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 Al 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

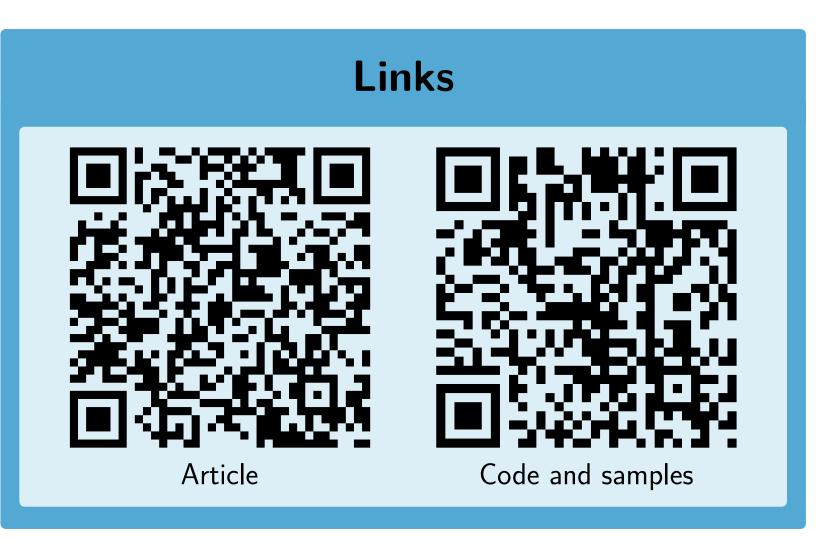
- We unify gradient flows, score-based diffusion models, and GANs by representing generated data as moving particles.
- A model is defined by:
- a gradient vector field that the particles follow;
- the possibility of incorporating a generator smoothing this movement.
- This suggests the existence of 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** \rightarrow Learns to progressively reverse a data degradation process.
- No generator (operates on the data space).
- SOTA performance.
- Easier to optimize.
- Slow inference.



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

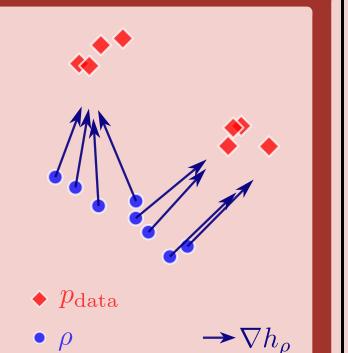
$$\alpha_t \nabla \log \left[p_{\text{data}} \star k_{\text{RBF}}^{\sigma(t)} \right] - \beta_t \nabla \log \rho_t$$

Discriminator Gradient

$$-
abla(c\circ f_{
ho_t})$$
 where $f_{
ho_t}$ discriminates ho_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 **gradi**ent ascent path on $h_{\rho_t}(x_t)$.
- h_{ρ} is usually a predefined functional approximated with neural networks. • p

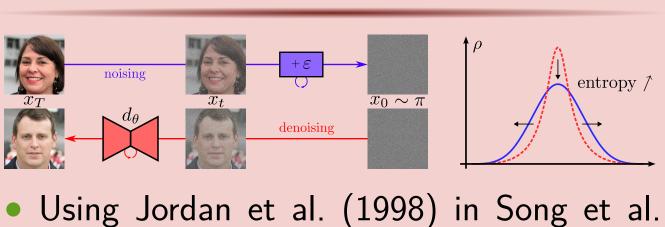


Wasserstein Gradient Flows

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

Objective $\mathcal{F}(ho)$		$h_{ ho}$
Forward KL	$\mathbb{E}_{ ho} \log ho \! / p_{ ext{data}}$	$-\log ho\!/p_{ m data}$
f-divergence	$\mathbb{E}_{p_{\mathrm{data}}} f\left(ho/p_{\mathrm{data}} ight)$	$-f'(ho/p_{ m data})$
Squared MMD w.r.t. kernel k	$\mathbb{E}_{\substack{x,x'\sim\rho\\y,y'\sim p_{\text{data}}}} \begin{bmatrix} k(x,x')\\+k(y,y')\\-2k(x,y) \end{bmatrix}$	$\mathbb{E}_{y \sim p_{\text{data}}} [k(y, \cdot)] \\ -\mathbb{E}_{x \sim \rho} [k(x, \cdot)]$
Entropy	$\mathbb{E}_{ ho}\log ho$	$-\log \rho$

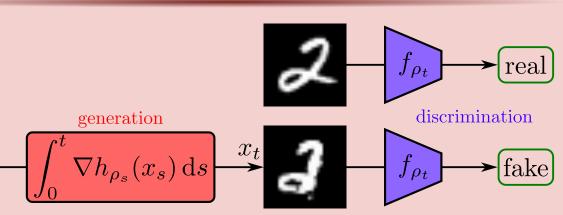
Score-Based Diffusion



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

(2019) and Karras et al. (2022):

Discriminator Flows

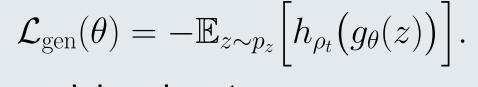


- Particles directly follow the discriminator gradient.
- The discriminator is simultaneously trained and used to generate data.

GANs

Interacting Particle Models (Generator)

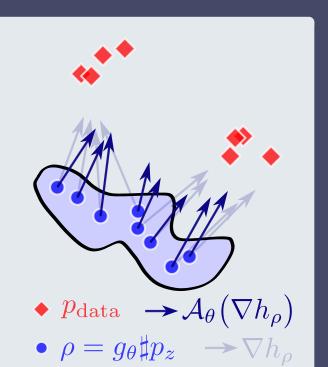
• Training with the same loss:



• At training time *t*:

 $\mathrm{d}g_{\theta_t}(z) = \eta \left[\mathcal{A}_{\theta_t}(z) \right] \left(\nabla h_{\rho_t} \right) \mathrm{d}t.$ Generalization of PMs where parti-

cles interact with each other.



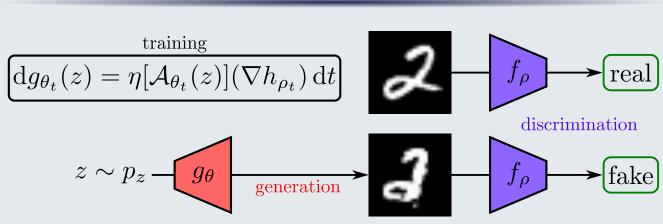
Stein Gradient Flows

• Stein gradient flows (Liu, 2017) are kernelized Wasserstein gradient flows: $dx_t = \mathbb{E}_{x_t' \sim \rho_t} \left| k(x_t, x_t') \nabla h_{\rho_t}(x_t') \right| dt.$

- Int-PMs under mild hypotheses (generalization of Durr et al. (2022)).
- ullet Hint towards the same $h_{
 ho}$ 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): $s_{\psi}^{p_{\mathrm{data}}}(\cdot,\sigma) \equiv \nabla \log |p_{\mathrm{data}} \star k_{\mathrm{RBF}}^{\sigma}|$. Generated distribution score (continuously
- updated like a discriminator): $s^{\rho}_{\phi} \equiv \nabla \log \rho$. • $\nabla h_{\rho}(x,\sigma) = s_{\psi}^{p_{\text{data}}}(x,\sigma) - s_{\phi}^{\rho}(x)$.



Gradient descent-ascent on the min-max objective yields the generator loss:

$$\mathcal{L}_{\mathrm{GAN}}(g_{ heta}) = \mathbb{E}_{z \sim p_z} \Big[\big(c \circ f_{
ho} \big) \big(g_{ heta}(z) \big) \Big].$$

Smoothing Operator

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

$$\begin{split} \big[\mathcal{A}_{\theta_t}(z) \big](V) &\triangleq \mathbb{E}_{z' \sim p_z} \bigg[k_{g_{\theta_t}} \big(z, z' \big) V \Big(g_{\theta_t} \big(z' \big) \Big) \bigg], \\ k_{g_{\theta_t}} \big(z, z' \big) &\triangleq \nabla_{\theta_t} g_{\theta_t} \big(z' \big)^\top \nabla_{\theta_t} g_{\theta_t} \big(z \big). \end{split}$$

- $k_{g_{\theta_t}}$ 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: $[\mathcal{A}_{\theta_t}(z)](V) = V(g_{\theta}(z))$.
- $\mathrm{d}g_{\theta_t}(z) = \nabla h_{\rho_t}(g_{\theta_t}(z)) \, \mathrm{d}t$: 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

- scheduling σ s w.r.t. training time t.
- ullet We randomly sample σ and also noise the particles:
- Generator update:

port the theory.

- few-step training of s^{ρ}_{ϕ} with denoising score matching;
- gradient descent step:

$$\theta \leftarrow \theta + \eta \underset{\sigma \sim p_{\sigma}, \varepsilon \sim \mathcal{N}(0, \sigma I_{D}), z \sim p_{z}}{\mathbb{E}} \left[\nabla_{\theta} g_{\theta}(z) \widetilde{\nabla h}_{\rho} (g_{\theta}(z) + \varepsilon, \sigma) \right].$$

Discriminator Flows in Practice

• Discriminator loss:

 $\mathcal{L}_{\mathrm{d}}(f;\rho,p_{\mathrm{data}}) = \mathbb{E}_{\rho}[a \circ f] - \mathbb{E}_{p_{\mathrm{data}}}[b \circ f] + \mathcal{R}(f;\rho,p_{\mathrm{data}}).$

- Naive training: successive f_{ρ_t} trainings and ρ_t updates.
- For efficiency purposes, we simultaneously learn all timeparameterized discriminators: $f_{\rho_t} = f_{\phi}(\cdot, t)$.
- Training step:
- sample $t \sim \mathcal{U}([0,1])$, $x_0 \sim \pi$;
- compute $x_t = -\eta \int_0^t \nabla (c \circ f_\phi(\cdot, s))(x_s) \, \mathrm{d}s;$
- 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:

$$\mathcal{L}_{\mathrm{gen}}(heta) = -\mathbb{E}_{z \sim p_z} \Big[h_{
ho_t} ig(g_{ heta}(z) ig) \Big].$$

- We do not take into account the dependency of ρ_t w.r.t. θ_t , to mimic PMs: $\rho = \text{StopGradient}(g_{\theta} \sharp p_z).$
- Continuous-time gradient descent:

$$\frac{\mathrm{d}\theta_t}{\mathrm{d}t} = -\eta \nabla_{\theta_t} \mathcal{L}_{\mathrm{gen}}(\theta_t) = \eta \nabla_{\theta_t} \mathbb{E}_{z \sim p_z} \Big[h_{\rho_t} \big(g_{\theta_t}(z) \big) \Big]
= \eta \mathbb{E}_{z \sim p_z} \Big[\nabla_{\theta_t} g_{\theta_t}(z) \nabla h_{\rho_t} \big(g_{\theta_t}(z) \big) \Big].$$

• Evolution of particles:

$$\frac{\mathrm{d}g_{\theta_t}(z)}{\mathrm{d}t} = \nabla_{\theta_t} g_{\theta_t}(z)^\top \frac{\mathrm{d}\theta_t}{\mathrm{d}t} = \eta \mathbb{E}_{z' \sim p_z} \left[\nabla_{\theta_t} g_{\theta_t}(z)^\top \nabla_{\theta_t} g_{\theta_t}(z') \nabla h_{\rho_t} \Big(g_{\theta_t}(z') \Big) \right].$$

Other Models & Flows

- Int-PMs and Stein (generalization of Durr et al. (2022)): $k(g_{\theta_t}(z), g_{\theta_t}(z')) = 0$ $k_{g_{\theta_{\star}}}(z,z')$ 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).

- Two practical issues:
- sliced score matching to train s^{ρ}_{ϕ} ;

PMs (no generator) Int-PMs (generator)

Hybrid models are viable, and sup-

EDM: diffusion (Karras et al., 2022).

GAN

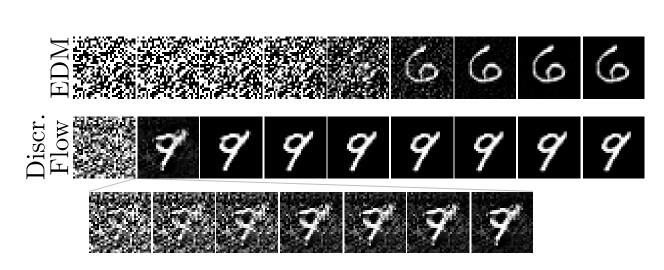
$$egin{aligned}
abla h_{
ho} &=
abla \log[p_{ ext{data}} \star k_{ ext{RBF}}^{\sigma}] -
abla \log[
ho_t \star k_{ ext{RBF}}^{\sigma}], \\ &\equiv \widetilde{\nabla} h_{
ho}(\cdot, \sigma) = s_{\psi}^{p_{ ext{data}}}(\cdot, \sigma) - s_{\phi}^{
ho}(\cdot, \sigma). \end{aligned}$$

Score GAN

Experimental Results



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.