

Unifying GANs and Score-Based Diffusion as Generative Particle Models

Jean-Yves Franceschi,¹ Mike Gartrell,¹ Ludovic Dos Santos,^{1,*} Thibaut Issenhuth,^{1,2,*} Emmanuel de Bézenac,^{3,*} Mickaël Chen,^{4,*} Alain Rakotomamonjy^{1,*}
¹Criteo AI Lab, Paris, France ²LIGM, Ecole des Ponts, Univ Gustave Eiffel, CNRS, Marne-la-Vallée, France ³SAM, D-MATH, ETH Zürich, Zürich-8092, Switzerland ⁴Valeo.ai, Paris, France

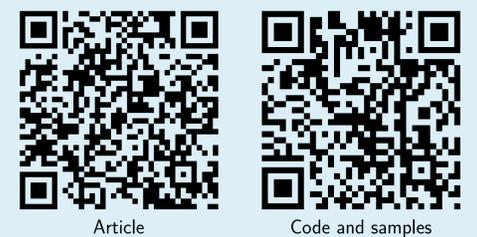
Our Contributions

- We unify gradient flows, score-based diffusion models, and GANs by representing generated data as moving particles.
- A model is defined by:
 - a gradient vector field that the particles follow;
 - the possibility of incorporating a generator smoothing this movement.
- This suggests the existence of hybrid models:
 - a generator trained with diffusion guidance (Score GANs);
 - a GAN trained without a generator (Discriminator Flows).

GANs vs Diffusion

- Traditional opposition in the literature.
- GANs** → Generator trained by discriminating true vs fake data.
- Generator (manifold learning / hypothesis).
 - Close to SOTA performance.
 - Harder to optimize.
 - Fast inference.
- Diffusion** → Learns to progressively reverse a data degradation process.
- No generator (operates on the data space).
 - SOTA performance.
 - Easier to optimize.
 - Slow inference.

Links



Particle-Based Framework

Generated particles $x_t \sim \rho_t$ follow a gradient vector field ∇h_{ρ_t} , i.e. optimize an objective h_{ρ_t} .

Wasserstein Gradient

$$-\nabla_W \mathcal{F}(\rho_t) = -\nabla \frac{\partial \mathcal{F}(\rho_t)}{\partial \rho_t}$$

Log Ratio Gradient

$$\alpha_t \nabla \log [p_{\text{data}} \star k_{\text{RBF}}^{\sigma(t)}] - \beta_t \nabla \log \rho_t$$

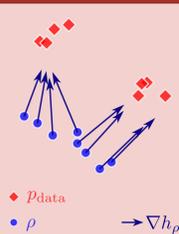
Discriminator Gradient

$$-\nabla (c \circ f_{\rho_t})$$

where f_{ρ_t} discriminates ρ_t from p_{data}

Particle Models (No Generator)

- At generation / inference time t :
 $x_0 \sim \pi = \rho_0$, $dx_t = \nabla h_{\rho_t}(x_t) dt$.
- Independently moving particles.
- Each x_t individually follows a **gradient ascent** path on $h_{\rho_t}(x_t)$.
- h_{ρ_t} is usually a predefined functional approximated with neural networks.

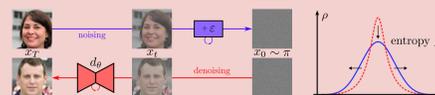


Wasserstein Gradient Flows

- Gradient descent for functionals over distributions \mathcal{F} (Santambrogio, 2017).

	Objective $\mathcal{F}(\rho)$	h_{ρ}
Forward KL	$\mathbb{E}_{\rho} \log \rho / p_{\text{data}}$	$-\log \rho / p_{\text{data}}$
f -divergence	$\mathbb{E}_{p_{\text{data}}} f(\rho / p_{\text{data}})$	$-f'(\rho / p_{\text{data}})$
Squared MMD w.r.t. kernel k	$\mathbb{E}_{x, x' \sim \rho} [k(x, x')] + k(y, y')$	$\mathbb{E}_{y \sim p_{\text{data}}} [k(y, \cdot)] - \mathbb{E}_{x \sim \rho} [k(x, \cdot)]$
Entropy	$\mathbb{E}_{\rho} \log \rho$	$-\log \rho$

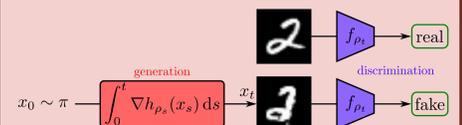
Score-Based Diffusion



- Using Jordan et al. (1998) in Song et al. (2019) and Karras et al. (2022):

$$dx_t = \alpha_t \nabla \log [p_{\text{data}} \star k_{\text{RBF}}^{\sigma(t)}](x_t) dt + \sqrt{2\beta_t} dW_t \Leftrightarrow -\beta_t \nabla \log \rho_t dt.$$

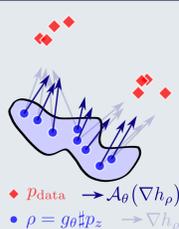
Discriminator Flows



- Particles directly follow the discriminator gradient.
- The discriminator is simultaneously trained and used to generate data.

Interacting Particle Models (Generator)

- Training with **the same loss**:
 $\mathcal{L}_{\text{gen}}(\theta) = -\mathbb{E}_{z \sim p_z} [h_{\rho_t}(g_{\theta}(z))]$.
- At training time t :
 $dg_{\theta_t}(z) = \eta [\mathcal{A}_{\theta_t}(z)] (\nabla h_{\rho_t}) dt$.
- Generalization of PMs where particles interact with each other.



Stein Gradient Flows

- Stein gradient flows (Liu, 2017) are kernelized Wasserstein gradient flows:

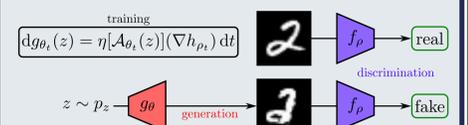
$$dx_t = \mathbb{E}_{x' \sim \rho_t} [k(x_t, x') \nabla h_{\rho_t}(x')] dt.$$

- Int-PMs under mild hypotheses (generalization of Durr et al. (2022)).
- Hint towards the same h_{ρ} being used in a PM and an Int-PM.

Score GANs

- Estimate ∇h_{ρ} with two score matching networks and use it in parameter update equation of generator training.
- Data score (pretrained like diffusion):
 $s_{\psi}^{\text{data}}(\cdot, \sigma) \equiv \nabla \log [p_{\text{data}} \star k_{\text{RBF}}^{\sigma}]$.
- Generated distribution score (continuously updated like a discriminator): $s_{\phi}^{\rho} \equiv \nabla \log \rho$.
- $\widehat{\nabla} h_{\rho}(x, \sigma) = s_{\psi}^{\text{data}}(x, \sigma) - s_{\phi}^{\rho}(x)$.

GANs



- Gradient descent-ascent on the min-max objective yields the generator loss:
 $\mathcal{L}_{\text{GAN}}(g_{\theta}) = \mathbb{E}_{z \sim p_z} [(c \circ f_{\rho})(g_{\theta}(z))]$.

Smoothing Operator

- $\mathcal{A}_{\theta_t}(z)$ is a linear operator on vector fields (kernel integral operator):

$$[\mathcal{A}_{\theta_t}(z)](V) \triangleq \mathbb{E}_{z' \sim p_z} [k_{g_{\theta_t}}(z, z') V(g_{\theta_t}(z'))],$$

$$k_{g_{\theta_t}}(z, z') \triangleq \nabla_{\theta_t} g_{\theta_t}(z')^{\top} \nabla_{\theta_t} g_{\theta_t}(z).$$

- $k_{g_{\theta_t}}$ is the generator's Neural Tangent Kernel (NTK, Jacot et al., 2018).
- Special case: $k_{g_{\theta_t}}(z, z') = \delta_{z-z'} I_d$ (generator with infinite capacity).
 - No interaction between particles: $[\mathcal{A}_{\theta_t}(z)](V) = V(g_{\theta_t}(z))$.
 - $dg_{\theta_t}(z) = \nabla h_{\rho_t}(g_{\theta_t}(z)) dt$: we retrieve PMs.
- General case: \mathcal{A}_{θ_t} represents the parameterization of ρ as a manifold.
 - \mathcal{A}_{θ_t} smooths the original vector field ∇h_{ρ_t} by convolving it with k .
 - Particles interact with each other through generator parameterization.

Score GANs in Practice

- Two practical issues:
 - sliced score matching to train s_{ϕ}^{ρ} ;
 - scheduling σ s w.r.t. training time t .
- We randomly sample σ and also noise the particles:
 $\nabla h_{\rho} = \nabla \log [p_{\text{data}} \star k_{\text{RBF}}^{\sigma}] - \nabla \log [\rho_t \star k_{\text{RBF}}^{\sigma}]$
 $\equiv \widehat{\nabla} h_{\rho}(\cdot, \sigma) = s_{\psi}^{\text{data}}(\cdot, \sigma) - s_{\phi}^{\rho}(\cdot, \sigma)$.
- Generator update:
 - few-step training of s_{ϕ}^{ρ} with denoising score matching;
 - gradient descent step:
 $\theta \leftarrow \theta + \eta \mathbb{E}_{\sigma \sim p_{\sigma}, \varepsilon \sim \mathcal{N}(0, \sigma I_D), z \sim p_z} [\nabla_{\theta} g_{\theta}(z) \widehat{\nabla} h_{\rho}(g_{\theta}(z) + \varepsilon, \sigma)]$.

Discriminator Flows in Practice

- Discriminator loss:
 $\mathcal{L}_d(f; \rho, p_{\text{data}}) = \mathbb{E}_{\rho} [a \circ f] - \mathbb{E}_{p_{\text{data}}} [b \circ f] + \mathcal{R}(f; \rho, p_{\text{data}})$.
- Naive training: successive f_{ρ_t} trainings and ρ_t updates.
- For efficiency purposes, we simultaneously learn all time-parameterized discriminators: $f_{\rho_t} = f_{\phi}(\cdot, t)$.
- Training step:
 - sample $t \sim \mathcal{U}([0, 1])$, $x_0 \sim \pi$;
 - compute $x_t = -\eta \int_0^t \nabla (c \circ f_{\phi}(\cdot, s))(x_s) ds$;
 - train $f_{\phi}(\cdot, t)$ to discriminate between x_t and p_{data} .
- Generalization of some gradient flows.

From PMs to Int-PMs

- We assign to each generated particle $x = g_{\theta}(z)$ the same loss as in PMs:
 $\mathcal{L}_{\text{gen}}(\theta) = -\mathbb{E}_{z \sim p_z} [h_{\rho_t}(g_{\theta}(z))]$.
- We do not take into account the dependency of ρ_t w.r.t. θ_t , to mimic PMs: $\rho = \text{StopGradient}(g_{\theta} \# p_z)$.
- Continuous-time gradient descent:

$$\frac{d\theta_t}{dt} = -\eta \nabla_{\theta_t} \mathcal{L}_{\text{gen}}(\theta_t) = \eta \nabla_{\theta_t} \mathbb{E}_{z \sim p_z} [h_{\rho_t}(g_{\theta_t}(z))]$$

$$= \eta \mathbb{E}_{z \sim p_z} [\nabla_{\theta_t} g_{\theta_t}(z) \nabla h_{\rho_t}(g_{\theta_t}(z))].$$

- Evolution of particles:

$$\frac{dg_{\theta_t}(z)}{dt} = \nabla_{\theta_t} g_{\theta_t}(z)^{\top} \frac{d\theta_t}{dt} = \eta \mathbb{E}_{z' \sim p_z} [\nabla_{\theta_t} g_{\theta_t}(z)^{\top} \nabla_{\theta_t} g_{\theta_t}(z') \nabla h_{\rho_t}(g_{\theta_t}(z'))].$$

Other Models & Flows

- Int-PMs and Stein (generalization of Durr et al. (2022)): $k(g_{\theta_t}(z), g_{\theta_t}(z')) = k_{g_{\theta_t}}(z, z')$ in the NTK regime.
- Langevin diffusion (Song et al., 2019) is a KL flow.
- Under some hypotheses, GANs are Stein flows (Franceschi et al., 2022; Yi et al., 2023): KL flow for f -divergence GANs, squared MMD for IPM GANs.
- As a consequence, under similar hypotheses, Discriminator Flows with the same losses are Wasserstein flows.
- Many methods use neural networks to approximate the flow (Alvarez-Melis et al., 2022; Heng et al., 2023).

Experimental Results

Dataset	PMs (no generator)		Int-PMs (generator)	
	EDM	Discr. Flow	GAN	Score GAN
MNIST	3	4	3	15
CelebA	10	41	19	35

- Hybrid models are viable, and support the theory.
EDM: diffusion (Karras et al., 2022).

Properties



- Discriminator flows** learn a path to the data distribution, unlike **diffusion**.
- PMs vs Int-PMs**: Int-PMs are prone to mode collapse but are faster than PMs at inference and have better latent space properties.

Perspectives

- Our work paves the way for new hybrid models.
- Model improvements: Score GANs for score distillation, Discriminator Flows for generation efficiency.
- Framework improvements: convergence guarantees, second-order and discrete-time optimization, more accurate GAN modeling.