

---

# Noise-conditional Maximum Likelihood Estimation with Score-based Sampling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We introduce a simple yet effective modification to the standard maximum likeli-  
2 hood estimation (MLE) framework for autoregressive generative models. Rather  
3 than maximizing a single unconditional likelihood of the data under the model, we  
4 maximize a family of *noise-conditional* likelihoods consisting of the data perturbed  
5 by a continuum of noise levels. We find that models trained this way are more  
6 robust to noise, obtain higher test likelihoods, and generate higher quality images.  
7 They can also be sampled from via a novel score-based sampling scheme which  
8 combats the classical *covariate shift* problem that occurs during sample generation  
9 in autoregressive models. Applying this augmentation to autoregressive image  
10 models, we obtain 3.32 bits per dimension on the ImageNet 64x64 dataset, and  
11 substantially improve the quality of generated samples in terms of the Fréchet  
12 Inception distance (FID) — from 37.50 to 13.50 on the CIFAR-10 dataset.

## 13 1 Introduction

14 Likelihood maximization models, *i.e.*, models trained by maximizing log-likelihood, are a leading  
15 class of modern generative models. Of these, autoregressive models boast state-of-the-art performance  
16 in many domains, including images Salimans et al. [2017], Child et al. [2019], text Vaswani et al.  
17 [2017], and audio Oord et al. [2016]. These architectures also show great promise for modeling long  
18 range dependencies Tay et al. [2020], Gu et al. [2021].

19 However, while log-likelihood is broadly agreed upon as one of the most rigorous metrics for  
20 goodness-of-fit in statistical and generative modeling, models with high likelihoods do not necessarily  
21 produce samples of high visual quality. This phenomenon has been discussed at length by Theis et al.  
22 [2015], Huszár [2015], and corroborated in empirical studies Grover et al. [2018], Kim et al. [2022].

23 Autoregressive models have an additional affliction: they have notoriously unstable dynamics during  
24 sample generation Bengio et al. [2015] due to their sequential sampling algorithm, which can cause  
25 errors to compound across time steps. Such errors cannot usually be corrected *ex post facto* due to  
26 the autoregressive structure of the model, and can substantially affect downstream steps as we find  
27 that their likelihoods are sensitive to even the most minor of perturbations.

28 Score-based diffusion models Song et al. [2020], Ho et al. [2020] offer a different perspective on the  
29 matter. Even though sampling is also sequential, diffusion models are more robust to perturbations  
30 because, in essence, they are trained as denoising functions Ho et al. [2020]. Moreover, the update  
31 direction in each step is unconstrained (unlike token-wise autoregressive models, which can only  
32 update one token at a time, and only once), meaning the model can correct errors from previous  
33 steps. However, diffusion models are poor likelihood models, as they cannot be trained via maximum  
34 likelihood, and density evaluations are inexact and require solving ODEs involving hundreds to  
35 thousands of function evaluations. Thus we wonder: is there a conceptual middle ground?

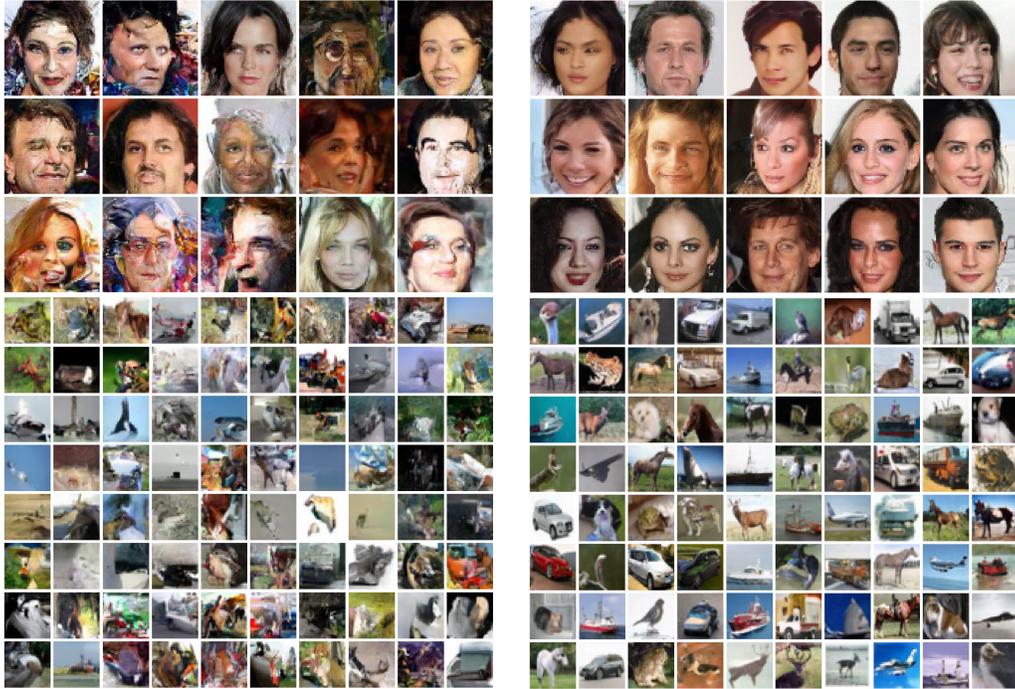


Figure 1: Generated samples on CelebA 64x64 (above) and CIFAR-10 (below). Autoregressive models trained via vanilla maximum likelihood (left) are brittle to sampling errors and can quickly diverge, producing nonsensical results. Those trained via our proposed algorithm (right) are more robust, which can significantly increase the coherence of the generated images.

36 In this paper, we offer such a framework. We further analyze the likelihood-sample quality mismatch  
 37 in autoregressive models, and propose techniques inspired by diffusion models to alleviate it. In  
 38 particular, we leverage the fact that the score function is naturally learned as a byproduct of maximum  
 39 likelihood estimation. This allows a novel two-part sampling strategy with noisy sampling and  
 40 score-based refinement.

41 Our contributions are threefold. 1) We investigate the pitfalls of training and inference under the  
 42 log-likelihood maximization scheme, particularly regarding sensitivity to noise corruptions. 2) We  
 43 propose a simple sanity test for checking the noise-robustness of likelihood models. 3) We introduce  
 44 a novel framework for the training and sampling of likelihood maximization models that improves  
 45 noise-robustness and substantially boosts the sample quality of the resulting model. Ultimately, we  
 46 obtain a model that can generate samples at a quality approaching that of diffusion models, without  
 47 losing the maximum likelihood framework and  $\mathcal{O}(1)$  likelihood evaluation speed of likelihood  
 48 maximization models.

## 49 2 The Pitfalls of Maximum Likelihood

50 We first show that density models trained to maximize the standard log-likelihood are surprisingly  
 51 sensitive to minor perturbations. We then discuss why this is bad for generative modeling performance.

### 52 2.1 A Simple Sanity Test

53 Consider the class of minimally corrupted probability densities we call  $p_\pi$ , where

$$p_\pi = p_{data} * p_{mult_{\{-1,0,1\}}(\pi/2, 1-\pi, \pi/2)}, \quad \pi \in [0, 1]. \quad (1)$$

54 Here,  $*$  denotes the convolution operator, and  $p_{mult_{\{a,b,c\}}(\alpha,\beta,\gamma)}$  is the density a  $d$ -dimensional  
 55 multinomial distribution taking on  $a$ ,  $b$ , and  $c$  with probabilities  $\alpha$ ,  $\beta$ , and  $\gamma$ , respectively.  $p_\pi$  is  
 56 *minimally corrupted* in the sense that, if  $p_{data}$  is an integer-discretized distribution (say, 8-bit images),

57  $p_\pi$  describes the distribution of points  $p_{data}$  that have had their least significant bit incremented or  
58 decremented with probability  $\pi$ .

59 To the human eye, the difference between samples drawn from  $p_\pi$  and  $p_{data}$  is almost imperceptible,  
60 even for  $\pi = 1$  (see Fig 2). However, for likelihood models, this perturbation drastically increases  
61 the negative log-likelihood of the data under the model (see Table 1), to the point that it significantly  
62 undermines (if not outright nullifies) any recent advances in density estimation. This basic inconsis-  
63 tency suggests that the learned density of many standard likelihood models is brittle and overly  
64 emphasizes bit-level statistics that have little influence on the inherent content of the image.

## 65 2.2 Why We Should Care

66 We provide two reasons for why failing this test is problematic, especially for autoregressive models.

67 First, noise is natural — and being less robust to noise also means being a poorer fit to natural data.  
68 Outside of the log-likelihood, measures of generative success in generative models fall under two  
69 categories: qualitative assessments (*i.e.*, the no-reference perceptual quality assessment Wang et al.  
70 [2002] or ‘eyeballing’ it) and quantitative heuristics (*i.e.*, computing statistics of hidden activations of  
71 pretrained CNNs Salimans et al. [2016], Heusel et al. [2017], Sajjadi et al. [2018]). Both strategies  
72 either rely directly on the human visual system, or are known to be closely related to it Güçlü and  
73 van Gerven [2015], Yamins et al. [2014], Khaligh-Razavi and Kriegeskorte [2014], Eickenberg et al.  
74 [2017], Cichy et al. [2016]. Therefore, implicit in the use of these criteria is the existence of a human  
75 (or human-like) model of images  $q_{human}$ , where  $q_{human} \approx p_{data}$  Huszár [2015]. The fact that  
76 we find samples from  $p_\pi$  nearly indistinguishable from  $p_{data}$ , whereas  $p_\theta$  finds them very different  
77 suggests that  $p_{data} \approx q_{human} \neq p_\theta$ .

78 Second, sample quality suffers. This holds for general likelihood models, given what we argue  
79 in the first point — namely  $p_\theta \neq q_{human}$ . However, noise-sensitivity is doubly problematic in  
80 autoregressive models. Due to the sequential nature of autoregressive sampling and the fact that  
81 models are trained entirely on data from the *true* distribution, any sampling error can drastically  
82 affect the sampling trajectory. This is related to the well-known *covariate shift* phenomenon Bengio  
83 et al. [2015], Shimodaira [2000] Moreover, such errors compound quickly. Table 1 shows that  
84 mis-sampling pixels by even a single bit can cause drastic changes to the overall likelihood. This can  
85 explain why standard autoregressive models commonly produce nonsensical results (Fig 1).

## 86 3 Noise-conditional Maximum Likelihood

87 To alleviate the problems discussed in Section 2, we propose a simple modification to the standard  
88 objective in maximum likelihood estimation. Rather than evaluating a single likelihood as in the  
89 vanilla formulation, we consider a family of noise-conditional likelihoods

$$\mathbb{E}_{\sigma \sim \mu} \mathbb{E}_{\mathbf{x} \sim p_\sigma} \log p_{\theta, \sigma}(\mathbf{x}), \quad (2)$$

90 where  $p_\sigma$  is a stochastic process indexed by noise scales  $\sigma$  describing a noise-corrupted version of  
91  $p_{data}$ , and  $\mu$  is a distribution over such scales. We call this the noise-conditional maximum likelihood  
92 (NCML) framework. In general, (2) is an all-purpose plug-in objective that can be used with any  
93 likelihood model adapted to accept a noise conditioning vector, though a continuous likelihood (e.g.  
94 Salimans et al. [2017], Li and Kluger [2022]) is necessary for computation of the score function.

95 Letting  $\sigma$  be the time index of a diffusion process, our approach becomes closely related to score-based  
96 diffusion models Song et al. [2020], albeit with two crucial differences.

97 First, instead of merely estimating the noise-conditional score function of the perturbed data density  
98  $p_\sigma$  for  $\sigma \in [0, T]$ , we directly estimate  $p_\sigma$  itself. However, we still learn the noise-conditional  
99 score function as a by-product of NCML. Moreover, we may access the score function by simply  
100 differentiating the log likelihood. Therefore, we can refine sampled points via Langevin dynamics.  
101 This provides an alternative strategy for sampling from  $p_{\theta, \sigma}$ , which we explore in 3.1.

102 Second, we need not design our diffusion so that  $p_T$  approximates the limiting stationary distribution  
103 of the process. This is necessary in diffusion models as the limiting prior is the only tractable  
104 distribution to initialize the sampling algorithm with. Since we have learned the density itself for all  
105  $\sigma \in [0, T]$ , we may initialize from any point of the diffusion, which increases the flexibility of the  
106 sampling strategy, and can drastically reduce the steps required to solve the reverse diffusion.

Model	CIFAR-10			ImageNet 64x64			
	FID	NLL $\pi = 0^*$	NLL $\pi = 0.5$	NLL $\pi = 1$	NLL $\pi = 0^*$	NLL $\pi = 0.5$	NLL $\pi = 1$
<b>ELBO</b>							
VDM	7.41	<b>2.49</b>	<b>3.75</b>	<b>3.97</b>	<b>3.40</b>	3.76	3.88
ScoreFlow	<b>5.40</b>	2.90	3.82	3.99	-	-	-
VDVAE	-	2.84	3.90	4.10	3.52	<b>3.66</b>	<b>3.82</b>
<b>Likelihood</b>							
Flow++	-	3.09	3.86	4.08	3.69	3.82	3.99
DenseFlow	48.15	2.98	3.80	4.02	3.35	3.68	3.85
PixelCNN++	55.72	2.92	3.84	4.01	3.52	3.84	4.00
PixelSNAIL	36.62	2.85	3.83	3.99	-	-	-
Sparse Transformer	37.50	<b>2.80</b>	3.82	3.98	3.44	3.73	3.89
NCPN (NCML-VE)	32.71	2.90	3.77	3.98	<b>3.32</b>	3.67	3.85
NCPN (NCML-subVP)	23.42	2.95	3.69	3.94	3.36	3.66	3.82
NCPN (NCML-VP)	<b>13.50</b>	3.42	<b>3.65</b>	<b>3.91</b>	3.50	<b>3.64</b>	<b>3.80</b>

Table 1: Results on CIFAR-10 and ImageNet 64x64. Negative log-likelihood (NLL) is in bits per dimension. Lower is better. \*NLL with  $\pi = 0$  is equivalent to NLL of the original data.

### 107 3.1 Sampling with Autoregressive NCML Models

108 The NCML framework allows for two sampling strategies. The first is to draw directly from the noise-  
109 free distribution  $p_{\theta,0}$ , in which case the conditional likelihood simplifies to a standard (unconditional)  
110 likelihood, and sampling is identical to that for a standard autoregressive model.

111 However, as discussed in Section 2, this strategy is unstable and tends to quickly accumulate errors.  
112 This motivates an alternative sampling strategy, which involves drawing from  $p_{\theta,\sigma}$  for  $\sigma > 0$ , then  
113 solving a reverse diffusion process back to  $\sigma = 0$ . The latter is possible due to the fact that the reverse  
114 diffusion is itself a diffusion process that depends on the score function Anderson [1982], which we  
115 have access to. This is identical to the sampling procedure in score-based diffusion models Song et al.  
116 [2020], except for the key difference that we need not initialize with the prior distribution.

## 117 4 Experiments and Discussion

118 In all experiments in Table 1, we choose  $p_\sigma$  to be the variance exploding (VE), variance preserving  
119 (VP), and sub-variance preserving (sub-VP) SDEs, respectively. Due to space constraints, we refer  
120 to Song et al. [2020] for more details. For our architecture, we introduce the noise-conditional  
121 pixel-wise network (NCPN), which consists of a PixelCNN backbone with added attention layers.  
122 We evaluate all models on minimally perturbed transformations (see 2.1) of CIFAR-10 and ImageNet  
123 64x64 for  $\pi \in \{0, \frac{1}{2}, 1\}$ , where we note that  $p_{\pi=0} = p_{data}$ . All noise-conditional models, *i.e.*, ours,  
124 VDM Kingma et al. [2021], and ScoreFlow Song et al. [2021], are evaluated at  $t = 0$ .

125 While it is clear that all models have reduced likelihoods when evaluated on the perturbed distribution  
126  $p_\pi$ ,  $\pi \in \{\frac{1}{2}, 1\}$ , we note that our models are more robust to such transformations, even though they  
127 are evaluated under the noiseless condition, and trained on a different class of noise, *i.e.*, the marginal  
128 likelihoods of diffusion processes. Furthermore, sample quality across all models correlates better  
129 with likelihoods on the perturbed distributions than likelihoods on the base distribution.

## 130 5 Conclusion

131 We proposed a simple sanity test for checking the robustness of likelihoods to minor perturbations.  
132 We found that most likelihood models are not robust under this test, and developed a new framework  
133 that improves performance in this setting, with substantial improvements in training and sampling.

## 134 References

- 135 Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their*  
136 *Applications*, 12(3):313–326, 1982.
- 137 Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence  
138 prediction with recurrent neural networks. *Advances in neural information processing systems*, 28,  
139 2015.
- 140 Xi Chen, Nikhil Mishra, Mostafa Rohaninejad, and Pieter Abbeel. Pixelsnail: An improved autore-  
141 gressive generative model. In *International Conference on Machine Learning*, pages 864–872.  
142 PMLR, 2018.
- 143 Rewon Child. Very deep vaes generalize autoregressive models and can outperform them on images.  
144 *arXiv preprint arXiv:2011.10650*, 2020.
- 145 Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse  
146 transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- 147 Radoslaw Martin Cichy, Aditya Khosla, Dimitrios Pantazis, Antonio Torralba, and Aude Oliva.  
148 Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object  
149 recognition reveals hierarchical correspondence. *Scientific reports*, 6(1):1–13, 2016.
- 150 Michael Eickenberg, Alexandre Gramfort, Gaël Varoquaux, and Bertrand Thirion. Seeing it all:  
151 Convolutional network layers map the function of the human visual system. *NeuroImage*, 152:  
152 184–194, 2017.
- 153 Matej Grčić, Ivan Grubišić, and Siniša Šegvić. Densely connected normalizing flows. *Advances in*  
154 *Neural Information Processing Systems*, 34:23968–23982, 2021.
- 155 Aditya Grover, Manik Dhar, and Stefano Ermon. Flow-gan: Combining maximum likelihood and  
156 adversarial learning in generative models. In *Proceedings of the AAAI conference on artificial*  
157 *intelligence*, volume 32, 2018.
- 158 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured  
159 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 160 Umut Güçlü and Marcel AJ van Gerven. Deep neural networks reveal a gradient in the complexity of  
161 neural representations across the ventral stream. *Journal of Neuroscience*, 35(27):10005–10014,  
162 2015.
- 163 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochre-  
164 iter. Gans trained by a two time-scale update rule converge to a local nash equilibrium.  
165 In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and  
166 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Cur-  
167 ran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper/2017/file/](https://proceedings.neurips.cc/paper/2017/file/8a1d694707eb0fefef65871369074926d-Paper.pdf)  
168 [8a1d694707eb0fefef65871369074926d-Paper.pdf](https://proceedings.neurips.cc/paper/2017/file/8a1d694707eb0fefef65871369074926d-Paper.pdf).
- 169 Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. Flow++: Improving flow-  
170 based generative models with variational dequantization and architecture design. In *International*  
171 *Conference on Machine Learning*, pages 2722–2730. PMLR, 2019.
- 172 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
173 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 174 Ferenc Huszár. How (not) to train your generative model: Scheduled sampling, likelihood, adversary?  
175 *arXiv preprint arXiv:1511.05101*, 2015.
- 176 Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsupervised,  
177 models may explain it cortical representation. *PLoS computational biology*, 10(11):e1003915,  
178 2014.
- 179 Dongjun Kim, Seungjae Shin, Kyungwoo Song, Wanmo Kang, and Il-Chul Moon. Soft truncation: A  
180 universal training technique of score-based diffusion model for high precision score estimation. In  
181 *International Conference on Machine Learning*, pages 11201–11228. PMLR, 2022.

- 182 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances*  
183 *in neural information processing systems*, 34:21696–21707, 2021.
- 184 Henry Li and Yuval Kluger. Neural inverse transform sampler. In *International Conference on*  
185 *Machine Learning*, pages 12813–12825. PMLR, 2022.
- 186 Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves,  
187 Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw  
188 audio. *arXiv preprint arXiv:1609.03499*, 2016.
- 189 Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing  
190 generative models via precision and recall. *Advances in neural information processing systems*, 31,  
191 2018.
- 192 Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and  
193 Xi Chen. Improved techniques for training gans. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon,  
194 and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Cur-  
195 ran Associates, Inc., 2016. URL [https://proceedings.neurips.cc/paper/2016/file/](https://proceedings.neurips.cc/paper/2016/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf)  
196 [8a3363abe792db2d8761d6403605aeb7-Paper.pdf](https://proceedings.neurips.cc/paper/2016/file/8a3363abe792db2d8761d6403605aeb7-Paper.pdf).
- 197 Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the  
198 pixelcnn with discretized logistic mixture likelihood and other modifications. *arXiv preprint*  
199 *arXiv:1701.05517*, 2017.
- 200 Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-  
201 likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 202 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
203 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*  
204 *arXiv:2011.13456*, 2020.
- 205 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of  
206 score-based diffusion models. *Advances in Neural Information Processing Systems*, 34:1415–1428,  
207 2021.
- 208 Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao,  
209 Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient  
210 transformers. *arXiv preprint arXiv:2011.04006*, 2020.
- 211 Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative  
212 models. *arXiv preprint arXiv:1511.01844*, 2015.
- 213 Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks.  
214 In *International conference on machine learning*, pages 1747–1756. PMLR, 2016.
- 215 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz  
216 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*  
217 *systems*, 30, 2017.
- 218 Zhou Wang, Hamid R Sheikh, and Alan C Bovik. No-reference perceptual quality assessment of  
219 jpeg compressed images. In *Proceedings. International conference on image processing*, volume 1,  
220 pages I–I. IEEE, 2002.
- 221 Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J  
222 DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual  
223 cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.

## 224 **A Appendix**

### 225 **A.1 Additional Experimental Details**

226 We compare against Kingma et al. [2021], Song et al. [2021], Child [2020], Ho et al. [2019], Grcić  
227 et al. [2021], Van Den Oord et al. [2016], Chen et al. [2018], Child et al. [2019]. Some results could  
228 not be included due to the irreproducibility of the techniques.

229 For our NCML-trained models, the diffusion times of the VE, VP, and sub-VP SDEs were chosen to  
230 be  $t = 0.5$ ,  $t = 0.1$ , and  $t = 0.025$ , respectively. The values are somewhat arbitrary, but chosen such  
231 that the standard deviation of the per-pixel differences between samples in  $p_{data}$  and their noised  
232 counterparts in  $p_T$  was  $\approx 10$ . We suspect that further improvements can be made to the empirical  
233 results if these numbers were chosen more judiciously.

### 234 **A.2 Additional Figures**

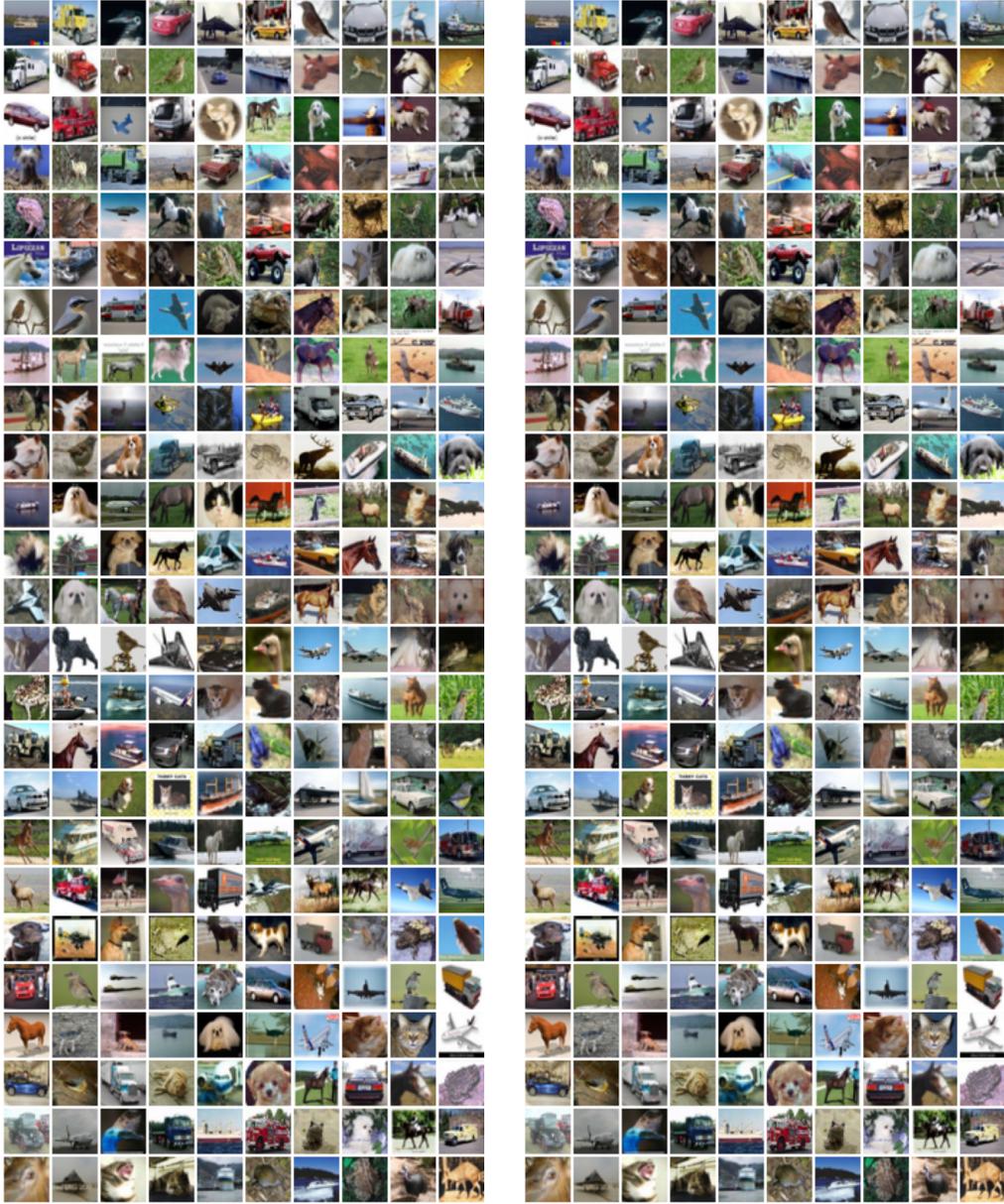


Figure 2: Images from  $p_{data}$  (left) versus their corresponding noisy counterparts from  $p_{\pi}$  (right), where  $\pi = 1$ . Differences are almost, if not entirely, imperceptible to the human eye.