A) for a matrix $A : \mathbb{R}^{m \times n}$, $m < n$, that is distributed according to a known density $p(A)$. The techniques in this work were later generalized to accommodate additive Gaussian Noise [15, 17, 1] in the measurements. More recently there have been efforts to further broaden the family of degradation models considered through Expectation-Maximization approaches that involve multiple training runs [52, 5].

Recent work from [17] has shown that, at least for the Gaussian corruption model, leveraging the low-quality data can tremendously increase the performance of the trained generative models. In particular, the authors consider the setting where we have access to a few samples from p_0 , let's denote them $\mathcal{D}_0\{x_0^{(i)}\}_{i=1}^{N_1}$ and many samples from p_{t_n} , let's denote them $\mathcal{D}_{t_n}\{x_{t_n}^{(i)}\}_{i=1}^{N_2}$, where $p_{t_n} = p_0 \otimes \mathcal{N}(0, \sigma^2(t_n)\mathbf{I})$ is a smoothed version of p_0 at a known noise level t_n . The clean samples are used to learn denoisers for all noise levels $t \in [0, T]$ while the noisy samples are used to learn denoisers only for $t \geq t_n$, using the training objective:

$$J_{\text{ambient}}(\theta) = \mathbb{E}_{t \in \mathcal{U}(t_n, T)} \sum_{i=1}^{N_2} \mathbb{E}_{x_t | x_{t_n}^{(i)}} \left[\left\| \alpha(t) h_\theta(x_t, t) + (1 - \alpha(t)) x_t - x_{t_n}^{(i)} \right\|^2 \right], \quad (2.2)$$

with $\alpha(t) = \frac{\sigma^2(t) - \sigma^2(t_n)}{\sigma^2(t)}$. Note that the objective of equation 2.2 only requires samples from p_{t_n} (instead of p_0) and can be used to train for all times $t \geq t_n$. This algorithm uses $N_1 + N_2$ datapoints to learn denoisers for $t > t_n$ and only N_1 datapoints to learn denoisers for $t \leq t_n$. The authors show that even for $N_1 \ll N_2$, the model performs similarly to the setting of training with $(N_1 + N_2)$ clean datapoints. The main limitation of this method and its related works is that the degradation process needs to be known. However, in many applications, we have data from heterogeneous sources and various qualities, but there is no analytic form or any prior on the corruption model.

Data filtering. One of the most crude, but widely used, approaches for dealing with heterogeneous data sources is to remove the low-quality data and train only the high-quality subset [43, 25, 23]. While this yields better results than naively training on the entire distribution, it leads to a decrease in diversity and relies on heuristics for optimizing the filtering. An alternative strategy is to train on the entire distribution and then fine-tune on high-quality data [13, 54]. This approach better trades the quality-diversity trade-off but still incurs a loss of diversity and is hard to calibrate.

Training with synthetic data. A lot of recent works have shown that synthetic data can improve the generative capabilities of diffusion models when mixed properly with real data from the distribution of interest [24, 3, 4]. In this work, we show that it helps significantly to view synthetic data as corrupted versions of the samples from the real distribution and incorporate this perspective into the training objective.

3 Method

We propose a new framework that extends beyond [17] to enable training generative models directly from arbitrarily corrupted and out-of-distribution data, without requiring prior knowledge of the degradation process. We begin by formalizing the setting of interest.

Problem Setting. We are given a dataset $\mathcal{D} = \{w_0^{(i)}\}_{i=1}^N$ consisting of N datapoints. Each point in \mathcal{D} is drawn from a mixture distribution \tilde{p}_0 , which mixes p_0 (the distribution of interest) and an alternative distribution q_0 that may contain various forms of degradation or out-of-distribution content. We assume access to two labeled subsets, S_G, S_B , where points in S_G are known to come from the clean distribution p_0 , and points in S_B from the corrupted distribution q_0 . While this assumption simplifies the initial exposition, we relax it in Section G.1. We focus on the practically relevant regime where $|S_G| \ll |\mathcal{D}|$ —i.e., access to high-quality data is severely limited. The objective is to learn a generative model that (approximately) samples from the clean distribution p_0 , leveraging both clean and corrupted samples in its training.

We now describe how degraded and out-of-distribution samples can be effectively leveraged during training in both the high-noise and low-noise regimes of the diffusion process.

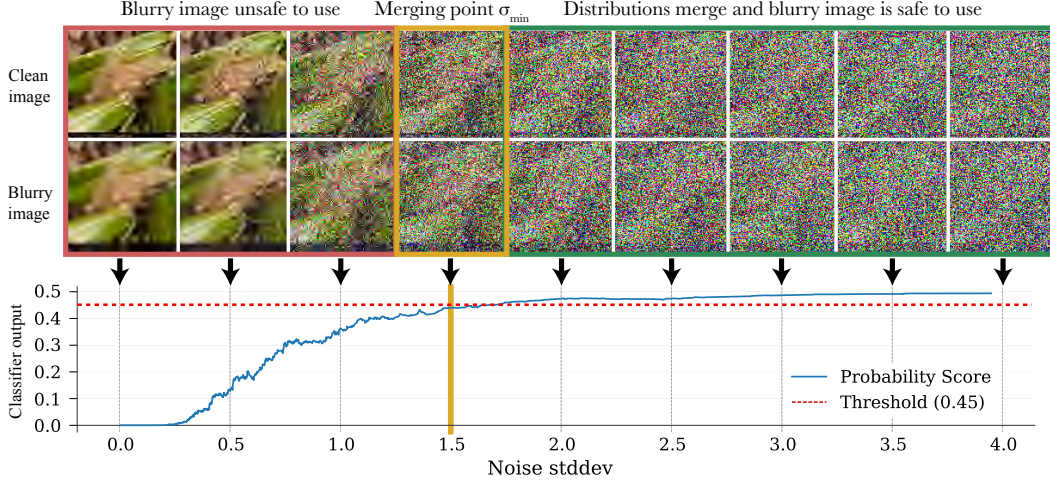


Figure 2: A time-dependent classifier trained to distinguish noisy clean and blurry images (blur kernel standard deviation $\sigma_B = 0.6$). At low noise the classifier is able to perfectly identify the blurry images, and outputs a probability close to 0. As the noise increases and the information in the image is destroyed, the clean and blurry distributions converge and the classifier outputs a prediction close to 0.5. The red line plots the threshold (selected at $\tau = 0.45$), which is crossed at $\sigma_t = 1.64$.

3.1 Learning in the high-noise regime (leveraging low-quality data)

Addition of gaussian noise contracts distribution distances. The first key idea of our method is that, at high diffusion times t , the noised target distribution p_t and the noised corrupted distribution \tilde{p}_t become increasingly similar (Theorem B.2), effectively attenuating the discrepancy introduced by corruption. This effect is illustrated in Figure 2 (top), where we compare a clean image and its degraded counterpart (in this case, corrupted by Gaussian blur). As the diffusion time t increases, the noised versions of both samples become visually indistinguishable. Consequently, samples from \tilde{p}_0 can be leveraged to learn (the score of) p_t , for $t > t_n^{\min}$. We formalize this intuition in Section B, and we also quantify that for large t there are statistical efficiency benefits for using a large sample from \tilde{p}_0 versus a small sample from p_0 .

Heuristic selection of the noise level. From the discussion so far, it follows that to use samples from \tilde{p}_0 , we need to assign them to a noise level t_n^{\min} . One can select this noise level empirically, i.e. we can ablate this parameter by training different models and selecting the one that maximizes the generative performance. However, this approach requires multiple trainings, which can be costly. Instead, we can find the desired noise level in a principled way as detailed below.

Training a classifier under additive Gaussian noise. To identify the appropriate noise level, we train a time-conditional classifier to distinguish between the noised distributions p_t and q_t across various diffusion times t . We use a single neural network $c_\theta^{\text{noise}}(x_t, t)$ that is conditioned on the diffusion time t , following the approach of time-aware classifiers used in classifier guidance [21]. The classifier is trained using labeled samples from S_G (clean) and S_B (corrupted) via the following objective:

$$J_{\text{noise}}(\theta) = \sum_{x_0 \in S_G} \mathbb{E}_{x_t|x_0} [-\log c_\theta^{\text{noise}}(x_t, t)] + \sum_{y_0 \in S_B} \mathbb{E}_{y_t|y_0} [-\log(1 - c_\theta^{\text{noise}}(y_t, t))] \quad (3.1)$$

Annotation. Once the classifier is trained, we use it to determine the minimal level of noise that must be added to the low-quality distribution q_0 so that it closely approximates a smoothed version of the high-quality distribution p_0 . Formally, we compute:

$$t_n^{\min} = \inf \left\{ t \in [0, T] : \frac{1}{|S_B|} \sum_{y_0 \in S_B} \mathbb{E}_{y_t|y_0} [c_\theta^{\text{noise}}(y_t, t)] > \tau \right\}, \quad (3.2)$$

for $\tau = 0.5 - \epsilon$ and for some $\epsilon > 0$. Subsequently, we form the annotated dataset $\mathcal{D}_{\text{annot}} = \{(w_0^{(i)} + \sigma_{t_n^{\min}} Z^{(i)}, t_n^{\min})\}_{i=1}^N \cup \{(x_0, 0) | x_0 \in S_G\}$, where the random variables $Z^{(i)}$ are i.i.d. standard

normals. In particular, our annotated dataset indicates that we should only use the samples from \mathcal{D} for diffusion times $t \geq t_n^{\min}$, for which the distributions have approximately merged and hence it is safe to use them. In fact, the optimal classifier assigns time t_n that corresponds to the first time for which $d_{\text{TV}}(p_t, q_t) \leq \epsilon$.

From arbitrary corruption to additive Gaussian noise. The afore-described approach reduces our problem of learning from data with arbitrary corruption to the setting of learning from data corrupted with additive Gaussian noise. The price we pay for this reduction is the information loss due to the extra noise we add to the samples during the annotation stage. We can now extend the objective function (2.2) to train our diffusion model. Suppose our annotated dataset is comprised of samples $\{(x_{t_i^{\min}}^{(i)}, t_i^{\min})\}$. Then our objective becomes:

$$J_{\text{ambient-o}}(\theta) = \mathbb{E}_{t \in \mathcal{U}[0, T]} \sum_{i: t_i^{\min} < t} \mathbb{E}_{x_t | x_{t_i^{\min}}^{(i)}} \left[\left\| \alpha(t, t_i^{\min}) h_{\theta}(x_t, t) + (1 - \alpha(t, t_i^{\min})) x_t - x_{t_i^{\min}}^{(i)} \right\|^2 \right],$$

where $\alpha(t, t_i^{\min}) = \frac{\sigma^2(t) - \sigma^2(t_i^{\min})}{\sigma^2(t)}$.

Moreover, the method is particularly well-suited to certain types of corruptions but is less effective for others. Because the addition of Gaussian noise suppresses high-frequency components—due to the spectral power law of natural images—our approach is most effective for corruptions that primarily degrade high frequencies (e.g., blur). In contrast, degradations that affect low-frequency content—such as color shifts, contrast reduction, or fog-like occlusions—are more challenging. This limitation is illustrated in Figure 9: masked images, for example, require significantly more noise to become usable compared to high-frequency corruptions like blur. In the extreme, the method reduces to a filtering approach, as infinite noise nullifies all information in the corrupted samples.

3.2 Learning in the low-noise regime (synthetic and out-of-distribution data)

So far, our algorithm implicitly results in varying amounts of training data across diffusion noise levels. At high noise, the model can leverage abundant low-quality data, whereas at low noise levels, it must rely solely on the limited set of high-quality samples. We now extend the algorithm to enable the use of synthetic and out-of-distribution data for learning denoisers at low-noise diffusion times.

To achieve this, we leverage another fundamental property of natural images: *locality*. At low diffusion times, the denoising task can be solved using only a small local region of the image, without requiring full spatial context. We validate this hypothesis experimentally in the Experiments Section (Figures 11, 12, 13, 14), where we show that there is a mapping between diffusion time t and the crop size needed to perform the denoising optimally at this diffusion time. Intuitively, the higher the noise, the more context is required to accurately reconstruct the image. Conversely, for lower noise, the local information within a small neighborhood suffices to achieve effective denoising. We use $\text{crop}(t)$ to denote the minimal crop size needed to perform optimal denoising at time t . If there are two distributions p_0 and \tilde{p}_0 that agree on their marginals (i.e. crops), they can be used interchangeably for low-diffusion times. Note that the distributions don't have to agree globally, they only have to agree on a local (patch) level. Formally, let $A(t)$ be a random patch selector of size $\text{crop}(t)$. Let also p_0, \tilde{p}_0 two distributions that satisfy:

$$A(t) \# p_0 = A(t) \# \tilde{p}_0, \quad (3.3)$$

where $A(t) \# p_0$ denotes the pushforward measure¹ of p_0 under $A(t)$. Then, the cropped portions of the tilted distributions provide equivalent information to the crops of the original distribution for denoising.

Training a crops classifier. Note that the condition of Equation (3.3) can be trivially satisfied if $A(t)$ masks all the pixels or even if $A(t)$ just selects a single pixel. We are interested in finding what is the maximum crop size for which this condition is approximately true. Once again, we can use a classifier to solve this task. The input to the classifier, $c_{\theta}^{\text{crops}}$, is a crop of an image that either arises from p_0 or \tilde{p}_0 , and the classifier needs to classify between these two cases.

¹Given measure spaces (X_1, Σ_1) and (X_2, Σ_2) , a measurable function $f : X_1 \rightarrow X_2$, and a probability measure $p : \Sigma_1 \rightarrow [0, \infty)$, the pushforward measure $f \# p$ is defined as $(f \# p)(B) := p(f^{-1}(B)) \forall B \in \Sigma_2$.

186 **Annotation and training using the trained classifier.** Once the classifier is trained, we are now
 187 interested in finding the biggest crop size for which the distributions p_0, \tilde{p}_0 cannot be confidently
 188 distinguished. Formally,

$$t_n^{\max} = \sup \left\{ t \in [0, T] : \frac{1}{|S_B|} \sum_{y_0 \in S_B} [c_\theta^{\text{crops}}(A(t)(y_t))] > \tau \right\}, \quad (3.4)$$

189 for $\tau = 0.5 - \epsilon$ and for some small $\epsilon > 0$ ². For times $t \leq t_n^{\max}$, the out-of-distribution images from
 190 \tilde{p}_0 can be used with the regular diffusion objective as images from p_0 , as for these times the denoiser
 191 only looks at crops and at the crop level the distributions have converged.

192 **The donut paradox.** Each sample can be used for $t \geq t_i^{\min}$ and for $t \leq t_i^{\max}$, but not for $t \in$
 193 (t_i^{\max}, t_i^{\min}) . We call this the *donut paradox* as there is a hole in the middle of the diffusion trajectory
 194 for which we have fewer available data. These times do not have enough noise for the distributions
 195 to merge globally, but also the required receptive field for denoising is big enough so that there are
 196 differences on a crop level. We show an example of this effect in Figure 10.

Table 1: ImageNet results with and without classifier-free guidance.

ImageNet-512	Train FID ↓				Test FID ↓				Model Size	
	FID		FIDv2		FID		FIDv2		Mparams	NFE
	no CFG	w/ CFG	no CFG	w/ CFG	no CFG	w/ CFG	no CFG	w/ CFG		
EDM2-XS	3.57	2.91	103.39	79.94	3.77	3.68	115.16	93.86	125	63
Ambient-o-XS	3.59	2.89	107.26	79.56	3.69	3.58	115.02	92.96	125	63
EDM2-XXL	1.91 (1.93)	1.81	42.84	33.09	2.88	2.73	56.42	46.22	1523	63
Ambient-o-XXL	1.99	1.87	43.38	33.34	2.81	2.68	56.40	46.02	1523	63
Ambient-o-XXL+crops	1.91	1.80	42.84	32.63	2.78	2.53	56.39	45.78	1523	63

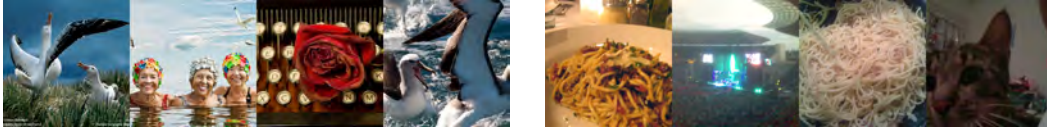


Figure 3: Results using CLIP to obtain the high-quality and the low-quality sets of ImageNet.

197 4 Experiments

198 **Controlled experiments to show utility from low-quality data.** To verify our method, we first
 199 do synthetic experiments on artificially corrupted data. We use EDM [34] as our baseline, and we
 200 train networks on CIFAR-10 and FFHQ. For the first experiments, we only use the high-noise part of
 201 our Ambient-o method (Section 3.1). We underline that for all of our experiments, we only change
 202 the way we use the data, and we keep all the optimization and network hyperparameters as is. We
 203 compare against using all the data as equal (despite the corruption) and the filtering strategy of
 204 only training on the clean samples. For evaluation, we measure FID [29] with respect to the full
 205 uncorrupted dataset (which is not available during training).

206 For the blurring experiments, we use a Gaussian kernel with standard deviation $\sigma_B =$
 207 0.4, 0.6, 0.8, 1.0, and we corrupt 90% of the data. We show some corrupted images in Appendix
 208 Figure 5a. To obtain annotations, we train a blurry vs clean image classifier under noise, as explained
 209 in Section 3.1. For the experiments in the main paper, we use a balanced dataset for the training
 210 of the classifier. We ablate the effect of having fewer training samples for the classifier training in
 211 Appendix Section F where we show that reducing the number of clean samples available for classifier
 212 training leads to a small drop in performance. Once equipped with the trained classifier, each sample
 213 is annotated on its own based on the amount of noise that is needed to confuse the classifier (sample
 214 dependent annotation). We present our results in Table 2a. As shown, for all corruption strengths,
 215 Ambient Omni, significantly outperforms the two baseline methods. In the one to the last column of
 216 Table 2a, we further show the average annotation of the classifier. As expected, the average assigned
 217 noise level increases as the corruption intensifies.

²We subtract an ϵ to allow for approximate mixing of the two distributions and hence smaller annotation times.

Table 2: In a controlled experiment with restricted access only to 10% of the clean dataset, our method of Ambient-o uses corrupted and out-of-distribution data to improve performance.

(a) Gaussian blurred data at different levels.

(b) Additional out-of-distribution data.

Method	Parameters Values (σ_B)	$\bar{\sigma}_n^{\min}$	FID	Source Data	Additional Data	Method	$\bar{\sigma}_n^{\max}$	FID
Only Clean (10%)	-	-	8.79	Dogs (10%)	None	-	-	12.08
All data	1.0	0	45.32		Cats	Fixed σ	0.2	11.14
	0.8		28.26		Cats	Fixed σ	0.1	<u>9.85</u>
	0.6		11.42		Cats	Fixed σ	0.05	10.66
	0.4		2.47		Cats	Fixed σ	0.025	12.07
Ambient-o	1.0	2.84	6.16		Cats	Classifier	0.09	8.92
	0.8	1.93	6.00		Procedural	Classifier	0.042	<u>10.98</u>
	0.6	1.38	5.34	Cats (10%)	None	-	-	5.20
	0.4	0.22	2.44		Dogs	Classifier	0.13	5.11
				Wildlife	Classifier	0.08	4.89	

Controlled experiments to show utility from out-of-distribution images. We now want to validate the method developed in Section 3.2 for leveraging out-of-distribution data. To start with, we want to find the mapping between diffusion times and the size of the receptive field required for an optimal denoising prediction. To do so, we take a pre-trained denoising diffusion model and measure the denoising loss at a given location as we increase the size of the context. We provide the corresponding plot in the Supplemental Figures 13, 11. The main finding is that while providing more context always leads to a decrease in the average loss, for sufficiently small noise levels, the loss nearly plateaus before the full image context is provided. That implies that the perfect denoiser for a given noise level only needs to look at a localized part of the image.

Equipped with the mapping between diffusion times and crop sizes, we now proceed to a fun experiment. Specifically, we show that it is possible to use images of cats to improve a generative model for dogs (!) and vice-versa. The cats here represent out-of-distribution data that can be used to improve the performance in the distribution of interest (in our toy example, dogs distribution). To perform this experiment, we train a classifier that discriminates between cats and dog images by looking at crops of various sizes (Section 3.2). Figure 18 shows the predictions of an 8×8 crops-classifier for an image of a cat, illustrating that there are a number of crops that are misclassified as crops from a dog image. We report results for this experiment in Table 2b and we observe improvements in FID arising from using out-of-distribution data. Beyond natural images, we show that it is even possible to use procedurally generated data from Shaders [6] to (slightly) improve the performance. Figure 19 shows an example of such an image and the corresponding predictions of a crops classifier. Table 2b contains more results and ablations between annotating all the out-of-distribution at a single noise level vs. sample-dependent annotations.

Corruptions of natural datasets – ImageNet results. Up to this point, our corrupted data has been artificially constructed to study our method in a controlled setting. However, it turns out that even in real datasets such as ImageNet, there are images with significant degradations such as heavy blur, low lighting, and low contrast, and also images with fantastic detail, clear lightning, and sharp contrast. Here, the high-quality and the low-quality sets are not given and hence we have to estimate them. We opt to use the CLIP-IQA quality metric [66] to separate ImageNet into high-quality (top 10% CLIP-IQA) and low-quality (bottom 90% CLIP-IQA) sets. Figure 3 shows some of the top and bottom quality images according to our metric. Given the high-quality and low-quality sets, we are now back to the previous setting where we can use the developed Ambient-o methodology. We underline that there is a rich literature regarding quality-assessment methods [69, 68, 49, 67].

We use Ambient-o to refer to our method that uses low-quality data at high diffusion times (Section 4) and Ambient-o+crops to refer to the extended version of our method that uses crops from potentially low-quality images at low-diffusion times. Perhaps surprisingly, there are images in ImageNet that have lower global quality but high-quality crops that we can use for low-noise. We present results in Table 1, where we show the best FID [29] and FD_{DINOv2} obtained by different methods. We show the highest and lowest quality crops, alongside their associated full images, of ImageNet according to CLIP in Figure 15.

As shown in the Table, our method leads to state-of-the-art FID scores, improving over the previous state-of-the-art baseline EDM-2 [35] at both the low and high parameter count settings. The benefits are more pronounced when we measure test FID as our method memorizes significantly less due to the addition of noise during the annotation stage of our pipeline (Section 3.1). Beyond FID, we

provide qualitative results in Figure 1 (bottom) and Appendix Figures 7, 8. We further show that the quality of the generated images measured by CLIP increased compared to the baseline in Appendix Table 5. The observed improvements are proof that the ability to learn from data with heterogeneous qualities can be truly impactful for realistic settings beyond synthetic corruptions typically studied in prior work.

Text-to-image results. For our final set of experiments, we show how Ambient-o can be used to improve the performance of text-to-image diffusion models. We use the code-base of MicroDiffusion [54], as it is open-data and trainable with modest compute (≈ 2 days on 8-H100 GPUs). Schwag et al. [54] use four main datasets to train their model: Conceptual Captions (12M) [56], Segment Anything (11M) [41], JourneyDB (4.2M) [63], and DiffusionDB (10.7M) [70]. Of these four, DiffusionDB is of significantly lower quality than the others as it contains solely synthetic data from an outdated diffusion model. This presents an opportunity for the use of our method. Can we use this lower-quality data and improve the performance of the trained network?

We set $\sigma_{\min} = 2$ for all samples from DiffusionDB and $\sigma_{\min} = 0$ for all other datasets and we train a diffusion model with Ambient-o. We note that we did not ablate this hyperparameter and it is quite likely that improved results would be obtained by tuning it or by training a high-quality vs low-quality data classifier for the annotation. Despite that, our trained model achieves a remarkable FID of **10.61** in COCO, significantly improving the baseline FID of 12.37 (Table 4). We present qualitative results in Figure 1 and GPT-4o evaluations on DrawBench and PartiPrompt in Figure 23. Ambient-o and baseline generations for different prompts can be found in Figure 1.

As an additional ablation, we compared our method with the recipe of doing a final fine-tuning on the highest-quality subset, as done in the works of [54, 13]. Compared to this baseline, our method obtained slightly worse COCO FID (10.61 vs 10.27) but obtained much greater diversity, as seen visually in Figure 29 and quantitatively through $> 13\%$ increases in DINO Vendi Diversity on prompts from DiffDB (3.22 vs 3.65.). This corroborates our intuition that data filtration leads to decreased diversity. Ambient-o uses all the data but can strike a fine balance between high-quality and diverse generation.

(a) Measuring fidelity and prompt alignment of generated images on COCO dataset.

Method	FID-30K (\downarrow)	Clip-FD-30K (\downarrow)	Clip-score (\uparrow)
Baseline	12.37	10.07	0.345
Ambient-o	10.61	9.40	0.348

(b) Measuring performance on the GenEval benchmark.

Method	Overall	Objects		Counting	Colors	Position	Color attribution
		Single	Two				
Baseline	0.44	0.97	0.33	0.35	0.82	0.06	0.14
Ambient-o	0.47	0.97	0.40	0.36	0.82	0.11	0.14

Figure 4: Quantitative benefits of Ambient-o on COCO [44] zero-shot generation and GenEval [26].

5 Conclusion

Is it possible to get good generative models from bad data? Our framework extracts value from low-quality, synthetic, and out-of-distribution sources. At a time when the ever-growing data demands of GenAI are at odds with the need for quality control, Ambient-o lights a path for both to be achieved simultaneously.

References

- [1] Asad Aali, Marius Arvinte, Sidharth Kumar, and Jonathan I Tamir. “Solving Inverse Problems with Score-Based Generative Priors learned from Noisy Data”. In: *arXiv preprint arXiv:2305.01166* (2023) (cit. on pp. 1, 3).
- [2] Asad Aali, Giannis Daras, Brett Levac, Sidharth Kumar, Alex Dimakis, and Jon Tamir. “Ambient Diffusion Posterior Sampling: Solving Inverse Problems with Diffusion Models Trained on Corrupted Data”. In: *The Thirteenth International Conference on Learning Representations*. 2025. URL: <https://openreview.net/forum?id=qeXcMutEZY> (cit. on p. 1).
- [3] Sina Alemohammad, Josue Casco-Rodriguez, Lorenzo Luzi, Ahmed Imtiaz Humayun, Hossein Babaei, Daniel LeJeune, Ali Siahkoobi, and Richard G Baraniuk. “Self-consuming generative models go mad”. In: *arXiv preprint arXiv:2307.01850* 4 (2023), p. 14 (cit. on p. 3).

- [4] Sina Alemohammad, Ahmed Imtiaz Humayun, Shruti Agarwal, John Collomosse, and Richard Baraniuk. “Self-improving diffusion models with synthetic data”. In: *arXiv preprint arXiv:2408.16333* (2024) (cit. on p. 3).
- [5] Weimin Bai, Yifei Wang, Wenzheng Chen, and He Sun. “An Expectation-Maximization Algorithm for Training Clean Diffusion Models from Corrupted Observations”. In: *arXiv preprint arXiv:2407.01014* (2024) (cit. on pp. 1, 3).
- [6] Manel Baradad, Chun-Fu Chen, Jonas Wulff, Tongzhou Wang, Rogerio Feris, Antonio Torralba, and Phillip Isola. “Procedural Image Programs for Representation Learning”. In: *Advances in Neural Information Processing Systems*. Ed. by Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho. 2022. URL: <https://openreview.net/forum?id=wJwHTgIoEOP> (cit. on p. 7).
- [7] Ashish Bora, Eric Price, and Alexandros G Dimakis. “AmbientGAN: Generative models from lossy measurements”. In: *International conference on learning representations*. 2018 (cit. on p. 1).
- [8] Zdravko I Botev, Joseph F Grotowski, and Dirk P Kroese. “Kernel density estimation via diffusion”. In: (2010) (cit. on p. 13).
- [9] Sitan Chen, Sinho Chewi, Jerry Li, Yuanzhi Li, Adil Salim, and Anru R Zhang. “Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions”. In: *arXiv preprint arXiv:2209.11215* (2022) (cit. on p. 1).
- [10] Sitan Chen, Giannis Daras, and Alex Dimakis. “Restoration-Degradation Beyond Linear Diffusions: A Non-Asymptotic Analysis For DDIM-type Samplers”. In: *Proceedings of the 40th International Conference on Machine Learning*. Ed. by Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett. Vol. 202. Proceedings of Machine Learning Research. PMLR, 23–29 Jul 2023, pp. 4462–4484. URL: <https://proceedings.mlr.press/v202/chen23e.html> (cit. on p. 1).
- [11] Tianyu Chen, Yasi Zhang, Zhendong Wang, Ying Nian Wu, Oscar Leong, and Mingyuan Zhou. *Denoising Score Distillation: From Noisy Diffusion Pretraining to One-Step High-Quality Generation*. 2025. arXiv: 2503.07578 [cs.LG]. URL: <https://arxiv.org/abs/2503.07578> (cit. on p. 1).
- [12] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. “StarGAN v2: Diverse image synthesis for multiple domains”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020, pp. 8188–8197 (cit. on p. 21).
- [13] Xiaoliang Dai et al. *Emu: Enhancing Image Generation Models Using Photogenic Needles in a Haystack*. 2023. arXiv: 2309.15807 [cs.CV] (cit. on pp. 1, 3, 8).
- [14] Giannis Daras, Yeshwanth Cherapanamjeri, and Constantinos Costis Daskalakis. “How Much is a Noisy Image Worth? Data Scaling Laws for Ambient Diffusion.” In: *The Thirteenth International Conference on Learning Representations*. 2025. URL: <https://openreview.net/forum?id=qZwtPEw2qN> (cit. on p. 1).
- [15] Giannis Daras, Yuval Dagan, Alexandros G Dimakis, and Constantinos Daskalakis. “Consistent diffusion models: Mitigating sampling drift by learning to be consistent”. In: *arXiv preprint arXiv:2302.09057* (2023) (cit. on pp. 1, 3).
- [16] Giannis Daras, Mauricio Delbracio, Hossein Talebi, Alex Dimakis, and Peyman Milanfar. “Soft Diffusion: Score Matching with General Corruptions”. In: *Transactions on Machine Learning Research* (2023). ISSN: 2835-8856. URL: <https://openreview.net/forum?id=W98rebBxlQ> (cit. on p. 1).
- [17] Giannis Daras, Alexandros G Dimakis, and Constantinos Daskalakis. “Consistent Diffusion Meets Tweedie: Training Exact Ambient Diffusion Models with Noisy Data”. In: *arXiv preprint arXiv:2404.10177* (2024) (cit. on pp. 1, 3).
- [18] Giannis Daras, Kulin Shah, Yuval Dagan, Aravind Gollakota, Alex Dimakis, and Adam Klivans. “Ambient Diffusion: Learning Clean Distributions from Corrupted Data”. In: *Thirty-seventh Conference on Neural Information Processing Systems*. 2023. URL: <https://openreview.net/forum?id=wBJBLy9kBY> (cit. on pp. 1, 3).
- [19] Mauricio Delbracio and Peyman Milanfar. “Inversion by direct iteration: An alternative to denoising diffusion for image restoration”. In: *arXiv preprint arXiv:2303.11435* (2023) (cit. on p. 1).

- [20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. “ImageNet: A large-scale hierarchical image database”. In: *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 2009, pp. 248–255. DOI: [10.1109/CVPR.2009.5206848](https://doi.org/10.1109/CVPR.2009.5206848) (cit. on p. 21).
- [21] Prafulla Dhariwal and Alexander Nichol. “Diffusion models beat gans on image synthesis”. In: *Advances in neural information processing systems* 34 (2021), pp. 8780–8794 (cit. on p. 4).
- [22] Sander Dieleman. *Diffusion is spectral autoregression*. 2024. URL: <https://sander.ai/2024/09/02/spectral-autoregression.html> (cit. on p. 2).
- [23] Logan Engstrom, Andrew Ilyas, Benjamin Chen, Axel Feldmann, William Moses, and Aleksander Madry. “Optimizing ml training with metagradient descent”. In: *arXiv preprint arXiv:2503.13751* (2025) (cit. on p. 3).
- [24] Damien Ferbach, Quentin Bertrand, Avishek Joey Bose, and Gauthier Gidel. “Self-consuming generative models with curated data provably optimize human preferences”. In: *arXiv preprint arXiv:2407.09499* (2024) (cit. on p. 3).
- [25] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. “DataComp: In search of the next generation of multimodal datasets”. In: *arXiv preprint arXiv:2304.14108* (2023) (cit. on pp. 1, 3).
- [26] Dhruva Ghosh, Hanna Hajishirzi, and Ludwig Schmidt. *GenEval: An Object-Focused Framework for Evaluating Text-to-Image Alignment*. 2023. arXiv: [2310.11513](https://arxiv.org/abs/2310.11513) [cs.CV]. URL: <https://arxiv.org/abs/2310.11513> (cit. on p. 8).
- [27] Sachin Goyal, Pratyush Maini, Zachary C Lipton, Aditi Raghunathan, and J Zico Kolter. “Scaling Laws for Data Filtering—Data Curation cannot be Compute Agnostic”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2024, pp. 22702–22711 (cit. on p. 1).
- [28] Dan Hendrycks and Thomas Dietterich. “Benchmarking neural network robustness to common corruptions and perturbations”. In: *arXiv preprint arXiv:1903.12261* (2019) (cit. on p. 1).
- [29] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. “Gans trained by a two time-scale update rule converge to a local nash equilibrium”. In: *Advances in neural information processing systems* 30 (2017) (cit. on pp. 6, 7).
- [30] Jonathan Ho, Ajay Jain, and Pieter Abbeel. “Denoising diffusion probabilistic models”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 6840–6851 (cit. on p. 1).
- [31] Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. “Adaptive Mixtures of Local Experts”. In: *Neural Computation* 3.1 (Mar. 1991). eprint: <https://direct.mit.edu/neco/article-pdf/3/1/79/812104/neco.1991.3.1.79.pdf>, pp. 79–87. ISSN: 0899-7667. DOI: [10.1162/neco.1991.3.1.79](https://doi.org/10.1162/neco.1991.3.1.79). URL: <https://doi.org/10.1162/neco.1991.3.1.79> (cit. on p. 22).
- [32] Yiding Jiang, Allan Zhou, Zhili Feng, Sadhika Malladi, and J Zico Kolter. “Adaptive data optimization: Dynamic sample selection with scaling laws”. In: *arXiv preprint arXiv:2410.11820* (2024) (cit. on p. 1).
- [33] Mason Kamb and Surya Ganguli. “An analytic theory of creativity in convolutional diffusion models”. In: *arXiv preprint arXiv:2412.20292* (2024) (cit. on p. 13).
- [34] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. “Elucidating the design space of diffusion-based generative models”. In: *arXiv preprint arXiv:2206.00364* (2022) (cit. on pp. 6, 21).
- [35] Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. “Analyzing and Improving the Training Dynamics of Diffusion Models”. In: *Proc. CVPR*. 2024 (cit. on pp. 2, 7, 19–21).
- [36] Tero Karras, Samuli Laine, and Timo Aila. “A style-based generator architecture for generative adversarial networks”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 4401–4410 (cit. on p. 21).
- [37] Sergey Kastryulin, Dzhamil Zakirov, and Denis Prokopenko. *PyTorch Image Quality: Metrics and Measure for Image Quality Assessment*. Open-source software available at <https://github.com/photosynthesis-team/piq>. 2019. URL: <https://github.com/photosynthesis-team/piq> (cit. on p. 18).
- [38] Sergey Kastryulin, Jamil Zakirov, Denis Prokopenko, and Dmitry V. Dylov. *PyTorch Image Quality: Metrics for Image Quality Assessment*. 2022. DOI: [10.48550/ARXIV.2208.14818](https://doi.org/10.48550/ARXIV.2208.14818). URL: <https://arxiv.org/abs/2208.14818> (cit. on p. 18).

- [39] Varun A Kelkar, Rucha Deshpande, Arindam Banerjee, and Mark A Anastasio. “Ambient-Flow: Invertible generative models from incomplete, noisy measurements”. In: *arXiv preprint arXiv:2309.04856* (2023) (cit. on p. 1).
- [40] Diederik P Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *arXiv preprint arXiv:1412.6980* (2014) (cit. on pp. 21, 22).
- [41] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. *Segment Anything*. 2023. arXiv: 2304.02643 [cs.CV]. URL: <https://arxiv.org/abs/2304.02643> (cit. on pp. 8, 21).
- [42] Alex Krizhevsky and Geoffrey Hinton. “Learning multiple layers of features from tiny images”. In: (2009) (cit. on p. 21).
- [43] Jeffrey Li et al. *DataComp-LM: In search of the next generation of training sets for language models*. 2024. arXiv: 2406.11794 [cs.LG] (cit. on pp. 1, 3).
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. *Microsoft COCO: Common Objects in Context*. 2015. arXiv: 1405.0312 [cs.CV]. URL: <https://arxiv.org/abs/1405.0312> (cit. on p. 8).
- [45] Yvette Y Lin, Angela F Gao, and Katherine L Bouman. “IMAGING AN EVOLVING BLACK HOLE BY LEVERAGING SHARED STRUCTURE”. In: *ICASSP* (2024) (cit. on p. 1).
- [46] Zeyuan Liu, Zhihe Yang, Jiawei Xu, Rui Yang, Jiafei Lyu, Baoxiang Wang, Yunjian Xu, and Xiu Li. “ADG: Ambient Diffusion-Guided Dataset Recovery for Corruption-Robust Offline Reinforcement Learning”. In: *arXiv preprint arXiv:2505.23871* (2025) (cit. on p. 1).
- [47] Haoye Lu, Qifan Wu, and Yaoliang Yu. “SFBD: A Method for Training Diffusion Models with Noisy Data”. In: *Frontiers in Probabilistic Inference: Learning meets Sampling*. 2025. URL: <https://openreview.net/forum?id=6HN14zuHRb> (cit. on p. 1).
- [48] Haoye Lu, Qifan Wu, and Yaoliang Yu. *Stochastic Forward-Backward Deconvolution: Training Diffusion Models with Finite Noisy Datasets*. 2025. arXiv: 2502.05446 [cs.LG]. URL: <https://arxiv.org/abs/2502.05446> (cit. on p. 1).
- [49] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. “Making a “completely blind” image quality analyzer”. In: *IEEE Signal processing letters* 20.3 (2012), pp. 209–212 (cit. on p. 7).
- [50] William Peebles and Saining Xie. *Scalable Diffusion Models with Transformers*. 2023. arXiv: 2212.09748 [cs.CV]. URL: <https://arxiv.org/abs/2212.09748> (cit. on p. 22).
- [51] Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, Thomas Wolf, et al. “The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale”. In: *arXiv preprint arXiv:2406.17557* (2024) (cit. on p. 1).
- [52] François Rozet, G r me Andry, Fran ois Lanusse, and Gilles Louppe. “Learning Diffusion Priors from Observations by Expectation Maximization”. In: *arXiv preprint arXiv:2405.13712* (2024) (cit. on pp. 1, 3).
- [53] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. “Laion-5b: An open large-scale dataset for training next generation image-text models”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 25278–25294 (cit. on p. 1).
- [54] Vikash Sehwal, Xianghao Kong, Jingtao Li, Michael Spranger, and Lingjuan Lyu. “Stretching Each Dollar: Diffusion Training from Scratch on a Micro-Budget”. In: *arXiv preprint arXiv:2407.15811* (2024) (cit. on pp. 2, 3, 8, 21, 22).
- [55] Kulin Shah, Alkis Kalavasis, Adam R. Klivans, and Giannis Daras. *Does Generation Require Memorization? Creative Diffusion Models using Ambient Diffusion*. 2025. arXiv: 2502.21278 [cs.LG] (cit. on p. 1).
- [56] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. “Conceptual Captions: A Cleaned, Hypernymed, Image Alt-text Dataset For Automatic Image Captioning”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by Iryna Gurevych and Yusuke Miyao. Melbourne, Australia: Association for Computational Linguistics, July 2018, pp. 2556–2565. DOI: 10.18653/v1/P18-1238. URL: <https://aclanthology.org/P18-1238/> (cit. on pp. 8, 21).

- [57] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. *Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer*. 2017. arXiv: [1701.06538](https://arxiv.org/abs/1701.06538) [cs.LG]. URL: <https://arxiv.org/abs/1701.06538> (cit. on p. 22).
- [58] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. “Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models”. In: *arXiv preprint arXiv:2212.03860* (2022) (cit. on p. 1).
- [59] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. “Understanding and Mitigating Copying in Diffusion Models”. In: *arXiv preprint arXiv:2305.20086* (2023) (cit. on p. 1).
- [60] Jiaming Song, Chenlin Meng, and Stefano Ermon. “Denoising diffusion implicit models”. In: *arXiv preprint arXiv:2010.02502* (2020) (cit. on p. 21).
- [61] Yang Song and Stefano Ermon. “Generative modeling by estimating gradients of the data distribution”. In: *Advances in Neural Information Processing Systems* 32 (2019) (cit. on p. 1).
- [62] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. “Score-based generative modeling through stochastic differential equations”. In: *arXiv preprint arXiv:2011.13456* (2020) (cit. on p. 1).
- [63] Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, Jifeng Dai, Yu Qiao, Limin Wang, and Hongsheng Li. *JourneyDB: A Benchmark for Generative Image Understanding*. 2023. arXiv: [2307.00716](https://arxiv.org/abs/2307.00716) [cs.CV]. URL: <https://arxiv.org/abs/2307.00716> (cit. on pp. 8, 21).
- [64] Ayush Tewari, Tianwei Yin, George Cazenavette, Semon Rezchikov, Josh Tenenbaum, Frédo Durand, Bill Freeman, and Vincent Sitzmann. “Diffusion with forward models: Solving stochastic inverse problems without direct supervision”. In: *Advances in Neural Information Processing Systems* 36 (2023), pp. 12349–12362 (cit. on p. 1).
- [65] Antonio Torralba, Phillip Isola, and William T Freeman. *Foundations of computer vision*. MIT Press, 2024 (cit. on p. 2).
- [66] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. “Exploring CLIP for Assessing the Look and Feel of Images”. In: *AAAI*. 2023 (cit. on pp. 7, 21).
- [67] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. “Exploring clip for assessing the look and feel of images”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 37. 2. 2023, pp. 2555–2563 (cit. on p. 7).
- [68] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. “Image quality assessment: from error visibility to structural similarity”. In: *IEEE transactions on image processing* 13.4 (2004), pp. 600–612 (cit. on p. 7).
- [69] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. “Multiscale structural similarity for image quality assessment”. In: *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*. Vol. 2. Ieee. 2003, pp. 1398–1402 (cit. on p. 7).
- [70] Zijie J Wang, Evan Montoya, David Munechika, Haoyang Yang, Benjamin Hoover, and Duen Horng Chau. “DiffusionDB: A Large-scale Prompt Gallery Dataset for Text-to-Image Generative Models”. In: *arXiv preprint arXiv:2210.14896* (2022) (cit. on pp. 8, 21).
- [71] Yasi Zhang, Tianyu Chen, Zhendong Wang, Ying Nian Wu, Mingyuan Zhou, and Oscar Leong. “Restoration Score Distillation: From Corrupted Diffusion Pretraining to One-Step High-Quality Generation”. In: *arXiv preprint arXiv:2505.13377* (2025) (cit. on p. 1).

514 A Limitations and Future Work

515 Our work opens several avenues for improvement. On the theoretical side, we aim to establish
 516 matching lower bounds to demonstrate that learning from the mixture distribution becomes provably
 517 optimal beyond a certain noise threshold. Algorithmically, while our method performs well under
 518 high-frequency corruptions, it remains an open question whether more effective training strategies
 519 could be used for different types of corruptions (e.g., masking). Moreover, real-world datasets often
 520 exhibit patch-wise heterogeneity—for example, facial regions are frequently blurred for privacy,
 521 leading to uneven corruption across image crops. We plan to investigate patch-level noise annotations
 522 to better capture this structure in future work. Computationally, the full-version of our algorithm
 523 requires the training of classifiers for annotations that increases the runtime. This overhead can be

avoided by using hand-picked annotation times based on quality proxies as done in our synthetic data experiment. Finally, we believe the true potential of Ambient-o lies in scientific applications, where data often arises from heterogeneous measurement processes.

B Theory

We study the 1-d case, but all our claims easily extend to any dimension. We compare two algorithms:

Algorithm 1. Algorithm 1 trains a diffusion model using access to n_1 samples from a target density p_0 , assumed to be supported in $[0, 1]$ and be λ_1 -Lipschitz.

Algorithm 2. Algorithm 2 trains a diffusion model using access to $n_1 + n_2$ samples from a density \tilde{p}_0 that is a mixture of the a target density p_0 and another density q_0 , assumed to be supported in $[0, 1]$ and be λ_2 -Lipschitz: $\tilde{p}_0 = \frac{n_1}{n_1+n_2}p_0 + \frac{n_2}{n_1+n_2}q_0$.

We want to compare how well these algorithms estimate the distribution $p_t := p_0 \otimes \mathcal{N}(0, \sigma_t^2)$. We use $\hat{p}_t^{(1)}, \hat{p}_t^{(2)}$ to denote the estimates obtained for p_t by Algorithms 1 and 2 respectively.

Diffusion modeling is Gaussian kernel density estimation. We start by making a connection between the optimal solution to the diffusion modeling objective and kernel density estimation. Given a finite dataset $\{W^{(i)}\}_{i=1}^n$, the optimal solution to the diffusion modeling objective should match the empirical density at time t , which is:

$$\hat{p}_t(x) = \frac{1}{n\sigma_t} \sum_{i=1}^n \phi\left(\frac{W^{(i)} - x}{\sigma_t}\right), \quad (\text{B.1})$$

where $\phi(u) = \frac{1}{\sqrt{2\pi}}e^{-u^2/2}$ is the Gaussian kernel. We observe that equation B.1 is identical to a Gaussian kernel density estimate, given samples $\{W^{(i)}\}_{i=1}^n$ ³.

We establish the following result for Gaussian kernel density estimation.

Theorem B.1 (Gaussian Kernel Density Estimation). *Let $\{W^{(i)}\}_{i=1}^n$ be a set of n independent samples from a λ -Lipschitz density p . Let \hat{p} be the empirical density, $p_\sigma := p \otimes \mathcal{N}(0, \sigma^2)$ and $\hat{p}_\sigma = \hat{p} \otimes \mathcal{N}(0, \sigma^2)$. Then, with probability at least $1 - \delta$ with respect to the sample randomness,*

$$\text{d}_{\text{TV}}(p_\sigma, \hat{p}_\sigma) \lesssim \frac{1}{n} + \frac{1}{\sigma^2 n} + \sqrt{\frac{\log n + \log(1 \vee \lambda) + \log 2/\delta}{\sigma^2 n}}. \quad (\text{B.2})$$

The proof of this result is given in the Appendix.

Comparing the performance of Algorithms 1 and 2. Applying Theorem B.1 directly to the p_0 density, we immediately get that the estimate $\hat{p}_t^{(1)}(x)$ obtained by Algorithm 1 satisfies:

$$\text{d}_{\text{TV}}(p_t, \hat{p}_t^{(1)}) \lesssim \frac{1}{n_1} + \frac{1}{\sigma_t^2 n_1} + \sqrt{\frac{\log n_1 + \log(1 \vee \lambda_1) + \log 2/\delta}{\sigma_t^2 n_1}}. \quad (\text{B.3})$$

Let us now see what we get by applying Theorem B.1 to Algorithm 2, which uses samples from the tilted distribution \tilde{p}_0 . Since this distribution is $\left(\frac{n_1}{n_1+n_2}\lambda_1 + \frac{n_2}{n_1+n_2}\lambda_2\right)$ -Lipschitz, we get that:

$$\text{d}_{\text{TV}}(\tilde{p}_t, \hat{p}_t^{(2)}) \lesssim \frac{1}{(n_1 + n_2)} + \frac{1}{\sigma_t^2 (n_1 + n_2)} + \sqrt{\frac{\log(n_1 + n_2) + \log(1 \vee \frac{n_1}{n_1+n_2}\lambda_1 + \frac{n_2}{n_1+n_2}\lambda_2) + \log 2/\delta}{\sigma_t^2 (n_1 + n_2)}},$$

where $\tilde{p}_t := \tilde{p}_0 \otimes \mathcal{N}(0, \sigma_t^2)$.

Further, we have that: $\text{d}_{\text{TV}}(p_t, \hat{p}_t^{(2)}) \leq \text{d}_{\text{TV}}(\tilde{p}_t, p_t) + \text{d}_{\text{TV}}(\tilde{p}_t, \hat{p}_t^{(2)})$. We already have a bound for the second term. To bound the first term, we prove the following theorem.

³This connection has been observed in prior works too, e.g., see [33, 8].

Theorem B.2 (Distance contraction under noise). *Consider distributions P and Q supported on a subset of \mathbb{R}^d with diameter D . Then*

$$d_{\text{TV}}(P \otimes \mathcal{N}(0, \sigma^2 \mathbf{I}), Q \otimes \mathcal{N}(0, \sigma^2 \mathbf{I})) \leq d_{\text{TV}}(P, Q) \cdot \frac{D}{2\sigma}.$$

554 Applying this theorem we get that: $d_{\text{TV}}(\tilde{p}_t, p_t) \leq \frac{1}{2\sigma_t} d_{\text{TV}}(\tilde{p}_0, p_0) \leq \frac{1}{2\sigma_t} \cdot \frac{n_2}{n_1+n_2} d_{\text{TV}}(p_0, q_0)$, where
 556 for the second inequality we used that $d_{\text{TV}}(p_0, \tilde{p}_0) \leq \frac{n_2}{n_1+n_2} d_{\text{TV}}(p_0, q_0)$.

557 Putting everything together, Algorithm (2) achieves an estimation error:

$$d_{\text{TV}}(p_t, \hat{p}_t^{(2)}) \lesssim \frac{1}{(n_1+n_2)} + \frac{1}{\sigma_t^2(n_1+n_2)} + \sqrt{\frac{\log(n_1+n_2) + \log(1 \vee \frac{n_1}{n_1+n_2}\lambda_1 + \frac{n_2}{n_1+n_2}\lambda_2) + \log 2/\delta}{\sigma_t^2(n_1+n_2)}} + \frac{n_2}{\sigma_t(n_1+n_2)} d_{\text{TV}}(p_0, q_0).$$

558 Comparing this with the bound obtained in Equation B.3, we see that if n_2 is sufficiently larger than
 559 n_1 or if $\lambda_2 \leq \lambda_1$, there is a t_n^{\min} such that for any $t \geq t_n^{\min}$, the upper-bound obtained by Algorithm
 560 2 is better than the upper-bound obtained by Algorithm 1. That implies that for high-diffusion times,
 561 using biased data might be helpful for learning, as the bias term (final term) decays with the amount
 562 of noise. Going back to equation B, note that the switching point $t \geq t_n^{\min}$ depends on the distance
 563 $d_{\text{TV}}(\tilde{p}_t, p_t)$ that decays as shown in Theorem B.2. Once this distance becomes small enough, our
 564 computations above suggest that we benefit from biased data. The classifier of Section 3.1, if optimal,
 565 exactly tracks the distance $d_{\text{TV}}(\tilde{p}_t, p_t)$ and, as a result, tracks the switching point.

566 C Theoretical Results

567 C.1 Kernel Estimation

568 **Assumption C.1.** The density p is λ lipschitz.

569 Let $\{X^{(i)}\}_{i=1}^n$ a set of n independent samples from a density p that satisfies Assumption C.1. Let \hat{p}
 570 be the empirical density on those samples.

571 We are interested in bounding the total variation distance between $p_\sigma := p \otimes \mathcal{N}(0, \sigma^2)$ and $\hat{p}_\sigma =$
 572 $\hat{p} \otimes \mathcal{N}(0, \sigma^2)$. In particular,

$$\hat{p}_\sigma(x) = \frac{1}{n\sigma} \sum_{i=1}^n \phi\left(\frac{X^{(i)} - x}{\sigma}\right), \quad (\text{C.1})$$

573 where $\phi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}$ is the Gaussian kernel. We want to argue that the TV distance between
 574 p_σ and \hat{p}_σ is small given sufficiently many samples n . For simplicity, let's fix the support of p to be
 575 $[0, 1]$. We have:

$$d_{\text{TV}}(p_\sigma, \hat{p}_\sigma) = \frac{1}{2} \int_0^1 |p_\sigma(x) - \hat{p}_\sigma(x)| dx = \sum_{l=0}^{L-1} \int_{l/L}^{(l+1)/L} |p_\sigma(x) - \hat{p}_\sigma(x)| dx \quad (\text{C.2})$$

576 Now let us look at one of the terms of the summation.

$$\int_{l/L}^{(l+1)/L} |p_\sigma(x) - \hat{p}_\sigma(x)| dx = \int_{l/L}^{(l+1)/L} |p_\sigma(x) - p_\sigma(l/L) + p_\sigma(l/L) - \hat{p}_\sigma(x)| dx \quad (\text{C.3})$$

$$\leq \int_{l/L}^{(l+1)/L} |p_\sigma(x) - p_\sigma(l/L)| dx + \int_{l/L}^{(l+1)/L} |p_\sigma(l/L) - \hat{p}_\sigma(x)| dx. \quad (\text{C.4})$$

577 We first work on the first term. Using Lemma C.6:

$$\int_{l/L}^{(l+1)/L} |p_\sigma(x) - p_\sigma(l/L)| dx \leq \lambda \int_{l/L}^{(l+1)/L} |x - l/L| dx \quad (\text{C.5})$$

$$= \frac{\lambda}{2L^2}. \quad (\text{C.6})$$

578 Next, we work on the second term.

$$\int_{l/L}^{(l+1)/L} |p_\sigma(l/L) - \hat{p}_\sigma(x)| dx = \int_{l/L}^{(l+1)/L} |p_\sigma(l/L) - \hat{p}_\sigma(l/L) + \hat{p}_\sigma(l/L) - \hat{p}_\sigma(x)| dx \quad (\text{C.7})$$

$$\leq \int_{l/L}^{(l+1)/L} |p_\sigma(l/L) - \hat{p}_\sigma(l/L)| dx + \int_{l/L}^{(l+1)/L} |\hat{p}_\sigma(l/L) - \hat{p}_\sigma(x)| dx. \quad (\text{C.8})$$

579 According to Lemma C.5, we have that \hat{p}_σ is $\hat{\lambda} = \frac{1}{\sigma^2 \sqrt{2\pi e}}$ Lipschitz. Then, the second term becomes:

$$\int_{l/L}^{(l+1)/L} |\hat{p}_\sigma(l/L) - \hat{p}_\sigma(x)| dx \leq \hat{\lambda} \int_{l/L}^{(l+1)/L} |l/L - x| dx = \frac{\hat{\lambda}}{2L^2}. \quad (\text{C.9})$$

580 It remains to bound the following term

$$\int_{l/L}^{(l+1)/L} |p_\sigma(l/L) - \hat{p}_\sigma(l/L)| dx = \frac{|p_\sigma(l/L) - \hat{p}_\sigma(l/L)|}{L} \quad (\text{C.10})$$

581 We will be applying Hoeffding's Inequality, stated below:

582 **Theorem C.2** (Hoeffding's Inequality). *Let Y_1, \dots, Y_n be independent random variables in $[a, b]$ with*
 583 *mean μ . Then,*

$$\Pr \left(\left| \frac{1}{n} \sum_{i=1}^n Y_i - \mu \right| \geq t \right) \leq 2 \exp(-2nt^2/(b-a)^2). \quad (\text{C.11})$$

584 Recall that \hat{p}_σ can be written as

$$\hat{p}_\sigma(x) = \frac{1}{n} \sum_{i=1}^n \frac{\phi((X^{(i)} - x)/\sigma)}{\sigma} = \frac{1}{n} \sum_{i=1}^n Y_i, \quad (\text{C.12})$$

585 in terms of the random variables $Y_i := \frac{\phi((X^{(i)} - x)/\sigma)}{\sigma}$. These random variables are supported in
 586 $\left[0, \frac{1}{\sqrt{2\pi\sigma^2}}\right]$. So, for any x , we have that:

$$\Pr(|\hat{p}_\sigma(x) - \mathbb{E}[\hat{p}_\sigma(x)]| \geq t) \leq 2 \exp(-4\pi\sigma^2 nt^2). \quad (\text{C.13})$$

587 Taking $t = \sqrt{\frac{\log(2L/\delta)}{4\pi\sigma^2 n}}$ and using the above inequality and the union bound, we have that, with
 588 probability at least $1 - \delta$, for all $l \in \{0, 1, \dots, L-1\}$:

$$|\hat{p}_\sigma(l/L) - \mathbb{E}[\hat{p}_\sigma(l/L)]| \leq \sqrt{\frac{\log(2L/\delta)}{4\pi\sigma^2 n}}. \quad (\text{C.14})$$

589 Let us now compute the expected value of $\hat{p}_\sigma(x)$.

$$\mathbb{E}[\hat{p}_\sigma(x)] = \mathbb{E} \left[\frac{1}{n\sigma} \sum_{i=1}^n \phi \left(\frac{X^{(i)} - x}{\sigma} \right) \right] \quad (\text{C.15})$$

$$= \frac{1}{n\sigma} \sum_{i=1}^n \mathbb{E} \left[\phi \left(\frac{X^{(i)} - x}{\sigma} \right) \right] \quad (\text{C.16})$$

$$= \frac{1}{\sigma} \int p(u) \phi \left(\frac{x-u}{\sigma} \right) du \equiv (p \otimes \mathcal{N}(0, \sigma^2))(x) = p_\sigma(x). \quad (\text{C.17})$$

590 Combining equation C.14 and equation C.17, we get:

$$|\hat{p}_\sigma(l/L) - p_\sigma(x)| \leq \sqrt{\frac{\log(2L/\delta)}{4\pi\sigma^2n}}. \quad (\text{C.18})$$

Putting everything together we have:

$$d_{\text{TV}}(p_\sigma, \hat{p}_\sigma) \leq \frac{\lambda}{2L} + \frac{1}{2L\sigma^2\sqrt{2\pi e}} + \sqrt{\frac{\log(2L/\delta)}{4\pi\sigma^2n}}.$$

Choosing $L = n \cdot \max\{\lambda, 1\}$ we get that:

$$d_{\text{TV}}(p_\sigma, \hat{p}_\sigma) \lesssim \frac{1}{n} + \frac{1}{\sigma^2n} + \sqrt{\frac{\log n + \log(1 \vee \lambda) + \log 2/\delta}{\sigma^2n}}.$$

591 C.2 Evolution of parameters under noise

592 *Proof of theorem B.2:* We will use the following facts:

593 *Fact 1* (Direct corollary of the optimal coupling theorem). There exists a coupling γ of P and Q ,
 594 which samples a pair of random variables $(X, Y) \sim \gamma$ such that $\Pr_\gamma[X \neq Y] = d_{\text{TV}}(P, Q)$.

595 *Fact 2.* For any $x, y \in \mathbb{R}^d$: $d_{\text{TV}}(\mathcal{N}(x, \sigma^2\text{I}), \mathcal{N}(y, \sigma^2\text{I})) \leq \|x - y\|/2\sigma$

Proof. The KL divergence between $\mathcal{N}(\mu_1, \Sigma_1)$ and $\mathcal{N}(\mu_2, \Sigma_2)$ is

$$\text{KL}(\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)) = \frac{1}{2} \left(\text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)\Sigma_2^{-1}(\mu_2 - \mu_1) - d + \log \frac{|\Sigma_2|}{|\Sigma_1|} \right).$$

Applying this general result to our case:

$$\text{KL}(\mathcal{N}(x, \sigma^2\text{I}), \mathcal{N}(y, \sigma^2\text{I})) = \frac{1}{2} \left(\frac{\|x - y\|^2}{\sigma^2} \right).$$

596 We conclude by applying Pinsker's inequality. □

597 A corollary of Fact 2 and the optimal coupling theorem is the following:

598 *Fact 3.* Fix arbitrary $x, y \in \mathbb{R}^d$. There exists a coupling $\gamma_{x,y}$ of $\mathcal{N}(0, \sigma^2\text{I})$ and $\mathcal{N}(0, \sigma^2\text{I})$, which
 599 samples a pair of random variables $(Z, Z') \sim \gamma_{x,y}$ such that $\Pr_{\gamma_{x,y}}[x + Z \neq y + Z'] = \|x - y\|/2\sigma$.

600 Now let us denote by $\tilde{P} = P \otimes \mathcal{N}(0, \sigma^2\text{I})$ and $\tilde{Q} = Q \otimes \mathcal{N}(0, \sigma^2\text{I})$. To establish our claim in the
 601 theorem statement, it suffices to exhibit a coupling $\tilde{\gamma}$ of \tilde{P} and \tilde{Q} which samples a pair of random
 602 variables $(\tilde{X}, \tilde{Y}) \sim \tilde{\gamma}$ such that: $\Pr_{\tilde{\gamma}}[\tilde{X} \neq \tilde{Y}] \leq d_{\text{TV}}(P, Q) \cdot \frac{D}{2\sigma}$. We define coupling $\tilde{\gamma}$ as follows:

- 603 1. Sample $(X, Y) \sim \gamma$ (as specified in Fact 1); then
- 604 2. sample $(Z, Z') \sim \gamma_{X,Y}$ (as specified in Fact 3); then
- 605 3. output $(\tilde{X}, \tilde{Y}) := (X + Z, Y + Z')$.

606 Let us argue the following:

607 **Lemma C.3.** *The afore-described sampling procedure $\tilde{\gamma}$ is a valid coupling of \tilde{P} and \tilde{Q} .*

608 *Proof.* We need to establish that the marginals of $\tilde{\gamma}$ are \tilde{P} and \tilde{Q} . We will only show that for
 609 $(\tilde{X}, \tilde{Y}) \sim \tilde{\gamma}$ according to the afore-described sampling procedure, the marginal distribution of \tilde{X} is
 610 \tilde{P} , as the proof for \tilde{Y} is identical. Since γ is a coupling of P and Q , for $(X, Y) \sim \gamma$, the marginal
 611 distribution of X is P . By Fact 3, conditioning on any value of X and Y , the marginal distribution of
 612 Z is $\mathcal{N}(0, \sigma^2\text{I})$. Thus, $\tilde{X} = X + Z$, where $X \sim P$ and independently $Z \sim \mathcal{N}(0, \sigma^2\text{I})$, and thus the
 613 distribution of \tilde{X} is \tilde{P} . □

614 **Lemma C.4.** *Under the afore-described coupling $\tilde{\gamma}$: $\Pr_{\tilde{\gamma}}[\tilde{X} \neq \tilde{Y}] \leq d_{\text{TV}}(P, Q) \cdot \frac{D}{2\sigma}$.*

615 *Proof.* Notice that, when $X = Y$, by Fact 3, $Z = Z'$ with probability 1, and therefore $\tilde{X} = \tilde{Y}$. So
 616 for event $\tilde{X} \neq \tilde{Y}$ to happen, it must be that $X \neq Y$ happens and, conditioning on this event, that
 617 $X + Z \neq Y + Z'$ happens. By Fact 1, $\Pr_\gamma[X \neq Y] = d_{TV}(P, Q)$. By Fact 3, for any realization of
 618 (X, Y) , $\Pr_{\gamma_{X,Y}}[X + Z \neq Y + Z'] = \frac{\|X - Y\|}{2\sigma} \leq \frac{D}{2\sigma}$, where we used that P and Q are supported on
 619 a set with diameter D . Putting the above together, the claim follows. \square

620 \square

621 C.3 Auxiliary Lemmas

622 **Lemma C.5** (Lipschitzness of the empirical density). *For a collection of points $X^{(1)}, \dots, X^{(n)}$
 623 consider the function $\hat{p}_\sigma(x) = \frac{1}{n\sigma} \sum_{i=1}^n \phi\left(\frac{X^{(i)} - x}{\sigma}\right)$, where $\phi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}$ is the Gaussian
 624 kernel. Then p_σ is $\left(\frac{1}{\sigma^2 \sqrt{2\pi e}}\right)$ -Lipschitz.*

625 *Proof.* Let us compute the derivative of \hat{p}_σ :

$$\hat{p}'_\sigma(x) = \frac{1}{n\sigma} \sum_{i=1}^n \frac{d}{dx} \phi\left(\frac{x - X^{(i)}}{\sigma}\right) \quad (\text{C.19})$$

$$= \frac{1}{\sqrt{2\pi}n\sigma} \sum_{i=1}^n \exp\left(-\frac{(X^{(i)} - x)^2}{2\sigma^2}\right) \frac{X^{(i)} - x}{\sigma^2} \quad (\text{C.20})$$

$$\leq \frac{1}{\sqrt{2\pi}\sigma^2} \max_u \exp(-u^2/2) u \quad (\text{C.21})$$

$$\leq \frac{1}{\sigma^2 \sqrt{2\pi e}}. \quad (\text{C.22})$$

626 \square

627 **Lemma C.6** (Lipschitzness of a density convolved with a Gaussian). *Let p be a density that is
 628 λ -Lipschitz. Let $p_\sigma = p \otimes \mathcal{N}(0, \sigma^2 I)$. Then, p_σ is also λ -Lipschitz.*

629 *Proof.* Let us denote with $\phi_\sigma(\cdot)$ the Gaussian density with variance σ^2 . We have that:

$$p_\sigma(x) - p_\sigma(y) = \int (p(x - \tau) - p(y - \tau)) \phi_\sigma(\tau) d\tau \Rightarrow \quad (\text{C.23})$$

$$|p_\sigma(x) - p_\sigma(y)| \leq \int |p(x - \tau) - p(y - \tau)| \phi_\sigma(\tau) d\tau \quad (\text{C.24})$$

$$\leq \lambda |x - y| \cdot \int \phi_\sigma(\tau) d\tau \quad (\text{C.25})$$

$$= \lambda |x - y|. \quad (\text{C.26})$$

630 \square

631 D Additional Results

632 D.1 CIFAR-10 controlled corruptions

633 Figures 5a, 5b and 6 show gaussian blur, motion blur, and JPEG corrupted CIFAR-10 images
 634 respectively at different levels of severity. Appendix Table 3 shows results for JPEG compressed
 635 data at different levels of compression. We also tested our method for motion blurred data with high
 636 severity, visualized in the last row of Appendix Figure 6, obtaining a best FID of 5.85 (compared to
 637 8.79 of training on only the clean data).

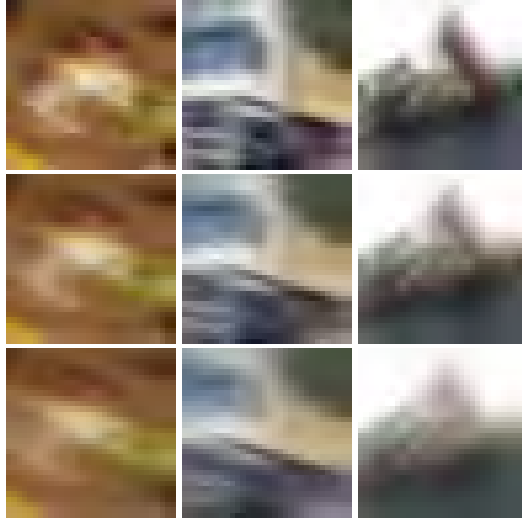
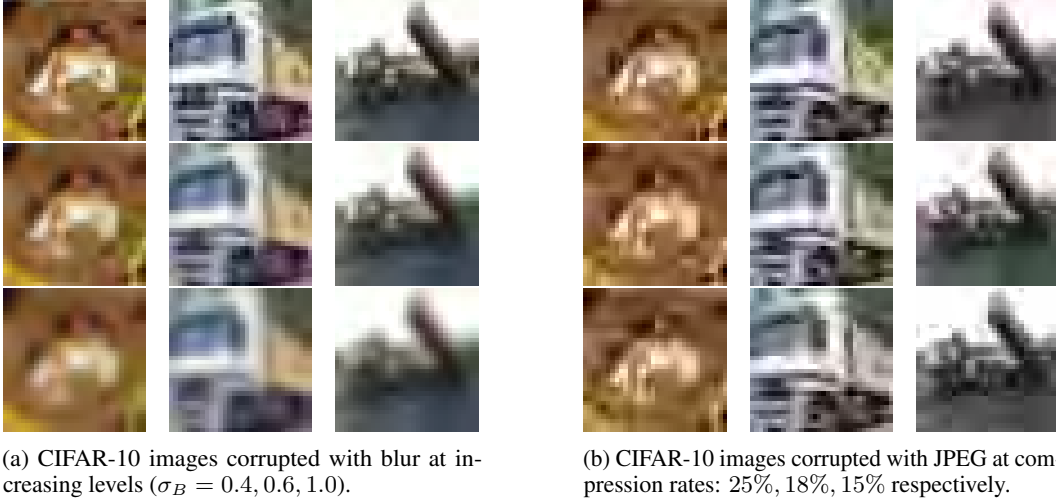


Figure 6: CIFAR-10 images corrupted with motion blur at increasing levels of corruption.

638 D.2 FFHQ-64x64 controlled corruptions

639 In Appendix 4 we show additional results for learning from blurred data on the FFHQ dataset.
 640 Similarly to the main paper, we observe that our Ambient-o algorithm leads to improvements over
 641 just using the high-quality data that are inversely proportional to the corruption level.

642 D.3 ImageNet results

643 In the main paper, we used FID as a way to measure the quality of generated images. However, FID
 644 is computed with respect to the test dataset that might also have samples of poor quality. Further,
 645 during FID computation, quality and diversity are entangled. To disentangle the two, we generate
 646 images using the EDM-2 baseline and our Ambient-o model and we use CLIP to evaluate the quality
 647 of the generated image (through the CLIP-IQA metric implemented in the PIQ package [38, 37]). We
 648 present results and win-rates in Table 5. As shown, Ambient-o achieves a better per-image quality
 649 compared to the baseline despite using exactly the same model, hyperparameters, and optimization
 650 algorithm. The difference comes solely from better use of the available data.

Table 3: Results for learning from JPEG compressed data on CIFAR-10.

Method	Dataset	Clean (%)	Corrupted (%)	JPEG Compression (Q)	$\bar{\sigma}_{t_n}^{\min}$	FID
Only Clean	Cifar-10	10	0	–	–	8.79
				15%	1.60	6.67
				18%	1.40	6.43
Ambient Omni	Cifar-10	10	90	25%	1.27	6.34
				50%	1.03	5.94
				75%	0.81	5.57
				90%	0.63	4.72

Table 4: Results for learning from blurred data, FFHQ.

Method	Dataset	Clean (%)	Corrupted (%)	Parameters Values (σ_B)	$\bar{\sigma}_{t_n}^{\min}$	FID
Only Clean	FFHQ	10	0	-	-	5.12
Ambient Omni	FFHQ	10	90	0.8	2.89	4.95
		10	90	0.6	2.12	4.65
		10	90	0.4	0.63	3.32

E Ambient diffusion implementation details and loss ablations

Similar to the EDM-2 [35] paper, we use a pre-condition weight to balance the importance of different diffusion times. Specifically, we modulate the EDM2 weight $\lambda(\sigma)$ by a factor:

$$\lambda_{\text{amb}}(\sigma, \sigma_{\min}) = \sigma^4 / (\sigma^2 - \sigma_{\min}^2)^2 \quad (\text{E.1})$$

for our ambient loss based on a similar analysis to [35]. We further use a buffer zone around the annotation time of each sample to ensure that the loss doesn't have singularities due to divisions by 0. We ablate the precondition term and the buffer size in Appendix Table 6.

For our ablations, we focus on the setting of training with 10% clean data and 90% corrupted data with Gaussian blur of $\sigma_B = 0.6$. Using no ambient pre-conditioning and no buffer, we obtain an FID of 5.56. In the same setting, adding the ambient pre-conditioning weight $\lambda_{\text{amb}}(\sigma, \sigma_{\min})$ improves FID by 0.13 points. Next, we ablate two strategies to mitigate the impact of the singularity of $\lambda_{\text{amb}}(\sigma, \sigma_{\min})$ at $\sigma = \sigma_{\min}$. The first strategy clips the ambient pre-conditioning weight at a specified maximum value $\lambda_{\text{amb}}^{\text{MAX}}$, but still trains for σ arbitrarily close to σ_{\min} . The second strategy also specifies a maximum value, but imposes a buffer

$$\sigma > \sqrt{1 + \frac{1}{\lambda_{\text{amb}}^{\text{MAX}} - 1}} \sigma_{\min} \quad (\text{E.2})$$

that restricts training to noise levels σ such that $\lambda_{\text{amb}}(\sigma, \sigma_{\min}) \leq \lambda_{\text{amb}}^{\text{MAX}}$. Clipping the ambient weight to $\lambda_{\text{amb}}^{\text{MAX}} = 2.0$ minimally improves FID to 5.35, but clipping to 4.0 significantly worsens it to 5.69. Adding a buffer at $\lambda_{\text{amb}}^{\text{MAX}} = 2.0$ slightly worsens FID to 5.40, but slackening the buffer to 4.0 minimally improves FID to 5.34. We opt for the buffering strategy in favor of the clipping strategy since performance appears convex in the buffer parameter, and because it obtains the best FID.

F Annotation ablations

We ablate the choice of using a fixed annotation vs sample-adaptive annotations in Appendix Table 7. We find that using sample-adaptive annotations achieves improved results. Nevertheless, both annotation methods yield improvements over the training on filtered data and the training on everything baselines. To show that our method works for more corruption types, we perform an equivalent experiment with JPEG compressed data at different compression ratios and we achieve similar results, presented in Appendix Table 3. We ablate the impact of the amount of training data and the number of training iterations on the classifier annotations in Appendix Section F. We show results for motion blur (Figure 6 and Section D.1) and for the FFHQ dataset (Table 4).

Table 5: Additional comparison between EDM-2 XXL and our Ambient-o model using the CLIP IQA metric for image quality assesment. Ambient-o leads to improved scores despite using the exact same architecture, data and hyperparameters. For this experiment, we use the models with guidance optimized for DINO FD since they are the ones producing the higher quality images.

Metric	EDM-2 [35] XXL	Ambient-o XXL crops
Average CLIP IQA score	0.69	0.71
Median CLIP IQA score	0.79	0.80
Win-rate	47.98%	52.02%

Table 6: Ablation study of ambient weight and stability buffer on Cifar-10 with 10% clean data and 90% corrupted data with blur of 0.6.

Method	FID ↓
<i>No ambient preconditioning weight and no buffer:</i>	
$\lambda_{\text{amb}}(\sigma, \sigma_{\min}) = 1 \ \& \ \sigma > \sigma_{\min}$	5.49
<i>Adding ambient preconditioning weight:</i>	
+ Weight $\lambda_{\text{amb}}(\sigma, \sigma_{\min}) = \sigma^4 / (\sigma^2 - \sigma_{\min}^2)^2$	5.36
<i>Adding stability buffer/clipping:</i>	
+ Clip $\lambda_{\text{amb}}(\sigma, \sigma_{\min})$ at 2.0	5.35
+ Clip $\lambda_{\text{amb}}(\sigma, \sigma_{\min})$ at 4.0	5.69
+ Buffer $\lambda_{\text{amb}}(\sigma, \sigma_{\min})$ at 2.0 i.e. $\sigma > \sqrt{2}\sigma_{\min}$	5.40
+ Buffer $\lambda_{\text{amb}}(\sigma, \sigma_{\min})$ at 4.0 i.e. $\sigma > (2/\sqrt{3})\sigma_{\min}$	5.34

678 **Balanced vs unbalanced data:** We ablate the impact of classifier training data on the setting of
679 CIFAR-10 with 10% clean data and 90% corrupted data with gaussian blur with $\sigma_B = 0.6$. When
680 annotating with a classifier trained on the same unbalanced dataset we train the diffusion model on
681 we obtained a best FID of 6.04, compared to the 5.34 obtained if we train on a balanced dataset.

682 **Training iterations:** We ablate the impact of classifier training iterations on the setting of CIFAR-10
683 with 10% clean data and 90% corrupted data with JPEG compression at compression rate of 18%,
684 training the classifier with a balanced dataset. We report minute variations in the best FID, obtaining
685 6.50, 6.58, and 6.49 when training the classifier for 5e6, 10e6, and 15e6 images worth of training
686 respectively.

Table 7: Comparison with baselines for training with data corrupted by Gaussian Blur at different levels. The dataset used in this experiment is CIFAR-10.

Method	Clean (%)	Corrupted (%)	Parameters Values (σ_B)	$\bar{\sigma}_{t_n}^{\min}$	FID
Only Clean	10	0	-	-	8.79
No annotations	10	90	1.0	0	45.32
			0.8		28.26
			0.4		2.47
Single annotation	10	90	1.0	2.32	6.95
			0.8	1.89	6.66
			0.4	0.00	2.47
Classifier annotations	10	90	1.0	2.84	6.16
	10	90	0.8	1.93	6.00
	10	90	0.4	0.22	2.44

G Training Details

G.1 Formation of the high-quality and low-quality sets.

In the theoretical problem setting we assumed the existence of a good set S_G from the clean distribution and a bad set S_B from the corrupted distribution. In practice, we do not actually possess these sets initially, but we can construct them so long as we have access to a measure of "quality". Given a function on images which tells us whether its good enough to generate or not e.g. CLIP-IQA quality [66] greater than some threshold, we can define our good set S_G as the good enough images and S_B as the complement. From this point on we can apply the methodology of ambient-o as developed, either employing classifier annotations as in our pixel diffusion experiments, or fixed annotations as in our large scale ImageNet and text-to-image experiments.

G.2 Datasets

CIFAR-10. CIFAR-10 [42] consists of 60,000 32x32 images of ten classes (airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck).

FFHQ. FFHQ [36] consists of 70,000 512x512 images of faces from Flickr. We used the dataset at 64x64 resolution for our experiments.

AFHQ. AFHQ [12] consists of 5,653 images of cats, 5,239 images of dogs and 5,000 images of wildlife, for a total of 15,892 images.

ImageNet. ImageNet [20] consists of 1,281,167 images of variable resolution from 1000 classes.

Conceptual Captions. Conceptual Captions [56] consists of 12M (image url, caption) pairs.

Segment Anything. Segment Anything [41] consists of 11.1M high-resolution images annotated with segmentation masks. Since the original dataset did not have real captions, we use the same LLaVA generated captions created by the MicroDiffusion [54] paper.

JourneyDB. JourneyDB consists of 4.4M synthetic image-caption pairs from Midjourney [63].

DiffusionDB. DiffusionDB consists of 14M synthetic image-caption pairs, mostly generated from Stable Diffusion models [70]. We use the same 10.7M quality-filtered subset created by the MicroDiffusion paper [54].

G.3 Diffusion model training

CIFAR-10. We use the EDM [34] codebase as a reference to train class-conditional diffusion models on CIFAR-10. The architecture is a Diffusion U-Net [60] with ~55M parameters. We use the Adam optimizer [40] with learning rate 0.001, batch size 512, and no weight decay. While the original EDM paper trained for 200×10^6 images worth of training, when training with corrupted data we saw best results around 20×10^6 images. On a single 8xV100 node we achieved a throughput of 0.8s per 1k images, for an average of 4.4h per training run.

FFHQ. Same as for CIFAR-10, except learning was set to $2e - 4$, we trained for a maximum of 100×10^6 images worth of training, and saw best results around 30×10^6 images worth.

AFHQ. Same as FFHQ.

ImageNet. We use the EDM2 [35] codebase as a reference to train class-conditional diffusion models on ImageNet. The architecture is a Diffusion U-Net [60] with ~125M parameters. We use the Adam optimizer [40] with reference learning rate 0.012, batch size 2048, and no weight decay. Same as the original codebase, we trained for ~2B worth of images. On 32 H200 GPUs, XS models took ~3 days to train, while XXL models took ~7 days.

728 **MicroDiffusion.** We use the MicroDiffusion codebase [54] as a reference to train text-to-image
729 models on an academic budget. We follow their recipe exactly, changing only the standard denoising
730 diffusion loss to the ambient diffusion loss. The architecture is a Diffusion Transformer [50] utilizing
731 Mixture-of-Experiments (MoE) feedforward layers [57, 31], with $\sim 1.1\text{B}$ parameters. We use the
732 AdamW optimizer [40] with reference learning rates $2.4e - 4/8e - 5/8e - 5/8e - 5$ for each of the
733 four phases and batch size 2048 for all phases. On 8 H200 GPUs, training takes ~ 2 days to train.

734 G.4 Classifier training

735 Classifier training is done using the same optimization recipe (optimizer, learning rate, batch size,
736 etc.) as diffusion model training, except we change the architecture to an encoder-only "Half-Unet",
737 simply by removing the decoder half of the original UNet architecture. The training of the classifier
738 is substantially shorter compared to the diffusion training since classification is task is easier than
739 generation.

740 H Additional Figures

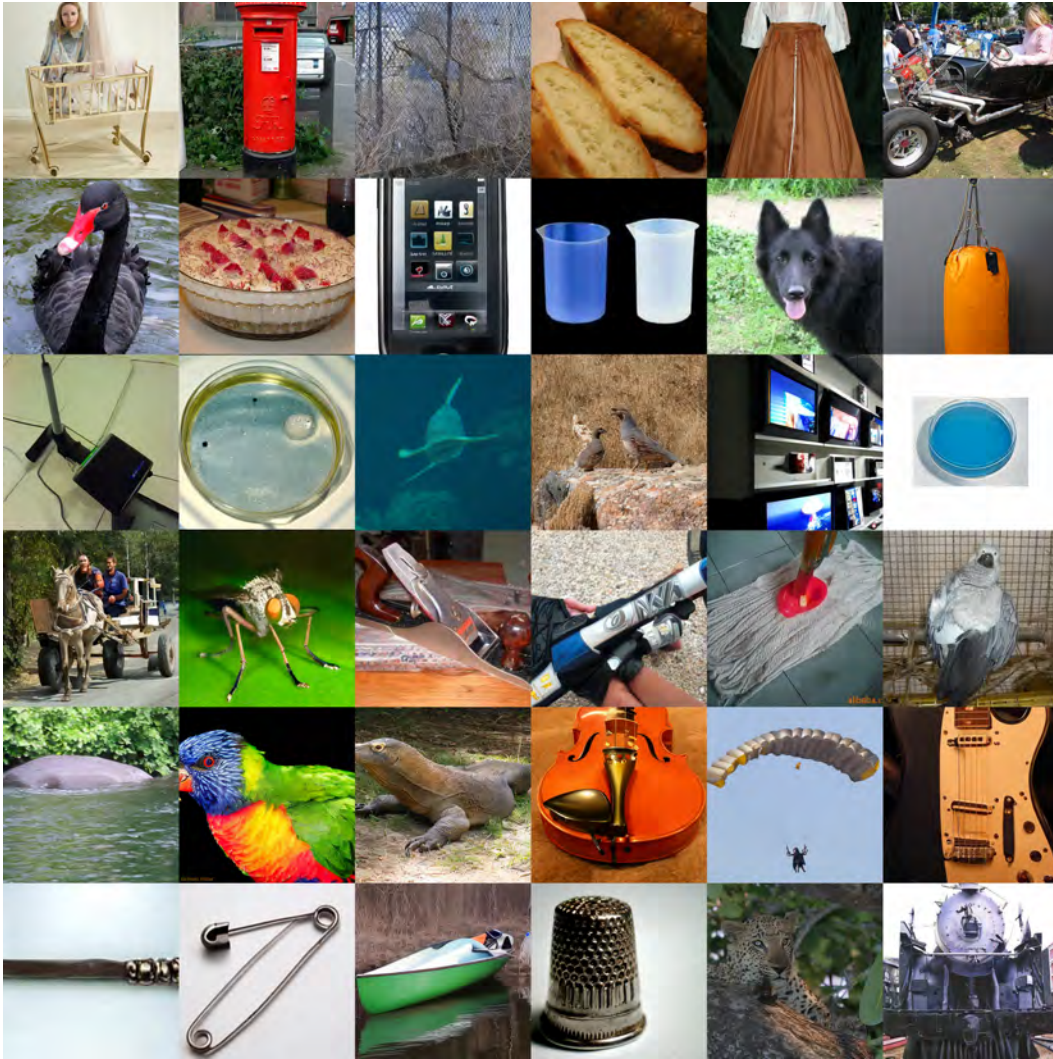


Figure 7: Uncurated generations from our Ambient-o XXL model trained on ImageNet.

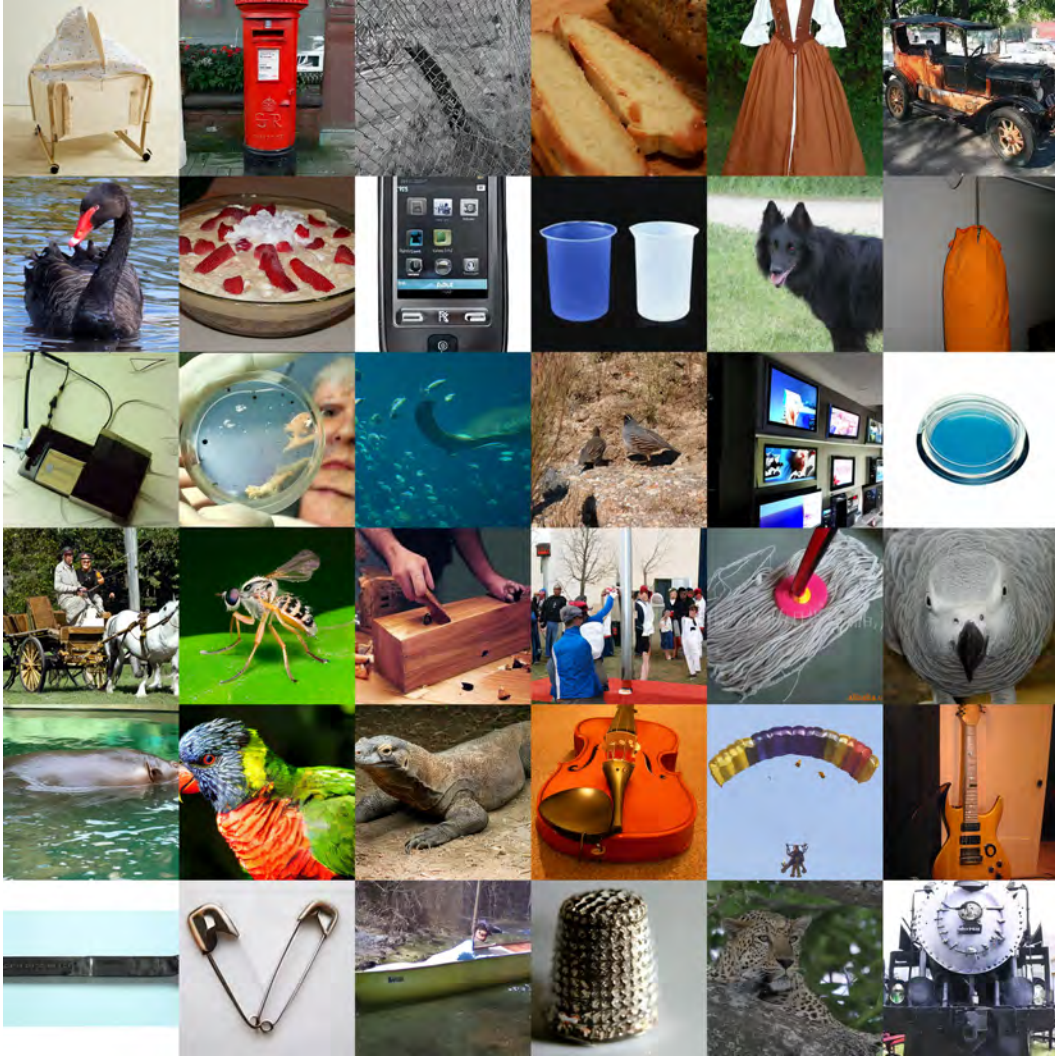


Figure 8: Uncurated generations from our Ambient-o+crops XXL model trained on ImageNet.

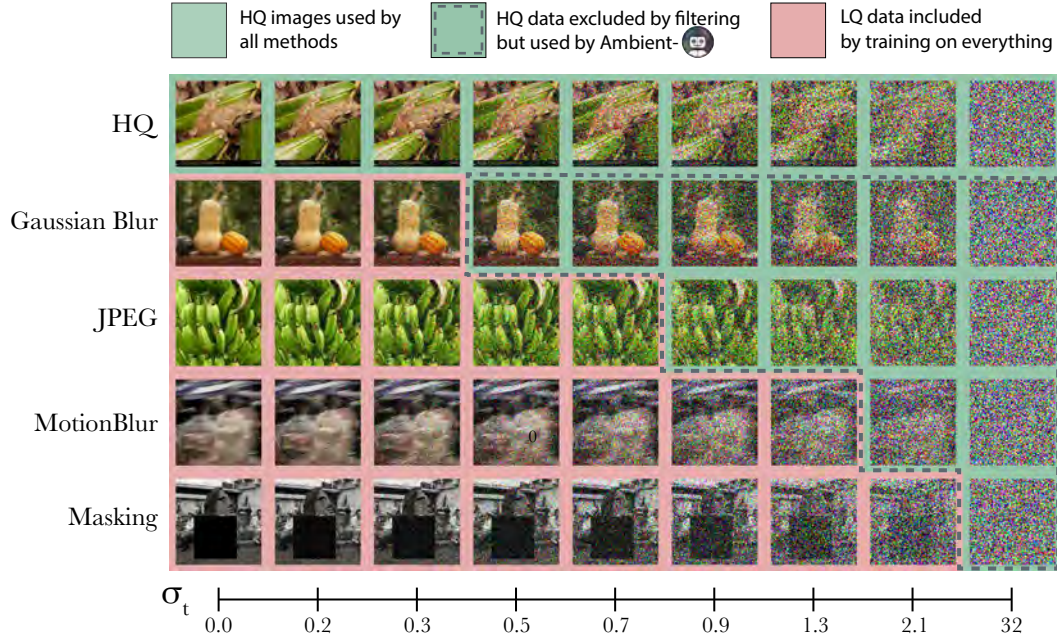


Figure 9: **Visual summary of our method for using low-quality data at high-noise.** We see how the various corrupted images become indistinguishable from the High Quality (HQ) after a minimum noise level. These noisy versions of Low Quality (LQ) images are actually high-quality data, which filtering approaches discard, but Ambient Omni uses.

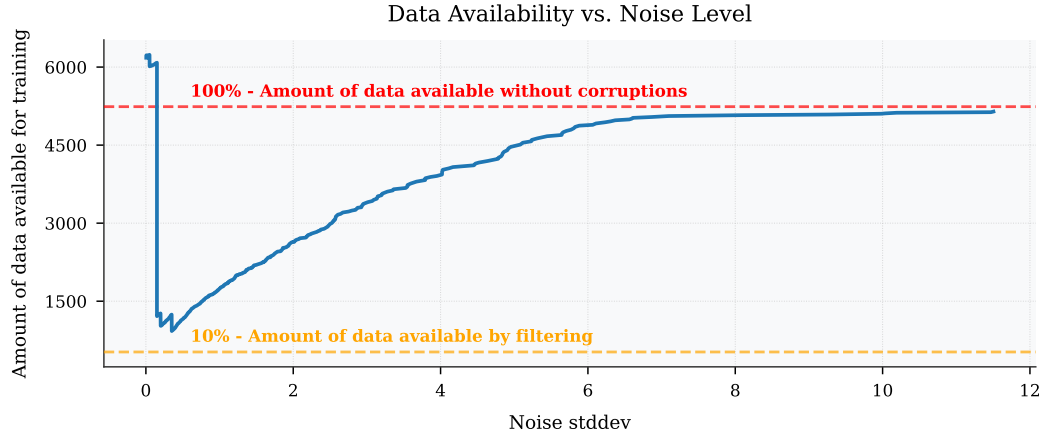


Figure 10: Amount of samples available at each noise level when training a generative model for dogs in the following setting: (1) we have 10% of the dogs dataset uncorrupted, (2) we have the other 90% of the dogs dataset corrupted with gaussian blur with $\sigma_B = 0.6$, and (3) we have 100% of the clean dataset of cats. At low noise levels, we can train on both the high quality dogs and a lot of the cats, resulting in $> 100\%$ of samples available relative to the original dogs dataset size. As the noise level starts to increase, we stop being able to use the out-of-distribution cat samples, but start gaining some blurry dog samples. As the noise level approaches the maximum all the blurry dogs become available for training, such that the amount of data available approaches 100%.

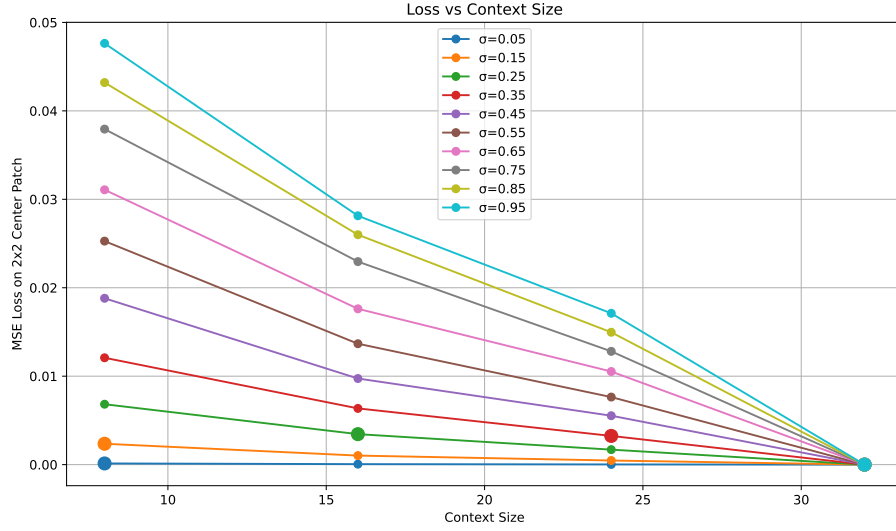


Figure 11: ImageNet-512x512: denoising loss of an optimally trained model, measured at 2×2 center patch, as we increase the context size given to the model (horizontal axis) and the noise level (different curves). As expected, for higher noise, more context is needed for optimal denoising. The large dot on each curve marks the point where the loss nearly plateaus.

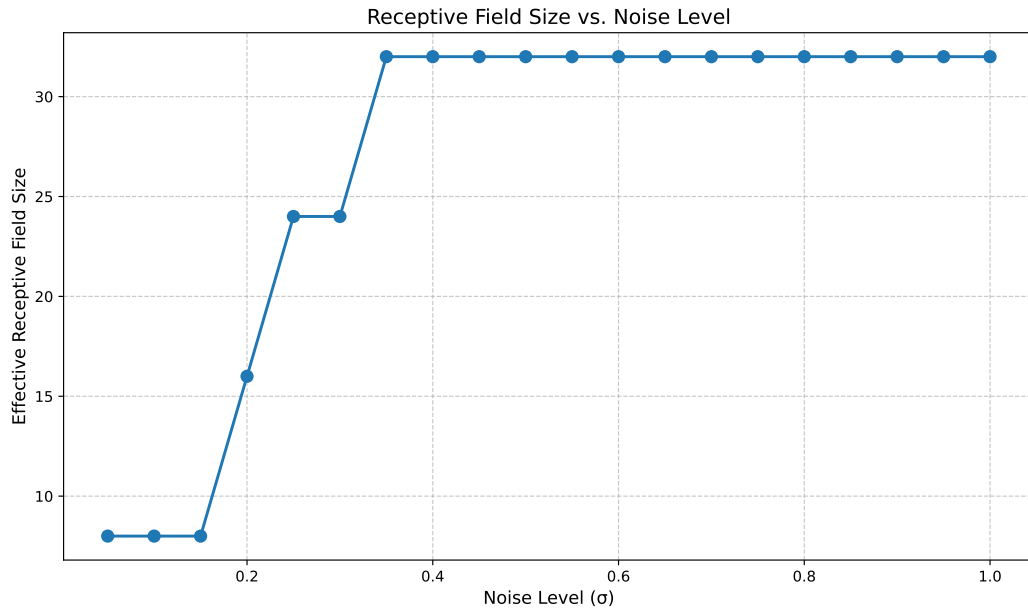


Figure 12: ImageNet-512x512: context size needed to be within $\epsilon = 1e - 3$ of the optimal loss for different noise levels. As expected, for higher noise, more context is needed for optimal denoising.

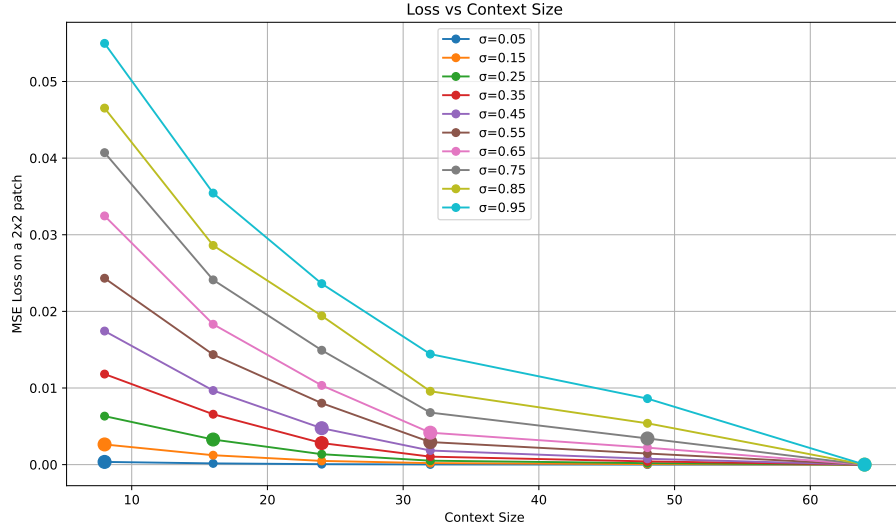


Figure 13: FFHQ: denoising loss of an optimally trained model, measured at 2×2 center patch, as we increase the context size given to the model (horizontal axis) and the noise level (different curves). As expected, for higher noise, more context is needed for optimal denoising. The large dot on each curve marks the point where the loss nearly plateaus.

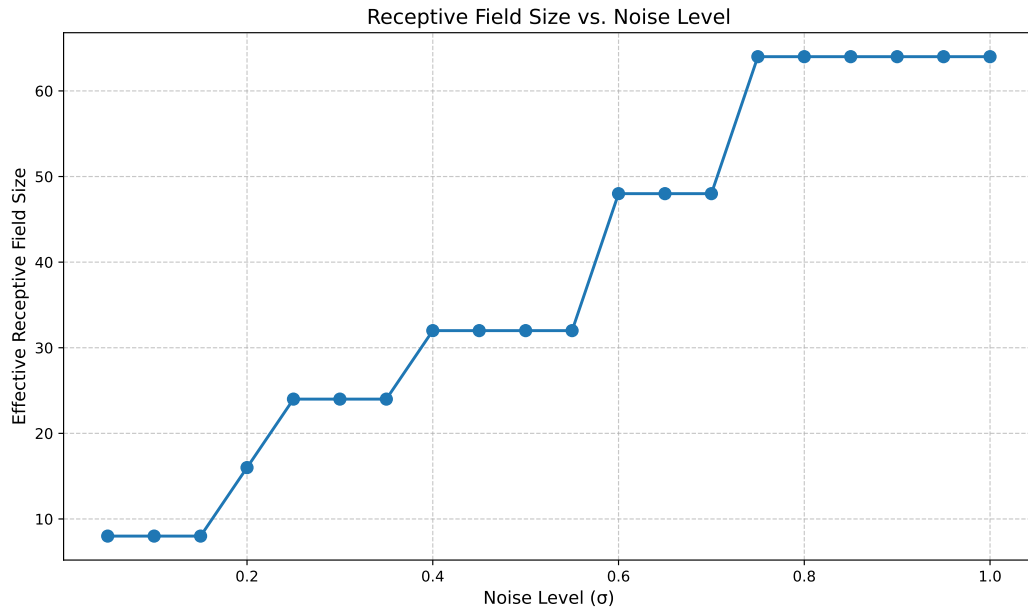
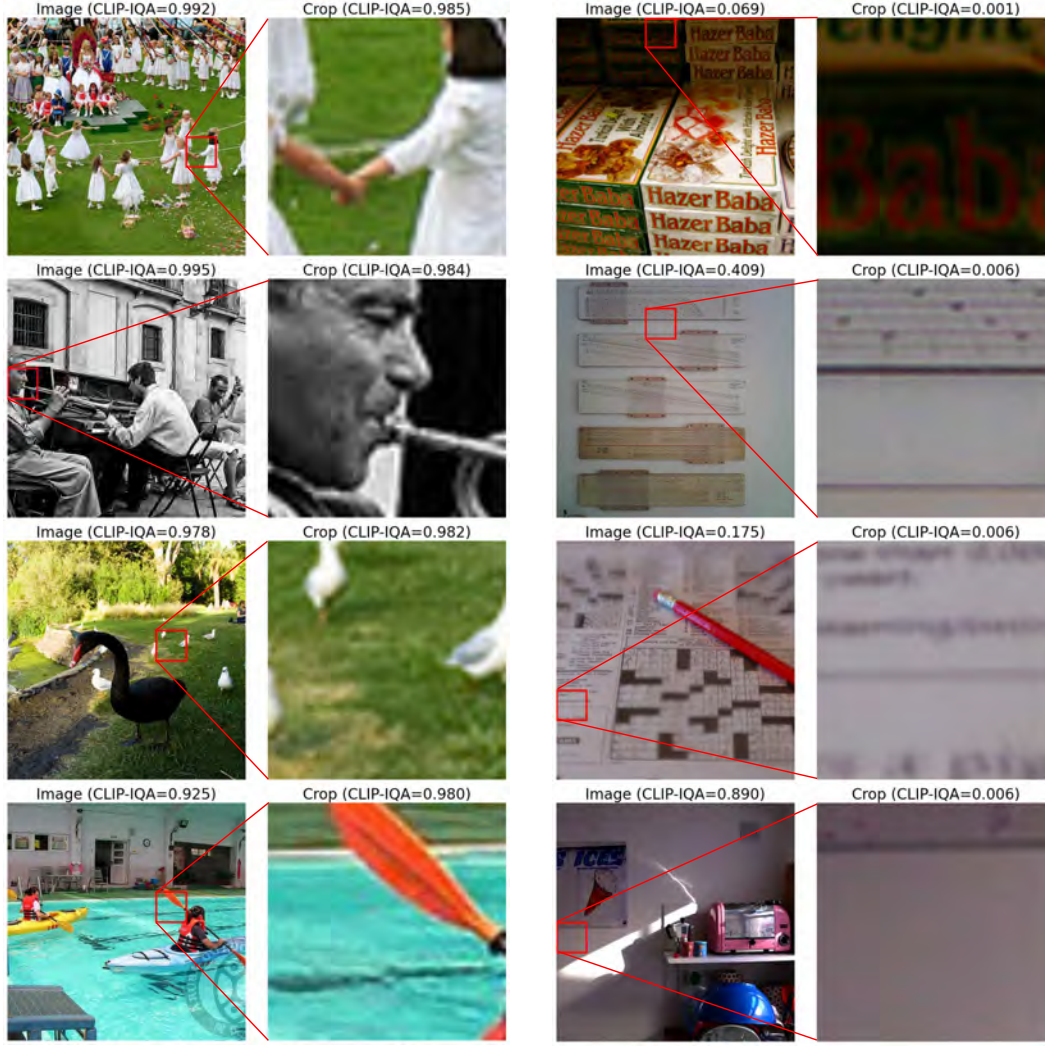


Figure 14: FFHQ: context size needed to be within $\epsilon = 1e - 3$ of the optimal loss for different noise levels. As expected, for higher noise, more context is needed for optimal denoising.



(a) High quality crops

(b) Low quality crops

Figure 15: Results using CLIP to find (a) high-quality and (b) low-quality crops on ImageNet.

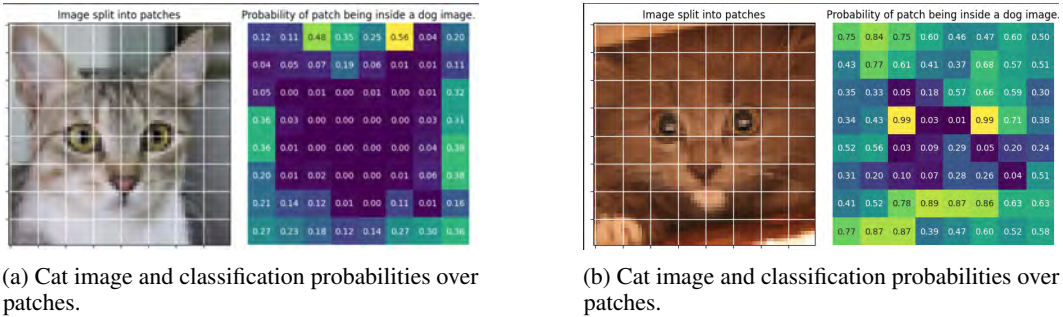


Figure 16: Two examples of cats from the AFHQ dataset. We partition each cat into non overlapping patches and we compute the probabilities of the patch belonging to an image of a dog using a cats vs dogs classifier trained on patches. The cat on the right has a lot more patches that could belong to a dog image according to the classifier, possibly due to the color or the texture of the fur.

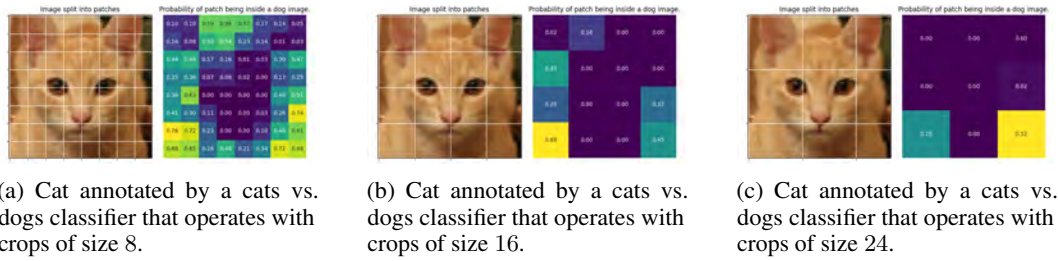


Figure 17: Patch-based annotations of a cat image from AFHQ using cats vs. dogs classifiers trained on different patch sizes.

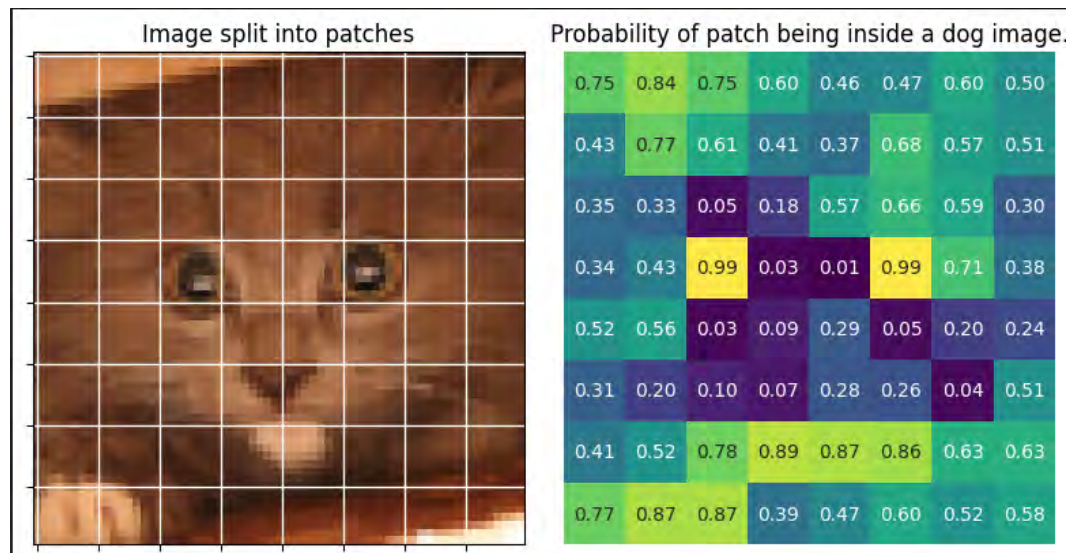


Figure 18: Patch level probabilities for dogness in a cat image.

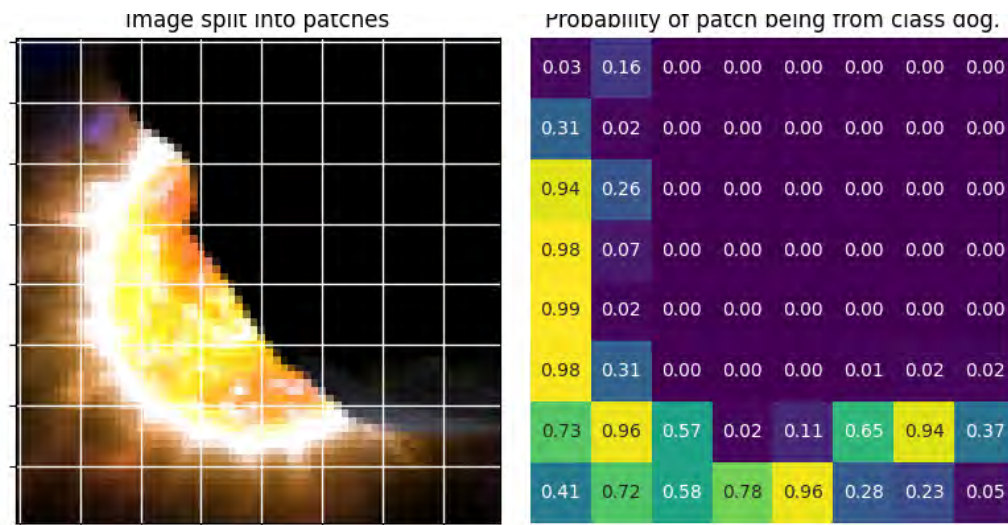
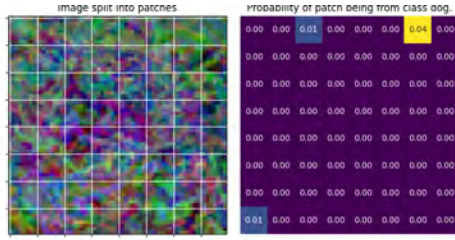
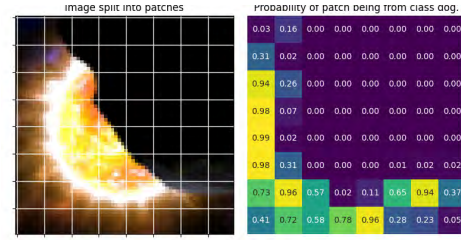


Figure 19: Patch level probabilities for dogness in a synthetic image (procedural program). The cat has more useful patches than this non-realistic procedural program.

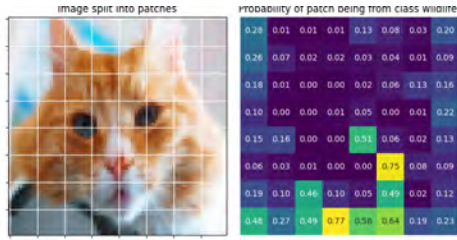


(a) Synthetic image and classification probabilities over patches.

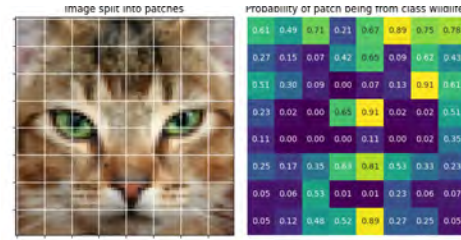


(b) Synthetic image and classification probabilities over patches.

Figure 20: Two examples of procedurally generated images. We partition each image into non overlapping patches and we compute the probabilities of the patch belonging to an image of a dog using a synthetic image vs dogs classifier trained on patches. The image on the right has a lot more patches that could belong to a dog image according to the classifier, possibly due to the color or the texture.



(a) Cat image and classification probabilities over patches.



(b) Cat image and classification probabilities over patches.

Figure 21: Two examples of cat images. We partition each image into nonoverlapping patches and we compute the probabilities of the patch belonging to an image of wildlife using a cats vs wildlife classifier trained on patches. The image on the right has a lot more patches that could belong to a wildlife image according to the classifier, possibly due to the color or the texture.

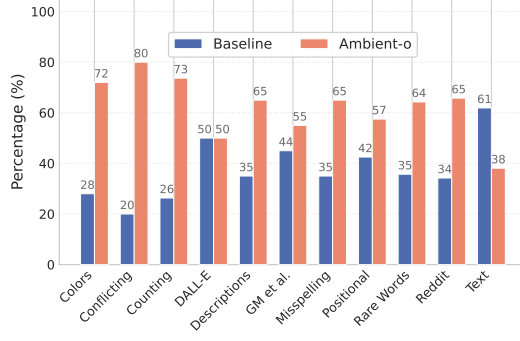


(a) Example batch.

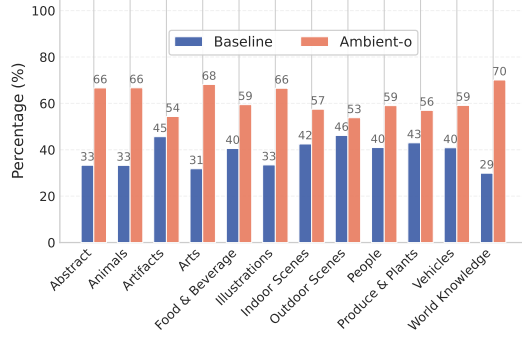


(b) Noisy batch.

Figure 22: Example batch.



(a) DrawBench

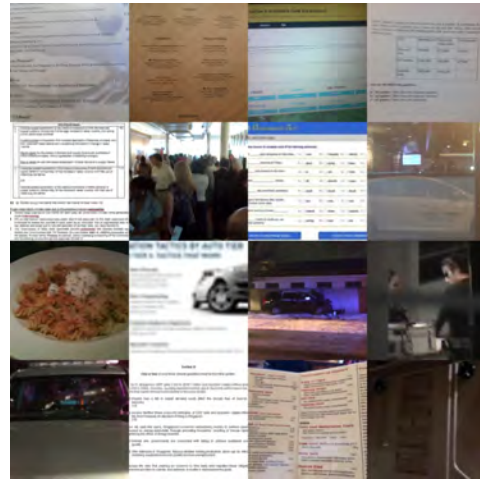


(b) PartiPrompts

Figure 23: Assessing image quality with GPT-4o on DrawBench and PartiPrompts.



(a) Highest quality images from CC12M according to CLIP.

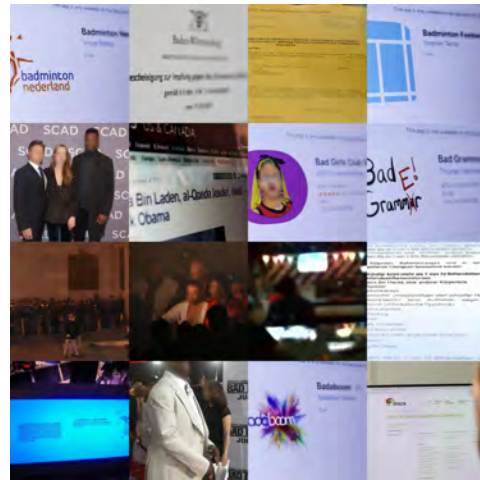


(b) Lowest quality images from CC12M according to CLIP.

Figure 24: CLIP annotations for quality of images from CC12M.



(a) Highest quality images from SA1B according to CLIP.

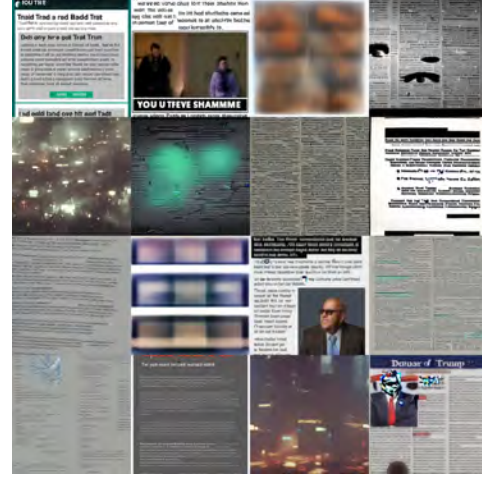


(b) Lowest quality images from SA1B according to CLIP.

Figure 25: CLIP annotations for quality of images from SA1B.



(a) Highest quality images from DiffDB according to CLIP.

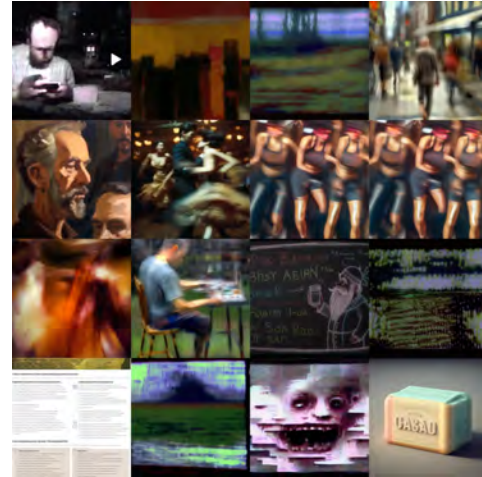


(b) Lowest quality images from DiffDB according to CLIP.

Figure 26: CLIP annotations for quality of images from DiffDB.



(a) Highest quality images from JDB according to CLIP.



(b) Lowest quality images from JDB according to CLIP.

Figure 27: CLIP annotations for quality of images from JDB.

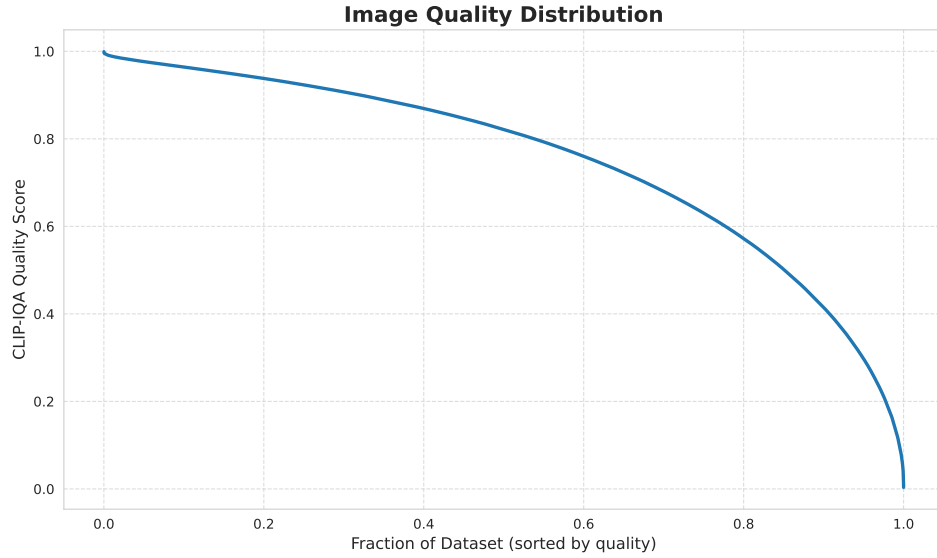


Figure 28: Distribution of image qualities according to CLIP for ImageNet-512.



Figure 29: **Examples of mode collapse.** Left: baseline model finetuned on a high-quality subset. Right: Ambient-o model using all the data. As shown, finetuning decreases output diversity.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: Our method does not use any information about the type of corruption, and our experiments show it generalizes to low quality data found in the wild, not just a few artificially controlled corruptions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: We openly discuss the limitations of our approach, such as:

- (a) The high and low quality distributions never *perfectly* merge, so our method always introduces a (small) distribution error compared to filtering.
- (b) Our method does not work well with certain corruption types, such as masking. These "ill-suited" corruptions require a very large amount of noise to merge, such that they are effectively never used during training and our method reduces to filtering in these cases.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best

judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: Our theorems include all premises and assumptions used to prove the result. Informal proofs are found in the main text, referencing formal proofs in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: All the information on the algorithm and the training recipe needed to reproduce our experiments is included in the paper (either in the main text or the appendix). Additionally, we will make the training and evaluation code public after acceptance of the paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: All data used is publically accessible. We will release the full training and evaluation code upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide the core elements in the main text and the full details in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: Obtaining error bars would require extremely computationally expensive retraining of diffusion models.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: GPU type and number and compute time is provided in the appendix for all experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The work does not use human trials, and all data used is publically available. We analyse the potential negative impacts of improving generative model abilities in ??.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: See ??.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We are not releasing any datasets. We will be releasing the models upon paper acceptance, but there has already been a model trained and open-sourced from the same dataset. Moreover, our work is far away from state-of-the-art text-to-image generation, and thus does not introduce extra risks that do not already exist.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

1002 Justification: Prior work has already trained and made public models trained on the same
1003 data we use to train. Moreover, all datasets are publically available and were introduced by
1004 prior research work, which we explicitly state and cite.

1005 Guidelines:

- 1006 • The answer NA means that the paper does not use existing assets.
- 1007 • The authors should cite the original paper that produced the code package or dataset.
- 1008 • The authors should state which version of the asset is used and, if possible, include a
1009 URL.
- 1010 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1011 • For scraped data from a particular source (e.g., website), the copyright and terms of
1012 service of that source should be provided.
- 1013 • If assets are released, the license, copyright information, and terms of use in the
1014 package should be provided. For popular datasets, paperswithcode.com/datasets
1015 has curated licenses for some datasets. Their licensing guide can help determine the
1016 license of a dataset.
- 1017 • For existing datasets that are re-packaged, both the original license and the license of
1018 the derived asset (if it has changed) should be provided.
- 1019 • If this information is not available online, the authors are encouraged to reach out to
1020 the asset's creators.

1021 13. New assets

1022 Question: Are new assets introduced in the paper well documented and is the documentation
1023 provided alongside the assets?

1024 Answer: [NA]

1025 Justification: We do not release any new datasets.

1026 Guidelines:

- 1027 • The answer NA means that the paper does not release new assets.
- 1028 • Researchers should communicate the details of the dataset/code/model as part of their
1029 submissions via structured templates. This includes details about training, license,
1030 limitations, etc.
- 1031 • The paper should discuss whether and how consent was obtained from people whose
1032 asset is used.
- 1033 • At submission time, remember to anonymize your assets (if applicable). You can either
1034 create an anonymized URL or include an anonymized zip file.

1035 14. Crowdsourcing and research with human subjects

1036 Question: For crowdsourcing experiments and research with human subjects, does the paper
1037 include the full text of instructions given to participants and screenshots, if applicable, as
1038 well as details about compensation (if any)?

1039 Answer: [NA]

1040 Justification: No research with human subjects.

1041 Guidelines:

- 1042 • The answer NA means that the paper does not involve crowdsourcing nor research with
1043 human subjects.
- 1044 • Including this information in the supplemental material is fine, but if the main contribu-
1045 tion of the paper involves human subjects, then as much detail as possible should be
1046 included in the main paper.
- 1047 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1048 or other labor should be paid at least the minimum wage in the country of the data
1049 collector.

1050 15. Institutional review board (IRB) approvals or equivalent for research with human 1051 subjects

1052 Question: Does the paper describe potential risks incurred by study participants, whether
 1053 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
 1054 approvals (or an equivalent approval/review based on the requirements of your country or
 1055 institution) were obtained?

1056 Answer: [NA]

1057 Justification: No research with human subjects.

1058 Guidelines:

- 1059 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1060 human subjects.
- 1061 • Depending on the country in which research is conducted, IRB approval (or equivalent)
- 1062 may be required for any human subjects research. If you obtained IRB approval, you
- 1063 should clearly state this in the paper.
- 1064 • We recognize that the procedures for this may vary significantly between institutions
- 1065 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
- 1066 guidelines for their institution.
- 1067 • For initial submissions, do not include any information that would break anonymity (if
- 1068 applicable), such as the institution conducting the review.

1069 **16. Declaration of LLM usage**

1070 Question: Does the paper describe the usage of LLMs if it is an important, original, or

1071 non-standard component of the core methods in this research? Note that if the LLM is used

1072 only for writing, editing, or formatting purposes and does not impact the core methodology,

1073 scientific rigor, or originality of the research, declaration is not required.

1074 Answer: [NA]

1075 Justification: No important, original, or non-standard usage of LLMs in the paper.

1076 Guidelines:

- 1077 • The answer NA means that the core method development in this research does not
- 1078 involve LLMs as any important, original, or non-standard components.
- 1079 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
- 1080 for what should or should not be described.