



HAL
open science

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 Jul 2022

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

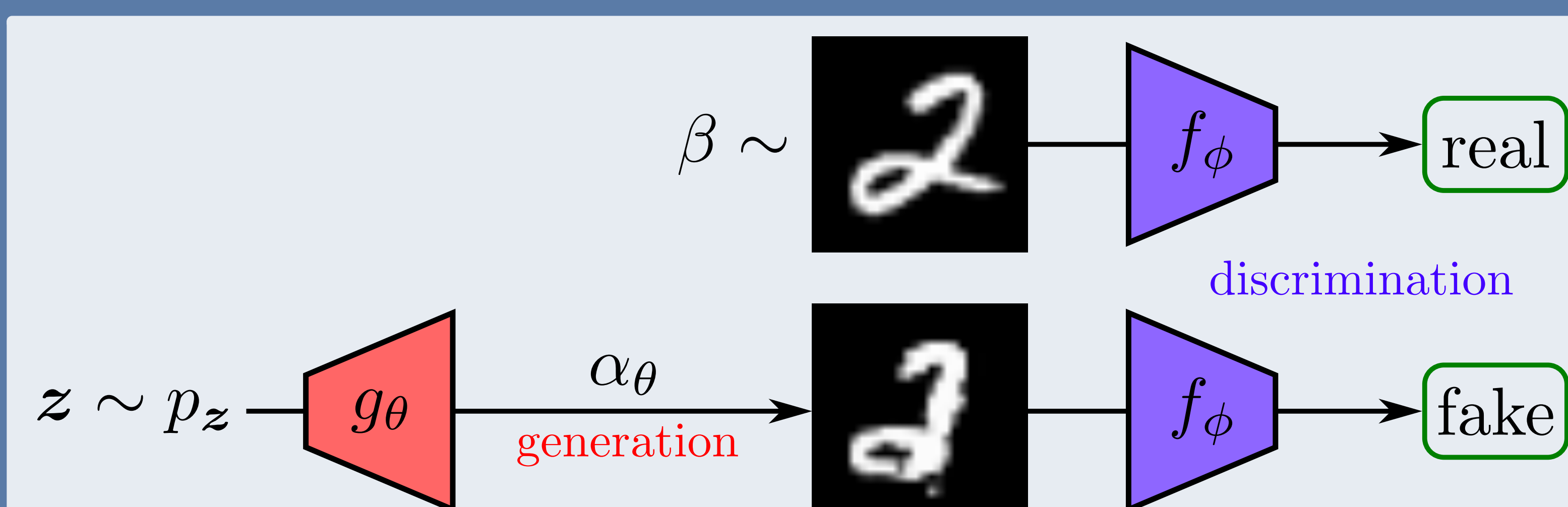
J.-Y. Franceschi,^{*,1,2} E. de Bézenac,^{*,3,2} I. Ayed,^{*,2,4} M. Chen,⁵ S. Lamprier,² P. Gallinari^{2,1}

¹Criteo AI Lab ²Sorbonne Université, CNRS, ISIR ³SAM, D-MATH, ETH Zürich ⁴ThereSIS Lab, Thales ⁵Valeo.ai

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.

Generative Adversarial Networks



$$\inf_{\theta} \sup_{\phi} \mathcal{L}(g_{\theta}, f_{\phi}) = \inf_{\theta} \mathcal{L}(g_{\theta}, f_{\phi}^*) \approx \inf_{\theta} \mathcal{C}(\alpha_{\theta}, \beta),$$

$$\mathcal{L}(g_{\theta}, f_{\phi}) = \mathbb{E}_{x \sim \alpha_{\theta}} [(a \circ f_{\phi})(x)] - \mathbb{E}_{y \sim \beta} [(b \circ f_{\phi})(y)].$$

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

Issue 1: Alternating Optimization

- In practice, GANs are optimized alternatingly, f_{ϕ} and g_{θ} being considered to be independent of each other.

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

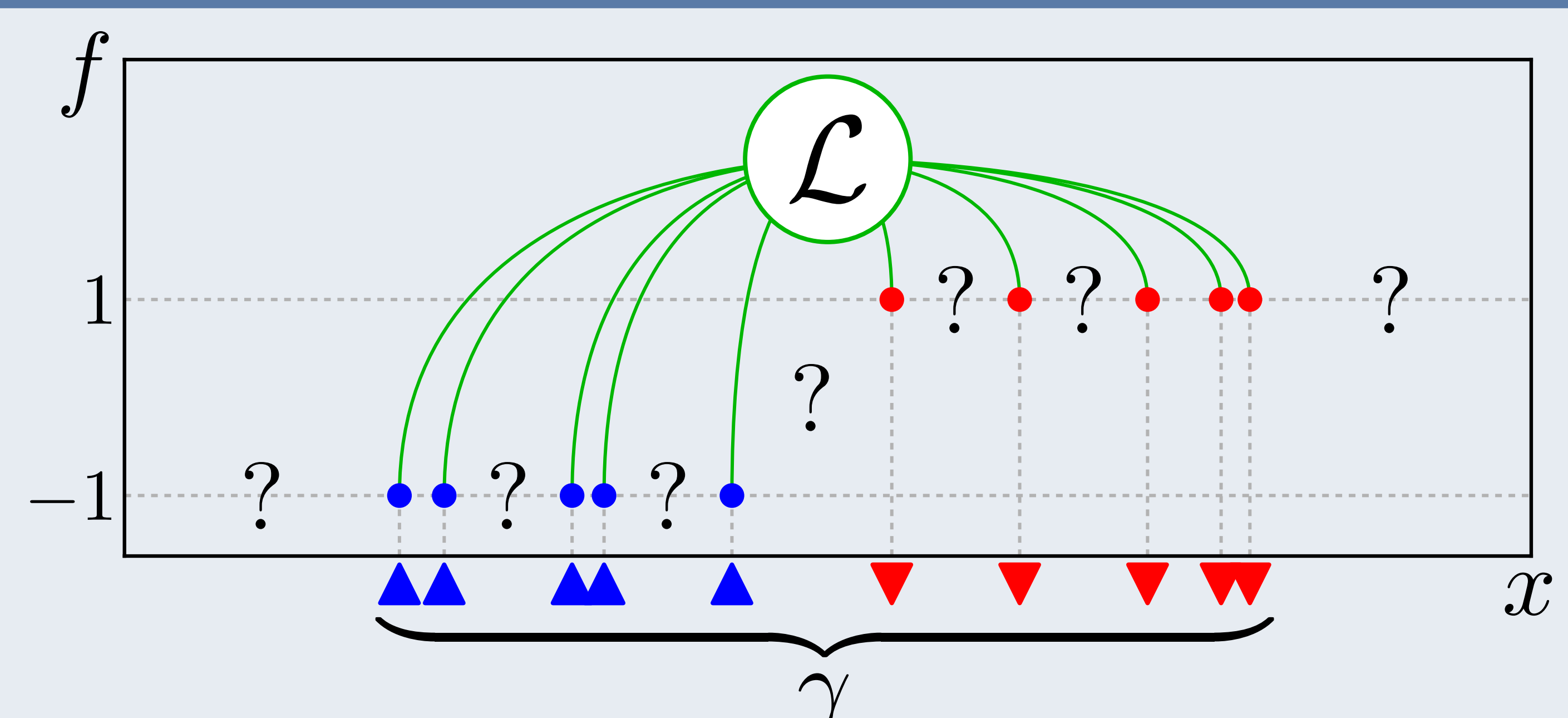
$$\phi \leftarrow \phi + \lambda \nabla_{\phi} \mathcal{L}(g_{\theta}, f_{\phi}).$$

- Gradient received by g_{θ} :

$$\nabla_{\theta} \mathcal{L}(g_{\theta}, f_{\phi}^*) \Rightarrow \nabla_{\theta} \mathcal{L}(g_{\theta}, f_{\phi}).$$

- This changes the true generator loss \mathcal{C} . Hence, sound studies should model alternating optimization.

Issue 2: Ill-Defined Discr. Gradients



In an alternating optimization setting:

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

- The gradient of the generator requires ∇f (chain rule).
- Without any assumption on the structure of f , as \mathcal{L} is only defined on training points γ , ∇f is not defined.
- The parameter gradient of the generator is thus also ill-defined.
- 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_{ϕ} with parameters ϕ , 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:

$$\mathcal{T}_{k, \gamma}: L^2(\gamma) \rightarrow \mathcal{H}_k^{\gamma}, h \mapsto \int_x k(\cdot, x) h(x) d\gamma(x),$$

where $\mathcal{H}_k^{\gamma} \subseteq L^2(\Omega)$ is the RKHS of k generated by γ .

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 :

$$\partial_t f_t = \mathcal{T}_{k, \gamma}(\nabla^{\gamma} \mathcal{L}_{\alpha}(f_t)).$$

Discriminator Structure:

Under mild assumptions, f_t is uniquely defined and:

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

- $\mathcal{T}_{k, \gamma}$ smooths out gradients over the whole input space by sending them into \mathcal{H}_k^{γ} .
- \mathcal{H}_k^{γ} 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 ∇f_t is well-defined.
- The parameter gradient of the generator $\nabla_{\theta} \mathcal{L}(g_{\theta}, f_t)$ 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.),
- a standard activation function (tanh, softplus, ReLU-like, sigmoid, 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 .

Resulting Convergence Results

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

Gradient Flow of Generated Distribution:

$$\partial_t \alpha_t^z = -\nabla_{\mathcal{G}_{k_g}} \mathcal{C}(\alpha_t)$$

$$= \nabla_x \cdot \left(\alpha_t^z \mathcal{T}_{k_g, p_z} \left(z \mapsto \nabla_x (a \circ f_{\alpha_t})(x)|_{x=g_{\ell}(z)} \right) \right),$$

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

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

$$\partial_t \alpha_t = -\nabla_{\mathcal{G}_{k_g}} \mathcal{C}(\alpha_t) = -\nabla_{\mathcal{W}} \mathcal{C}(\alpha_t).$$

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

IPM GANs ($a = b = \text{id}$) and NTK MMD:

We find $f_t = f_0 + t \left(\mathbb{E}_{x \sim \alpha} [k(x, \cdot)] - \mathbb{E}_{y \sim \beta} [k(y, \cdot)] \right)$, hence, if $f_0 = 0$:

$$\mathcal{C}(\alpha) = \mathcal{L}_{\alpha}(f_t) \propto \text{MMD}_k^2(\alpha, \beta).$$

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 ($\mathcal{T}_{k_g, p_z} = \text{id}$).
- We discover the singular performance of ReLU architectures for generative modeling and explain it by studying generator gradients with our framework.

Toolkit GAN(TK)²

- GAN analysis toolkit based on our framework.
- Written in JAX with the Neural Tangents library.
- Link to paper and code \rightarrow

