

---

# Local Coverage Governs Memorization in Diffusion Models

---

Anonymous Authors<sup>1</sup>

## Abstract

Diffusion models are known to memorize training data, but which samples are most likely to be memorized? While memorization is often treated as a global property, in practice diffusion models simultaneously generate both memorized and novel samples. In this work, we show that memorization is governed by *local data coverage*. Leveraging the connection between diffusion models and kernel density estimation (KDE), we derive a theoretical criterion that predicts whether a point is memorized or generalized based on the density of training data in its neighborhood and the overall sample complexity. In the high-dimensional limit, this leads to a sharp local transition: regions of low coverage are dominated by isolated training samples (memorization), while dense regions support interpolation and generalization. We validate these predictions empirically, showing that memorization increases with local sparsity and that diffusion models exhibit a coexistence of memorized and novel samples within the same model. Extending this framework to multi-class settings, we further show that classes with higher intra-class diversity (and thus lower local coverage) are more strongly memorized. Our results provide a unified, local view of memorization in diffusion models, explaining when and where memorization occurs in terms of data geometry.

## 1. Introduction

**How does memorization arise *locally* in diffusion models?** Diffusion models are powerful generative methods capable of learning complex high-dimensional data distributions from finite datasets. Yet when trained with limited data, they may reproduce training examples rather than

generate novel samples, a phenomenon commonly referred to as *memorization* (Somepalli et al., 2022; Carlini et al., 2023; Kadkhodaie et al., 2023). Empirically, memorization decreases as the number of training samples increases, suggesting a fundamental relationship between finite-sample effects and generative generalization.

Theoretical work has begun to analyze this phenomenon through the close connection between diffusion models and kernel density estimation (KDE) (Pham et al., 2024; Ambrogioni, 2023; Biroli & Mézard, 2024; Achilli et al., 2025a; Lucibello & Mézard, 2024). In this picture, each training sample contributes a local kernel, and generalisation arises from the superposition of many such kernels, see Figure 1a for a sketch. When kernels overlap strongly, the model assigns finite probability to the space between examples; when overlap is weak, generation can collapse onto individual training points. These analyses predict that a critical sample complexity, there is a phase transitions from memorization to generalization that is *global*: the whole generative model goes from memorizing to generating new samples at once.

Recent work has also gone beyond this global perspective. First, Achilli et al. (2026) showed that memorization may emerge progressively through the loss of manifold dimensions, leading to a form of *geometric memorization* in which some directions of variability are lost before exact copying occurs. Second, Garnier-Brun et al. (2026) showed a coexistence of memorization and generalization in early-stopped diffusion models. This suggests that memorization is richer than a single abrupt transition.

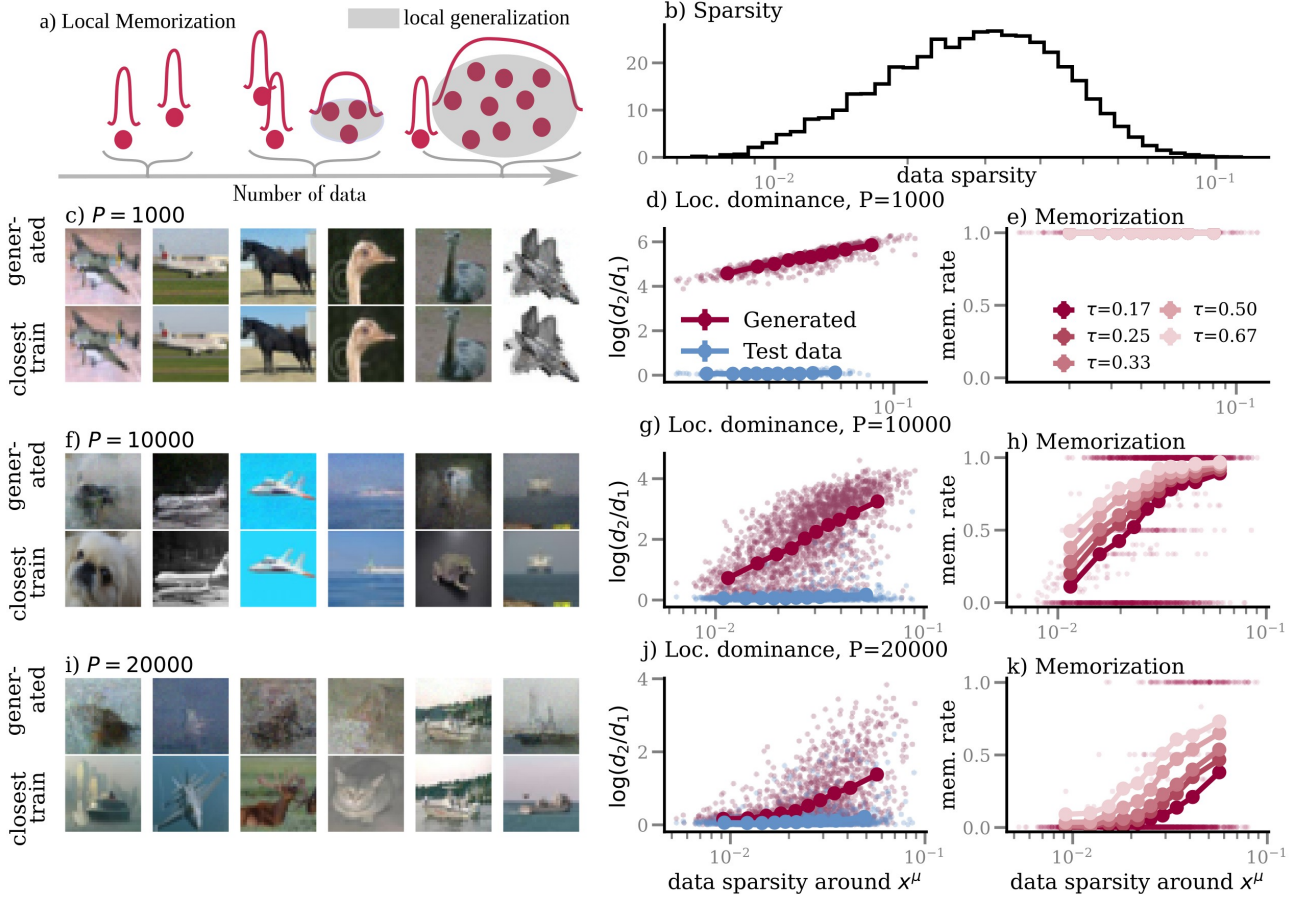
Additionally, empirical observations indicate an additional level of heterogeneity that remains theoretically unexplained. Memorizing diffusion models often reproduce only a subset of their training data rather than the entire dataset (Fang et al., 2025; Carlini et al., 2023). Memorization is also known to depend on guidance strength, conditioning, architecture, and dataset structure (Somepalli et al., 2023; Kim et al., 2025; Gu et al., 2025; Yoon et al., 2023). Thus, at finite sample size, a diffusion model often generates a mixture of copied and novel samples see Figure 1 c) f) and i). This raises a natural question:

**Which regions of data space, classes, or individual samples are most likely to be memorized?**

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.



**Figure 1. Local memorization.** a) Sketch of kernel density approximation and local memorization phenomenon. b) Distribution of data sparsity for 20,000 samples from CIFAR-10. c) Generated samples + closest training image (measured by cosine similarity) for diffusion model trained on  $P = 1000$  training examples. d) Local dominance as a function of sparsity. For each generated sample  $\hat{x}$ , we compute distances  $d_1$  and  $d_2$  to its nearest and second-nearest training samples, and define dominance via their ratio. Each point corresponds to a training sample  $x^\mu$ , aggregating generated samples assigned to it. Scatter shows raw values, markers indicate binned averages. Blue points show the same quantity computed using test data, providing a baseline without memorization. e) Memorization rate per training sample  $x^\mu$ , defined as the fraction of generated samples satisfying  $d_1/d_2 < \tau$ . Scatter shows raw values for  $\tau = 0.1$ , markers indicate binned averages across thresholds. f)–k) Same measurements for models trained on larger datasets. As the number of training samples increases, both dominance and memorization decrease, and their dependence on sparsity weakens. Overall, sparse regions exhibit strong single-sample dominance and higher memorization, while dense regions promote interpolation across multiple training points.

In this work, we give a concrete criterion for memorization by introducing a *local* theory of memorization in diffusion models. Our central hypothesis is simple: memorization is governed not only by the total number of training samples, but by their *local coverage*. Regions of data space containing many nearby examples support interpolation and generalization, whereas isolated regions are prone to sample retrieval. In Figure 1 we show an example of this behavior for diffusion models with U-net architecture (Ronneberger et al., 2015) trained on subsets of the CIFAR-10 image dataset (Krizhevsky, 2009). We find that more isolated training data are preferentially memorized, whereas data points in regions of higher local density (lower local sparsity in Figure 1) are memorized less.

To formalize this intuition, we extend the KDE framework (Lucibello & Mézard, 2024; Biroli & Mézard, 2024; Achilli et al., 2025b) to a local setting. We quantify the local coverage around a point  $x \in \mathbb{R}^N$  by the probability mass contained in a ball of radius  $h$  around  $x$ ,

$$p_{\text{in}}(x) = \int_{B_h(x)} d\rho(x'), \quad (1)$$

where  $B_h(x)$  denotes a  $N$ -dimensional hypersphere centered at  $x$  and  $\rho$  the density of the data from which the training points are drawn. This quantity measures how densely the data distribution  $\rho$  populates the neighborhood of  $x$ . In the high-dimensional regime in which the number of samples scales as

$$P = e^{\alpha N}, \quad (2)$$

we define

$$\nu_{\text{in}}(x) = \lim_{N \rightarrow \infty} \frac{1}{N} \ln p_{\text{in}}(x). \quad (3)$$

Our main result shows that memorization at a point  $x$  is controlled by the pair  $(\nu_{\text{in}}(x), \alpha)$  with a sharp transition:

$$\begin{cases} x \in A_{\text{mem}} & \text{if } \nu_{\text{in}}(x) < -\alpha \\ x \in A_{\text{gen}} & \text{otherwise} \end{cases} \quad (4)$$

where we refer to  $A_{\text{gen}}$  as the region where the model correctly interpolates, and  $A_{\text{mem}}$  is a region where the learned density is characterized by isolated peaks centered on training examples. This result shows that depending on local coverage, the same diffusion model may exhibit retrieval-like behavior in some regions of space and generative interpolation in others. Points in low-coverage regions ( $\nu_{\text{in}}(x)$  small) are memorized, while points in high-coverage regions are generalized. This local theory explains several empirical phenomena:

- **Coexistence of memorization and generalization:** copied and novel samples can arise simultaneously from the same model.
- **Gradual disappearance of memorization:** as  $P$  increases, memorization vanishes first in dense regions and persists longest in sparse regions.
- **Class-dependent memorization:** classes with larger intra-class diversity (and therefore lower local coverage) are memorized more strongly.

We validate these predictions on diffusion models trained on standard image datasets such as CIFAR-10 and CelebA (Liu et al., 2015). Consistent with theory, training examples originating from regions of lower density, as well as classes with larger intra-class diversity are memorized more strongly. More broadly, our results suggest that memorization in diffusion models is not a global binary property, but a local finite-sample phenomenon governed by the geometry and coverage of the training distribution.

## 2. Related Work

A growing body of empirical and theoretical work has established that diffusion models can memorize training data, particularly in low-data regimes. Memorization is not a uniform phenomenon: it depends strongly on properties of the training setup. In particular, it is exacerbated by data duplication (Carlini et al., 2023), prompt conditioning and classifier free guidance (Somepalli et al., 2023; Wen et al., 2024; Kim et al., 2025; Gu et al., 2025), model capacity (Yoon et al., 2023; George et al., 2025) and training duration (Bonnaire et al., 2025; Favero et al., 2025), (Garnier-Brun et al., 2026) whereas it decreases with learning rate

(Wu et al., 2025) and dataset size (Kadkhodaie et al., 2024; Somepalli et al., 2023). These findings suggest that memorization is a structured and predictable effect, rather than a rare failure mode.

Strategies to measure and mitigate memorization include determining the local intrinsic dimension of generated points (Ross et al., 2025), and estimating how "peaked" the estimated density is around a given point (Jeon et al., 2025) as well as identifying (Wen et al., 2024; Kim et al., 2025) memorized prompts. Importantly, these strategies share a common intuition: isolated data points produce a highly peaked local density, and produce samples of low intrinsic dimension. This perspective is naturally consistent with a kernel density estimation (KDE) view of diffusion models. Isolated training points induce highly concentrated local density estimates and generate low-dimensional samples, while dense regions support smoother interpolation. Similarly, conditioning can be interpreted as restricting the effective sample set contributing to the density estimate, thereby increasing local sparsity and promoting memorization. In this sense, existing empirical methods already implicitly rely on a local density perspective on memorization.

On the theoretical side, several works have established a connection between diffusion models and kernel density estimation (Pidstrigach, 2022; Ambrogioni, 2024; Li et al., 2024). Within this framework, one can determine a global "collapse phase" in diffusion models as a function of effective noise in the backward/sampling process (Biroli et al., 2024), where the model begins to reproduce training data. More generally, the sample complexity required for smooth interpolation is known to scale exponentially with the ambient or manifold dimension (Lucibello & Mézard, 2024; Achilli et al., 2025b) of the data. Recent work has further shown that memorization can emerge gradually through a loss of manifold dimensions (Achilli et al. (2026) and that this process may be spatially inhomogeneous. However, existing theoretical analyses remain largely global: they characterize when a model as a whole memorizes or generalizes, but do not predict which specific regions or samples are more likely to be memorized. In contrast, we develop a local theory of memorization that explicitly links memorization to data coverage at the level of individual points. This provides a principled explanation for the heterogeneous memorization patterns observed in practice.

## 3. Background: The equivalence between diffusion models and kernel density estimation

We briefly recall a key observation: diffusion models trained on finite data behave like KDE. Diffusion models define a forward noising process that gradually corrupts a data point

165  $x_0$  from  $\rho$  through the addition of noise

$$166 \quad x_t(\epsilon_t) = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, \text{Id}), \quad (5)$$

168 where  $\bar{\alpha}_t \in (0, 1)$  is a decreasing function in  $t$ . As  $t$  increases, the original signal  $x_0$  is gradually suppressed compared to the isotropic Gaussian noise, until one obtains  $x_T$  whose distribution is close to  $\mathcal{N}(0, \text{Id})$ . A neural network  $\epsilon_\theta$  with parameters  $\theta$  is trained to predict the noise  $\epsilon_t$ , which corresponds to learning the score  $\nabla \ln \rho_t(x_t)$ , where  $\rho_t$  is the distribution of  $x_t$ . This is achieved via the denoising score matching objective

$$177 \quad L = \frac{1}{NTP} \sum_t \sum_{x^\mu \in \mathcal{D}} \mathbb{E}_{\epsilon_t} \|\epsilon_t - \epsilon_\theta(x_t^\mu(\epsilon_t), t)\|^2, \quad (6)$$

178 where  $\mathcal{D}$  is a training set of size  $P$  points  $\{x^\mu\}_{\mu=1}^P$  drawn i.i.d. from  $\rho$ .

181 To make the connection to KDE explicit, consider an idealized setting in which the model has infinite capacity and is trained to optimality. In this case, for each  $t$ , the learned function  $\epsilon_\theta(\cdot, t)$  can be treated as an arbitrary function  $\epsilon(\cdot, t)$  and the loss can be minimized functionally. Setting this functional derivative to zero yields:

$$187 \quad \epsilon(x, t) \propto \nabla_x \ln \sum_\mu \mathcal{N}(x; \bar{\alpha}_t x^\mu, (1 - \bar{\alpha}_t)\text{Id}). \quad (7)$$

189 This identity reveals that, when trained to convergence on a finite number of samples, diffusion models learn a Gaussian mixture centered at the training data points. In other words, diffusion models trained on a finite dataset recover the score of the empirical distribution convolved with Gaussian noise, rather than the true underlying density  $\rho$ . In the low-noise limit ( $\bar{\alpha}_t \rightarrow 1$ ), this mixture becomes sharply peaked around individual training samples, and the model effectively reproduces them. At finite noise, there is a direct correspondence between Equation (7) and KDE. KDE attempts to approximate a density  $\rho$  using a mixture of kernel functions  $K$  centered at the data points:

$$202 \quad Z(x) = \frac{1}{P} \sum_\mu K\left(\frac{|x - x^\mu|}{h}\right). \quad (8)$$

204 The correspondence with KDE becomes explicit by identifying the kernel as Gaussian with bandwidth  $h = \sqrt{1 - \bar{\alpha}_t}$  up to a rescaling of the variables  $x, x^\mu$  by  $\sqrt{\bar{\alpha}_t}$ . In the ideal case, one has sufficient number of samples  $P$  such that  $Z$  becomes a smooth approximation of  $\rho$ . This behavior with  $P$  is one hypothesis to explain memorization and generalization in diffusion models.

211 This perspective suggests that generalization is not uniform across the data space. Even when the total number of samples is large, the quality of the KDE approximation depends on local sample density. Regions with high data density are well approximated, while sparse regions remain dominated by individual kernels. In the following, we formalize this intuition by deriving conditions under which local regions exhibit memorization or generalization.

## 4. Theory: Kernel density estimation with hard spheres

To study kernel density estimation (KDE), we will now characterize the behavior of the log-density

$$z(x) = \frac{1}{N} \ln Z(x), \quad (9)$$

where  $Z(x)$  is the mixture of kernel functions  $K$  centered at the data points defined in Equation (8). The difficulty in establishing the behavior of  $z$  is that it is a random quantity that depends on the draw of the  $P$  training points from  $\rho$ . Our goal is to compute its distribution over such draws to derive general properties of high dimensional KDE.

To this end, we use a method introduced by Dotsenko (2011), which allows us to compute the cumulative distribution function (CDF) of  $z(x)$  via an exponential transform. We define

$$W_N(y, x) = \frac{\overline{\exp(-\exp(-N(y - z(x))))}}{\xrightarrow{N \rightarrow \infty} \overline{\Theta(y - z(x))}} \quad (10)$$

where the  $\overline{f(z(x))}$  denotes the average over draws of the  $P$  data points from  $\rho$  and  $\Theta$  is the Heaviside function with convention  $\Theta(0) = e^{-1}$ . In the limit  $N \rightarrow \infty$ , this quantity converges to the cumulative distribution function of  $z(x)$ . If  $K$  is bounded from above (e.g. Gaussian  $K$ ), we can expand the outer exponential function in powers of  $Z^n(x)$  and exchange the order of averaging and summation. This reduces the problem to computing moments of the form  $\overline{Z^n(x)}$ , which can be evaluated exactly using the Fourier transforms of  $\rho$  and  $K$  (see Section A). We then obtain the final expression

$$W_N(y, x) = \exp\{e^{\alpha N} \ln[1 - I_\circ(y, x)]\} \quad (11)$$

where

$$I_\circ(y, x) = \int d\tilde{x} \rho(\tilde{x}) \left(1 - e^{-K\left[\frac{x - \tilde{x}}{h}\right]} \exp[-N(\alpha + y)]\right) \leq 1. \quad (12)$$

This expression eliminates the explicit dependence on the random dataset, reducing the problem to an integral over the data distribution  $\rho$ . Although this integral remains high-dimensional, its structure depends only on the distribution of kernel values  $K\left(\frac{x - \tilde{x}}{h}\right)$ . In particular, if the logarithm of the kernel concentrates, this integral simplifies significantly. Even without such assumptions, the distribution of kernel values can be estimated empirically from data. This formulation makes explicit that the behavior of  $z(x)$  is controlled by the statistics of distances between  $x$  and samples drawn from  $\rho$ , providing a direct link between KDE performance and local data geometry.

We now analyze the limit  $N \rightarrow \infty$ . Depending on the scaling of  $\lim_{N \rightarrow \infty} e^{\alpha N} I_\circ(y, x) =: i_\circ(y, x)$ , we obtain

three regimes:

$$\lim_{N \rightarrow \infty} W_N(y, x) = \begin{cases} 0 & \text{if } i_o(y, x) = \infty \\ e^{-i_o(y, x)} & \text{if } i_o(y, x) = \mathcal{O}(1) \\ 1 & \text{if } i_o(y, x) = 0 \end{cases} \quad (13)$$

As a function of  $y$ , we find that  $I_o$  decreases and thus  $W_N$  increases. A typical scenario is that the distribution of  $z$  concentrates with  $N \rightarrow \infty$  thus  $\Pr(\{y > z(x)\})$  becomes a step function, i.e. the measure of points  $y$  where the second condition in Equation (13) is fulfilled must shrink to zero. However, our results could in principle be used to compute fluctuations in  $z(x)$ , corresponding to a non-vanishing interval in  $y$  where the second condition is fulfilled.

**Hard spheres.** We now restrict ourselves to one particular kernel, namely

$$K(x) = \Theta(R^2 - x^2)/V_R \quad (14)$$

where  $V_R$  is the volume of the  $N$ -dimensional sphere with radius  $R$ . As we will see, this choice of kernel allows us to express Equation (13) in a particularly simple form.

In high dimensions, we expect that this choice of kernel is qualitatively equivalent to a Gaussian one with standard deviation  $R$ , as the mass of both distributions concentrates near a thin shell at radius  $R$  when  $N \rightarrow \infty$ . We fix  $R = 2\pi e$ , for which the volume  $V_R$  remains  $\mathcal{O}(1)$  as  $N \rightarrow \infty$ . The effective scale of the kernel relative to the data is controlled by the bandwidth  $h$ , which can be interpreted as a rescaling of the data points by a factor  $1/h$ . Unless otherwise specified, we set  $h = 1$ .

Under these assumptions, the expression for  $I_o$  simplifies significantly. Assuming that  $\nu_{\text{in}}(x)$  remains  $\mathcal{O}(1)$  as  $N \rightarrow \infty$ , we obtain the following sharp characterization of  $z(x)$ :

$$z(x) \rightarrow \begin{cases} \nu_{\text{in}}(x) & \text{if } \nu_{\text{in}}(x) \geq -\alpha \\ -\infty & \text{if } \nu_{\text{in}}(x) < -\alpha \end{cases}, \quad (15)$$

see Section A for a derivation.

This result shows that the behavior of  $z(x)$  is controlled by the comparison between the local coverage  $\nu_{\text{in}}(x)$  and the sample complexity  $\alpha$ . Regions with sufficiently high local density behave as smooth KDE estimates, while regions with low density are dominated by the absence of nearby samples.

The divergence  $z(x) \rightarrow -\infty$  for  $\nu_{\text{in}}(x) < -\alpha$  reflects the fact that, in these regions, the kernel density estimate vanishes with high probability. Indeed,

$$\begin{aligned} \Pr(\{Z(x) = 0\}) &= (1 - p_{\text{in}}(x))^P = e^{P \ln(1 - e^{-N\nu_{\text{in}}(x)})} \\ &\rightarrow \Theta(-[\alpha + \nu_{\text{in}}(x)]) \end{aligned}$$

On the other hand, when  $\alpha$  is large enough, then  $Z(x)$  approaches  $p_{\text{in}}(x)$ , hence the kernel density estimate converges to a local average of  $\rho$  over a sphere.

This establishes a sharp dichotomy between two regimes. In regions where  $\nu_{\text{in}}(x) \geq -\alpha$ , the kernel density estimate provides a smooth approximation of the underlying distribution, corresponding to a generalization regime. In contrast, in regions where  $\nu_{\text{in}}(x) < -\alpha$ , the estimate vanishes with high probability, indicating that no nearby samples are present.

In the context of diffusion models, such regions are dominated by isolated training points. If probability mass is assigned to these locations, it must concentrate on individual samples, leading to memorization. This establishes the coexistence of memorized and generalized regions, as described in Eq. (4).

## 5. Experiments

### 5.1. Local sparsity and dominance

We now empirically investigate how local data density influences memorization behavior in diffusion models. Our goal is to test the central prediction of the KDE-based theory that sparse regions are more prone to memorization than dense ones. This implies a **coexistence of memorization and generalization**: copied and novel samples can arise simultaneously from the same model. Furthermore, it predicts a **gradual, rather than abrupt disappearance of memorization**: as  $P$  increases, memorization vanishes first in dense regions and persists longest in sparse regions. Consequently, regions of low data density lead to more localized, sample-specific behavior, while dense regions promote interpolation across multiple training points. While our theory predicts a sharp transition as a function of local coverage, here we test its qualitative consequences using empirical proxies for density.

We train diffusion models with a U-net architecture on randomly chosen subsets of CIFAR-10 and CelebA. We then generate samples from these models and assess their closeness to training examples using different measures outlined below. Details on the experimental procedure are given in Section C. We first define how we measure sparsity and memorization and then describe the outcome of the experiments shown in Figure 1 and Figure 3.

**Local sparsity.** For each training sample  $x^\mu$ , we quantify local sparsity using nearest-neighbor distances, which serve as a proxy for inverse local coverage. Concretely, we define

$$s(x^\mu) = \frac{1}{Nk} \sum_{i=1}^k \|x^\mu - x^{\mu_i}\|^2, \quad (16)$$

where  $\{x^{\mu_i}\}$  are the  $k$  nearest neighbors of  $x^\mu$ . Larger values of  $s(x^\mu)$  correspond to sparser regions of the data distribution. We also validate that the results remain consis-

275 tent for different choices of  $k$  in  $\{5, 10, 20, 50\}$ . The results  
276 reported in the figures correspond to  $k = 10$ .

277 **Local dominance.** Given a generated sample  $\hat{x}$ , we compute  
278 its distances  $d_1(\hat{x})$  and  $d_2(\hat{x})$  to the nearest and second-  
279 nearest training samples, and define the *dominance*  
280

$$281 \Delta(\hat{x}) = \log \frac{d_2(\hat{x})}{d_1(\hat{x})}. \quad (17)$$

282 Large values of  $\Delta$  indicate that a single training point domi-  
283 nates the local score estimation around the generated sample,  
284 while  $\Delta \approx 0$  corresponds to interpolation between multiple  
285 neighbors. We assign each generated sample to its near-  
286 est training point and compute, for each  $x^\mu$ , the average  
287 dominance over all generated samples assigned to it. To  
288 disentangle model-specific effects from dataset geometry,  
289 we construct a baseline by replacing generated samples with  
290 held-out test data, processed in the same way.

291 **Memorization rate.** We further report a binary memo-  
292 rization metric used in several previous studies (Wu et al.,  
293 2025; Yoon et al., 2023; Bonnaire et al., 2025) based on the  
294 condition  $d_1/d_2 < \tau$ .

295 **Results.** The middle column of Figure 1 shows the average  
296 dominance as a function of local sparsity for different train-  
297 ing set sizes. We observe a clear **monotonic increase of**  
298 **dominance with sparsity** for generated samples: sparse re-  
299 gions exhibit strong single-sample dominance, while dense  
300 regions show low dominance and interpolation. In contrast,  
301 the test-data baseline remains nearly flat, indicating that  
302 this effect is not explained by nearest-neighbor geometry  
303 alone. As the number of training samples increases, the  
304 overall level of dominance decreases and the dependence  
305 on sparsity weakens, consistent with the expectation that  
306 higher sample density reduces isolated regions. The cor-  
307 responding results for the memorization rate are shown in  
308 the right column of Figure 1. While we do not directly esti-  
309 mate the threshold  $\nu_{\text{in}}(x) = -\alpha$ , the observed dependence  
310 on local sparsity is consistent with a transition controlled  
311 by local coverage. For the canonical value  $\tau = 0.33$  com-  
312 monly used in the literature, we see a clear **coexistence of**  
313 **memorization and generalization** that depends on data  
314 sparsity. Moreover, while the precise behavior depends on  
315  $\tau$ , the same qualitative trend is observed: **memorization**  
316 **increases with sparsity**. In the appendix Section B we  
317 show the outcome of the same experiment on downsampled  
318 CelebA, an image dataset of celebrity portraits. The results  
319 of these experiments are also consistent with the hypothesis.  
320 These results support a local version of the KDE picture:  
321 diffusion models exhibit a continuous transition from inter-  
322 polation in dense regions to single-sample dominance in  
323 sparse regions.  
324  
325  
326  
327  
328  
329

## 5.2. Class-dependent Memorization

We now study memorization at a coarser level, focusing  
on **class-dependent effects**. In a multi-class setting, the  
KDE picture predicts that classes with more concentrated  
support (i.e. lower local sparsity) should be memorized less  
than classes with higher intra-class variability. A schematic  
illustration is shown in Figure 2a. To formalize this intuition,  
consider a mixture of class-conditional densities  $\rho_c$ , such  
that

$$\rho(x) = \sum_c w_c \rho_c(x),$$

where the weights  $w_c$  sum to one, and we assume that  
all weights are order one in  $N$ . We define the local in-  
distribution log-density for class  $c$  as

$$\nu_{\text{in},c}(x) := \frac{1}{N} \ln p_{\text{in},c}(x) = \frac{1}{N} \ln \int_{B_h(x)} d\rho_c(x'), \quad (18)$$

and assume that  $\nu_{\text{in},c}(x)$  remains  $\mathcal{O}(1)$  in the large- $N$  limit.  
Then the total in-distribution density satisfies

$$\begin{aligned} \nu_{\text{in}}(x) &= \frac{1}{N} \ln \sum_c e^{N(\nu_{\text{in},c}(x) + N^{-1} \ln w_c)} \\ &\rightarrow \max_c \nu_{\text{in},c}(x), \end{aligned}$$

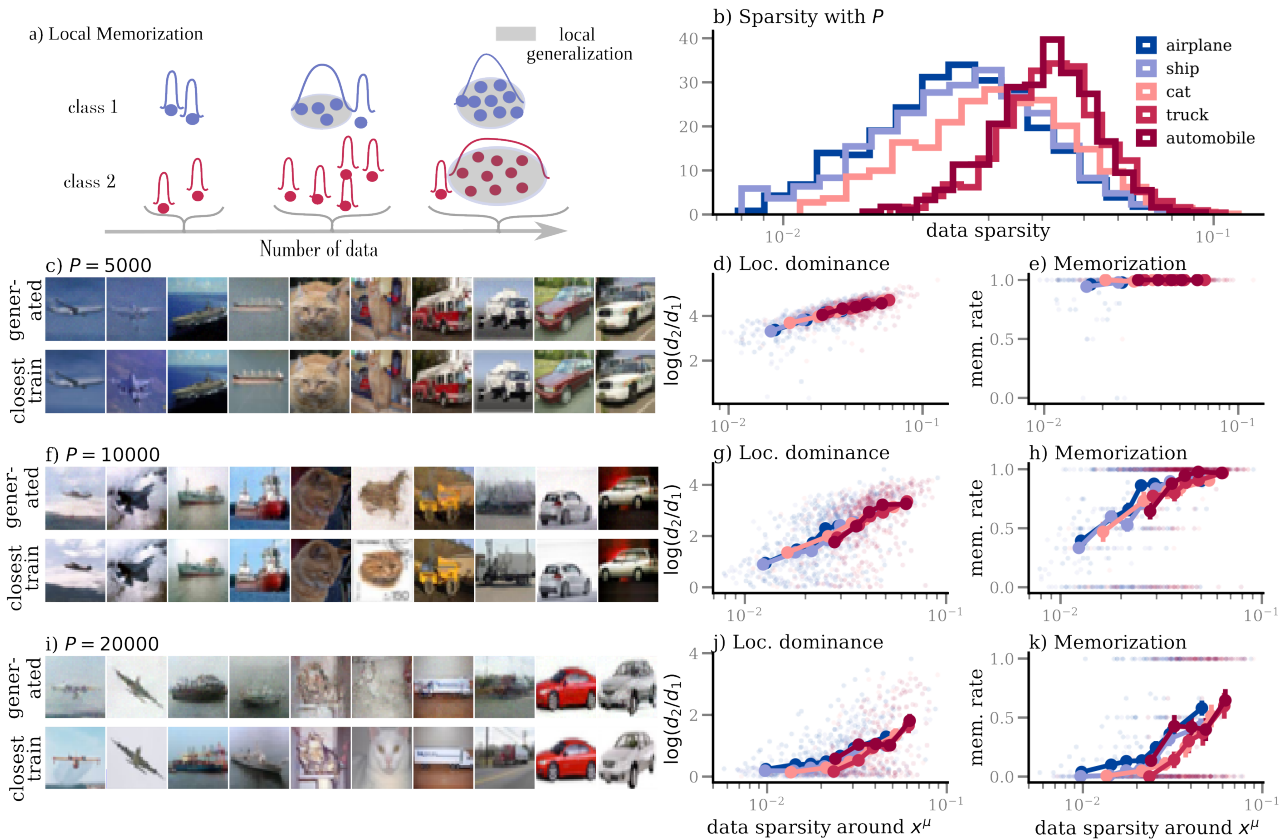
i.e. it is dominated by the locally densest class. Conse-  
quently, we obtain the same scaling behavior as in Equa-  
tion (15):

$$z(x) \rightarrow \begin{cases} \max_c \nu_{\text{in},c}(x) & \text{if } \max_c \nu_{\text{in},c}(x) \geq -\alpha, \\ -\infty & \text{otherwise,} \end{cases} \quad (19)$$

showing that memorization is governed by the class with  
the highest local coverage around  $x$ . This leads to a clear  
prediction: classes with higher intra-class variability (and  
thus lower local coverage) should be more prone to memo-  
rization.

We now demonstrate empirically that such class-dependent  
density differences naturally arise. As a concrete exam-  
ple, we isolate the memorization behavior of different sub-  
classes of CIFAR-10. In Figure 2b we show that the classes  
"airplane" and "ship" in CIFAR-10 have, on average, con-  
siderably lower intra-class diversity than the classes "truck"  
and "automobile", likely due to more homogeneous blue  
background colors in the former two classes. Correspond-  
ingly, we observe both higher local dominance, and higher  
memorization in more diverse classes.

In the appendix, Section B, we show that the same trends  
emerge for two additional experiments. First, when one  
sorts the CelebA datasets into classes according to whether  
images have the attribute "wearing hat" or "blond", im-  
ages where "wearing hat" is true are more diverse than  
those where it is false. Correspondingly, portraits featur-  
ing hats are memorized more. Second, we train diffusion



**Figure 2. Class-wise memorization.** a) Schematic illustration of KDE for two classes with different local sparsities. b) Distribution of per-class sparsity around training points for different subclasses of CIFAR-10, measured from 20,000 training samples. c) Generated samples together with their closest training example (measured by cosine similarity) for a diffusion model trained on  $P = 1000$  samples, sorted to contain 2 nearest neighbors for each of the four classes shown in b). d) Local dominance as a function of sparsity. For each generated sample  $\hat{x}$ , we compute distances  $d_1$  and  $d_2$  to its nearest and second-nearest training samples, and define dominance via their ratio. Each point corresponds to a training sample  $x^\mu$ , aggregating generated samples assigned to it. Scatter shows raw values, markers indicate binned averages. Different colors correspond to different classes. e) Memorization rate of each sample  $x^\mu$ , defined as the fraction of generated samples satisfying  $d_1/d_2 < \tau = 1/3$ . f) - k) report the same measures, but for diffusion models trained on larger datasets. "Airplane" and "ship" samples lie in denser regions and exhibit lower dominance and memorization, while "truck" and "automobile" samples are sparser and more frequently memorized, consistent with the prediction that local coverage controls class-dependent memorization.

models on a mixture of MNIST (Deng, 2012) and CIFAR-10 images. MNIST, consisting of handwritten digits, is structurally simpler and occupies a much more concentrated region of image space, whereas CIFAR-10 exhibits substantially higher variability. Treating MNIST and CIFAR-10 as two distinct classes, we observe that MNIST samples consistently exhibit lower local sparsity than CIFAR-10 samples. As predicted by the theory, this translates into reduced memorization. These results support the hypothesis that classes with higher intra-class diversity are memorized more strongly than classes with lower intra-class diversity.

## 6. Discussion

**Summary.** In this work, we test the predictive power of kernel density estimation to explain local memorization in diffusion models. Kernel density estimation in high dimensions exhibits a rich phenomenology that has not been fully explored. While earlier works (Lucibello & Mézard, 2024; Biroli & Mézard, 2024; Achilli et al., 2025b) focused on establishing *global* phase transitions, focusing on memorization or generalization of *typical samples from the data distribution*, its extension to anisotropic distributions revealed a more refined picture, that can e.g. be associated with gradual loss of manifold dimensions (Achilli et al., 2026). In this work, we explore yet another crucial property of KDE: memorization is governed by *local coverage*. We show that, in the high-dimensional limit, each point  $x$  undergoes a local transition as a function of sample complexity, separating regions where the density is smoothly estimated from regions dominated by isolated training samples. This leads to a coexistence of *generalization* and *memorization* within the same model. Our analysis further predicts that, in multi-class settings, the behavior at a point  $x$  is controlled by the class with the highest local coverage. Empirically, we find consistent evidence for these predictions: memorization increases with local sparsity and intra-class variability, and different classes exhibit systematically different memorization behavior.

Taken together, these results show that memorization in diffusion models is not a global failure mode, but a *localized phenomenon* driven by the geometry of the data distribution. In particular, isolated regions of the data space can remain memorized even when the majority of generated samples are novel.

**Limitations.** Our analysis relies on the assumption that  $\frac{1}{N} \ln p_{\text{in}}(x)$  satisfies a large deviation principle, analogous to concentration assumptions used in prior work (Lucibello & Mézard, 2024; Achilli et al., 2025a). While this enables a sharp theoretical characterization, it predicts an abrupt transition between memorization and generalization, whereas empirically we observe a more gradual dependence on local sparsity. We therefore hypothesize that fluctuations are

at the core of this more gradual increase, which could be taken into account through Equation (13). Incorporating such effects is an important direction for future work.

A second limitation stems from the choice of the metric space underlying the KDE approximation. Our theory assumes that diffusion models operate in the ambient space. In practice, however, models may implicitly operate in a lower-dimensional representation, for example through projection onto a data manifold (Achilli et al., 2025a; 2026) or by exploiting symmetries in the data (Kamb & Ganguli, 2025). In such cases, the effective kernel should be understood as acting in this learned metric space, which depends on model architecture and may help explain the observed dependence of memorization on model capacity (Yoon et al., 2023).

**Outlook.** Our results suggest that controlling memorization requires shaping the *local geometry* of the data representation. In particular, learning representations that increase local coverage, by mapping data to lower-dimensional spaces or by exploiting invariance may reduce memorization but reduces model expressivity, presenting a tradeoff between increased local coverage and preserving generative performance.

Another important direction concerns training dynamics. Memorization is known to depend on training time and optimization hyperparameters (Bonnaire et al., 2025; Favero et al., 2025; Wu et al., 2025), and recent work suggests that memorization and generalization can emerge simultaneously during training (Garnier-Brun et al., 2026). Understanding how convergence to the KDE solution in the generalized and memorized regimes depends on optimization dynamics is a promising direction for future research.

**Broader impacts.** This work is foundational, but memorization in generative models is directly related to privacy and copyright risks. Understanding which examples are likely to be memorized may help audit and mitigate training-data reproduction, although the same insights could potentially inform extraction attempts.

## References

- Achilli, B., Ambrogioni, L., Lucibello, C., Mézard, M., and Ventura, E. The Capacity of Modern Hopfield Networks under the Data Manifold Hypothesis, March 2025a. URL <http://arxiv.org/abs/2503.09518>. arXiv:2503.09518 [cond-mat].
- Achilli, B., Ambrogioni, L., Lucibello, C., Mézard, M., and Ventura, E. Memorization and generalization in generative diffusion under the manifold hypothesis. *Journal of Statistical Mechanics: Theory and Experiment*, 2025(7):073401, July 2025b. ISSN 1742-5468. doi: 10.

- 1088/1742-5468/ade136. URL <https://doi.org/10.1088/1742-5468/ade136>.
- Achilli, B., Ventura, E., Silvestri, G., Pham, B., Raya, G., Krotov, D., Lucibello, C., and Ambrogioni, L. Losing dimensions: Geometric memorization in generative diffusion, March 2026. URL <http://arxiv.org/abs/2410.08727>. arXiv:2410.08727 [stat].
- Ambrogioni, L. In search of dispersed memories: Generative diffusion models are associative memory networks, November 2023. URL <http://arxiv.org/abs/2309.17290>. arXiv:2309.17290 [cs, stat].
- Ambrogioni, L. In Search of Dispersed Memories: Generative Diffusion Models Are Associative Memory Networks. *Entropy*, 26(5):381, May 2024. ISSN 1099-4300. doi: 10.3390/e26050381. URL <https://www.mdpi.com/1099-4300/26/5/381>.
- Biroli, G. and Mézard, M. Kernel Density Estimators in Large Dimensions, October 2024. URL <http://arxiv.org/abs/2408.05807>. arXiv:2408.05807 [cs].
- Biroli, G., Bonnaire, T., de Bortoli, V., and Mézard, M. Dynamical regimes of diffusion models. *Nature Communications*, 15(1):9957, November 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-54281-3. URL <https://www.nature.com/articles/s41467-024-54281-3>.
- Bonnaire, T., Urfin, R., Biroli, G., and Mézard, M. Why Diffusion Models Don't Memorize: The Role of Implicit Dynamical Regularization in Training, May 2025. URL <http://arxiv.org/abs/2505.17638>. arXiv:2505.17638 [cs].
- Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwal, V., Tramèr, F., Balle, B., Ippolito, D., and Wallace, E. Extracting Training Data from Diffusion Models, January 2023. URL <http://arxiv.org/abs/2301.13188>. arXiv:2301.13188 [cs].
- Deng, L. The MNIST Database of Handwritten Digit Images for Machine Learning Research [Best of the Web]. *IEEE Signal Processing Magazine*, 29(6):141–142, November 2012. ISSN 1558-0792. doi: 10.1109/MSP.2012.2211477. URL <https://ieeexplore.ieee.org/document/6296535>.
- Dotsenko, V. Replica solution of the random energy model. *EPL (Europhysics Letters)*, 95(5):50006, September 2011. ISSN 0295-5075, 1286-4854. doi: 10.1209/0295-5075/95/50006. URL <https://iopscience.iop.org/article/10.1209/0295-5075/95/50006>.
- Durrett, R. Probability: Theory and Examples.
- Fang, Z., Jiang, Z., Chen, H., Zhang, X., Tang, K., Li, X., and Li, J. A Closer Look on Memorization in Tabular Diffusion Model: A Data-Centric Perspective, August 2025. URL <http://arxiv.org/abs/2505.22322>. arXiv:2505.22322 [cs].
- Favero, A., Sclocchi, A., and Wyart, M. Bigger Isn't Always Memorizing: Early Stopping Overparameterized Diffusion Models, September 2025. URL <http://arxiv.org/abs/2505.16959>. arXiv:2505.16959 [cs].
- Garnier-Brun, J., Biggio, L., Beltrame, D., Mézard, M., and Saglietti, L. Biased Generalization in Diffusion Models, March 2026. URL <http://arxiv.org/abs/2603.03469>. arXiv:2603.03469 [cs].
- George, A. J., Veiga, R., and Macris, N. Denoising Score Matching with Random Features: Insights on Diffusion Models from Precise Learning Curves, February 2025. URL <http://arxiv.org/abs/2502.00336>. arXiv:2502.00336 [cs].
- Gu, X., Du, C., Pang, T., Li, C., Lin, M., and Wang, Y. On Memorization in Diffusion Models, February 2025. URL <http://arxiv.org/abs/2310.02664>. arXiv:2310.02664 [cs].
- Jeon, D., Kim, D., and No, A. Understanding and Mitigating Memorization in Generative Models via Sharpness of Probability Landscapes, August 2025. URL <http://arxiv.org/abs/2412.04140>. arXiv:2412.04140 [cs].
- Kadkhodaie, Z., Guth, F., Simoncelli, E., and Mallat, S. Generalization in diffusion models arises from geometry-adaptive harmonic representation. *ArXiv*, abs/2310.02557:null, 2023. doi: 10.48550/arXiv.2310.02557. URL <https://www.semanticscholar.org/paper/a8724abaf519ab9113cf9dcc4c6d17f984de52cf>.
- Kadkhodaie, Z., Guth, F., Simoncelli, E. P., and Mallat, S. Generalization in diffusion models arises from geometry-adaptive harmonic representations, April 2024. URL <http://arxiv.org/abs/2310.02557>. arXiv:2310.02557 [cs].
- Kamb, M. and Ganguli, S. An analytic theory of creativity in convolutional diffusion models. June 2025. URL <https://openreview.net/forum?id=ilpL2qAcla>.
- Kim, J., Kim, S., and Lee, J.-S. How Diffusion Models Memorize, September 2025. URL <http://arxiv.org/abs/2509.25705>. arXiv:2509.25705 [cs].
- Krizhevsky, A. Learning Multiple Layers of Features from Tiny Images. 2009. URL <https://www.semanticscholar.org/paper/>

- 495 Learning-Multiple-Layers-of-Features-from-Wu, Y. H., Marion, P., Biau, G., and Boyer, C. Taking a Big  
496 5d90f06bb70a0a3dced62413346235c02b1aa086. Step: Large Learning Rates in Denoising Score Match-  
497 ing Prevent Memorization. In *Proceedings of Thirty*  
498 *Eighth Conference on Learning Theory*, pp. 5718–5756.  
499 PMLR, July 2025. URL <https://proceedings.mlr.press/v291/wu25a.html>.
- 500 Li, S., Chen, S., and Li, Q. A Good Score  
501 Does not Lead to A Good Generative Model, Jan-  
502 uary 2024. URL <http://arxiv.org/abs/2401.04856>. arXiv:2401.04856 [cs].
- 503 Liu, Z., Luo, P., Wang, X., and Tang, X. Deep Learn-  
504 ing Face Attributes in the Wild. pp. 3730–3738.  
505 IEEE Computer Society, December 2015. ISBN  
506 978-1-4673-8391-2. doi: 10.1109/ICCV.2015.425.  
507 URL <https://www.computer.org/csdl/proceedings-article/iccv/2015/8391d730/12OmNzGlRCR>.
- 508 Lucibello, C. and Mézard, M. Exponential Capac-  
509 ity of Dense Associative Memories. *Physical*  
510 *Review Letters*, 132(7):077301, February 2024.  
511 doi: 10.1103/PhysRevLett.132.077301. URL  
512 [https://link.aps.org/doi/10.1103/](https://link.aps.org/doi/10.1103/PhysRevLett.132.077301)  
513 [PhysRevLett.132.077301](https://link.aps.org/doi/10.1103/PhysRevLett.132.077301).
- 514 Pham, B., Raya, G., Negri, M., Zaki, M. J., Ambrogioni,  
515 L., and Krotov, D. Memorization to Generalization: The  
516 Emergence of Diffusion Models from Associative Mem-  
517 ory. November 2024. URL <https://openreview.net/forum?id=zVMMaVy2BY>.
- 518 Pidstrigach, J. Score-Based Generative Models Detect Man-  
519 ifolds. October 2022. URL <https://openreview.net/forum?id=AiNrnIrDfd9>.
- 520 Ronneberger, O., Fischer, P., and Brox, T. U-Net: Convo-  
521 lutional Networks for Biomedical Image Segmentation,  
522 May 2015. URL <http://arxiv.org/abs/1505.04597>. arXiv:1505.04597 [cs].
- 523 Ross, B. L., Kamkari, H., Wu, T., Hosseinzadeh, R., Liu,  
524 Z., Stein, G., Cresswell, J. C., and Loaiza-Ganem, G. A  
525 Geometric Framework for Understanding Memorization  
526 in Generative Models, March 2025. URL <http://arxiv.org/abs/2411.00113>. arXiv:2411.00113  
527 [stat].
- 528 Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and  
529 Goldstein, T. Diffusion Art or Digital Forgery? Investi-  
530 gating Data Replication in Diffusion Models, Decem-  
531 ber 2022. URL <http://arxiv.org/abs/2212.03860>. arXiv:2212.03860 [cs].
- 532 Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and  
533 Goldstein, T. Understanding and Mitigating Copying  
534 in Diffusion Models. November 2023. URL <https://openreview.net/forum?id=HtMXRgUMt>.
- 535 Wen, Y., Liu, Y., Chen, C., and Lyu, L. DETECTING, EX-  
536 PLAINING, AND MITIGATING MEMO- RIZATION  
537 IN DIFFUSION MODELS. 2024.
- 538 Yoon, T., Choi, J. Y., Kwon, S., and Ryu, E. K. Diffu-  
539 sion Probabilistic Models Generalize when They Fail to  
540 Memorize. July 2023. URL <https://openreview.net/forum?id=shciCbSk9h#all>.

## A. Average over draws of samples

In this appendix, we give a detailed derivation of the results presented in Section 4. Our starting point is the expansion

$$W_N(y, x) = \overline{\exp(-\exp(-N(y - z(x))))} = \sum_{n=0}^{\infty} \frac{(-1)^n}{n!} \exp(-nN(y - z(x))).$$

Recall that the notation  $\overline{f(z)}$  denotes the average over the data set.

As a first step, we clarify the conditions under which we can exchange the order of summation and average on the right hand side of the expression. Let us first assume that  $K$  is bounded from above by  $M \in \mathbb{R}$ . Almost all kernels will fulfill this constraint. Then we find that likewise  $Z$  is bounded from above and below,  $Z \in [0, M]$ . Moreover, since we must have  $z(x) \leq 0$ , we may also assume that  $y \leq 0$ , because we wish to evaluate the cumulative distribution function  $\Pr(\{z(x) \leq y\})$  only at valid values of  $y$ .

Let us now define

$$S_{m,N}(y, x) = \sum_{n=0}^m \frac{(-1)^n}{n!} \exp(nNy) Z^n(x)$$

which trivially fulfills

$$\overline{S_{m,N}(y, x)} = \sum_{n=0}^m \frac{(-1)^n}{n!} \exp(nNy) \overline{Z^n(x)}.$$

Moreover, we have that  $|S_{m,N}(y, x)|$  is also bounded, which follows directly from

$$|S_{m,N}(y, x)| \leq \sum_{n=0}^m \frac{1}{n!} |(-1)^n \exp(nNy) Z^n(x)| \leq \sum_{n=0}^m \frac{M^n}{n!} \leq \exp(M).$$

Then we find that  $S_{m,N}(y, x)$  fulfills all conditions of the dominated convergence theorem (see e.g. section 1.6 in (Durrett)) which states that if  $S_{m,N}(y, x) \rightarrow W_N(y, x)$  as  $m \rightarrow \infty$  and  $|S_{m,N}(y, x)|$  is bounded then  $\overline{S_{m,N}(y, x)} \rightarrow \overline{W_N(y, x)}$ . It follows that we may switch the order of summation and average.

Hence, if the kernel is bounded from above, we find

$$W_N(y, x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{n!} \exp(nNy) \overline{Z^n(x)}. \quad (20)$$

We now proceed to find an expression for  $\overline{Z^n(x)}$ . As a first step, we define the log-Fourier transform

$$\kappa(\xi) = \ln \frac{1}{2\pi} \int dx \exp(ix^T \xi) K(x)$$

$$K(x) = \int d\xi \exp(-ix^T \xi + \kappa(\xi))$$

which is also known of the cumulant generating function of  $K$ . Making use of this expression, we reformulate

$$Z_N(x) = \frac{1}{P} \sum_{\mu} \int d\xi^{\mu} \exp(-i(x - x^{\mu})^T \xi^{\mu} + \kappa(\xi^{\mu})).$$

Although this operation initially to complicated  $Z$  by introducing the integrals over  $\xi^{\mu}$ , observe that the samples  $x^{\mu}$  now only appear in the exponent in a linear fashion, which will simplify taking their average. We now want to evaluate

$$\begin{aligned} \overline{Z^n(x)} &= \int \prod_{\mu} d\rho(x^{\mu}) Z^n(x) \\ &= P^{-n} \int \prod_{\mu=1}^P d\rho(x^{\mu}) \sum_{\mu_1, \dots, \mu_P} \int \left( \prod_{i=1}^n d\xi_i \right) \exp \left( -i \sum_{i=1}^n (x - x^{\mu_i})^T \xi_i + \kappa(\xi_i) \right) \end{aligned}$$

The problem with evaluating this integral is now the following: Since we are summing over  $\mu_1, \dots, \mu_P$ , there occur cases where two indices are equal  $\mu_i = \mu_j$  which couple the corresponding integrals in  $\xi_i$ . We now group integrals of all equal indices together, if there are  $m$  such indices, this leads to the following integral

$$\bar{I}(x, m) = P^{-m} \int d\rho(x^\mu) \int \left( \prod_{i=1}^m d\xi_i \right) \exp \left( -i \sum_{i=1}^m (x - x^\mu)^T \xi_i + \kappa(\xi_i) \right),$$

which we will now simplify. First, we define the cumulant generating function of  $\rho$ ,

$$\phi(\xi) = \ln \int dx \rho(x) \exp(i\xi) \quad (21)$$

We can then compute the average over  $x^\mu$  as

$$\bar{I}(x, m) = P^{-m} \int \left( \prod_{i=1}^m d\xi_i \right) \exp \left( \phi \left[ \sum_{i=1}^m \xi_i \right] - i \sum_{i=1}^m [x \xi_i + \kappa(\xi_i)] \right),$$

where now the sum over the  $\xi_i$  appears in the cumulant generating function  $\phi$ . The calculation can therefore be further simplified by defining an auxiliary variable  $r = \sum_{i=1}^m \xi_i$  and enforcing this definition with a  $\delta$ -constraint, which amounts to introducing a factor  $1 = \frac{1}{(2\pi)^d} \int d\tilde{r} dr \exp[i\tilde{r}^T (r - \sum_{i=1}^m \xi_i)]$  into the integral, which simplifies as follows:

$$\begin{aligned} \bar{I}(x, m) &= \frac{1}{(2\pi)^N P^m} \int d\tilde{r} dr \int \left( \prod_{i=1}^m d\xi_i \right) \exp \left( i(\tilde{r} - x)^T \left( \sum_{i=1}^m \xi_i \right) - i\tilde{r}r + \phi(r) + \sum_{i=1}^m \kappa(\xi_i) \right) \\ &= \frac{1}{(2\pi)^N P^m} \int d\tilde{r} dr \exp(-i\tilde{r}r + \phi(r)) \left[ \int d\xi \exp(-i(x - \tilde{r})^T \xi + \kappa(\xi)) \right]^m \\ &= \frac{1}{(2\pi)^N P^m} \int d\tilde{r} dr \exp(-i\tilde{r}r + \phi(r)) [K(x - \tilde{r})]^m \end{aligned}$$

where in the first to the second line, we exploited that all the integrals in  $\xi_i$  decouple, and in the second to the third, we used the fact that the reverse Fourier transform of the Kernel appears through  $\exp(\kappa(\xi))$ . Similarly, we can now use that the backtransform of  $\phi$  also naturally appears to find

$$\bar{I}(x, m) = \int d\tilde{r} \rho(\tilde{r}) \exp(m \{ \ln K(x - \tilde{r}) - \ln N \})$$

we now make use of this expression to find  $W_N$ . First, note that

$$\bar{Z}^n(x) = \sum_{\substack{m_1, \dots, m_k \\ \sum_i m_i = n \\ k \leq P}} \binom{n}{m_1, \dots, m_k} \bar{I}(x, m_1) \dots \bar{I}(x, m_k) \quad (22)$$

The binomial factor  $\binom{n}{m_1, \dots, m_k}$  accounts for all ways of obtaining  $k$  groups of identical indices  $\mu_{j_1} = \dots = \mu_{j_{m_i}}$ . The condition  $\sum_i m_i = n$  ensures that all powers sum up to  $n$ , and the condition  $k \leq P$  comes from the fact that we can have at

most  $k = P$  distinct terms (corresponding to no repetitions, each sample appears only once). Inserting this into  $W_N$ , we find

$$\begin{aligned}
 W_N(y, x) &= \sum_{n=0}^{\infty} \exp(-ndy) (-1)^n \frac{1}{n!} \sum_{\substack{m_1, \dots, m_k \\ \sum_i m_i = n \\ k \leq P}} \frac{n!}{m_1! \dots m_k!} \bar{I}(x, m_1) \dots \bar{I}(x, m_k), \\
 &= \sum_{n=0}^{\infty} \exp(-ndy) (-1)^n \sum_{\substack{m_1, \dots, m_k \\ \sum_i m_i = n \\ k \leq P}} \frac{1}{m_1! \dots m_k!} \bar{I}(x, m_1) \dots \bar{I}(x, m_k), \\
 &= \left[ \sum_{m=0}^{\infty} \frac{\exp(-dy)^m (-1)^m}{m!} \bar{I}(x, m) \right]^P
 \end{aligned}$$

The exponent in  $P$  comes from the fact that we cannot have more than  $P$  different powers  $m_i$ s at the same time. Defining

$$G_N(\tilde{r}, x) = \frac{1}{N} \ln K(x - \tilde{r}) - \alpha$$

we find that

$$\bar{I}(x, m) = \int d\tilde{r} \rho(\tilde{r}) \exp(m N G_N(\tilde{r}, x))$$

which inserted into the expression for  $W_N$  finally yields

$$\begin{aligned}
 W_N(y, x) &= \left[ \int d\tilde{r} \rho(\tilde{r}) \sum_{m=0}^{\infty} \frac{[-\exp(-Ny + N G_N(\tilde{r}, x))]^m}{m!} \right]^P \\
 &= \left[ \int d\tilde{r} \rho(\tilde{r}) \exp(-\exp\{-N[y - G_N(\tilde{r}, x)]\}) \right]^P
 \end{aligned}$$

which, when relabeling  $\tilde{r} \rightarrow \tilde{x}$ , finally yields Equation (11). This completes the derivation of our results for general kernels  $K$  and general data distributions  $\rho$ .

### A.1. Hard spherical shells

We now outline how to derive the result for kernels that are hard shells. Using the definition of  $p_{\text{in}}$ , we find that the integral  $I_o$  defined in Equation (12) decomposes into two areas: those where  $K = 0$  and those where  $K = 1$ . Using this distinction, we then find

$$I_o(y, x) = p_{\text{in}}(x) (1 - \exp(-\exp\{-N(y + \alpha)\})) \quad (23)$$

now additionally assuming that

$$\nu(x) = \frac{1}{N} \ln p_{\text{in}}(x) \quad (24)$$

scales as  $\mathcal{O}_N(1)$ , we find that

$$W_N(x, y) \rightarrow \begin{cases} \Theta(y - \nu_{\text{in}}(x)) & y + \alpha \geq 0 \\ \Theta(-(\nu_{\text{in}}(x) + \alpha)) & \text{else} \end{cases}$$

which yields the probability that  $z \geq y$  in the limit. The first line shows that when  $z$  is larger than  $\alpha$ , we recover  $z \rightarrow \nu_{\text{in}}(x)$ .

The second line must be treated with care. If  $\nu(x) + \gamma > 0$ , then this implies  $y < \nu_{\text{in}}(x)$ , meaning  $y$  is strictly smaller than  $z$ , i.e.  $y$  is just not large enough to have had the jump in the cumulative probability. On the other hand, when  $\nu_{\text{in}}(x) + \gamma < 0$ , then we have that  $P_N(\{z_N(x) < y\}) = 1$  everywhere, meaning that  $z$  is smaller than any finite value of  $y$ . In this case, the jump happens at  $z \rightarrow -\infty$ , or zero local density.

Summarizing these findings we find that the cumulative distribution has the shape of a step function. Hence the value of  $z(x)$  concentrates on the "jump location", given by

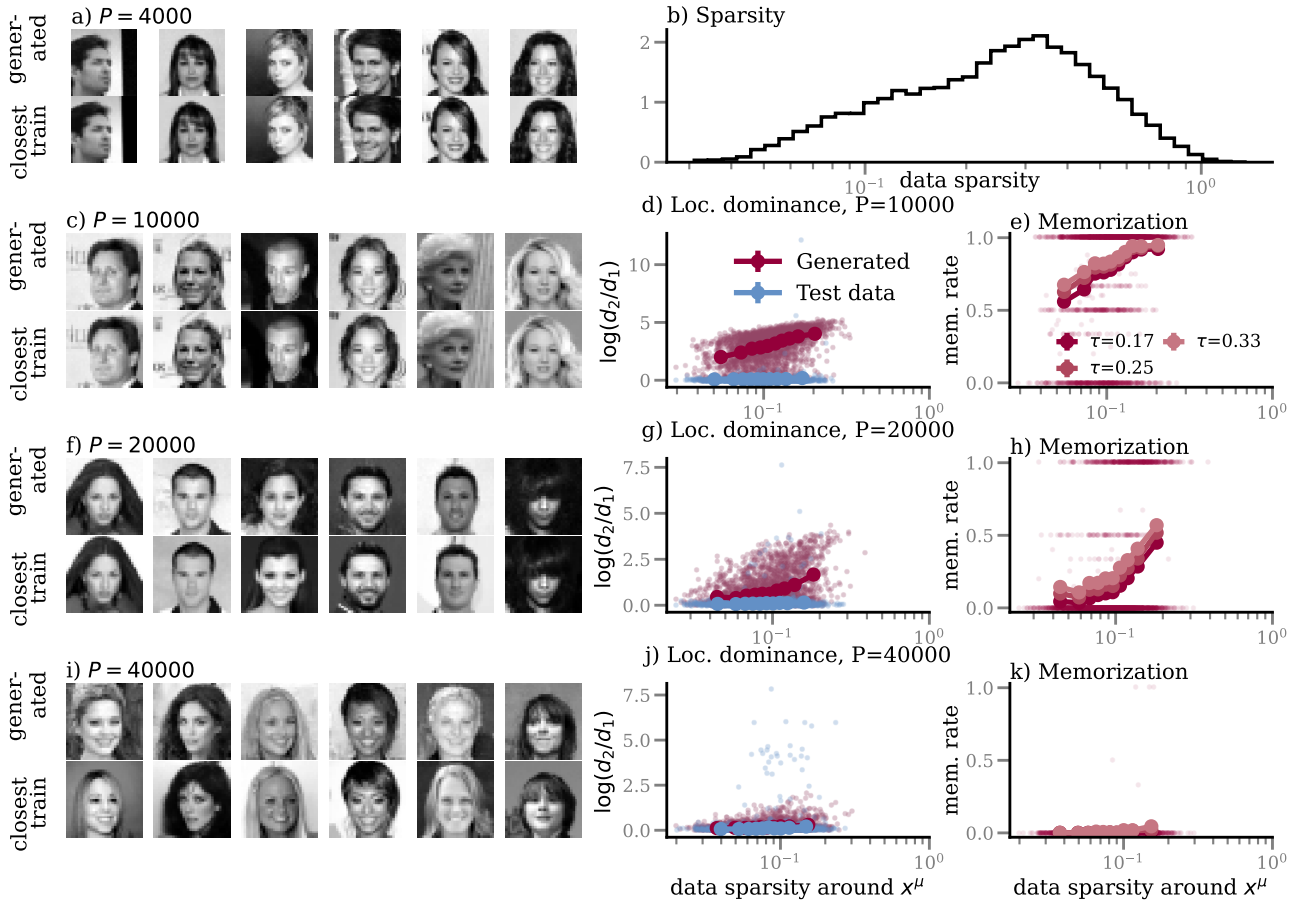
$$z(x) \rightarrow \begin{cases} \nu_{\text{in}} & \text{if } \nu_{\text{in}} \geq -\alpha \\ -\infty & \text{if } \nu_{\text{in}} < -\alpha \end{cases} \quad (25)$$

which exactly the expression we report in Section 4.

## B. Additional Experiments

### B.1. Local sparsity and dominance

In Figure 3, we report the analogous experiment to Figure 1 for CelebA data, which we downsample to  $32 \times 32$  greyscale pixels. We find that both local dominance and memorization increase with local sparsity. In comparison to CIFAR-10, we observe that the CelebA dataset appears to have fewer samples in very low density regions ( compare logarithmic scale of Figure 3 d) to Figure 1 d).



**Figure 3. Local memorization in CelebA data.** a) Sketch of kernel density approximation and local memorization phenomenon. b) Distribution of data sparsity for 40,000 samples from CelebA. c) Generated samples + closest training image (measured by cosine similarity) for diffusion model trained on  $P = 4000$  training examples. d) Local dominance as a function of sparsity. For each generated sample  $\hat{x}$ , we compute distances  $d_1$  and  $d_2$  to its nearest and second-nearest training samples, and define dominance via their ratio. Each point corresponds to a training sample  $x^\mu$ , aggregating generated samples assigned to it. Scatter shows raw values, markers indicate binned averages. Blue points show the same quantity computed using test data, providing a baseline without memorization. e) Memorization rate per training sample  $x^\mu$ , defined as the fraction of generated samples satisfying  $d_1/d_2 < \tau$ . Scatter shows raw values for  $\tau = 0.1$ , markers indicate binned averages across thresholds. f)–k) Same measurements for models trained on larger datasets. As the number of training samples increases, both dominance and memorization decrease, and their dependence on sparsity weakens. Overall, sparse regions exhibit strong single-sample dominance and higher memorization, while dense regions promote interpolation across multiple training points.

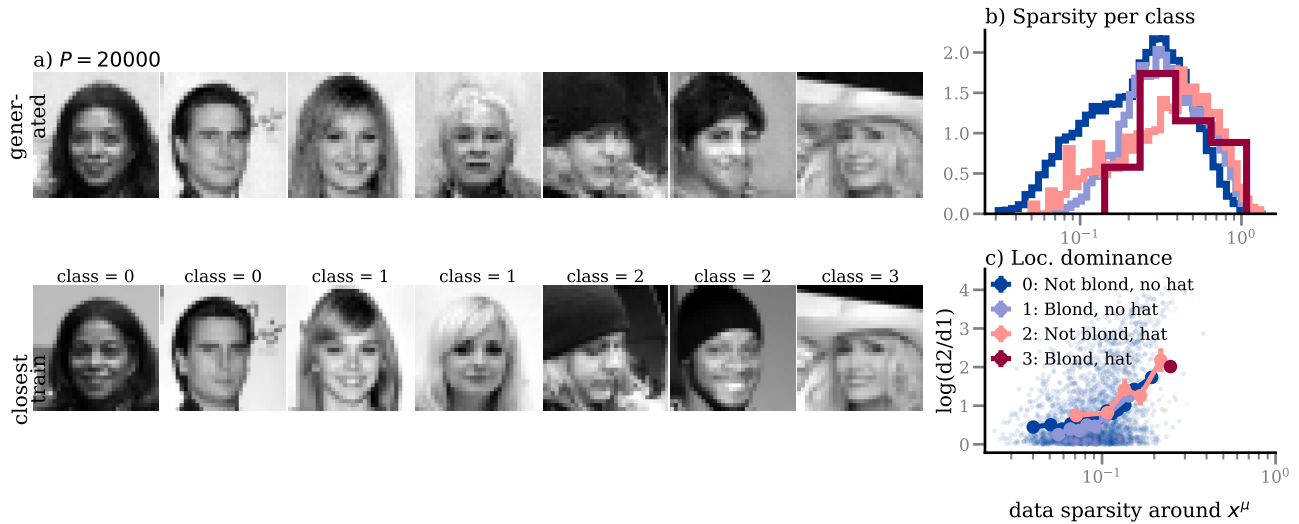


Figure 4. **Attribute-conditioned memorization in CelebA.** a) Generated samples + closest training image (measured by cosine similarity) for diffusion model trained on  $P = 20000$  training examples, sorted by classes. b) Distribution of data sparsity for 40000 samples from CelebA. d) Local dominance of training sample  $x^\mu$  against local data sparsity around  $x^\mu$  for a diffusion model trained on  $P = 20000$  training examples from CelebA.

## B.2. Attribute-dependent Memorization, CelebA

The CelebA dataset consists of celebrity portraits that are annotated with "attributes" such as hair color or accessories. We construct four classes from these attributes, conditioning on "blond" and "not blond" as well as "wearing hat" and the opposite. Again, we find that local dominance (and thus memorization) increases with data sparsity, see Figure 4. Samples from classes where "wearing hat" is true are typically more diverse. Therefore this class has a higher average sparsity per class, and is therefore more likely to be memorized. These classes are not balanced: approximately one sixth of images has attribute "blond" and approximately 5% of images have attribute "wearing hat". Consequently the diffusion model generates fewer samples that are memorized data points with these attributes, and even fewer samples are closest to training data where both attributes are true. This is reflected in the lower bin resolution in Figure 4b) and fewer such points appearing in Figure 4 c).

## B.3. Class-dependent Memorization, CIFAR-10+MNIST

In a separate experiment, we train diffusion models on a mixture of MNIST (Deng, 2012) and CIFAR-10 images. MNIST, consisting of handwritten digits, is structurally simpler and occupies a much more concentrated region of image space, whereas CIFAR-10 exhibits substantially higher variability. Treating MNIST and CIFAR-10 as two distinct classes, we observe that MNIST samples consistently exhibit lower local sparsity than CIFAR-10 samples. As predicted by the theory, this translates into reduced memorization: in Figure 5, both local dominance and memorization rates increase with sparsity. For the comparison of CIFAR and MNIST, all distances are computed on  $\ell_2$ -normalized samples, so that sparsity reflects relative geometric structure (angular similarity) rather than differences in overall scale. We verified that using unnormalized samples yields qualitatively similar results. On average, MNIST points lie in denser regions and are therefore memorized less than CIFAR-10 points.

## C. Experimental Details

**Model.** We train denoising diffusion models using a standard discrete-time formulation with  $T = 1000$  timesteps. The score network is parameterized by a U-Net architecture (Ronneberger et al., 2015) with convolutional filters. All models operate directly in pixel space.

**Training procedure.** Models are trained using the standard diffusion objective with mean squared error loss. Optimization is performed using Adam with learning rate  $10^{-4}$  and batch size 100. Training proceeds for a fixed number of  $10^6$  gradient

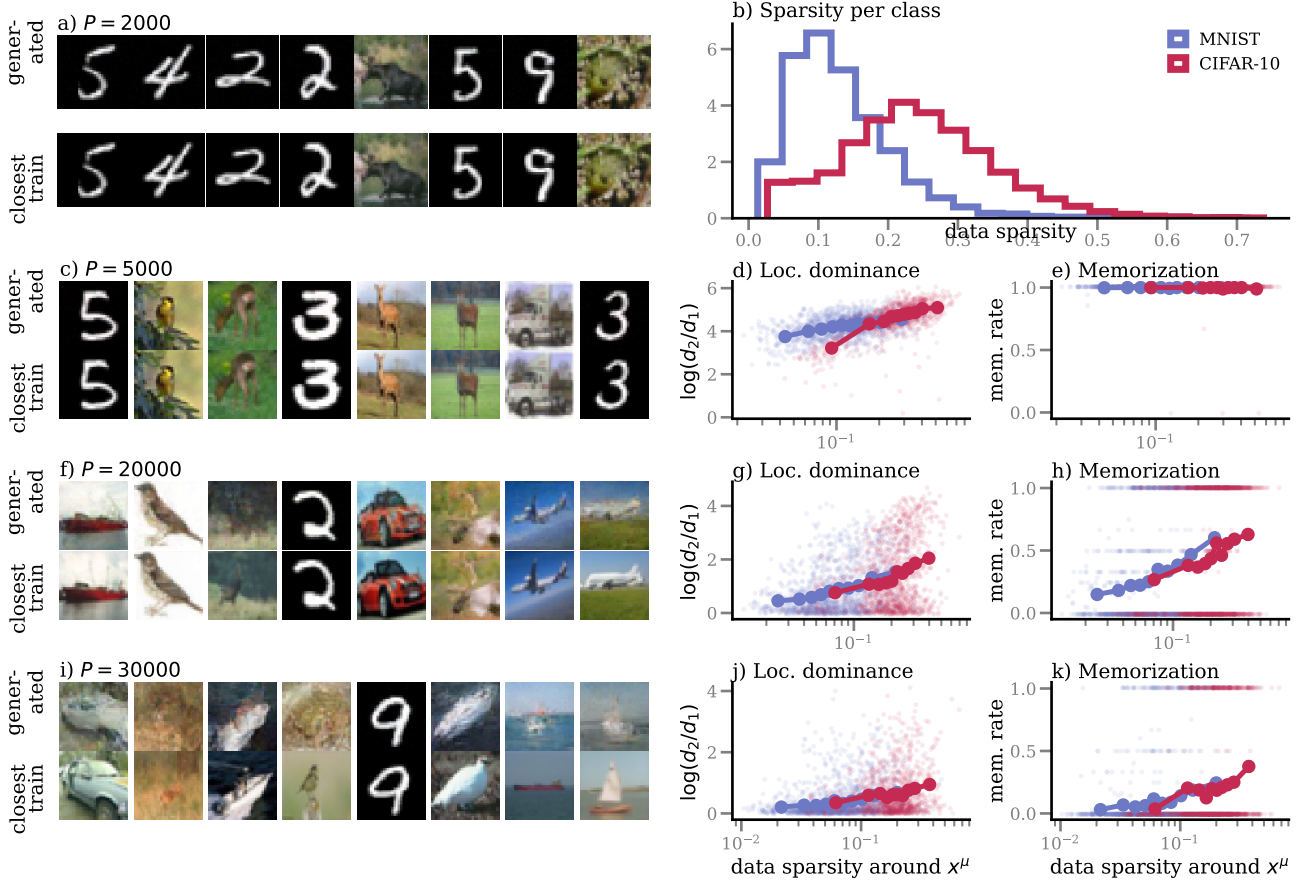


Figure 5. **Class-wise memorization on combined MNIST and CIFAR-10.** a) Generated samples together with their closest training example (measured by cosine similarity) for a diffusion model trained on  $P = 2000$  samples. b) Distribution of per-class sparsity around training points for four classes in CIFAR-10 measured from 20,000 training samples. c) Generated samples together with their closest training example (measured by cosine similarity) for a diffusion model trained on  $P = 5000$  samples. d) Class-wise local dominance of a training point  $x^\mu$  as a function of local data sparsity. Local dominance is defined via the ratio between the distance  $d_1$  to the closest generated sample and the distance  $d_2$  between that generated sample and the second-closest training point. MNIST points are shown in blue, CIFAR-10 points in red. e) Memorization rate of each sample  $x^\mu$ , defined as the fraction of generated samples satisfying  $d_1/d_2 < \tau = 1/3$ . Different classes are shown as different colors. f)–k) Same measurements for models trained on larger datasets. As the number of training samples increases, both dominance and memorization decrease, and their dependence on sparsity weakens.

880 steps for all models.

881  
 882 **Evaluation** For all experiments, the dataset is deterministically split into a training set of size  $P$  and a held-out test set.  
 883 The same split is reused for all evaluations to ensure comparability across checkpoints and metrics. At evaluation time,  
 884 we generate  $N_{\text{gen}} = 4000$  samples from the trained diffusion model. We compare generated samples to training and test  
 885 data using either cosine similarity or normalized  $\ell_2$  distance. In the cosine case, all inputs are flattened and  $\ell_2$ -normalized  
 886 before computing similarities. Memorization is quantified at the level of individual generated samples. For each generated  
 887 sample  $\tilde{x}$ , we compute its two nearest training neighbors and define distances  $d_1$  and  $d_2$ . A sample is considered memorized  
 888 if  $\frac{d_1}{d_2} < \tau$  for a range of thresholds  $\tau \in \{1/6, 1/4, 1/3, 1/2, 2/3\}$ . This yields both per-sample and per-training-point  
 889 memorization statistics. Local data density is estimated using nearest-neighbor statistics computed solely on the training set.  
 890 For each training sample, we compute the distances to the  $k$  nearest neighbors for  $k \in \{2, 5, 10, 20, 50\}$  and average over  
 891 these distances as defined in Equation (16) These statistics serve as proxies for local sparsity and are later correlated with  
 892 memorization behavior.

### 893 C.1. CIFAR-10 Experiments

894 We construct the CIFAR-10 dataset by combining the original training and test splits (60,000 images total), and then  
 895 randomly partitioning them into three equal subsets (train/validation/test), each containing approximately 20,000 images.  
 896 All images are represented as  $32 \times 32$  RGB tensors normalized to  $[0, 1]$ . No additional preprocessing or augmentation is  
 897 applied.  
 898  
 899

### 900 C.2. MNIST + CIFAR-10 Experiments

901 To study memorization under controlled differences in local data density, we construct a combined dataset from CIFAR-10  
 902 and MNIST, which exhibit markedly different intrinsic complexities. CIFAR-10 images are highly variable, while MNIST  
 903 digits occupy a much more concentrated region of image space. This setup allows us to induce systematic differences in  
 904 local sparsity across classes.  
 905

906 We merge the original training and test splits of both datasets. To ensure compatibility, MNIST images are resized to  $32 \times 32$   
 907 and converted to three channels by duplicating the grayscale channel, resulting in RGB tensors consistent with CIFAR-10.  
 908 CIFAR-10 images are used without modification. All images are represented as  $32 \times 32$  RGB tensors normalized to  $[0, 1]$ ,  
 909 and no additional preprocessing or augmentation is applied.  
 910

911 To isolate dataset-level effects, we ignore the original class labels and instead assign a binary label indicating the dataset of  
 912 origin (CIFAR-10 or MNIST). The two datasets are then balanced by subsampling the larger dataset such that both contribute  
 913 an equal number of samples. The resulting dataset therefore consists of an equal mixture of CIFAR-10 and MNIST images.  
 914

### 915 C.3. CelebA Experiments

916 This dataset consists of high-resolution face images together with 40 binary attributes per image. All images are resized to  
 917  $32 \times 32$  pixels, and converted to grayscale. The resulting images are represented as single-channel tensors normalized to  
 918  $[0, 1]$ . No additional data augmentation is applied.  
 919  
 920  
 921  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934