

Geometric Deep Learning and Neural Tangent Kernels

Jan E. Gerken



CHALMERS
UNIVERSITY OF TECHNOLOGY



UNIVERSITY OF
GOTHENBURG

WASP | WALLENBERG AI
AUTONOMOUS SYSTEMS
AND SOFTWARE PROGRAM

WASP AI/Math supervisor workshop 2024

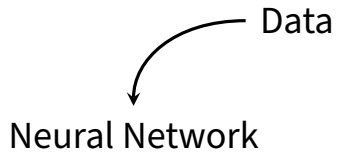
Part I: Geometric Deep Learning

Deep learning

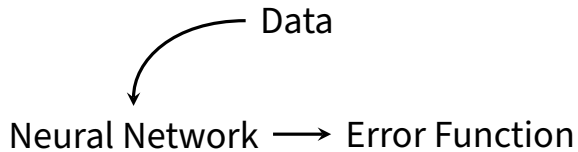
Deep learning

Data

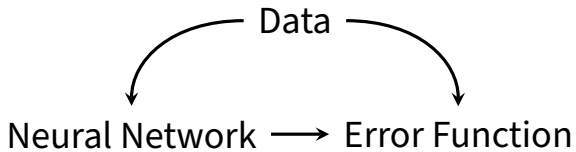
Deep learning



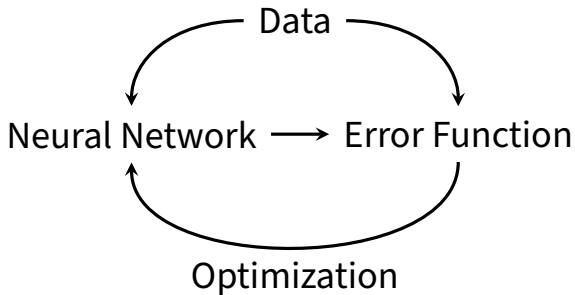
Deep learning



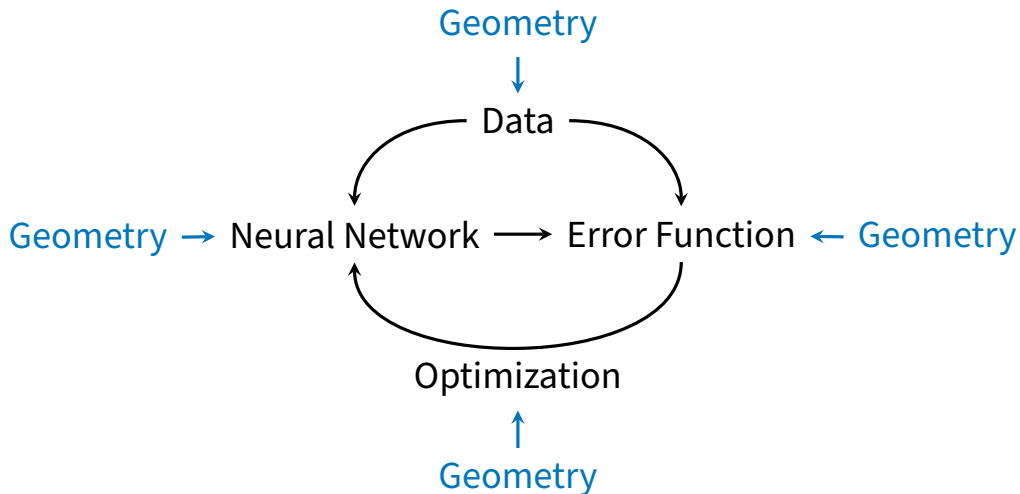
Deep learning



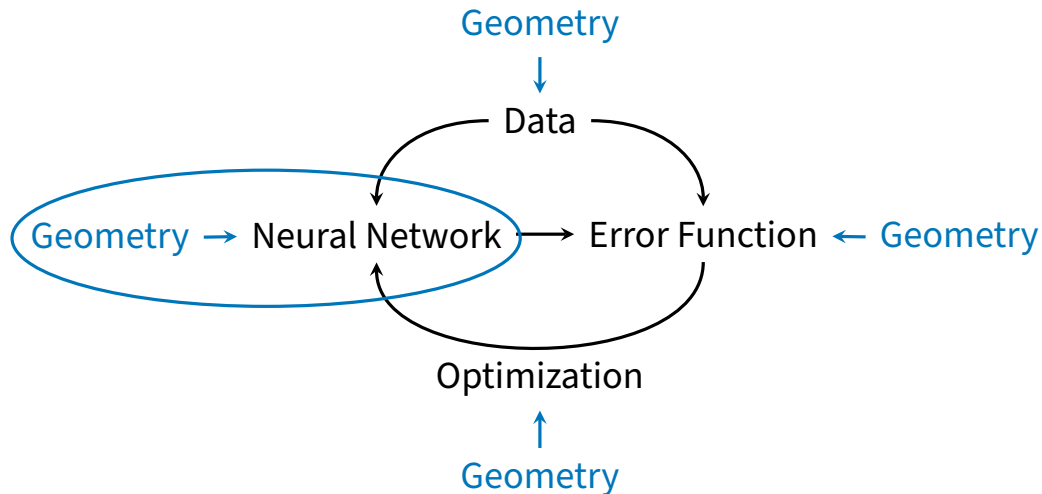
Deep learning



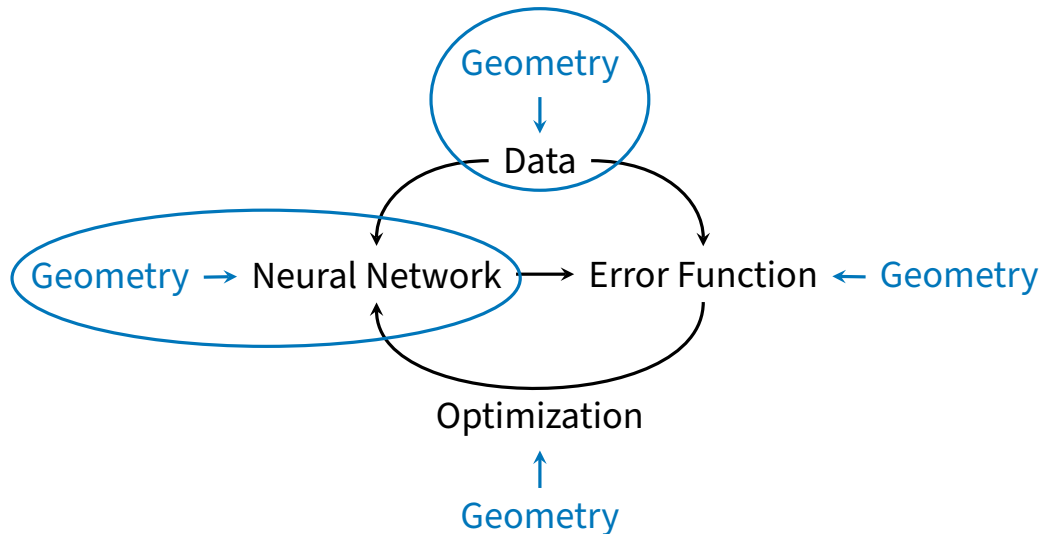
Geometric deep learning



Geometric deep learning



Geometric deep learning



HEAL-SWIN

in collaboration with



Oscar Carlsson



Hampus Linander



Heiner Spieß



Fredrik Ohlsson



Christoffer Petersson

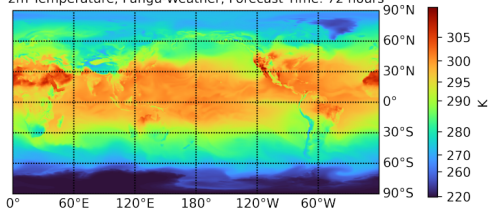


Daniel Persson

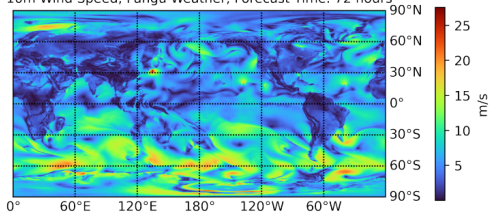
Spherical data

Spherical data

2m Temperature, Pangu-Weather, Forecast Time: 72 hours

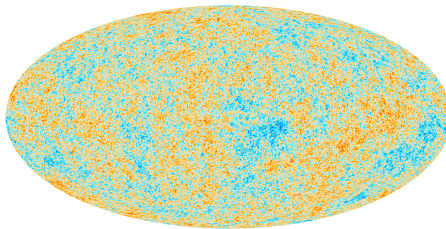
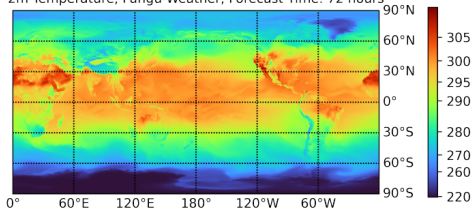


10m Wind Speed, Pangu-Weather, Forecast Time: 72 hours

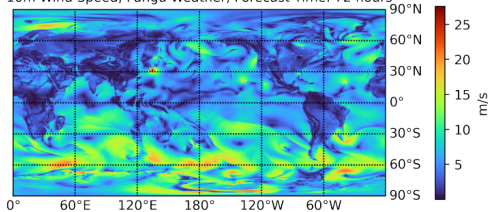


Spherical data

2m Temperature, Pangu-Weather, Forecast Time: 72 hours

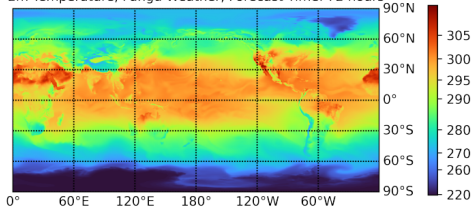


10m Wind Speed, Pangu-Weather, Forecast Time: 72 hours

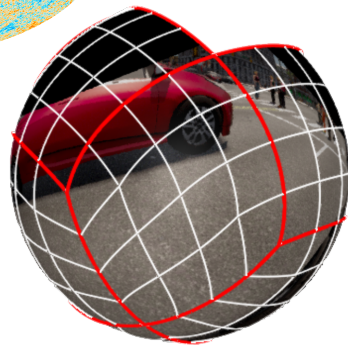
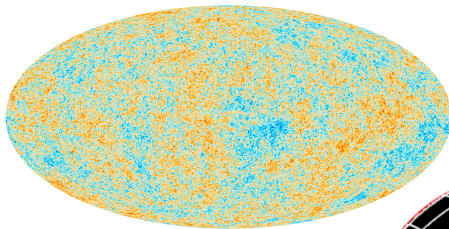
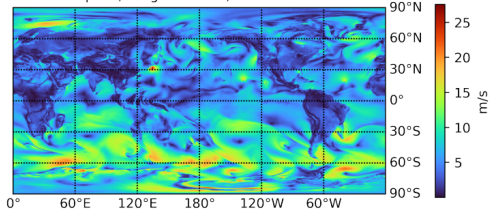


Spherical data

2m Temperature, Pangu-Weather, Forecast Time: 72 hours



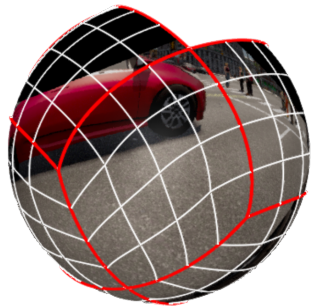
10m Wind Speed, Pangu-Weather, Forecast Time: 72 hours



Fisheye images as spherical data



project



SWIN transformer

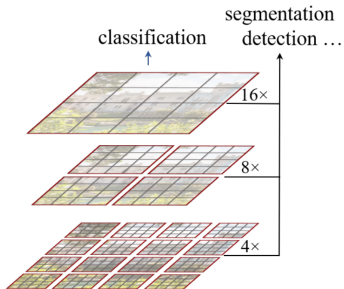
[Liu et al. 2021]

SWIN = Shifting Windows

SWIN transformer

[Liu et al. 2021]

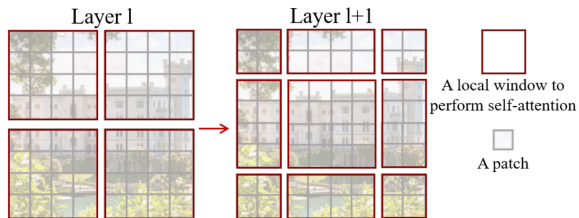
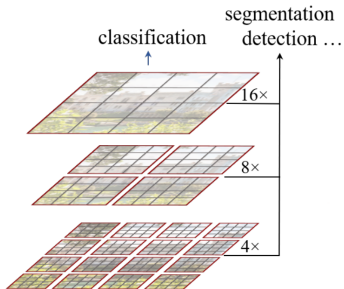
SWIN = Shifting Windows



SWIN transformer

[Liu et al. 2021]

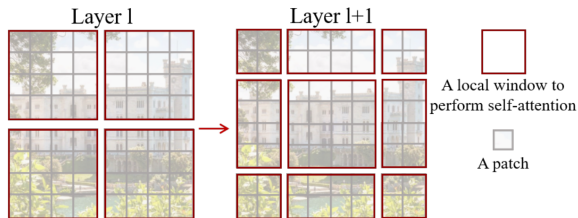
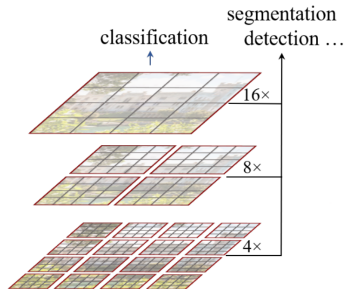
SWIN = Shifting Windows



SWIN transformer

[Liu et al. 2021]

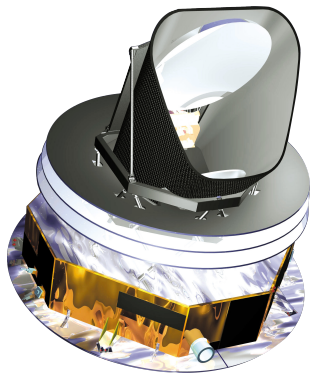
SWIN = Shifting Windows



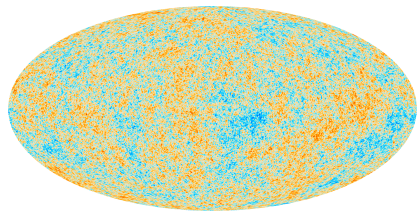
Goal: Construct SWIN transformer for spherical data

Sampling on the sphere

Sampling on the sphere



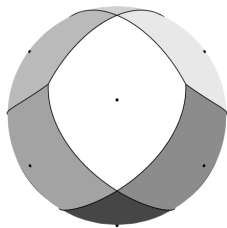
sample



HEALPix

[Gorski et al. 1998]

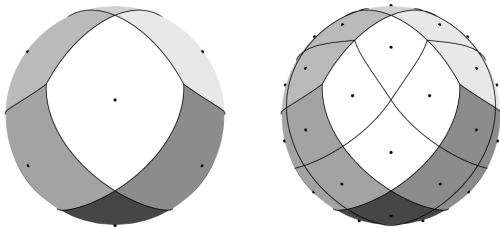
HEALPix = Hierarchical Equal Area iso-Latitude Pixelisation



HEALPix

[Gorski et al. 1998]

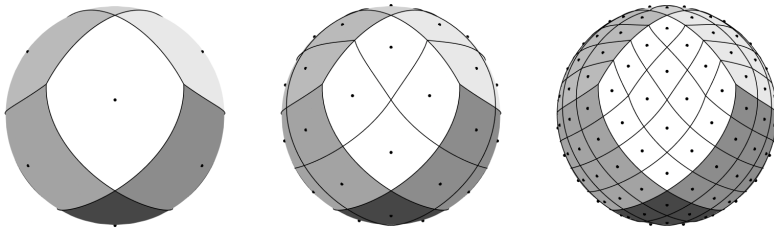
HEALPix = Hierarchical Equal Area iso-Latitude Pixelisation



HEALPix

[Gorski et al. 1998]

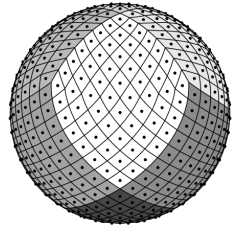
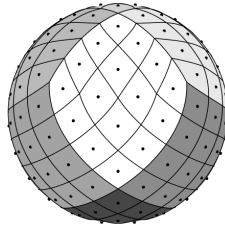
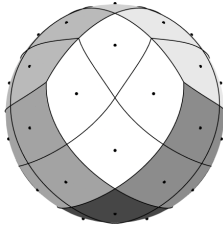
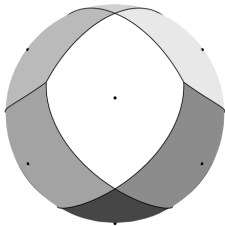
HEALPix = Hierarchical Equal Area iso-Latitude Pixelisation



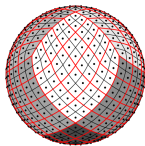
HEALPix

[Gorski et al. 1998]

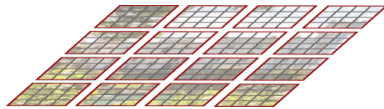
HEALPix = Hierarchical Equal Area iso-Latitude Pixelisation



HEAL-SWIN: Windowing

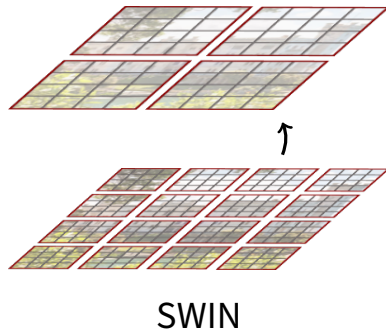
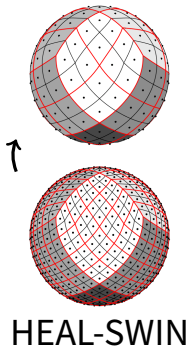


HEAL-SWIN

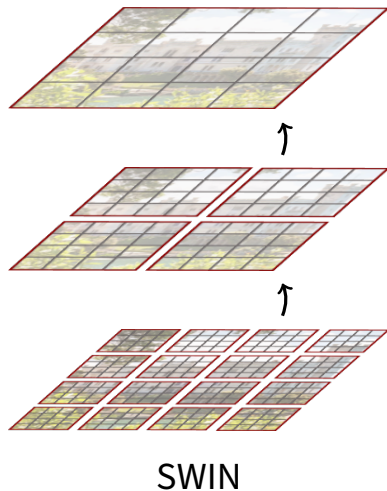
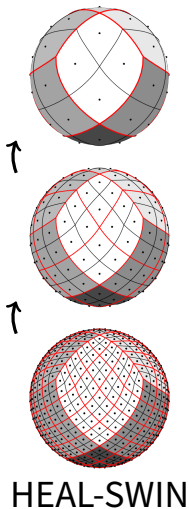


SWIN

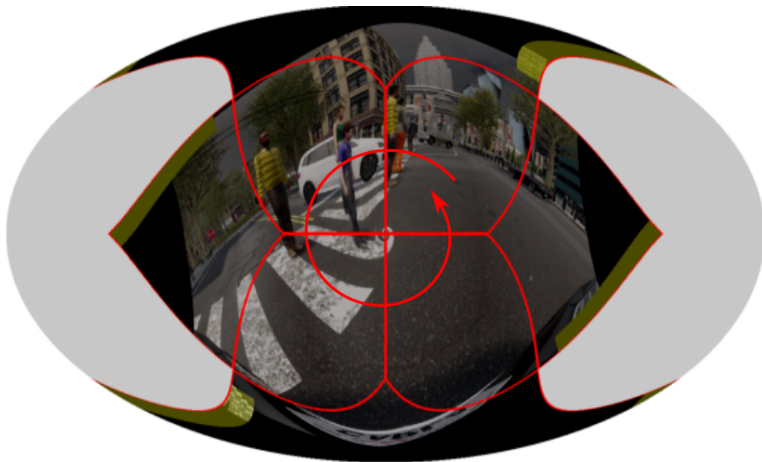
HEAL-SWIN: Windowing



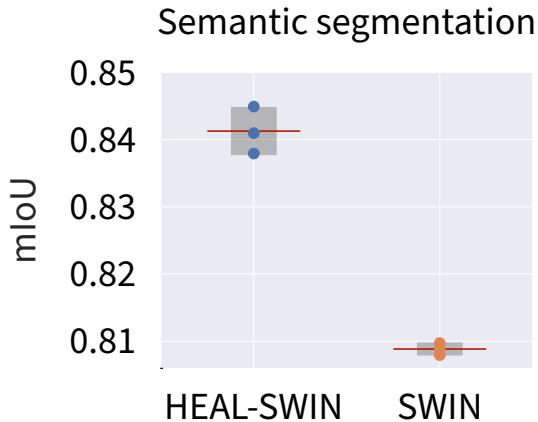
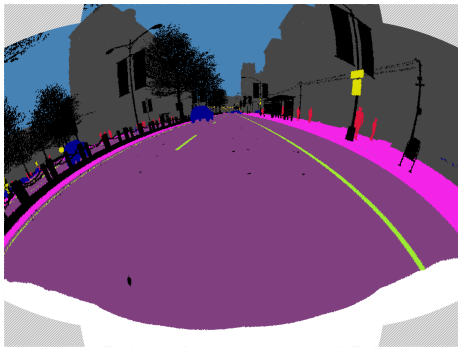
HEAL-SWIN: Windowing



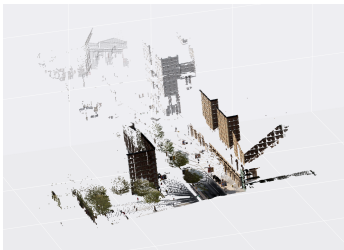
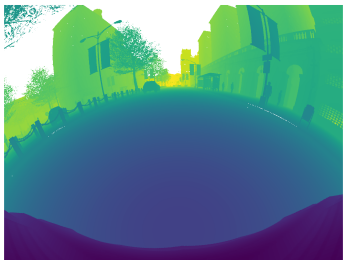
HEAL-SWIN: Shifting



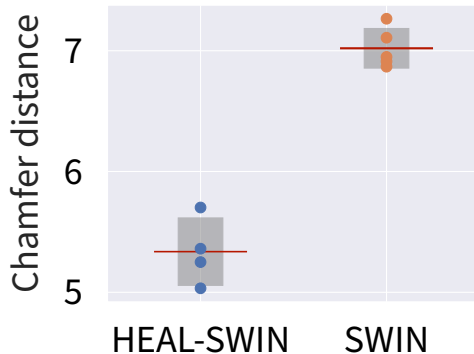
Semantic segmentation



Depth estimation



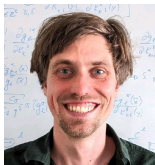
Depth estimation error



Part II: Neural Tangent Kernels

Emergent Equivariance in Deep Ensembles

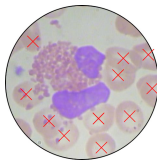
in collaboration with



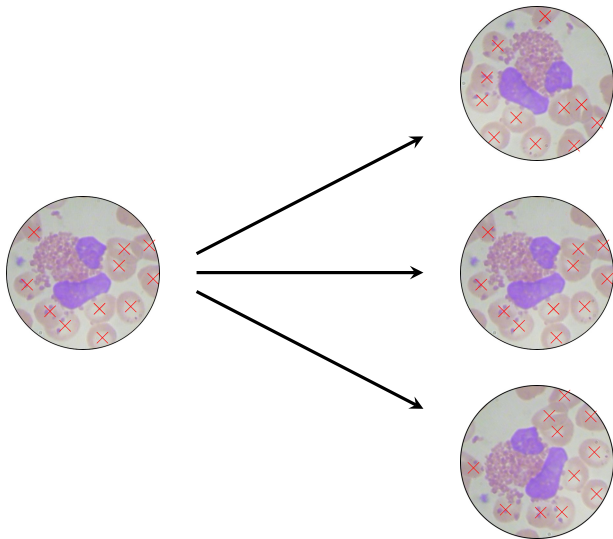
Pan Kessel

Data augmentation

Data augmentation



Data augmentation



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?

Neural Tangent Kernel

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 η

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 η (indicated by a blue arrow pointing to the fraction)

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

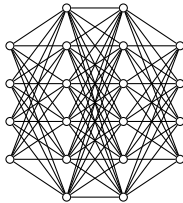
training sample x_i (indicated by a blue arrow pointing to the kernel argument)

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

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

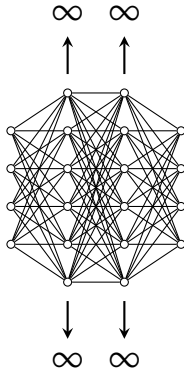
Infinite width limit

[Jacot et al. 2018]



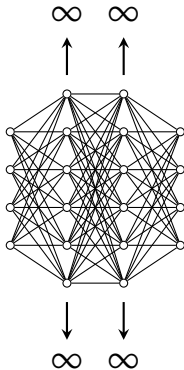
Infinite width limit

[Jacot et al. 2018]



Infinite width limit

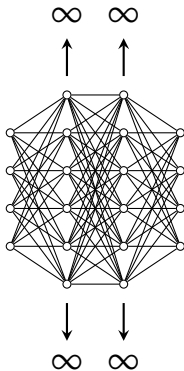
[Jacot et al. 2018]



👍 NTK becomes independent of initialization

Infinite width limit

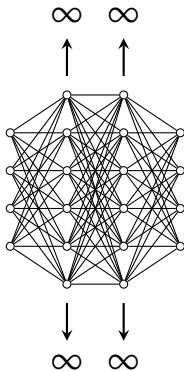
[Jacot et al. 2018]



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

Infinite width limit

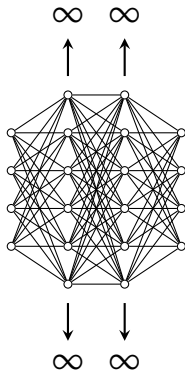
[Jacot et al. 2018]



- 👍 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]

Ⓢ 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

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

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

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

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. A single blue arrow also points from "augmented labels" to the term t in the exponent of the exponential function.

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 $\Theta(X, X)$ and $\Theta(X, X)^{-1}$ terms. A final blue arrow points from "augmented labels" to the $\rho(g)$ 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
 - at infinite width

Main conclusion

Deep ensembles trained with data augmentation are equivariant.

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

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.

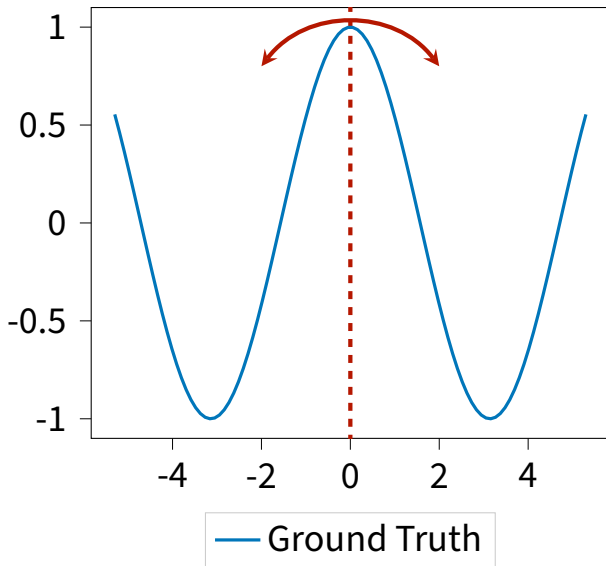
Intuitive explanation

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

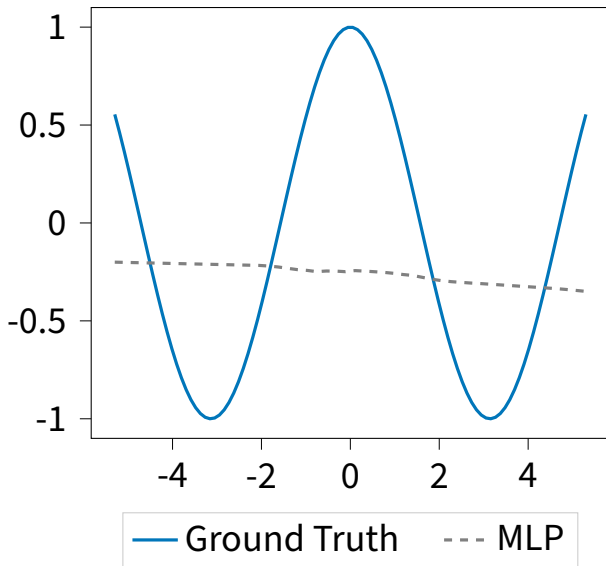
⊙ At infinite width, the mean output at initialization is zero everywhere.

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

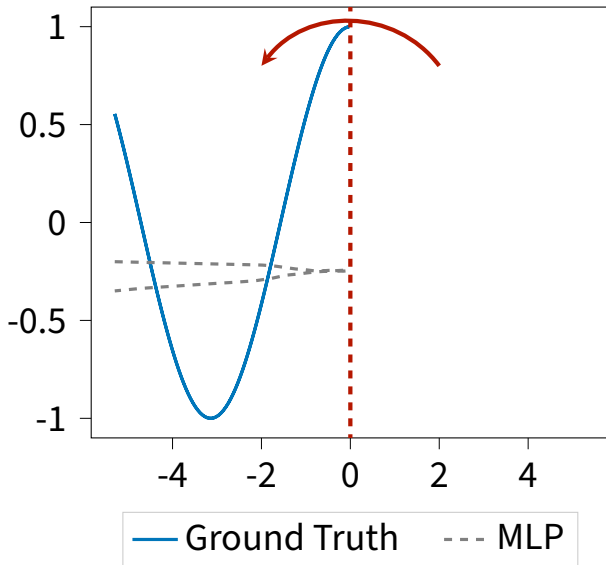
Toy example



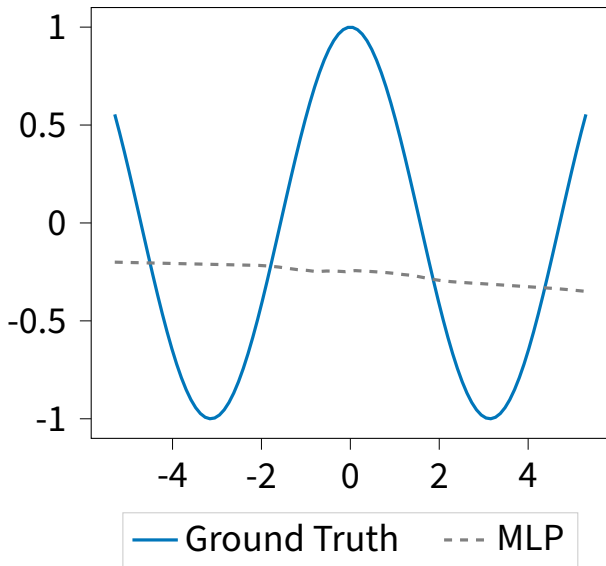
Initialization



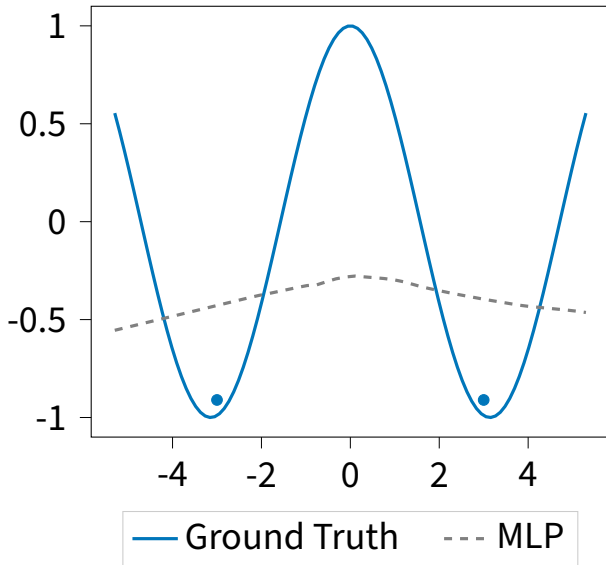
Initialization



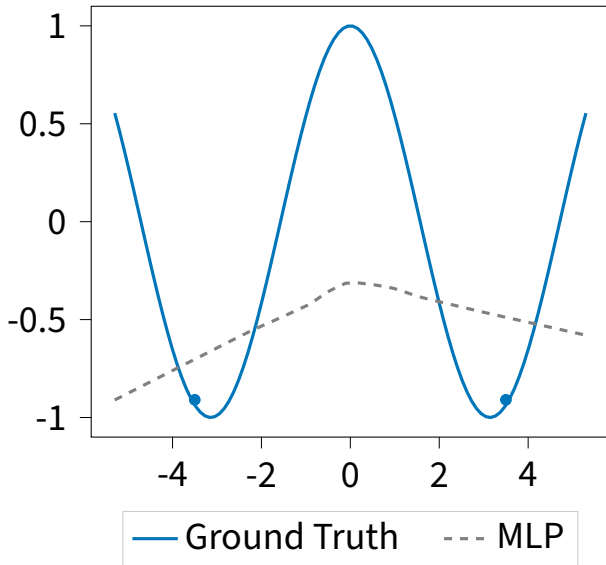
Initialization



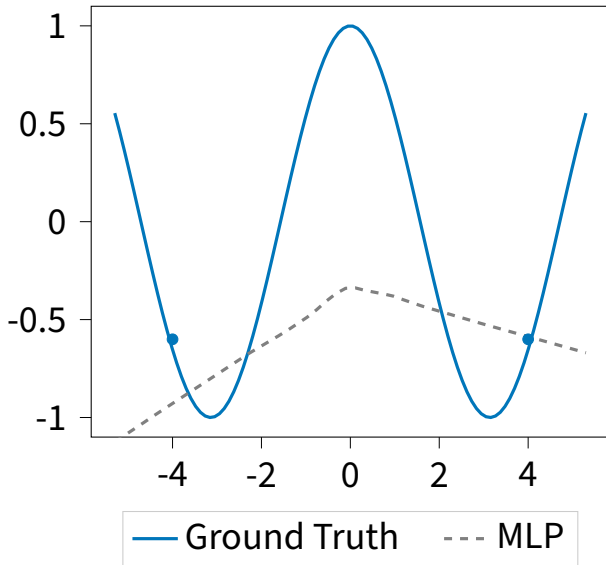
After 1 Training Step



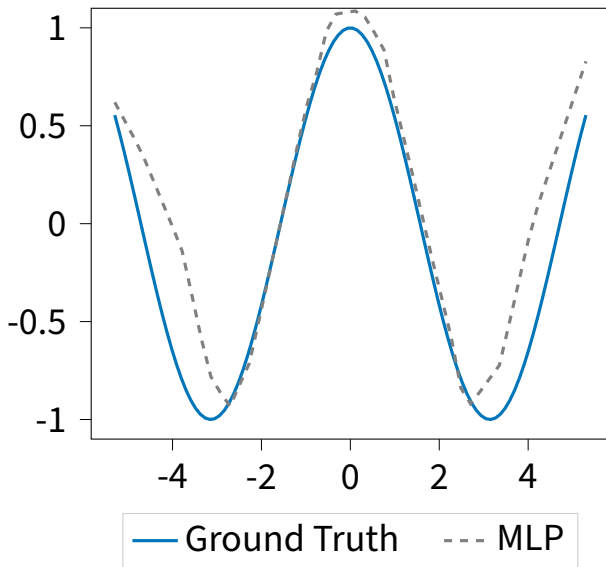
After 2 Training Steps



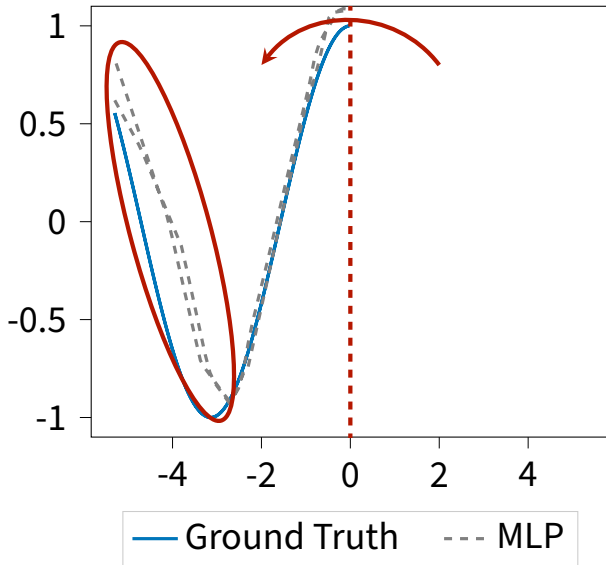
After 3 Training Steps



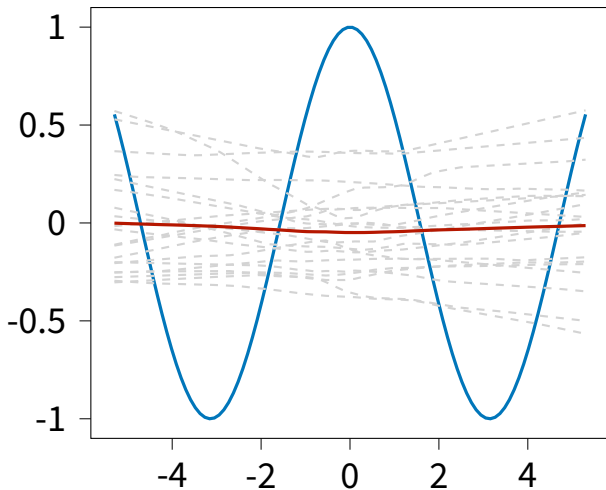
After 2000 Training Steps



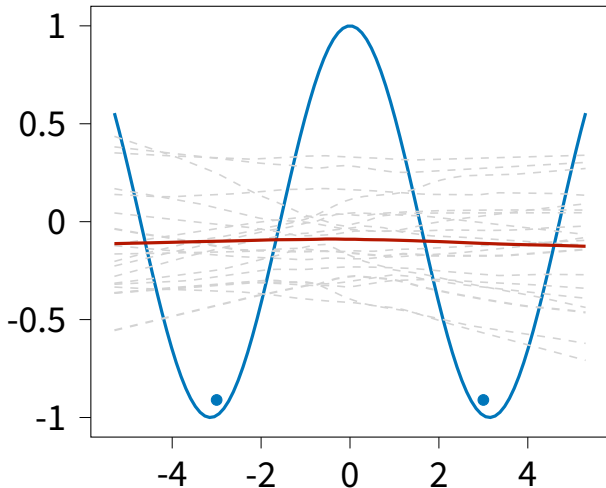
After 2000 Training Steps



Initialization

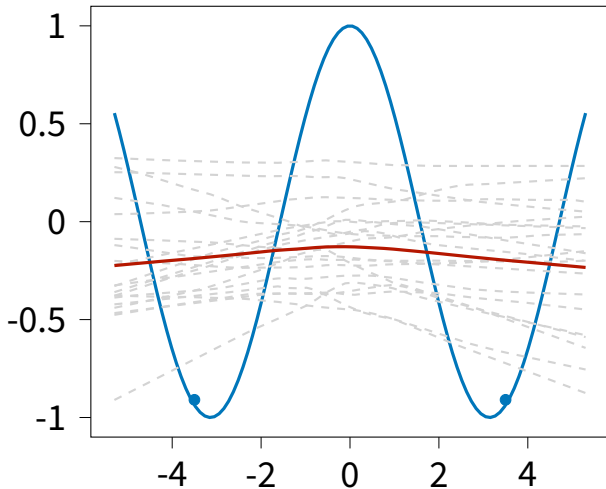


After 1 Training Step



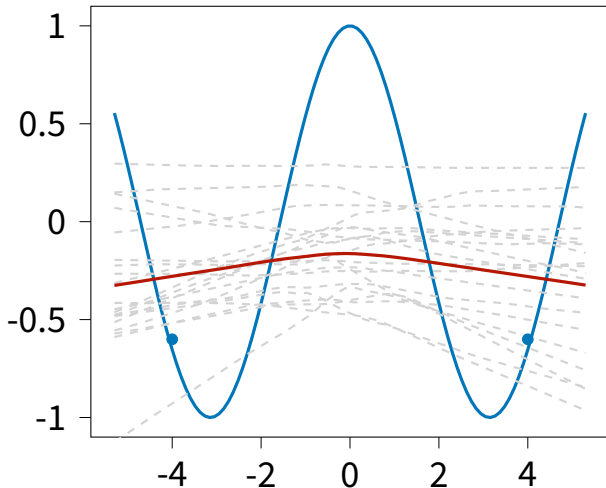
— Ground Truth - - - MLP — Ensemble Mean

After 2 Training Steps

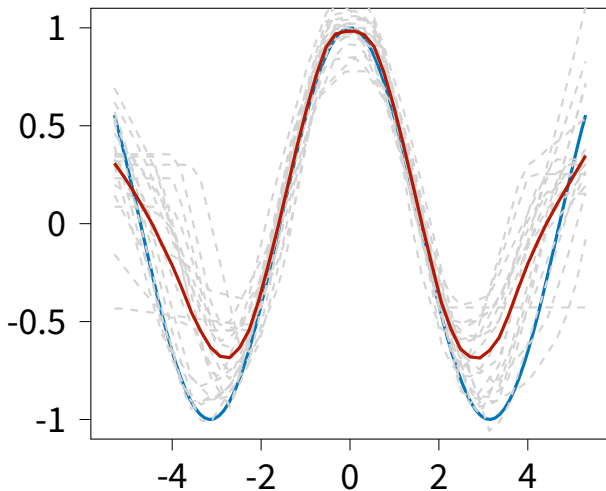


— Ground Truth - - - MLP — Ensemble Mean

After 3 Training Steps

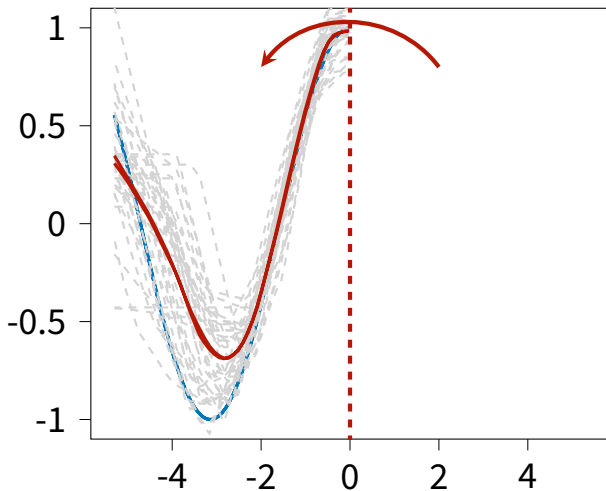


After 2000 Training Steps



— Ground Truth - - - MLP — Ensemble Mean

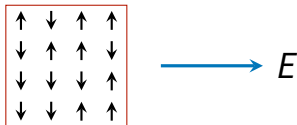
After 2000 Training Steps



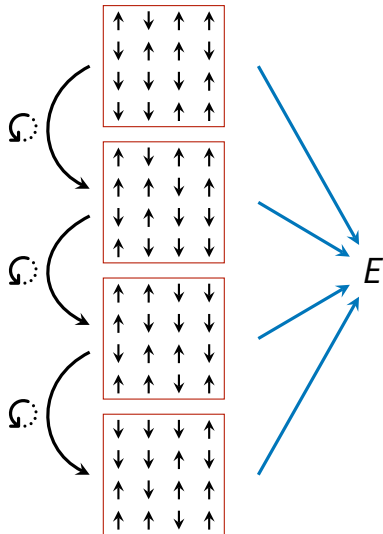
— Ground Truth - - - MLP — Ensemble Mean

Experiments

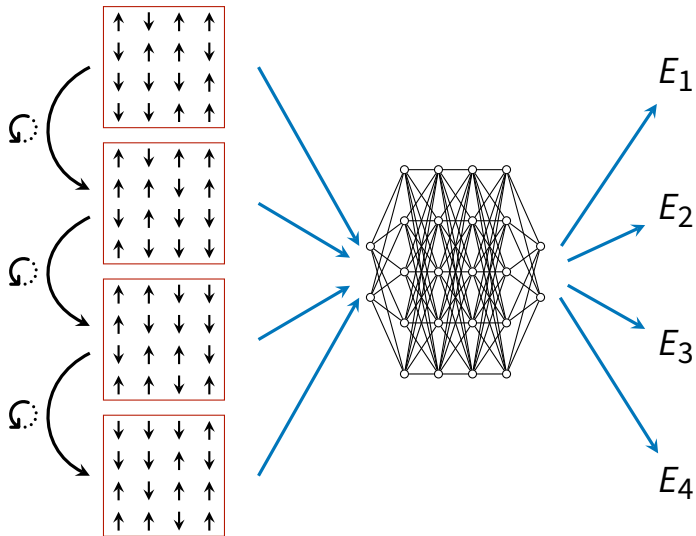
Ising model



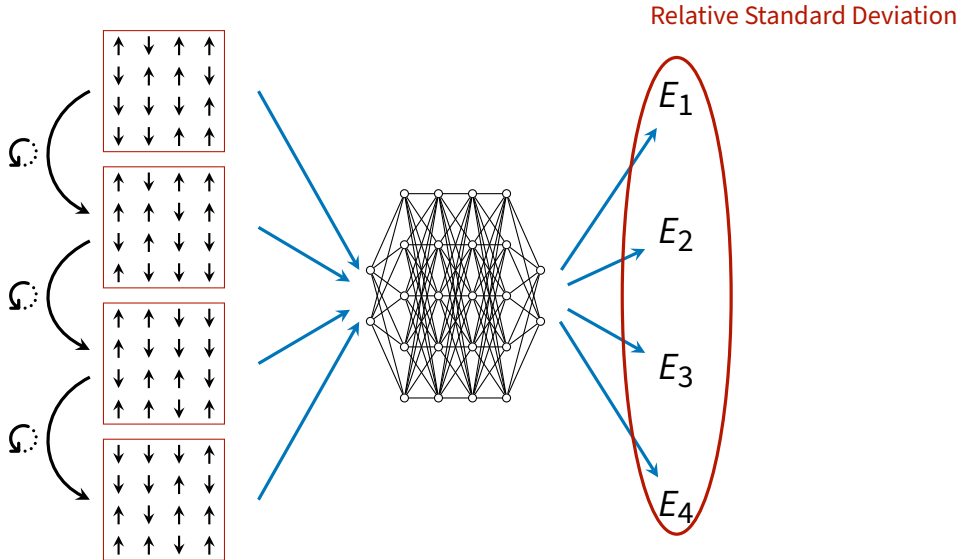
Ising model

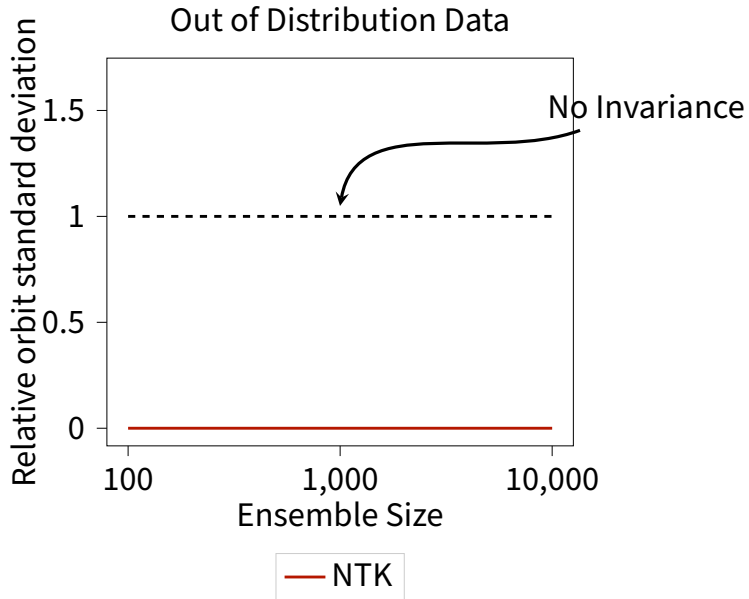


Ising model

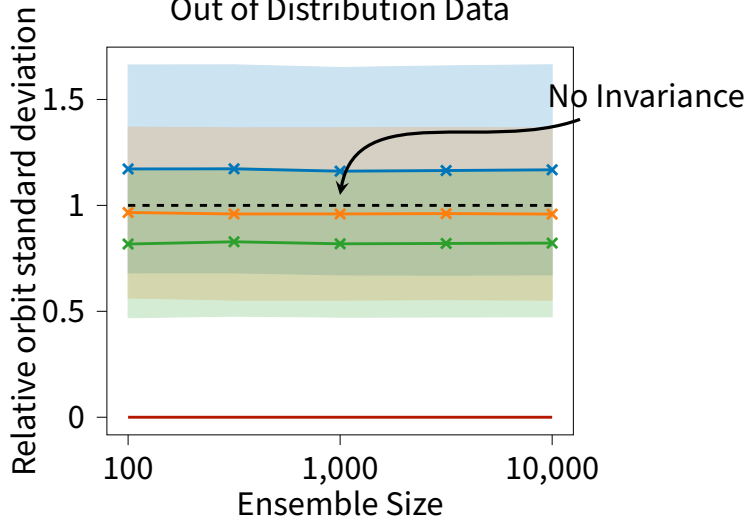


Ising model

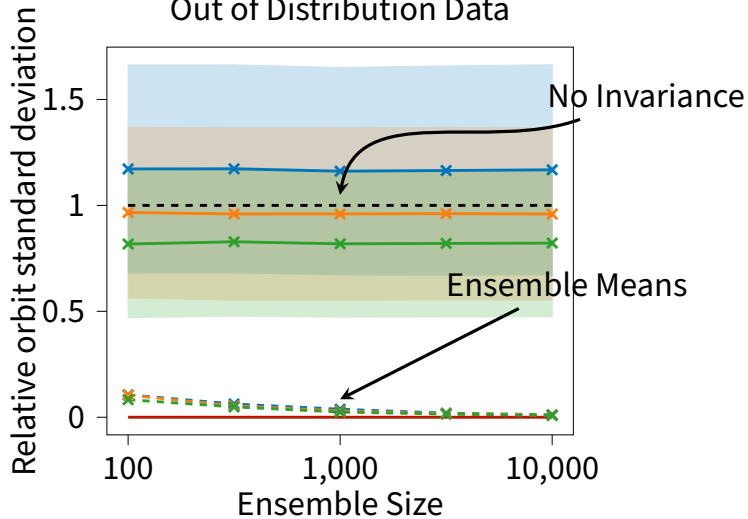




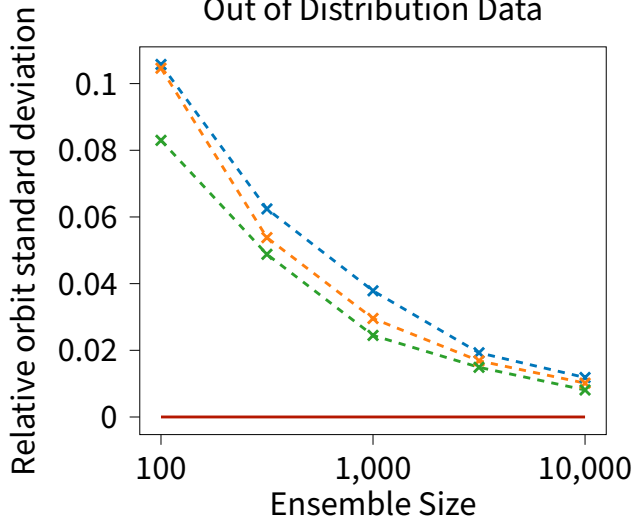
Out of Distribution Data



Out of Distribution Data



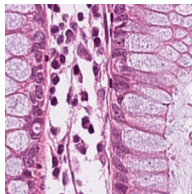
Out of Distribution Data



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

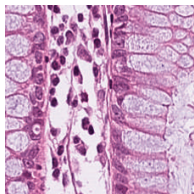
Histological slices

[Kather et al. 2018]



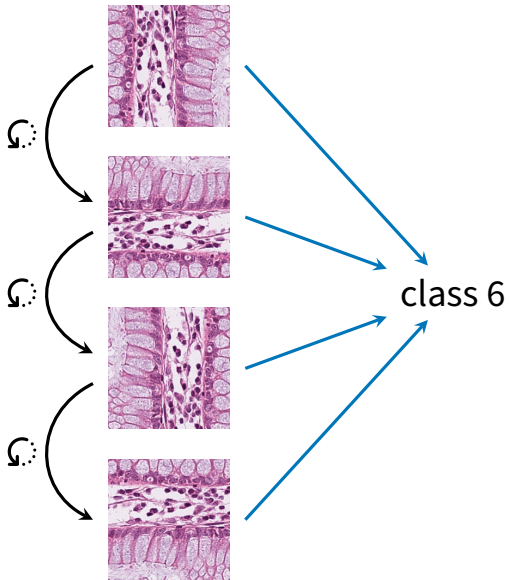
Histological slices

[Kather et al. 2018]

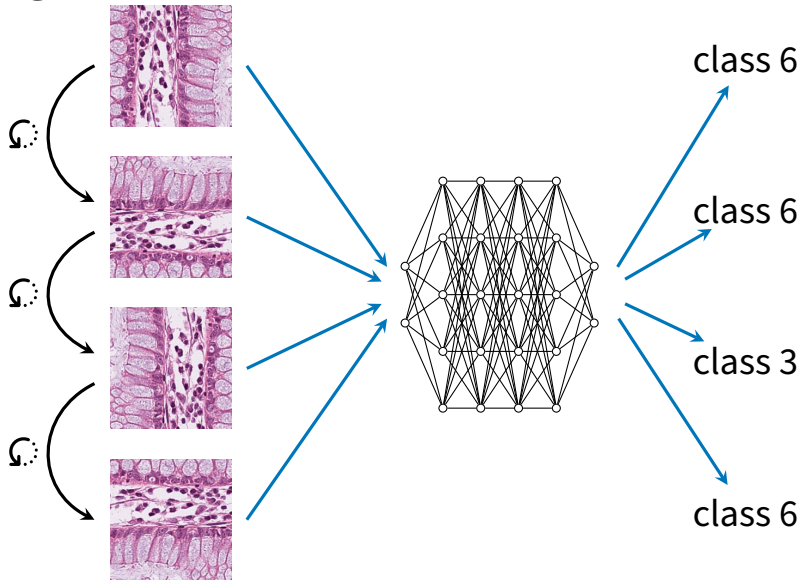


→ class 6

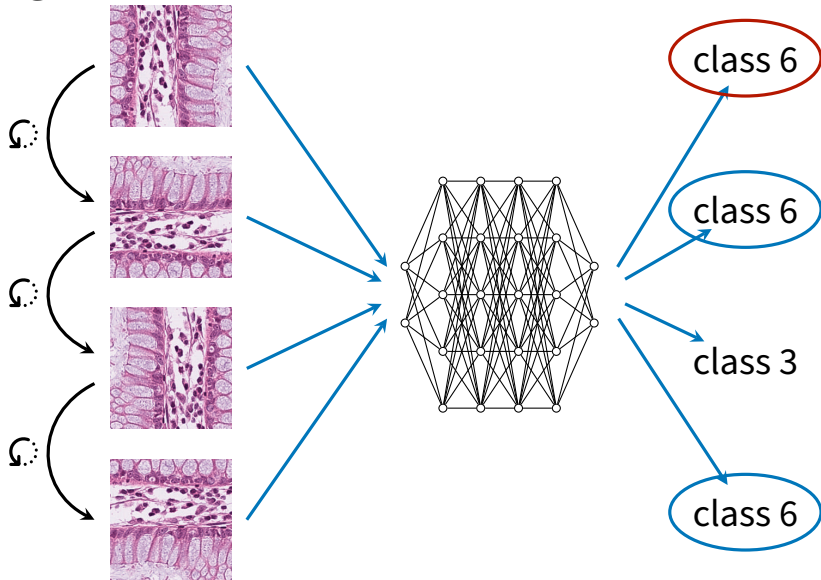
Histological slices



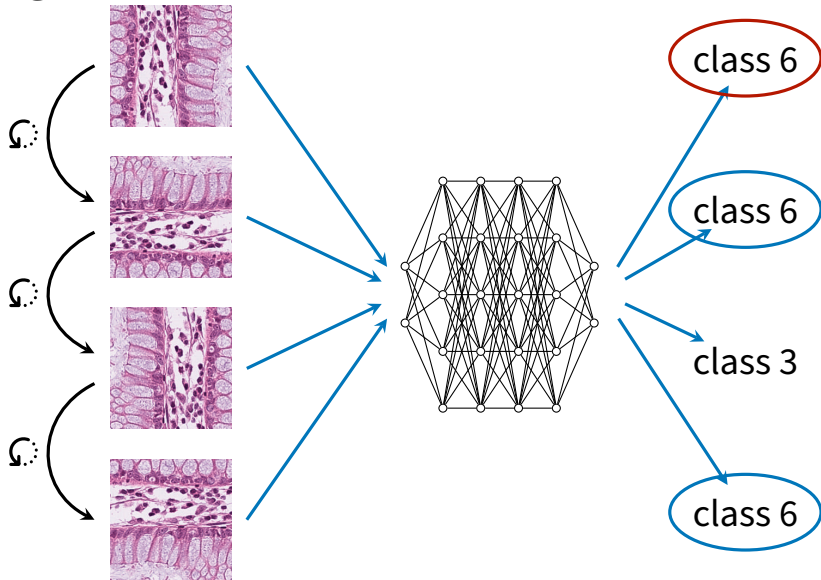
Histological slices



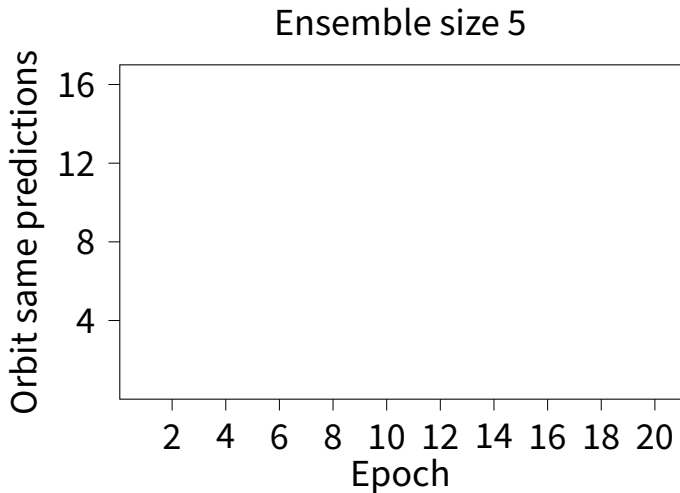
Histological slices



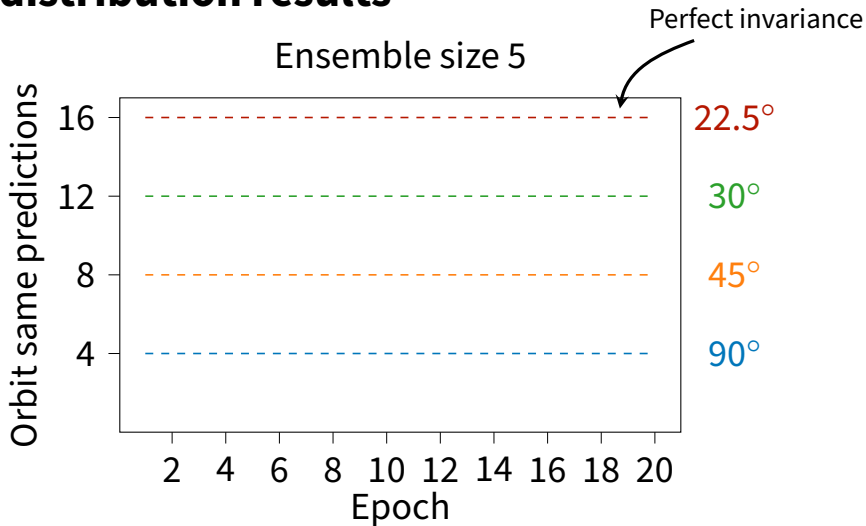
Histological slices



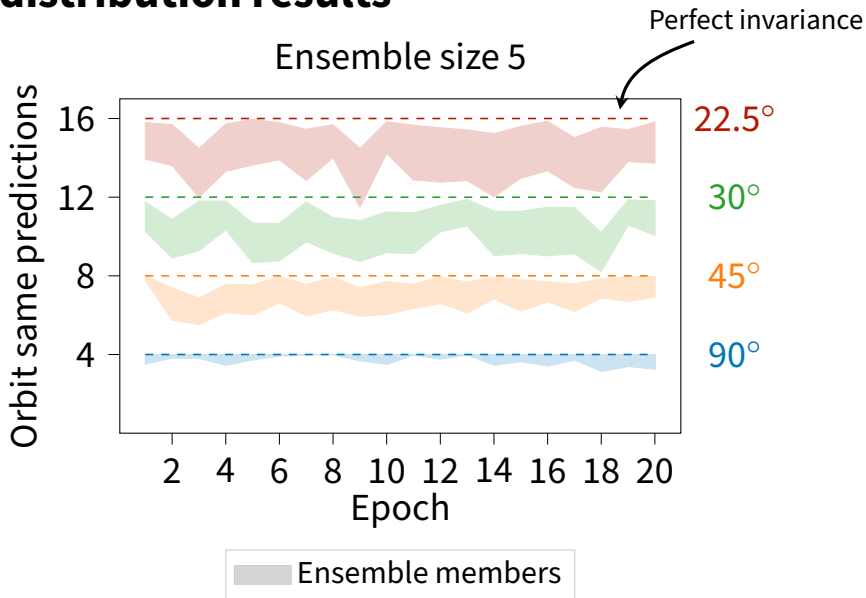
Out of distribution results



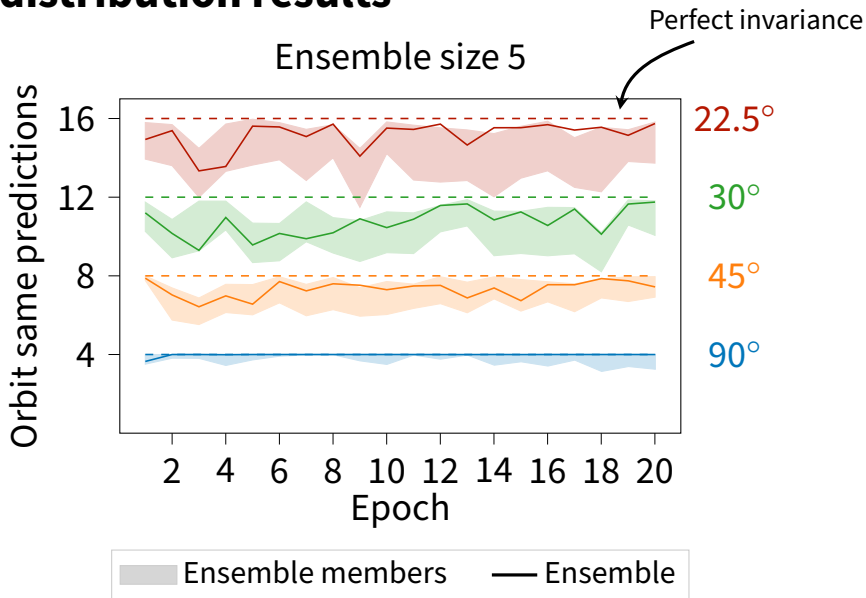
Out of distribution results



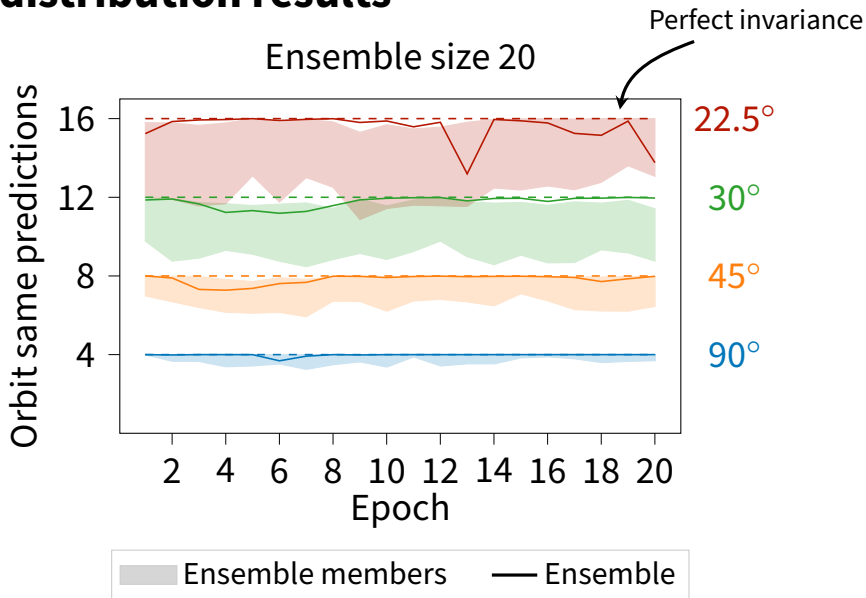
Out of distribution results



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 ± 0.12	7.72 ± 0.34	15.24 ± 0.69
only DA	3.41 ± 0.18	6.73 ± 0.24	12.77 ± 0.71
E2CNN ¹	4 ± 0.0	7.71 ± 0.21	15.08 ± 0.34
Canon ²	4 ± 0.0	7.45 ± 0.14	12.41 ± 0.85

¹[Weiler et al. 2019], ²[Kaba et al. 2022]

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

- *HEAL-SWIN: A Vision Transformer On The Sphere*

Oscar Carlsson^{*}, Jan E. Gerken^{*}, Hampus Linander, Heiner Spieß, Fredrik Ohlsson, Christoffer Petersson, Daniel Persson

CVPR 2024

- *Emergent Equivariance in Deep Ensembles*

Jan E. Gerken^{*}, Pan Kessel^{*}

ICML 2024 (Oral)

^{*} Equal contribution



Group Website