

# Emergent Equivariance in Deep Ensembles

Jan E. Gerken



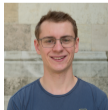
UNIVERSITY OF  
GOTHENBURG

WASP | WALLENBERG AI  
AUTONOMOUS SYSTEMS  
AND SOFTWARE PROGRAM

in collaboration with



Pan Kessel



Philipp Misof

# Data augmentation

👍 Easy to implement

👍 No specialized architecture necessary

# Data augmentation

- 👍 Easy to implement
- 👍 No specialized architecture necessary
- 👎 No exact equivariance

# Data augmentation

- 👍 Easy to implement
- 👍 No specialized architecture necessary
- 👎 No exact equivariance

Can we understand data augmentation theoretically?

# Empirical NTK

Training dynamics under continuous gradient descent:

$$\frac{d\mathcal{N}_{\theta}(x)}{dt} = -\frac{\eta}{N} \sum_{i=1}^N \Theta_{\theta}(x, x_i) \frac{\partial L}{\partial \mathcal{N}(x_i)}$$

learning rate  $\eta$

loss  $L$

training sample  $x_i$

# Empirical NTK

Training dynamics under continuous gradient descent:

$$\frac{d\mathcal{N}_\theta(x)}{dt} = -\frac{\eta}{N} \sum_{i=1}^N \Theta_\theta(x, x_i) \frac{\partial L}{\partial \mathcal{N}(x_i)}$$

learning rate  $\eta$  (indicated by a blue arrow pointing to the fraction)

loss  $L$  (indicated by a blue arrow pointing to the derivative)

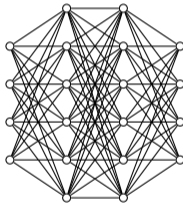
training sample  $x_i$  (indicated by a blue arrow pointing to the kernel argument)

with the **empirical neural tangent kernel (NTK)**

$$\Theta_\theta(x, x') = \sum_{\mu} \frac{\partial \mathcal{N}(x)}{\partial \theta_{\mu}} \frac{\partial \mathcal{N}(x')}{\partial \theta_{\mu}}$$

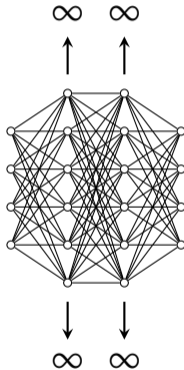
# Infinite width limit

[Jacot et al. 2018]



# Infinite width limit

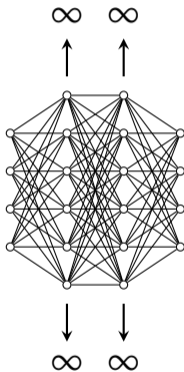
[Jacot et al. 2018]





# Infinite width limit

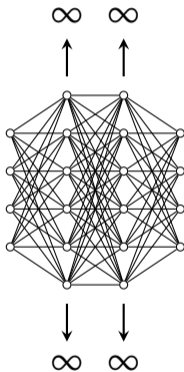
[Jacot et al. 2018]



👍 NTK becomes independent of initialization

# Infinite width limit

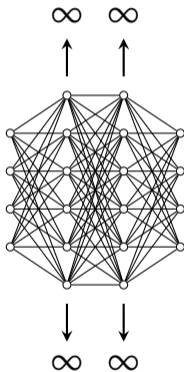
[Jacot et al. 2018]



- 👍 NTK becomes independent of initialization
- 👍 NTK becomes constant in training

# Infinite width limit

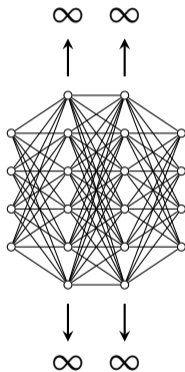
[Jacot et al. 2018]



- 👍 NTK becomes independent of initialization
- 👍 NTK becomes constant in training
- 👍 NTK can be computed for most networks

# Infinite width limit

[Jacot et al. 2018]



- 👍 NTK becomes independent of initialization
- 👍 NTK becomes constant in training
- 👍 NTK can be computed for most networks
- ✓ Training dynamics can be solved

# Mean prediction from NTK

[Jacot et al. 2018]

ⓘ At infinite width, the mean prediction is given by

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{I} - e^{-\eta \Theta(X, X) t}) \gamma$$

# Mean prediction from NTK

[Jacot et al. 2018]

① At infinite width, the mean prediction is given by

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{I} - e^{-\eta \Theta(X, X) t}) \gamma$$

neural tangent kernel

# Mean prediction from NTK

[Jacot et al. 2018]

ⓘ At infinite width, the mean prediction is given by

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{I} - e^{-\eta \Theta(X, X) t}) Y$$

neural tangent kernel

train data

# Mean prediction from NTK

[Jacot et al. 2018]

ⓘ At infinite width, the mean prediction is given by

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{I} - e^{-\eta \Theta(X, X) t}) \gamma$$

neural tangent kernel

learning rate

train data



# Mean prediction from NTK

[Jacot et al. 2018]

ⓘ At infinite width, the mean prediction is given by

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{I} - e^{-\eta \Theta(X, X) t}) Y$$

neural tangent kernel

train labels

learning rate

train data

# **Data augmentation**

## Data augmentation at infinite width

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{1} - e^{-\eta \Theta(X, X) t}) \gamma$$

# Data augmentation at infinite width

$$\mu_t(x) = \Theta(x, X) \Theta(X, X)^{-1} (\mathbb{1} - e^{-\eta \Theta(X, X) t}) Y$$

The diagram illustrates the components of the equation. The text "augmented data" is positioned to the left of the equation, with three blue arrows pointing to the terms  $\Theta(x, X)$ ,  $\Theta(X, X)^{-1}$ , and  $\Theta(X, X)$  in the equation. The text "augmented labels" is positioned below the equation, with two blue arrows pointing to the terms  $\mathbb{1}$  and  $Y$  in the equation.

# Data augmentation at infinite width

group transformation

$$\mu_t(\rho(g)x) = \Theta(\rho(g)x, X)\Theta(X, X)^{-1}(\mathbb{I} - e^{-\eta\Theta(X, X)t})Y$$

augmented data

augmented labels

# Data augmentation at infinite width

group transformation for augmented data

$$\mu_t(\rho(g)x) = \Theta(\rho(g)x, X) \Theta(X, X)^{-1} (\mathbb{I} - e^{-\eta \Theta(X, X)t}) Y$$

augmented data augmented labels

# Data augmentation at infinite width

group transformation

$$\mu_t(\rho(g)x) = \Theta(x, X)\Theta(X, X)^{-1}(\mathbb{I} - e^{-\eta\Theta(X, X)t})\rho(g)Y$$

augmented data

augmented labels

The diagram illustrates the equation for data augmentation at infinite width. The equation is  $\mu_t(\rho(g)x) = \Theta(x, X)\Theta(X, X)^{-1}(\mathbb{I} - e^{-\eta\Theta(X, X)t})\rho(g)Y$ . A blue arrow labeled "group transformation" points to the  $\rho(g)$  term. Below the equation, the text "augmented data" has three blue arrows pointing to the  $\Theta(x, X)$ ,  $\Theta(X, X)^{-1}$ , and  $\Theta(X, X)$  terms. The text "augmented labels" has two blue arrows pointing to the  $\Theta(X, X)$  term and the  $\rho(g)$  term. A final blue arrow points from "augmented labels" to the  $Y$  term.

# Data augmentation at infinite width

group transformation

augmented labels

$$\mu_t(\rho(g)x) = \Theta(x, X)\Theta(X, X)^{-1}(\mathbb{I} - e^{-\eta\Theta(X, X)t})\underbrace{\rho(g)Y}_{=Y}$$

for invariance



# Data augmentation at infinite width

group transformation

$$\begin{aligned}\mu_t(\rho(g)x) &= \Theta(x, X)\Theta(X, X)^{-1}(\mathbb{I} - e^{-\eta\Theta(X, X)t})\underbrace{\rho(g)Y}_{=Y} \\ &= \mu_t(x)\end{aligned}$$

for invariance

# Mean prediction

$$\mu_t(x)$$

# Mean prediction

$$\mu_t(x) = \mathbb{E}_{\theta_0 \sim \text{initializations}}[\mathcal{N}_{\theta_t}(x)]$$

# Mean prediction

$$\mu_t(x) = \mathbb{E}_{\theta_0 \sim \text{initializations}}[\mathcal{N}_{\theta_t}(x)] = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{\theta_0 = \text{init}_1}^{\text{init}_n} \mathcal{N}_{\theta_t}(x)$$

# Mean prediction

$$\mu_t(x) = \mathbb{E}_{\theta_0 \sim \text{initializations}}[\mathcal{N}_{\theta_t}(x)] = \lim_{n \rightarrow \infty} \underbrace{\frac{1}{n} \sum_{\theta_0 = \text{init}_1}^{\text{init}_n} \mathcal{N}_{\theta_t}(x)}_{\text{mean prediction of deep ensemble}}$$

## **Main conclusion**

Deep ensembles trained with data augmentation are equivariant.

# Main conclusion

Deep ensembles trained with data augmentation are equivariant.

- ✓ Proof of exact equivariance for
  - full data augmentation
  - infinite ensembles

# Main conclusion

Deep ensembles trained with data augmentation are equivariant.

- ✓ Proof of exact equivariance for
  - full data augmentation
  - infinite ensembles
- ✓ Equivariance holds for all training times



## Main conclusion

Deep ensembles trained with data augmentation are equivariant.

- ✓ Proof of exact equivariance for
  - full data augmentation
  - infinite ensembles
- ✓ Equivariance holds for all training times
- ✓ Equivariance holds away from the training data

# Main conclusion

Deep ensembles trained with data augmentation are equivariant.

- ✓ Proof of exact equivariance for
  - full data augmentation
  - infinite ensembles
- ✓ Equivariance holds for all training times
- ✓ Equivariance holds away from the training data
- ✓ Holds also for finite-width networks

[Nordenfors, Flinth 2024]

# Intuitive explanation

- ✓ Equivariance holds for all training times
- ✓ Equivariance holds away from the training data

# Intuitive explanation

- ✓ Equivariance holds for all training times
- ✓ Equivariance holds away from the training data

⊙ At infinite width, the mean output at initialization is zero everywhere.

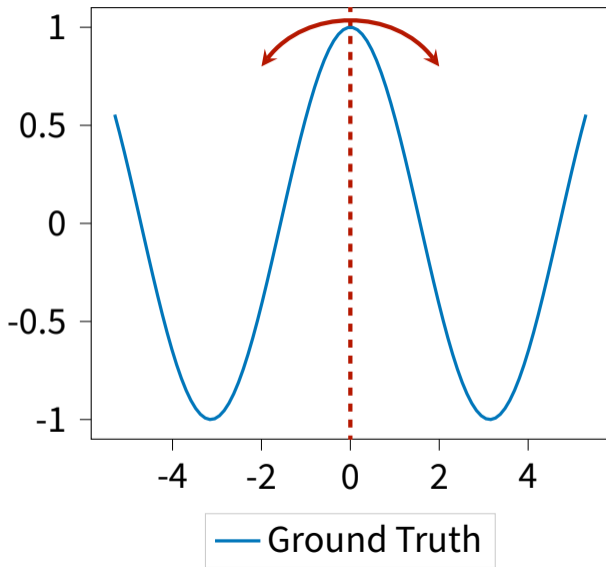
# Intuitive explanation

- ✓ Equivariance holds for all training times
- ✓ Equivariance holds away from the training data

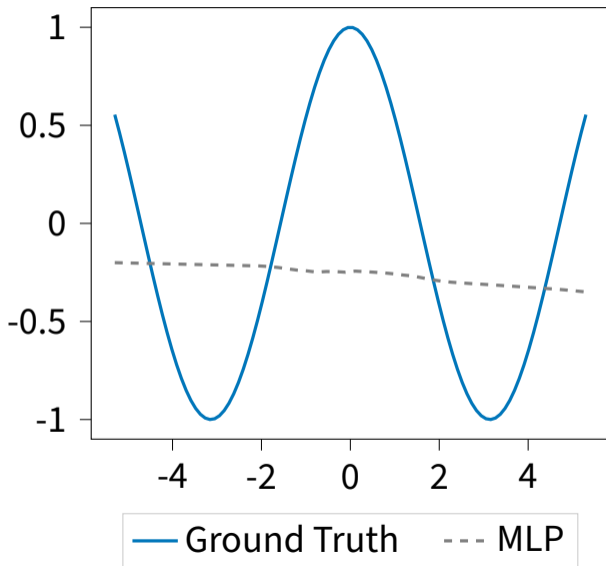
⊙ At infinite width, the mean output at initialization is zero everywhere.

⇒ Training with full data augmentation leads to an equivariant function.

# Toy example

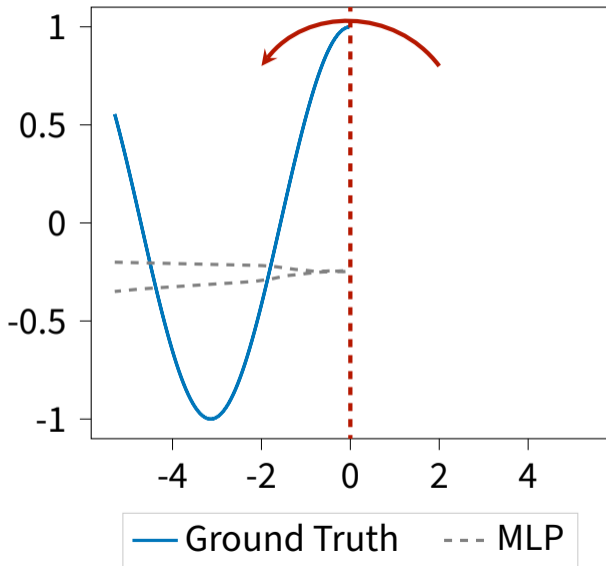


## Initialization

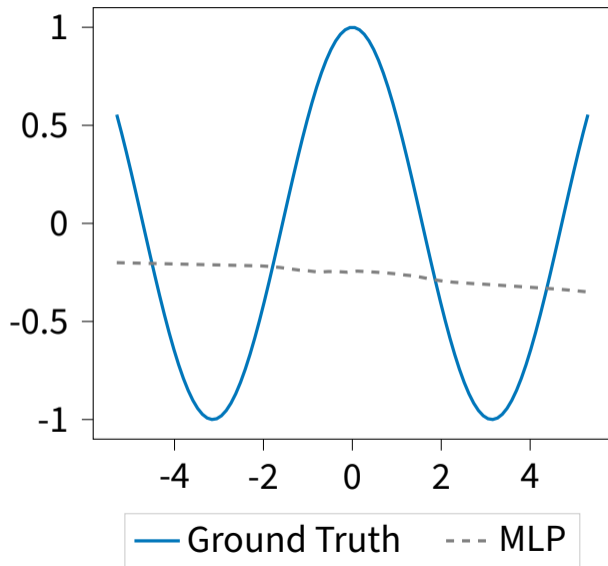




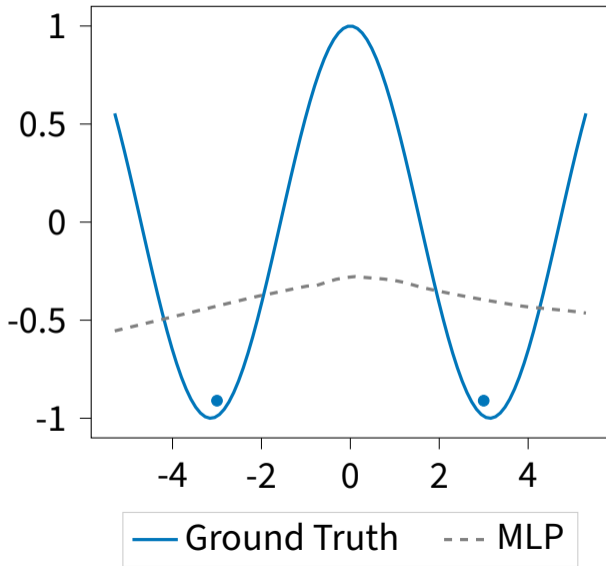
## Initialization



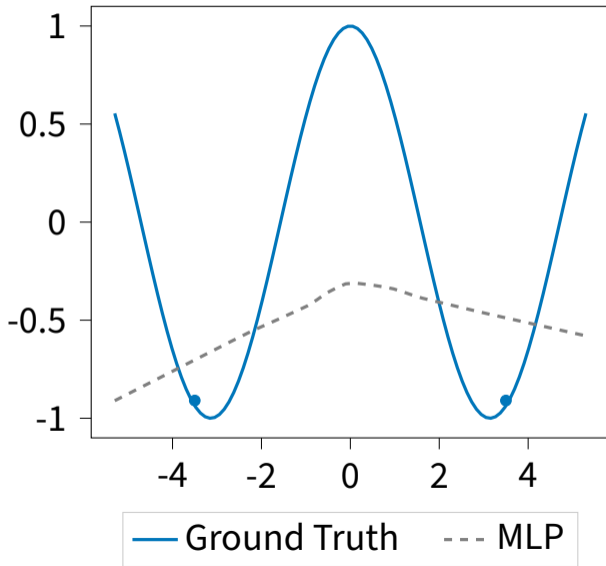
## Initialization



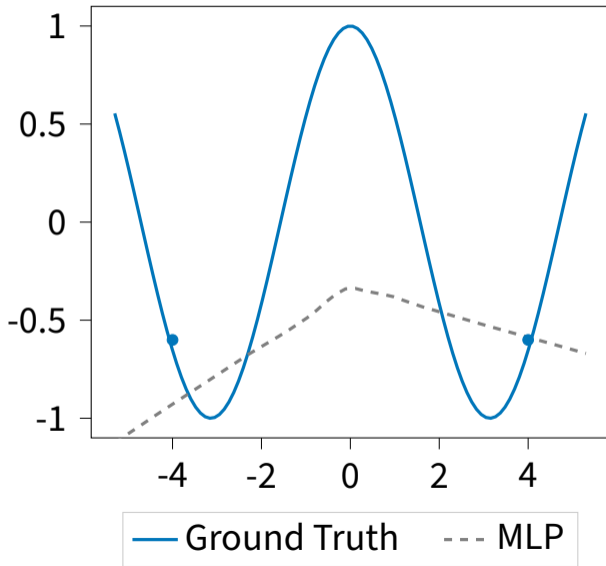
## After 1 Training Step



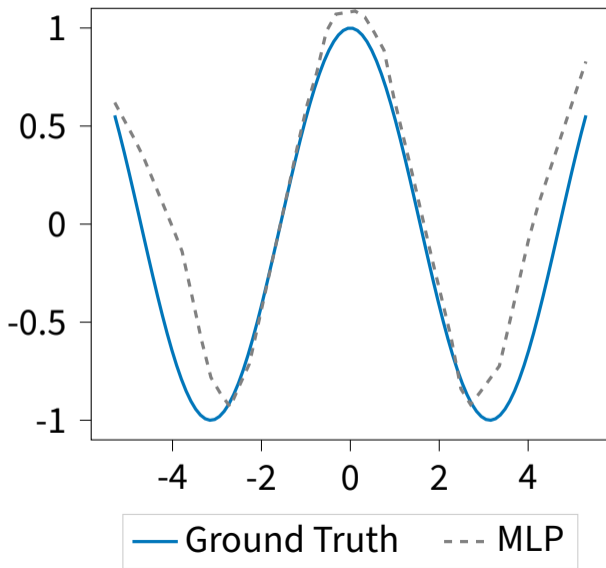
## After 2 Training Steps



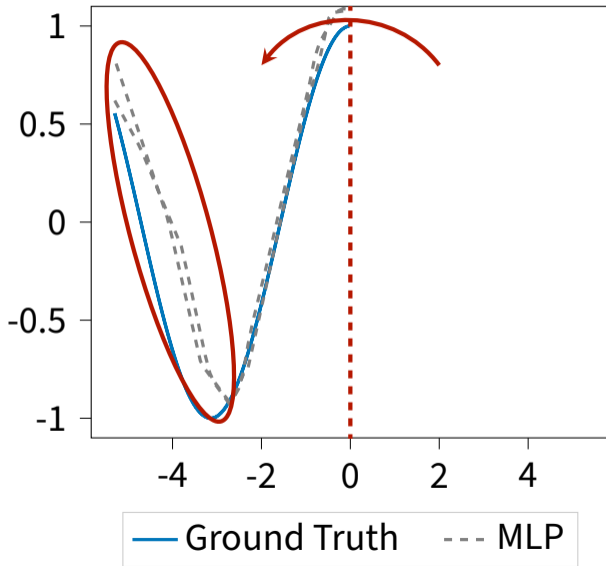
## After 3 Training Steps



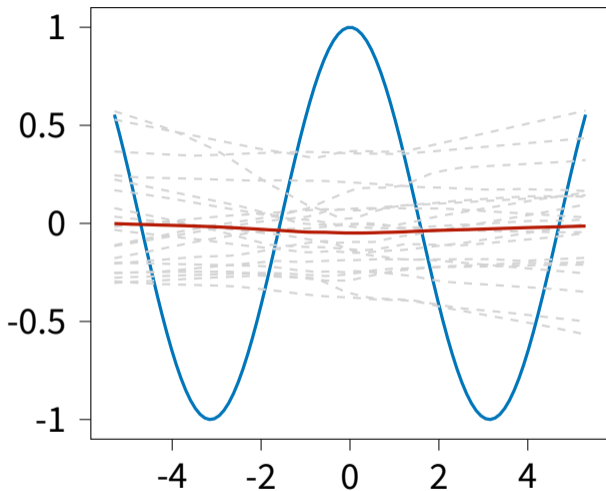
## After 2000 Training Steps



# After 2000 Training Steps



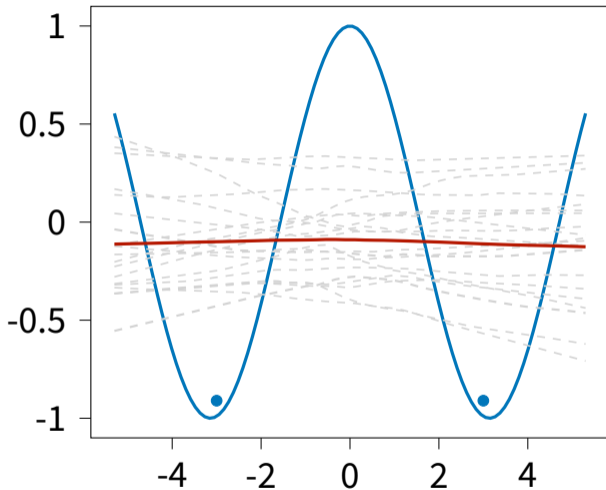
## Initialization



— Ground Truth    - - - MLP    — Ensemble Mean

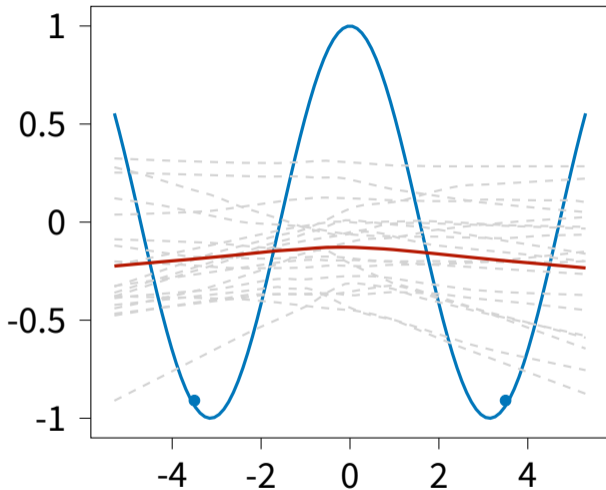


## After 1 Training Step



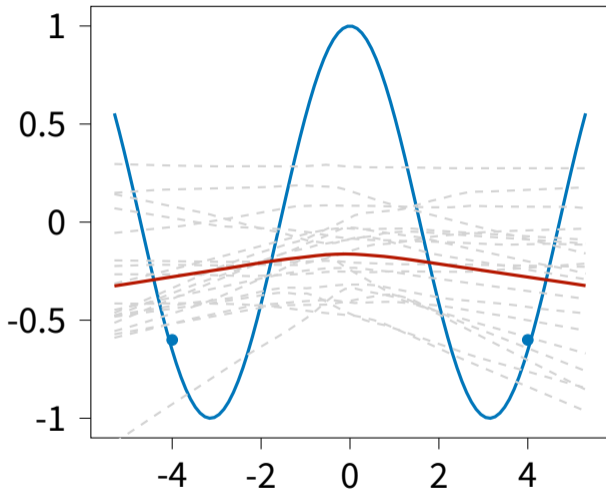
— Ground Truth    - - - MLP    — Ensemble Mean

## After 2 Training Steps



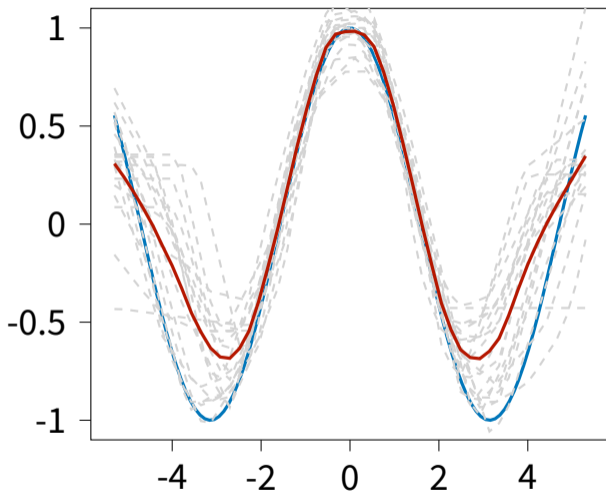
— Ground Truth    - - - MLP    — Ensemble Mean

## After 3 Training Steps



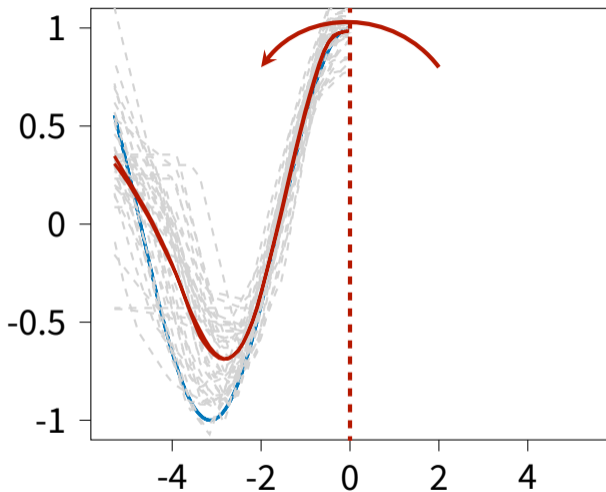
— Ground Truth    - - - MLP    — Ensemble Mean

## After 2000 Training Steps



— Ground Truth    - - - MLP    — Ensemble Mean

## After 2000 Training Steps

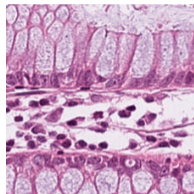
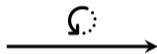
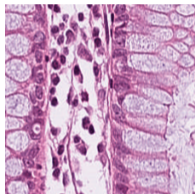


— Ground Truth    - - - MLP    — Ensemble Mean

# **What Does An Augmented Ensemble Converge To?**

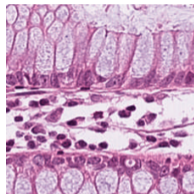
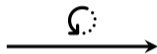
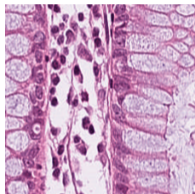
# Rotating images

# Rotating images





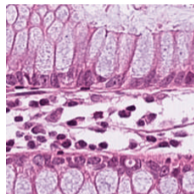
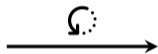
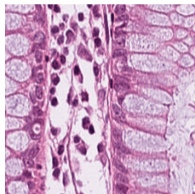
# Rotating images



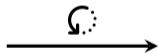
$f(x)$

$f : \text{pixels} \rightarrow \text{colors}$

# Rotating images



$f(x)$



$f : \text{pixels} \rightarrow \text{colors}$

$f(\rho(g^{-1})x)$

$= [\rho_{\text{reg}}(g)f](x)$

# Data augmentation and NTKs

# Data augmentation and NTKs

Consider two ensembles:

trained without data augmentation

trained with data augmentation

# Data augmentation and NTKs

Consider two ensembles:

trained without data augmentation

trained with data augmentation

If

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

# Data augmentation and NTKs

Consider two ensembles:

trained without data augmentation

trained with data augmentation

If

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

Then

$$\mu_t^{\text{non-aug}}(x) = \mu_t^{\text{aug}}(x)$$

at infinite width.

# Data augmentation and NTKs

Consider two ensembles:

trained without data augmentation

trained with data augmentation

If

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

Then

$$\mu_t^{\text{non-aug}}(x) = \mu_t^{\text{aug}}(x) \quad \forall t$$

at infinite width.

# Data augmentation and NTKs

Consider two ensembles:

trained without data augmentation

trained with data augmentation

If

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

Then

$$\mu_t^{\text{non-aug}}(x) = \mu_t^{\text{aug}}(x) \quad \forall t \quad \forall x$$

at infinite width.



## Data augmentation and NTKs

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

## Data augmentation and NTKs

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

- Ⓢ Given an architecture with NTK  $\Theta^{\text{aug}}$ ,  
find an architecture with NTK  $\Theta^{\text{non-aug}}$

# Group convolutions

[Cohen, Welling 2016]

# Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt  $\rho_{\text{reg}}$

# Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt  $\rho_{\text{reg}}$

- Ordinary convolutions

$$f'(y) = \int_X dx \kappa(x - y) f(x)$$

# Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt  $\rho_{\text{reg}}$

- Ordinary convolutions

$$f'(y) = \int_X dx \kappa(x - y) f(x)$$

- Group convolutions

$$f'(g) = \int_X dx \kappa(\rho(g^{-1})x) f(x)$$

lifting

# Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt  $\rho_{\text{reg}}$

- Ordinary convolutions

$$f'(y) = \int_X dx \kappa(x - y) f(x)$$

- Group convolutions

$$f'(g) = \int_X dx \kappa(\rho(g^{-1})x) f(x) \quad \text{lifting}$$

$$f'(g) = \int_G dg \kappa(g^{-1}h) f(h) \quad \text{group convolution}$$

# Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt  $\rho_{\text{reg}}$

- Ordinary convolutions

$$f'(y) = \int_X dx \kappa(x - y) f(x)$$

- Group convolutions

$$f'(g) = \int_X dx \kappa(\rho(g^{-1})x) f(x) \quad \text{lifting}$$

$$f'(g) = \int_G dg \kappa(g^{-1}h) f(h) \quad \text{group convolution}$$

$$f' = \frac{1}{\text{vol}(G)} \int_G dg f(g) \quad \text{group pooling}$$



# GCNNs

Stack GConv-layers to obtain an invariant network

# GCNNs

Stack GConv-layers to obtain an invariant network



# GCNNs

Stack GConv-layers to obtain an invariant network



→ lifting

# GCNNs

Stack GConv-layers to obtain an invariant network



→ lifting →  $\sigma$

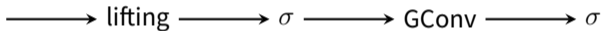
# GCNNs

Stack GConv-layers to obtain an invariant network



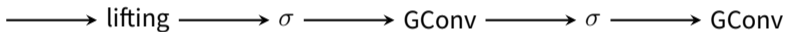
# GCNNs

Stack GConv-layers to obtain an invariant network



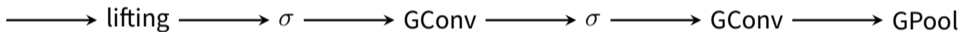
# GCNNs

Stack GConv-layers to obtain an invariant network



# GCNNs

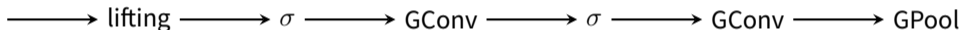
Stack GConv-layers to obtain an invariant network





# NTKs for GCNNs

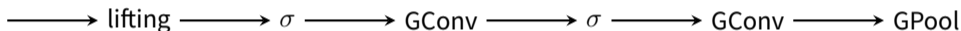
Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network

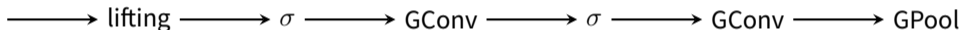


Compute NTK with layer-wise recursion

0

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network

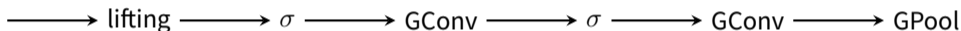


Compute NTK with layer-wise recursion

$$0 \longrightarrow \Theta_{g,g'}^{(1)}(f,f')$$

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network

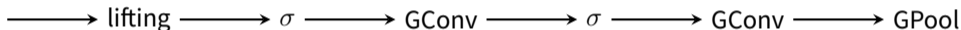


Compute NTK with layer-wise recursion

$$0 \longrightarrow \Theta_{g,g'}^{(1)}(f,f') \longrightarrow \Theta_{g,g'}^{(2)}(f,f')$$

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network

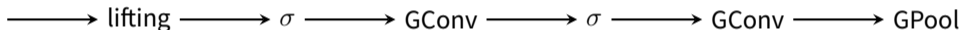


Compute NTK with layer-wise recursion

$$0 \longrightarrow \Theta_{g,g'}^{(1)}(f,f') \longrightarrow \Theta_{g,g'}^{(2)}(f,f') \longrightarrow \Theta_{g,g'}^{(3)}(f,f')$$

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network

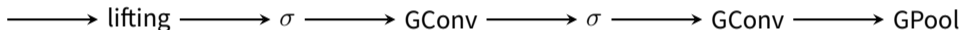


Compute NTK with layer-wise recursion

$$0 \longrightarrow \Theta_{g,g'}^{(1)}(f,f') \longrightarrow \Theta_{g,g'}^{(2)}(f,f') \longrightarrow \Theta_{g,g'}^{(3)}(f,f') \longrightarrow \Theta_{g,g'}^{(4)}(f,f')$$

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network

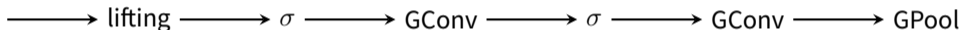


Compute NTK with layer-wise recursion

$$0 \longrightarrow \Theta_{g,g'}^{(1)}(f,f') \longrightarrow \Theta_{g,g'}^{(2)}(f,f') \longrightarrow \Theta_{g,g'}^{(3)}(f,f') \longrightarrow \Theta_{g,g'}^{(4)}(f,f') \longrightarrow \Theta_{g,g'}^{(5)}(f,f')$$

# NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

$$0 \longrightarrow \Theta_{g,g'}^{(1)}(f,f') \longrightarrow \Theta_{g,g'}^{(2)}(f,f') \longrightarrow \Theta_{g,g'}^{(3)}(f,f') \longrightarrow \Theta_{g,g'}^{(4)}(f,f') \longrightarrow \Theta_{g,g'}^{(5)}(f,f') \longrightarrow \Theta(f,f')$$



# NTKs of MLPs and GCNNs

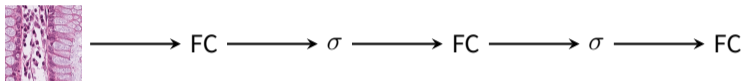
# NTKs of MLPs and GCNNs

- Consider two neural networks

# NTKs of MLPs and GCNNs

- Consider two neural networks

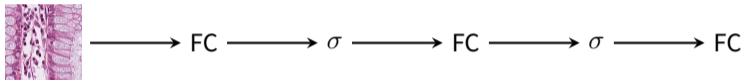
An MLP



# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



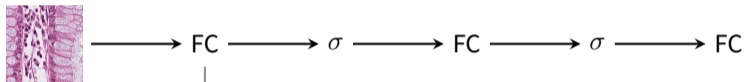
A GCNN



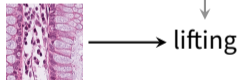
# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



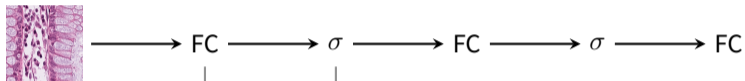
A GCNN



# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



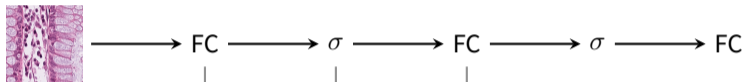
A GCNN



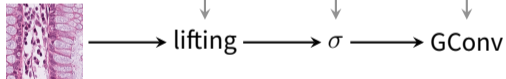
# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



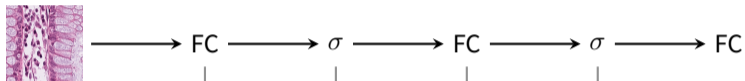
A GCNN



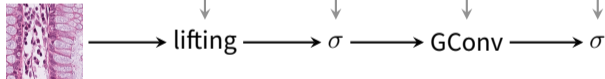
# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



A GCNN





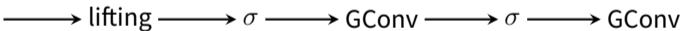
# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



A GCNN



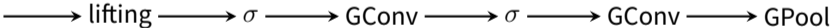
# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



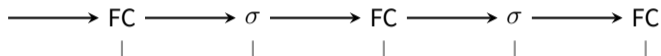
A GCNN



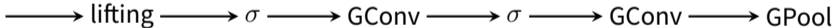
# NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



A GCNN



- Then


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

## Data augmentation of MLPs

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

# Data augmentation of MLPs

before: non-aug


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

# Data augmentation of MLPs

before: non-aug

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

before: aug

# Data augmentation of MLPs

before: non-aug

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

before: aug

⇒ training the MLP on  
 $G$ -augmented data

# Data augmentation of MLPs

before: non-aug

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

before: aug

⇒ training the MLP on  $G$ -augmented data = training the GCNN on unaugmented data



# Data augmentation of MLPs

before: non-aug

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

before: aug

⇒ training the MLP on  $G$ -augmented data = training the GCNN on unaugmented data  
in the ensemble mean

# Data augmentation of MLPs

before: non-aug

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

before: aug

⇒ training the MLP on  
G-augmented data

=

training the GCNN on  
unaugmented data

in the ensemble mean,  $\forall t, \forall x$

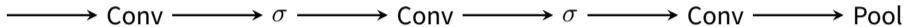
# Data augmentation of CNNs

## Data augmentation of CNNs

- Consider a CNN and a GCNN invariant wrt. roto-translations

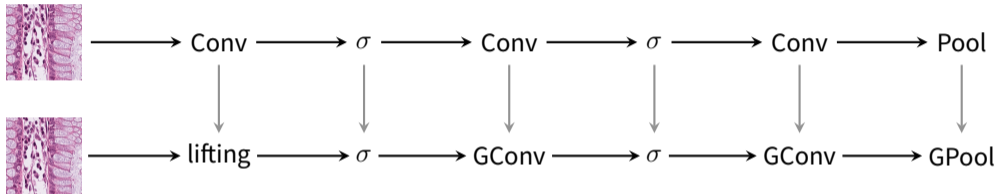
# Data augmentation of CNNs

- Consider a CNN and a GCNN invariant wrt. roto-translations



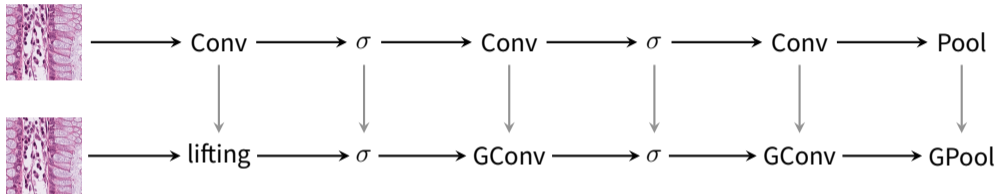
# Data augmentation of CNNs

- Consider a CNN and a GCNN invariant wrt. roto-translations



# Data augmentation of CNNs

- Consider a CNN and a GCNN invariant wrt. roto-translations

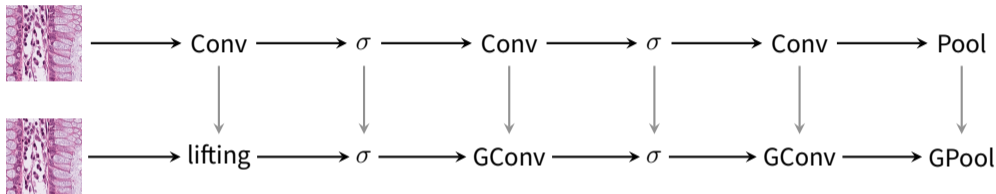


- Then

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{n} \sum_{r \in C_n} \Theta^{\text{CNN}}(f, \rho_{\text{reg}}(r)f')$$

# Data augmentation of CNNs

- Consider a CNN and a GCNN invariant wrt. roto-translations



- Then

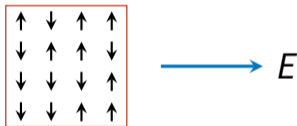
$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{n} \sum_{r \in C_n} \Theta^{\text{CNN}}(f, \rho_{\text{reg}}(r)f')$$

⇒ By training the CNN on rotated images, one obtains a roto-translation invariant GCNN

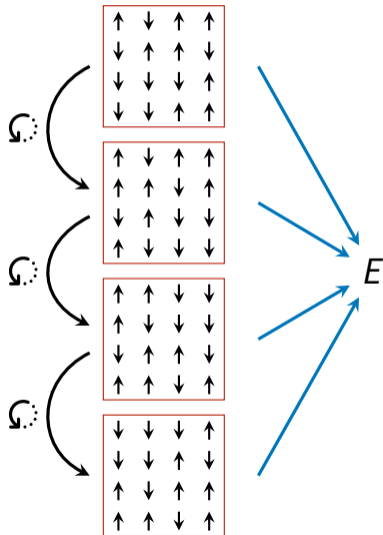


# Experiments

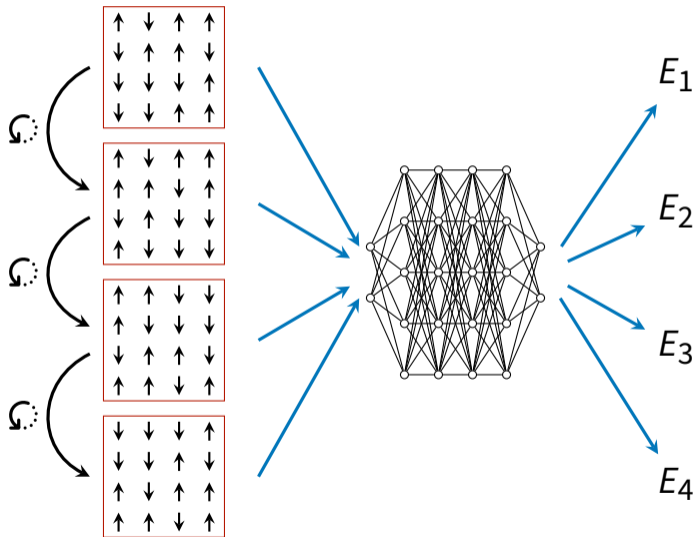
# Ising model



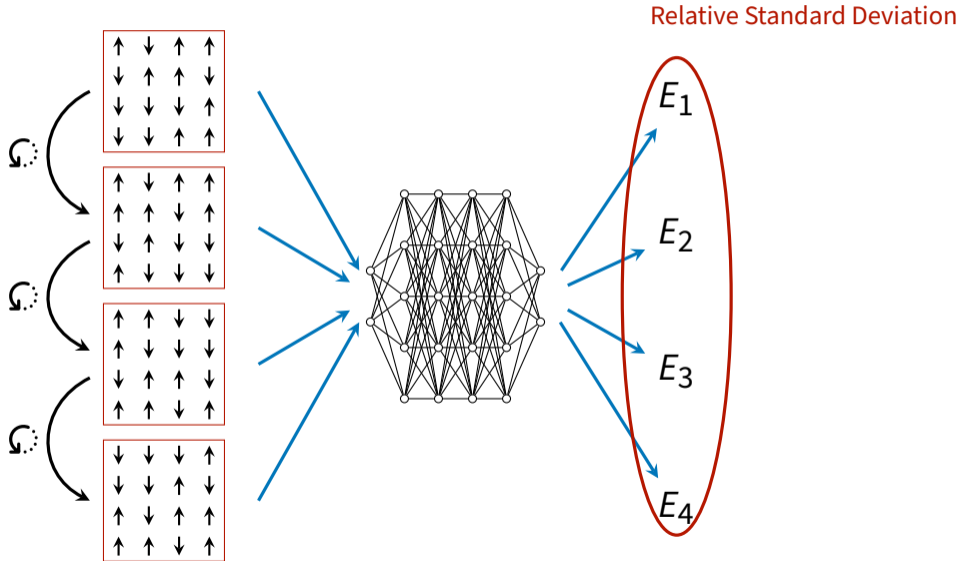
# Ising model

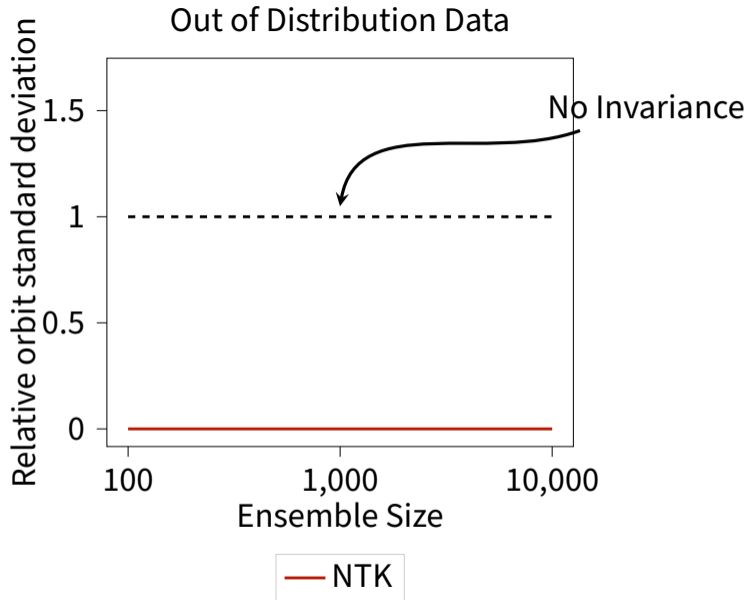


# Ising model

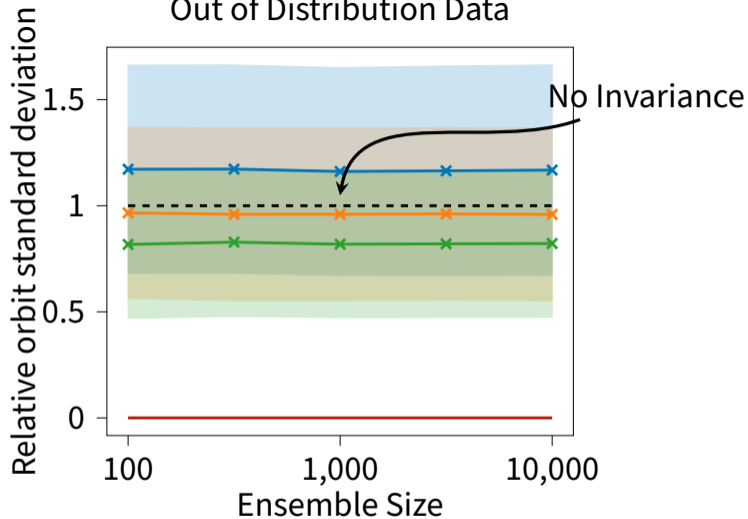


# Ising model



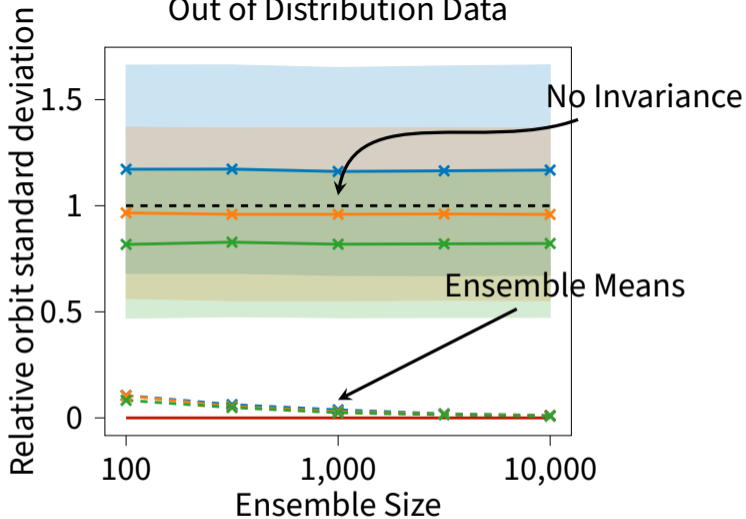


## Out of Distribution Data



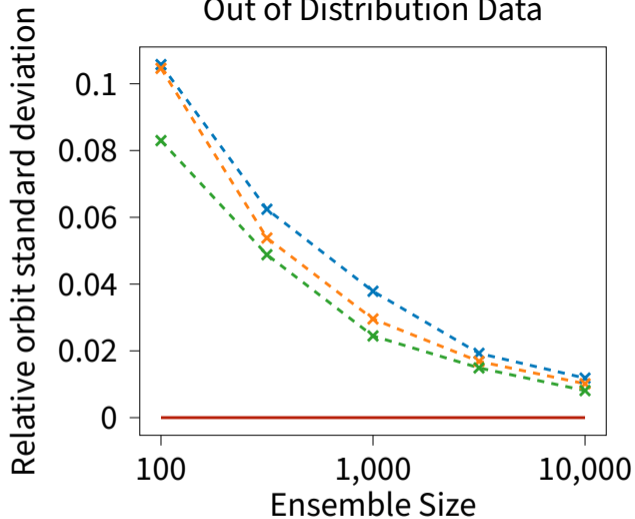
— NTK    × Width 512    × Width 1024    × Width 2048

### Out of Distribution Data





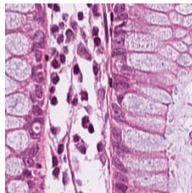
## Out of Distribution Data



— NTK    -x- Width 512    -x- Width 1024    -x- Width 2048

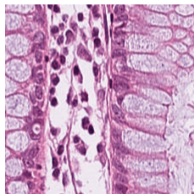
# Histological slices

[Kather et al. 2018]



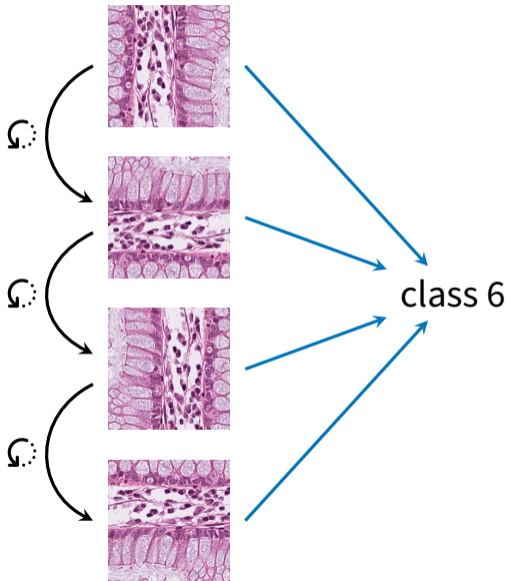
# Histological slices

[Kather et al. 2018]

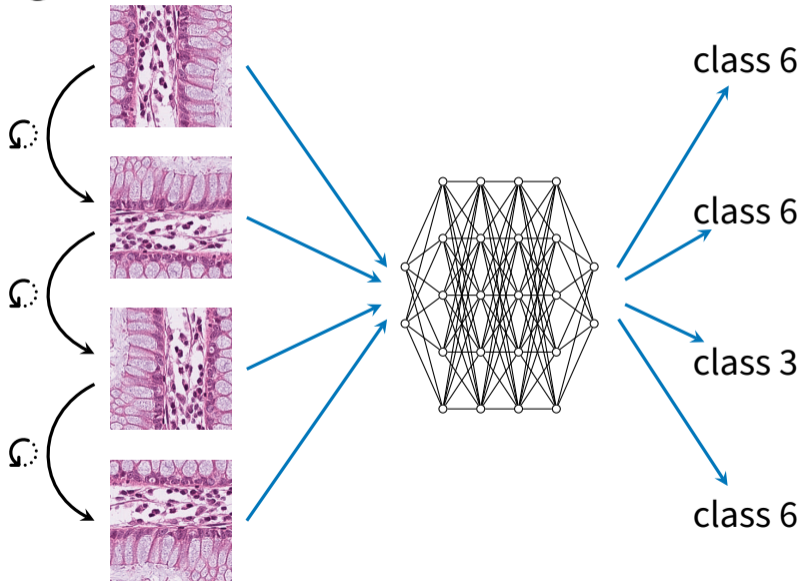


→ class 6

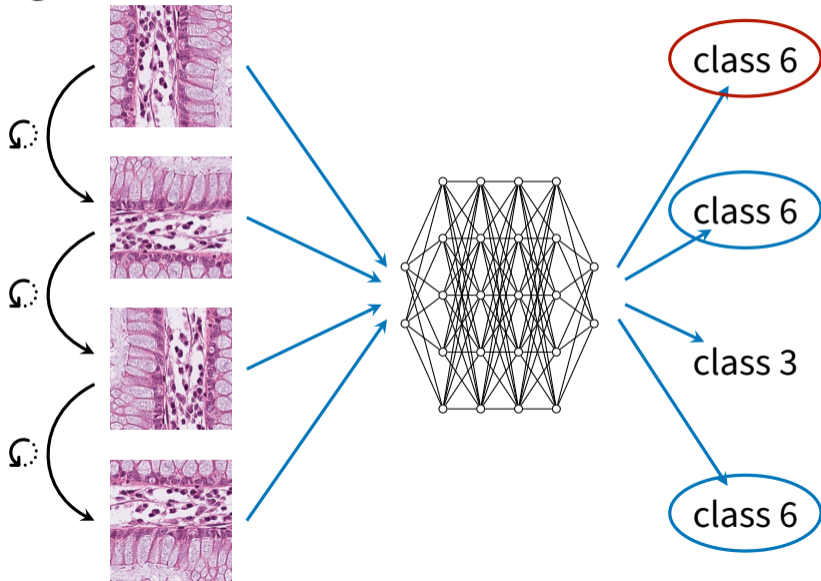
# Histological slices



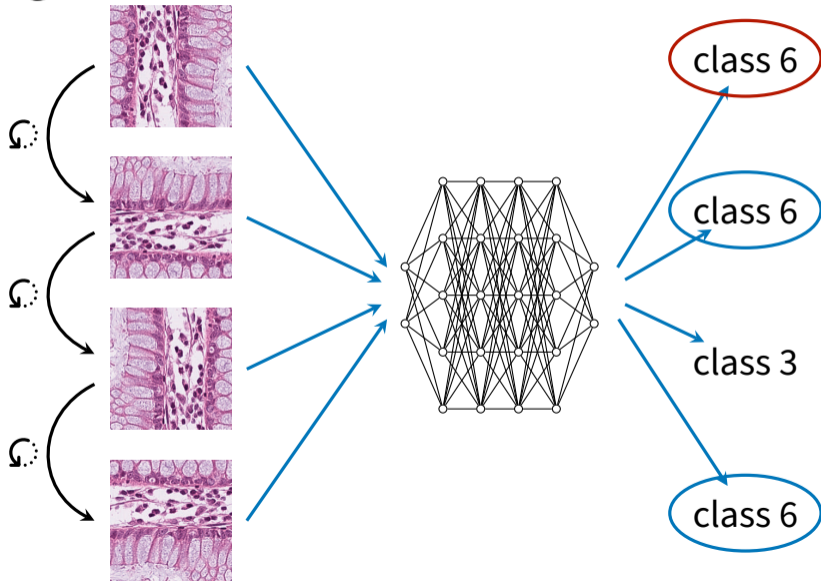
# Histological slices



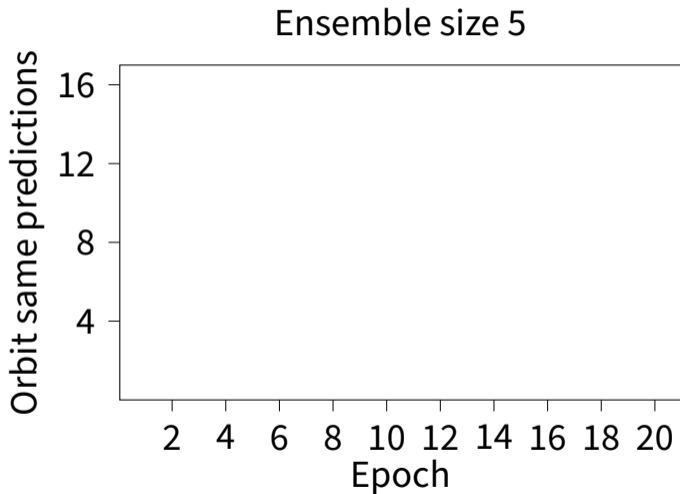
# Histological slices



# Histological slices

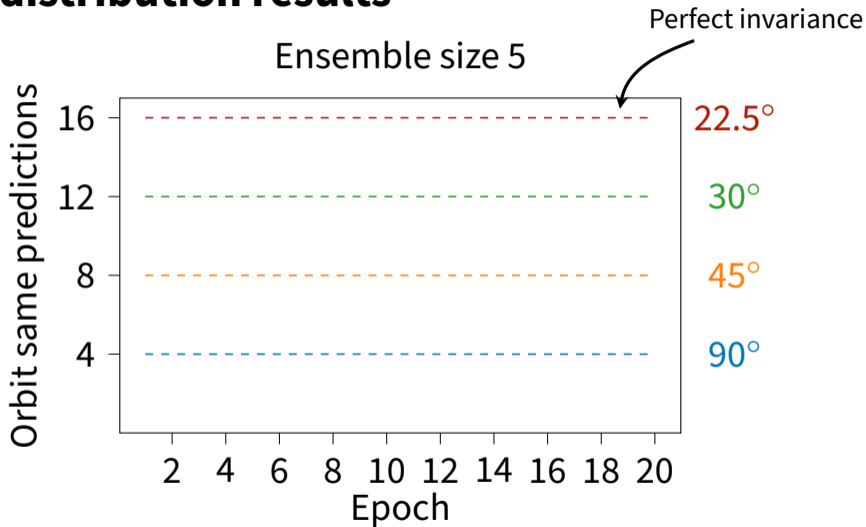


# Out of distribution results

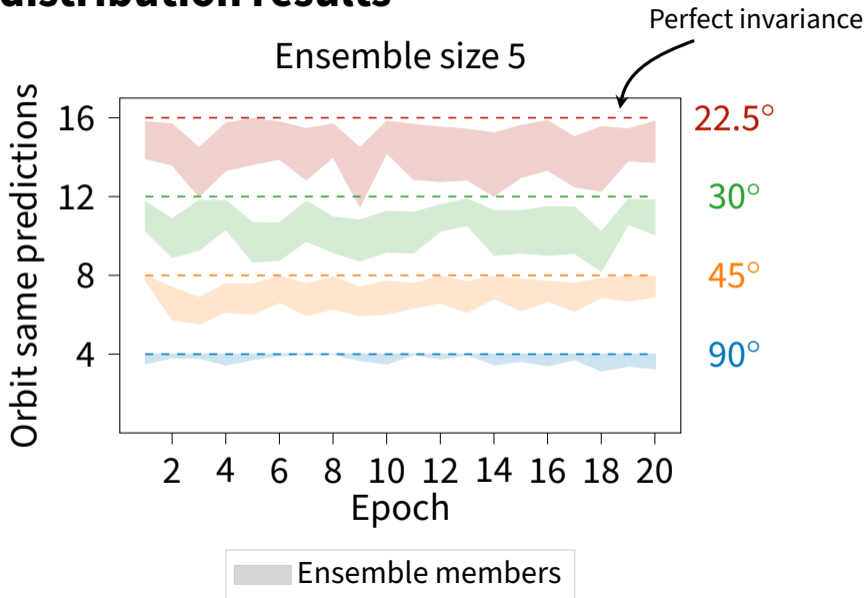




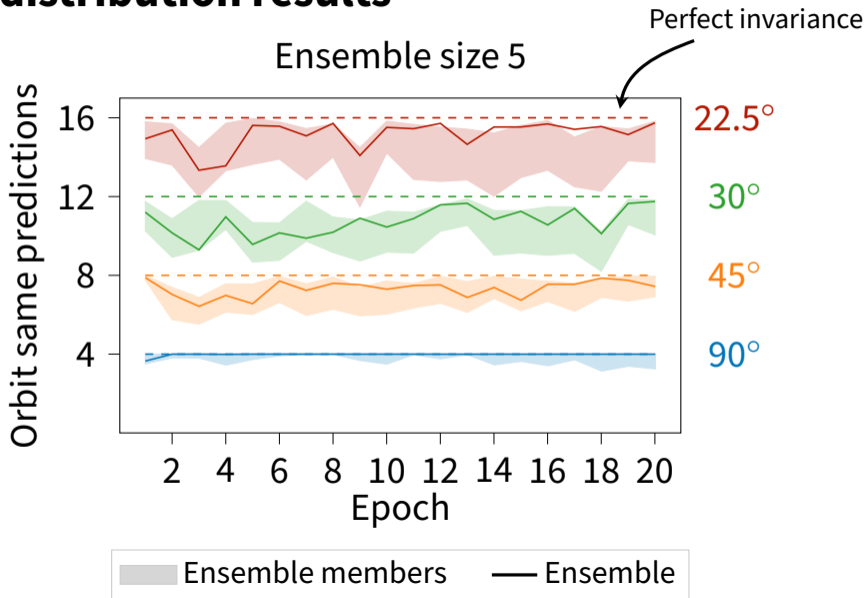
# Out of distribution results



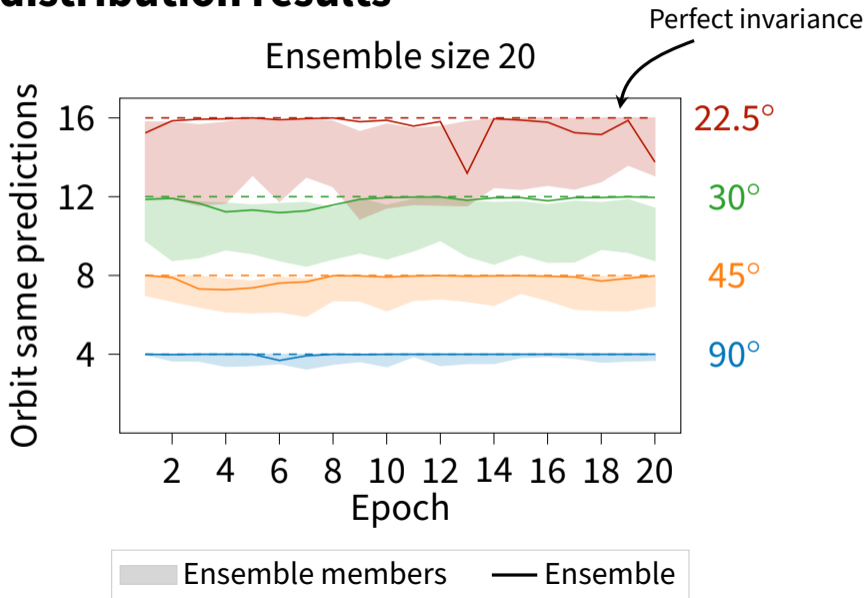
# Out of distribution results



# Out of distribution results



# Out of distribution results



# Further experimental results

## Further experimental results

- ✓ Emergent invariance for rotated FashionMNIST

## Further experimental results

- ✓ Emergent invariance for rotated FashionMNIST
- ✓ Partial augmentation for continuous symmetries

## Further experimental results

- ✓ Emergent invariance for rotated FashionMNIST
- ✓ Partial augmentation for continuous symmetries
- ✓ Emergent equivariance (as opposed to invariance)



# Comparison to other methods

## Comparison to other methods

⇒ Models trained on rotated FashionMNIST

## Comparison to other methods

⇒ Models trained on rotated FashionMNIST

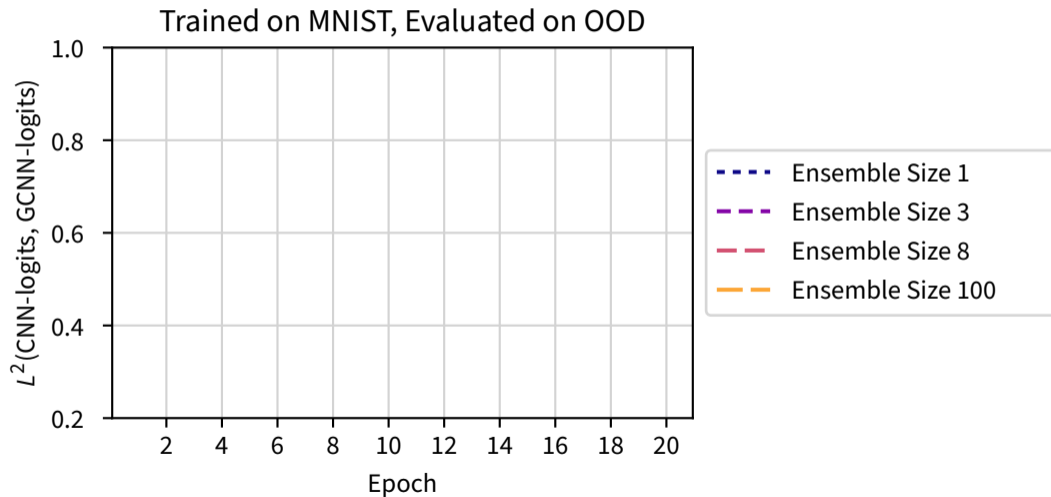
Orbit same predictions out of distribution:

	$C_4$	$C_8$	$C_{16}$
DeepEns+DA	$3.85 \pm 0.12$	<b><math>7.72 \pm 0.34</math></b>	<b><math>15.24 \pm 0.69</math></b>
only DA	$3.41 \pm 0.18$	$6.73 \pm 0.24$	$12.77 \pm 0.71$
E2CNN <sup>1</sup>	<b><math>4 \pm 0.0</math></b>	<b><math>7.71 \pm 0.21</math></b>	<b><math>15.08 \pm 0.34</math></b>
Canon <sup>2</sup>	<b><math>4 \pm 0.0</math></b>	<b><math>7.45 \pm 0.14</math></b>	$12.41 \pm 0.85$

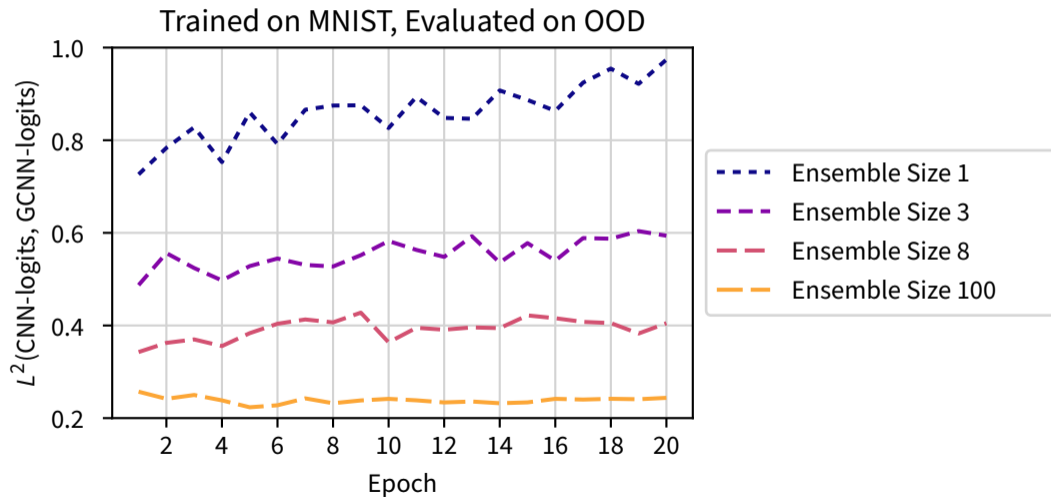
<sup>1</sup>[Weiler et al. 2019], <sup>2</sup>[Kaba et al. 2022]

# Convergence of augmented CNNs to GCNNs

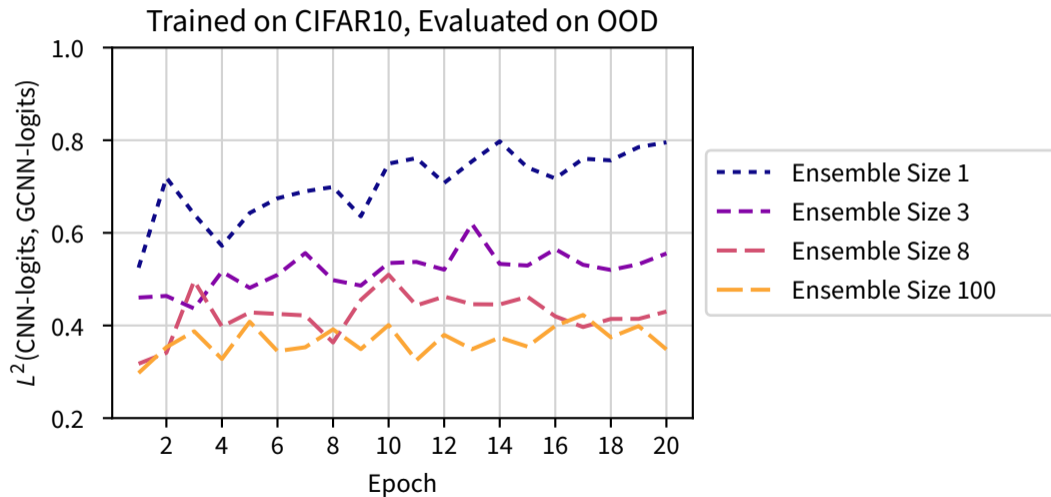
# Convergence of augmented CNNs to GCNNs



# Convergence of augmented CNNs to GCNNs



# Convergence of augmented CNNs to GCNNs



# Key takeaways



## Key takeaways

If you need ensembles

- 👍 use data augmentation to obtain an equivariant model.

## Key takeaways

If you need ensembles

👍 use data augmentation to obtain an equivariant model.

If you need data augmentation

👍 use an ensemble to boost the equivariance.

## Key takeaways

If you need ensembles

👍 use data augmentation to obtain an equivariant model.

If you need data augmentation

👍 use an ensemble to boost the equivariance.

Analysis of neural tangent kernel can lead to powerful practical insights!

# Papers

- [Emergent Equivariance in Deep Ensembles](#)

Jan E. Gerken\*, Pan Kessel\*

ICML 2024 (Oral)

\* Equal contribution

- [Equivariant Neural Tangent Kernels](#)

Philipp Misof, Pan Kessel, Jan E. Gerken

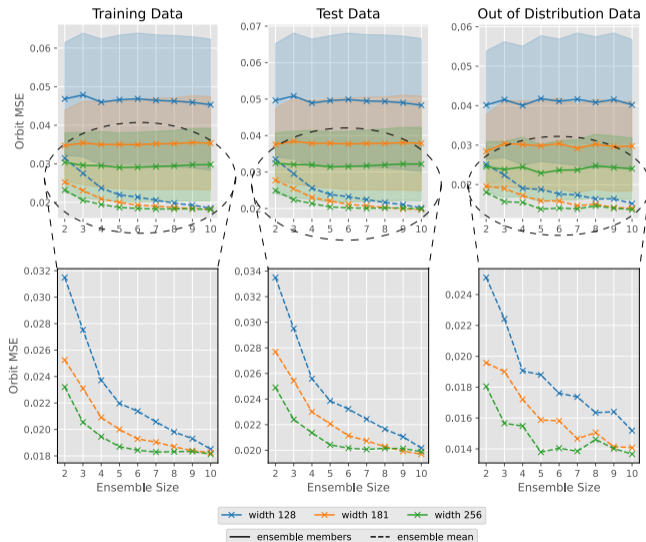
arXiv: 2406.06504



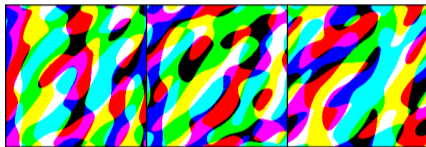
# Thank you

# Backup

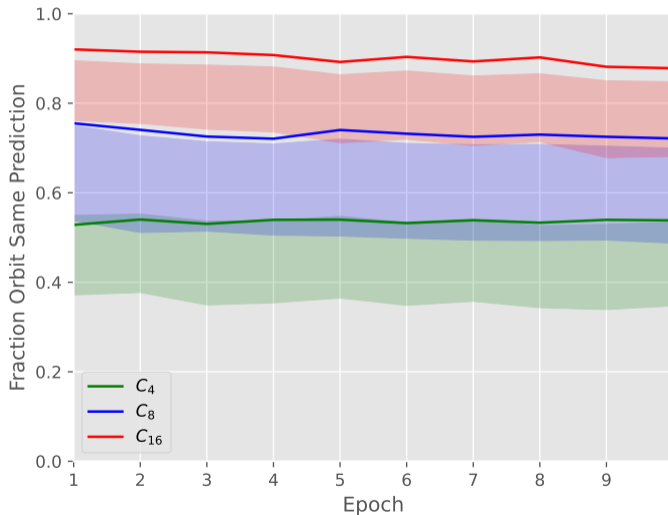
# Emergent equivariance of cross products



# Histological Data – OOD samples



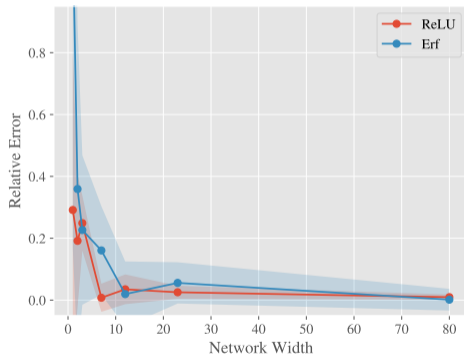
# Emergent continuous symmetry on FashionMNIST



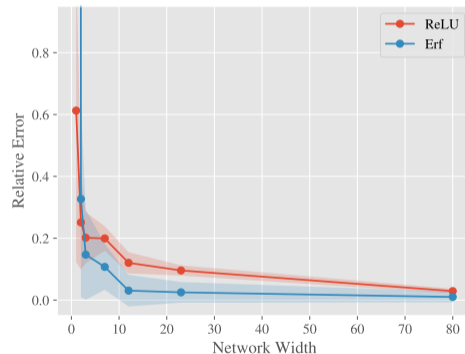


# Kernel convergence

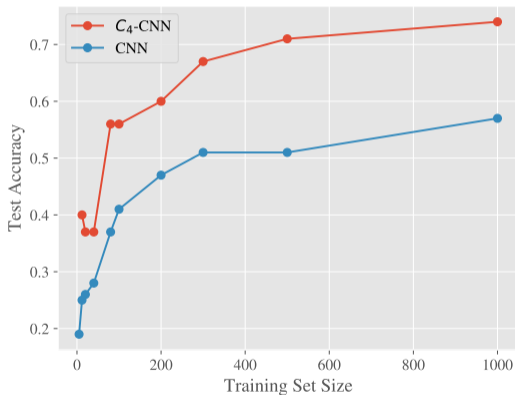
## NNGP



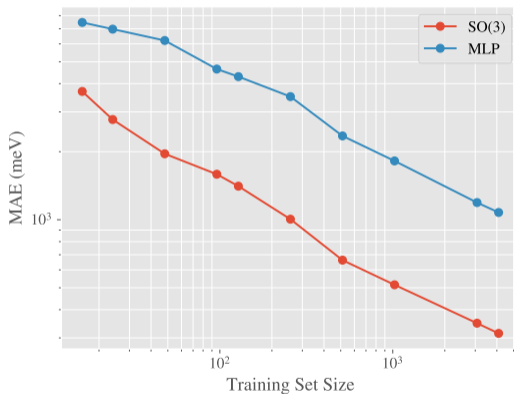
## NTK



# Equivariant NTKs for medical image classification

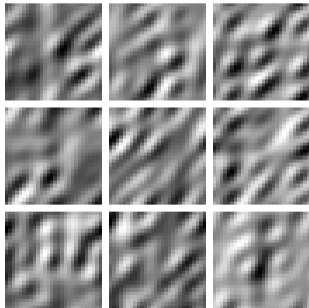


# Equivariant NTKs for molecular property regression



# OOD samples for CNN to GCNN convergence

MNIST



CIFAR10

