# GENERATIVE UNCERTAINTY IN DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

## Abstract

Diffusion and flow matching models have recently driven significant breakthroughs in generative modeling. While state-of-the-art models produce high-quality samples on average, individual samples can still be low quality. Detecting such samples without human inspection remains a challenging task. To address this, we propose a Bayesian framework for estimating the *generative uncertainty* of synthetic samples. We outline how to make Bayesian inference practical for large, modern generative models and introduce a new semantic likelihood to address the challenges posed by high-dimensional sample spaces. Through our experiments, we demonstrate that the proposed generative uncertainty effectively identifies poor-quality samples and significantly outperforms existing uncertainty-based methods. Notably, our Bayesian framework can be applied *post-hoc* to any pretrained diffusion or flow matching model (via the Laplace approximation), and we propose simple yet effective techniques to minimize its computational overhead during sampling.

023 024

025

000

001 002 003

004

005 006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

Diffusion and flow-matching models (Sohl-Dickstein et al., 2015; Song et al., 2020a;b; Lipman et al., 2022) have recently pushed the boundaries of generative modeling due to their strong theoretical underpinnings and easy-to-scale nature. Across various domains, they have enabled the generation of increasingly realistic samples (Rombach et al., 2022; Esser et al., 2024; Li et al., 2024). Despite impressive progress, state-of-the-art models can still generate low-quality images that contain artifacts or fail to align with the provided conditioning information. As a result, users deploying these models may need multiple generations to manually find a high-quality sample. This raises a key question: *how can we detect poor generations?* 

034 Bayesian inference has long been applied to detect poor-quality predictions in predictive models (Gal et al., 2016; Wilson, 2020; Arbel et al., 2023). By capturing the uncertainty of the model 035 parameters due to limited training data, each prediction can be assigned a predictive uncertainty, 036 which, when high, serves as a warning that the prediction may be unreliable. Despite its widespread 037 use for principled uncertainty quantification in predictive models, Bayesian methodology has been 038 far less commonly applied to detecting poor generations in generative modeling. One notable exception is BayesDiff (Kou et al., 2024), which uses Bayesian uncertainty to filter out low-quality 040 samples in diffusion models for natural images. However, BayesDiff's uncertainty estimates are 041 not particularly indicative of a sample's actual quality and instead serve primarily as a detector of 042 images with 'cluttered' backgrounds (Fig. 5). 043

In this work, we propose a Bayesian framework for estimating generative uncertainty in modern 044 generative models, such as diffusion. To scale Bayesian inference for large diffusion models, we 045 employ the (last-layer) Laplace approximation (Daxberger et al., 2021a). Additionally, to address 046 the challenge posed by the high-dimensional sample spaces of data such as natural images, we 047 introduce a semantic likelihood, where we leverage pretrained image encoders (such as CLIP 048 (Radford et al., 2021)) to compute variability in a latent, *semantic* space instead. Through our experiments, we demonstrate that generative uncertainty is an effective tool for detecting low-quality 050 samples and propose simple strategies to minimize the sampling overhead introduced by Bayesian 051 inference. In particular, we make the following contributions:

- 052
- 1. We formalize the notion of *generative uncertainty* and propose a method to estimate it for modern generative models (Section 3). Analogous to how predictive uncertainty helps

identify unreliable predictions in predictive models, generative uncertainty can be used to detect low-quality generations in generative models.

- 2. We show that our generative uncertainty *strongly outperforms* previous uncertainty-based approaches for filtering out poor samples (Kou et al., 2024; De Vita & Belagiannis, 2025). Additionally, we achieve competitive performance with non-uncertainty-based methods, such as realism scores (Kynkäänniemi et al., 2019) and rarity scores (Han et al., 2023) (Section 4.1).
  - 3. We propose effective strategies to reduce the sampling overhead of Bayesian uncertainty (Section 4.2) and demonstrate the applicability of our framework beyond diffusion models by applying it to a (latent) flow matching model (Section D.7).

### 2 BACKGROUND

055

056

057

058

059

060

061

062

063

064 065

066 067

068

074

083 084

085

087

105 106

107

#### 2.1 GENERATIVE MODELING

**Sampling in Generative Models** Modern deep generative models like variational autoencoders (VAEs), generative adversarial networks (GANs), and diffusion models differ in their exact probabilistic frameworks and training schemes, yet share a common sampling recipe: start with random noise and transform it into a new data sample (Tomczak, 2022). Specifically, let  $x \in \mathcal{X}$ denote a data sample and  $z \in \mathcal{Z}$  an initial noise. A new sample is generated by:

$$oldsymbol{z} \sim p(oldsymbol{z}) \,, \ \ \hat{oldsymbol{x}} = g_{ heta}(oldsymbol{z})$$

where  $p(\mathbf{z})$  is an initial noise (prior) distribution, typically a standard Gaussian  $\mathcal{N}(0, I)$ , and  $g_{\theta}: \mathcal{Z} \to \mathcal{X}$  is a generator function with model parameters  $\theta \in \mathbb{R}^{P}$ .

**Diffusion Models** The primary focus of this work is on diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020b).<sup>1</sup> These models operate by progressively corrupting data into Gaussian noise and learning to reverse this process. For a data sample  $x_0 \sim q(x)$ , the forward (noising) process is defined as

$$\boldsymbol{x}_t = \sqrt{\bar{lpha}_t} \boldsymbol{x}_0 + \sqrt{1 - \bar{lpha_t}} \epsilon, \ \epsilon \sim \mathcal{N}(0, I)$$

where  $\bar{\alpha}_t = \prod_{s=1}^t (1-\beta_s)$  and  $\{\beta_s\}_{s=1}^T$  is a noise schedule chosen such that  $x_T \sim \mathcal{N}(0, I)$  (approximately). In the backward process, a denoising network,  $f_{\theta}$ , is learned via a simplified regression objective (among various possible parameterizations, see Song et al. (2020b) or Karras et al. (2022)):

$$\mathcal{L}(\theta; \mathcal{D}) = \mathbb{E}_{t, \boldsymbol{x}_0, \epsilon} \left[ \left| f_{\theta}(\sqrt{\bar{\alpha}_t} \boldsymbol{x}_0 + \sqrt{1 - \bar{\alpha_t}} \epsilon, t) - \epsilon \right| \right| _2^2 \right]$$
(1)

where  $\mathcal{D} = \{x_n\}_{n=1}^N$  denotes a training dataset of images. After training, diffusion models generate new samples via a generator function,  $g_{\hat{\theta}}$ , which consists of sequentially applying the learned denoiser,  $f_{\hat{\theta}}$ , and following specific transition rules from samplers such as DDPM (Ho et al., 2020) or DDIM (Song et al., 2020a).

## 094 2.2 BAYESIAN DEEP LEARNING

096 Bayesian neural networks (BNNs) go beyond point predictions and allow for principled uncertainty quantification (MacKay, 1995; Kendall & Gal, 2017; Jospin et al., 2022). Let  $h_{\psi}: \mathcal{X} \to \mathcal{Y}$  denote 097 a predictive model with parameters  $\psi \in \mathbb{R}^O$  and  $\mathcal{D} = \{(\bm{x}_n, \bm{y}_n)\}_{i=1}^N$  denote training data. A 098 point-prediction model would find a single fixed set of parameters,  $\hat{\psi} = \arg \max \mathcal{L}(\psi; \mathcal{D})$ , that 099 maximizes an objective function, L. Meanwhile, a Bayesian Neural Network (BNN) specifies a 100 prior,  $p(\psi)$ , over model parameters and defines a likelihood,  $p(\boldsymbol{y}|h_{\psi}(\boldsymbol{x}))$ , which together yield a 101 posterior distribution via Bayes' rule:  $p(\psi|\mathcal{D}) \propto p(\psi) \prod_{n=1}^{N} p(\boldsymbol{y}_n | h_{\psi}(\boldsymbol{x}_n))$ . Under this Bayesian view, a predictive model for a new test point  $\boldsymbol{x}_*$  is then obtained via the posterior predictive 102 103 distribution (Murphy, 2022): 104

$$p(\mathbf{y}|\boldsymbol{x}_*, \mathcal{D}) = \mathbb{E}_{p(\psi|\mathcal{D})} [p(\mathbf{y}|h_{\psi}(\boldsymbol{x}_*))].$$

<sup>&</sup>lt;sup>1</sup>Our framework extends beyond diffusion models and can be applied to other generative model families. See Figure TODO for a demonstration on flow matching models (Lipman et al., 2022).

108 For large models, finding the exact posterior distribution is computationally intractable, hence an 109 approximate posterior  $q(\psi|\mathcal{D})$  is used instead. Popular approaches for approximate inference in-110 clude deep ensembles (Lakshminarayanan et al., 2017), variational inference (Blundell et al., 2015; 111 Zhang et al., 2018), SWAG (Mandt et al., 2017; Maddox et al., 2019), and Laplace approximation 112 Daxberger et al. (2021a). Moreover, to alleviate computational overhead, it is common to give a "Bayesian treatment" only to a subset of parameters (Kristiadi et al., 2020; Daxberger et al., 113 2021b; Sharma et al., 2023). Finally, the intractable expectation integral in the posterior predictive 114 is approximated via Monte-Carlo (MC) sampling: 115

116

118

125

126 127

128

129

130

131

132

133

134

136

139

140

144

145

 $p(\mathbf{y}|\boldsymbol{x}_*, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^{M} p(\mathbf{y}|h_{\psi_m}(\boldsymbol{x}_*)), \ \psi_m \sim q(\psi|\mathcal{D}),$ (2)

with M denoting the number of MC samples. By measuring the variability of the posterior predictive distribution, e.g., its entropy, one can obtain an estimate of the model's predictive uncertainty for a given test point  $u(x_*)$ . The utility of such uncertainties has been demonstrated on a wide range of tasks like out-of-distribution (OOD) detection (Daxberger et al., 2021a) and active learning (Gal et al., 2017).

## 3 GENERATIVE UNCERTAINTY VIA BAYESIAN INFERENCE

While Bayesian neural networks (BNNs) have traditionally been applied to predictive models to estimate *predictive uncertainty*, we demonstrate how to apply them to diffusion to estimate *generative uncertainty* in this section (see Figure 4 for an overview of our method). Later in Section 4, we show that generative uncertainty can be used to detect poor-quality samples. Our focus is on generative models for natural images, where  $x \in \mathbb{R}^{H \times W \times C}$ . For ease of exposition, we consider unconditional generation in this section, though our methodology can also be applied directly to conditional models.

## 135 3.1 BAYESIAN DIFFUSION

As in traditional Bayesian predictive models (cf. Section 2.2), the central principle for obtaining a
 Bayesian notion of uncertainty in diffusion models is the posterior predictive distribution:

$$p(\mathbf{x}|\boldsymbol{z}, \mathcal{D}) = \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} \big[ p(\mathbf{x}|g_{\boldsymbol{\theta}}(\boldsymbol{z})) \big].$$
(3)

Here, we use z (with a slight abuse of notation) to denote the entire randomness involved in the diffusion sampling process.<sup>2</sup> Generative uncertainty is then defined as the variability of the posterior predictive:

$$u(\boldsymbol{z}) := H(p(\mathbf{x}|\boldsymbol{z}, \mathcal{D})) \tag{4}$$

where  $H(\cdot)$  denotes the entropy. We propose a tractable estimator of entropy later in Eq. 8. We choose entropy as the variability measure due to its simplicity and widespread use in quantifying predictive uncertainty. However, we note that alternative measures of variability, such as pairwise-distance estimators (PAiDEs) (Berry & Meger, 2023), can also be employed.

In the same way that the predictive uncertainty  $u(x_*)$ , of a predictive model  $h_{\psi}$  provides insight into the quality of its prediction for a new test point  $x_*$ , the generative uncertainty u(z) of a diffusion model  $g_{\theta}$  should offer information about the quality of the generation  $g_{\theta}(z)$  for a "new" random noise sample z. We demonstrate this relationship experimentally in Section 4. Next, we discuss how to make Bayesian inference on (large) diffusion models computationally tractable.

3.2 LAST-LAYER LAPLACE APPROXIMATION

State-of-the-art diffusion models are extremely large (100M to 1B+ parameters) and can take weeks to train. Consequently, the computational overhead of performing Bayesian inference on such large models is a significant concern. To address this, we adopt the Laplace approximation (MacKay,

160 161

155

156

<sup>2</sup>For example, in DDIM (Song et al., 2020a) and ODE sampling (Song et al., 2020b),  $\boldsymbol{z} = \epsilon_T$ , whereas in DDPM (Ho et al., 2020) and SDE sampling (Song et al., 2020b),  $\boldsymbol{z} = \{\epsilon_T, \ldots, \epsilon_1\}$  with  $\epsilon_t \sim \mathcal{N}(0, I)$ .

167 168 The Laplace approximation of the posterior is given by:

$$q(\theta|\mathcal{D}) = \mathcal{N}(\theta|\hat{\theta}, \Sigma), \ \Sigma = \left(\nabla_{\theta}^{2} \mathcal{L}(\theta; \mathcal{D})\big|_{\hat{\theta}}\right)^{-1},\tag{5}$$

where  $\theta$  represents the parameters of a pre-trained diffusion model, and  $\Sigma$  is the inverse Hessian of the diffusion (regression) loss from Eq. 1. To reduce the computational cost further, we apply a "Bayesian" treatment only to the last layer of the denoising network  $f_{\theta}$ .

174 We note that the use of last-layer Laplace approximation for diffusion models has been previously 175 proposed in BayesDiff (Kou et al., 2024). While our implementation of the Laplace approximation 176 closely follows theirs, there are significant differences in how we utilize the approximate posterior,  $q(\theta|\mathcal{D})$ . Specifically, in our approach, we use it within the traditional Bayesian framework (Eq. 3) to 177 sample new diffusion model parameters, leaving the diffusion sampling process,  $g_{\theta}$ , unchanged. In 178 contrast, BayesDiff resamples new weights from  $q(\theta|D)$  at every diffusion sampling step t, which 179 necessitates substantial modifications to the diffusion sampling process through their variance 180 propagation approach. We later demonstrate in Section 4 that modifications such as variance 181 propagation are unnecessary for obtaining Bayesian generative uncertainty and staying closer to the 182 traditional Bayesian setting leads to the best empirical performance. 183

## 3.3 SEMANTIC LIKELIHOOD

169 170

189 190

202

203

We next discuss the choice of likelihood for estimating generative uncertainty in diffusion models.
 Since the denoising problem in diffusion is modeled as a (multi-output) regression problem, the most straightforward approach is to place a simple Gaussian distribution over the generated sample:

$$p(\mathbf{x}|g_{\theta}(\boldsymbol{z})) = \mathcal{N}(\mathbf{x}|g_{\theta}(\boldsymbol{z}), \sigma^2 I), \tag{6}$$

191 where  $\sigma^2$  represents the observation noise.

However, as we will demonstrate in Section 4, this likelihood leads to non-informative estimates 193 of generative uncertainty (Eq. 4). The primary issue is that the sample space of natural images 194 is high-dimensional (i.e.,  $|\mathcal{X}| = HWC$ ). Consequently, placing the likelihood directly in the 195 sample space causes the variability of the posterior predictive distribution to be based on pixel-level 196 differences. This is problematic because it is well-known that two images can appear nearly 197 identical to the human eye while exhibiting a large  $L_2$ -norm difference in pixel space  $\mathcal{X}$  (see, for example, the literature on adversarial examples (Szegedy, 2013)). To get around this, we propose 199 to map the generated samples to a "semantic" latent space, S, via a pre-trained feature extractor,  $c_{\phi}: \mathcal{X} \to \mathcal{S}$  (e.g., an inception-net (Szegedy et al., 2016) or a CLIP encoder (Radford et al., 2021)). 200 The resulting *semantic likelihood* has the form 201

$$p(\mathbf{x}|g_{\theta}(\boldsymbol{z});\phi) = \mathcal{N}(\mathbf{e}(\mathbf{x})|c_{\phi}(g_{\theta}(\boldsymbol{z})),\sigma^{2}I)$$
(7)

204 where  $\mathbf{e}(\mathbf{x}) \in S$  is the vector of extracted semantic features.

By combining the (last-layer) Laplace approximate posterior and the semantic likelihood, we can now approximate the posterior predictive (Eq. 3) as

$$p(\mathbf{x}|\boldsymbol{z}, \mathcal{D}) pprox \mathcal{N}ig(\mathbf{e}(\mathbf{x}) \mid ar{m{e}}, \operatorname{Diag}ig(rac{1}{M}\sum_{m=1}^M m{e}_m^2 - ar{m{e}}^2ig) + \sigma^2ig),$$

$$\bar{\boldsymbol{e}} = \frac{1}{M} \sum_{m=1}^{M} \boldsymbol{e}_{m}, \quad \boldsymbol{e}_{m} = c_{\phi} \big( g_{\theta_{m}}(\boldsymbol{z}) \big), \ \theta_{m} \sim q(\theta|\mathcal{D}) , \tag{8}$$

where *M* denotes the number of Monte Carlo samples. Additionally, we approximate the posterior predictive with a single Gaussian via moment matching here, a common practice in Bayesian neural networks for regression problems (Lakshminarayanan et al., 2017; Antorán et al., 2020). 216 Unlike in the posterior predictive for predictive models (Eq. 2), where it is used to obtain both the 217 prediction and the associated uncertainty, the generative posterior predictive (Eq. 8) is used solely 218 to estimate the generative uncertainty u(z). The actual samples  $\hat{x}$  are still generated using the 219 pre-trained diffusion model  $g_{\hat{\theta}}$  (see Algorithm 1).

220 221

222

3.4 EPISTEMIC UNCERTAINTY

Bayesian uncertainty is commonly decomposed into two components: *aleatoric* and *epistemic* 223 uncertainty (Hüllermeier & Waegeman, 2021; Smith et al., 2024). Aleatoric uncertainty represents 224 the irreducible uncertainty inherent in the data-generating process, while epistemic uncertainty 225 arises from observing only a limited amount of training data. In our framework, we fix the 226 observation noise in the semantic likelihood (Eq. 7) to a small constant value (e.g.,  $\sigma = 0.001$ ). As 227 a result, the generative uncertainty we capture is primarily epistemic in nature, reflecting uncertainty 228 about the model parameters  $\theta$  due to limited training data via  $q(\theta | \mathcal{D})$ . Extending our framework to 229 additionally capture the aleatoric uncertainty of a generative process presents an interesting avenue 230 for future research. Furthermore, since the parameters  $\phi$  of the semantic feature extractor  $c_{\phi}$  are 231 kept fixed in the semantic likelihood, the resulting generative uncertainty u(z) continues to reflect 232 the epistemic uncertainty of the diffusion model parameters  $\theta$ .

233 234

235

242

243 244

245

246

247

248

249

250 25

4 **EXPERIMENTS** 

236 In our experiments, we demonstrate that generative uncertainty is an effective method for detecting 237 poor samples in diffusion models (Section 4.1). We also discuss the sampling overhead introduced 238 by our Bayesian approach and show that it can be effectively minimized (Section 4.2). Finally, we 239 extend our Bayesian framework beyond diffusion by applying it to detect low-quality samples in a (latent) flow matching model (Appendix D.7). Code to replicate all our experiments and figures is 240 publicly available at GITHUB REPO. 241

## 4.1 DETECTING POOR GENERATIONS

Table 1: Image generation results for 10K filtered samples (out of 12K) based on various metrics. Our generative uncertainty outperforms previously proposed uncertainty-based approaches in terms of image quality (AU (De Vita & Belagiannis, 2025), BayesDiff (Kou et al., 2024)), as indicated by higher FID and precision scores, and is competitive with non-uncertainty methods (Realism (Kynkäänniemi et al., 2019), Rarity (Han et al., 2023)). We report mean values along with standard deviations over 3 runs with different random seeds.

251		ADM (DI	DIM), ImageNe	ImageNet 128×128		UViT (DPM), ImageNet 256×256		
252		<b>FID</b> $(\downarrow)$	<b>Precision</b> $(\uparrow)$	Recall $(\uparrow)$	<b>FID</b> $(\downarrow)$	<b>Precision</b> $(\uparrow)$	Recall $(\uparrow)$	
254	Random	$11.31\pm0.07$	$58.90 \pm 0.36$	$70.68 \pm 0.38$	$9.46 \pm 0.12$	$60.94 \pm 0.24$	$73.82 \pm 0.33$	
255	BayesDiff	$11.20\pm0.05$	$58.80 \pm 0.05$	$70.62\pm0.32$	$9.16 \pm 0.17$	$61.77\pm0.19$	$73.72\pm0.38$	
56	AŬ	$11.39\pm0.05$	$58.82 \pm 0.42$	$70.70\pm0.38$	$9.20 \pm 0.12$	$61.80 \pm 0.33$	$73.46\pm0.24$	
57	Ours	$10.14\pm0.08$	$61.26\pm0.26$	$69.60 \pm 0.49$	$7.89 \pm 0.12$	$64.14\pm0.17$	$71.92 \pm 0.35$	
58	Realism	$9.76 \pm 0.04$	$67.95 \pm 0.19$	$66.32 \pm 0.40$	$8.24 \pm 0.09$	$70.29 \pm 0.15$	$69.12 \pm 0.32$	
259	Rarity	$10.09\pm0.02$	$64.99 \pm 0.16$	$67.73 \pm 0.47$	$8.37 \pm 0.11$	$67.21 \pm 0.10$	$67.76 \pm 0.48$	
260								

261 To evaluate whether our newly introduced generative uncertainty can be used to detect low-quality 262 generations, we follow the experimental setup from prior work on uncertainty-based filtering (Kou 263 et al., 2024; De Vita & Belagiannis, 2025). Specifically, we generate 12K samples using a given 264 diffusion model and compute the uncertainty estimate for each sample. We then select the 10K 265 samples with the *lowest* uncertainty. If uncertainty reliably reflects the visual quality of generated 266 samples, filtering based on it should yield greater improvements in population-level metrics (such as FID) compared to selecting a random subset of 10K images. 267

268

**Implementation Details** To ensure a fair comparison with BayesDiff (Kou et al., 2024), we 269 adopt their proposed implementation of the last-layer Laplace approximation. Specifically, we use



Figure 1: Images with the highest (*left*) and the lowest (*right*) generative uncertainty (Eq. 8) among 12K generations using a UViT diffusion model (Bao et al., 2023). Generative uncertainty correlates with visual quality, as high-uncertainty samples exhibit numerous artefacts, whereas low-uncertainty samples resemble canonical images of their respective conditioning class.



Figure 2: Images with the highest (*bottom*) and the lowest (*top*) generative uncertainty among 128 generations using a UViT diffusion model (Bao et al., 2023) for 2 classes: black swan (*left*) and Tibetan terrier (*right*).

an Empirical Fisher approximation of the Hessian with a diagonal factorization (Daxberger et al., 2021a). When computing the posterior predictive distribution (Eq. 8), we use M = 5 Monte Carlo samples. For the semantic feature extractor  $c_{\phi}$ , we leverage a pretrained CLIP encoder (Radford et al., 2021). Additional implementation details are provided in Appendix E.

**Baselines** We first compare our proposed generative uncertainty to existing uncertainty-based approaches for detecting low-quality samples: BayesDiff (Kou et al., 2024) and the aleatoric uncertainty (AU) approach proposed by De Vita & Belagiannis (2025). BayesDiff estimates epistemic uncertainty in diffusion models using a last-layer Laplace approximation and tracks this uncertainty throughout the entire sampling process via their variance propagation method. In contrast, De Vita & Belagiannis (2025) computes aleatoric uncertainty by measuring the sensitivity of intermediate diffusion scores to random perturbations. Unlike our approach, both methods estimate uncertainty directly in pixel space. 

Importantly, we also compare our method against non-uncertainty-based sample-level metrics, such as the *realism* score (Kynkäänniemi et al., 2019) and the *rarity* score (Han et al., 2023). These metrics work by measuring the distance of a generated sample from the data manifold (derived from a reference dataset) in a semantic space spanned by the inception-net features (Szegedy et al., 2016).
Notably, prior work (Kou et al., 2024; De Vita & Belagiannis, 2025) has not considered such comparisons, which we believe are essential for assessing the practical utility of uncertainty-based filtering.

324 **Evaluation Metrics** In addition to the widely used Fréchet Inception Distance (FID) (Heusel 325 et al., 2017) for evaluating the quality of a filtered set of images, we also report *precision* and 326 recall metrics (Sajjadi et al., 2018; Kynkäänniemi et al., 2019). To compute these quantities we 327 fit two manifolds in feature space: one for the generated images and another for the reference 328 (training) images. Precision is the proportion of generated images that lie in the reference image manifold while recall is the proportion of reference images that lie in the generated image manifold. 329 Precision measures the quality (or fidelity) of generated samples, whereas recall quantifies their 330 diversity (or coverage over the reference distribution). 331

332

Results We present our main results on the ImageNet dataset in Table 1. We first observe that
 existing uncertainty-based approaches (BayesDiff and AU) result in little to no improvement in
 metrics that assess sample quality (FID and precision). In contrast, our generative uncertainty
 method leads to significant improvements in terms of both FID and precision. For example, on
 the UViT model (Bao et al., 2023), a subset of images selected based on our uncertainty measure
 achieves an FID of 7.89, significantly outperforming both the Random baseline (9.45) and existing
 uncertainty-based methods (BayesDiff 9.16, AU 9.20).

Next, in order to qualitatively demonstrate the effectiveness of our approach, we show 25 samples with the highest and lowest generative uncertainty (out of the original 12K samples) according to our method in Figure 1. High-uncertainty samples exhibit numerous artefacts, and in most cases, it is difficult to determine what exactly they depict. Combined with the quantitative results in Table 1, this supports our hypothesis that (Bayesian) generative uncertainty is an effective metric for identifying low-quality samples. Conversely, the lowest-uncertainty samples are of high quality, with most appearing as 'canonical' examples of their respective (conditioning) class.

- 347 For comparison, in Figure 5 we also depict the 25 "worst" and "best" samples according to the 348 uncertainty estimate from BayesDiff (Kou et al., 2024). It is evident that their uncertainty is less 349 informative for sample quality than ours. Moreover, their uncertainty measure appears to be very 350 sensitive to the background pixels. Most images with the highest uncertainty have a 'cluttered' 351 background, whereas most images with the lowest uncertainty have a 'clear' background. We attribute this issue to the fact that in BayesDiff the uncertainty is computed directly in the pixel 352 space, unlike in our approach where we use the semantic likelihood (Section 3.3) to move away 353 from the (high-dimensional) sample space. To further verify the importance of the semantic 354 likelihood, in Figure 7 we perform an ablation where we compute our generative uncertainty 355 directly in the pixel-space. It is clear that without semantic likelihood, our uncertainty becomes 356 overly sensitive to the background pixels in the same way as in BayesDiff. 357
- Returning to Table 1, we observe that filtering based on our generative uncertainty results in some 358 loss of sample diversity, as evidenced by lower recall scores (e.g., 73.82 for Random vs. 71.92 for 359 our method on the UViT model). We attribute this to the fact that, in our main experiment, 12K 360 images are generated unconditionally.<sup>3</sup> As a result, all 1000 ImageNet classes are represented. Since 361 certain classes produce images with higher uncertainty (see Appendix D.6 for a detailed analysis), 362 filtering based on uncertainty inevitably alters the class distribution among the selected samples. 363 We expect this issue to be less pronounced in conditional generation (see Figure 2). Moreover, the 364 trade-off between improving sample quality (precision) and reducing diversity (recall) has been 365 observed before, see for example the literature on classifier-free guidance (Ho & Salimans, 2022). 366

Lastly, we compare our proposed method with non-uncertainty-based approaches—a compar-367 ison missing in prior literature (Kou et al., 2024; De Vita & Belagiannis, 2025). For realism 368 (Kynkäänniemi et al., 2019), we retain the 10K images with the highest scores, whereas for rarity 369 (Han et al., 2023), we keep those with the lowest scores. As shown in Table 1, our generative 370 uncertainty is the only uncertainty-based method that approaches realism and rarity in terms of FID 371 (e.g., 7.89 for ours vs. 8.24 for realism and 8.37 for rarity on UViT). However, a large gap remains 372 in precision (e.g., 64.14 for ours vs. 70.29 for realism and 67.21 for rarity on UViT). Notably, 373 realism and rarity sacrifice the most sample diversity, as indicated by their lowest recall scores (e.g., 374 69.12 for realism and 67.76 for rarity on UViT).

375 376

<sup>&</sup>lt;sup>3</sup>Following Kou et al. (2024), we actually still use class-conditional diffusion models but randomly sample a class for each of the 12K generated samples.

Furthermore, Table 2 shows that our score can be effectively combined with realism or rarity scores. Specifically, combining our score with realism yields an FID of 7.60 on UViT, compared to 8.26 when combining realism and rarity. We attribute higher benefits from ensembling our score to the fact that, while realism and rarity exhibit a strong negative Spearman correlation (-0.85), our uncertainty measure is less correlated with them (-0.27 with realism, 0.38 with rarity), as shown in Figure 10.

384 385

386

#### 4.2 IMPROVING SAMPLING EFFICIENCY

We next examine the sampling costs associ-387 ated with Bayesian inference in diffusion sam-388 pling. As shown in Algorithm 1, obtaining an 389 uncertainty estimate u(z) for a generated sam-390 ple  $\hat{x}_0 = g_{\theta}(z)$  requires generating M addi-391 tional samples, resulting in MT additional net-392 work function evaluations (NFEs). For the re-393 sults presented in Table 1, we use M = 5 and 394 the default number of sampling steps T = 50, 395 leading to an additional 250 NFEs for uncer-396 tainty estimation-on top of the 50 NFEs required to generate the original sample. Since 397 this overhead may be prohibitively expensive in 398 certain deployment scenarios, we next explore 399 strategies to reduce the sampling cost associ-400 ated with our generative uncertainty. 401

The most straightforward approach is to reduce the number of Monte Carlo samples M. Encouragingly, reducing M to as low as 1 still achieves highly competitive performance (see Figure 3). Further efficiency gains can be



Figure 3: FID results for 10K ImageNet-filtered images using our generative uncertainty on ADM model (Dhariwal & Nichol, 2021). We vary the number of Monte Carlo samples M and diffusion sampling steps T (see Algorithm 1). By default, we use M = 5 with T = 50, incurring an additional 250 NFEs for uncertainty estimation. Encouragingly, setting M = 1 and T = 25 still achieves competitive performance while reducing the sampling overhead by 10x.

achieved by reducing the number of sampling steps T, leveraging the flexibility of diffusion models to adjust T on the fly. Importantly, we lower T only for the additional M samples used for uncertainty assessment while keeping the default T for the original sample  $\hat{x}_0$  to ensure that the generation quality is not compromised. Taken together, reducing M and T significantly improves the efficiency of our generative uncertainty. Concretely, using the ADM model (Dhariwal & Nichol, 2021), our generative uncertainty method with M = 1 and T = 25 achieves an FID of 10.36, which still strongly outperforms both the Random (11.31) and BayesDiff (11.20) baselines while requiring only 25 additional NFEs.

414 415

416

## 5 CONCLUSION

We introduced generative uncertainty and demonstrated how to estimate it in modern generative models such as diffusion. Our experiments showed the effectiveness of generative uncertainty in filtering out low-quality samples. For future work, it would be interesting to explore broader applications of Bayesian principles in generative modeling beyond detecting poor-quality generations.
Promising directions include guiding synthetic data generation and optimizing diffusion hyperparameters via marginal likelihood using the Laplace approximation.

423

- 424
- 425
- 426
- 427 428
- 429

430

432	REFERENCES
433	Itel Bitel (CEB

434	Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying
435	framework for flows and diffusions. arXiv preprint arXiv:2303.08797, 2023.

- Anastasios N Angelopoulos, Amit Pal Kohli, Stephen Bates, Michael Jordan, Jitendra Malik, Thayer Alshaabi, Srigokul Upadhyayula, and Yaniv Romano. Image-to-image regression with distribution-free uncertainty quantification and applications in imaging. In *International Conference on Machine Learning*, pp. 717–730. PMLR, 2022.
- Javier Antorán, James Allingham, and José Miguel Hernández-Lobato. Depth uncertainty in neural networks. *Advances in neural information processing systems*, 33:10620–10634, 2020.
- Julyan Arbel, Konstantinos Pitas, Mariia Vladimirova, and Vincent Fortuin. A primer on bayesian
   neural networks: review and debates. *arXiv preprint arXiv:2309.16314*, 2023.
- Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22669–22679, 2023.
- Omer Belhasin, Yaniv Romano, Daniel Freedman, Ehud Rivlin, and Michael Elad. Principal uncertainty quantification with spatial correlation for image restoration problems. *IEEE Transactions* on Pattern Analysis and Machine Intelligence, 2023.
- Lucas Berry and David Meger. Escaping the sample trap: Fast and accurate epistemic uncertainty
   estimation with pairwise-distance estimators. *arXiv preprint arXiv:2308.13498*, 2023.
- Lucas Berry, Axel Brando, and David Meger. Shedding light on large generative networks: Estimating epistemic uncertainty in diffusion models. In *The 40th Conference on Uncertainty in Artificial Intelligence*, 2024.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in
   neural network. In *International conference on machine learning*, pp. 1613–1622. PMLR, 2015.
- Matthew Albert Chan, Maria J Molina, and Christopher Metzler. Estimating epistemic and aleatoric uncertainty with a single model. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- 464 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary 465 differential equations. *Advances in neural information processing systems*, 31, 2018.
- François Cornet, Grigory Bartosh, Mikkel N Schmidt, and Christian A Naesseth. Equivariant neural diffusion for molecule generation. In *38th Conference on Neural Information Processing Systems*, 2024.
- 470 Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space. *arXiv preprint* 471 *arXiv:2307.08698*, 2023.
- Erik Daxberger and José Miguel Hernández-Lobato. Bayesian variational autoencoders for unsupervised out-of-distribution detection. *arXiv preprint arXiv:1912.05651*, 2019.
- Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and
   Philipp Hennig. Laplace redux-effortless bayesian deep learning. Advances in Neural Information
   *Processing Systems*, 34:20089–20103, 2021a.
- Erik Daxberger, Eric Nalisnick, James U Allingham, Javier Antorán, and José Miguel HernándezLobato. Bayesian deep learning via subnetwork inference. In *International Conference on Machine Learning*, pp. 2510–2521. PMLR, 2021b.
- Michele De Vita and Vasileios Belagiannis. Diffusion model guided sampling with pixel-wise aleatoric uncertainty estimation. *IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2025.
- 485 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

486 487 488 489 490	Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In <i>Forty-first International Conference on Machine Learning</i> , 2024. URL https://openreview.net/forum?id=FPnUhsQJ5B.
491 492 493 494	Gianni Franchi, Dat Nguyen Trong, Nacim Belkhir, Guoxuan Xia, and Andrea Pilzer. Towards understanding and quantifying uncertainty for text-to-image generation, 2024. URL https://arxiv.org/abs/2412.03178.
495 496	Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In <i>International conference on machine learning</i> , pp. 1183–1192. PMLR, 2017.
497 498	Yarin Gal et al. Uncertainty in deep learning. 2016.
499 500	Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and LingPeng Kong. Diffuseq: Sequence to sequence text generation with diffusion models. <i>arXiv preprint arXiv:2210.08933</i> , 2022.
501 502 503	Shuyang Gu, Jianmin Bao, Dong Chen, and Fang Wen. Giqa: Generated image quality assessment. In <i>Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020,</i> <i>Proceedings, Part XI 16</i> , pp. 369–385. Springer, 2020.
504 505 506 507	Jiyeon Han, Hwanil Choi, Yunjey Choi, Junho Kim, Jung-Woo Ha, and Jaesik Choi. Rarity score : A new metric to evaluate the uncommonness of synthesized images. In <i>The Eleventh International Conference on Learning Representations (ICLR)</i> , 2023.
508 509 510	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. <i>Advances in neural information processing systems</i> , 30, 2017.
511 512	Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. <i>arXiv preprint arXiv:2207.12598</i> , 2022.
513 514 515	Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.
516 517 518	Emiel Hoogeboom, V1ctor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3d. In <i>International conference on machine learning</i> , pp. 8867–8887. PMLR, 2022.
519 520 521	Eyke Hüllermeier and Willem Waegeman. Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods. <i>Machine learning</i> , 110(3):457–506, 2021.
522 523 524	Jaehui Hwang, Junghyuk Lee, and Jong-Seok Lee. Anomaly score: Evaluating generative models and individual generated images based on complexity and vulnerability. In <i>Proceedings of the</i> <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 8754–8763, 2024.
525 526 527	Alexander Immer, Matthias Bauer, Vincent Fortuin, Gunnar Rätsch, and Khan Mohammad Emtiyaz. Scalable marginal likelihood estimation for model selection in deep learning. In <i>International Conference on Machine Learning</i> , pp. 4563–4573. PMLR, 2021.
529 530 531	Laurent Valentin Jospin, Hamid Laga, Farid Boussaid, Wray Buntine, and Mohammed Bennamoun. Hands-on bayesian neural networks—a tutorial for deep learning users. <i>IEEE Computational</i> <i>Intelligence Magazine</i> , 17(2):29–48, 2022.
532 533 534	Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion- based generative models. <i>Advances in neural information processing systems</i> , 35:26565–26577, 2022.
535 536 537	Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? <i>Advances in neural information processing systems</i> , 30, 2017.
538 539	Siqi Kou, Lei Gan, Dequan Wang, Chongxuan Li, and Zhijie Deng. Bayesdiff: Estimating pixel- wise uncertainty in diffusion via bayesian inference. In <i>The Twelfth International Conference on</i> <i>Learning Representations (ICLR)</i> , 2024.

551

559

565

566

567

568

570

571

572

573

- 540 Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes 541 overconfidence in relu networks. In International conference on machine learning, pp. 5436-542 5446. PMLR, 2020.
- Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent 544 space. arXiv preprint arXiv:2210.10960, 2022.
- 546 Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved 547 precision and recall metric for assessing generative models. Advances in neural information 548 processing systems, 32, 2019.
- 549 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive 550 uncertainty estimation using deep ensembles. Advances in neural information processing systems, 30, 2017. 552
- Zehui Li, Yuhao Ni, Guoxuan Xia, William Beardall, Akashaditya Das, Guy-Bart Stan, and Yiren 553 Zhao. Absorb & escape: Overcoming single model limitations in generating heterogeneous ge-554 nomic sequences. In The Thirty-eighth Annual Conference on Neural Information Processing 555 Systems, 2024. URL https://openreview.net/forum?id=XHTl2k1LYk. 556
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching 558 for generative modeling. arXiv preprint arXiv:2210.02747, 2022.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and 560 transfer data with rectified flow. arXiv preprint arXiv:2209.03003, 2022. 561
- 562 Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast 563 ode solver for diffusion probabilistic model sampling in around 10 steps. Advances in Neural Information Processing Systems, 35:5775–5787, 2022. 564
  - Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. Diffusion hyperfeatures: Searching through time and space for semantic correspondence. Advances in Neural Information Processing Systems, 36, 2024.
- 569 David JC MacKay. Bayesian interpolation. Neural computation, 4(3):415–447, 1992.
  - David JC MacKay. Bayesian neural networks and density networks. Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, 354(1):73-80, 1995.
- Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. 574 A simple baseline for bayesian uncertainty in deep learning. Advances in neural information 575 processing systems, 32, 2019. 576
- 577 D Mandt, Matthew D Hoffman, and David M Blei. Stochastic gradient descent as approximate 578 bayesian inference. Journal of Machine Learning Research, 18(134):1-35, 2017.
- 579 Kevin P Murphy. Probabilistic machine learning: an introduction. MIT press, 2022. 580
- 581 Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do 582 deep generative models know what they don't know? arXiv preprint arXiv:1810.09136, 2018.
- Koichi Namekata, Amirmojtaba Sabour, Sanja Fidler, and Seung Wook Kim. Emerdiff: Emerging 584 pixel-level semantic knowledge in diffusion models. arXiv preprint arXiv:2401.11739, 2024. 585
- 586 Katharina Ott, Michael Tiemann, and Philipp Hennig. Uncertainty and structure in neural ordinary differential equations. arXiv preprint arXiv:2305.13290, 2023.
- 588 William Peebles and Saining Xie. Scalable diffusion models with transformers. In Proceedings of 589 the IEEE/CVF International Conference on Computer Vision, pp. 4195–4205, 2023. 590
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 592 models from natural language supervision. In International conference on machine learning, pp. 8748-8763. PMLR, 2021.

594 595 596	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High- resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF Con-</i> <i>ference on Computer Vision and Pattern Recognition (CVPR)</i> , pp. 10684–10695, June 2022.
598 599	Yunus Saatci and Andrew G Wilson. Bayesian gan. Advances in neural information processing systems, 30, 2017.
600 601 602	Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. <i>Advances in neural information processing systems</i> , 31, 2018.
603 604 605	Swami Sankaranarayanan, Anastasios Angelopoulos, Stephen Bates, Yaniv Romano, and Phillip Isola. Semantic uncertainty intervals for disentangled latent spaces. In <i>NeurIPS</i> , 2022.
606 607 608	Mrinank Sharma, Sebastian Farquhar, Eric Nalisnick, and Tom Rainforth. Do bayesian neural net- works need to be fully stochastic? In <i>International Conference on Artificial Intelligence and</i> <i>Statistics</i> , pp. 7694–7722. PMLR, 2023.
609 610 611	Zhenming Shun and Peter McCullagh. Laplace approximation of high dimensional integrals. <i>Journal of the Royal Statistical Society Series B: Statistical Methodology</i> , 57(4):749–760, 1995.
612 613 614	Freddie Bickford Smith, Jannik Kossen, Eleanor Trollope, Mark van der Wilk, Adam Foster, and Tom Rainforth. Rethinking aleatoric and epistemic uncertainty. <i>arXiv preprint arXiv:2412.20892</i> , 2024.
615 616 617	Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In <i>International conference on machine learning</i> , pp. 2256–2265. PMLR, 2015.
618 619 620	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. <i>arXiv</i> preprint arXiv:2010.02502, 2020a.
621 622 623	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. <i>arXiv preprint arXiv:2011.13456</i> , 2020b.
624 625	C Szegedy. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199, 2013.
626 627 628	Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethink- ing the inception architecture for computer vision. In <i>Proceedings of the IEEE conference on</i> <i>computer vision and pattern recognition</i> , pp. 2818–2826, 2016.
629 630 631	Jacopo Teneggi, Matthew Tivnan, Web Stayman, and Jeremias Sulam. How to trust your diffusion model: A convex optimization approach to conformal risk control. In <i>International Conference on Machine Learning</i> , pp. 33940–33960. PMLR, 2023.
632 633 634	Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models, 2016. URL https://arxiv.org/abs/1511.01844.
635 636	Jakub M. Tomczak. <i>Deep Generative Modeling</i> . Springer, Germany, February 2022. ISBN 978-3-030-93157-5. doi: 10.1007/978-3-030-93158-2.
637 638 639 640	Ba-Hien Tran, Babak Shahbaba, Stephan Mandt, and Maurizio Filippone. Fully bayesian autoen- coders with latent sparse gaussian processes. In <i>International Conference on Machine Learning</i> , pp. 34409–34430. PMLR, 2023.
641 642 643	Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In <i>International conference on machine learning</i> , pp. 9690–9700. PMLR, 2020.
644 645 646	Andrew Gordon Wilson. The case for bayesian deep learning. <i>arXiv preprint arXiv:2001.10995</i> , 2020.
	Oinhand Vi Vien for Chan Channel Zhang Zahai Zhang Linen Zhu and Vien Sie Kang Diffusion

647 Qiuhua Yi, Xiangfan Chen, Chenwei Zhang, Zehai Zhou, Linan Zhu, and Xiangjie Kong. Diffusion models in text generation: a survey. *PeerJ Computer Science*, 10:e1905, 2024.

648 649 650	Cheng Zhang, Judith Bütepage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 41(8):2008–2026, 2018.
651 652 653 654	Zicheng Zhang, Haoning Wu, Chunyi Li, Yingjie Zhou, Wei Sun, Xiongkuo Min, Zijian Chen, Xiaohong Liu, Weisi Lin, and Guangtao Zhai. A-bench: Are lmms masters at evaluating ai- generated images? <i>arXiv preprint arXiv:2406.03070</i> , 2024.
655 656 657 658	Ganning Zhao, Vasileios Magoulianitis, Suya You, and C-C Jay Kuo. A lightweight generalizable evaluation and enhancement framework for generative models and generated samples. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pp. 450–459, 2024.
659	
660	
661	
662	
663	
664	
665	
666	
667	
668	
669	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

<ul> <li>The supplementary material is organized as follows:</li> <li>In Appendix A, we supplement our methods (Section 3) with an algorithm and a diagram explaining our generative uncertainty.</li> <li>In Appendix B, we describe related literature.</li> <li>In Appendix C, we point out limitations of our work.</li> <li>In Appendix D.1, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>In Appendix D.3, we perform ablations on our semantic likelihood (Section 3.3).</li> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2022).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	702 703	Appendix
<ul> <li>In Appendix A, we supplement our methods (Section 3) with an algorithm and a diagram explaining our generative uncertainty.</li> <li>In Appendix B, we describe related literature.</li> <li>In Appendix D.1, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>In Appendix D.2, we perform ablations on our semantic likelihood (Section 3.3).</li> <li>In Appendix D.3, we demonstrate how to use our generative uncertainty for pixel-wise uncertainty.</li> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	704	The supplementary material is organized as follows:
<ul> <li>explaining our generative uncertainty.</li> <li>in Appendix B, we describe related literature.</li> <li>in Appendix D I, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>in Appendix D I, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>in Appendix D I, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>in Appendix D I, we qualitatively compare our semantic likelihood (Section 3.3).</li> <li>in Appendix D A, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D 5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D 5, we further compare our generative uncertainty to looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D 7, we apply our generative uncertainty to detect low-quality samples in a faterin flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	705 706	• In Appendix A, we supplement our methods (Section 3) with an algorithm and a diagram
<ul> <li>In Appendix D., we describe related interaute.</li> <li>In Appendix D.1, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>In Appendix D.2, we perform ablations on our semantic likelihood (Section 3.3).</li> <li>In Appendix D.3, we demonstrate how to use our generative uncertainty for pixel-wise uncertainty.</li> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	707	explaining our generative uncertainty.
<ul> <li>In Appendix D., we point out imitations of our work.</li> <li>In Appendix D.1, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>In Appendix D.2, we perform ablations on our semantic likelihood (Section 3.3).</li> <li>In Appendix D.3, we demonstrate how to use our generative uncertainty for pixel-wise uncertainty.</li> <li>In Appendix D.5, we what diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.7, we provide implementation and experimental details.</li> </ul>	709	• In Appendix B, we describe related literature.
<ul> <li>In Appendix D.1, we qualitatively compare our method with BayesDiff (Kou et al., 2024).</li> <li>In Appendix D.2, we perform ablations on our semantic likelihood (Section 3.3).</li> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	710	• In Appendix C, we point out limitations of our work.
<ul> <li>In Appendix D.2, we perform ablations on our semantic likelihood (Section 3.3).</li> <li>In Appendix D.3, we demonstrate how to use our generative uncertainty for pixel-wise uncertainty.</li> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	711	• In Appendix D.1, we qualitatively compare our method with BayesDiff (Kou et al., 2024).
<ul> <li>In Appendix D.3, we demonstrate how to use our generative uncertainty for pixel-wise uncertainty.</li> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	712	• In Appendix D.2, we perform ablations on our semantic likelihood (Section 3.3).
<ul> <li>In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.</li> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkiänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	713 714	• In Appendix D.3, we demonstrate how to use our generative uncertainty for pixel-wise uncertainty.
<ul> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores.</li> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	715 716	• In Appendix D.4, we show that diffusion's own likelihood is not useful for filtering out poor samples.
<ul> <li>In Appendix D.5, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	717 718	<ul> <li>In Appendix D.5, we further compare our generative uncertainty to realism (Kynkäänniemi et al. 2019) and rarity (Han et al. 2023) scores.</li> </ul>
<ul> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	719 720	<ul> <li>In Appendix D.6, we investigate the drop in sample diversity by looking at the average generative uncertainty per conditioning class.</li> </ul>
<ul> <li>In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al., 2023).</li> <li>In Appendix E, we provide implementation and experimental details.</li> </ul>	721	generative uncertainty per conditioning class.
<ul> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimental details.</li> <li>In Appendix E, we provide implementation and experimentation and experimental details.</li> <li>In Appendix E, we provide implementation and experimentation and experimentation</li></ul>	722	• In Appendix D.7, we apply our generative uncertainty to detect low-quality samples in a latent flow matching model (Dao et al. 2023)
<ul> <li>In Appendix E, we provide implementation allockperimentation allockperimentat</li></ul>	723	• In Appendix F, we provide implementation and experimental details
726         727         728         729         730         731         732         733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	725	• In Appendix E, we provide implementation and experimental details.
727         728         729         730         731         732         733         734         735         736         737         738         739         730         731         732         733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751          752	726	
728         729         730         731         732         733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         754	727	
729         730         731         732         733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	728	
730         731         732         733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	729	
731         732         733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	730	
733         734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753	732	
734         735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	733	
735         736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754         755	734	
736         737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	735	
737         738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754	736	
738         739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754         755	737	
739         740         741         742         743         744         745         746         747         748         749         750         751         752         753         754         755	738	
740         741         742         743         744         745         746         747         748         749         750         751         752         753         754         755	739	
742         743         744         745         746         747         748         749         750         751         752         753         754         755	740	
743         744         745         746         747         748         749         750         751         752         753         754         755	742	
744         745         746         747         748         749         750         751         752         753         754         755	743	
745         746         747         748         749         750         751         752         753         754         755	744	
746         747         748         749         750         751         752         753         754         755	745	
747         748         749         750         751         752         753         754         755	746	
748         749         750         751         752         753         754         755	747	
750 751 752 753 754 755	748	
751 752 753 754 755	750	
752 753 754 755	751	
753 754 755	752	
754 755	753	
755	754	
	755	



Figure 4: Demonstration of how we compute *generative uncertainty* for each random noise z. We sample M sets of model parameters from our posterior distribution  $q(\theta|D)$  and generate M im-ages. Then, we evaluate the semantic likelihood for each by computing feature embeddings with a pretrained encoder  $c_{\phi}$  (e.g., CLIP) and take the uncertainty (e.g. entropy) over these embeddings. Random noises z with low uncertainty (*left*) tend to lead to consistent, high quality generations while random noises with high uncertainty (right) lead to poor, discordant generations.

778	Algorithm 1: Diffusion Sampling with Generative Uncertainty
779	<b>Input</b> : random noise $z$ , pretrained diffusion model $g_{\hat{\theta}}$ , Laplace posterior $q(\theta \mathcal{D})$ (Eq. 5),
780	number of MC samples M, semantic feature extractor $c_{\phi}$ , semantic likelihood noise $\sigma$
781	<b>Output:</b> generated sample $\hat{x}_0$ , generative uncertainty estimate $u(z)$
782	1 Generate a sample $\hat{x}_0 = g_{\hat{\theta}}(z)$
783	<sup>2</sup> Get semantic features $\boldsymbol{e}_0 = c_{\phi}(\hat{\boldsymbol{x}}_0)$
784	$\mathfrak{s}$ for $m=1  ightarrow M$ do
785	4 $\mid  heta_m \sim q( heta \mid \mathcal{D})$
786	5 $\hat{m{x}}_m = g_{ heta_m}(m{z})$
787	6 $ig  egin{array}{c} m{e}_m = c_\phi(\hat{m{x}}_m) \end{array}$
788	7 end
789	8 Compute $p(\boldsymbol{x} \boldsymbol{z}, \mathcal{D})$ using $\{\boldsymbol{e}_m\}_{m=0}^M$ (Eq. 7)
790	• Compute the entropy $u(\mathbf{z}) = H(p(\mathbf{x} \mathbf{z}, \mathcal{D}))$
701	10 return $\hat{x}_0, u(z)$
791	
192	

#### В **RELATED WORK**

**Uncertainty quantification in diffusion** models has recently gained significant attention. Most related to our work are BayesDiff (Kou et al., 2024), which uses a Laplace approximation to track epistemic uncertainty throughout the sampling process, and De Vita & Belagiannis (2025), which captures aleatoric uncertainty via the sensitivity of diffusion score estimates. Our work extends both by proposing a more general (applicable beyond diffusion), simpler approach (requiring no sampling modifications), and a more effective (see Section 4.1) uncertainty framework. 

Also related is DECU (Berry et al., 2024), which employs an efficient variant of deep ensembles (Lakshminarayanan et al., 2017) to capture the epistemic uncertainty of conditional diffusion mod-els. However, DECU does not consider using uncertainty to detect poor-quality generations, as its framework provides uncertainty estimates at the level of the conditioning variable, whereas ours estimates uncertainty at the level of initial random noise. Similarly, in Chan et al. (2024) the use of hyper-ensembles is proposed to capture epistemic uncertainty in diffusion models for inverse problems such as super-resolution, but, as in DECU, their approach does not provide uncertainty es-timates in unconditional settings or in conditional settings with low-dimensional conditioning (such as class-conditional generation). Moreover, both DECU (Berry et al., 2024) and Chan et al. (2024)

require modifying and retraining diffusion model components, whereas our approach operates *post-hoc* with any pretrained diffusion model via the Laplace approximation (Daxberger et al., 2021a).
A recent approach, PUNC (Franchi et al., 2024), focuses specifically on text-to-image models. The uncertainty of image generation with respect to text conditioning is measured through the alignment between a caption generated from a generated image and the original prompt used to generate said image.

Additionally, a large body of work explores conformal prediction for uncertainty quantification in diffusion models (Angelopoulos et al., 2022; Sankaranarayanan et al., 2022; Teneggi et al., 2023;
Belhasin et al., 2023). However, these approaches are primarily designed for inverse problems (e.g., deblurring), and cannot be directly applied to detect low-quality samples in unconditional generation.

821 Bayesian inference in generative models has been explored previously outside the domain of dif-822 fusion models. Prominent examples include Saatci & Wilson (2017) where a Bayesian version 823 of a GAN is proposed, showing improvements for semi-supervised learning, and Daxberger & 824 Hernández-Lobato (2019), where a Bayesian VAE (Tran et al., 2023) is shown to provide more 825 informative likelihood estimates for the unsupervised out-of-distribution detection compared to the 826 non-Bayesian counterparts (Nalisnick et al., 2018). Since diffusion models can be interpreted as neu-827 ral ODEs (Song et al., 2020b), another relevant work is Ott et al. (2023), which employs a Laplace approximation to quantify uncertainty when solving neural ODEs (Chen et al., 2018). However, Ott 828 et al. (2023) focuses solely on low-dimensional regression problems. 829

830 Non-uncertainty based approaches for filtering out poor generations include the realism 831 (Kynkäänniemi et al., 2019), rarity (Han et al., 2023), and anomaly scores (Hwang et al., 2024). 832 Our work is the first to establish a connection between these scores and uncertainty-based meth-833 ods, which we hope will inspire the development of even better sample-level metrics in the future. Additionally, a large body of work focuses on specially designed sample-quality scoring models 834 (Gu et al., 2020; Zhao et al., 2024) or, alternatively, on leveraging large pretrained vision-language 835 models (VLMs) (Zhang et al., 2024) for scoring generated images. However, these approaches re-836 quire either access to sample-quality labels or rely on (expensive) external VLMs. In contrast, our 837 uncertainty-based method requires neither, making it a more accessible and scalable alternative.

838 839

## C LIMITATIONS

840 841

842 While we have demonstrated in Section 4 that semantic likelihood is essential for addressing the 843 over-sensitivity of prior work to background pixels (Kou et al., 2024), our reliance on a pretrained 844 image encoder like CLIP (Radford et al., 2021) limits the applicability of our diffusion uncertainty framework to natural images. Removing the dependence on such encoders would unlock the ap-845 plication our Bayesian framework to other modalities where diffusion models are used, such as 846 molecules (Hoogeboom et al., 2022; Cornet et al., 2024) or text (Gong et al., 2022; Yi et al., 2024). 847 Exploring whether insights from the literature on uncovering semantic features in diffusion models 848 (Kwon et al., 2022; Luo et al., 2024; Namekata et al., 2024) could help achieve this represents a 849 promising direction for future work. 850

Moreover, the large size of modern diffusion models necessitates the use of cheap and scalable Bayesian approximate inference techniques, such as the (diagonal) last-layer Laplace approximation employed in our work (following (Kou et al., 2024)). A more comprehensive comparison of available approximate inference methods could be valuable, as improving the quality of the posterior approximation may further enhance the detection of low-quality samples based on Bayesian generative uncertainty.

857 858

D ADDITIONAL RESULTS

## D.1 QUALITATIVE COMPARISON WITH BAYESDIFF

860 861

859

To further highlight the differences between our generative uncertainty and BayesDiff (Kou et al., 2024), we present samples with the highest and lowest uncertainty according to BayesDiff in Figure 5. These samples are drawn from the same set of 12K ImageNet "unconditional" images generated

using the UViT model (Bao et al., 2023) as in Figure 1. Notably, BayesDiff's uncertainty score appears highly sensitive to background pixels—images with high uncertainty tend to have cluttered backgrounds, while those with low uncertainty typically feature clear backgrounds. Furthermore, as reflected in BayesDiff's poor performance in terms of FID and precision (see Table 1), some low-uncertainty examples exhibit noticeable artefacts, whereas certain high-uncertainty samples are of rather high-quality. For example, the image of a dog in the bottom-right corner of the high-uncertainty grid in Figure 5 looks quite good despite being assigned (very) high uncertainty.

Similarly, in Figure 6, we show low- and high-uncertainty samples according to BayesDiff for the
same set of 128 images per class as in Figure 2. Once again, we observe that BayesDiff's uncertainty
metric is less informative regarding a sample's visual quality compared to our generative uncertainty.









Figure 5: Images with the highest (*left*) and the lowest (*right*) BayesDiff uncertainty (Kou et al., 2024) among 12K generations using a UViT diffusion model (Bao et al., 2023). BayesDiff uncertainty correlates poorly with visual quality and is overly sensitive to the background pixels. Same set of 12K generated images is used as in Figure 1 to ensure a fair comparison.



Figure 6: Images with the highest (*bottom*) and the lowest (*top*) BayesDiff uncertainty (Kou et al., 2024) among 128 generations using a UViT diffusion model (Bao et al., 2023) for 2 classes: black swan (*left*) and Tibetan terrier (*right*). Same set of 128 generated images per class is used as in Figure 2 to ensure a fair comparison.

## D.2 ABLATION ON SEMANTIC LIKELIHOOD

To highlight the importance of using a semantic likelihood (Section 3.3) when leveraging uncertainty to detect low-quality generations, we conduct an ablation study in which we replace it with a standard Gaussian likelihood applied directly in pixel space (Eq. 6). Figure 7 presents the highest and lowest uncertainty images according to this 'pixel-space' generative uncertainty. Notably, pixelspace uncertainty is overly sensitive to background pixels, mirroring the issue observed in BayesDiff (Kou et al., 2024) (see Appendix D.1). This highlights the necessity of using semantic likelihood to obtain uncertainty estimates that are truly informative about the visual quality of generated samples.



Figure 7: Images with the highest (*left*) and the lowest (*right*) 'pixel-space' generative uncertainty among 12K generations using a UViT diffusion model (Bao et al., 2023). Pixel-space uncertainty correlates poorly with visual quality and is overly sensitive to the background pixels. Same set of 12K generated images is used as in Figure 1 to ensure a fair comparison.



Figure 8: Pixel-wise uncertainty based on our generative uncertainty for 5 generated samples using UViT (Bao et al., 2023).

## D.3 PIXEL-WISE UNCERTAINTY

While not the primary focus of our work, we demonstrate how our generative uncertainty framework (Algorithm 1) can be adapted to obtain pixel-wise uncertainty estimates. This is achieved by
replacing our proposed semantic likelihood (Eq. 7) with a standard 'pixel-space' likelihood (Eq. 6).
Figure 8 illustrates pixel-wise uncertainty estimates for 5 generated samples.

Although pixel-wise uncertainty received significant attention in past work (Kou et al., 2024; Chan et al., 2024; De Vita & Belagiannis, 2025), there is currently no principled method for evaluating its quality. Most existing approaches rely on qualitative inspection, visualizing pixel-wise uncertainty for a few generated samples (as we do in Figure 8). This further motivates our focus on sample-wise uncertainty estimates, where more rigorous evaluation frameworks—such as improvements in FID and precision on a set of filtered images (see Table 1)—enable more meaningful comparisons between different approaches.

#### 968 D.4 COMPARISON WITH LIKELIHOOD

970 We compare our generative uncertainty filtering criterion with a likelihood selection approach on 971 the 12K images generated by ADM trained on ImageNet 128x128. In the same way as in our other comparisons, we retain the 10K generated images with highest likelihood. We utilize the implementation in Dhariwal & Nichol (2021) to compute the bits-per-dimension of each sample (one-to-one with likelihood). The 25 samples with lowest and highest likelihood are shown in Figure
Visually, the likelihood objective heavily prefers simple images with clean backgrounds and not necessarily image quality. Note that this is consistent with other works that have reported likelihood to be an inconsistent identifier of image quality (Theis et al., 2016).

Figure 9: The 25 "Worst" (*left*) and "Best" (*right*) samples generated by ADM trained on ImageNet 128x128 selected by lowest and highest likelihood among 12K generations.

## D.5 COMPARISON WITH REALISM & RARITY

To better understand the relationship between our generative uncertainty and non-uncertainty-based approaches such as realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) scores, we compute the Spearman correlation coefficient between different sample-level metrics on a set of 12K generated images from the experiment in Section 4.1. As shown in Figure 10, realism and rarity scores exhibit a strong correlation (< -0.8). This is unsurprising, as both scores are derived from the distance of a generated sample to a data manifold obtained using a reference dataset (e.g., a subset of training data or a separate validation dataset).<sup>4</sup>

In contrast, our generative uncertainty exhibits a weaker correlation (< 0.4) with both realism and rarity scores. We attribute this to the fact that our uncertainty primarily reflects the limited training data used in training diffusion models (i.e., epistemic uncertainty, see Section 3.4), rather than the distance to a reference dataset, as is the case for realism and rarity scores.

Next, we investigate whether combining different scores can improve the detection of low-quality generations. When combining two scores, we first rank the 12K images based on each score individually, then compute the combined ranking by summing the two rankings and re-ranking accordingly. The results, shown in Table 2, indicate that combining realism and rarity leads to minor or no improvements in FID (9.81 compared to 9.76 for realism alone on ADM (Dhariwal & Nichol, 2021)). However, combining our generative uncertainty with either realism or rarity achieves the best FID performance (9.54 on ADM). These results suggest that ensembling scores that capture different aspects of generated sample quality is a promising direction for future research.

1018 1019

995

996 997 998

999

1020

D.6 CLASS-AVERAGED GENERATIVE UNCERTAINTY

1021To better understand the drop in sample diversity (recall) when using our generative uncertainty to<br/>filter low-quality samples in Table 1, we analyze the distribution of average entropy per conditioning<br/>class. Specifically, for each of the 12K generated images, we randomly sample a conditioning class

<sup>1024</sup> 

<sup>&</sup>lt;sup>4</sup>Such distance-based approaches are also commonly used to estimate prediction's quality in predictive models; see, for example, Van Amersfoort et al. (2020).



Figure 10: Spearman correlation coefficient between different sample quality metrics for 12K ImageNet images generated using ADM (Dhariwal & Nichol, 2021) (*left*) and UViT (Bao et al., 2023) (*right*).

Table 2: Image generation results for 10K filtered samples (out of 12K) based on combined metrics.
Combining our generative uncertainty outperforms combining realism and recall in terms of FID.
We report mean values along with standard deviation over 3 runs with different random seeds.

	ADM (D	DIM), ImageNe	et 128×128	UViT (D	PM), ImageNet	256×256
	<b>FID</b> $(\downarrow)$	<b>Precision</b> $(\uparrow)$	Recall $(\uparrow)$	<b>FID</b> $(\downarrow)$	<b>Precision</b> $(\uparrow)$	Recall $(\uparrow)$
Realism + Rarity	$9.81\pm0.06$	$67.06 \pm 0.29$	$66.73 \pm 0.37$	$8.26\pm0.07$	$69.01 \pm 0.33$	$69.86 \pm 0.36$
Ours + Realism	$9.54\pm0.04$	$66.41 \pm 0.15$	$67.04 \pm 0.47$	$7.60\pm0.10$	$68.33 \pm 0.09$	$69.75 \pm 0.42$
<b>Ours + Rarity</b>	$9.56\pm0.06$	$65.44 \pm 0.26$	$67.36 \pm 0.54$	$7.56\pm0.12$	$67.48 \pm 0.18$	$70.18 \pm 0.40$

1040

1044

to mimic unconditional generation. As a result, all 1,000 ImageNet classes are represented among
the 12K generated samples. Next, we compute our generative uncertainty (i.e., entropy; see Eq. 8)
for each sample and then average the uncertainties within each class. A plot of class-averaged uncertainties is shown in Figure 11. Since class-averaged uncertainties exhibit considerable variance,
the class distribution in the 10K filtered samples deviates somewhat from that of the original 12K images, thereby explaining the reduction in diversity (recall).

While our primary focus in this work is on providing per-sample uncertainty estimates u(z), we can also obtain uncertainty estimates for the conditioning variable u(y) (e.g., a class label), by averaging over all samples corresponding to a particular  $y \in \mathcal{Y}$  as done in Figure 11. These estimates resemble the epistemic uncertainty scores proposed in DECU (Berry et al., 2024) and could be used to identify conditioning variables for which generated samples are likely to be of poor quality. We leave further exploration of generative uncertainty at the level of conditioning variables for future work.

1065

#### 1066 D.7 FLOW MATCHING

1067 To demonstrate that our generative uncertainty 1068 framework (Section 3) extends beyond diffu-1069 sion models, we apply it here to the recently 1070 popularized flow matching approach (Lipman 1071 et al., 2022; Liu et al., 2022; Albergo et al., 1072 2023). Specifically, we consider a latent flow 1073 matching formulation (Dao et al., 2023) with 1074 a DiT backbone (Peebles & Xie, 2023). For 1075 sampling, we employ a fifth-order Runge-Kutta 1076 ODE solver (dopri5). In Figure 12, we illustrate the samples with the highest and low-1077 est generative uncertainty among 12K gener-1078 ated samples. On a filtered set of 10K images, 1079 our generative uncertainty framework achieves



Figure 11: A histogram of class-averaged generative uncertainties for 12K generated samples using UViT (Bao et al., 2023).

1080 an FID of 10.48 and a precision of 64.71, significantly outperforming a random baseline, which yields an FID of 11.80 and a precision of 61.04.



Figure 12: Images with the highest (*left*) and the lowest (*right*) generative uncertainty (Eq. 8) among 12K generations using a latent flow matching model (Dao et al., 2023). Generative uncertainty correlates with visual quality, as high-uncertainty samples exhibit numerous artefacts, whereas low-1102 uncertainty samples resemble canonical images of their respective conditioning class.

1103 1104

1100

1101

E IMPLEMENTATION DETAILS

1105 1106

1110

All our experiments can be conducted on a single A100 GPU, including the fitting of the Laplace 1107 posterior (Section 3.2). Code for reproducing our experiments is publicly available at GITHUB 1108 REPO. 1109

Laplace Approximation When fit-1111 ting a last-layer Laplace approxima-1112 tion (Section 3.2), we closely fol-1113 low the implementation from Bayes-1114 Diff (Kou et al., 2024). Specifically, 1115 we use the empirical Fisher approxi-1116 mation with a diagonal factorization 1117 for Hessian computation. The prior precision parameter and observation 1118 noise are fixed at  $\gamma = 1$  and  $\sigma = 1$ , 1119 respectively. For Hessian computa-1120 tion, we utilize 1% of the training 1121 data for ImageNet 128×128 and 2% 1122

	All Params.	LL Params.	LL Name
ADM	$\sim 421 \times 10^6$	$\sim 14\times 10^3$	out.2
UViT	$\sim 500 \times 10^6$	$\sim 18 \times 10^3$	decoder_pred
DiT	$\sim 131 \times 10^6$	$\sim 1.2 \times 10^6$	final_layer

Table 3: Details of our last-layer (LL) Laplace approximation. The first column presents the total number of model parameters, while the second and third columns indicate the number of parameters in the last layer and its name, respectively

for ImageNet 256×256. Further details about the last layer of each diffusion model are provided in 1123 Table 3, where we observe that fewer than 1% of the parameters receive a 'Bayesian treatment'. For 1124 code implementation, we rely on the laplace<sup>5</sup> library (Daxberger et al., 2021a). 1125

As discussed in Section C, improving the quality of the Laplace approximation—such as incorpo-1126 rating both first and last layers instead of only the last layer (Daxberger et al., 2021b; Sharma et al., 1127 2023) or optimizing Laplace hyperparameters (e.g., prior precision and observation noise) (Immer 1128 et al., 2021)—could further enhance the quality of generative uncertainty and represents a promising 1129 direction for future work. 1130

1131 **Sampling with Generative Uncertainty** For our main experiment in Section 4.1, we generate 1132 12K images using the pretrained ADM model (Dhariwal & Nichol, 2021) for ImageNet 128×128 1133

<sup>&</sup>lt;sup>5</sup>https://github.com/aleximmer/Laplace

and the UViT model (Bao et al., 2023) for ImageNet 256×256. Following BayesDiff (Kou et al., 2024), we use a DDIM sampler (Song et al., 2020a) for the ADM model and a DPM-2 sampler (Lu et al., 2022) for the UViT model, both with T = 50 sampling steps.

To compute generative uncertainty (Algorithm 1), we first sample M = 5 sets of weights from the posterior  $q(\theta|\mathcal{D})$ . Then, for each of the initial 12K random seeds, we generate M additional samples. The same set of model weights  $\{\theta_m\}_{m=1}^M$  is used for all 12K samples for efficiency reasons. For semantic likelihood (Eq. 7), we use a pretrained CLIP encoder (Radford et al., 2021) and set the semantic noise to  $\sigma^2 = 0.001$ .

**Baselines** For all baselines, we use the original implementation provided by the respective papers, except for (De Vita & Belagiannis, 2025), which we reimplemented ourselves since we were unable to get their code to run. Moreover, we use the default settings (e.g., hyperparameters) recommended by the authors for all baselines. For realism (Kynkäänniemi et al., 2019) and rarity (Han et al., 2023) we use InceptionNet (Szegedy et al., 2016) as a feature extractor and a subset of 50K ImageNet training images as the reference dataset. For samples where the rarity score is undefined (i.e., those that lie outside the estimated data manifold), we set it to inf.