

Implicit geometric regularization in flow matching via density weighted Stein operators

Anonymous authors

Paper under double-blind review

Abstract

Flow Matching (FM) has emerged as a powerful paradigm for continuous normalizing flows, yet standard FM implicitly performs an unweighted L^2 regression over the entire ambient space. In high dimensions, this leads to a fundamental inefficiency: the vast majority of the integration domain consists of low-density “void” regions where the target velocity fields are often chaotic or ill-defined. In this paper, we propose γ -Flow Matching (γ -FM), a density-weighted variant that aligns the regression geometry with the underlying probability flow. While density weighting is desirable, naive implementations would require evaluating the intractable target density. We circumvent this by introducing a Dynamic Density-Weighting strategy that estimates the *target* density directly from training particles. This approach allows us to dynamically downweight the regression loss in void regions without compromising the simulation-free nature of FM. Theoretically, we establish that γ -FM minimizes the transport cost on a statistical manifold endowed with the γ -Stein metric. Spectral analysis further suggests that this geometry induces an implicit Sobolev regularization, effectively damping high-frequency oscillations in void regions. Empirically, γ -FM significantly improves vector field smoothness and sampling efficiency on high-dimensional latent datasets, while demonstrating intrinsic robustness to outliers.

1 Introduction

The Manifold Hypothesis posits that high-dimensional real-world data concentrate near a low-dimensional manifold embedded in the ambient space (Fefferman et al., 2016). While theoretical verification of this hypothesis remains a subject of active research (Pope et al., 2021), the remarkable success of deep generative models offers a constructive validation: if data were uniformly distributed in the high-dimensional void, efficient learning of the probability distribution would be computationally intractable (Bengio et al., 2013). Thus, the capability to accurately model and sample from the data distribution is, in itself, a testament to the existence of such low-dimensional structures.

Flow Matching (FM) has emerged as a powerful paradigm for capturing this distribution by unifying diffusion models and continuous normalizing flows (CNFs) (Lipman et al., 2023; Albergo & Vanden-Eijnden, 2023). It directly regresses a vector field that transports a simple base distribution to the data distribution. Unlike score-based diffusion models, FM avoids the need to estimate the score function or to solve reverse-time stochastic differential equations. Instead, it learns a deterministic ordinary differential equation (ODE) whose solution defines an invertible map between latent noise and data.

Despite its conceptual elegance, FM faces a fundamental challenge in high-dimensional settings: the curse of dimensionality and volume imbalance. In high dimensions, the data manifold occupies a negligible fraction of the ambient space, and the majority of the integration domain consists of low-density “voids.” In standard FM, the training objective uniformly integrates the regression error over the entire probability path. Consequently, the model is forced to solve the regression problem even in these vast void regions, where probability paths are sparse or effectively irrelevant to the final generation quality. Forcing a neural network to fit target velocities in these “don’t-care” regions results in a rough vector field that wastes model capacity and causes numerical stiffness during ODE integration.

In this work, we propose γ -Flow Matching (γ -FM) as a principled remedy for this issue. Our key idea is to reinterpret FM as a regression problem under a density-weighted geometry induced by the γ -divergence (Fujisawa & Eguchi, 2008). Instead of treating all spatial locations equally, we utilize the power density $p_t(x)^\gamma$ as an importance weight that naturally highlights the data manifold. This approach effectively “focuses” the learning process: the model prioritizes accurate vector field estimation where the data actually resides, while being allowed to remain smooth and simple in the empty ambient space. Latent flow matching has been independently explored by Lipman et al. (2023), who combine standard flow matching with pre-trained autoencoders and provide a Wasserstein-2 control for the resulting latent flows. Our work is complementary: we adopt a similar latent-flow setting, but instead of modifying the representation, we modify the *regression geometry* itself via the γ -weighted objective and analyse its effect through a γ -Stein and nonlinear Fokker-Planck viewpoint. While Chen & Lipman (2024) explicitly extend Flow Matching to Riemannian manifolds, our approach induces an implicit manifold geometry in the ambient space via density weighting, avoiding the need for explicit charts or geodesics.

Weighting Schemes in Diffusion and Flow Models In the realm of diffusion models, the choice of the weighting function plays a crucial role in balancing different signal-to-noise ratios across diffusion times. For instance, methods like EDM and variance-preserving (VP) SDEs employ carefully designed noise schedules to shape the training objective. However, these approaches typically rely on weights that depend solely on the time t or the noise schedule. Our proposed γ -Flow Matching differs fundamentally by introducing *spatially* varying weights based on the model density $p_t(x)$, thereby prioritizing regions where the model is confident while deprioritizing regions of low trust.

Density-weighted divergences and geometry The core motivation for our method roots in the geometry induced by density-powered divergences, in particular the γ -divergence, which defines an L^2 -type structure weighted by $p(x)^\gamma$. Classically, the γ -divergence has been used to downweight outliers and contaminated data, since regions where $p(x)$ is small automatically receive a small weight. In our setting we reinterpret this mechanism geometrically: low-density regions in the ambient space behave as geometric “voids” where the teacher signals are unstable and less informative, and the p^γ -weighting provides a way to organize regression on the data manifold while de-emphasizing these regions.

Flow Regularization Regularizing continuous flows to improve ODE solver stability is an active area of research. Typical strategies involve explicit Lipschitz penalties, spectral normalization, or Jacobian regularization to smooth the vector field. In contrast, our approach introduces an *implicit* regularization mechanism: by modifying the loss geometry via a density weight, the learned vector field naturally avoids wild oscillations in the voids. This can be seen as a form of geometric regularization that shapes the flow to follow the data manifold more faithfully, without needing to impose explicit gradient penalties.

Our contributions are summarized as:

1. **Dynamic Density-Weighting:** We formulate γ -FM as a weighted regression scheme that leverages the target density p_t . We propose a tractable, simulation-free implementation using particle-based density estimation.
2. **Variance Reduction Analysis:** We theoretically demonstrate that our weighting scheme minimizes the variance of the gradient estimator by suppressing high-noise signals in low-density regions.
3. **Geometric Regularization:** We show that γ -FM implicitly regularizes the vector field, reducing the Jacobian norm and enabling high-quality sampling with fewer function evaluations (NFE).
4. **Robustness & Efficiency:** Experiments on latent space CIFAR-10 confirm that γ -FM achieves a superior trade-off between generation quality and computational cost, while exhibiting strong robustness to outliers.

2 Flow Matching and Robust Divergences

2.1 Conditional Flow Matching

Avoiding the Likelihood Bottleneck Standard Continuous Normalizing Flows (CNFs) model the data distribution p_1 by transporting a simple base distribution p_0 through an ODE defined by a vector field v_θ . The change of variables formula for CNFs expresses

$$\log p_1(x) = \log p_0(\phi_1^{-1}(x)) - \int_0^1 \text{Tr}(\nabla_x v_\theta(x_t, t)) dt,$$

where ϕ_t is the flow map generated by the vector field v_θ , and the integral captures the accumulated divergence of v_θ along the trajectory. Maximizing such models typically involves maximizing the log-likelihood, which requires a stable and accurate computation of the Jacobian trace to account for the change in volume (the normalization constant).

Flow Matching (FM) circumvents this bottleneck by bypassing the explicit computation of the normalizing constant. Instead, it defines a probability path $(p_t)_{t \in [0,1]}$ between a known base distribution p_0 and the target p_1 , and learns a vector field whose associated continuity equation transports p_0 to p_1 . Rather than maximizing a likelihood, FM directly regresses a neural vector field v_θ to match a target vector field u_t that generates the desired probability path.

Formally, let $p_t(x)$ be a probability density path connecting p_0 and p_1 over $t \in [0, 1]$. This path satisfies the continuity equation:

$$\frac{\partial p_t}{\partial t} + \nabla \cdot (p_t u_t) = 0,$$

where $u_t(x)$ is the time-dependent vector field generating the flow. The goal is to regress the model vector field $v_\theta(x, t)$ to match the target vector field $u_t(x)$. The ideal marginal regression loss is defined as:

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \mathbb{E}_{x_t \sim p_t(x)} [\|v_\theta(x_t, t) - u_t(x_t)\|^2].$$

However, directly accessing the marginal vector field $u_t(x)$ and the marginal density $p_t(x)$ is generally intractable.

To solve this, [Lipman et al. \(2023\)](#) introduced *Conditional Flow Matching* (CFM), which uses a conditional probability path $p_t(x | x_1)$ given the data endpoint x_1 , and a corresponding conditional vector field $u_t(x | x_1)$. Crucially, it has been shown that the gradients of the intractable objective (2.1) are identical to those of the conditional objective:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \mathbb{E}_{x_1 \sim p_1} \mathbb{E}_{x_t \sim p_t(\cdot | x_1)} [\|v_\theta(x_t, t) - u_t(x_t | x_1)\|^2].$$

In this framework, the marginal vector field $u_t(x)$ implicitly emerges from the conditional field $u_t(x | x_1)$, allowing one to design $p_t(\cdot | x_1)$ in a convenient way (e.g., Gaussian interpolation) without computing the Jacobian trace or normalization constants.

2.2 Motivation: Geometric Focusing via γ -Divergence

From the information-geometric viewpoint, it is natural to interpret flow matching through the geometry induced by divergences and the associated Riemannian metrics on statistical manifolds ([Amari & Nagaoka, 2000](#); [Ay et al., 2017](#)). Standard Flow Matching implicitly minimizes the discrepancy between the target and model vector fields in an unweighted L^2 sense. From the perspective of Optimal Transport, this objective corresponds to minimizing the standard kinetic energy of the flow, $\mathcal{E}(v) = \int \|v(x)\|^2 p_t(x) dx$, which is associated with the Wasserstein-2 geometry and the linear heat equation. Recent works have explicitly targeted Wasserstein-optimal paths via minibatch optimal transport couplings ([Tong et al., 2024](#)), yet these approaches typically operate in the unweighted L^2 geometry. While statistically consistent, this "flat" geometry is inefficient in high dimensions because it assigns equal transport cost to the vast low-density "voids" as it does to the concentrated data manifold. To address this, we propose changing the underlying geometry

of the regression itself, moving from the L^2 regression (kinetic-energy) viewpoint to the robust γ -geometry, in analogy with the transition from the Fisher divergence to the γ -Fisher divergence; see Barp et al. (2019) for the diffusion score-matching divergence.

Formulation: Dynamic Density-Weighted Regression. We define the γ -Flow Matching (γ -FM) objective as the minimization of the *weighted* kinetic energy. Instead of the standard energy, we align the regression with the density-weighted geometry induced by the γ -divergence (Fujisawa & Eguchi, 2008). We define the weighted loss as:

$$\mathcal{L}_\gamma(\theta) := \mathbb{E}_{t \sim \mathcal{U}[0,1]} \mathbb{E}_{x_1 \sim p_1} \mathbb{E}_{x_t \sim p_t(\cdot | x_1)} [w_\gamma(x_t, t) \|v_\theta(x_t, t) - u_t(x_t | x_1)\|^2], \quad (1)$$

Theoretically, to align with the geometry of the γ -divergence, the weight should depend on the model density $q_{\theta(t)}$:

$$w_\gamma^{\text{ideal}}(x, t) \propto q_{\theta(t)}(x)^\gamma. \quad (2)$$

However, evaluating the model density $q_{\theta(t)}(x)$ during training is computationally prohibitive as it requires solving the ODE. Therefore, we adopt the target density $p_t(x)$ as a tractable proxy. Under the assumption that the model successfully tracks the target flow (i.e., $q_{\theta(t)} \approx p_t$), we define our practical weighting scheme as:

$$w_\gamma(x, t) := p_t(x)^\gamma. \quad (3)$$

This approximation allows us to compute weights solely based on the training data interpolation, preserving the simulation-free nature of Flow Matching. Here, we adopt the Conditional Flow Matching (CFM) framework, where $u_t(x | x_1)$ is the conditional vector field generating the probability path from noise to a specific data point x_1 . Crucially, since the training samples x_t are drawn from the target path p_t , this weighting naturally evolves over time. In the high-density manifold regions, the weight $p_t(x)^\gamma$ is significant, enforcing accurate vector field matching. Conversely, in the empty ambient space (voids) where $p_t(x) \approx 0$, the weight vanishes. This effectively removes the chaotic, ill-defined target signals in the voids from the optimization landscape.

Remark 2.1 (Geometric Interpretation via Weighted Transport). *Our weighting choice $w_\gamma \propto p_t^\gamma$ is not a heuristic modification; it fundamentally alters the metric structure of the transport problem. Standard FM minimizes the kinetic energy $\int \|v\|^2 p dx$, which underpins the Benamou-Brenier formula for optimal transport. In contrast, our objective \mathcal{L}_γ corresponds to minimizing the γ -weighted kinetic energy:*

$$\mathcal{E}_\gamma(v_t) = \int \|v_t(x)\|^2 p_t(x)^{1+\gamma} dx.$$

*This defines a Riemannian metric (specifically, the γ -weighted Fisher information metric) where the infinitesimal transport cost is scaled by the density power $p^{1+\gamma}$. In this geometry, distances in low-density regions are compressed to zero. Consequently, γ -FM does not simply "ignore" outliers; it solves the regression problem on a statistical manifold where the voids are geometrically insignificant. This modification naturally links the regression to the **Porous Medium Equation** (nonlinear diffusion) rather than the Heat Equation (linear diffusion), providing a theoretical guarantee for compact support preservation as discussed in Section 3, see Otto (2001) for extensive discussion.*

2.3 Tractability via Particle-Based Estimation

Evaluating the exact density $p_t(x)$ for the conditional probability path (e.g., Gaussian mixtures in CFM) can be computationally expensive or numerically unstable in high dimensions. Moreover, we seek a method that captures the *local* geometry of the batch without solving differential equations.

We circumvent this bottleneck by adopting a **particle-based estimation** strategy. Since the mini-batch samples $\mathcal{B}_t = \{x_t^{(i)}\}_{i=1}^B$ at any given time t are drawn from p_t , their spatial distribution provides a direct Monte Carlo estimate of the density. We employ a robust kernel-based proxy for the density. Specifically, we define the weight for a sample $x_t \in \mathcal{B}_t$ based on the distance to its k -nearest neighbors:

$$w_\gamma(x_t) \approx \exp\left(-\frac{\gamma}{\sigma} \bar{d}_k(x_t)\right), \quad (4)$$

where $\bar{d}_k(x_t) = \frac{1}{k} \sum_{j=1}^k \|x_t - x_t^{(j)}\|$ is the mean distance to the k nearest neighbors in the batch, and σ is a scaling constant. In effect, we set $\sigma = \text{median}\{\bar{d}_k(x_t^{(i)})\}_{i=1}^B$ within each minibatch. While a naive k -NN search can be computationally expensive, we show in Appendix B that the overhead is negligible for typical batch sizes and that the performance is robust to the choice of k . We emphasize that (4) is a monotone surrogate for the ideal escort weight $w_\gamma(x, t) \propto q_{\theta(t)}(x)^\gamma$: $\bar{d}_k(x_t)$ increases in locally low-density regions, hence \tilde{w}_γ down-weights updates in voids while preserving the standard FM objective structure.

This exponential weighting scheme has two key advantages:

- **Simulation-Free:** It requires only pairwise distance computations within the batch, preserving the efficiency of Standard FM.
- **Dynamic Adaptation:** It naturally adapts to the flow. At $t \approx 0$ (noise), particles are spread out, leading to uniform weights. At $t \rightarrow 1$ (data), particles concentrate on the manifold, creating a sharp weighting profile that isolates the data structure. In the high-dimensional voids surrounding the manifold, the effective density \hat{p}_t vanishes. Consequently, w_γ suppresses the regression loss in these empty regions, preventing the model from overfitting to unstable target signals where no data exists. By focusing the training budget solely on the populated regions of the probability path, γ -FM learns a vector field that is accurate on the manifold and smooth elsewhere, as evidenced by the reduced Jacobian norm in our experiments.

3 Theoretical Analysis

In this section, we analyze the theoretical properties of γ -FM. We establish that our density-weighting scheme is not merely a heuristic, but acts as a variance-optimal estimator and enforces physical constraints consistent with nonlinear diffusion. To provide a concrete basis for the following analysis, we summarize the practical training procedure of γ -FM in Algorithm 1. This algorithm implements the simulation-free density estimation discussed in Section 2. For simplicity we use linear interpolants in experiments.

Algorithm 1 Training γ -Flow Matching with Dynamic Density-Weighting**Require:** Training data \mathcal{D} , Batch size B , Weighting parameter $\gamma \geq 0$, Neighbors k **Ensure:** Trained vector field parameters θ

```

1: Initialize neural network parameters  $\theta$ 
2: while not converged do
3:                                      $\triangleright$  1. Sample flow matching variables
4:   Sample data batch  $x_1 \sim \mathcal{D}$ 
5:   Sample noise batch  $x_0 \sim p_0 = \mathcal{N}(0, I)$ 
6:   Sample time steps  $t \sim \mathcal{U}[0, 1]$ 
7:                                      $\triangleright$  2. Compute interpolants and targets
8:    $x_t \leftarrow (1 - t)x_0 + tx_1$ 
9:    $u_t \leftarrow x_1 - x_0$ 
10:                                      $\triangleright$  3. Dynamic Density-Weighting
11:   if  $\gamma > 0$  then
12:     Compute pairwise distances matrix for batch  $\{x_t\}$ 
13:     for  $i = 1$  to  $B$  do
14:        $\bar{d}_k(x_t^{(i)}) \leftarrow$  mean distance to  $k$ -nearest neighbors
15:        $w_i \leftarrow \exp\left(-\frac{\gamma}{\sigma} \bar{d}_k(x_t^{(i)})\right)$ 
16:     end for
17:     Normalize weights:  $w_i \leftarrow w_i / (\frac{1}{B} \sum_j w_j)$ 
18:   else
19:      $w_i \leftarrow 1$ 
20:   end if
21:                                      $\triangleright$  4. Optimization step
22:    $\mathcal{L}(\theta) \leftarrow \frac{1}{B} \sum_{i=1}^B w_i \|v_\theta(x_t^{(i)}, t) - u_t^{(i)}\|^2$ 
23:    $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta)$ 
24: end while

```

Although the modification in Algorithm 1 appears minimal, its theoretical implications are profound. In the remainder of this section, we rigorously justify this design choice, demonstrating that this simple weighting scheme induces a fundamental shift in the regression geometry.

3.1 Variance Reduction in High-Dimensional Voids

We first formalize the “Manifold Focusing” effect as a variance reduction problem. Recall that the conditional flow matching objective targets the individual paths $u_t(x|x_1)$. Let $\Sigma_t(x) := \text{Var}_{x_1|x}[u_t(x|x_1)]$ denote the intrinsic variance (ambiguity) of the target signal at location x .

Proposition 3.1 (Variance of the Weighted Estimator). *Assume that the gradient of the vector field model with respect to its parameters is bounded, i.e., there exists a constant $K > 0$ such that $\|\nabla_\theta v_\theta(x, t)\|_{\text{op}}^2 \leq K$ for all x, t . Then, the trace of the covariance matrix of the gradient estimator \hat{g}_γ satisfies the bound:*

$$\text{Tr}(\text{Var}[\hat{g}_\gamma]) \leq 4K \int_{\mathbb{R}^d} p_t(x) w_\gamma(x)^2 \text{Tr}(\Sigma_t(x)) dx + C_{\text{signal}}, \quad (5)$$

where C_{signal} represents the variance contribution from the learnable mean field.

Proof. Let the per-sample loss for a target path connecting x_0 to x_1 be $J_t(\theta; x_1) = w_\gamma(x_t) \|v_\theta(x_t) - u_t(x_t|x_1)\|^2$. The stochastic gradient is $\hat{g}_\gamma = \nabla_\theta J_t = 2w_\gamma(x_t)(\nabla_\theta v_\theta)^\top (v_\theta - u_t(x_t|x_1))$. Using the Law of Total Variance, we decompose the variance over the marginal density $p_t(x)$:

$$\text{Var}[\hat{g}_\gamma] = \mathbb{E}_{x \sim p_t}[\text{Var}(\hat{g}_\gamma | x)] + \text{Var}_{x \sim p_t}[\mathbb{E}[\hat{g}_\gamma | x]].$$

We focus on the first term (intrinsic noise). For a fixed location x , the conditional variance is due to the variability of the target $u_t(x|x_1)$:

$$\text{Var}(\hat{g}_\gamma | x) = \mathbb{E}_{x_1|x} \left[\|\hat{g}_\gamma - \mathbb{E}[\hat{g}_\gamma|x]\|^2 \right] \quad (6)$$

$$= 4w_\gamma(x)^2 \mathbb{E}_{x_1|x} \left[\left\| (\nabla_\theta v_\theta(x))^\top (u_t(x) - u_t(x|x_1)) \right\|^2 \right]. \quad (7)$$

Using the operator norm inequality $\|A^\top b\|^2 \leq \|A\|_{\text{op}}^2 \|b\|^2$, we have:

$$\text{Tr}(\text{Var}(\hat{g}_\gamma | x)) \leq 4w_\gamma(x)^2 \|\nabla_\theta v_\theta(x)\|_{\text{op}}^2 \text{Tr}(\Sigma_t(x)).$$

Applying the boundedness assumption $\|\nabla_\theta v_\theta(x)\|_{\text{op}}^2 \leq K$, we obtain:

$$\text{Tr}(\text{Var}(\hat{g}_\gamma | x)) \leq 4Kw_\gamma(x)^2 \text{Tr}(\Sigma_t(x)).$$

Integrating this with respect to $p_t(x)$ yields the first term of the bound in Eq. (5). The second term (C_{signal}) corresponds to $\text{Var}_x[\mathbb{E}[\hat{g}_\gamma|x]]$, which depends on the learnable signal $u_t(x)$ and is independent of the conditional noise variance $\Sigma_t(x)$. \square

In standard FM ($w_\gamma = 1$), the integral is dominated by the volume of the void space where $\Sigma_t(x)$ is large. By choosing $w_\gamma(x) \propto p_t(x)^\gamma$, the integrand becomes proportional to $p_t(x)^{1+2\gamma}\Sigma_t(x)$. Under the reasonable assumption that the signal ambiguity scales inversely with density (i.e., $\Sigma_t(x) \sim p_t(x)^{-\alpha}\Sigma_0$ for $\alpha > 0$), the γ -weighting with $\gamma \geq \alpha/2$ ensures that the noise contribution vanishes:

$$\lim_{p_t(x) \rightarrow 0} p_t(x)^{1+2\gamma}\Sigma_t(x) = 0.$$

This suggests that γ -FM suppresses gradient noise from the voids, concentrating the optimization budget on the high-density region.

When is the variance–density scaling plausible? The heuristic scaling $\Sigma_t(x) \propto p_t(x)^{-\alpha}$ is most plausible for conditional path designs in which the conditional target $u_t(x_t | x_1)$ becomes increasingly ill-conditioned in low-density regions of the marginal p_t . A representative example is Gaussian CFM, where x_t is obtained by adding Gaussian noise and u_t involves a score-like term of the intermediate marginal; in such settings the conditional variance of u_t given x_t typically increases as $p_t(x_t)$ decreases, reflecting amplification of estimation noise in “void” regions. More generally, for transport-noise interpolations that mix a deterministic drift toward x_1 with a stochastic perturbation, the signal-to-noise ratio of the conditional direction deteriorates away from the data manifold, so that $\text{tr} \Sigma_t(x)$ is larger where $p_t(x)$ is smaller. Our analysis in Proposition 3.1 should be read in this spirit: it formalizes how γ -weighting suppresses contributions from such high-variance, low-density regions, rather than requiring an exact power-law identity.

3.2 Physical Consistency with Liouville Dynamics

This subsection is intended as a geometric analogy, not as a derivation of the exact training dynamics of γ -FM. We use a classical model problem from optimal transport—the Wasserstein gradient flow of a Tsallis-type energy—to make precise a qualitative mechanism: density-power weighting suppresses motion in low-density (“void”) regions. In particular, the resulting continuum dynamics exhibit a degenerate diffusion whose effective diffusivity vanishes as $p \rightarrow 0$, which leads to a finite-speed propagation effect.

Consider the generalized γ -entropy functional:

$$\mathcal{F}_\gamma[p] = \frac{1}{\gamma} \int p(x)^{\gamma+1} dx.$$

Notice that this functional corresponds exactly to the *self-divergence term* in the definition of the γ -divergence (up to a sign). Just as minimizing γ -divergence statistically ignores outliers (Section 3), the gradient flow of its associated entropy \mathcal{F}_γ physically restricts the spread of probability mass.

The Wasserstein gradient flow is defined by the continuity equation driving the density along the gradient of the variation (see, e.g., [Jordan et al. \(1998\)](#); [Ambrosio et al. \(2008\)](#)).

$$\partial_t p = \nabla \cdot \left(p \nabla \frac{\delta \mathcal{F}_\gamma[p]}{\delta p} \right).$$

Calculating the variation $\frac{\delta \mathcal{F}_\gamma[p]}{\delta p} = \frac{\gamma+1}{\gamma} p^\gamma$ and substituting it gives the explicit dynamics:

$$\partial_t p(x, t) = \frac{\gamma+1}{\gamma} \nabla \cdot (p(x, t) \nabla p(x, t)^\gamma) = \Delta(p(x, t)^{1+\gamma}), \quad (8)$$

Equation (8) is the Wasserstein gradient flow of \mathcal{F}_γ and will be used here as an exactly analyzable toy model that captures the qualitative effect of γ -weighting. For $\gamma > 0$, the diffusion is *degenerate*: the effective diffusivity scales like p^γ and vanishes as $p \rightarrow 0$, which is the mechanism behind finite-speed propagation.

Proposition 3.2 (Preservation of Compact Support). *Let $p(x, t)$ be the unique weak solution to the PME (8) with $\gamma > 0$, subject to a non-negative initial condition $p_0 \in L^1(\mathbb{R}^d) \cap L^\infty(\mathbb{R}^d)$. If the initial support $\text{supp}(p_0)$ is compact, then the support $\text{supp}(p(\cdot, t))$ remains compact for all $t > 0$. Moreover, there exists a constant C depending on the initial mass such that the support is contained in a ball $\mathcal{B}(0, R(t))$ with*

$$R(t) \leq C(1+t)^\beta, \quad \text{where } \beta = \frac{1}{d\gamma + 2}.$$

Proof. The proof relies on the Comparison Principle for the Porous Medium Equation ([Vázquez, 2007](#)). The principle states that if two solutions u and v satisfy $u(x, 0) \leq v(x, 0)$ everywhere, then $u(x, t) \leq v(x, t)$ for all $t > 0$. We construct an explicit supersolution using the Barenblatt–Pattle solution $B(x, t; M)$, which represents the diffusion of a Dirac mass M . Let $B(x, \tau; M)$ be the Barenblatt profile centered at the origin with mass M at a time shift $\tau > 0$:

$$B(x, \tau) = \tau^{-\alpha} \left(C(M) - \kappa \tau^{-2\beta} |x|^2 \right)_+^{1/\gamma},$$

where $(\cdot)_+ = \max(\cdot, 0)$. Since the initial data p_0 is bounded and has compact support, we can choose a sufficiently large mass M and a time shift $\tau > 0$ such that the Barenblatt profile covers the initial data:

$$p_0(x) \leq B(x, \tau) \quad \text{for all } x \in \mathbb{R}^d.$$

By the Comparison Principle, this ordering is preserved for all subsequent times $t > 0$:

$$p(x, t) \leq B(x, t + \tau).$$

The support of the Barenblatt solution $B(\cdot, t + \tau)$ is explicitly known to be a ball of radius

$$R_B(t) = \sqrt{\frac{C(M)}{\kappa}} (t + \tau)^\beta.$$

Since $0 \leq p(x, t) \leq B(x, t + \tau)$, the support of p must be contained within the support of B . Thus,

$$\text{supp}(p(\cdot, t)) \subseteq \mathcal{B}(0, R_B(t)),$$

which implies that the support remains compact and expands at a rate of at most $O(t^\beta)$. In the limit $\gamma \rightarrow 0$, the exponent $\beta \rightarrow 1/2$, but the Barenblatt profile converges to a Gaussian which is strictly positive everywhere. Thus, the strict containment within a finite ball holds if and only if $\gamma > 0$. \square

In the toy model (8), the limit $\gamma \rightarrow 0$ reduces to the heat equation, whose solutions become instantly positive everywhere, reflecting infinite-speed propagation. By contrast, for $\gamma > 0$ the degeneracy at $p \approx 0$ yields an evolving interface and a finite-propagation behavior, as formalized in Proposition 3.2. While γ -FM does *not* literally solve (8), the same mechanism provides a useful intuition: when the weight behaves like $w_\gamma(x, t) \propto p_t(x)^\gamma$, updates are strongly down-weighted in low-density regions, and empirically this reduces spurious mass placed in “void” areas (a “void rejection” effect).

Remark 3.3 (Generalized entropy and maximum entropy principle). *The functional*

$$\mathcal{F}_\gamma[p] = \frac{1}{\gamma} \int p(x)^{\gamma+1} dx$$

is, up to an affine rescaling, the negative of the Tsallis q -entropy with $q = 1 + \gamma$ (Tsallis, 1988). It is well known that maximizing the Tsallis entropy under mass and second-moment constraints,

$$\int p(x) dx = 1, \quad \int \|x\|^2 p(x) dx = m_2,$$

yields generalized Gaussian (or q -Gaussian) densities of the form

$$p_q(x) = Z^{-1} \left(1 - (1 - q)\beta \|x - \mu\|^2 \right)_+^{\frac{1}{1-q}},$$

for suitable parameters $\beta > 0$ and $\mu \in \mathbb{R}^d$. These q -Gaussians are precisely the self-similar Barenblatt profiles of the porous medium equation (8) for an appropriate correspondence between q and the nonlinearity exponent $1 + \gamma$; see, e.g., Malacarne et al. (2001); Takatsu (2012). In particular, they exhibit compact support when $\gamma > 0$, providing concrete examples of the finite-support behaviour described in Proposition 3.2.

3.3 Motivating Example: 1D Double-Well Potential

To visualize the macroscopic effect of our weighting scheme, it is instructive to consider a one-dimensional toy problem where the dynamics can be analyzed exactly. Consider a target distribution defined by a double-well potential $V(x) = \frac{1}{4}x^4 - \frac{3}{4}x^2$, where the target density satisfies $p_{data}(x) \propto \exp(-V(x))$. We analyze the probability flow transporting a standard Gaussian noise $\mathcal{N}(0, 1)$ to this target.

The γ -weighted continuity equation can be mapped to a nonlinear Fokker–Planck equation with density-dependent diffusion:

$$\partial_t p_t = \nabla \cdot (p_t \nabla V) + D \nabla \cdot (p_t \nabla p_t^\gamma), \quad (9)$$

where D is a diffusion constant. The behavior of this system critically depends on γ :

- **Case $\gamma = 0$ (Heat Equation):** The diffusion term becomes linear ($D \Delta p_t$). A fundamental property of the Heat Equation is its *infinite speed of propagation*. As illustrated in Figure 1 (Standard FM), probability mass instantaneously leaks into the high-potential barrier regions (voids) between the wells. This corresponds to the model "hallucinating" paths where no data exists.
- **Case $\gamma > 0$ (Porous Medium Equation):** Eq. (9) becomes the Porous Medium Equation (PME). A hallmark of the PME is its *finite speed of propagation*. The effective diffusivity $D_{\text{eff}} \propto p_t^\gamma$ vanishes in low-density regions. Consequently, the diffusion physically stops at the boundaries of the wells. As shown in Figure 1 (γ -FM), this creates a sharp geometric barrier that confines the flow to the main modes, preventing leakage into the void.

This analytical example provides a rigorous justification for the "Manifold Focusing" effect: our dynamic weighting w_γ empirically emulates the finite-propagation physics of the PME, naturally enforcing compact support for the learned distribution.

3.4 Geometric Foundation: The γ -Stein Viewpoint

The physical behavior described by the PME (Section 3.2) is not accidental; it arises from the intrinsic geometry induced by the density weighting. While standard Flow Matching minimizes kinetic energy in a flat Euclidean geometry, we argue that γ -FM minimizes energy on a statistical manifold endowed with a density-dependent metric.

The γ -Stein Metric. From the perspective of Information Geometry (Eguchi, 2009), the γ -divergence generates a Riemannian metric structure on the space of probability densities. Following standard conventions

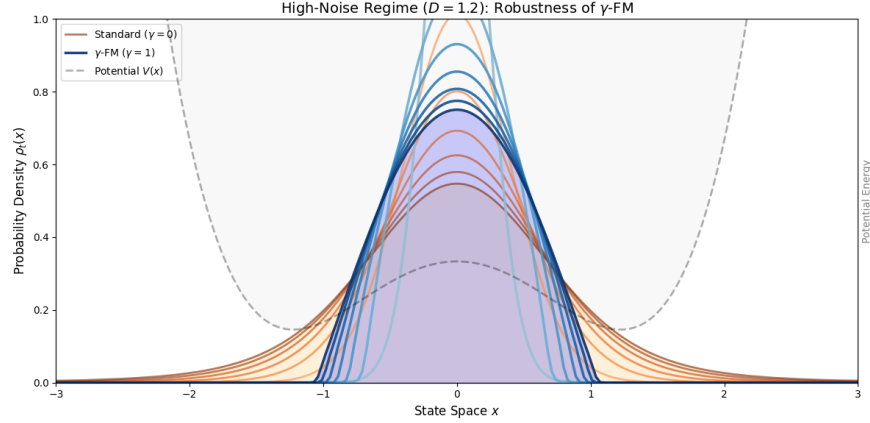


Figure 1: **Finite vs. Infinite Propagation.** Evolution of density under a double-well potential. **(Left) Standard FM ($\gamma = 0$):** Corresponds to linear diffusion (Heat Equation), causing probability mass to leak continuously into the low-density barrier (void). **(Right) γ -FM ($\gamma = 1$):** Corresponds to nonlinear diffusion (Porous Medium Equation) with finite propagation speed. The flow respects the potential barrier, keeping the mass tightly confined to the modes. This illustrates the mechanism of void rejection.

in information geometry, we hereinafter denote the parametric family of model densities by q_θ . The weighting $w_\gamma \approx q_\theta^\gamma$ in our objective (Eq. 1) naturally corresponds to the γ -weighted Fisher information metric $g^{(\gamma)}$ (Fujisawa & Eguchi, 2008; Matsuzoe, 2017):

$$g_{ij}^{(\gamma)}(\theta) = \int_{\mathbb{R}^d} \partial_{\theta_i} \log q_\theta(x) \partial_{\theta_j} \log q_\theta(x) q_\theta(x)^{1+\gamma} dx. \quad (10)$$

This metric measures distance based on the *escort measure* $d\mu_\gamma \propto q^{1+\gamma} dx$. Crucially, regions with low density $q(x) \approx 0$ make negligible contribution to the metric tensor. Geometrically, this means that "distances" in the void regions are compressed to zero, effectively removing them from the optimization landscape. The detailed discussion is given in Appendix A.

Flow Matching as Geodesic Optimization. Let $\{q_\theta : \theta \in \Theta\}$ be a parametric family of densities on \mathbb{R}^d , and let $t \mapsto \theta(t)$ be a time-dependent curve in parameter space with associated path of densities $q_{\theta(t)}$. A probability flow $(q_{\theta(t)}, v_t)_{t \in [0,1]}$ satisfies the continuity equation

$$\partial_t q_{\theta(t)}(x) + \nabla \cdot (q_{\theta(t)}(x) v_t(x)) = 0.$$

For a fixed θ , we equip the space of vector fields with the γ -weighted inner product

$$\langle u, v \rangle_{\theta, \gamma} := \int_{\mathbb{R}^d} u(x)^\top v(x) q_\theta(x)^{1+\gamma} dx,$$

and denote by $L_\gamma^2(q_\theta)$ the corresponding Hilbert space. The γ -Stein operator associated with q_θ is defined by

$$\mathcal{A}_{q_\theta}^{(\gamma)} f(x) := q_\theta(x)^{-1} \nabla \cdot (q_\theta(x)^{\gamma+1} f(x)), \quad f : \mathbb{R}^d \rightarrow \mathbb{R}^d.$$

We regard the tangent space of the statistical manifold at θ as the closure (in $L_\gamma^2(q_\theta)$) of the range of $\mathcal{A}_{q_\theta}^{(\gamma)}$:

$$T_\theta \mathcal{M}_\gamma := \overline{\{\mathcal{A}_{q_\theta}^{(\gamma)} f : f \in C_c^\infty(\mathbb{R}^d, \mathbb{R}^d)\}}.$$

In other words, $T_\theta \mathcal{M}_\gamma$ consists of all γ -Stein-type velocity fields that preserve total mass under the weighted geometry.

The γ -weighted flow-matching objective can be written as

$$\mathcal{L}_\gamma = \int_0^1 \|v_t - v_t^*\|_{L_\gamma^2(q_{\theta(t)})}^2 dt,$$

where v_t^* is the ideal velocity field induced by the target transport. For each t , the minimizer u_t of \mathcal{L}_γ over $T_{\theta(t)}\mathcal{M}_\gamma$ is the orthogonal projection of v_t^* onto the tangent space with respect to $\langle \cdot, \cdot \rangle_{\theta(t), \gamma}$. This motivates the definition of the projection operator

$$\Pi_{\theta, \gamma} : L_\gamma^2(q_\theta) \rightarrow T_\theta \mathcal{M}_\gamma, \quad \Pi_{\theta, \gamma}[f] := \arg \min_{u \in T_\theta \mathcal{M}_\gamma} \|f - u\|_{L_\gamma^2(q_\theta)}^2.$$

In particular, $u_t = \Pi_{\theta(t), \gamma}[v_t^*]$ is the γ -Stein projection of the ideal velocity field onto the statistical manifold.

Following Appendix A, the γ -Stein connection $\nabla^{(\gamma)}$ is characterised by the requirement that the covariant derivative of a time-dependent velocity field u_t along a curve $t \mapsto \theta(t)$ is obtained by projecting the ordinary time derivative back onto the tangent space:

$$D_t^{(\gamma)} u_t := \Pi_{\theta(t), \gamma}[\partial_t u_t].$$

A curve $\theta(t)$ is a (parametric) geodesic of the γ -Stein geometry if its associated velocity field $u_t \in T_{\theta(t)}\mathcal{M}_\gamma$ is covariantly constant, that is,

$$D_t^{(\gamma)} u_t = \Pi_{\theta(t), \gamma}[\partial_t u_t] = 0 \quad \text{for all } t \in [0, 1].$$

Therefore, minimizing \mathcal{L}_γ over admissible flows $(q_{\theta(t)}, u_t)_{t \in [0, 1]}$ can be interpreted as searching for geodesic curves on the statistical manifold endowed with the γ -Stein metric and connection. The case $\gamma = 0$ reduces to the flat L^2 geometry, where the tangent space coincides with $L^2(q_{\theta(t)})$ and $\Pi_{\theta(t), 0}$ is the identity, so that geodesics correspond to straight lines and the velocity field must be defined everywhere. For $\gamma > 0$, the metric is weighted by $q_{\theta(t)}^{1+\gamma}$, so that directions supported in low-density regions have negligible norm. As a consequence, the optimal velocity field produced by γ -FM concentrates on the high-density manifold of the data, not because it “ignores” data, but because the intrinsic geometry itself is dominated by the manifold structure through the γ -Stein projection.

3.5 Implicit Geometric Regularization

We now formalize the implicit-regularization intuition stated earlier by deriving a Sobolev-type roughness functional induced by the γ -weighted objective. Let $J_\theta(x, t) = \nabla_x v_\theta(x, t)$ be the Jacobian of the model. The stiffness of the ODE solver is controlled by the Lipschitz constant $L(t) \approx \sup_x \|J_\theta(x, t)\|_F$.

In unweighted regression, minimizing the loss in voids (where the target u_t is chaotic) requires $v_\theta(x, t)$ to change rapidly, driving $\|J_\theta(x, t)\|_F$ to be large. We formalize the regularization effect as a bound on the weighted Sobolev norm:

$$\mathcal{L}_{\text{roughness}}(t) := \int_{\mathbb{R}^d} q_{\theta(t)}(x) w_\gamma(x, t) \|\nabla_x v_\theta(x, t)\|_F^2 dx. \quad (11)$$

By downweighting the voids ($w_\gamma(x, t) \rightarrow 0$), γ -FM relaxes the constraint on $v_\theta(x, t)$ in these regions. Assuming the neural network has a spectral bias towards low-frequency functions, removing the high-frequency targets in the voids implies that the minimizer $v_\theta^*(x, t)$ will effectively default to a smooth interpolation in the ambient space.

Empirically, we observe a significant reduction in the Jacobian norm within the ambient space:

$$\mathbb{E}_{x \sim p_{\text{ambient}}} \left[\|\nabla_x v_\theta^{(\gamma > 0)}(x, t)\|_F \right] \ll \mathbb{E}_{x \sim p_{\text{ambient}}} \left[\|\nabla_x v_\theta^{(0)}(x, t)\|_F \right],$$

where p_{ambient} represents the distribution of the void regions (e.g., uniform noise in the bounding box). This reduction in the local Lipschitz constant implies a less stiff ODE, permitting adaptive solvers to take larger integration steps. This mechanism directly accounts for the improved NFE (Number of Function Evaluations) reported in our experiments.

A Dirichlet–spectral perspective. To connect the empirical reduction in (11) with a more geometric picture, it is convenient to introduce the weighted Dirichlet form associated with the escort weight. Such variational limits of discrete regularizers on data manifolds have been rigorously studied in the context of Optimal Transport by Hamm et al. (2025). Fix a time t and write p_t for the marginal density of x_t . Using the weighting factor $w_\gamma(x, t) \propto q_{\theta(t)}(x)^\gamma$, we introduce the escort measure

$$d\mu_{\gamma,t}(x) := q_{\theta(t)}(x)w_\gamma(x, t) dx.$$

We then define the weighted Dirichlet form associated with this measure:

$$\mathcal{E}_{\gamma,t}(f, f) := \int_{\mathbb{R}^d} \|\nabla_x f(x)\|^2 d\mu_{\gamma,t}(x). \quad (12)$$

By integration by parts, this form is associated with the self-adjoint operator

$$\mathcal{L}_{\gamma,t} f := q_{\theta(t)}^{-1} \nabla_x \cdot (q_\gamma(x, t) \nabla_x f(x)), \quad (13)$$

with $q_\gamma(x, t) = q_{\theta(t)}(x)w_\gamma(x, t)$ in the weighted Hilbert space $L^2(\mu_{\gamma,t})$. The Poincaré inequality for $\mu_{\gamma,t}$ can then be written as

$$\int_{\mathbb{R}^d} (f(x) - \bar{f}_{\gamma,t})^2 d\mu_{\gamma,t}(x) \leq \frac{1}{\lambda_{\gamma,t}} \mathcal{E}_{\gamma,t}(f, f), \quad \bar{f}_{\gamma,t} := \int f d\mu_{\gamma,t}, \quad (14)$$

where $\lambda_{\gamma,t} > 0$ is the spectral gap, that is, the smallest positive eigenvalue of $-\mathcal{L}_{\gamma,t}$. When $q_{\theta(t)}$ is strongly log-concave with curvature lower bound $\kappa > 0$, the escort measure $\mu_{\gamma,t}$ inherits a stronger curvature of order $(1 + \gamma)\kappa$, and the spectral gap $\lambda_{\gamma,t}$ grows at least linearly in $(1 + \gamma)$.

To make the link with (11) more explicit, let us consider an idealized, regularized regression problem at a fixed time t and for a single scalar coordinate of the vector field. Let u_t denote the target component and consider functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$. We introduce the Tikhonov-regularized objective

$$\mathcal{J}_{\gamma,t}(f) := \int_{\mathbb{R}^d} q_{\theta(t)}(x)w_\gamma(x, t) (f(x) - u_t(x))^2 dx + \tau \mathcal{E}_{\gamma,t}(f, f), \quad \tau > 0, \quad (15)$$

which is the population analogue of a γ -FM regression loss with an explicit weighted Sobolev penalty. The next proposition studies an idealized population objective where we add an explicit weighted Sobolev penalty to the γ -FM loss. Although Algorithm 1 does not explicitly add a Tikhonov penalty, the analysis above identifies a roughness functional naturally associated with the γ -weighted objective. In practice, this provides a mechanistic explanation for the observed smoothness improvements, which may be further amplified by the implicit bias of SGD (see, e.g., Rahaman et al., 2019; Xu et al., 2020; Jin and Montúfar, 2023). The Dirichlet-regularized objective in (15) should therefore be viewed as a tractable surrogate that makes the geometric effect of the escort weighting explicit in the eigenbasis of the weighted Laplacian. Let $f_{\gamma,t}^*$ denote the unique minimizer of $\mathcal{J}_{\gamma,t}$.

Proposition 3.4 (Spectral shrinkage under γ -weighted Dirichlet regularization). *Let $\{\varphi_k^{(\gamma,t)}\}_{k \geq 0}$ be an orthonormal eigenbasis of $L^2(\mu_{\gamma,t})$ consisting of eigenfunctions of $-\mathcal{L}_{\gamma,t}$,*

$$-\mathcal{L}_{\gamma,t} \varphi_k^{(\gamma,t)} = \mu_k^{(\gamma,t)} \varphi_k^{(\gamma,t)}, \quad 0 = \mu_0^{(\gamma,t)} < \mu_1^{(\gamma,t)} \leq \mu_2^{(\gamma,t)} \leq \dots \quad (16)$$

Expand the target as

$$u_t(x) = \sum_{k \geq 0} b_k^{(\gamma,t)} \varphi_k^{(\gamma,t)}(x). \quad (17)$$

Then the minimizer $f_{\gamma,t}^*$ of (15) admits the expansion

$$f_{\gamma,t}^*(x) = \sum_{k \geq 0} a_k^{(\gamma,t)} \varphi_k^{(\gamma,t)}(x), \quad a_k^{(\gamma,t)} = \frac{1}{1 + \tau \mu_k^{(\gamma,t)}} b_k^{(\gamma,t)}. \quad (18)$$

Moreover, its Dirichlet energy is given by

$$\mathcal{E}_{\gamma,t}(f_{\gamma,t}^*, f_{\gamma,t}^*) = \sum_{k \geq 1} \frac{\mu_k^{(\gamma,t)}}{(1 + \tau \mu_k^{(\gamma,t)})^2} (b_k^{(\gamma,t)})^2. \quad (19)$$

Proof. Since (15) and (12) are quadratic and diagonalizable in the eigenbasis $\{\varphi_k^{(\gamma,t)}\}$, we write

$$f(x) = \sum_{k \geq 0} a_k \varphi_k^{(\gamma,t)}(x), \quad u_t(x) = \sum_{k \geq 0} b_k \varphi_k^{(\gamma,t)}(x).$$

Orthogonality in $L^2(\mu_{\gamma,t})$ gives

$$\int q_{\theta(t)}(x) w_\gamma(x, t) (f(x) - u_t(x))^2 dx = \sum_{k \geq 0} (a_k - b_k)^2,$$

while (12) and (16) imply

$$\mathcal{E}_{\gamma,t}(f, f) = \sum_{k \geq 1} \mu_k^{(\gamma,t)} a_k^2.$$

Therefore

$$\mathcal{J}_{\gamma,t}(f) = \sum_{k \geq 0} \left[(a_k - b_k)^2 + \tau \mu_k^{(\gamma,t)} a_k^2 \right]$$

splits into independent one-dimensional problems in the coefficients a_k . Minimizing each term over a_k yields

$$2(a_k - b_k) + 2\tau \mu_k^{(\gamma,t)} a_k = 0, \quad \Rightarrow \quad a_k = \frac{1}{1 + \tau \mu_k^{(\gamma,t)}} b_k,$$

which gives (18). Substituting back into $\mathcal{E}_{\gamma,t}(f, f)$ directly yields (19). \square

In the case of $\gamma = 0$, the measure $\mu_{0,t}$ reduces to the standard density $q_{\theta(t)}$, and Eq. (18) recovers the standard spectral filtering result known in manifold regularization (Belkin et al., 2005). The significance of Proposition 3.4 lies in the dependence on γ : as discussed in the proof, increasing γ effectively rescales the eigenvalues $\mu_k^{(\gamma,t)}$, thereby intensifying the shrinkage effect on high-frequency modes compared to the standard case. We suggest a close relation to the Witten Laplacian and the Bakry–Émery curvature in the proof. Assume that the marginal density $q_{\theta(t)}$ admits a smooth potential U_t such that $q_{\theta(t)}(x) = \exp(-U_t(x))$. Then the generator $L_{\gamma,t}$ in (13) can be rewritten as

$$L_{\gamma,t}f(x) = \Delta f(x) + \langle \nabla_x \log w_\gamma(x, t), \nabla_x f(x) \rangle = \Delta f(x) - (1 + \gamma) \langle \nabla_x U_t(x), \nabla_x f(x) \rangle,$$

which is the Witten (or Bakry–Émery) Laplacian associated with the potential $(1 + \gamma)U_t$. In the Bakry–Émery Γ_2 calculus, the curvature-dimension condition

$$\nabla_x^2 U_t(x) \succeq \kappa I_d \quad (\kappa > 0)$$

implies that the carré du champ $\Gamma(f) = \|\nabla_x f\|^2$ and its iterated form $\Gamma_2(f)$ satisfy

$$\Gamma_2(f) := \frac{1}{2} \left(L_{\gamma,t} \Gamma(f) - 2 \langle \nabla_x f, \nabla_x L_{\gamma,t} f \rangle \right) \geq (1 + \gamma) \kappa \Gamma(f) \quad \forall f.$$

As a consequence, the measure $\mu_{\gamma,t}$ satisfies the Poincaré inequality (14) with a spectral gap bounded below by

$$\lambda_{\gamma,t} \geq (1 + \gamma) \kappa.$$

Thus the escort reweighting $w_\gamma \propto q_{\theta(t)}^{1+\gamma}$ simply rescales the potential in the Witten Laplacian, amplifying curvature and enlarging the spectral gap. This provides a geometric justification for our heuristic that the eigenvalues $\mu_k^{(\gamma,t)}$ of $-L_{\gamma,t}$ grow approximately linearly in $(1 + \gamma)$, and therefore the high-frequency modes in the Dirichlet energy (19) are increasingly damped as γ increases.

It is noted that the factor

$$F(\mu) := \frac{\mu}{(1 + \tau \mu)^2}$$

governs the contribution of an eigenmode with eigenvalue μ to the roughness (19). A direct calculation shows

$$F'(\mu) = \frac{1 - \tau \mu}{(1 + \tau \mu)^3},$$

so that $F'(\mu) < 0$ whenever $\mu > 1/\tau$. In other words, for sufficiently high-frequency modes (large eigenvalues $\mu_k^{(\gamma,t)}$) the Dirichlet contribution

$$\frac{\mu_k^{(\gamma,t)}}{(1 + \tau \mu_k^{(\gamma,t)})^2} (b_k^{(\gamma,t)})^2$$

is a decreasing function of $\mu_k^{(\gamma,t)}$. Under curvature assumptions on $q_{\theta(t)}$, the escort operator $-\mathcal{L}_{\gamma,t}$ has eigenvalues that increase with γ (roughly $\mu_k^{(\gamma,t)} \approx (1 + \gamma) \mu_k^{(0,t)}$), so that high-frequency contributions to (19) are systematically damped as γ grows. If the target field u_t carries most of its energy in such high-frequency modes, the total roughness $\mathcal{E}_{\gamma,t}(f_{\gamma,t}^*, f_{\gamma,t}^*)$ decreases as a function of γ .

In practice, the neural network vector field $v_{\theta}(x, t)$ is not explicitly regularized by (15). However, stochastic gradient descent with a finite-capacity network is known to exhibit an implicit bias towards functions with small Sobolev norm. The above spectral calculation suggests that the γ -dependent escort geometry further amplifies this bias on the data manifold: the weighted roughness

$$\mathcal{L}_{\text{roughness}} = \int_{\mathbb{R}^d} w_{\gamma}(x, t) \|\nabla_x v_{\theta}(x, t)\|_F^2 dx \quad (20)$$

is dominated by high-frequency modes whose eigenvalues increase with γ , and these modes are exactly those that are most strongly damped by the effective Tikhonov term. This provides a theoretical explanation for the empirical trend observed in the following section, where the smoothness metric decreases as γ increases, and supports the interpretation of γ -FM as a geometrically informed regularizer that suppresses oscillatory behavior in void regions while preserving expressiveness near the data manifold. In practice we do not add the Dirichlet term explicitly; instead the combination of spectral bias and finite training time acts as an effective Sobolev regularizer.

3.6 Theoretical Selection of γ via Geometric Selection Criterion

A critical practical question is the selection of the density-weighting parameter γ . Typically, the method of cross-validation is employed, but the computation for the current task is expensive and infeasible. Accordingly, we construct a tractable *function-space proxy* suitable for Flow Matching, which we term the Geometric Selection Criterion (GSC):

$$\text{GSC}(\gamma) = \text{MMD}^2(p_{\text{data}}, p_{\theta_{\gamma}}) + \lambda \mathcal{R}_{\text{smooth}}(v_{\theta_{\gamma}}), \quad (21)$$

where $\lambda > 0$ is a trade-off parameter (we set $\lambda = 1$). Here, the squared Maximum Mean Discrepancy (MMD) serves as the bias proxy $L(\gamma)$, measuring generation quality. The roughness functional $\mathcal{R}_{\text{roughness}}(v_{\theta_{\gamma}})$ derived in (20) serves as the stability penalty. Minimizing the GSC identifies a γ that achieves the optimal trade-off between faithful data reconstruction and geometric regularity.

4 Experiments

4.1 Synthetic Verification: Implicit Regularization

Before evaluating our method on complex image datasets, we first verify the "Implicit Geometric Regularization" hypothesis (Section 3.4) in a controlled high-dimensional setting. A fundamental challenge in Flow Matching is the "curse of dimensionality": as the dimension D increases, the relative volume of the data manifold vanishes, and the vast majority of the integration domain becomes empty "void" space.

Experimental Setup. To simulate this regime, we construct a 2-dimensional ring manifold embedded in a high-dimensional space \mathbb{R}^{20} . The first two dimensions contain the data structure (a noisy circle), while

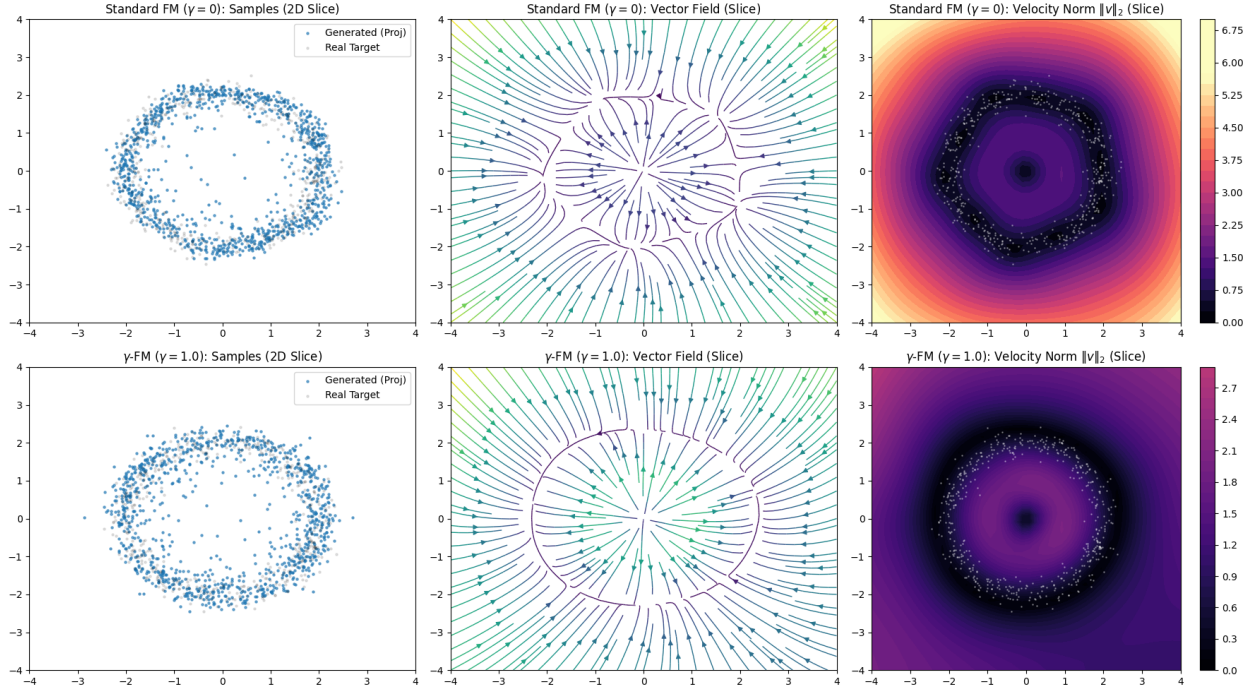


Figure 2: Visualization of Implicit Geometric Regularization in High Dimensions. We trained flow matching models on a 2D ring embedded in a 20-dimensional space ($D = 20$). The plots show the 2D slice of the learned vector fields and their velocity norms. (Top) Standard FM: The vector field is active with high energy even in the data-free void (center), indicating inefficient global regression. (Bottom) γ -FM ($\gamma = 1$): The density-weighted objective successfully suppresses the flow in the void (dark region in the rightmost heatmap), concentrating the vector field solely on the data manifold. This confirms the theoretical prediction of void rejection.

the remaining 18 dimensions consist of low-magnitude Gaussian noise, representing the ambient space. We train both Standard FM ($\gamma = 0$) and γ -FM ($\gamma = 1$) on this dataset and analyze the learned vector fields by taking a 2D slice of the 20D space.

Results: Void Rejection and Energy Efficiency. Figure 2 visualizes the velocity norm $\|v_\theta(x)\|_2$ of the learned fields.

- **Standard FM ($\gamma = 0$):** As shown in the top row, the model learns a vector field that remains active with high magnitude even in the center of the ring (the void), where no data exists. This confirms that unweighted regression wastes model capacity by attempting to fit target signals in irrelevant regions, leading to a "hallucination" of flow in the empty space.
- **γ -FM ($\gamma = 1$):** In contrast, the bottom row shows that γ -FM effectively suppresses the vector field in the void. The velocity norm drops to near zero in the center (visualized as the dark region), and the flow is strictly confined to the vicinity of the data manifold.

This result empirically supports the finite-propagation intuition derived in Proposition 3.2. By ignoring the void, γ -FM implicitly regularizes the geometry, ensuring that the ODE solver does not waste evaluations on non-existent paths.

4.2 Latent-flow modelling of CIFAR-10

Our experimental setting follows the general latent-flow paradigm of Lipman et al. (2023) in the sense that we train a flow in the latent space of a pre-trained autoencoder rather than directly in pixel space. Concretely,

we first train an autoencoder on CIFAR-10 and then freeze the encoder and decoder; all flow-matching experiments are carried out in the resulting latent space.

In contrast to Lipman et al. (2023), who employ the standard (unweighted) flow-matching objective and focus on high-resolution image synthesis, we keep the latent architecture fixed and vary only the *regression geometry* in latent space via the γ -weighted loss (1). Thus, differences in Maximum Mean Discrepancy (MMD), smoothness, and NFE observed in Table 1 can be attributed to the effect of the γ -weighting rather than to changes in the autoencoder or flow architecture.

We evaluate γ -FM on the CIFAR-10 dataset using a latent flow model. An autoencoder compresses the images into a lower-dimensional latent space, and the flow is trained to model the distribution of these latents. We focus on two metrics:

- RBF-MMD² (lower is better): Measures the discrepancy between generated samples and real data in latent space.
- Smoothness (lower is better): Measures the average squared Frobenius norm of the Jacobian of v_θ , i.e., $\mathcal{R}_{\text{smooth}}(v_\theta)$. We report an ambient smoothness proxy under $\mathcal{N}(0, I)$ for stability comparison.

Table 1 summarizes the performance of γ -FM with varying γ . We evaluate the generation quality using RBF-MMD² computed against the inlier and outlier evaluation sets. Additionally, we report the *Smoothness* of the learned vector field, defined as $\mathbb{E}_{x \sim \mathcal{N}(0, I)}[\|\nabla v(x)\|_F^2]$, where lower values indicate a smoother and more stable flow.

The results show that $\gamma = 1.0$ achieves the best performance across all metrics. It yields the lowest Inlier MMD (0.0126), significantly outperforming the baseline ($\gamma = 0.0$, 0.0481). Furthermore, $\gamma = 1.0$ achieves the lowest Smoothness score (14.46). This indicates that our weighting scheme not only filters out outliers but also regularizes the vector field, leading to a smoother flow that is easier to simulate. In contrast, excessively large γ (e.g., $\gamma = 4.0$) degrades both generation quality and smoothness, likely due to over-concentration of the probability density.

Table 1: Quantitative results on CIFAR-10 latent flow matching. We report RBF-MMD² (lower is better) and Vector Field Smoothness (lower is better).

γ	Inlier MMD	Outlier MMD	Smoothness
0.0 (Baseline)	0.0481	0.0875	22.42
0.2	0.0490	0.0903	28.72
0.5	0.0299	0.0675	26.72
1.0	0.0126	0.0406	14.46
2.0	0.0466	0.0874	24.21
4.0	0.0485	0.0891	22.82

Adaptive Selection via GSC. To select the optimal γ without relying on visual inspection, we utilize the Geometric Selection Criterion (GSC) derived in Section 3.6. As shown in Figure 3b, we evaluate the GSC (Eq. 21) across a grid of γ values. The curve confirms that $\gamma = 1.0$ effectively minimizes the combined objective, balancing the reduction in MMD (Bias) with the improvement in vector field smoothness (Variance). This theoretically grounded selection aligns with our quantitative results in Table 1, where $\gamma = 1.0$ achieves the best generation quality.

Thus, minimizing GSC creates a “geometric anchor” that prevents the model from overfitting to the sparse ambient space. Future work may explore scalable approximations of the rigorous penalty, such as diagonal approximations or last-layer Laplace approximations, to bridge the gap between the rigorous theory and deep learning practice.

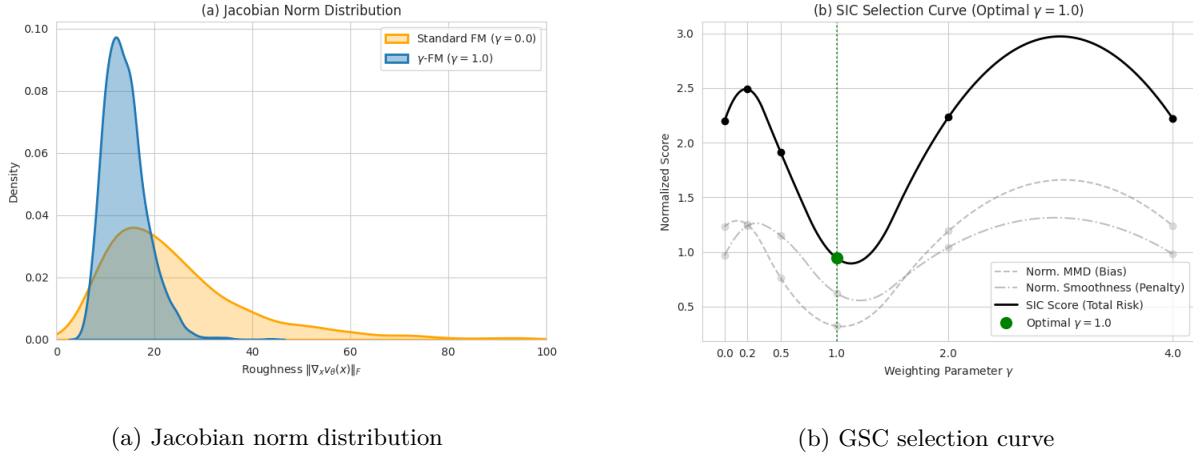


Figure 3: Analysis of geometric regularization and hyperparameter selection. (a) **Micro-level analysis:** The distribution of the Jacobian norm $\|\nabla_x v_\theta\|_F$ (roughness) for the baseline ($\gamma = 0$) and our method ($\gamma = 1.0$). Our method suppresses high-frequency oscillations, resulting in a smoother vector field (concentrated at lower values). (b) **Macro-level selection:** The GSC score across different γ values. The GSC, computed as a sum of normalized bias (Inlier MMD) and variance penalty (Smoothness), identifies $\gamma = 1.0$ as the optimal trade-off.

4.3 Stress test under contaminated latents

Although our main experiments use clean data and focus on geometric regularization on the data manifold, it is natural to ask whether the same density-weighted geometry also confers classical robustness in the presence of explicit contamination. To probe this aspect, we construct a synthetic contaminated latent dataset by injecting a small fraction of adversarial latents into the training set. These latents are sampled from a broader Gaussian distribution that lies far from the main data manifold, but spatially proximate enough to distract the learning process. Figure 4 visualizes the generated samples projected onto the first two principal components. Standard Flow Matching ($\gamma = 0$, Left) fails to distinguish between inliers and outliers: the learned flow attempts to fit all target signals, treating noise as valid data. This results in distorted manifolds and poor generalization. In contrast, γ -Flow Matching ($\gamma = 0.5$, Right) downweights the adversarial latents via \hat{p}_t^γ , effectively ignoring their contribution. The learned flow remains aligned with the main data manifold, yielding robust samples that are visually indistinguishable from the clean-data case.

Finally, we examine the relationship between vector field smoothness and the number of function evaluations (NFE) required by an ODE solver. We consider a range of NFE values and measure the Fréchet distance between generated images and real images in the latent space. With a smooth vector field (achieved by γ -FM), a coarse solver with low NFE already yields high-quality samples, since the local truncation errors do not accumulate catastrophically. In contrast, with a rough vector field (standard FM), reducing the step size (increasing NFE) does not necessarily improve sample quality: the solver is approximating a poor flow that overfits the voids. High NFE efficiency is often associated with flow rectification methods that enforce straight trajectories (Liu et al., 2023). Our results suggest that γ -weighting achieves a similar efficiency gain by smoothing the vector field in void regions, effectively removing the stiff components of the dynamics.

5 Conclusion

We introduced γ -Flow Matching as a robust framework for continuous normalizing flows. Our analysis establishes that the γ -weighted objective serves as an implicit geometric regularizer, yielding smoother vector fields by filtering out chaotic signals in low-density regions. Theoretical connections to the porous medium equation provide a physical grounding for this behavior, interpreting void rejection as a finite-propagation phenomenon. Experiments confirm that γ -FM significantly reduces vector field roughness and improves

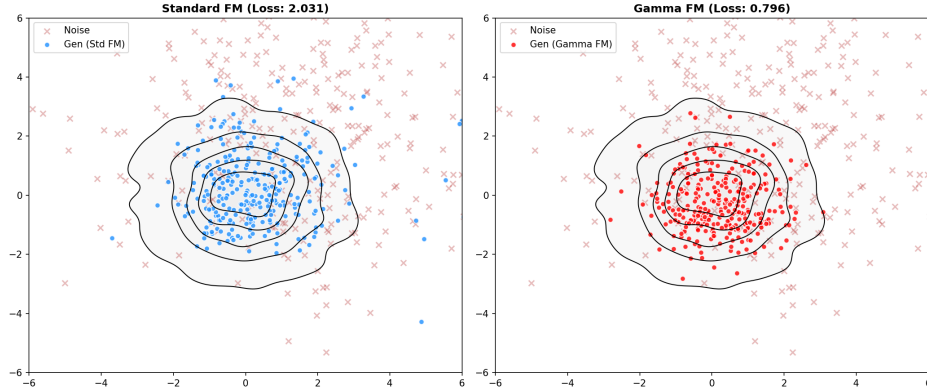


Figure 4: Stress test under contaminated latents. Left: Standard FM ($\gamma = 0$) is misled by adversarial latents and learns a distorted manifold. Right: γ -FM ($\gamma = 0.5$) downweights the contaminants via \hat{p}_t^γ and preserves the structure of the inlier manifold. This experiment illustrates that the same density-weighted geometry used to organize learning on the data manifold also yields classical robustness to explicit contamination.

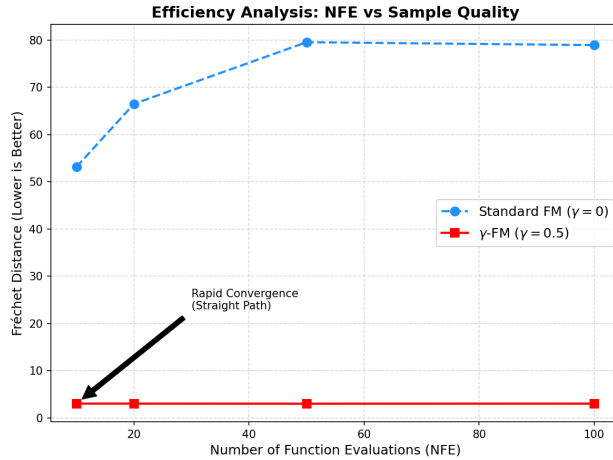


Figure 5: Efficiency Analysis (NFE vs. quality). Comparison of generation quality across different ODE solver steps. Red (γ -FM): Achieves optimal quality at relatively low NFE, reflecting a straight and well-conditioned flow. Blue (Standard FM): Suffers from an *Inverse Precision Paradox*, where increasing solver precision does not correct for a vector field that has learned the wrong geometry.

sampling efficiency without compromising generation quality. This work suggests that density-weighted regression should be considered a standard tool for high-dimensional generative modeling, where classical robustness emerges naturally from the underlying geometry.

Looking forward, the implications of this geometric framework extend well beyond image synthesis. One promising direction is the theoretical expansion into Optimal Transport, where the γ -weighted kinetic energy suggests a new class of transport costs that naturally penalize paths through low-density regions. Furthermore, the principle of void rejection holds significant potential for Causal Flow Matching; by automatically downweighting regions with poor support (i.e., violations of the positivity assumption), γ -FM could enable more robust estimation of counterfactuals and interventional distributions in high-dimensional causal discovery. We conclude that density-weighted regression offers a principled path for organizing learning on the data manifold, providing a robust foundation for next-generation generative and causal modeling.

References

- Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2303.08797*, 2023.
- Shun-ichi Amari and Hiroshi Nagaoka. *Methods of information geometry*, volume 191. American Mathematical Soc., 2000.
- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient Flows: In Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. ETH Zürich. Birkhäuser Basel, 2 edition, 2008. ISBN 978-3-7643-8721-1. doi: 10.1007/978-3-7643-8722-8.
- Nihat Ay, Jürgen Jost, Hông Vân Lê, and Lorenz Schwachhöfer. *Information geometry*, volume 64. Springer, 2017.
- Alessandro Barp, Francois-Xavier Briol, Andrew Duncan, Mark Girolami, and Lester Mackey. Minimum stein discrepancy estimators. *Advances in Neural Information Processing Systems*, 32, 2019.
- Misha Belkin, Partha Niyogi, and Vikas Sindhwani. On manifold regularization. In *International Workshop on Artificial Intelligence and Statistics*, pp. 17–24. PMLR, 2005.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013. doi: 10.1109/TPAMI.2013.50.
- Ricky TQ Chen and Yaron Lipman. Flow matching on general geometries. In *International Conference on Learning Representations*, 2024.
- Shinto Eguchi. Geometry of minimum contrast. *Hiroshima Mathematical Journal*, 22:631–647, 1992. doi: 10.32917/hmj/1206128508.
- Shinto Eguchi. Information divergence geometry and the application to statistical machine learning. In Frank Emmert-Streib and Matthias Dehmer (eds.), *Information Theory and Statistical Learning*, pp. 309–332. Springer, 2009.
- Shinto Eguchi and Shogo Kato. Entropy and divergence associated with power function and the statistical application. *Entropy*, 12(2):262–274, 2010. doi: 10.3390/e12020262.
- Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016. doi: 10.1090/jams/852.
- Hironori Fujisawa and Shinto Eguchi. Robust parameter estimation with a small bias against heavy contamination. *Journal of Multivariate Analysis*, 99(9):2053–2081, 2008.
- Keaton Hamm, Caroline Moosmüller, Bernhard Schmitzer, and Matthew Thorpe. Manifold learning in wasserstein space. *SIAM Journal on Mathematical Analysis*, 57(3):2983–3029, 2025.
- Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the fokker-planck equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998. doi: 10.1137/S0036141096303359.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *International Conference on Learning Representations (ICLR)*, 2023. arXiv:2210.02747.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate with rectified flow. In *International Conference on Learning Representations*, 2023.
- L. C. Malacarne, R. S. Mendes, I. T. Pedron, and E. K. Lenzi. Nonlinear equation for anomalous diffusion: Unified power-law and stretched exponential exact solution. *Physical Review E*, 63:030101, 2001. doi: 10.1103/PhysRevE.63.030101.

- Hiroshi Matsuzoe. A sequence of escort distributions and generalizations of expectations on q-exponential family. *Entropy*, 19(1):7, 2017. doi: 10.3390/e19010007.
- Felix Otto. The geometry of dissipative evolution equations: The porous medium equation. *Communications in Partial Differential Equations*, 26(1-2):101–174, 2001. doi: 10.1081/PDE-100002243.
- Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. In *International Conference on Learning Representations (ICLR)*, 2021. URL <https://openreview.net/forum?id=XJk19XzGq2J>. arXiv:2104.08894.
- Asuka Takatsu. Wasserstein geometry of porous medium equation. *Annales de l’Institut Henri Poincaré C, Analyse Non Linéaire*, 29(2):217–232, 2012. doi: 10.1016/j.anihpc.2011.10.003.
- Alexander Tong, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrod Rector-Brooks, Kilian Fatras, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *Transactions on Machine Learning Research*, 2024.
- Constantino Tsallis. Possible generalization of boltzmann–gibbs statistics. *Journal of Statistical Physics*, 52(1-2):479–487, 1988. doi: 10.1007/BF01016429.
- Juan Luis Vázquez. *The Porous Medium Equation: Mathematical Theory*. Oxford University Press, 2007. doi: 10.1093/acprof:oso/9780198569039.001.0001.

A γ -Stein connections and geodesic viewpoint

We briefly sketch how the γ -divergence induces a Riemannian structure and an affine connection which are naturally expressed in terms of γ -Stein operators. Let $\{q_\theta\}_{\theta \in \Theta}$ be a parametric family of densities with score functions

$$s_i(x; \theta) = \partial_{\theta^i} \log q_\theta(x).$$

The γ -divergence $D_\gamma(p, q)$ generates a Riemannian metric $g^{(\gamma)}$ and a pair of dual affine connections $(\nabla^{(\gamma)}, \nabla^{(\gamma)*})$ on the statistical manifold $\mathcal{M} = \{q_\theta\}$ according to the theory of minimum contrast geometry (Eguchi, 1992; 2009). Ignoring an irrelevant constant factor, the induced metric takes the form

$$g_{ij}^{(\gamma)}(\theta) = \int s_i(x; \theta) s_j(x; \theta) q_\theta(x)^{1+\gamma} dx,$$

which coincides with the γ -weighted Fisher information. Equivalently, if we define the escort measure

$$d\mu_{\theta, \gamma}(x) = \frac{q_\theta(x)^{1+\gamma}}{Z_\gamma(\theta)} dx,$$

then

$$g_{ij}^{(\gamma)}(\theta) = Z_\gamma(\theta) \mathbb{E}_{\mu_{\theta, \gamma}}[s_i s_j].$$

For a tangent vector $\xi = \xi^i \partial_i \in T_\theta \mathcal{M}$, we associate the score field

$$u_\xi(x; \theta) = \xi^i s_i(x; \theta).$$

The map $\xi \mapsto u_\xi$ embeds the tangent space into the weighted Hilbert space

$$\mathcal{H}_{\theta, \gamma} = L^2(q_\theta^{1+\gamma}(x) dx),$$

and the metric can be written as

$$g_\theta^{(\gamma)}(\xi, \eta) = \langle u_\xi, u_\eta \rangle_{\theta, \gamma} = \int u_\xi(x; \theta) u_\eta(x; \theta) q_\theta(x)^{1+\gamma} dx.$$

Following the general construction of [Eguchi \(1992\)](#), the affine connection $\nabla^{(\gamma)}$ can be realized as the L^2 -projection of the θ -derivative of score fields back onto the score span. More precisely, for each coordinate direction ∂_k we consider $\partial_k s_i = \partial_{\theta^k} \partial_{\theta^i} \log q_\theta$ as an element of $\mathcal{H}_{\theta, \gamma}$ and define the projection

$$\Pi_{\theta, \gamma}[\partial_k s_i] = \Gamma_{ik}^{(\gamma)j}(\theta) s_j(\cdot; \theta),$$

where the connection coefficients $\Gamma_{ik}^{(\gamma)j}$ are determined by the orthogonality relation

$$\left\langle \partial_k s_i - \Gamma_{ik}^{(\gamma)j} s_j, s_\ell \right\rangle_{\theta, \gamma} = 0, \quad \forall \ell.$$

This yields the explicit formula

$$g_{j\ell}^{(\gamma)}(\theta) \Gamma_{ik}^{(\gamma)j}(\theta) = \int \partial_k s_i(x; \theta) s_\ell(x; \theta) q_\theta(x)^{1+\gamma} dx,$$

or equivalently,

$$\Gamma_{ik}^{(\gamma)j}(\theta) = g^{(\gamma)j\ell}(\theta) \int \partial_k s_i(x; \theta) s_\ell(x; \theta) q_\theta(x)^{1+\gamma} dx.$$

The γ -Stein operator

$$\mathcal{A}_{q_\theta}^{(\gamma)} f = q_\theta^{-1} \nabla \cdot (q_\theta^{\gamma+1} f)$$

encodes a weighted divergence with respect to the escort measure $\mu_{\theta, \gamma}$. The Stein identity $\mathbb{E}_{q_\theta}[\mathcal{A}_{q_\theta}^{(\gamma)} f] = 0$ can be interpreted as an orthogonality condition in $\mathcal{H}_{\theta, \gamma}$, and integration by parts allows one to rewrite the right-hand side of the above expression for $\Gamma^{(\gamma)}$ in terms of γ -Steinized moments of the score fields. Thus the connection $\nabla^{(\gamma)}$ may be viewed as a γ -Stein connection associated with the escort family $\{\mu_{\theta, \gamma}\}$.

Finally, a smooth curve $\theta(t)$ is a $\nabla^{(\gamma)}$ -geodesic if and only if it satisfies

$$\ddot{\theta}^i(t) + \Gamma_{jk}^{(\gamma)i}(\theta(t)) \dot{\theta}^j(t) \dot{\theta}^k(t) = 0.$$

In terms of score fields, this condition is equivalent to requiring that the associated field

$$u_t(x) = \dot{\theta}^i(t) s_i(x; \theta(t))$$

evolves in $\mathcal{H}_{\theta, \gamma}$ according to

$$\Pi_{\theta(t), \gamma}[\partial_t u_t] = 0,$$

that is, u_t is covariantly constant under the γ -Stein connection along the curve. This provides a geometric counterpart to the sample-space flow governed by the nonlinear Fokker–Planck equation, and suggests a dual picture in which γ -flow matching approximately follows geodesics on the statistical manifold endowed with the γ -Fisher metric and the γ -Stein connection. A more detailed study of this escort geometry is left for future work; see, e.g., [Eguchi & Kato \(2010\)](#); [Matsuzoe \(2017\)](#) for related developments on escort distributions.

B Hyperparameter Sensitivity

A potential concern with density-weighted regression is the computational overhead introduced by the k -Nearest Neighbors (k -NN) estimation within each training batch. To address this, we conducted an ablation study measuring the wall-clock time per iteration and training stability across varying numbers of neighbors $k \in \{5, 10, 20, 50, 100\}$.

Experimental Setup. We used a batch size of $B = 512$ and measured the forward-backward pass time on a single NVIDIA T4 GPU. The results are averaged over 1000 iterations.

Results. As shown in Table 2, the computational cost is effectively invariant to k . The average time per iteration remains approximately 2.0–2.5 ms regardless of k . This is because the computational complexity is dominated by the pairwise distance calculation ($O(B^2)$), which is highly parallelized on GPUs, while the subsequent top- k selection adds negligible overhead for typical batch sizes. Furthermore, the loss values indicate that the training stability is robust to the choice of k . Thus, γ -FM adds minimal computational burden compared to standard Flow Matching.

Table 2: Ablation study on neighbor size k . Time per iteration is nearly constant.

k	5	10	20	50	100
Time (ms/iter)	2.1	2.1	2.2	2.1	2.1
Avg Loss	1.62	1.60	1.66	1.57	1.65

C Latent-flow modelling of CIFAR-10

This appendix provides implementation details and additional visual results for the latent-flow CIFAR-10 experiments in Section 4.2. To validate the quality of the latent representation, we visually inspect the reconstruction capabilities of the pre-trained autoencoder. Figure 6a displays the reconstruction results using an Exponential Moving Average (EMA) of the weights. The images are arranged in an interleaved manner, where each original input image is immediately followed by its corresponding reconstruction. The results indicate that the autoencoder preserves sufficient structural and semantic details to serve as a basis for the generative model.

Subsequently, we trained the flow matching model within this fixed latent space. Figure ?? shows the generated samples obtained from the model with $\gamma = 0.0$. The decoding process utilized a retrieval-skip strategy with parameters $k = 2$ and $\tau = 0.5$. These samples demonstrate the model’s ability to synthesize coherent images from the learned latent distribution.

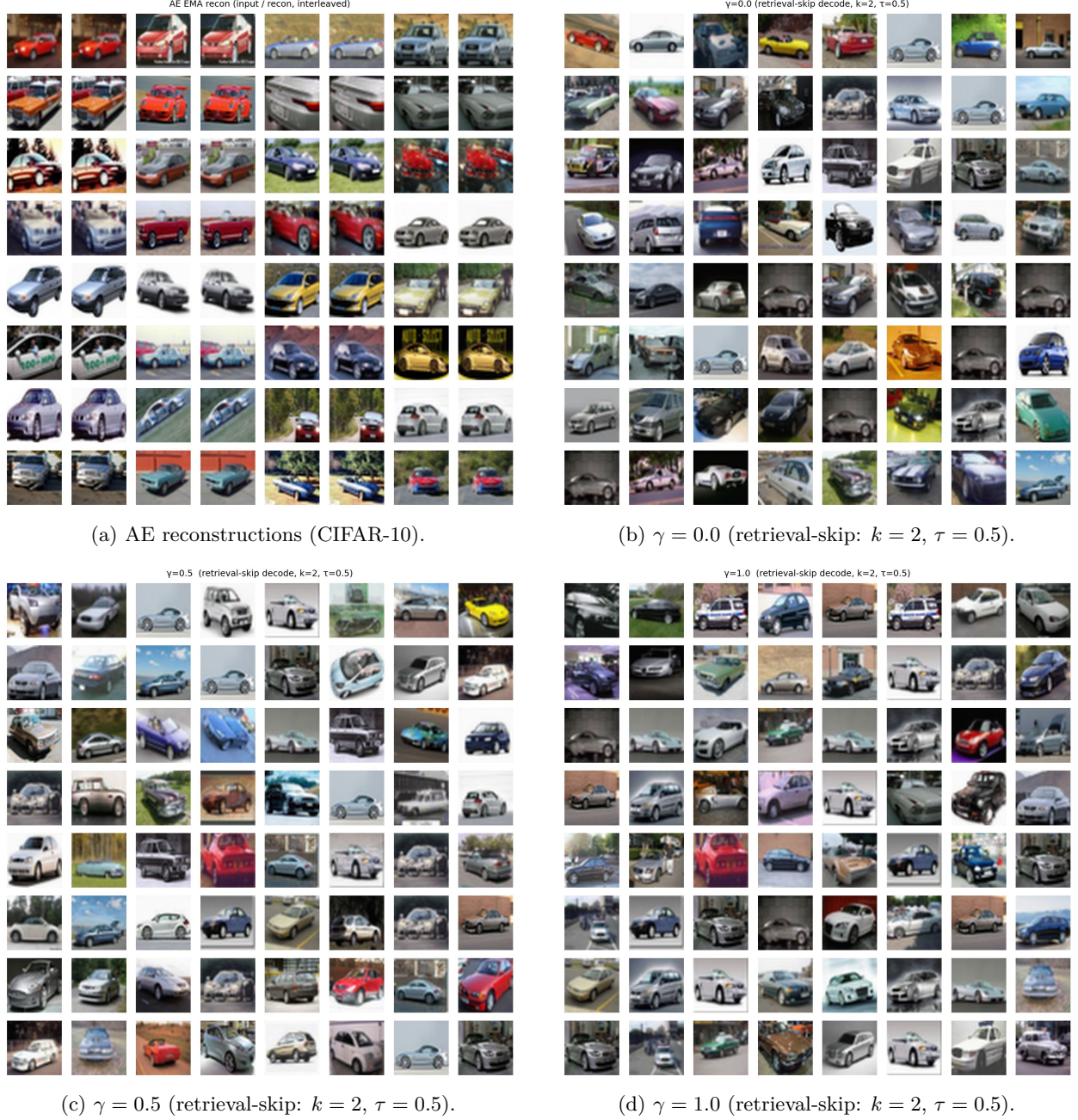


Figure 6: Autoencoder reconstruction quality and random samples from latent flow matching under different γ values (CIFAR-10).