

Equivariant Neural Tangent Kernels

by Philipp Misof

Department of Mathematical Sciences, Division of Algebra and Geometry

March 13, 2025



UNIVERSITY OF
GOTHENBURG



CHALMERS
UNIVERSITY OF TECHNOLOGY

WASP | WALLENBERG AI,
AUTONOMOUS SYSTEMS
AND SOFTWARE PROGRAM

Joint work with



Jan Gerken
(Chalmers)

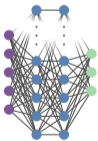
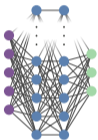


Pan Kessel
(Prescient Design, Switzerland)

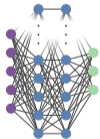
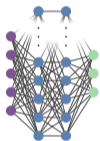
- 1 The Neural Tangent Kernel
- 2 NTK of Equivariant NNs
- 3 Concrete Examples
- 4 Data Augmentation vs GCNNs

What can we say about the **dynamics of NNs** without specifying their initial parameters?

What can we say about the **dynamics of NNs** without specifying their initial parameters?

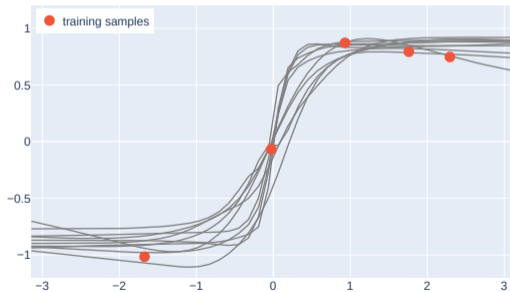


What can we say about the **dynamics of NNs** without specifying their initial parameters?

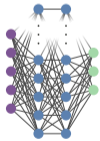
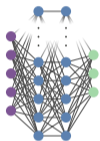


Let's learn $\sin(x)$

3-layer MLP with widths 8

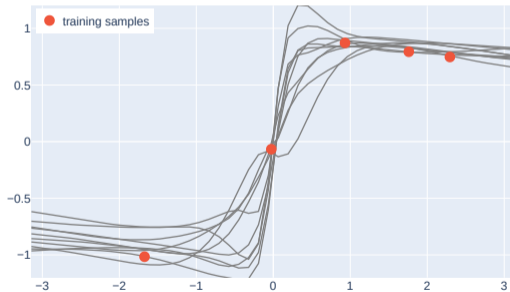


What can we say about the **dynamics of NNs** without specifying their initial parameters?

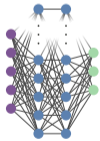
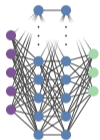


Let's learn $\sin(x)$

3-layer MLP with widths 64

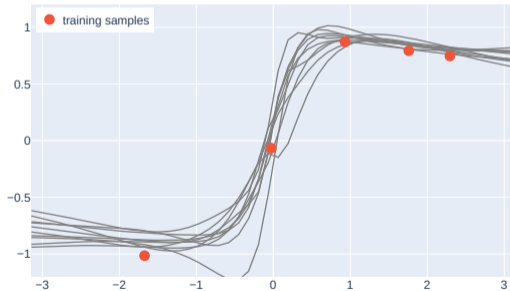


What can we say about the **dynamics of NNs** without specifying their initial parameters?

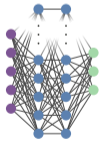
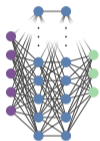


Let's learn $\sin(x)$

3-layer MLP with widths 256

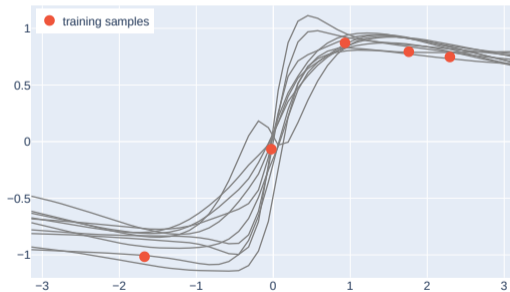


What can we say about the **dynamics of NNs** without specifying their initial parameters?

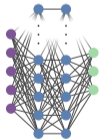
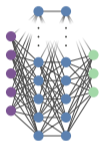


Let's learn $\sin(x)$

3-layer MLP with widths 1024

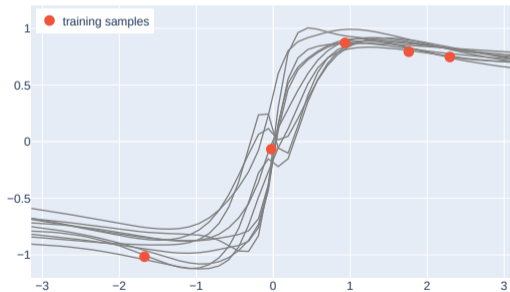


What can we say about the **dynamics of NNs** without specifying their initial parameters?

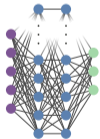
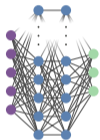


Let's learn $\sin(x)$

3-layer MLP with widths 4096

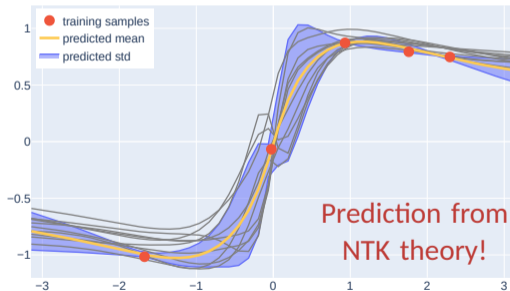


What can we say about the **dynamics of NNs** without specifying their initial parameters?



Let's learn $\sin(x)$

3-layer MLP with widths 4096



Training dynamics

(Jacot, Gabriel, and Hongler 2018)

Training dynamics

(Jacot, Gabriel, and Hongler 2018)

Assume **Gradient Flow** instead of discrete Gradient Descent

Training dynamics

(Jacot, Gabriel, and Hongler 2018)

Assume **Gradient Flow** instead of discrete Gradient Descent

$$\frac{d\mathcal{N}}{dt}(\mathbf{x}) = -\eta \sum_{i=1}^{n_{\text{train}}} \Theta_t(\mathbf{x}, \mathbf{x}_i) \frac{\partial \mathcal{L}}{\partial \mathcal{N}(\mathbf{x}_i)}$$

Training dynamics

(Jacot, Gabriel, and Hongler 2018)

Assume **Gradient Flow** instead of discrete Gradient Descent

The diagram shows the equation for Gradient Flow with several annotations:

- Neural Network output**: An arrow points from this text to the \mathcal{N} in the numerator of the derivative $\frac{d\mathcal{N}}{dt}$.
- Learning rate**: An arrow points from this text to the η in the equation.
- Loss**: An arrow points from this text to the \mathcal{L} in the numerator of the fraction $\frac{\partial \mathcal{L}}{\partial \mathcal{N}(x_i)}$.
- Training time**: An arrow points from this text to the t in the denominator of the derivative $\frac{d\mathcal{N}}{dt}$.
- (empirical) Neural Tangent Kernel**: An arrow points from this text to the $\Theta_t(x, x_i)$ term in the equation.

$$\frac{d\mathcal{N}}{dt}(x) = -\eta \sum_{i=1}^{n_{\text{train}}} \Theta_t(x, x_i) \frac{\partial \mathcal{L}}{\partial \mathcal{N}(x_i)}$$

Training dynamics

(Jacot, Gabriel, and Hongler 2018)

Assume **Gradient Flow** instead of discrete Gradient Descent

The diagram shows the equation for Gradient Flow: $\frac{d\mathcal{N}}{dt}(\mathbf{x}) = -\eta \sum_{i=1}^{n_{\text{train}}} \Theta_t(\mathbf{x}, \mathbf{x}_i) \frac{\partial \mathcal{L}}{\partial \mathcal{N}(\mathbf{x}_i)}$. Annotations include: 'Neural Network output' pointing to \mathcal{N} ; 'Learning rate' pointing to η ; 'Loss' pointing to \mathcal{L} ; 'Training time' pointing to t ; and '(empirical) Neural Tangent Kernel' pointing to $\Theta_t(\mathbf{x}, \mathbf{x}_i)$.

$$\frac{d\mathcal{N}}{dt}(\mathbf{x}) = -\eta \sum_{i=1}^{n_{\text{train}}} \Theta_t(\mathbf{x}, \mathbf{x}_i) \frac{\partial \mathcal{L}}{\partial \mathcal{N}(\mathbf{x}_i)}$$

(empirical) Neural Tangent Kernel

Training dynamics

(Jacot, Gabriel, and Hongler 2018)

Assume **Gradient Flow** instead of discrete Gradient Descent

Neural Network output

Learning rate

Loss

$$\frac{d\mathcal{N}}{dt}(\mathbf{x}) = -\eta \sum_{i=1}^{n_{\text{train}}} \Theta_t(\mathbf{x}, \mathbf{x}_i) \frac{\partial \mathcal{L}}{\partial \mathcal{N}(\mathbf{x}_i)}$$

Training time

(empirical) **Neural Tangent Kernel**

$$\Theta_t(\mathbf{x}, \mathbf{x}') = \sum_{\mu} \frac{\partial \mathcal{N}(\mathbf{x})}{\partial \theta_{\mu}} \left(\frac{\partial \mathcal{N}(\mathbf{x}')}{\partial \theta_{\mu}} \right)^{\top}$$

So far, nothing is gained

So far, nothing is gained

Θ_t is dependent on θ_t , which is **stochastic** and **time-dependent**

So far, nothing is gained

Θ_t is dependent on θ_t , which is **stochastic** and **time-dependent**

Freezing of the NTK

(Jacot, Gabriel, and Hongler 2018)

So far, nothing is gained

Θ_t is dependent on θ_t , which is **stochastic** and **time-dependent**

Freezing of the NTK

(Jacot, Gabriel, and Hongler 2018)

$$\Theta_t \xrightarrow[\text{layer width}]{\text{increasing}} \Theta = \mathbb{E}[\Theta_t]$$

So far, nothing is gained

Θ_t is dependent on θ_t , which is **stochastic** and **time-dependent**

Freezing of the NTK

(Jacot, Gabriel, and Hongler 2018)

$$\Theta_t \xrightarrow[\text{layer width}]{\text{increasing}} \Theta = \mathbb{E}[\Theta_t] \quad \rightarrow \quad \text{Simple ODE}$$

So far, nothing is gained

Θ_t is dependent on θ_t , which is **stochastic** and **time-dependent**

Freezing of the NTK

(Jacot, Gabriel, and Hongler 2018)

$$\Theta_t \xrightarrow[\text{layer width}]{\text{increasing}} \Theta = \mathbb{E}[\Theta_t] \quad \rightarrow \quad \text{Simple ODE}$$

In case of MLE loss, **closed-form solution** of the **mean**

$$\mu(x) = \Theta(x, \mathcal{X})\Theta(\mathcal{X}, \mathcal{X})^{-1}\mathcal{Y} \quad \text{as } t \rightarrow \infty$$

at ∞ width.

So far, nothing is gained

Θ_t is dependent on θ_t , which is **stochastic** and **time-dependent**

Freezing of the NTK

(Jacot, Gabriel, and Hongler 2018)

$$\Theta_t \xrightarrow[\text{layer width}]{\text{increasing}} \Theta = \mathbb{E}[\Theta_t] \quad \rightarrow \quad \text{Simple ODE}$$

In case of MLE loss, **closed-form solution** of the **mean**

$$\begin{array}{ccc} \text{Training inputs} & & \text{Training targets} \\ \downarrow & & \downarrow \\ \mu(x) = \Theta(x, \mathcal{X})\Theta(\mathcal{X}, \mathcal{X})^{-1}\mathcal{Y} & \text{as } t \rightarrow \infty & \end{array}$$

at ∞ width.

What do we gain?

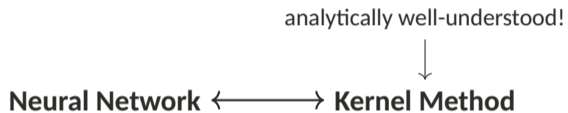
What do we gain?

- Useful correspondence

Neural Network \longleftrightarrow Kernel Method

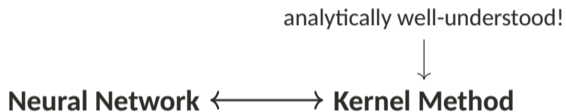
What do we gain?

- Useful correspondence



What do we gain?

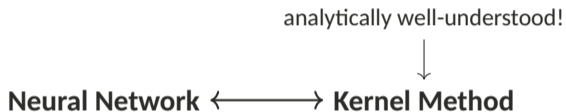
- Useful correspondence



- Non-linear **kernel method inspired by** empirical insights from **NNs** (*Arora et al. 2020*)

What do we gain?

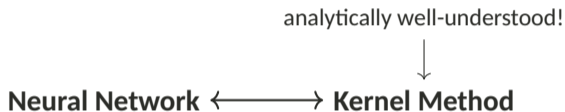
- Useful correspondence



- Non-linear **kernel method inspired by** empirical insights from **NNs** (*Arora et al. 2020*)
- Study **trainability vs generalization** regimes (*Xiao, Pennington, and Schoenholz 2020*)

What do we gain?

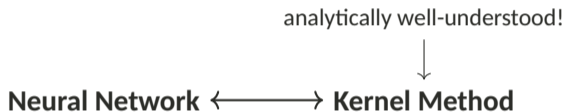
- Useful correspondence



- Non-linear **kernel method inspired by** empirical insights from **NNs** (*Arora et al. 2020*)
- Study **trainability vs generalization** regimes (*Xiao, Pennington, and Schoenholz 2020*)
- **Spectral bias** inside and outside the training set (*Bowman and Montufar 2022*)

What do we gain?

- Useful correspondence



- Non-linear **kernel method inspired by** empirical insights from **NNs** (*Arora et al. 2020*)
- Study **trainability vs generalization** regimes (*Xiao, Pennington, and Schoenholz 2020*)
- **Spectral bias** inside and outside the training set (*Bowman and Montufar 2022*)
- Tool for **dataset distillation** (*Nguyen, Chen, and Lee 2021*)

How is the NTK computed in practice?

How is the NTK computed in practice?

→ **layer by layer!** *(Jacot, Gabriel, and Hongler 2018)*

How is the NTK computed in practice?

→ **layer by layer!** *(Jacot, Gabriel, and Hongler 2018)*

Convenient: *Neural Network Gaussian Process Kernel (NNGP)*

$$K^{(\ell)}(\mathbf{x}, \mathbf{x}') = \mathbb{E} \left[\mathcal{N}^{(\ell)}(\mathbf{x}) \left(\mathcal{N}^{(\ell)}(\mathbf{x}') \right)^\top \right]$$

How is the NTK computed in practice?

→ **layer by layer!** (Jacot, Gabriel, and Hongler 2018)

Convenient: *Neural Network Gaussian Process Kernel (NNGP)*

$$K^{(\ell)}(\mathbf{x}, \mathbf{x}') = \mathbb{E} \left[\overset{\ell^{\text{th}} \text{ layer neurons}}{\mathcal{N}^{(\ell)}(\mathbf{x}) \left(\mathcal{N}^{(\ell)}(\mathbf{x}') \right)^{\text{T}}} \right]$$

How is the NTK computed in practice?

→ **layer by layer!** (Jacot, Gabriel, and Hongler 2018)

Convenient: *Neural Network Gaussian Process Kernel (NNGP)*

$$K^{(\ell)}(\mathbf{x}, \mathbf{x}') = \mathbb{E} \left[\overset{\ell^{\text{th}} \text{ layer neurons}}{\mathcal{N}^{(\ell)}(\mathbf{x}) \left(\mathcal{N}^{(\ell)}(\mathbf{x}') \right)^{\top}} \right]$$

Each NN layer then corresponds to a particular **recursion**

$$\begin{aligned} K^{(\ell+1)}(\mathbf{x}, \mathbf{x}') &= A^{(\ell)}(K^{(\ell)}(\mathbf{x}, \mathbf{x}')), \\ \Theta^{(\ell+1)}(\mathbf{x}, \mathbf{x}') &= B^{(\ell)}(\Theta^{(\ell)}(\mathbf{x}, \mathbf{x}'), K^{(\ell+1)}(\mathbf{x}, \mathbf{x}')) \end{aligned}$$

How is the NTK computed in practice?

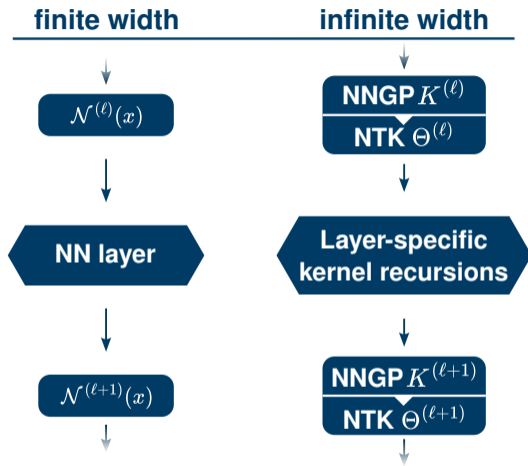
→ **layer by layer!** (Jacot, Gabriel, and Hongler 2018)

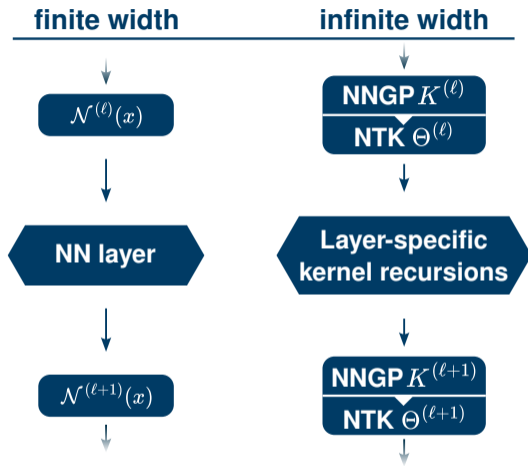
Convenient: *Neural Network Gaussian Process Kernel (NNGP)*

$$K^{(\ell)}(\mathbf{x}, \mathbf{x}') = \mathbb{E} \left[\overset{\ell^{\text{th}} \text{ layer neurons}}{\mathcal{N}^{(\ell)}(\mathbf{x}) \left(\mathcal{N}^{(\ell)}(\mathbf{x}') \right)^{\top}} \right]$$

Each NN layer then corresponds to a particular **recursion**

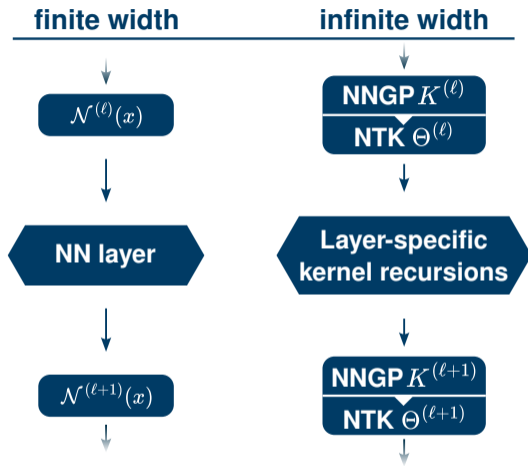
$$\begin{aligned} K^{(\ell+1)}(\mathbf{x}, \mathbf{x}') &= \overset{\text{layer-specific NNGP map}}{A^{(\ell)}}(K^{(\ell)}(\mathbf{x}, \mathbf{x}')), \\ \Theta^{(\ell+1)}(\mathbf{x}, \mathbf{x}') &= B^{(\ell)}(\Theta^{(\ell)}(\mathbf{x}, \mathbf{x}'), K^{(\ell+1)}(\mathbf{x}, \mathbf{x}')) \\ &\quad \overset{\text{layer-specific NTK map}}{\uparrow} \end{aligned}$$





→ Implemented in the neural-tangents library

(Novak et al. 2020)



→ Implemented in the neural-tangents library

(Novak et al. 2020)

Convention: We treat nonlinearities σ as individual layers

We extend the NTK to **equivariant NNs**

We extend the NTK to **equivariant NNs**

We cover



Group Convolution

We extend the NTK to **equivariant NNs**

We cover

Group Convolution

Group Pooling

We extend the NTK to **equivariant NNs**

We cover

Group Convolution

Group Pooling

Elementwise Non-linearity

- 1 The Neural Tangent Kernel
- 2 NTK of Equivariant NNs**
- 3 Concrete Examples
- 4 Data Augmentation vs GCNNs

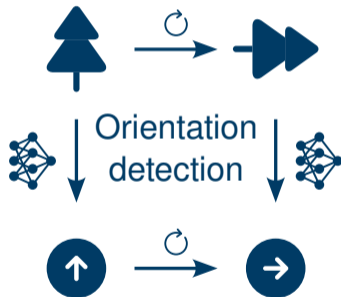
Recap on equivariant NNs

Recap on equivariant NNs

Want to enforce **symmetry** w.r.t a group G acting on the input signal f

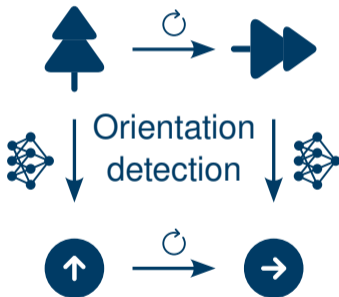
Recap on equivariant NNs

Want to enforce **symmetry** w.r.t a group G acting on the input signal f



Recap on equivariant NNs

Want to enforce **symmetry** w.r.t a group G acting on the input signal f



$$\begin{array}{ccc} f & \xrightarrow{\rho_{\text{in}}(g)} & \rho_{\text{in}}(g)(f) \\ \downarrow \mathcal{N} & & \downarrow \mathcal{N} \\ \mathcal{N}(f) & \xrightarrow{\rho_{\text{out}}(g)} & \rho_{\text{out}}(g)[\mathcal{N}(f)] \end{array}$$

$$\forall g \in G$$

Recap on Group Convolutional Neural Networks (GCNNs)

Recap on Group Convolutional Neural Networks (GCNNs)

Group convolutional layer

$$[\mathcal{N}^{(\ell+1)}(\mathbf{f})](\mathbf{g}) = \frac{1}{\sqrt{n_\ell |\mathcal{S}_\kappa|}} \int_G d\mathbf{h} \kappa(\mathbf{g}^{-1}\mathbf{h}) [\mathcal{N}^{(\ell)}(\mathbf{f})](\mathbf{h})$$

Recap on Group Convolutional Neural Networks (GCNNs)

Group convolutional layer

$$[\mathcal{N}^{(\ell+1)}(\mathbf{f})](\mathbf{g}) = \frac{1}{\sqrt{n_\ell |\mathcal{S}_\kappa|}} \int_G d\mathbf{h} \kappa(\mathbf{g}^{-1}\mathbf{h}) [\mathcal{N}^{(\ell)}(\mathbf{f})](\mathbf{h})$$

input signal \rightarrow

group element \rightarrow

filter \rightarrow

\uparrow # channels \uparrow filter support \leftarrow group

Recap on Group Convolutional Neural Networks (GCNNs)

Group pooling

$$\mathcal{N}^{(\ell+1)}(f) = \frac{1}{\text{vol}(G)} \int_G dg [\mathcal{N}^{(\ell)}(f)](g)$$

Recap on Group Convolutional Neural Networks (GCNNs)

Group pooling

$$\mathcal{N}^{(\ell+1)}(\mathbf{f}) = \frac{1}{\text{vol}(G)} \int_G dg [\mathcal{N}^{(\ell)}(\mathbf{f})](g)$$

Elementwise non-linearity σ

$$[\mathcal{N}^{(\ell+1)}(\mathbf{f})](g) = \sigma([\mathcal{N}^{(\ell)}(\mathbf{f})](g)), \quad \sigma : \mathbb{R} \rightarrow \mathbb{R} \text{ elementwise}$$

The Equivariant NTK

The Equivariant NTK

$$\mathbb{E} \left[\sum_{\mu} \frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f})](\mathbf{g})}{\partial\theta_{\mu}} \left(\frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f}')](\mathbf{g}')}{\partial\theta_{\mu}} \right)^{\top} \right]$$

The Equivariant NTK

$$\Theta_{g, g'}^{(\ell)}(\mathbf{f}, \mathbf{f}') = \mathbb{E} \left[\sum_{\mu} \frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f})](\mathbf{g})}{\partial \theta_{\mu}} \left(\frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f}')](\mathbf{g}')}{\partial \theta_{\mu}} \right)^{\top} \right]$$

\uparrow
Evaluation point in group space

The Equivariant NTK

$$\Theta_{g, g'}^{(\ell)}(\mathbf{f}, \mathbf{f}') = \mathbb{E} \left[\sum_{\mu} \frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f})](\mathbf{g})}{\partial \theta_{\mu}} \left(\frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f}')](\mathbf{g}')}{\partial \theta_{\mu}} \right)^{\top} \right]$$

\uparrow
Evaluation point in group space

∞ -width limit:

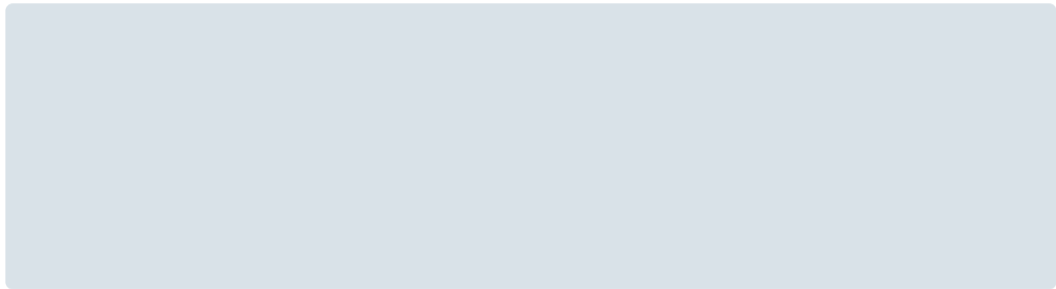
The Equivariant NTK

$$\Theta_{g, g'}^{(\ell)}(\mathbf{f}, \mathbf{f}') = \mathbb{E} \left[\sum_{\mu} \frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f})](\mathbf{g})}{\partial \theta_{\mu}} \left(\frac{\partial[\mathcal{N}^{(\ell)}(\mathbf{f}')](\mathbf{g}')}{\partial \theta_{\mu}} \right)^{\top} \right]$$

\uparrow
Evaluation point in group space

∞ -width limit: # channels $\rightarrow \infty$

Recursion of the **group convolutional layer**



Recursion of the group convolutional layer

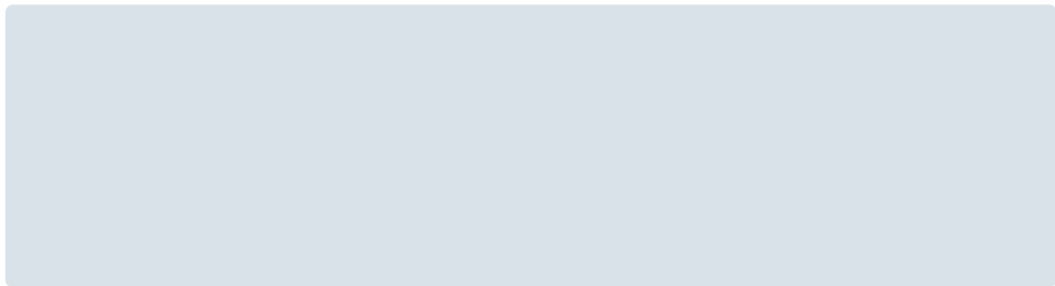
$$K_{g,g'}^{(\ell+1)}(f,f') = \frac{1}{\text{vol}(S_\kappa)} \int_{S_\kappa} dh K_{gh,g'h}^{(\ell)}(f,f')$$

Recursion of the group convolutional layer

$$K_{g,g'}^{(\ell+1)}(f, f') = \frac{1}{\text{vol}(\mathcal{S}_\kappa)} \int_{\mathcal{S}_\kappa} dh K_{gh,g'h}^{(\ell)}(f, f')$$

$$\Theta_{g,g'}^{(\ell+1)}(f, f') = K_{g,g'}^{(\ell+1)}(f, f') + \frac{1}{\text{vol}(\mathcal{S}_\kappa)} \int_{\mathcal{S}_\kappa} dh \Theta_{gh,g'h}^{(\ell)}(f, f')$$

Recursion of the **group pooling layer**



Recursion of the **group pooling layer**

$$K^{(\ell+1)}(f, f') = \frac{1}{\text{vol}(G)} \int_G dg dg' K_{g, g'}^{(\ell)}(f, f')$$

Recursion of the group pooling layer

$$K^{(\ell+1)}(f, f') = \frac{1}{\text{vol}(G)} \int_G dg dg' K_{g, g'}^{(\ell)}(f, f')$$

$$\Theta^{(\ell+1)}(f, f') = \frac{1}{\text{vol}(G)} \int_G dg dg' \Theta_{g, g'}^{(\ell)}(f, f')$$

Why is that interesting?

Why is that interesting?

- Can compare equivariant architectures with their non-equivariant counterparts **in a parameter-independent setting**

Why is that interesting?

- Can compare equivariant architectures with their non-equivariant counterparts **in a parameter-independent setting**
- Can compare **data augmentation** \leftrightarrow **equivariant architectures**

Why is that interesting?

- Can compare equivariant architectures with their non-equivariant counterparts **in a parameter-independent setting**
- Can compare **data augmentation** \leftrightarrow **equivariant architectures**
 - Ensembles of data-augmented networks are equivariant (*Gerken and Kessel 2024; Nordenfors and Flinth 2024*)

Why is that interesting?

- Can compare equivariant architectures with their non-equivariant counterparts **in a parameter-independent setting**
- Can compare **data augmentation** \leftrightarrow **equivariant architectures**
 - Ensembles of data-augmented networks are equivariant (*Gerken and Kessel 2024; Nordenfors and Flinth 2024*)
 - \rightarrow How does this form of equivariance differ?

Why is that interesting?

- Can compare equivariant architectures with their non-equivariant counterparts **in a parameter-independent setting**
- Can compare **data augmentation** \leftrightarrow **equivariant architectures**
 - Ensembles of data-augmented networks are equivariant (*Gerken and Kessel 2024; Nordenfors and Flinth 2024*)
 - \rightarrow How does this form of equivariance differ?
- Can study **biases** and the **influence of hyperparameters** on equivariant NNs through the NTK (*Xiao, Pennington, and Schoenholz 2020; Bowman and Montufar 2022*)

Why is that interesting?

- Can compare equivariant architectures with their non-equivariant counterparts **in a parameter-independent setting**
- Can compare **data augmentation** \leftrightarrow **equivariant architectures**
 - Ensembles of data-augmented networks are equivariant (*Gerken and Kessel 2024; Nordenfors and Flinth 2024*)
 - \rightarrow How does this form of equivariance differ?
- Can study **biases** and the **influence of hyperparameters** on equivariant NNs through the NTK (*Xiao, Pennington, and Schoenholz 2020; Bowman and Montufar 2022*)
- Obtain **equivariant Kernel methods**

- 1 The Neural Tangent Kernel
- 2 NTK of Equivariant NNs
- 3 Concrete Examples**
- 4 Data Augmentation vs GCNNs

On the practical side

On the practical side

Explicit expressions for

- **roto-translations** $G = C_4 \ltimes \mathbb{R}^2$ in the plane

On the practical side

Explicit expressions for

- **roto-translations** $G = C_4 \ltimes \mathbb{R}^2$ in the plane
- **3d-rotations** $G = SO(3)$ on the sphere S^2
via **spherical convolutions**

On the practical side

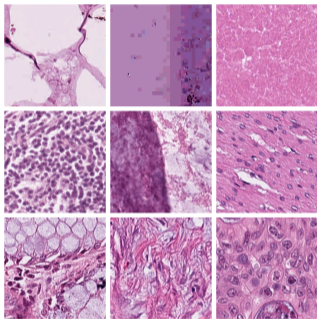
Explicit expressions for

- **roto-translations** $G = C_4 \times \mathbb{R}^2$ in the plane
- **3d-rotations** $G = SO(3)$ on the sphere S^2 via **spherical convolutions**

Provide **implementations** in the `neural-tangents` library (Novak et al. 2020)

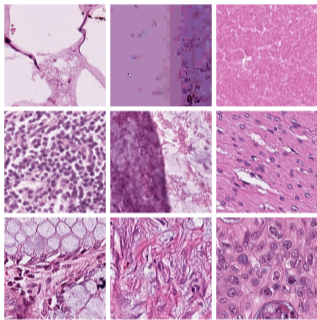
∞ -width limit in practice: Histological image classification

∞ -width limit in practice: Histological image classification



9 classes of microscopical
tissue images (*Kather, Halama,
and Marx 2018*)

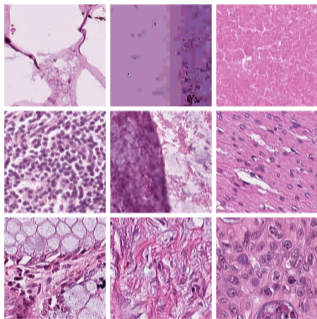
∞ -width limit in practice: Histological image classification



9 classes of microscopical
tissue images (*Kather, Halama,
and Marx 2018*)

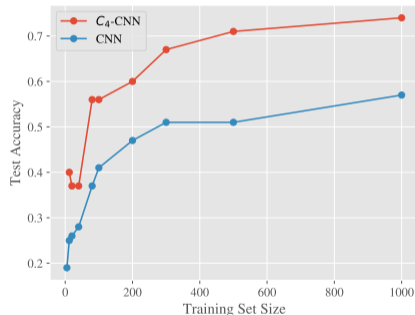
CNN
vs.
 $C_4 \times \mathbb{R}^2$ GCNN

∞ -width limit in practice: Histological image classification



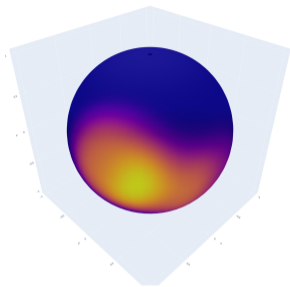
9 classes of microscopical tissue images (Kather, Halama, and Marx 2018)

CNN
vs.
 $C_4 \times \mathbb{R}^2$ GCNN



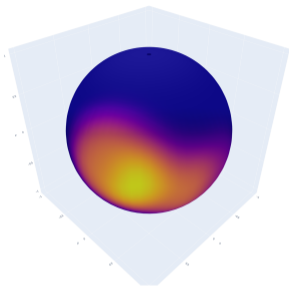
∞ -width limit in practice: Molecular Energy Prediction on QM9

∞ -width limit in practice: Molecular Energy Prediction on QM9



Atoms' environments are represented as signals on the sphere (Esteves, Slotine, and Makadia 2023)

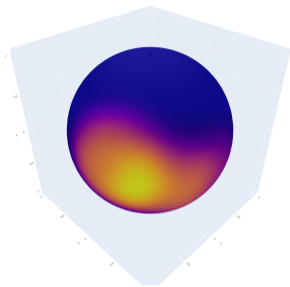
∞ -width limit in practice: Molecular Energy Prediction on QM9



MLP
vs.
SO(3) GCNN

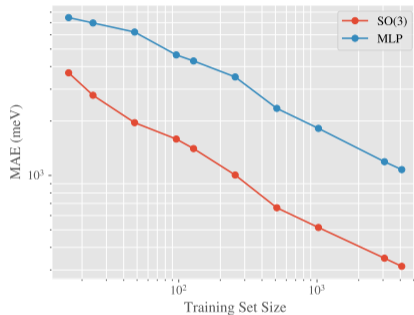
Atoms' environments are represented as signals on the sphere (Estevés, Slotine, and Makadia 2023)

∞ -width limit in practice: Molecular Energy Prediction on QM9

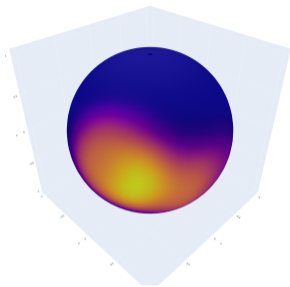


MLP
vs.
SO(3) GCNN

Atoms' environments are represented as signals on the sphere (Estevés, Slotine, and Makadia 2023)

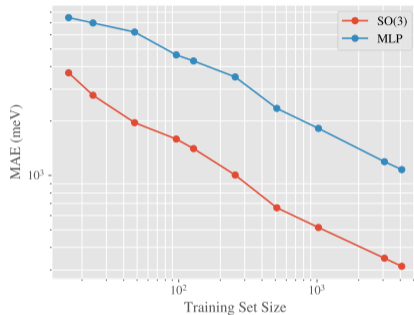


∞ -width limit in practice: Molecular Energy Prediction on QM9



MLP
vs.
SO(3) GCNN

Atoms' environments are represented as signals on the sphere (Estevés, Slotine, and Makadia 2023)



➡ Performance boost due to 3d-rotation invariance extends to the ∞ -width limit

- 1 The Neural Tangent Kernel
- 2 NTK of Equivariant NNs
- 3 Concrete Examples
- 4 Data Augmentation vs GCNNs**

Data Augmentation \leftrightarrow Group Convolutional (GC) NNs

Data Augmentation \leftrightarrow Group Convolutional (GC) NNs

Can construct a GCNN s.t.

$$\Theta^{\text{GC}}(f, f') = \frac{1}{\text{vol}(G)} \int_G dg \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

Data Augmentation \leftrightarrow Group Convolutional (GC) NNs

Can construct a GCNN s.t.

$$\Theta^{\text{GC}}(f, f') = \frac{1}{\text{vol}(G)} \int_G dg \underbrace{\Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')}_{\text{Effective MLP kernel under data augmentation}}$$

Data Augmentation \leftrightarrow Group Convolutional (GC) NNs

Can construct a GCNN s.t.

$$\Theta^{\text{GC}}(f, f') = \underbrace{\frac{1}{\text{vol}(\mathcal{G})} \int_{\mathcal{G}} dg \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')}_{\text{Effective MLP kernel under data augmentation}}$$

At ∞ -width and quadratic \mathcal{L} :



mean of an ensemble of data augmented MLPs equals the **mean** of an ensemble of GCNNs **at all training times t** .

Data Augmentation \leftrightarrow Group Convolutional (GC) NNs

Can construct a GCNN s.t.

$$\Theta^{\text{GC}}(f, f') = \underbrace{\frac{1}{\text{vol}(\mathcal{G})} \int_{\mathcal{G}} dg \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')}_{\text{Effective MLP kernel under data augmentation}}$$

At ∞ -width and quadratic \mathcal{L} :



mean of an ensemble of data augmented MLPs equals the **mean** of an ensemble of GCNNs **at all training times** t .

Data Augmentation \leftrightarrow Group Convolutional (GC) NNs

Can construct a GCNN s.t.

$$\Theta^{\text{GC}}(f, f') = \underbrace{\frac{1}{\text{vol}(\mathcal{G})} \int_{\mathcal{G}} dg \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')}_{\text{Effective MLP kernel under data augmentation}}$$

At ∞ -width and quadratic \mathcal{L} :

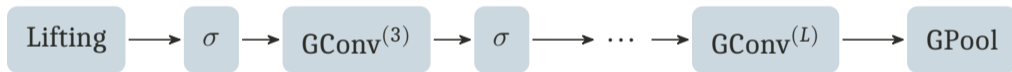
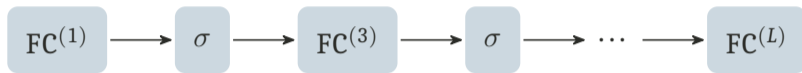


mean of an ensemble of data augmented MLPs equals the **mean** of an ensemble of GCNNs **at all training times t** .

Similar result for data augmented CNN $\leftrightarrow \mathcal{C}_4 \times \mathbb{R}^2$ GCNN

Architecture correspondence

Architecture correspondence



All group convolutions with global filter support $S_{\kappa}^{\ell} = G$ or $S_{\kappa}^1 = X$ for the lifting layer.

Data Augmentation vs. Group Convolutions **at finite width**

Data Augmentation vs. Group Convolutions **at finite width**

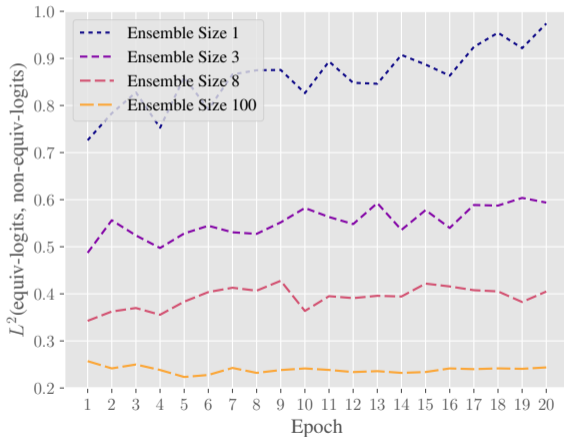
- Data augmented CNN vs $C_4 \times \mathbb{R}^2$
GCNN on **MNIST**

Data Augmentation vs. Group Convolutions **at finite width**

- Data augmented CNN vs $C_4 \times \mathbb{R}^2$
GCNN on **MNIST**
- Compare L_2 -difference of mean logits on **out-of-distribution** data

Data Augmentation vs. Group Convolutions **at finite width**

- Data augmented CNN vs $C_4 \times \mathbb{R}^2$ GCNN on **MNIST**
- Compare L_2 -difference of mean logits on **out-of-distribution** data



Summing up

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**
- Provide explicit expressions for **roto-translations** and **3d-rotations**

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**
- Provide explicit expressions for **roto-translations** and **3d-rotations**
- Showed how this analytic tool can be used to **analyze the connection between data augmentation and GCNNs**

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**
- Provide explicit expressions for **roto-translations** and **3d-rotations**
- Showed how this analytic tool can be used to **analyze the connection between data augmentation and GCNNs**

What's next?

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**
- Provide explicit expressions for **roto-translations** and **3d-rotations**
- Showed how this analytic tool can be used to **analyze the connection between data augmentation and GCNNs**

What's next?

- Relations to Quantum Field Theory (*Banta et al. 2024*)

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**
- Provide explicit expressions for **roto-translations** and **3d-rotations**
- Showed how this analytic tool can be used to **analyze the connection between data augmentation and GCNNs**

What's next?

- Relations to Quantum Field Theory (*Banta et al. 2024*)
- Trainability and generalization regimes of equivariant NNs (*Xiao, Pennington, and Schoenholz 2020*)

Summing up

- **Extended** the Neural Tangent Kernel theory **to equivariant neural networks**
- Provide explicit expressions for **roto-translations** and **3d-rotations**
- Showed how this analytic tool can be used to **analyze the connection between data augmentation and GCNNs**

What's next?

- Relations to Quantum Field Theory (*Banta et al. 2024*)
- Trainability and generalization regimes of equivariant NNs (*Xiao, Pennington, and Schoenholz 2020*)
- equivariant graph NNs

Do you want to know more?



arXiv:2406.06504