

#### A Neural Tangent Kernel Perspective of GANs

Jean-Yves Franceschi, Emmanuel de Bézenac, Ibrahim Ayed, Mickaël Chen, Sylvain Lamprier, Patrick Gallinari

#### ▶ To cite this version:

Jean-Yves Franceschi, Emmanuel de Bézenac, Ibrahim Ayed, Mickaël Chen, Sylvain Lamprier, et al.. A Neural Tangent Kernel Perspective of GANs. Thirty-ninth International Conference on Machine Learning, Jul 2022, Baltimore, MD, United States. . hal-03716574

#### HAL Id: hal-03716574 https://hal.science/hal-03716574

Submitted on 7 Jul2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



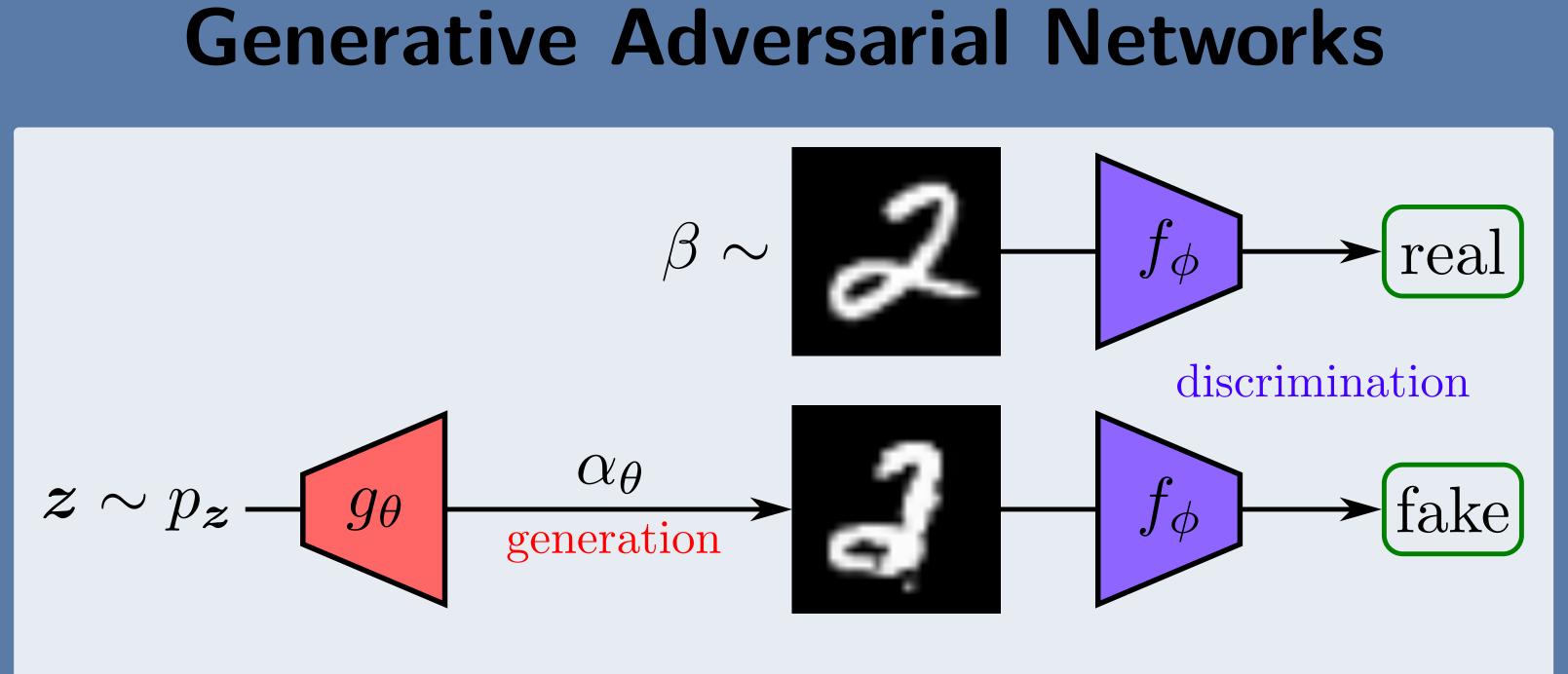
Distributed under a Creative Commons Attribution - NonCommercial - ShareAlike 4.0 International License



# J.-Y. Franceschi,<sup>\*,1,2</sup> E. de Bézenac,<sup>\*,3,2</sup> I. Ayed,<sup>\*,2,4</sup> M. Chen,<sup>5</sup> S. Lamprier,<sup>2</sup> P. Gallinari<sup>2,1</sup>

## Motivation & Outline

- Many analyses cannot explain GAN training as they fail to take into account alternating optimization and model the architecture and implicit biases of the discriminator.
- We propose a theoretical framework solving these issues using the theory of Neural Tangent Kernels.
- We deduce new insights about the flow and convergence of the generated distribution during training.



$$\inf_{\theta} \sup_{\phi} \mathcal{L}(g_{\theta}, f_{\phi}) = \inf_{\theta} \mathcal{L}(g_{\theta}, f_{\phi_{\theta}^{\star}}) \approx \inf_{\theta} \mathscr{C}(\alpha_{\theta}, \beta),$$
$$\mathcal{L}(g_{\theta}, f_{\phi}) = \mathbb{E}_{x \sim \alpha_{\theta}} \Big[ (a \circ f_{\phi})(x) \Big] - \mathbb{E}_{y \sim \beta} \Big[ (b \circ f_{\phi})(y) \Big].$$

E.g. for vanilla GAN  $\mathscr{C}$  is a Jensen-Shannon, for WGAN it is the 1-Wasserstein, for LSGAN it is a  $\chi^2$ -divergence.

## **Issue 1: Alternating Optimization**

• In practice, GANs are optimized alternatingly,  $f_{\phi}$  and  $g_{\theta}$  being considered to be independent of each other.

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}(g_{\theta}, f_{\phi});$$
  
$$\phi \leftarrow \phi + \lambda \nabla_{\phi} \mathcal{L}(g_{\theta}, f_{\phi}).$$

• Gradient received by  $g_{\theta}$ :

$$\nabla_{\theta} \mathcal{L} \Big( g_{\theta}, f_{\phi_{\theta}^{\star}} \Big) \qquad \Rightarrow \qquad \nabla_{\theta} \mathcal{L} \big( g_{\theta}, f_{\phi} \big).$$

• This changes the true generator loss  $\mathscr{C}$ . Hence, sound studies should model alternating optimization.

Thirty-ninth International Conference on Machine Learning — ICML 2022, Baltimore, MD, USA

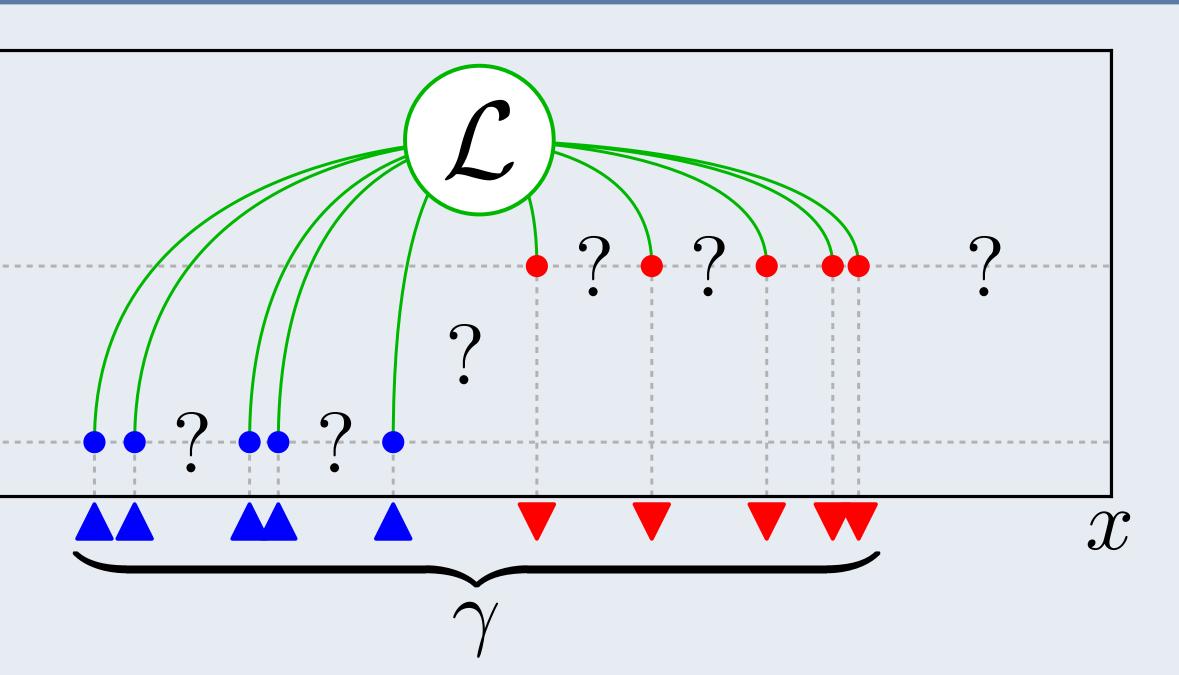
### $\nabla_{\theta} \mathcal{L}(q_{\theta}, .)$

- ill-defined.

 $\mathcal{T}_{k,\gamma}$ : L where  $\mathcal{H}_{k}^{\gamma}\subseteq$ 

# CRITEO ETHzürich Valeo A Neural Tangent Kernel THALES **Perspective of GANs** ALCO

## **Issue 2: III-Defined Discr. Gradients**



In an alternating optimization setting:

$$f) = \mathbb{E}_{z \sim p_z} \left[ \nabla_{\theta} g_{\theta}(z)^{\top} \nabla_x (a \circ f)(x) |_{x = g_{\theta}(z)} \right].$$

• The gradient of the generator requires  $\nabla f$  (chain rule). • Without any assumption on the structure of f, as  $\mathcal{L}$  is only defined on training points  $\gamma$ ,  $\nabla f$  is not defined. The parameter gradient of the generator is thus also

• Analyses need to take into account the structure of f.

Our solution: Modeling discriminator training as a neural network via its NTK.

### A Background on NTKs

The Neural Tangent Kernel: For a neural network  $f_{\phi}$  with parameters  $\phi$ , its NTK  $k_{f_{\phi}}$  is defined as:  $k_{f_{\phi}}(x,y) \triangleq \langle \partial_{\phi} f_{\phi}(x), \partial_{\phi} f_{\phi}(y) \rangle.$ 

In the infinite-width limit of f, during training:  $k_{f_{\phi}}(x,y) = k(x,y).$ 

### Kernel Integral Operator and RKHS:

$$\mathcal{L}^{2}(\gamma) \to \mathcal{H}_{k}^{\gamma}, \ h \mapsto \int_{x} k(\cdot, x) h(x) \,\mathrm{d}\gamma(x),$$
  
  $\subseteq L^{2}(\Omega)$  is the RKHS of  $k$  generated by  $\gamma$ 

### **Discriminator Inner Loop**

We consider the NNs in the NTK regime. This enables a theoretical study of their evolution w.r.t. training time t:

$$\partial_t f_t = \mathcal{T}_{k,\gamma} (\nabla$$

#### **Discriminator Structure:** Under mild assumptions, $f_t$ is uniquely defined and:

$$\forall t \in \mathbb{R}_+, f_t = f_0 + \mathcal{T}_{k,\gamma} \left( \int_0^t \nabla^{\gamma} \mathcal{L}_{\alpha}(f_s) \, \mathrm{d}s \right) \in f_0 + \mathcal{H}_k^{\gamma}$$

- $\mathcal{T}_{k,\gamma}$  smooths out gradients over the whole input space by sending them into  $\mathcal{H}_{k}^{\gamma}$ .
- $\mathcal{H}_{k}^{\gamma}$  depends on discriminator architecture.

### Differentiability of the Discriminator:

The discriminator trained with gradient descent is infinitely differentiable (almost) everywhere.

- The spatial gradient of the discriminator  $\nabla f_t$  is welldefined.
- is well-defined.

### **Underlying NTK Regularity Results**

To prove the above results, we establish novel regularity results on NTKs. Given, for the network f:

- a standard architecture (fully connected, convolutional, residual, etc.),
- Gaussian, etc.),

we prove that the NTK k of f is:

• smooth almost everywhere if the network has non-null bias terms. smooth everywhere if the activation is smooth.

These results, obtained from similar regularity results on the conjugate kernel of f, then transfer to f.

 $\nabla^{\gamma} \mathcal{L}_{\alpha}(f_t)$ ).

• The parameter gradient of the generator  $\nabla_{\theta} \mathcal{L}(g_{\theta}, f_t)$ 

• a standard activation function (tanh, softplus, ReLU-like, sigmoid,

## **Resulting Convergence Results**

Our finer-grained framework allows us to derive novel convergence insights, with highlighted results below.

#### **Gradient Flow of Generated Distribution:**

$$egin{aligned} &\mathcal{D}_{\ell} lpha_{\ell}^{z} = - 
abla_{k_{g}} \mathscr{C}(lpha_{\ell}) \ &= 
abla_{x} \cdot \left( lpha_{\ell}^{z} \mathcal{T}_{k_{g}, p_{z}} \Big( z \mapsto 
abla_{x} (a \circ f_{lpha_{\ell}})(x) \Big|_{x} \Big) \right) \end{aligned}$$

where  $\alpha_{\ell}^{z}$  is the distribution of  $(z, g_{\ell}(z))$  under  $z \sim p_{z}$ .

• In the non-interacting case, i.e.  $\mathcal{T}_{k_a,p_z} = id$ , this corresponds to a Wasserstein gradient flow:

$$\partial_{\ell} \alpha_{\ell} = -\nabla_{\mathscr{S}_{k_g}} \mathscr{C}(\alpha_{\ell}) = -\nabla_{\mathscr{W}} \mathscr{C}(\alpha_{\ell})$$

- In the general case, this is a gradient flow in a Stein geometry defined by the generator's NTK  $k_q$ .
- $\mathscr{C}(\alpha_{\ell})$  is automatically decreasing via this gradient flow, as fast as possible (locally).

IPM GANs (a = b = id) and NTK MMD:

We find  $f_t = f_0 + t \Big( \mathbb{E}_{x \sim \alpha} [k(x, \cdot)] - \mathbb{E}_{y \sim \beta} [k(y, \cdot)] \Big)$ , hence, if  $f_0 = 0$ :  $\mathscr{C}(\alpha) = \mathscr{L}_{\alpha}(f_t) \propto \mathrm{MMD}_k^2(\alpha, \beta).$ 

## **Empirical Study**

- We assess the adequacy of our framework by observing how close finite and infinite-width regimes are.
- We study the convergence of GANs on empirical distributions in the non-interacting case  $(\mathcal{T}_{k_a,p_z} = id)$ .
- We discover the singular performance of ReLU architectures for generative modeling and explain it by studying generator gradients with our framework.

# <sup>1</sup>Criteo Al Lab <sup>2</sup>Sorbonne Université, CNRS, ISIR <sup>3</sup>SAM, D-MATH, ETH Zürich <sup>4</sup>ThereSIS Lab, Thales <sup>5</sup>Valeo.ai

