# Symmetries in AI4Science

Jan E. Gerken

CHALMERS
UNIVERSITY OF TECHNOLOGY

UNIVERSITY OF
GOTHENBURG

WASP | WALLENBERG AI,
AUTONOMOUS SYSTEMS
AND SOFTWARE PROGRAM

Workshop on Machine Learning Based Sampling in Lattice Field Theory and Quantum Chemistry

TRA Colloquium

Bonn
22$^{th}$ October 2024

# Symmetries in physics

$$SU(2) \times SU(3) \times U(1) \longleftrightarrow$$

## Standard Model of Elementary Particles

# Symmetries in chemistry



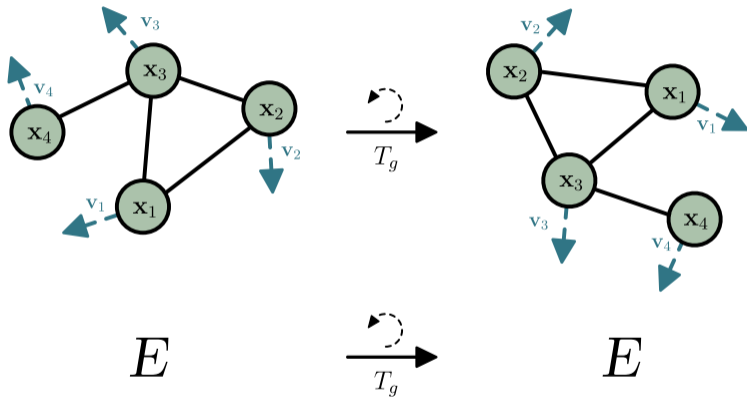Image from: [Satorras et al. 2021]

# Symmetries in prediction models

# Symmetries in prediction models

# Symmetries in prediction models



rotate

# Symmetries in prediction models

# Symmetries in prediction models

# Symmetries in prediction models



$$\iff \mathcal{N}(\rho_{\text{in}}(g)x) = \rho_{\text{out}}(g)\mathcal{N}(x)$$

# Symmetries in prediction models



$$\Longleftrightarrow \quad \mathcal{N}(\rho_{\text{in}}(g)x) = \rho_{\text{out}}(g)\mathcal{N}(x)$$

network

rotation

image

rotate

network

network

rotate

# Equivariance



$$\iff \mathcal{N}(\rho_{\text{in}}(g)x) = \rho_{\text{out}}(g)\mathcal{N}(x)$$

network

rotation   image

rotate

network   network

rotate

# Symmetries in generative models

- Sample from an invariant distribution

$$p(x) = p(\rho(g)x)$$

# Symmetries in generative models

- Sample from an invariant distribution

$$p(x) = p(\rho(g)x)$$

- For latent variable models:



invariant latent distribution

equivariant model

# Fundamental representation

- Groups act on vector spaces with representations

$$\rho : G \to \mathbb{R}^{n \times n}$$

# Fundamental representation

- Groups act on vector spaces with representations

$$\rho : G \to \mathbb{R}^{n \times n}$$

- Vectors transform in the defining representation of matrix Lie groups

$$f_x \to \rho(g) f_x$$

# Fundamental representation

- Groups act on vector spaces with representations

$$\rho : G \to \mathbb{R}^{n \times n}$$

- Vectors transform in the defining representation of matrix Lie groups

$$f_x \to \rho(g)f_x$$

- E.g. atom positions, force vectors

# Regular representation

- A group acts on itself by left multiplication $h \to gh$

# Regular representation

- A group acts on itself by left multiplication $h \to gh$

- This is leads to the <span style="color:red">regular representation</span> on $\mathbb{R}^{|G|}$

$$\rho_{\text{reg}}(g)e_h = e_{gh}$$

# Regular representation

- A group acts on itself by left multiplication $h \to gh$

- This is leads to the regular representation on $\mathbb{R}^{|G|}$

$$\rho_{\text{reg}}(g)e_h = e_{gh}$$

- Functions $G \to \mathbb{R}$ can be identified with $\mathbb{R}^{|G|}$

$$e_g \leftrightarrow \mathbb{I}_g$$

# Regular representation

- A group acts on itself by left multiplication $h \to gh$

- This is leads to the regular representation on $\mathbb{R}^{|G|}$

$$\rho_{\text{reg}}(g)e_h = e_{gh}$$

- Functions $G \to \mathbb{R}$ can be identified with $\mathbb{R}^{|G|}$

$$e_g \leftrightarrow \mathbb{I}_g$$

- The regular representation on functions $f : G \to \mathbb{R}$ is given by

$$(\rho_{\text{reg}}(g)f)(h) = f(g^{-1}h)$$

# Regular representation

- The regular representation on functions $f : G \to \mathbb{R}$ is given by

$$(\rho_{\mathsf{reg}}(g)f)(h) = f(g^{-1}h)$$

# Regular representation

- The regular representation on functions $f : G \to \mathbb{R}$ is given by

$$(\rho_{\text{reg}}(g)f)(h) = f(g^{-1}h)$$

- For functions on general domains, the regular representation is

$$(\rho_{\text{reg}}(g)f)(x) = f(\rho(g^{-1})x)$$

# Regular representation

- The regular representation on functions $f : G \to \mathbb{R}$ is given by

$$(\rho_{\text{reg}}(g)f)(h) = f(g^{-1}h)$$

- For functions on general domains, the regular representation is

$$(\rho_{\text{reg}}(g)f)(x) = f(\rho(g^{-1})x)$$

- This is how scalar fields transform

# Vector fields

- For functions $f : D \to \mathbb{R}^n$, combine defining- with regular representation

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

# Vector fields

- For functions $f : D \to \mathbb{R}^n$, combine defining- with regular representation

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

- This is how vector fields transform

# Vector fields

- For functions $f : D \to \mathbb{R}^n$, combine defining- with regular representation

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

- This is how vector fields transform

- $\pi$ can also be other representation, e.g. adjoint representation

$$\pi(g)f = \rho(g)\, f \,\rho^{-1}(g)$$

# Arbitrary representations

- Any finite-dimensional representation of a compact group can be written as a direct sum of irreducible representations

# Arbitrary representations

- Any finite-dimensional representation of a compact group can be written as a direct sum of irreducible representations

- In particular, tensor product representations can be decomposed into direct sums

$$\rho^{\ell} \otimes \rho^{m} = \bigoplus_{n} (\rho^{n})^{\oplus c_{n}^{\ell m}}$$

# Arbitrary representations

- Any finite-dimensional representation of a compact group can be written as a direct sum of irreducible representations

- In particular, tensor product representations can be decomposed into direct sums

$$\rho^{\ell} \otimes \rho^{m} = \bigoplus_{n} (\rho^{n})^{\oplus c_{n}^{\ell m}}$$

- Change of basis done via Clebsch–Gordan coefficients

# How to construct equivariant layers?

- Consider linear maps $M$ which are equivariant

# How to construct equivariant layers?

- Consider linear maps $M$ which are equivariant

- For vectors, $M$ needs to be an <span style="color:red">intertwiner</span>

$$M\rho_{\text{in}}(g) = \rho_{\text{out}}(g)M \qquad \forall g \in G \qquad (*)$$

# How to construct equivariant layers?

- Consider linear maps $M$ which are equivariant

- For vectors, $M$ needs to be an <span style="color:red">intertwiner</span>

$$M\rho_{\text{in}}(g) = \rho_{\text{out}}(g)M \qquad \forall g \in G \qquad (*)$$

- Schur's lemma: for complex, irreducible representations

$$M = \lambda\mathbb{1} \quad \text{if} \quad \rho_{\text{in}} = \rho_{\text{out}} \qquad \text{and} \qquad M = 0 \quad \text{otherwise}$$

# How to construct equivariant layers?

- Consider linear maps $M$ which are equivariant

- For vectors, $M$ needs to be an intertwiner

$$M\rho_{\text{in}}(g) = \rho_{\text{out}}(g)M \qquad \forall g \in G \qquad (*)$$

- Schur's lemma: for complex, irreducible representations

  $$M = \lambda \mathbb{I} \quad \text{if} \quad \rho_{\text{in}} = \rho_{\text{out}} \qquad \text{and} \qquad M = 0 \quad \text{otherwise}$$

- Hence, decompose $\rho_{in}$, $\rho_{out}$ into irreps to solve $(*)$

# How to construct equivariant layers?

- For the regular representation, linear equivariant layers are given by group convolutions [Cohen, Welling 2016]

$$[\psi * f](g) = \int_G \mathrm{d}h \, \psi(h^{-1}g) f(h)$$

# How to construct equivariant layers?

- For the regular representation, linear equivariant layers are given by group convolutions [Cohen, Welling 2016]

$$[\psi * f](g) = \int_G \mathrm{d}h \, \psi(h^{-1}g) f(h)$$

- For the translation group, these become the usual convolutions

# How to construct equivariant layers?

- For the regular representation, linear equivariant layers are given by red group convolutions [Cohen, Welling 2016]

$$[\psi * f](g) = \int_G dh\, \psi(h^{-1}g)f(h)$$

- For the translation group, these become the usual convolutions

- For combination with fundamental representation $(\pi(g)f(\rho^{-1}(g)x))$, convolution filter needs to be an intertwiner

[Review: Weiler et al. 2023]

12

# Equivariance in quantum chemistry

- Group:
  roto-translations of the molecule + permutations of identical atoms

# Equivariance in quantum chemistry

- Group:
  roto-translations of the molecule + permutations of identical atoms

- Use graph-NNs for permutation part

# Equivariance in quantum chemistry

- Group:
  roto-translations of the molecule + permutations of identical atoms

- Use graph-NNs for permutation part

- For SO(3), expand in irreps, use tensor products to combine features

[Review: Duval et al. 2023]

# Gauge symmetry

- In a gauge symmetry, group element can depend on position $x$

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

# Gauge symmetry

- In a gauge symmetry, group element can depend on position $x$

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

- Symmetry becomes local

# Gauge symmetry

- In a gauge symmetry, group element can depend on position $x$

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

- Symmetry becomes local

- In quantum chemistry only global symmetries

# Gauge symmetry

- In a <span style="color:red">gauge</span> symmetry, group element can depend on position $x$

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

- Symmetry becomes local

- In quantum chemistry only global symmetries

- Equivariance wrt local coordinate changes is also a gauge symmetry: Gauge CNNs [Cheng et al. 2019]

## Equivariance in lattice field theory

- In lattice field theory, typically combination of local and global symmetries: $G = SU(n) \times SE(3)$

$$\pi(g(x))f(\rho^{-1}(h)x) \qquad g(x) \in \mathsf{SU}(n) \qquad h \in \mathsf{SE}(3)$$

# Equivariance in lattice field theory

- In lattice field theory, typically combination of local and global symmetries: $G = SU(n) \times SE(3)$

  $$\pi(g(x))f(\rho^{-1}(h)x) \qquad g(x) \in \mathrm{SU}(n) \qquad h \in \mathrm{SE}(3)$$

- The gauge group acts in the adjoint representation

  $$\pi(g(x))f = \rho(g(x))\, f\, \rho^{\dagger}(g(x))$$

# Equivariance in lattice field theory

- In lattice field theory, typically combination of local and global symmetries: $G = SU(n) \times SE(3)$

$$\pi(g(x))f(\rho^{-1}(h)x) \qquad g(x) \in SU(n) \qquad h \in SE(3)$$

- The gauge group acts in the adjoint representation

$$\pi(g(x))f = \rho(g(x))\,f\,\rho^\dagger(g(x))$$

- By discretizing on the lattice, obtain links $U_\mu$ transforming as

$$U_\mu(x) \to \rho(g(x))U_\mu(x)\,\rho^\dagger(g(x+\hat{\mu}))$$

## How to construct gauge equivariant layers?

- Can build loops transforming as

$$W(x) \to \rho(g(x))W(x)\,\rho^{\dagger}(g(x)) \qquad (*)$$

## How to construct gauge equivariant layers?

- Can build loops transforming as

$$W(x) \rightarrow \rho(g(x))W(x)\,\rho^\dagger(g(x)) \qquad (*)$$

- Products of loops are equivariant, traces are invariant

$$W(x)\tilde{W}(x) \rightarrow \rho(g(x))W(x)\tilde{W}(x)\,\rho^\dagger(g(x))$$
$$\mathrm{tr}(W(x)) \rightarrow \mathrm{tr}(W(x))$$

# How to construct gauge equivariant layers?

- Can build loops transforming as

$$W(x) \to \rho(g(x))W(x)\,\rho^{\dagger}(g(x)) \qquad (*)$$

- Products of loops are equivariant, traces are invariant

$$W(x)\tilde{W}(x) \to \rho(g(x))W(x)\tilde{W}(x)\,\rho^{\dagger}(g(x))$$
$$\mathrm{tr}(W(x)) \to \mathrm{tr}(W(x))$$

- Use this together with convolutions to build gauge equivariant networks [Favoni et al. 2020]

# How to construct gauge equivariant layers?

- Can build loops transforming as

$$W(x) \to \rho(g(x))W(x)\rho^\dagger(g(x)) \qquad (*)$$

- Products of loops are equivariant, traces are invariant

$$W(x)\tilde{W}(x) \to \rho(g(x))W(x)\tilde{W}(x)\rho^\dagger(g(x))$$
$$\text{tr}(W(x)) \to \text{tr}(W(x))$$

- Use this together with convolutions to build gauge equivariant networks [Favoni et al. 2020]

- Can also manipulate invariants of $(*)$ [Boyda et al. 2021]

16

# How to construct gauge equivariant layers?

- Can build loops transforming as

$$W(x) \to \rho(g(x))W(x)\,\rho^\dagger(g(x)) \qquad (*)$$

- Products of loops are equivariant, traces are invariant

$$W(x)\tilde{W}(x) \to \rho(g(x))W(x)\tilde{W}(x)\,\rho^\dagger(g(x))$$
$$\mathrm{tr}(W(x)) \to \mathrm{tr}(W(x))$$

- Use this together with convolutions to build gauge equivariant networks [Favoni et al. 2020]

- Can also manipulate invariants of $(*)$ [Boyda et al. 2021]

- Can differentiate an invariant [Bacchio et al. 2023]

# Part II: Other ways of reaching equivariance

# (Frame) averaging

- Create exactly equivariant model by averaging over the group

$$\bar{f}(x) = \int_G \mathrm{d}h \, \pi(g) f(\rho^{-1}(g)x)$$

# (Frame) averaging

- Create exactly equivariant model by averaging over the group

$$\bar{f}(x) = \int_G \mathrm{d}h \, \pi(g) f(\rho^{-1}(g)x)$$

👍 It is sufficient to average over an equivariant subset $\mathcal{F}(x) \subset G$ (frame) [Puny et al. 2022]

# (Frame) averaging

- Create exactly equivariant model by averaging over the group

$$\bar{f}(x) = \int_G dh\, \pi(g) f(\rho^{-1}(g)x)$$

- 👍 It is sufficient to average over an equivariant subset $\mathcal{F}(x) \subset G$ (frame) [Puny et al. 2022]

- 👍 Works with any architecture

# (Frame) averaging

- Create exactly equivariant model by averaging over the group

$$\bar{f}(x) = \int_G \mathrm{d}h\, \pi(g) f(\rho^{-1}(g)x)$$

- 👍 It is sufficient to average over an equivariant subset $\mathcal{F}(x) \subset G$ (frame) [Puny et al. 2022]

- 👍 Works with any architecture

- 👎 Only approximate for continuous groups when sampling is necessary to evaluate the integral

# Canonicalization

- Use an equivariant map $D \to G$ to predict a canonicalizing transformation

# Canonicalization

[Kaba et al. 2022]

- Use an equivariant map $D \rightarrow G$ to predict a canonicalizing transformation

- Use non-equivariant network for prediction

# Canonicalization

- Use an equivariant map $D \to G$ to predict a canonicalizing transformation

- Use non-equivariant network for prediction

- 👍 Exactly equivariant

# Canonicalization

- Use an equivariant map $D \to G$ to predict a canonicalizing transformation

- Use non-equivariant network for prediction

👍 Exactly equivariant

👎 Still needs equivariant model

# Canonicalization

[Kaba et al. 2022]

- Use an equivariant map $D \rightarrow G$ to predict a canonicalizing transformation

- Use non-equivariant network for prediction

👍 Exactly equivariant

👎 Still needs equivariant model

👎 Equivariant function with codomain $G$ is hard to construct

# 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?

# Emergent Equivariance in Deep Ensembles

in collaboration with



Pan Kessel

# Empirical NTK

Training dynamics under continuous gradient descent:

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

learning rate

loss

training sample

# Empirical NTK

Training dynamics under continuous gradient descent:

learning rate

loss

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

training sample

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

# Infinite width limit

[Jacot et al. 2018]

# Infinite width limit

∞ ∞

👍 NTK becomes independent of initialization

∞ ∞

# Infinite width limit

👍 NTK becomes independent of initialization

👍 NTK becomes constant in training

# Infinite width limit

👍 NTK becomes independent of initialization

👍 NTK becomes constant in training

👍 NTK can be computed for most networks

# Infinite width limit

∞  ∞

↑  ↑



↓  ↓

∞  ∞

👍 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

⊙ At infinite width, the mean prediction is given by

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

# Mean prediction from NTK

[Jacot et al. 2018]

⊙ At infinite width, the mean prediction is given by

neural tangent kernel

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

# Mean prediction from NTK

⊙ At infinite width, the mean prediction is given by

neural tangent kernel

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

train data

# Mean prediction from NTK

[Jacot et al. 2018]

⊙ At infinite width, the mean prediction is given by

neural tangent kernel

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

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

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

# Data augmentation at infinite width

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

augmented data

augmented labels

# Data augmentation at infinite width

group transformation

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

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

# 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

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

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

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

# Mean prediction

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

## 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

Ground Truth

Initialization

Ground Truth ---- MLP

# Initialization



Ground Truth — — — MLP

Initialization

Ground Truth  - - - MLP

# After 1 Training Step



Ground Truth ----- MLP

# 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



Legend: Ground Truth, MLP, Ensemble Mean

After 2 Training Steps

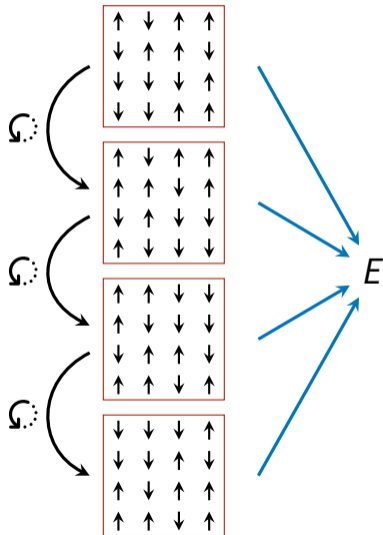Ground Truth — MLP — Ensemble Mean

# After 3 Training Steps



Legend: Ground Truth — MLP — Ensemble Mean

After 2000 Training Steps

Ground Truth — MLP — Ensemble Mean

# After 2000 Training Steps

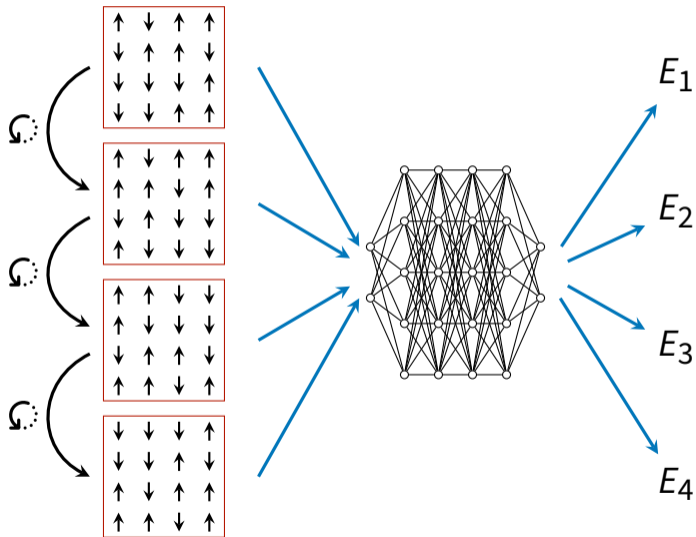

Legend: Ground Truth, MLP, Ensemble Mean

# Experiments

# Ising model



$$\longrightarrow E$$

# Ising model

# Ising model



$E_1$

$E_2$

$E_3$

$E_4$

# Ising model



Relative Standard Deviation

$E_1$

$E_2$

$E_3$

$E_4$

Out of Distribution Data

Relative orbit standard deviation vs. Ensemble Size

No Invariance

NTK

Out of Distribution Data

Relative orbit standard deviation vs Ensemble Size

No Invariance

Legend: NTK, Width 512, Width 1024, Width 2048

Out of Distribution Data

Relative orbit standard deviation vs Ensemble Size

No Invariance

Ensemble Means

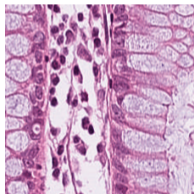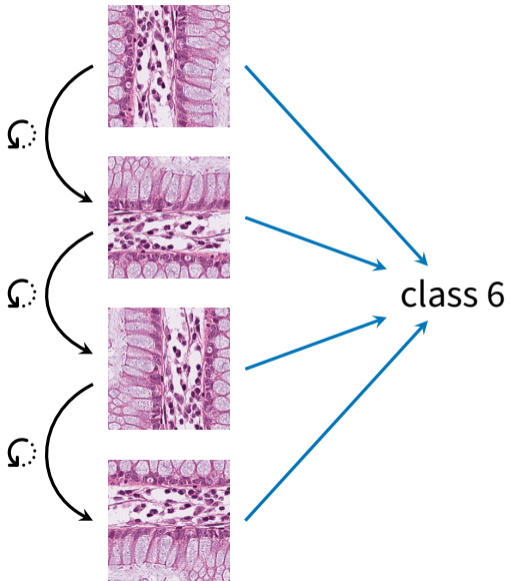Legend: NTK | Width 512 | Width 1024 | Width 2048

Out of Distribution Data

# Histological slices

[Kather et al. 2018]

# Histological slices

 → class 6

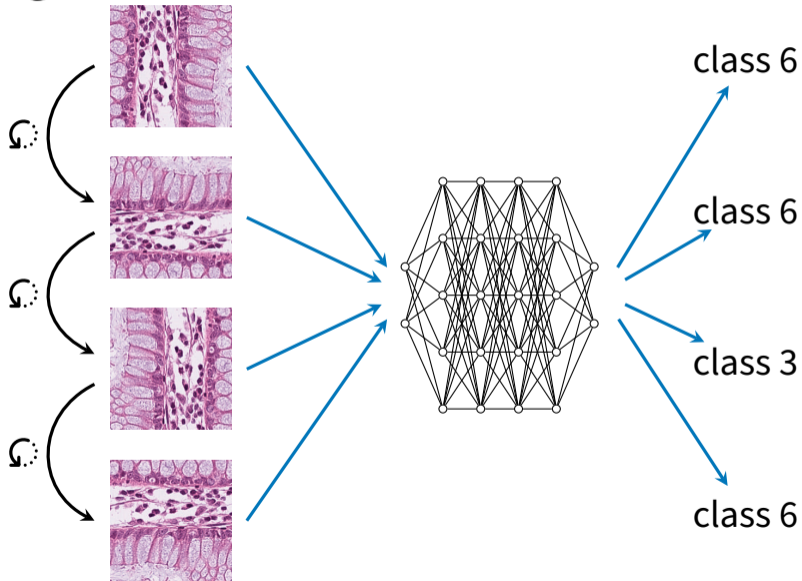# Histological slices
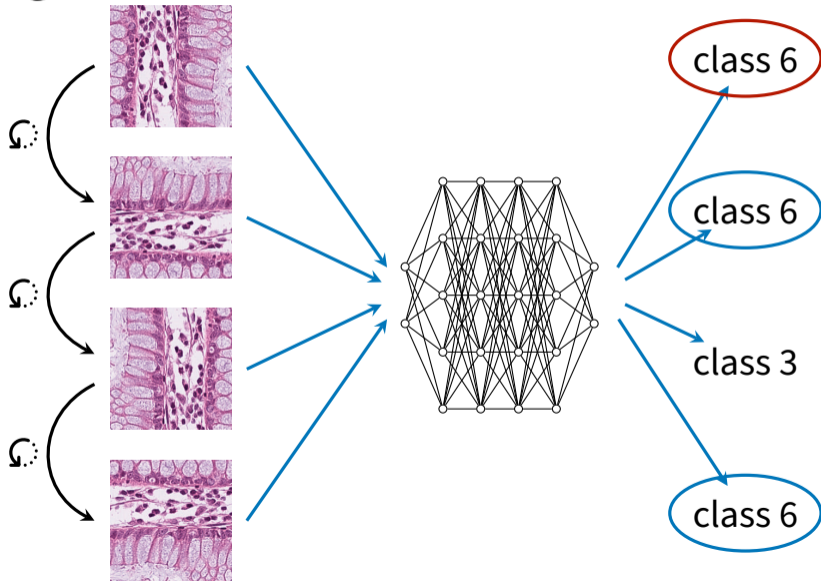


class 6

# Histological slices



class 6

class 6

class 3

class 6

# Histological slices

# Histological slices



Orbit Same Predictions = 3

class 6

class 6

class 3

class 6

38

# Out of distribution results



Ensemble size 5

# Out of distribution results



Ensemble size 5

Perfect invariance

Orbit same predictions vs Epoch

- 22.5° (red, at 16)
- 30° (green, at 12)
- 45° (orange, at 8)
- 90° (blue, at 4)

# Out of distribution results



Perfect invariance

Ensemble size 5

Orbit same predictions

22.5°
30°
45°
90°

Epoch

Ensemble members

# Out of distribution results



Ensemble size 5

Perfect invariance

# Out of distribution results



Ensemble size 20

Perfect invariance

22.5°
30°
45°
90°

Orbit same predictions

16
12
8
4

Epoch
2 4 6 8 10 12 14 16 18 20

Ensemble members — Ensemble

# Further experimental results

# Further experimental results

$\checkmark$ 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

- Geometric deep learning and equivariant neural networks
  Jan E. Gerken, Jimmy Aronsson, Oscar Carlsson, Hampus Linander,
  Fredrik Ohlsson, Christoffer Petersson, Daniel Persson

  Artificial Intelligence Review 2023

- Emergent Equivariance in Deep Ensembles
  Jan E. Gerken*, Pan Kessel*

  ICML 2024 (Oral)

  * Equal contribution



Group Website