

Emergent Equivariance in Deep Ensembles

Jan E. Gerken



CHALMERS
UNIVERSITY OF TECHNOLOGY



UNIVERSITY OF
GOTHENBURG



in collaboration with



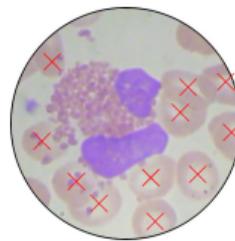
Pan Kessel



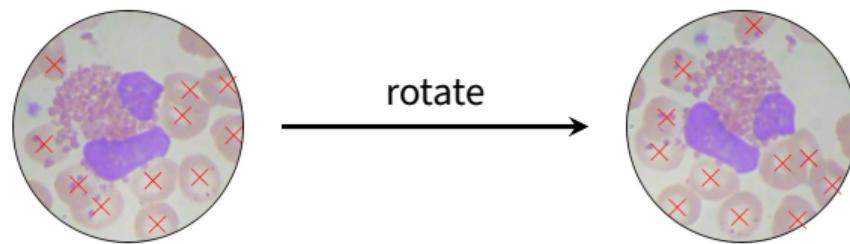
Philipp Misof

Symmetries in deep learning

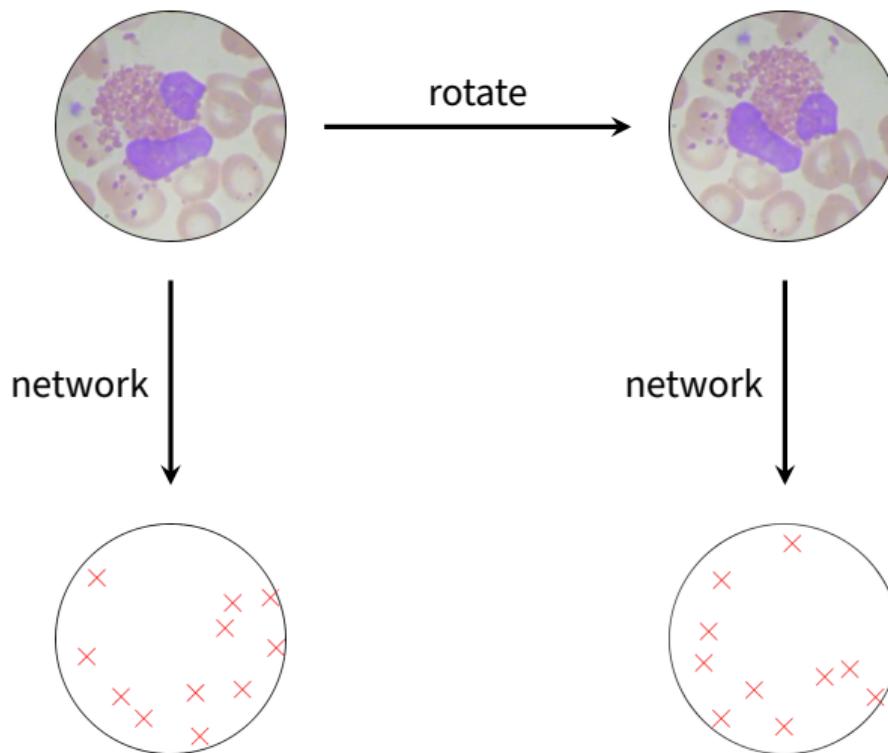
Symmetries in deep learning



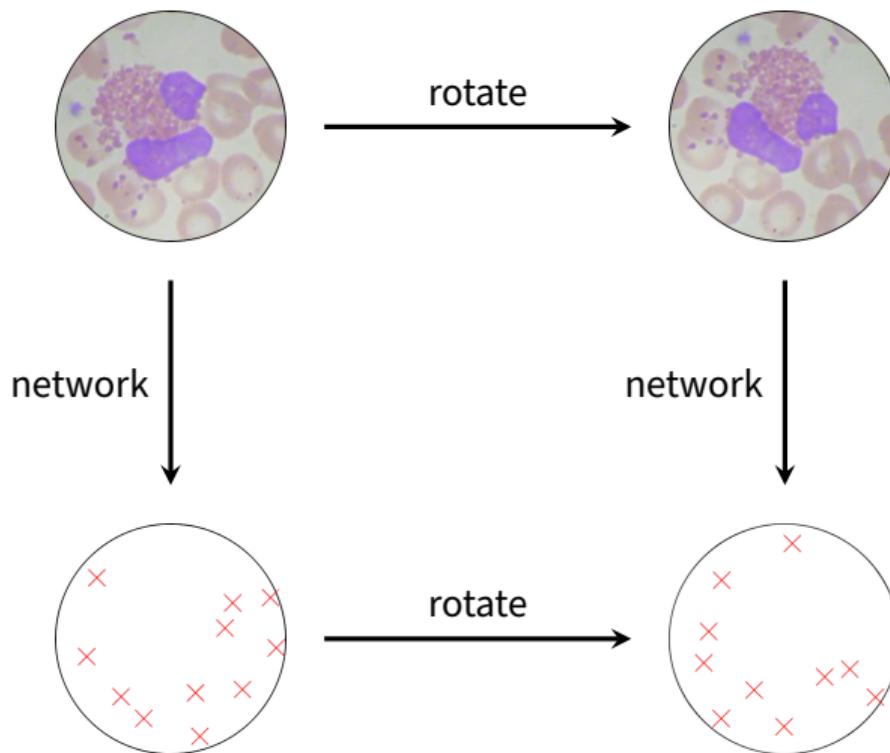
Symmetries in deep learning



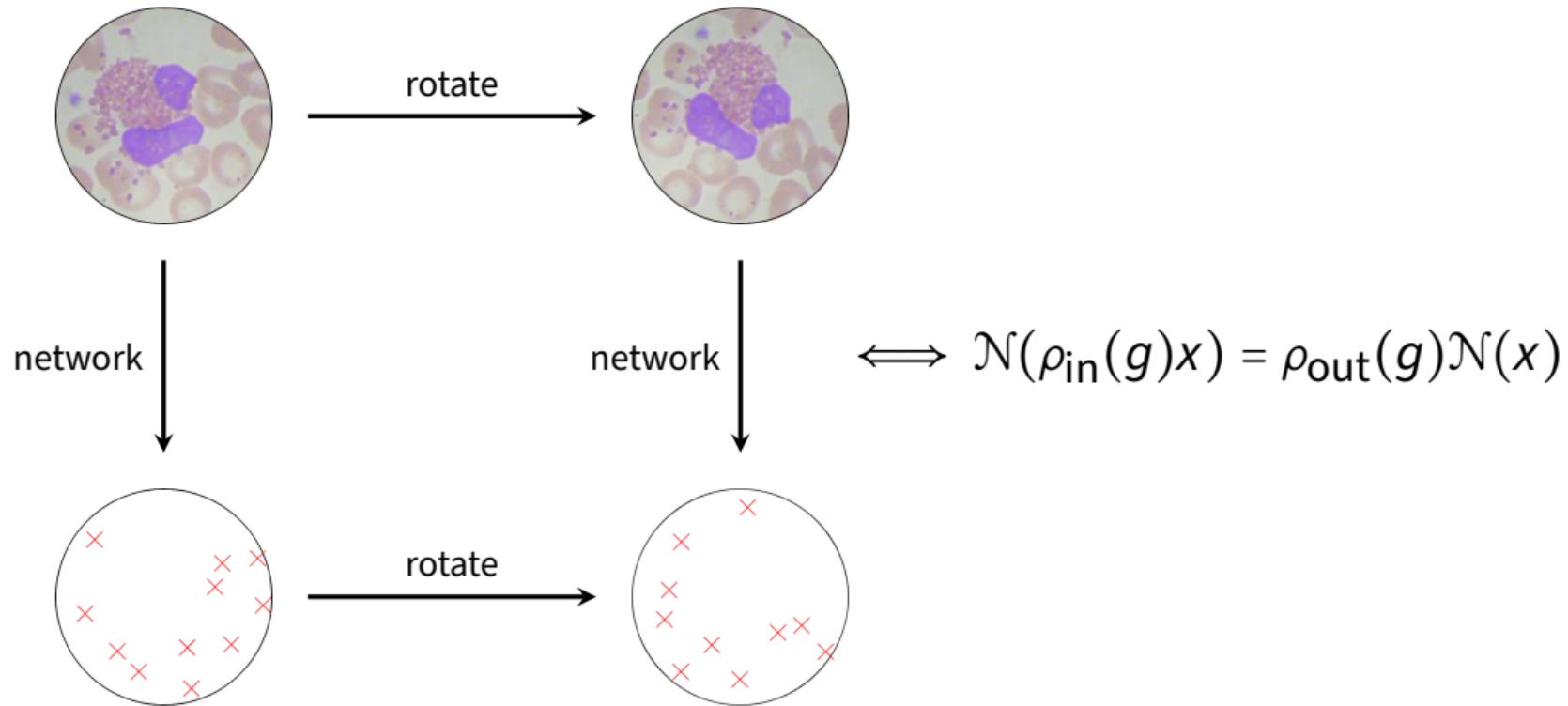
Symmetries in deep learning



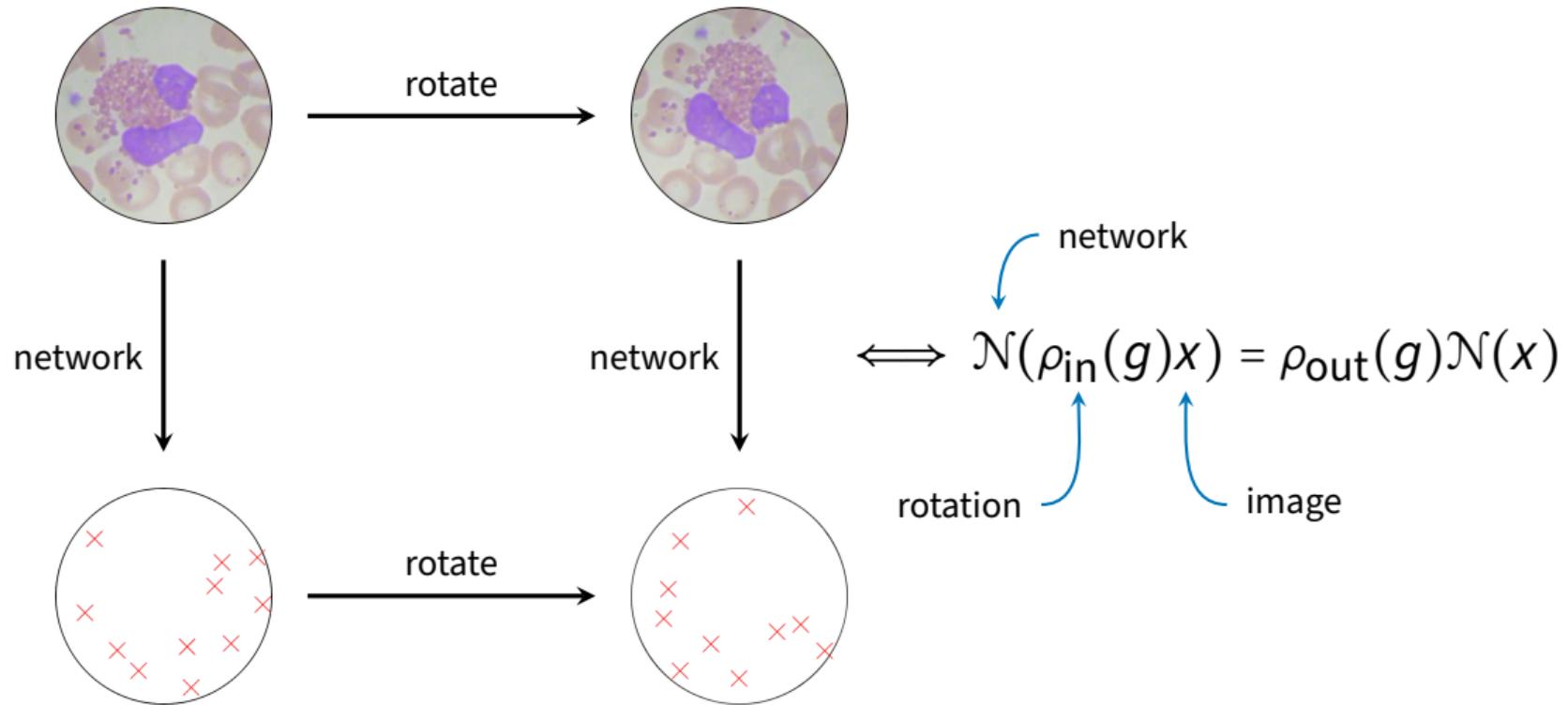
Symmetries in deep learning



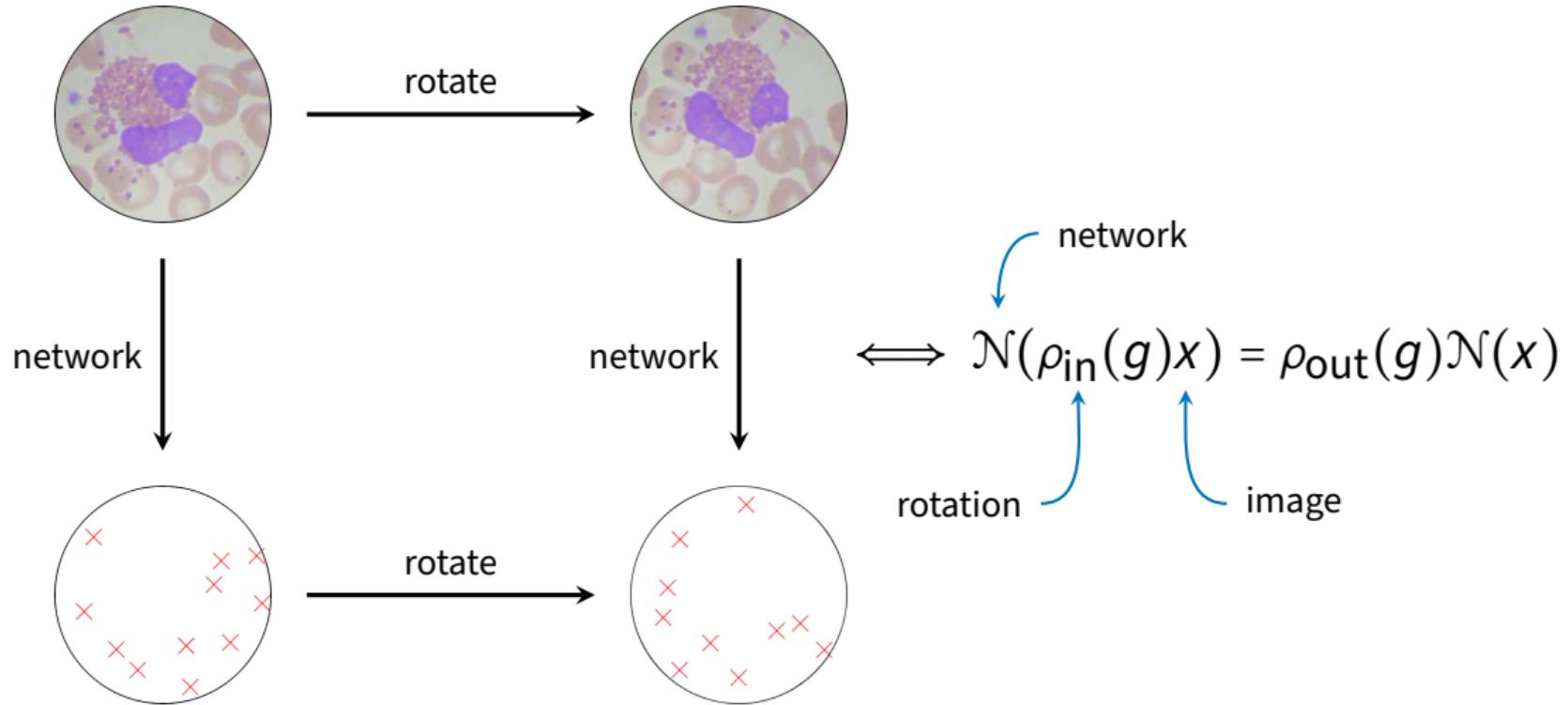
Symmetries in deep learning



Symmetries in deep learning



Equivariance



Equivariant neural networks

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen
University of Amsterdam

T.S.COHEN@UVA.NL

Max Welling
University of Amsterdam
University of California Irvine
Canadian Institute for Advanced Research

M.WELLING@UVA.NL

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

Convolution layers can be used effectively in a deep network because all the layers in such a network are *translation equivariant*: shifting the image and then feeding it through a number of layers is the same as feeding the original image through the same layers and then shifting the resulting feature maps (at least up to edge-effects). In

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen

University of Amsterdam

Max Welling

University of Amsterdam

University of California Irvine

Canadian Institute for Advanced Research

T.S.COHEN@UVA.NL

M.WELLING@UVA.NL

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

Convolution layers can be used effectively in a deep network because all the layers in such a network are *translation equivariant*: shifting the image and the feeding it through a number of layers is the same as feeding the original image through the same layers and then shifting the resulting feature maps (at least up to edge-effects). In

Equivariant Transformer Networks

Kai Sheng Tai¹ Peter Bailis¹ Gregory Valiant¹

Abstract

How can prior knowledge on the transformation invariance of a domain be incorporated into the architecture of a neural network? We propose Equivariant Transformers (ETs), a family of differentiable function-to-image mappings that preserve the robustness of neural networks to pre-defined continuous transformation groups. Through the use of specially-derived canonical coordinate systems, ETs incorporate functions that

scaling to each training image). While data augmentation typically helps reduce the test error of CNN-based models, there is no guarantee that transformation invariance will be enforced for data not seen during training.

In contrast to training time approaches like data augmentation, recent work on group equivariant CNNs (Cohen & Welling, 2016; Dierkes et al., 2016; Marcus et al., 2017; Woern et al., 2017; Tancik et al., 2017; Alabd, 2017; Cohen et al., 2018) has explored new CNN architectures that are designed to encode specific-to-domain invariances.

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen

University of Amsterdam

Max Welling

University of Amsterdam

University of California Irvine

Canadian Institute for Advanced Research

T.S.COHEN@UVA.NL

M.WELLING@UVA.NL

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

Convolution layers can be used effectively in a deep network because all the layers in such a network are *translation equivariant*: shifting the image and the feeding it through a number of layers is the same as feeding the original image through the same layers and then shifting the resulting feature maps (at least up to edge-effects). In

Equivariant Transformer Networks

Kai Sheng Tai¹ Peter Bailis¹ Gregory Valiant¹

Abstract

How can prior knowledge on the transformation invariance of a domain be incorporated into the architecture of a neural network? We propose Equivariant Transformers (ETs), a family of differentiable end-to-image mappings that preserve the robustness of neural networks to pre-defined continuous transformation groups. Through the use of specially-derived canonical coordinate systems, ETs implement functions that

scaling to each training image). While data augmentation typically helps reduce the test error of CNN-based models, there is no guarantee that transformation invariance will be enforced for data not seen during training.

In contrast to training time approaches like data augmentation, recent work on group equivariant CNNs (Cohen & Welling, 2016; Diefenbach et al., 2016; Marcus et al., 2017; Worrn et al., 2017; Mallya et al., 2017; Mallya, 2017; Cohen et al., 2018) has explored new CNN architectures that are designed to encode modifiable invariance constraints.

Theory for Equivariant Quantum Neural Networks

Quynh T. Nguyen,^{1,2} Louis Schatzki,^{3,4} Paolo Branca,^{1,5} Michael Rapone,^{1,6} Patrick J. Coles,³ Frédéric Sauvage,⁴ Martin Loeffelholz,^{1,7} and M. Cirne,³

¹Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA
²School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA

³Information Sciences, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

⁴Department of Electrical and Computer Engineering, University of Illinois Urbana-Champaign, Urbana, Illinois 61801, USA

⁵Department of Mathematics, University of Southern California, Los Angeles, California 90089, USA

⁶Department of Mathematics, University of California Davis, Davis, California 95616, USA
⁷Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

Quantum neural network architectures that have little-to-no inductive biases are known to face trainability and generalization issues. Inspired by a similar problem, recent breakthroughs in machine learning address this challenge by creating models encoding the symmetries of the learning task. This is materialized through the usage of equivariant neural networks whose action commutes with that of the symmetry. In this work we extend these ideas to the quantum regime by

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen

University of Amsterdam

Max Welling

University of Amsterdam

University of California Irvine

Canadian Institute for Advanced Research

T.S.COHEN@UVA.NL

M.WELLING@UVA.NL

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

Convolution layers can be used effectively in a deep network because all the layers in such a network are *translation equivariant*: shifting the image and feeding it through a number of layers is the same as feeding the original image through the same layers and then shifting the resulting feature maps (at least up to edge-effects). In

Equivariant Transformer Networks

Kai Sheng Tai¹ Peter Bailis¹ Gregory Valiant¹

Abstract

How can prior knowledge on the transformation invariance of a domain be incorporated into the architecture of a neural network? We propose Equivariant Transformers (ETs), a family of differentiable end-to-image mappings that preserve the robustness of neural networks to pre-defined continuous transformation groups. Through the use of specially-derived canonical coordinate systems, ETs implement functions that

scaling to each training image). While data augmentation typically helps reduce the test error of CNN-based models, there is no guarantee that transformation invariance will be enforced for data not seen during training.

In contrast to training time approaches like data augmentation, recent work on group equivariant CNNs (Cohen & Welling, 2016; Diefenbach et al., 2016; Marcus et al., 2017; Worrall et al., 2017; Mallya et al., 2017; Cohen et al., 2018) has explored new CNN architectures that are designed to encode modifiable invariance constraints.

Theory for Equivariant Quantum Neural Networks

Quynh T. Nguyen,^{1,2} Louis Schatzki,^{3,4} Paolo Branca,^{1,5} Michael Rapone,^{1,6} Patrick J. Coles,³ Frédéric Sauvage,⁴ Martin Lucca,^{1,7} and M. Cirone³

¹Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA
²School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA

³Information Sciences, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

⁴Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA
⁵Department of Physics, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

⁶Department of Mathematics, University of California Davis, Davis, California 95616, USA
⁷Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

Quantum neural network architectures that have little-to-no inductive biases are known to face trainability and generalization issues. Inspired by a similar problem, recent breakthroughs in machine learning address this challenge by creating models encoding the symmetries of the learning task. This is materialized through the usage of equivariant neural networks whose action commutes with that of the symmetry. In this work we extend these ideas to the quantum setting by

An Efficient Lorentz Equivariant Graph Neural Network for Jet Tagging

Shiqi Gong^{a,c,1} Qi Meng^b Jue Zhang^b Huilin Qu^c Congqiao Li^c Sitian Qian^d Weitao Du^a Zhi-Ming Ma^a Tie-Yan Liu^b

^aAcademy of Mathematics and Systems Science, Chinese Academy of Sciences,
Zhongguancun East Road, Beijing 100190, China

^bMicrosoft Research Asia,
Duning Street, Beijing 100089, China

^cCERN, EP Department,
CH-1211 Geneva 23, Switzerland

^dSchool of Physics, Peking University,
Chenfu Road, Beijing 100871, China

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen

University of Amsterdam

Max Welling

University of Amsterdam

University of California Irvine

Canadian Institute for Advanced Research

T.S.COHEN@UVA.NL

M.WELLING@UVA.NL

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

Convolution layers can be used effectively in a deep network because all the layers in such a network are *translation equivariant*: shifting the image and feeding it through a number of layers is the same as feeding the original image through the same layers and then shifting the resulting feature maps (at least up to edge-effects). In

An Efficient Lorentz Equivariant Graph Neural Network for Jet Tagging

Shiqi Gong^{a,1} Qi Meng^b Jue Zhang^b Huilin Qu^c Congqiao Li^c Sitian Qian^d Weitao Du^a Zhi-Ming Ma^a Tie-Yan Liu^b

^aAcademy of Mathematics and Systems Science, Chinese Academy of Sciences,
Zhongguancun East Road, Beijing 100190, China

^bMicrosoft Research Asia,
Daming Street, Beijing 100089, China

^cCERN, EP Department,
CH-1211 Geneva 23, Switzerland

^dSchool of Physics, Peking University,
Chenfu Road, Beijing 100871, China

Equivariant Transformer Networks

Kai Sheng Tai¹ Peter Bailis¹ Gregory Valiant¹

Abstract

How can prior knowledge on the transformation invariance of a domain be incorporated into the architecture of a neural network? We propose Equivariant Transformers (ETs), a family of differentiable end-to-image mappings that preserve the robustness of neural networks to pre-defined continuous transformation groups. Through the use of symmetry-derived canonical coordinate systems, ETs incorporate functions that

scaling to each training image). While data augmentation typically helps reduce the test error of CNN-based models, there is no guarantee that transformation invariance will be enforced for data not seen during training.

In contrast to training time approaches like data augmentation, recent work on group equivariant CNNs (Cohen & Welling, 2016; Diefenbach et al., 2016; Marcus et al., 2017; Worrall et al., 2017; Mallya et al., 2017; Cohen et al., 2018) has explored new CNN architectures that are

Theory for Equivariant Quantum Neural Networks

Quynh T. Nguyen,^{1,2} Louis Schatzki,^{3,4} Paolo Branca,^{1,5} Michael Rapone,^{1,6}
Patrick J. Coles,³ Frédéric Sauvage,⁴ Martin Lucca,^{1,7} and M. Cirone³

¹Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA
²School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA

³Information Sciences, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

⁴Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

⁵Department of Physics, Aerospace and Mechanical Engineering, Santa Fe Institute, Santa Fe, New Mexico 87501, USA

⁶Department of Mathematics, University of California Davis, Davis, California 95616, USA

⁷Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

Quantum neural network architectures that have little-to-no inductive biases are known to face quantum stability and generalization issues. Inspired by a similar problem, recent breakthroughs in machine learning address this challenge by creating models encoding the symmetries of the learning task. This is materialized through the usage of equivariant neural networks whose action commutes with that of the symmetry. In this work we present theory for the quantum analogues for

E(3)-Equivariant Graph Neural Networks for Data-Efficient and Accurate Interatomic Potentials

Simon Batzner^{a,1} Albert Musoelian,¹ Lixin Sun,¹ Mario Geiger,² Jonathan P. Maillet,³
Mordechai Kornblith,² Nicola Molinari,¹ Tess E. Smith,^{4,5} and Boris Kozinsky^{a,1,3}

¹John A. Paulson School of Engineering and Applied Sciences,
Harvard University, Cambridge, MA 02138, USA

²École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland
³Robert Bosch Research and Technology Center, Cambridge, MA 02120, USA

⁴Computational Research Division and Center for Advanced Mathematics for Energy Research Applications,
Lawrence Berkeley National Laboratory, Berkeley, CA 94730, USA
⁵Massachusetts Institute of Technology, Department of Electrical
Engineering and Computer Science, Cambridge, MA 02139, USA

This work presents Neural Equivariant Interatomic Potentials (NeqIP), an E(3)-equivariant neural network approach for learning interatomic potentials from *ab-initio* calculations for molecular dynamics simulations. While most contemporary symmetry-aware models use invariant convolutions and only act on scalars, NeqIP employs E(3)-equivariant convolutions for interactions of geometric tensors, resulting in a more information-rich and faithful representation of atomic environments. The method achieves state-of-the-art accuracy on a challenging and diverse set of molecules and

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen

University of Amsterdam

Max Welling

University of Amsterdam

University of California Irvine

Canadian Institute for Advanced Research

T.S.COHEN@UVA.NL

M.WELLING@UVA.NL

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

Convolution layers can be used effectively in a deep network because all the layers in such a network are *translation equivariant*: shifting the image and feeding it through a number of layers is the same as feeding the original image through the same layers and then shifting the resulting feature maps (at least up to edge-effects). In

Equivariant Transformer Networks

Kai Sheng Tai¹ Peter Bailis¹ Gregory Valiant¹

Abstract

How can prior knowledge on the transformation invariance of a domain be incorporated into the architecture of a neural network? We propose Equivariant Transformers (ETs), a family of differentiable end-to-image mappings that preserve the robustness of neural networks to pre-defined continuous transformations. Through the use of specially-derived canonical coordinate systems, ETs implement functions that

scaling to each training image). While data augmentation typically helps reduce the test error of CNN-based models, there is no guarantee that transformation invariance will be enforced for data not seen during training.

In contrast to training time approaches like data augmentation, recent work on group equivariant CNNs (Cohen & Welling, 2016; Dilemmas et al., 2016; Marcus et al., 2017; Worrn et al., 2017; Mallya et al., 2017; Cohen et al., 2018) has explored new CNN architectures that are designed to encode modularity via invariant mechanisms

Theory for Equivariant Quantum Neural Networks

Quynh T. Nguyen,^{1,2} Louis Schatzki,^{3,4} Paolo Branca,^{1,5} Michael Rapone,^{1,6} Patrick J. Coles,¹ Frédéric Sauvage,⁴ Martin Laločka,^{1,7} and M. Černý,³

¹Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA
²School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA

³Information Sciences, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

⁴Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

⁵Department of Mathematics, University of Southern California, Los Angeles, California 90089, USA

⁶Department of Mathematics, University of California Davis, Davis, California 95616, USA
⁷Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

Quantum neural network architectures that have little-to-no inductive biases are known to face trainability and generalization issues. Inspired by a similar problem, recent breakthroughs in machine learning address this challenge by creating models encoding the symmetries of the learning task. This is materialized through the usage of equivariant neural networks whose action commutes with that of the symmetry. In this work we present those to the quantum setting

HIERARCHICAL, ROTATION-EQUIVARIANT NEURAL NETWORKS TO SELECT STRUCTURAL MODELS OF PROTEIN COMPLEXES

Stephan Eismann*

Department of Applied Physics
Stanford University
seimann@stanford.edu

Raphael J.L. Townsend*

Department of Computer Science
Stanford University
raphael@cs.stanford.edu

Nathaniel Thomas*

Department of Physics
Stanford University
nthomas103@gmail.com

Mihail Jagota

Department of Electrical Engineering
Stanford University
mjagota@stanford.edu

Bowen Jing

Department of Computer Science
Stanford University
bjing@cs.stanford.edu

Ron O. Dror

Department of Computer Science
Stanford University
rondror@cs.stanford.edu

ABSTRACT

Predicting the structure of multi-protein complexes is a grand challenge in biochemistry, with major implications for basic science and drug discovery. Computational structure prediction methods generally leverage pre-defined structural features to distinguish accurate structural models from less accurate ones. This raises the question of whether it is possible to learn characteristics of accurate models directly from atomic coordinates of protein complexes, with no prior assumptions. Here we introduce a machine learning method that learns directly from the 3D positions of all atoms to

An Efficient Lorentz Equivariant Graph Neural Network for Jet Tagging

Shiqi Gong^{a,1} Qi Meng^b Jue Zhang^b Huilin Qu^c Congqiao Li^c Sitian Qian^d Weitao Du^a Zhi-Ming Ma^a Tie-Yan Liu^b

^aAcademy of Mathematics and Systems Science, Chinese Academy of Sciences, Zhongguancun East Road, Beijing 100190, China

^bMicrosoft Research Asia, Daming Street, Beijing 100089, China

^cCERN, EP Department, CH-1211 Geneva 23, Switzerland

^dSchool of Physics, Peking University, Chenfu Road, Beijing 100871, China

E(3)-Equivariant Graph Neural Networks for Data-Efficient and Accurate Interatomic Potentials

Simon Batzner^{a,1} Albert Musoelian¹ Lixin Sun¹ Mario Geiger² Jonathan P. Mallon³ Mordechai Kornbluth² Nicola Molinari¹ Tess E. Smith^{4,5} and Boris Kozinsky^{a,1,3}

¹John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138, USA

²École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland

³Robert Bosch Research and Technology Center, Cambridge, MA 02130, USA

⁴Computational Research Division, and Center for Advanced Computing for Energy Research Applications, Lawrence Berkeley National Laboratory, Berkeley, CA 94730, USA

⁵Massachusetts Institute of Technology, Department of Electrical Engineering and Computer Science, Cambridge, MA 02132, USA

This work presents Neural Equivariant Interatomic Potentials (NeqIP), an E(3)-equivariant neural network approach for learning interatomic potentials from *ab initio* calculations for molecular dynamics simulations. While most contemporary symmetry-aware models use invariant convolutions and only act on scalars, NeqIP employs E(3)-equivariant convolutions for interactions of geometric tensors, resulting in a more information-rich and faithful representation of atomic environments. The method achieves state-of-the-art accuracy on a challenging and diverse set of molecules and

Equivariant neural networks

Group Equivariant Convolutional Networks

Taco S. Cohen
University of Amsterdam
Max Welling
University of Amsterdam
University of California Irvine
Canadian Institute for Advanced Research

Abstract

We introduce Group equivariant Convolutional Neural Networks (G-CNNs), a natural generalization of convolutional neural networks that reduces sample complexity by exploiting symme-

T.S.COHEN@UVA.NL
M.WELLING@UVA.NL

M.WELLING@UVA.NL

An Efficient Lorentz Equivariant Graph Neural Network for Jet Tagging

Shiqi Gong^{a,b,1} Qi Meng^b Jue Zhang^b Huilin Qu^c Congqiao Li^c Sitian Qian^d Weitao Du^a Zhi-Ming Ma^a Tie-Yan Liu^b

^aAcademy of Mathematics and Systems Science, Chinese Academy of Sciences,
Zhongguancun East Road, Beijing 100190, China

^bMicrosoft Research Asia,
Daming Street, Beijing 100089, China

^cCERN, EP Department,
CH-1211 Geneva 23, Switzerland

^dSchool of Physics, Peking University,
Chenfu Road, Beijing 100871, China

Equivariant Transformer Networks

Karl Kötter¹, Pratik Jain¹, Balaji Raghav¹, Caglar Gulcehre², Myle Ott¹

Geometric Deep Learning and Equivariant Neural Networks

JAN E. GERKEN¹, JIMMY ARONSSON^{1*}, OSCAR CARLSSON^{1*}, HÄMUS LINANDER²,
FREDRIK OHLSSON³, CHRISTOFER PETERSSON^{1,4} AND DANIEL PERSSON¹

¹ Chalmers University of Technology, Department of Mathematical Sciences
SE-412 96 Gothenburg, Sweden

² Gothenburg University, Department of Physics
SE-412 96, Gothenburg, Sweden

³ Umeå University, Department of Mathematics and Mathematical Statistics
SE-901 87, Umeå, Sweden

⁴ Zenseact
SE-417 56, Gothenburg, Sweden

* equal contribution

Neural Equivariant Interatomic Potentials (NeqIP)
Engineering and Computer Science, Cambridge, MA 02139, USA

This work presents Neural Equivariant Interatomic Potentials (NeqIP), an E(3)-equivariant neural network approach for learning interatomic potentials from *ab initio* calculations for molecular dynamics simulations. While most contemporary symmetry-aware models use invariant convolutions and only act on scalars, NeqIP employs E(3)-equivariant convolutions for interactions of geometric tensors, resulting in a more information-rich and faithful representation of atomic environments. The method achieves state-of-the-art accuracy on a challenging and diverse set of molecules and

Theory for Equivariant Quantum Neural Networks

Quynh T. Nguyen^{1,2}, Louis Schatzki^{3,4}, Paolo Branca^{1,5}, Michael Rapone^{1,6},
Patrick J. Coles¹, Frédéric Sauvage¹, Martin Läzcano^{1,7} and M. Cirne^{1,8}

¹ Theoretical Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA
² School of Engineering and Applied Sciences, Harvard University, Cambridge, Massachusetts 02138, USA

³ Information Sciences, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

⁴ Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

⁵ Department of Mathematics, University of Southern California, Los Angeles, California 90089, USA

⁶ Department of Mathematics, University of California Davis, Davis, California 95616, USA
⁷ Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA

Quantum neural network architectures that have little-to-no inductive biases are known to face trainability and generalization issues. Inspired by a similar problem, recent breakthroughs in machine learning address this challenge by creating models encoding the symmetries of the learning task. This is materialized through the usage of equivariant neural networks whose action commutes with that of the symmetry. In this work we present theory close to the quantum scale for

ERARCHICAL, ROTATION-EQUIVARIANT NEURAL NETWORKS TO SELECT STRUCTURAL MODELS OF PROTEIN COMPLEXES

Stephan Eismann^{*}
Department of Applied Physics
Stanford University
seismann@stanford.edu

Raphael J.L. Townsend^{*}
Department of Computer Science
Stanford University
raphael@cs.stanford.edu

Nathaniel Thomas^{*}
Department of Physics
Stanford University
nthomas103@gmail.com

Mihir Jagota
Department of Electrical Engineering
Stanford University
mjagota@stanford.edu

Bowen Jing
Department of Computer Science
Stanford University
bjing@cs.stanford.edu

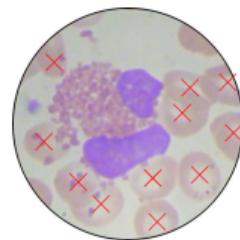
Ron O. Dror
Department of Computer Science
Stanford University
rondror@cs.stanford.edu

ABSTRACT

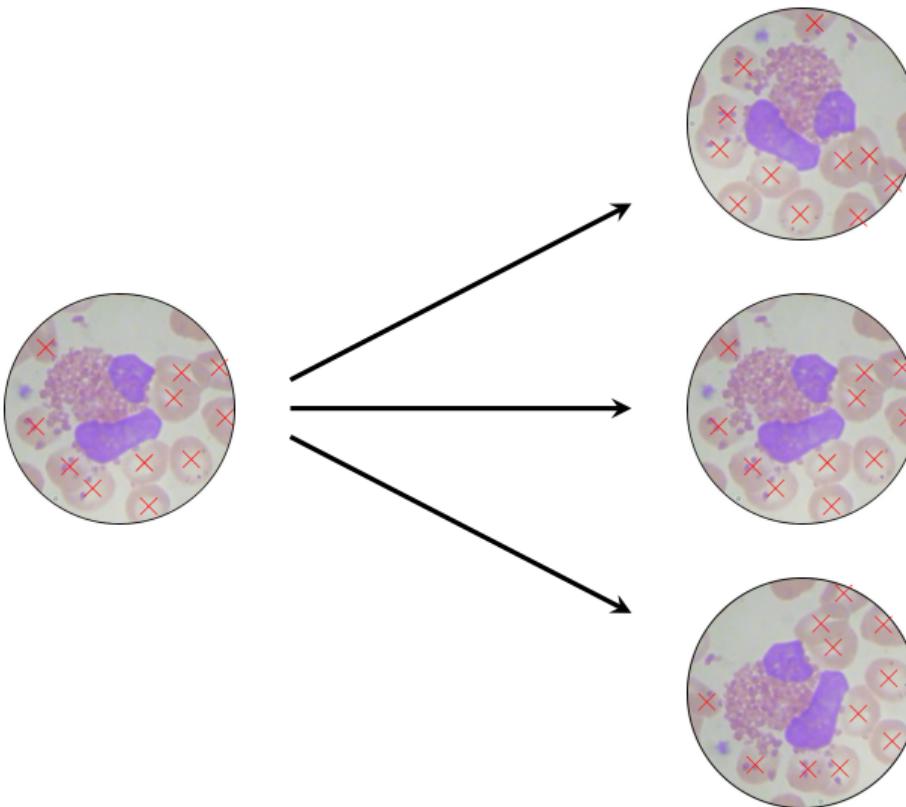
Predicting the structure of multi-protein complexes is a grand challenge in biochemistry, with major implications for basic science and drug discovery. Computational structure prediction methods generally leverage pre-defined structural features to distinguish accurate structural models from less accurate ones. This raises the question of whether it is possible to learn characteristics of accurate models directly from atomic coordinates of protein complexes, with no prior assumptions. Here we introduce a machine learning method that learns directly from the 3D positions of all atoms to

Data augmentation

Data augmentation



Data augmentation



Data augmentation

Data augmentation

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

<https://doi.org/10.1038/s43998-024-07487-w>

Received: 19 December 2023

Accepted: 29 April 2024

Published online: 8 May 2024

Open access

Check for updates

Josh Abramson^{1,2}, Jonas Adler³, Jack Danger^{1,4}, Richard Evans⁵, Tim Green^{2,6}, Alexander Pritzel^{2,7}, Olaf Ronneberger⁸, Lindsay Williams⁹, Andrew J. Bellard¹⁰, Jennifer C. Boudelle¹¹, Sebastian V. Bräuer¹², David A. Brown¹³, Chaitanya Chakraborty¹⁴, Michael O’Neill¹⁵, Daniel Riedl¹⁶, Kathryn Sawyer¹⁷, Zucheng Shao¹⁸, Nalini Sengupta¹⁹, Estri Aravant²⁰, Charles Boettig²¹, Ottavia Bartoli²², Alex Bridgland²³, Aleney Chempurac²⁴, Miles Congreve²⁵, Alexander L. Cowen-Rivers²⁶, Andrew Cowie²⁷, Michael Figariros²⁸, Michael Gromiha²⁹, Michael D. Hargreaves³⁰, Daniel Young J. Khor³¹, Christopher M. R. Lew³², Koko Perlin³³, Anil Pratap³⁴, Pouya Rayv³⁵, Sulabh Ray³⁶, Adrien Reinaud³⁷, Ashwak Thivarusundaram³⁸, Catherine Tong³⁹, Sergei Vakser⁴⁰, Ellen S. Zhang⁴¹, Michal Zelma⁴², Augustin Žíšek⁴³, Victor Kapit⁴⁴, Puthrenet Kochi⁴⁵, Max Jaderberg^{46,47}, Dennis Hesselbar^{48,49} & John M. Jumper^{1,2,50}

The introduction of AlphaFold 2 has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design^{1–4}. Here we describe the AlphaFold 3 model with a substantially updated architecture that can predict the structures of proteins and complexes of complex systems including proteins, nucleic acids, small molecules, ions and modified residues. The new AlphaFold model demonstrates substantially improved accuracy

Data augmentation

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

<https://doi.org/10.1101/244098>

Received: 19 December 2023

Accepted: 29 April 2024

Published online: 8 May 2024

Open access

Check for updates

Josh Abramson^{1,2}, Jonas Adler³, Jack Danger^{1,4}, Richard Evans⁵, Tim Green², Alexander Grillet², Olaf Ronneberger^{1,6}, Lindsay Willmott², Andrew J. Bellard¹, Jennifer Cao¹, Sebastian V. Engel¹, David A. Fiser¹, Michael G. Clark^{1,7}, Michael O'Neill¹, Daniel Reiter¹, Kathryn Sawyer¹, Alexander Shokhob¹, Zachary S. Smith¹, Natasja Sengulay¹, Ester Avantaggiati¹, Charles Boettig¹, Ottavia Bartoli^{1,8}, Alex Bridgland¹, Alexey Cherezov¹, Miles Congreve¹, Alexander L. Cowen-Rivers¹, Andrew Cowie¹, Michael Figari¹, Miles Figari¹, Michael G. Fiser¹, Daniel Gitter¹, Michael H. Gitter¹, Yannick K. Hildebrand¹, Christopher M. H. Lew¹, Koko Perner¹, Anil Purohit¹, Pouyan Raveh¹, Sulabh Ray¹, Michael Rhee¹, Adrien Rhee¹, Ashvak Thivaisundararao¹, Catherine Tong¹, Sergei Vakser¹, Ellen S. Zhang¹, Michal Zelma¹, Augustin Žíšek¹, Victor Bapst¹, Pudvezet Koch¹, Max Jaderberg^{1,9}, Dennis Hesselbar^{1,10} & John M. Jumper^{1,11}

The introduction of AlphaFold 2 has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design^{1–4}. Here we describe the AlphaFold 3 model with a substantially updated architecture and characteristics, including improved performance on a wide range of complex systems including proteins, nucleic acids, small molecules, ions and modified residues. The new AlphaFold model demonstrates substantially improved accuracy

The Importance of Being Scalable: Improving the Speed and Accuracy of Neural Network Interatomic Potentials Across Chemical Domains

Eric Qu¹
UC Berkeley
ericqu@berkeley.edu

Aditi S. Krishnapriyan¹
UC Berkeley, LBNL
aditik1@berkeley.edu

Abstract

Scaling has been a critical factor in improving model performance and generalization across various fields of machine learning. It involves how a model's performance changes with increases in model size or input data, as well as how efficiently computational resources are utilized to support this growth. Despite successes in scaling other types of machine learning models, the study of scaling in Neural Network Interatomic Potentials (NNIPs) remains limited. NNIPs act as

Data augmentation

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

<https://doi.org/10.1038/s43242-024-07487-w>

Received: 19 December 2023

Accepted: 29 April 2024

Published online: 8 May 2024

Open access

Check for updates

Josh Abramson^{1,2}, Jonas Adler³, Jack Danger^{1,4}, Richard Evans⁵, Tim Green², Alexander Grütz⁶, Olaf Ronneberger⁷, Lindsay Willmott⁸, Andrew J. Bellard⁹, Jennifer Cao¹⁰, Sebastian V. Engel¹¹, David A. Fajardo-Chavez¹², Michael O’Neill¹³, Daniel Riedl¹⁴, Kathryn Raynor¹⁵, Zachary S. Rosenblatt¹⁶, Paula Sungayal¹⁷, Estri Aravant¹⁸, Charles Boettig¹⁹, Ottavia Bartoli²⁰, Alex Brigland²¹, Aleney Chompanac²², Miles Congreve²³, Alexander L. Cowen-Rivers²⁴, Andrew Cowie²⁵, Michael Figariro²⁶, Miles Figariro²⁷, Michael Gromadam²⁸, Daniel H. Hwang²⁹, Yannick K. Kihne³⁰, Catherine M. R. Lew³¹, Koko Peris³², Anil Purohit³³, Praveen Ray³⁴, Sulabh Ray³⁵, Aditi S. Ray³⁶, Ashvika Thivakaranarayanan³⁷, Catherine Tong³⁸, Sergei Vakser³⁹, Ellen S. Zhang⁴⁰, Michal Zelma⁴¹, Augustin Žídek⁴², Victor Bapst⁴³, Pudhvezh Kath⁴⁴, Max Jaderberg⁴⁵, Dennis Hesselbar⁴⁶ & John M. Jumper^{1,2}

The introduction of AlphaFold 2 has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design^{1–4}. Here we describe the AlphaFold 3 model with a substantially updated architecture that achieves state-of-the-art performance in the prediction of complex interactions, including protein–protein, protein–nucleic acid, protein–small molecule, and protein–modified residue. The new AlphaFold model demonstrates substantially improved accuracy

The Importance of Being Scalable: Improving the Speed and Accuracy of Neural Network Interatomic Potentials Across Chemical Domains

Eric Qu¹
UC Berkeley
ericqu@berkeley.edu

Aditi S. Krishnapriyan¹
UC Berkeley, LBNL
aditik1@berkeley.edu

Abstract

Scaling has been a critical factor in improving model performance and generalization across various fields of machine learning. It involves how a model's performance changes with increases in model size or input data, as well as how efficiently computational resources are utilized to support this growth. Despite successes in scaling other types of machine learning models, the study of scaling in Neural Network Interatomic Potentials (NNIPs) remains limited. NNIPs act as

Swallowing the Bitter Pill: Simplified Scalable Conformer Generation

Yuyang Wang¹, Ahmed A. Elbag^{1,2}, Navdeep Jolly³, Joshua M. Suskind¹, Miguel Ángel Bautista¹

Abstract

We present a novel way to predict molecular conformers through a simple formulation that sidesteps many of the heuristics of prior works and achieves state of the art results by using the advantages of scale. By training a diffusion generative model directly on 3D atomic positions without any constraints about the implicit structure of molecules (e.g., minimum bond angles) we are able to radically simplify structure generation and predict the number of other

is the vast complexity of the 3D structure space, encompassing factors such as bond lengths and torsional angles. Designing the molecular conformer space with no constraints, such as bond types and spatial arrangements determined by chiral centers, the conformational space experiences exponential growth with the expansion of the graph size and the number of rotatable bonds (Aszkenasy & Gomez-Bombarelli, 2022). This complicates brute force and exhaustive approaches, making them virtually unusable for even moderately small molecules. Systematic methods, like DMTGA (Hawkins et al., 2018),

Data augmentation

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

<https://doi.org/10.1038/s43998-024-07487-w>
Received: 19 December 2023
Accepted: 29 April 2024
Published online: 8 May 2024
Open access
Check for updates

The introduction of AlphaFold 2 has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design^{1–3}. Here we describe the AlphaFold 3 model with a substantially updated architecture that achieves state-of-the-art performance across a wide range of complex systems including proteins, nucleic acids, small molecules, ions and modified residues. The new AlphaFold model demonstrates substantially improved accuracy

Probing the effects of broken symmetries in machine learning

Marco F. Langer¹, Sergey N. Prodnikov² and Michele Ceriotti^{1*}

Laboratory of Computational Science and Modelling and National Centre for Computational Design and Discovery of Novel Materials MARVEL, Institute of Materials, École Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland

* Author to whom any correspondence should be addressed.

E-mail: michele.ceriotti@epfl.ch

Keywords: machine learning, symmetry-constrained models, atomistic modeling, molecular simulations

Supplementary material for this article is available [online](#)

Abstract

Symmetry is one of the most central concepts in physics, and it is no surprise that it has also been widely adopted as an inductive bias for machine-learning models applied to the physical sciences. This is especially true for models targeting the properties of matter at the atomic scale. Both established and state-of-the-art approaches, with almost no exceptions, are built to be exactly equivariant to translations, permutations, and rotations of the atoms. Incorporating symmetries—rotations in particular—constraints the model design space and implies more complicated architectures that are often also computationally demanding. There are indications

The Importance of Being Scalable: Improving the Speed and Accuracy of Neural Network Interatomic Potentials Across Chemical Domains

Eric Qu¹
UC Berkeley
ericqu@berkeley.edu

Aditi S. Krishnapriyan¹
UC Berkeley, LBNL
aditik1@berkeley.edu

Abstract

Scaling has been a critical factor in improving model performance and generalization across various fields of machine learning. It involves how a model's performance changes with increases in model size or input data, as well as how efficiently computational resources are utilized to support this growth. Despite successes in scaling other types of machine learning models, the study of scaling in Neural Network Interatomic Potentials (NNIPs) remains limited. NNIPs act as

Swallowing the Bitter Pill: Simplified Scalable Conformer Generation

Yuyang Wang¹, Ahmed A. Elbag^{1,2}, Navdeep Jolly³, Joshua M. Suskind¹, Miguel Ángel Bautista¹

Abstract

We present a novel way to predict molecular conformers through a simple formulation that sidesteps many of the heuristics of prior works and achieves state of the art results by using the advantages of scale. By training a diffusion generative model directly on 3D atomic positions without any constraints about the explicit structure of molecules (e.g., minimum bond angles) we are able to radically simplify structure generation and make it much more efficient.

is the vast complexity of the 3D structure space, encompassing factors such as bond lengths and torsional angles. Despite the molecular complexity, the number of conformers that satisfy specific constraints, such as bond types and spatial arrangements determined by chiral centers, the conformational space experiences exponential growth with the expansion of the graph size and the number of rotatable bonds (Aszkenasy & Gomez-Bombarelli, 2022). This complicates brute force and exhaustive approaches, making them virtually unusable for even moderately small molecules. Systematic methods, like DMFGA (Hawkins et al., 2018),

Data augmentation

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

<https://doi.org/10.1038/s43247-024-07487-w>
Received: 19 December 2023
Accepted: 29 April 2024
Published online: 8 May 2024
Open access
Check for updates

The introduction of AlphaFold 2 has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design^{1–3}. Here we describe the AlphaFold 3 model with a substantially updated architecture that achieves state-of-the-art performance in predicting the structure of complex systems including proteins, nucleic acids, small molecules, ions and modified residues. The new AlphaFold model demonstrates substantially improved accuracy

Probing the effects of broken symmetries in machine learning

Marc F Langer¹, Sergey N Prochnikov² and Michele Ceriotti¹

Laboratory of Computational Science and Modelling and National Centre for Computational Design and Discovery of Novel Materials MARVEL, Institute of Materials, Ecole Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland

¹ Author to whom any correspondence should be addressed.

E-mail: michele.ceriotti@epfl.ch

Keywords: machine learning, symmetry-constrained models, atomistic modeling, molecular simulations

Supplementary material for this article is available [online](#)

Abstract

Symmetry is one of the most central concepts in physics, and it is no surprise that it has also been widely adopted as an inductive bias for machine-learning models applied to the physical sciences. This is especially true for models targeting the properties of matter at the atomic scale. Both established and state-of-the-art approaches, with almost no exceptions, are built to be exactly equivariant to translations, permutations, and rotations of the atoms. Incorporating symmetries—rotations in particular—constraints the model design space and implies more complicated architectures that are often also computationally demanding. There are indications

The Importance of Being Scalable: Improving the Speed and Accuracy of Neural Network Interatomic Potentials Across Chemical Domains

Eric Qu¹
UC Berkeley
ericqu@berkeley.edu

Aditi S. Krishnapriyan¹
UC Berkeley, LBNL
aditik1@berkeley.edu

Abstract

Scaling has been a critical factor in improving model performance and generalization across various fields of machine learning. It involves how a model's performance changes with increases in model size or input data, as well as how efficiently computational resources are utilized to support this growth. Despite successes in scaling other types of machine learning models, the study of scaling in Neural Network Interatomic Potentials (NNIPs) remains limited. NNIPs act as

Two for One: Diffusion Models and Force Fields for Coarse-Grained Molecular Dynamics

Marioes Arts,^{1,1,2,3} Victor García Satorras,^{1,4,5} Chin-Wei Huang,¹ Daniel Zügner,⁵ Marco Federici,^{1,3} Cecilia Clementi,^{5,6} Frank Noé,¹ Robert Pinsler,¹ and Rianne van den Berg¹

¹ Work done during an internship at Microsoft Research (Amsterdam).

² University of Copenhagen, Department of Computer Science, Universitetsparken 1, Copenhagen, 2100, Denmark.

³ AI4Science, Microsoft Research, Evert van de Beekstraat 354, Amsterdam, 1118 CZ, The Netherlands.

⁴ AI4Science, Microsoft Research, Karl-Liebknecht-Straße 32, Berlin, 10178, Germany.

⁵ University of Amsterdam, Information Institute, Science Park 904, Amsterdam, 1098 XH, The Netherlands.

⁶ Freie Universität Berlin, Department of Physics, Arnimallee 12, Berlin, 14195, Germany.

#AI4Science, Microsoft Research, 21 Station Road, Cambridge, CB1 2FB, United Kingdom.

^{1,2} Equal contribution.

E-mail: ma@di.ku.dk; victorgar@microsoft.com

Abstract

Course-grained (CG) molecular dynamics enables the study of biological processes at temporal and spatial scales that would be intractable at an atomistic resolution. However, accurately learning a CG force field remains a challenge. In this work, we leverage connections between score-based generative models, force fields and molecular

Swallowing the Bitter Pill: Simplified Scalable Conformer Generation

Yuyang Wang¹, Ahmed A. Elbag^{1,2}, Navdeep Jolly¹, Joshua M. Susskind¹, Miguel Ángel Bautista¹

Abstract

We present a novel way to predict molecular conformers through a simple formulation that sidesteps many of the heuristics of prior works and achieves state of the art results by using the advantages of scale. By training a diffusion generative model directly on 3D atomic positions without any constraints about the explicit structure of molecules (e.g., minimum bond angles) we are able to radically simplify structure generation and make it tractable for molecules up to 100 atoms.

is the vast complexity of the 3D structure space, encompassing factors such as bond lengths and torsional angles. Designing the molecular conformer space to satisfy these constraints, such as bond types and spatial arrangements determined by chiral centers, the conformational space experiences exponential growth with the expansion of the graph size and the number of rotatable bonds (Aszkenasy & Gomez-Bombarelli, 2022). This complicates brute force and exhaustive approaches, making them virtually unusable for even moderately small molecules. Systematic methods, like DMFGA (Hawkins et al., 2018),

Data augmentation

Article

Accurate structure prediction of biomolecular interactions with AlphaFold 3

<https://doi.org/10.1038/s43998-024-02487-w>
Received: 19 December 2023
Accepted: 29 April 2024
Published online: 8 May 2024
Open access
Check for updates

Josh Abrahams^{1,2}, Jonas Adler¹, Jack Danger^{1,3}, Richard Evans¹, Tim Green², Alexander Grütz¹, Olaf Henselberger¹, Lindsay Hiller¹, Andrew J. Bellard¹, Jennifer Johnson¹, Sebastian Kipperwinkel¹, David A. Lomax¹, Chaitanya Rangwala¹, Michael O’Neill¹, Daniel Reiter¹, Kathryn Raynor¹, Alexander Skarshewski¹, Zohreh Soltani¹, Nisha Srivastava¹, Estri Avanté¹, Charles Boettig¹, Ottavia Bartoli¹, Alex Bridgland¹, Alenay Chempur¹, Miles Congreve¹, Alexander L. Cowen-Rivers¹, Andrew Cowie¹, Michael Figari¹, Miles Gough¹, Daniel Hirschmann¹, Yannick J. Kihl¹, Catherine K. Kihl¹, Christopher M. R. Lee¹, Koko Lerner¹, Daniel Pogorelsky¹, Pouya Raveh¹, Sulabh Ray¹, Michael Riedel¹, Adrien Rossouw¹, Anikó Tóth-Lakatos¹, Catherine Tong¹, Sergei Vakser¹, Ellen S. Zhang¹, Michael Zeldes¹, Augustin Zidki¹, Victor Bapst¹, Pushmeet Kohli¹, Max Jaderberg^{1,2}, Dennis Hesselbarth^{1,2} & John M. Jumper^{1,2}

The introduction of AlphaFold 2 has spurred a revolution in modelling the structure of proteins and their interactions, enabling a huge range of applications in protein modelling and design^{1–3}. Here we describe the AlphaFold 3 model with a substantially updated architecture that increases its capacity to predict the complex structures of complexes including proteins, nucleic acids, small molecules, ions and modified residues. The new AlphaFold model demonstrates substantially improved accuracy

Probing the effects of broken symmetries in machine learning

Marc F. Langer¹, Sergey N. Prochnikov² and Michele Ceriotti¹

Laboratory of Computational Science and Modelling and National Centre for Computational Design and Discovery of Novel Materials MARVEL, Institute of Materials, Ecole Polytechnique Fédérale de Lausanne, 1015 Lausanne, Switzerland

¹ Author to whom any correspondence should be addressed.

E-mail: michele.ceriotti@epfl.ch

Keywords: machine learning, symmetry-constrained models, atomistic modeling, molecular simulations

Supplementary material for this article is available [online](#)

Abstract

Symmetry is one of the most central concepts in physics, and it is no surprise that it has also been widely adopted as an inductive bias for machine-learning models applied to the physical sciences. This is especially true for models targeting the properties of matter at the atomic scale. Both established and state-of-the-art approaches, with almost no exceptions, are built to exactly equivariant to translations, permutations, and rotations of the atoms. Incorporating symmetries—rotations in particular—constraints the model design space and implies more complicated architectures that are often also computationally demanding. There are indications

The Importance of Being Scalable: Improving the Speed and Accuracy of Neural Network Interatomic Potentials Across Chemical Domains

Eric Qu¹
UC Berkeley
ericqu@berkeley.edu

Aditi S. Krishnapriyan¹
UC Berkeley, LBNL
aditik1@berkeley.edu

Abstract

Scaling has been a critical factor in improving model performance and generalization across various fields of machine learning. It involves how a model’s performance changes with increases in model size or input data, as well as how efficiently computational resources are utilized to support this growth. Despite successes in scaling other types of machine learning models, the study of scaling in Neural Network Interatomic Potentials (NNIPs) remains limited. NNIPs act as

Two for One: Diffusion Models and Force Fields for Coarse-Grained Molecular Dynamics

Marioes Arts,^{1,1,2} Victor García Satorras,^{1,3,4} Chin-Wei Huang,¹ Daniel Zügner,⁵ Marco Federici,^{1,1} Cecilia Clementi,^{5,6} Frank Noé,¹ Robert Pinsler,¹ and Rianne van den Berg¹

¹ Work done during an internship at Microsoft Research (Amsterdam).
² University of Copenhagen, Department of Computer Science, Universitetsparken 1, Copenhagen, 2100, Denmark.

³ AI4Science, Microsoft Research, Evert van de Beekstraat 354, Amsterdam, 1118 CZ, The Netherlands.
⁴ AI4Science, Microsoft Research, Karl-Liebknecht-Straße 32, Berlin, 10178, Germany.

⁵ University of Amsterdam, Information Institute, Science Park 904, Amsterdam, 1098 XH, The Netherlands.

⁶ Freie Universität Berlin, Department of Physics, Arnimallee 12, Berlin, 14195, Germany.
⁷ AI4Science, Microsoft Research, 21 Station Road, Cambridge, CB1 2FB, United Kingdom.

¹ Equal contribution.

* E-mail: ma@di.ku.dk; victorgar@microsoft.com

Abstract

Coarse-grained (CG) molecular dynamics enables the study of biological processes at temporal and spatial scales that would be intractable at an atomistic resolution. However, accurately learning a CG force field remains a challenge. In this work, we leverage connections between score-based generative models, force fields and molecular

Swallowing the Bitter Pill: Simplified Scalable Conformer Generation

Yuyang Wang¹, Ahmed A. Elbag^{1,2}, Navdeep Jolly¹, Joshua M. Susskind¹, Miguel Ángel Bautista¹

Abstract

We present a novel way to predict molecular conformers through a simple formulation that sidesteps many of the heuristics of prior works and achieves state of the art results by using the advantages of scale. By training a diffusion generative model directly on 3D atomic positions without any constraints about the chemical structure of molecules (or even the bond and angles) we are able to radically simplify structure generation and make it much more efficient.

is the vast complexity of the 3D structure space, encompassing factors such as bond lengths and torsional angles. Despite the molecular specific constraints, such as bond types and spatial arrangements determined by chiral centers, the conformational space experiences exponential growth with the expansion of the graph size and the number of rotatable bonds (Aszled & Gomez-Bombarelli, 2022). This complicates brute force and exhaustive approaches, making them virtually unfeasible for even moderately small molecules. Systematic methods, like DMFGA (Hawkins et al., 2018),

DOES EQUIVARIANCE MATTER AT SCALE?

Johann Bremer¹, Sönke Behrends¹, Pim de Haan², Taco Cohen¹
Quacquarelli AI Research¹
pim1@johannbremer.de

ABSTRACT

Given large data sets and sufficient compute, is it beneficial to design neural architectures for the structure and symmetries of each problem? Or is it more efficient to learn them from data? We study empirically how equivariant and non-equivariant networks scale with compute and training samples. Focusing on a benchmark problem of rigid-body interactions and on general-purpose transformer architectures, we perform a series of experiments, varying the model size, training steps, and dataset size. We find that non-equivariant models with data augmentation are less efficient, but training non-equivariant models with data augmentation can close this gap given sufficient epochs. Second, scaling with compute follows a power law, with equivariant models outperforming non-equivariant ones at each tested compute budget. Finally, the optimal allocation of a compute budget onto model size and training duration differs between equivariant and non-equivariant models.

Data augmentation

- thumb-up Easy to implement
- thumb-up 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

loss

training sample

The diagram illustrates the components of the training dynamics equation. It shows a blue arrow pointing from the text 'learning rate' to the term $-\frac{\eta}{N}$. Another blue arrow points from the text 'loss' to the term $\frac{\partial L}{\partial \mathcal{N}(x_i)}$. A third blue arrow points from the text 'training sample' to the summation term $\sum_{i=1}^N \Theta_\theta(x, x_i)$.

Empirical NTK

Training dynamics under continuous gradient descent:

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

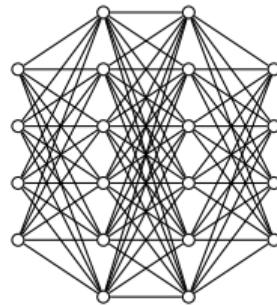
↑
learning rate ↑
↑
training sample ↑
loss

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

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

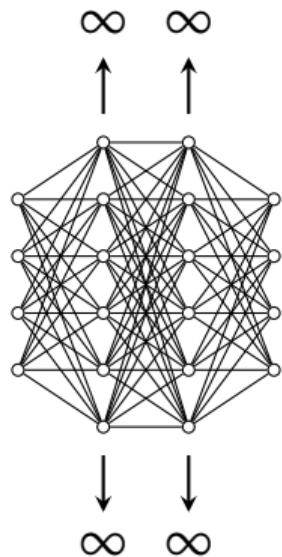
Infinite width limit

[Jacot et al. 2018]



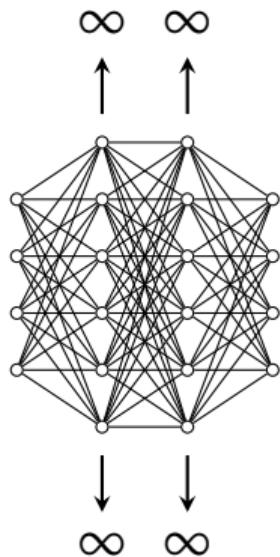
Infinite width limit

[Jacot et al. 2018]



Infinite width limit

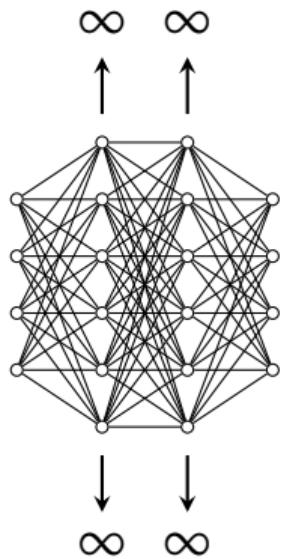
[Jacot et al. 2018]



👍 NTK becomes independent of initialization

Infinite width limit

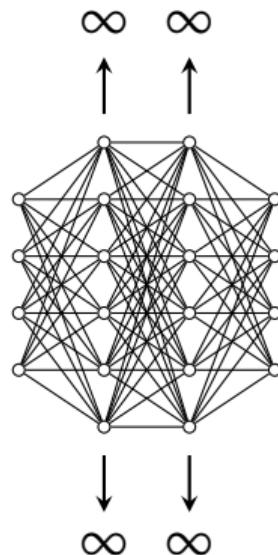
[Jacot et al. 2018]



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

Infinite width limit

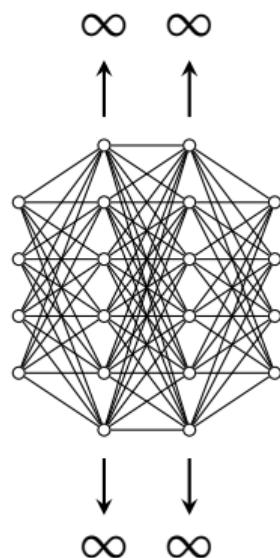
[Jacot et al. 2018]



- 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]

- ① At infinite width, the mean prediction is given by

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

Mean prediction from NTK

[Jacot et al. 2018]

- ① At infinite width, the mean prediction is given by

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

neural tangent kernel



Mean prediction from NTK

[Jacot et al. 2018]

- ① At infinite width, the mean prediction is given by

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

neural tangent kernel

train data

Mean prediction from NTK

[Jacot et al. 2018]

- ① At infinite width, the mean prediction is given by

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

neural tangent kernel

learning rate

train data

Mean prediction from NTK

[Jacot et al. 2018]

- ① At infinite width, the mean prediction is given by

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

Diagram illustrating the components of the mean prediction formula:

- neural tangent kernel**: Points to the term $\Theta(x, X)$.
- train labels**: Points to the term Y .
- learning rate**: Points to the term $e^{-\eta\Theta(X, X)t}$.
- train data**: Points to the term $\Theta(X, X)^{-1}$.

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

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

The diagram shows the mathematical expression for augmented data $\mu_t(x)$. It consists of several blue arrows pointing from labels to terms in the equation:

- An arrow points from the label "augmented data" to the term $\Theta(x, X)$.
- Two arrows point from the label "augmented data" to the term $\Theta(X, X)^{-1}$.
- Two arrows point from the label "augmented labels" to the term $(\mathbb{I} - e^{-\eta\Theta(X, X)t})Y$.
- A single arrow points from the label "augmented labels" to the label "augmented data".

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

The diagram illustrates the mathematical expression for data augmentation. A blue curved arrow labeled "group transformation" points downwards to the equation. Another blue curved arrow labeled "augmented data" points upwards from the left side of the equation to the term $\rho(g)x$. A third blue curved arrow labeled "augmented labels" points upwards from the right side of the equation to the term Y .

Data augmentation at infinite width

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

group transformation

for augmented data

augmented data

augmented labels

The diagram illustrates the components of the data augmentation formula. At the top, a blue curved arrow labeled "group transformation" points from the left towards the term $\Theta(\rho(g)x, X)$. Another blue curved arrow labeled "for augmented data" points from the right towards the term $\mathbb{I} - e^{-\eta}\Theta(X, X)t$. Below the equation, two blue arrows point upwards from the labels "augmented data" and "augmented labels" to the terms $\Theta(\rho(g)x, X)$ and $\Theta(X, X)^{-1}$ respectively.

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 group transformation equation. A blue arrow points from the term $\rho(g)$ in the equation to the label "group transformation" above it. Another blue arrow points from the term $\Theta(X, X)^{-1}$ to the label "augmented data" below it. A third blue arrow points from the term $e^{-\eta\Theta(X, X)t}$ to the label "augmented labels" below it.

Data augmentation at infinite width

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

group transformation

augmented labels

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

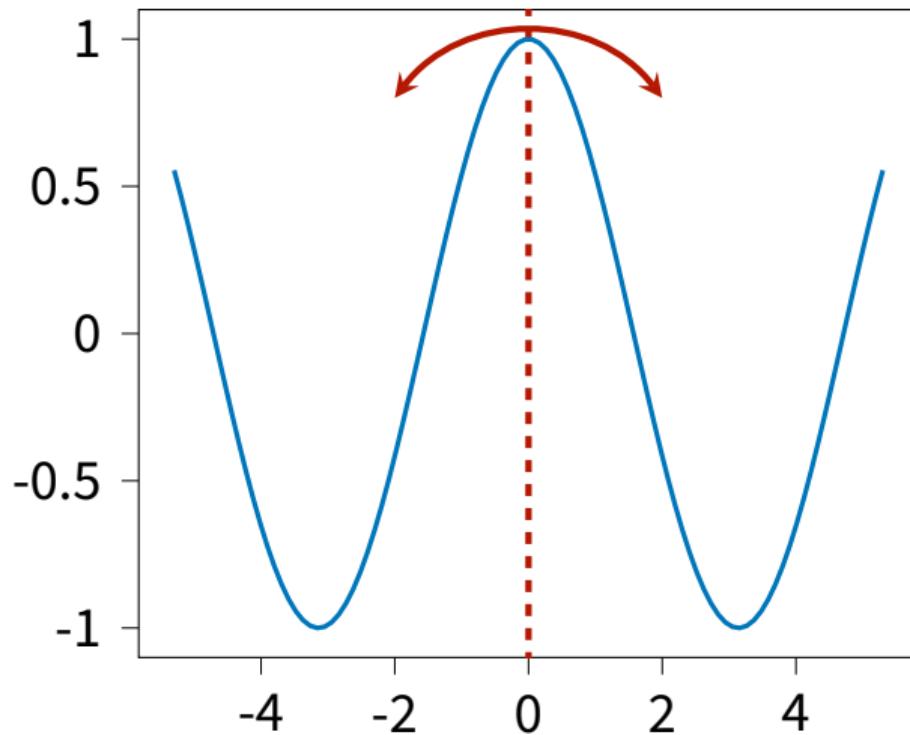
Intuitive explanation

- ✓ Equivariance holds for all training times
 - ✓ Equivariance holds away from the training data
-
- ➊ At infinite width, the mean output at initialization is zero everywhere.

Intuitive explanation

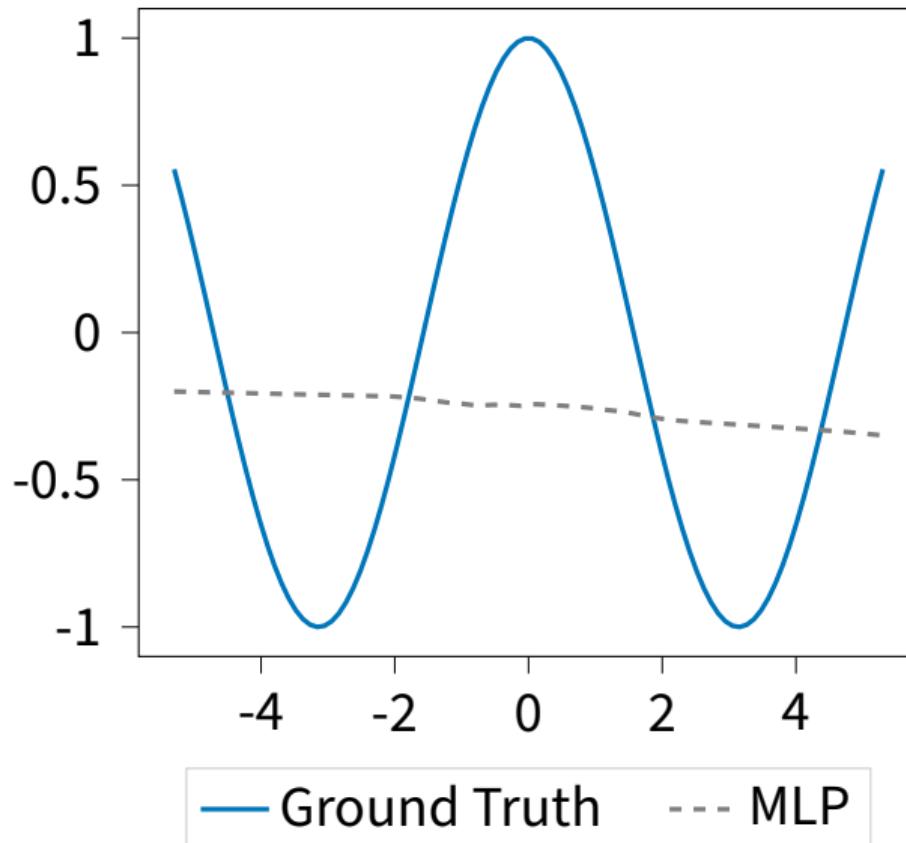
- ✓ Equivariance holds for all training times
 - ✓ Equivariance holds away from the training data
-
- ➊ At infinite width, the mean output at initialization is zero everywhere.
 - ⇒ Training with full data augmentation leads to an equivariant function.

Toy example

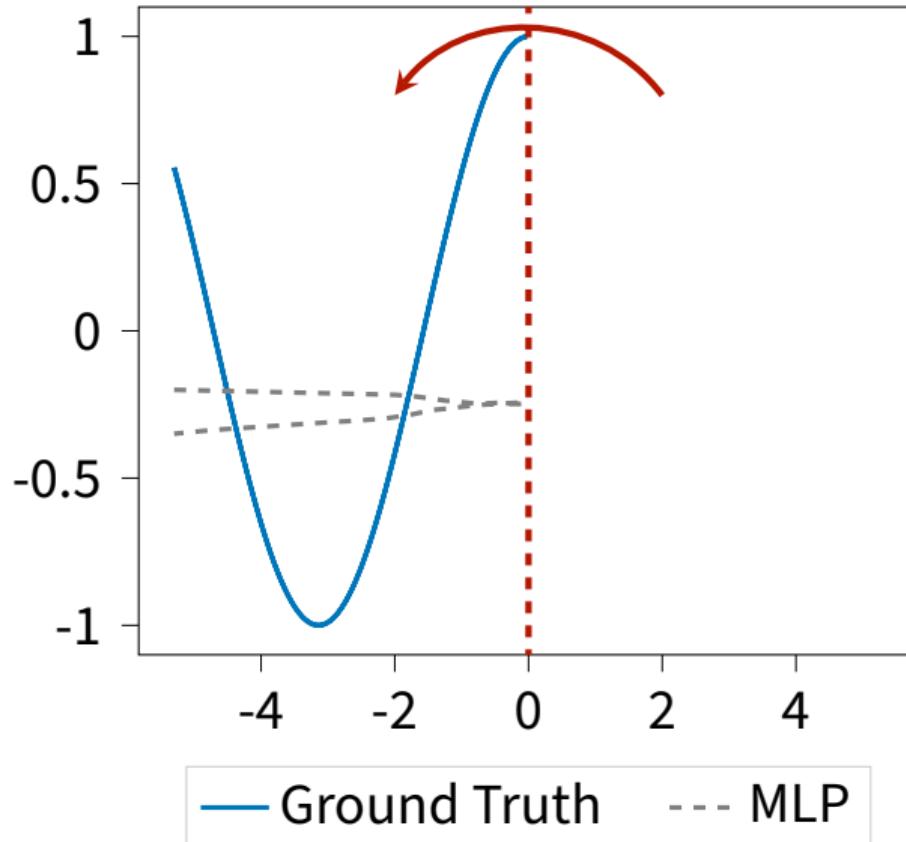


— Ground Truth

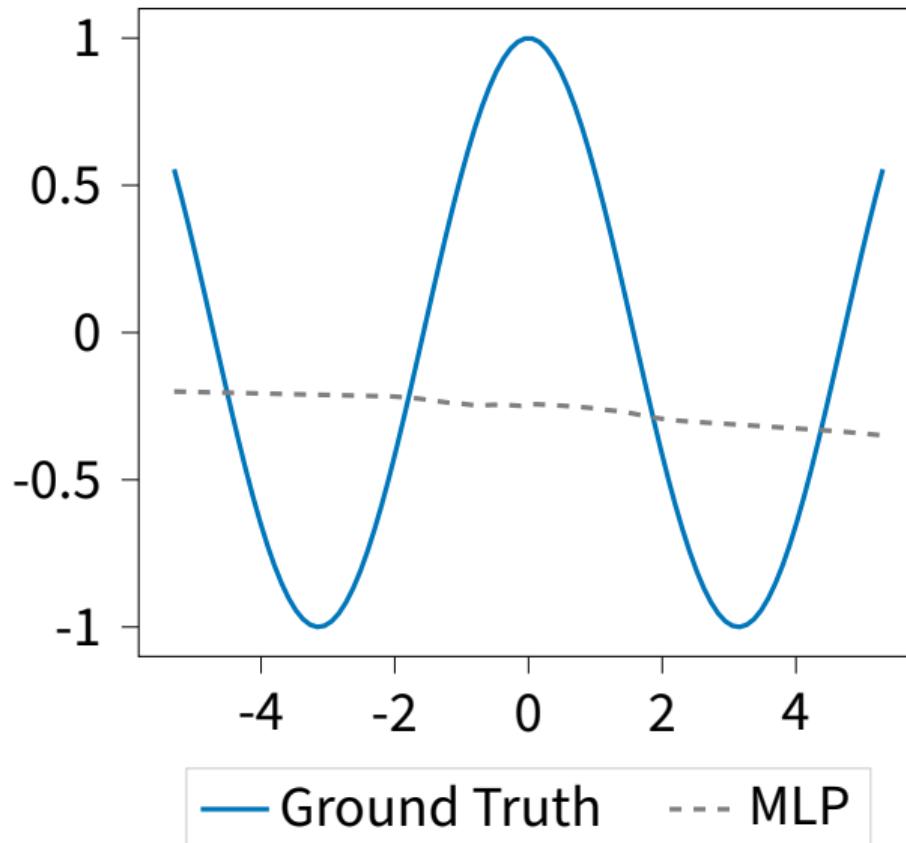
Initialization



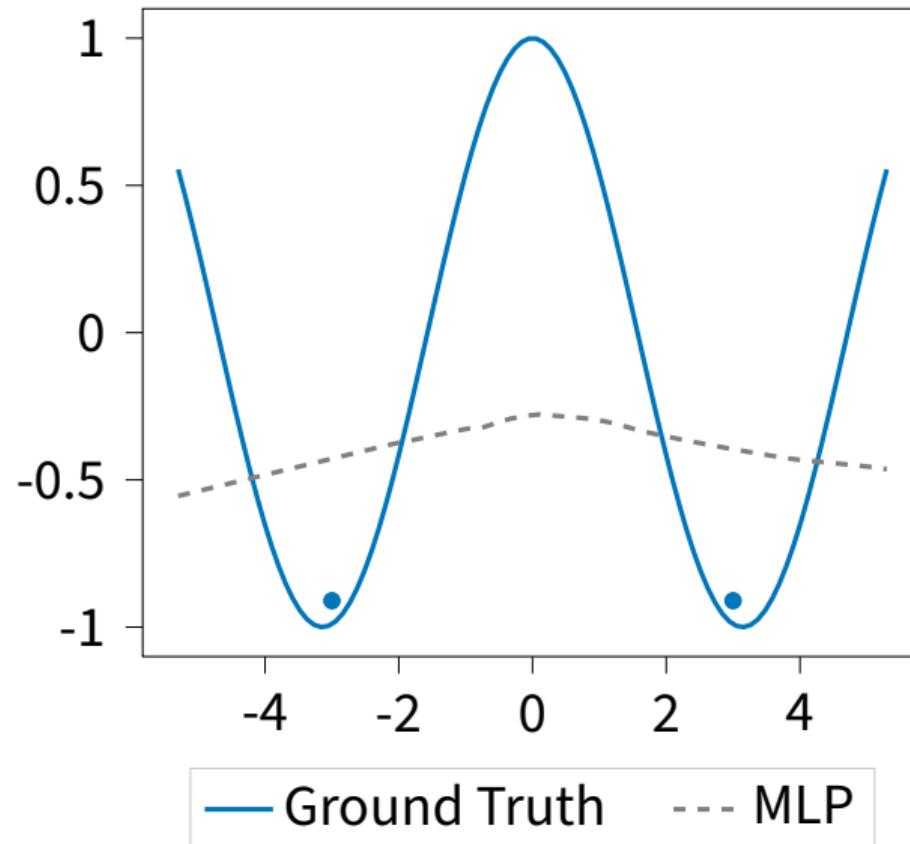
Initialization



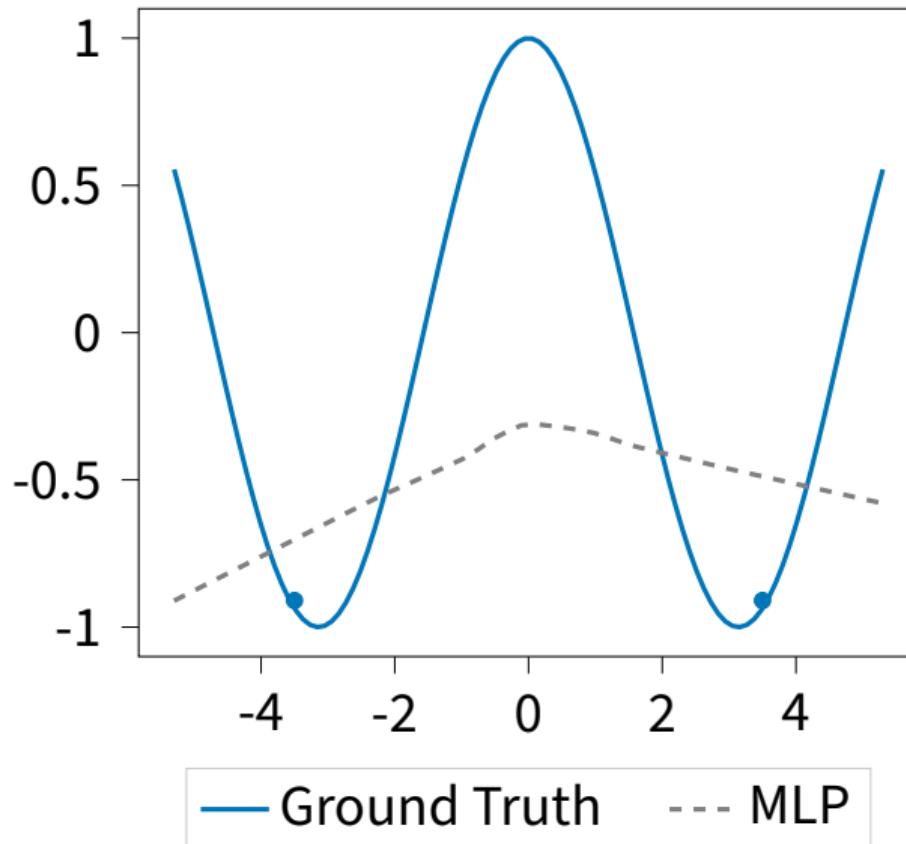
Initialization



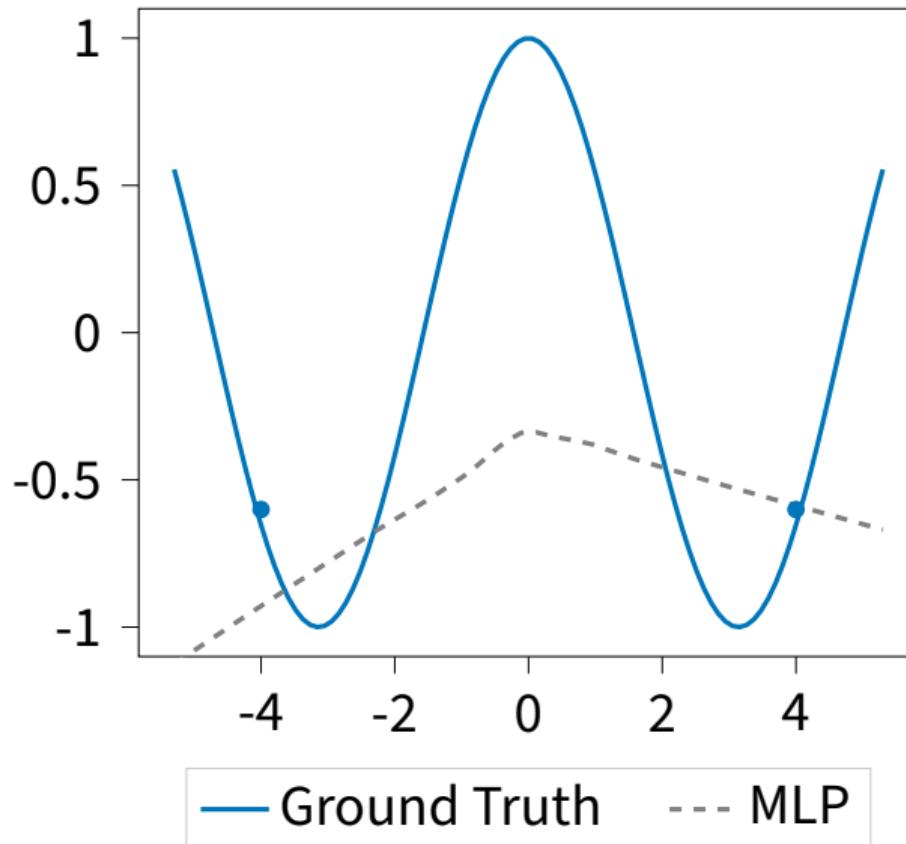
After 1 Training Step



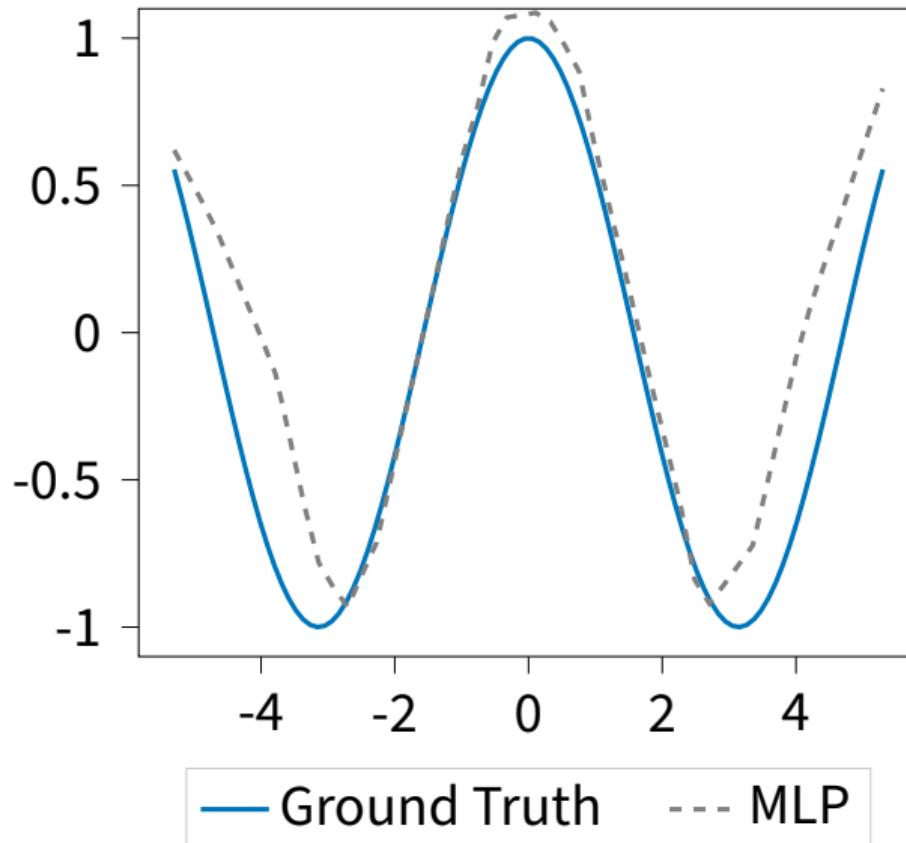
After 2 Training Steps



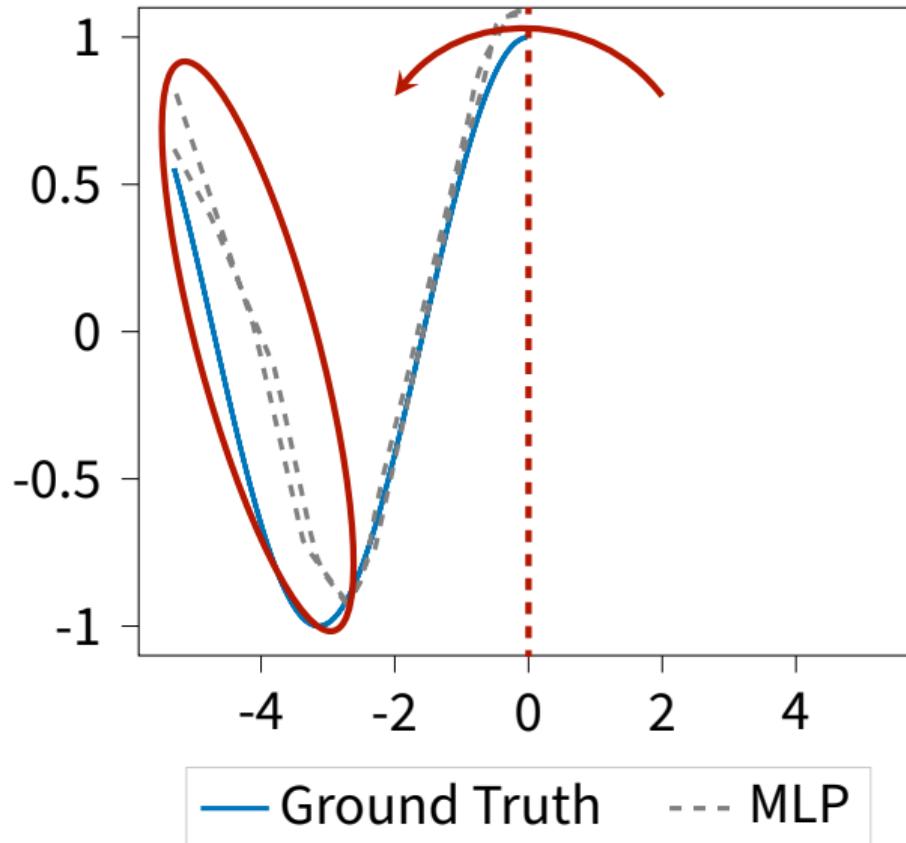
After 3 Training Steps



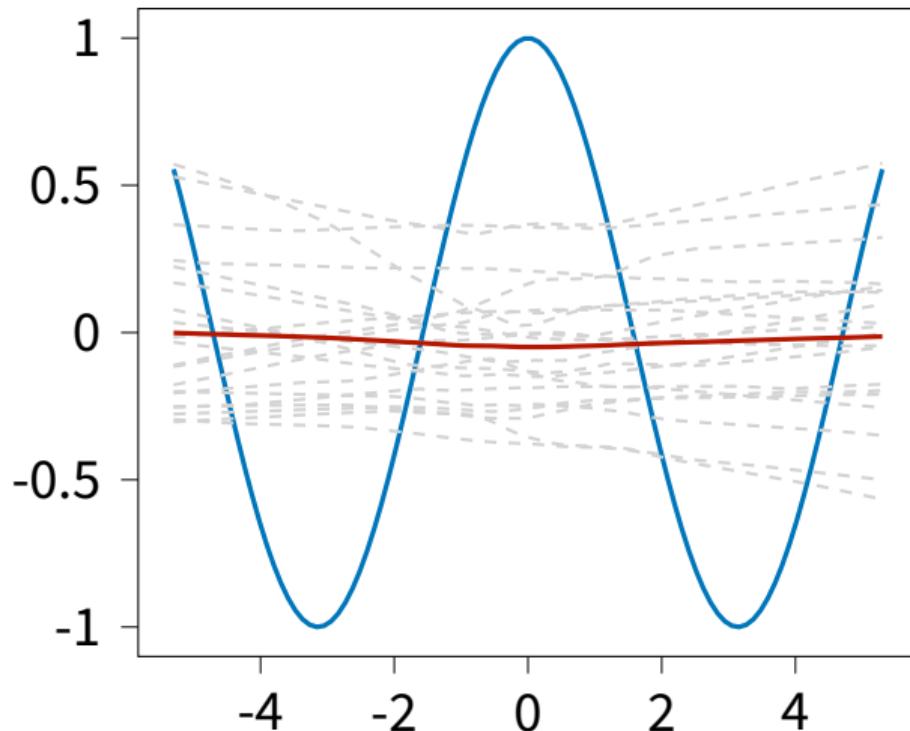
After 2000 Training Steps



After 2000 Training Steps

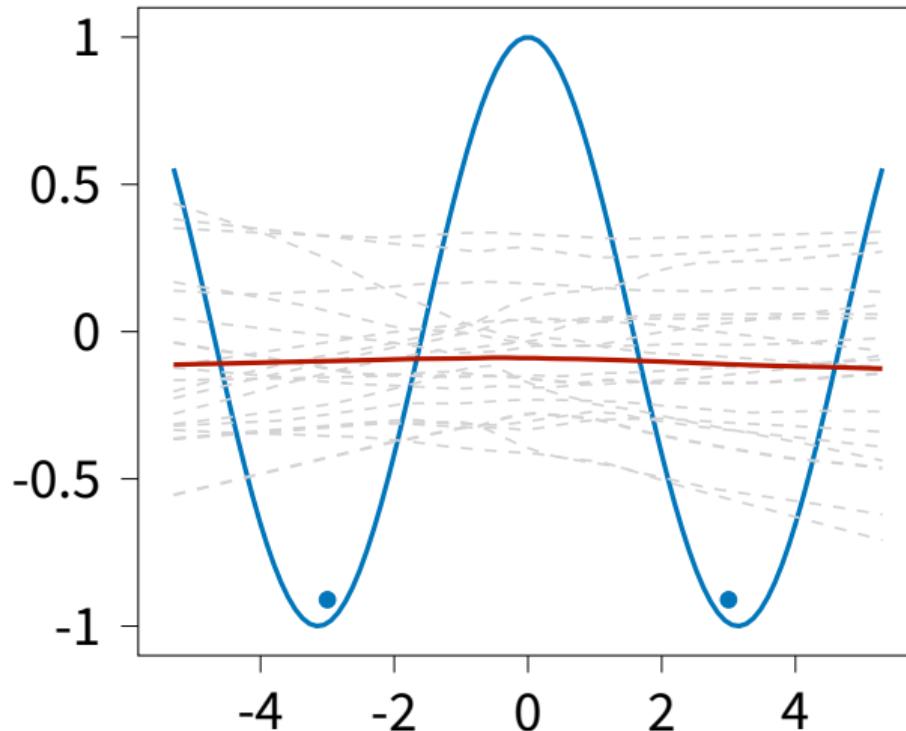


Initialization



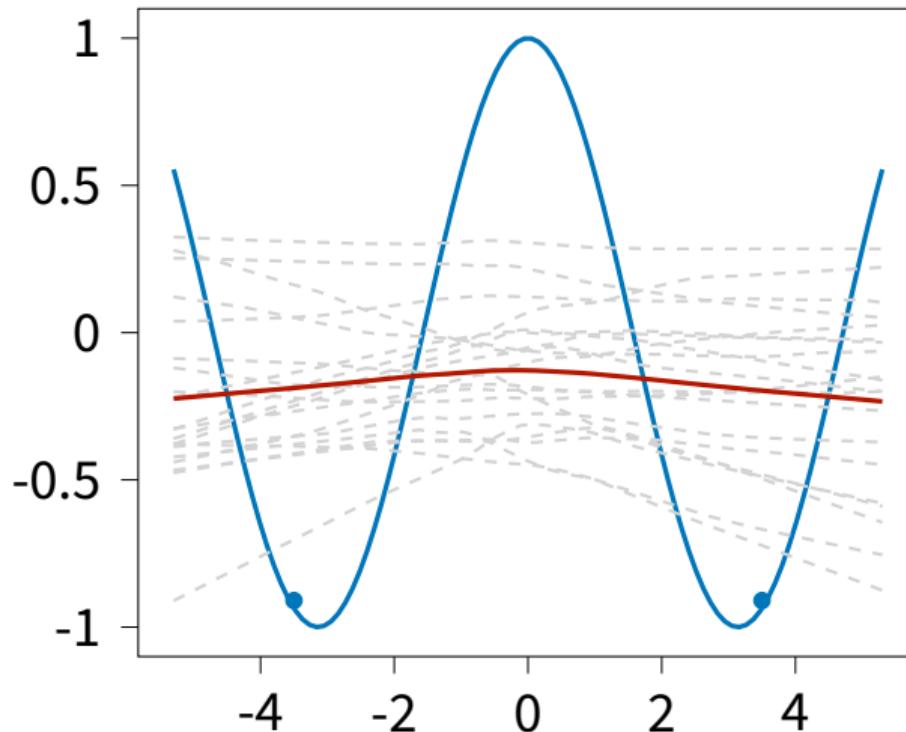
— Ground Truth - - - MLP — Ensemble Mean

After 1 Training Step



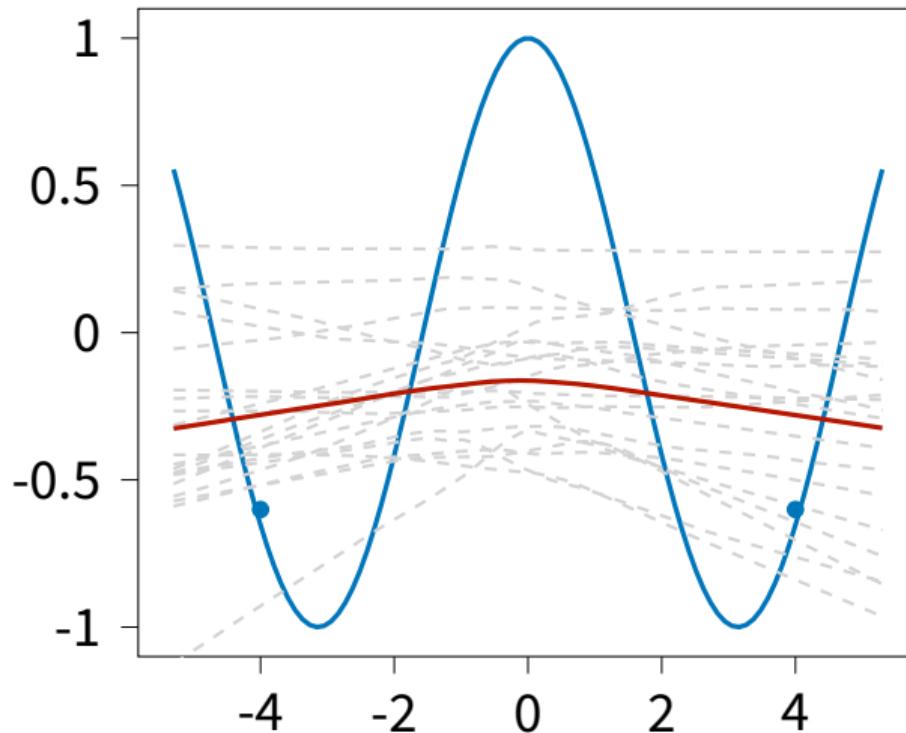
— Ground Truth - - - MLP — Ensemble Mean

After 2 Training Steps



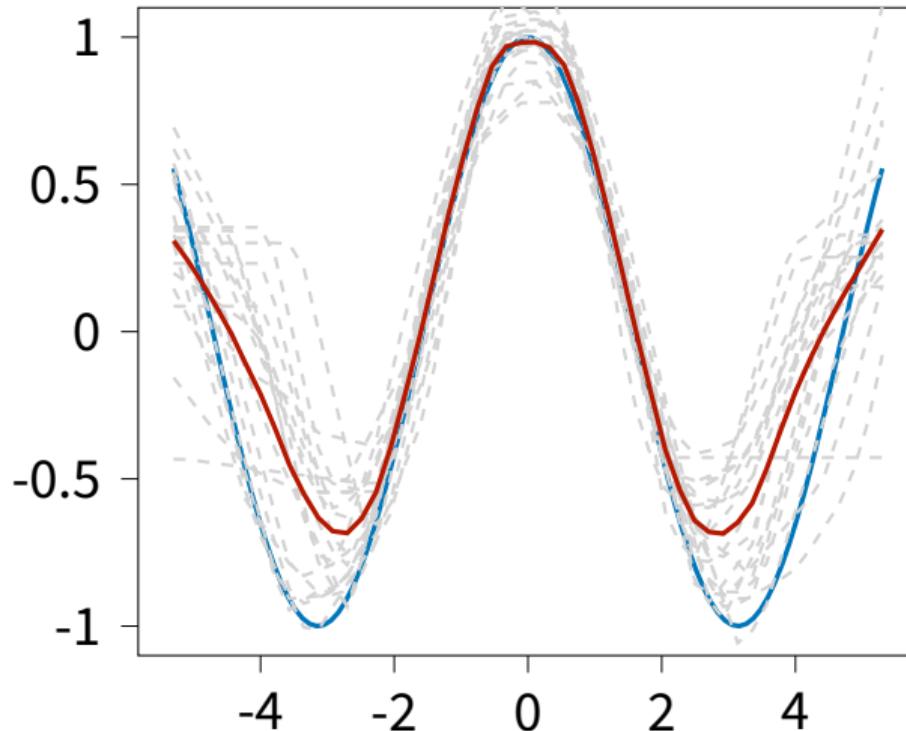
— Ground Truth - - - MLP — Ensemble Mean

After 3 Training Steps



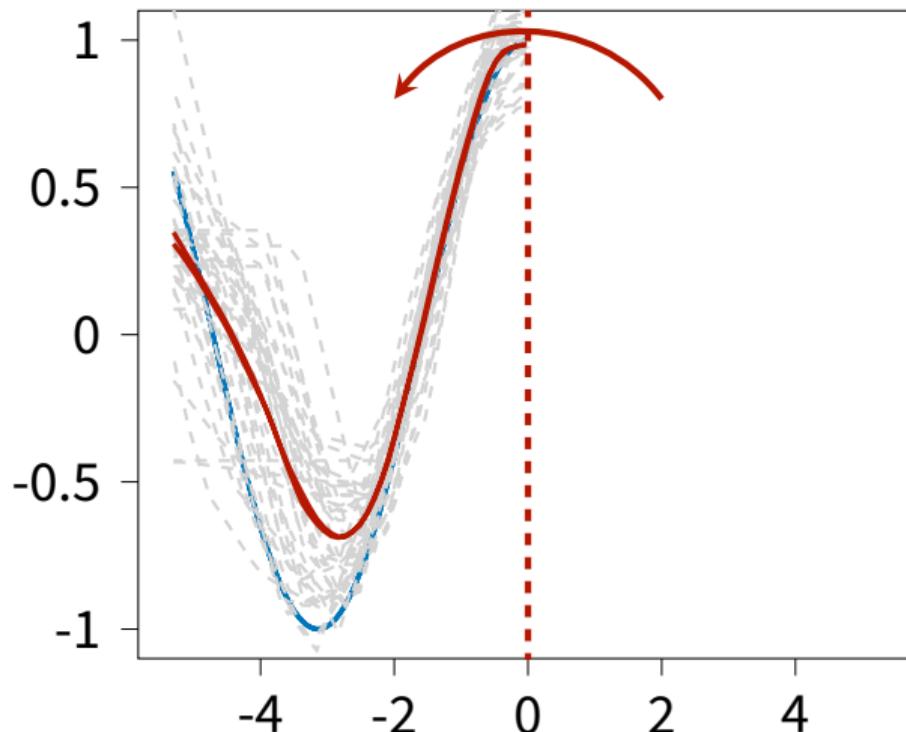
— Ground Truth - - - MLP — Ensemble Mean

After 2000 Training Steps



— Ground Truth - - - MLP — Ensemble Mean

After 2000 Training Steps

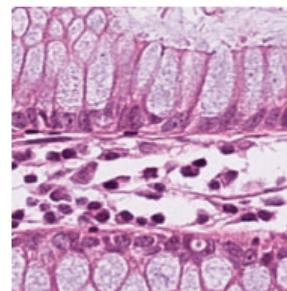
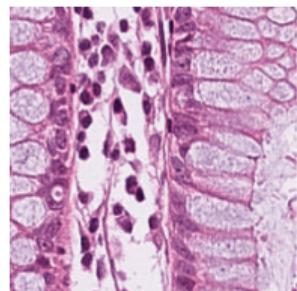


— Ground Truth - - - MLP — Ensemble Mean

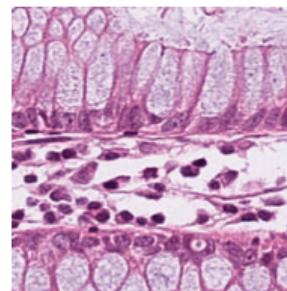
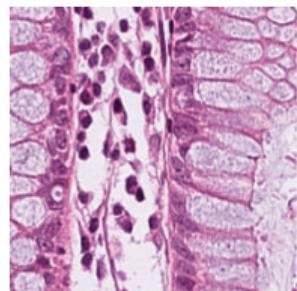
What Does An Augmented Ensemble Converge To?

Rotating images

Rotating images



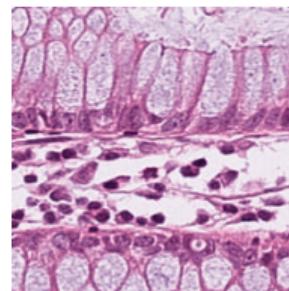
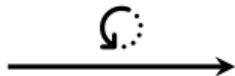
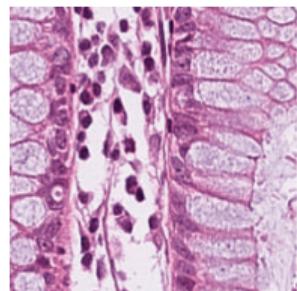
Rotating images



$$f(x)$$

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

Rotating images



$$\begin{aligned} f(x) \\ f : \text{pixels} \rightarrow \text{colors} \end{aligned}$$



$$\begin{aligned} f(\rho(g^{-1})x) \\ = [\rho_{\text{reg}}(g)f](x) \end{aligned}$$

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

$$\Theta^{\text{non-aug}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{aug}}(f, \rho_{\text{reg}}(g)f')$$

Data augmentation and NTKs

Consider two ensembles:

trained without data augmentation

trained with data augmentation

If

$$\Theta^{\text{non-aug}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{aug}}(f, \rho_{\text{reg}}(g)f')$$

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

$$\Theta^{\text{non-aug}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{aug}}(f, \rho_{\text{reg}}(g)f')$$

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

$$\Theta^{\text{non-aug}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{aug}}(f, \rho_{\text{reg}}(g)f')$$

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

$$\Theta^{\text{non-aug}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{aug}}(f, \rho_{\text{reg}}(g)f')$$

Data augmentation and NTKs

$$\Theta^{\text{non-aug}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{aug}}(f, \rho_{\text{reg}}(g)f')$$

- ① Given an architecture with NTK Θ^{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 ρ_{reg}

Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt ρ_{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 ρ_{reg}

- Ordinary convolutions

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

- Group convolutions

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

Group convolutions

[Cohen, Welling 2016]

Group conv's are the (unique) linear layers equivariant wrt ρ_{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 ρ_{reg}

- Ordinary convolutions

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

- Group convolutions

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

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

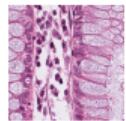
$$f' = \frac{1}{\text{vol}(G)} \int_G dg f(g) \quad \text{group pooling}$$

GCNNs

Stack GConv-layers to obtain an invariant network

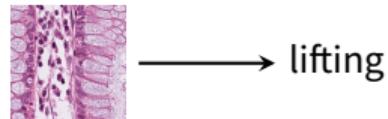
GCNNs

Stack GConv-layers to obtain an invariant network



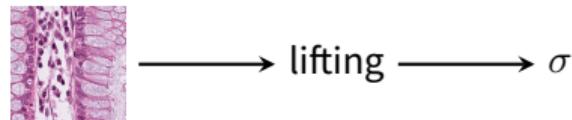
GCNNs

Stack GConv-layers to obtain an invariant network



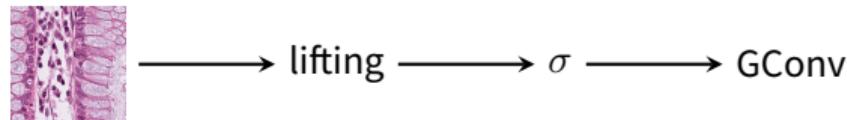
GCNNs

Stack GConv-layers to obtain an invariant network



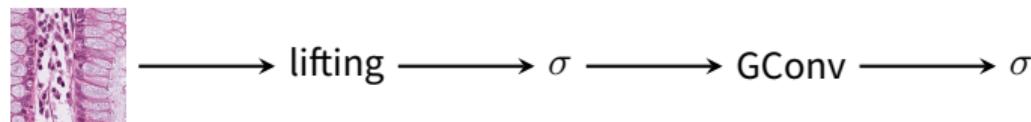
GCNNs

Stack GConv-layers to obtain an invariant network



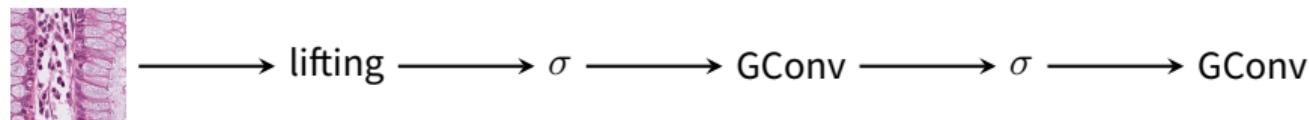
GCNNs

Stack GConv-layers to obtain an invariant network



GCNNs

Stack GConv-layers to obtain an invariant network



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

[Misof, Kessel, JG, 2024]

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

0

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

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

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

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

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

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

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

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

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

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

NTKs for GCNNs

Stack GConv-layers to obtain an invariant network



Compute NTK with layer-wise recursion

[Misof, Kessel, JG, 2024]

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

NTKs of MLPs and GCNNs

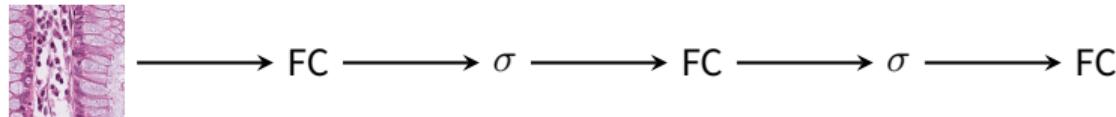
NTKs of MLPs and GCNNs

- Consider two neural networks

NTKs of MLPs and GCNNs

- Consider two neural networks

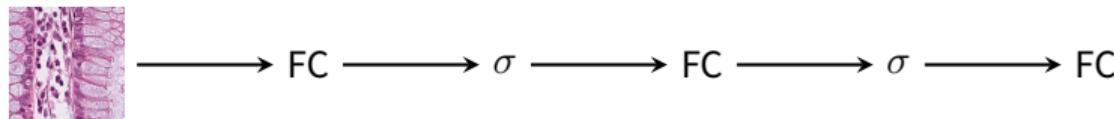
An MLP



NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



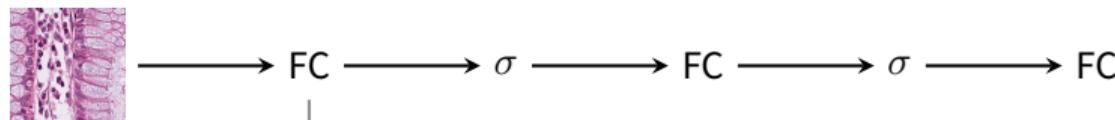
A GCNN



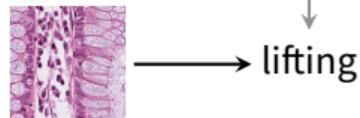
NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



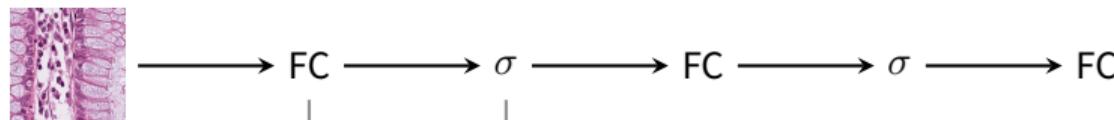
A GCNN



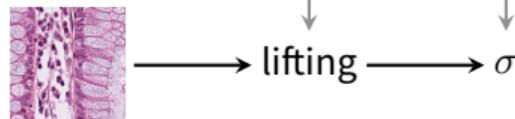
NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



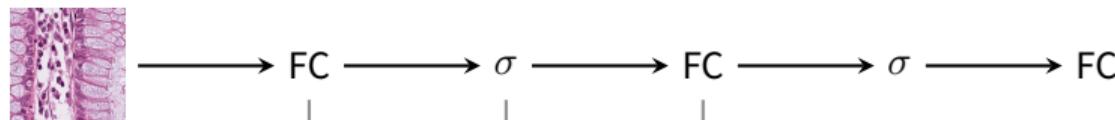
A GCNN



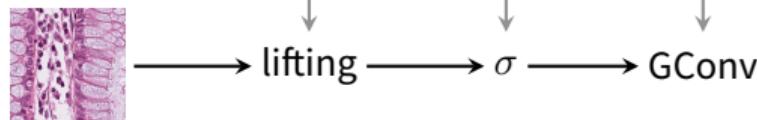
NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



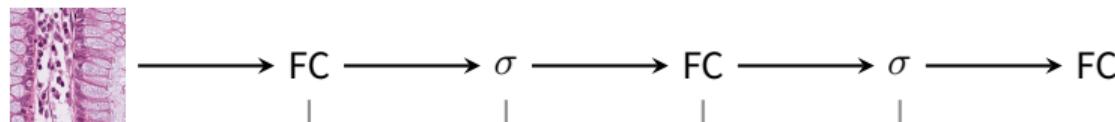
A GCNN



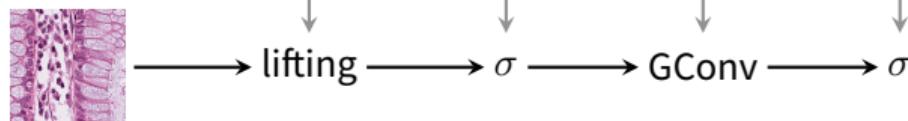
NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



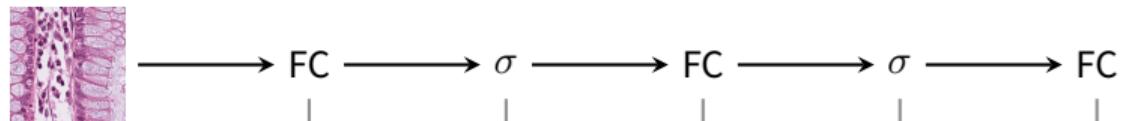
A GCNN



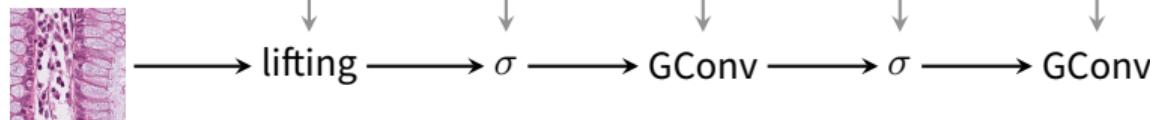
NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



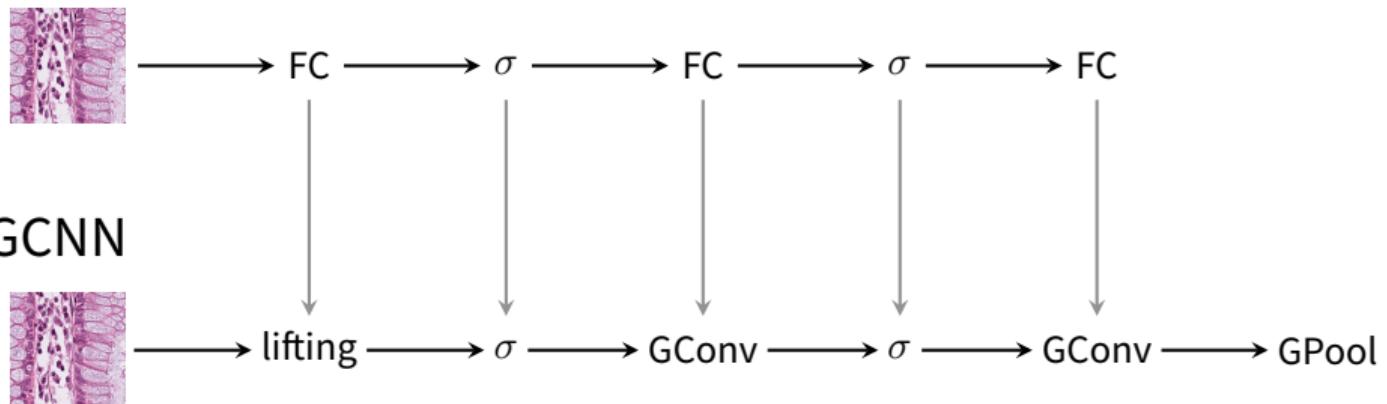
A GCNN



NTKs of MLPs and GCNNs

- Consider two neural networks

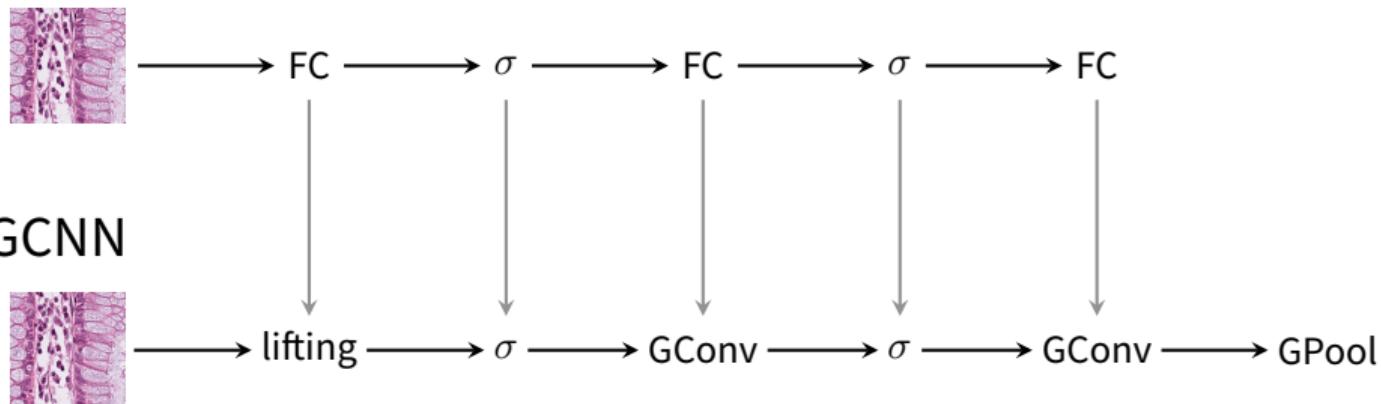
An MLP



NTKs of MLPs and GCNNs

- Consider two neural networks

An MLP



- Then

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

Data augmentation of MLPs

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

Data augmentation of MLPs

before: non-aug

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

Data augmentation of MLPs

before: non-aug

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

before: aug

Data augmentation of MLPs

before: non-aug

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

before: aug

- ⇒ training the MLP on
G-augmented data

Data augmentation of MLPs

before: non-aug

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

before: aug

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

Data augmentation of MLPs

before: non-aug

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

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

$$\Theta^{\text{GCNN}}(f, f') = \frac{1}{|G|} \sum_{g \in G} \Theta^{\text{MLP}}(f, \rho_{\text{reg}}(g)f')$$

before: aug

⇒ training the MLP on
G-augmented data

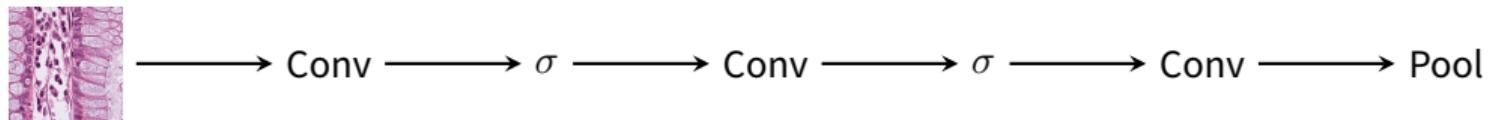
= training the GCNN on
unaugmented data

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

Data augmentation of CNNs

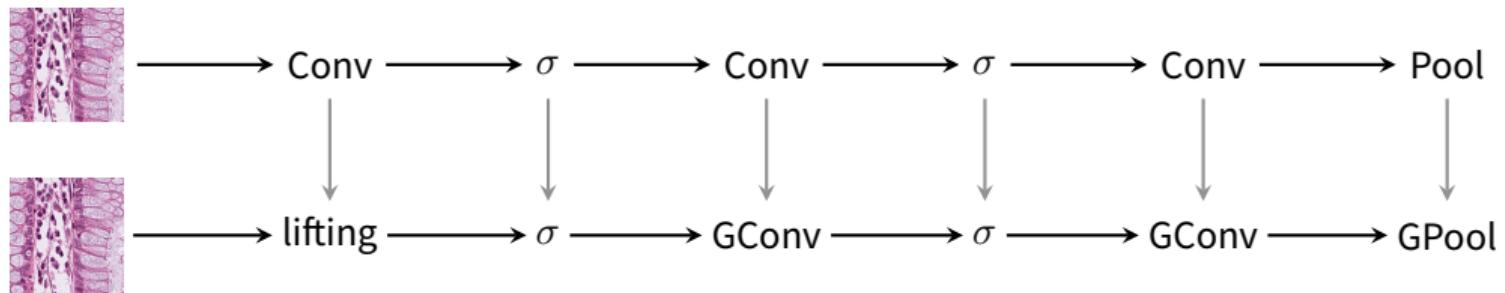
Data augmentation of CNNs

- Consider a CNN



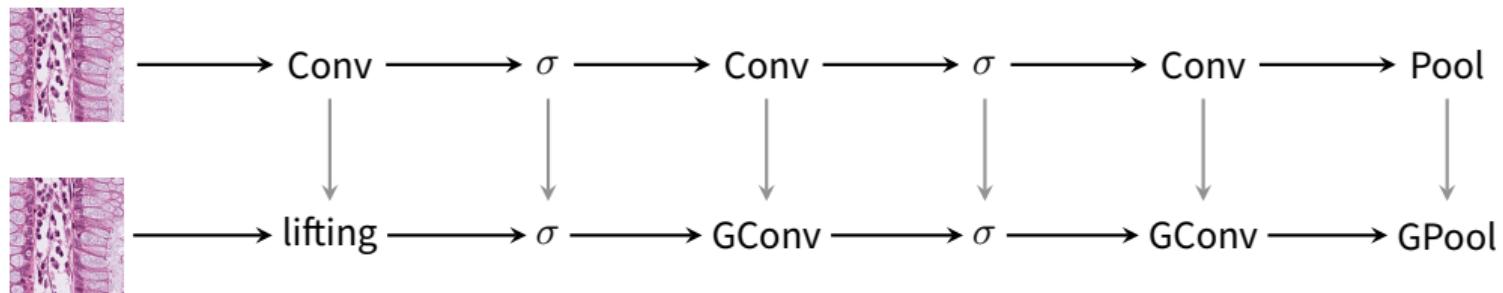
Data augmentation of CNNs

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



Data augmentation of CNNs

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

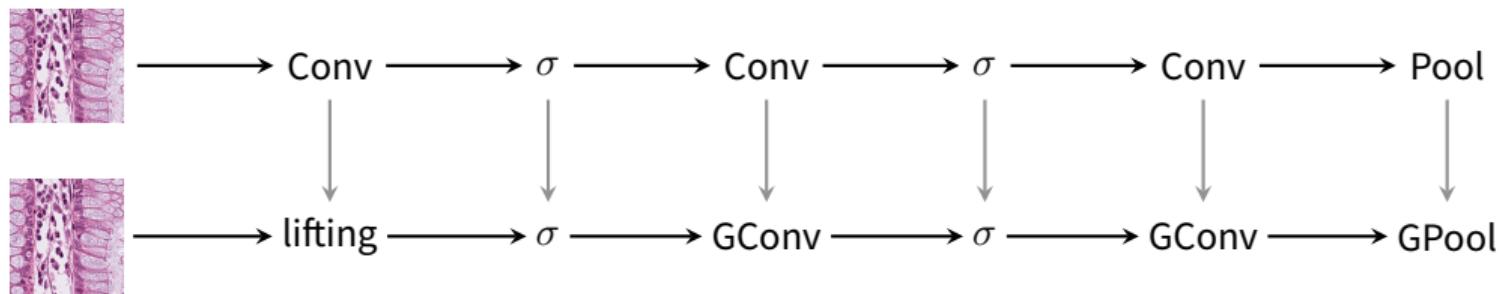


- Then

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

Data augmentation of CNNs

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



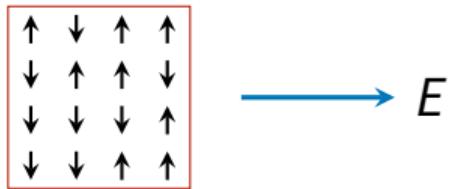
- Then

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

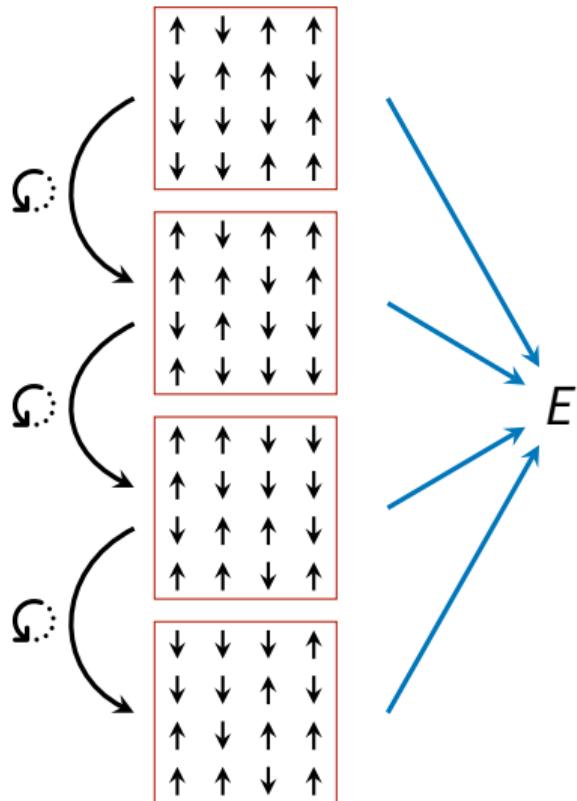
- By training the CNN on rotated images, one obtains a roto-translation invariant GCNN

Experiments

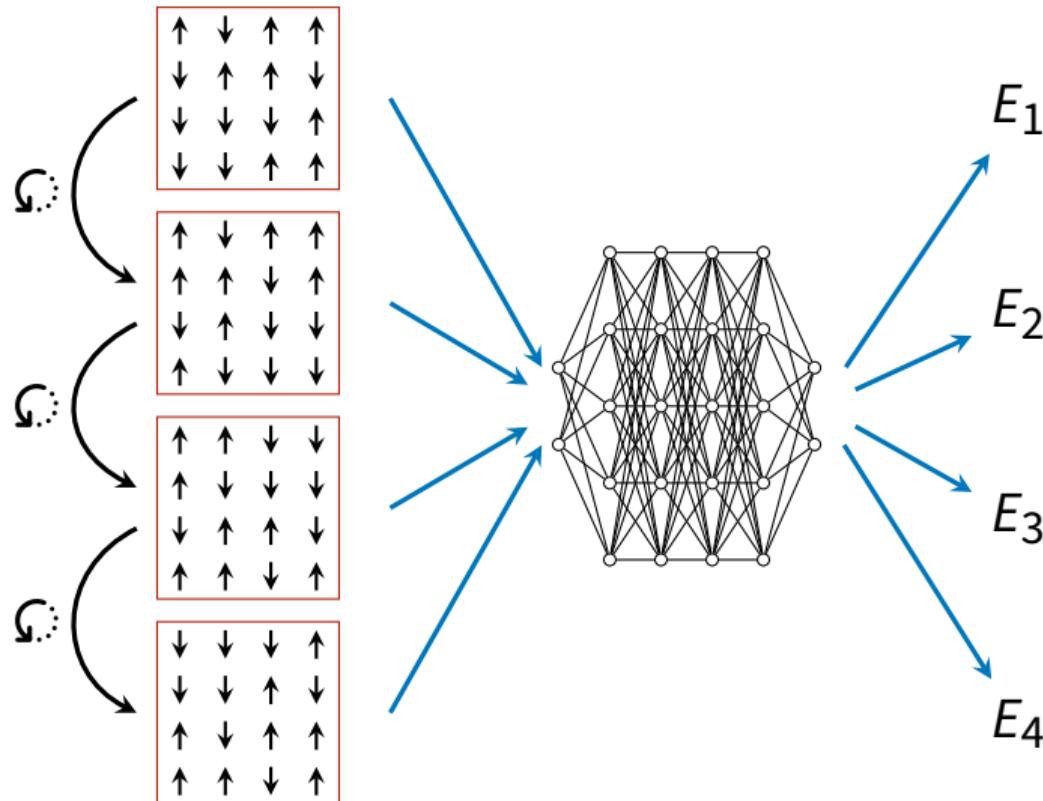
Ising model



Ising model

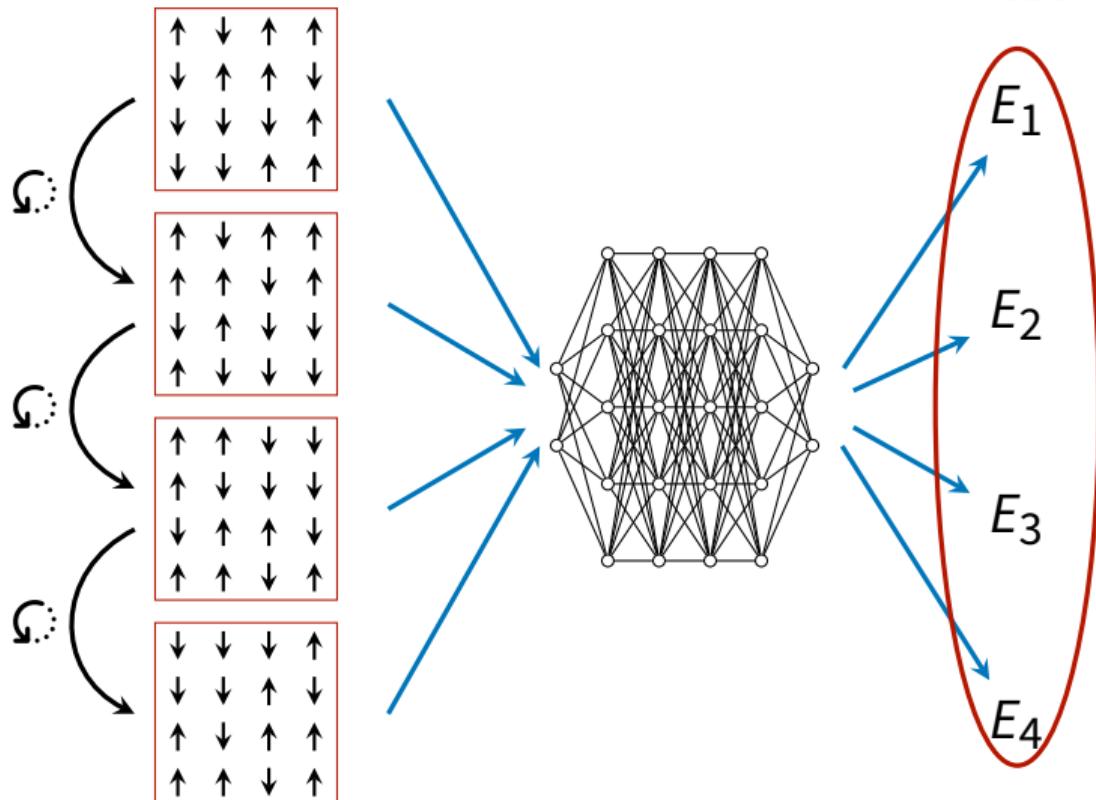


Ising model

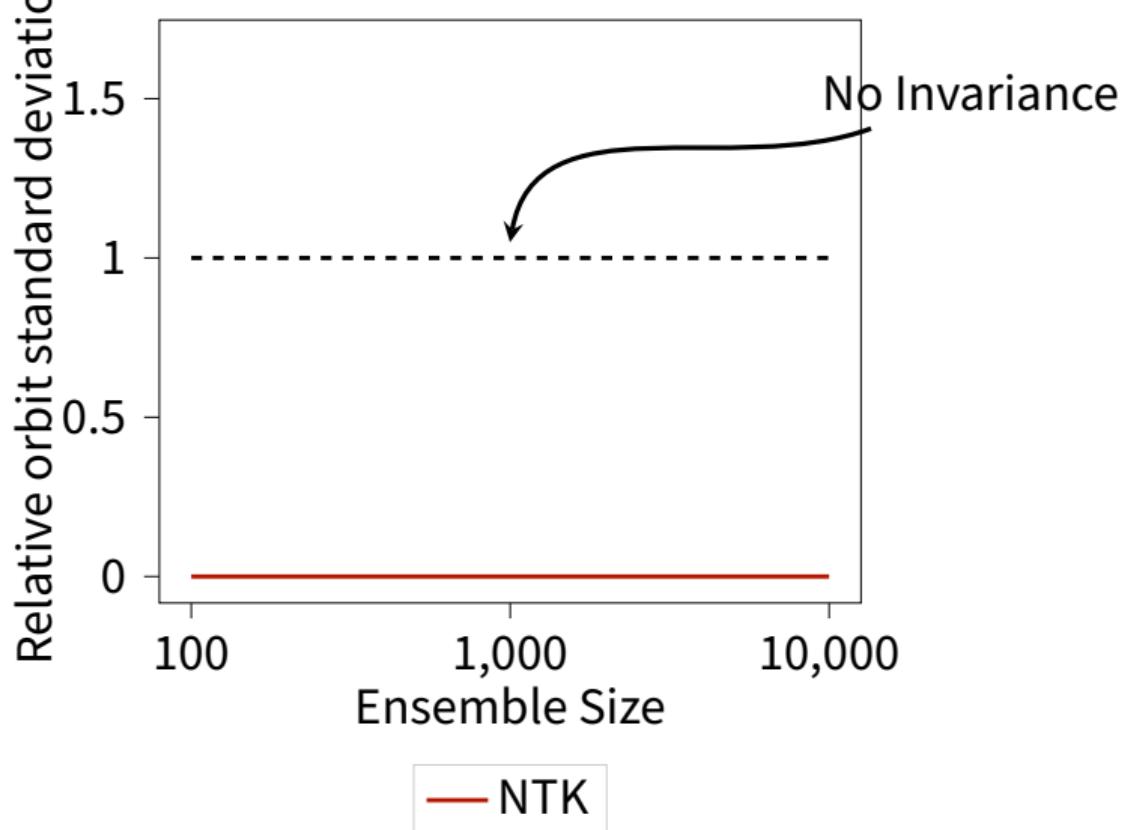


Ising model

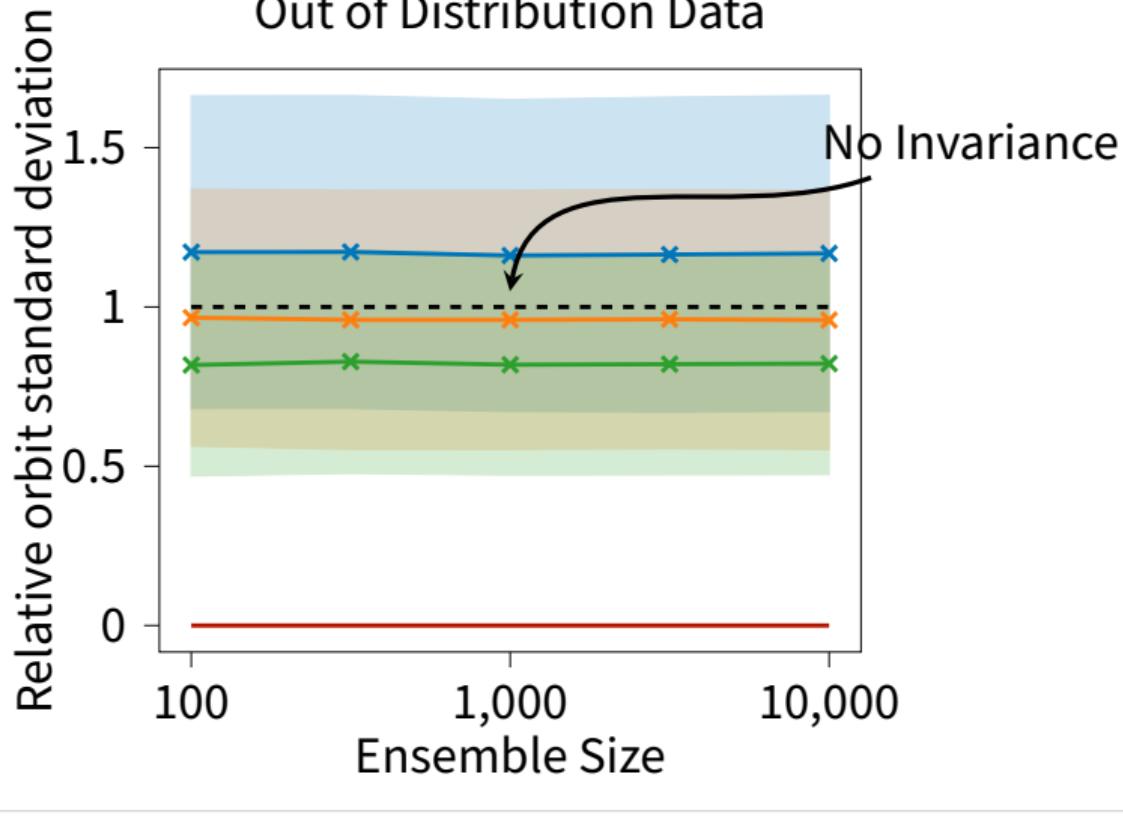
Relative Standard Deviation



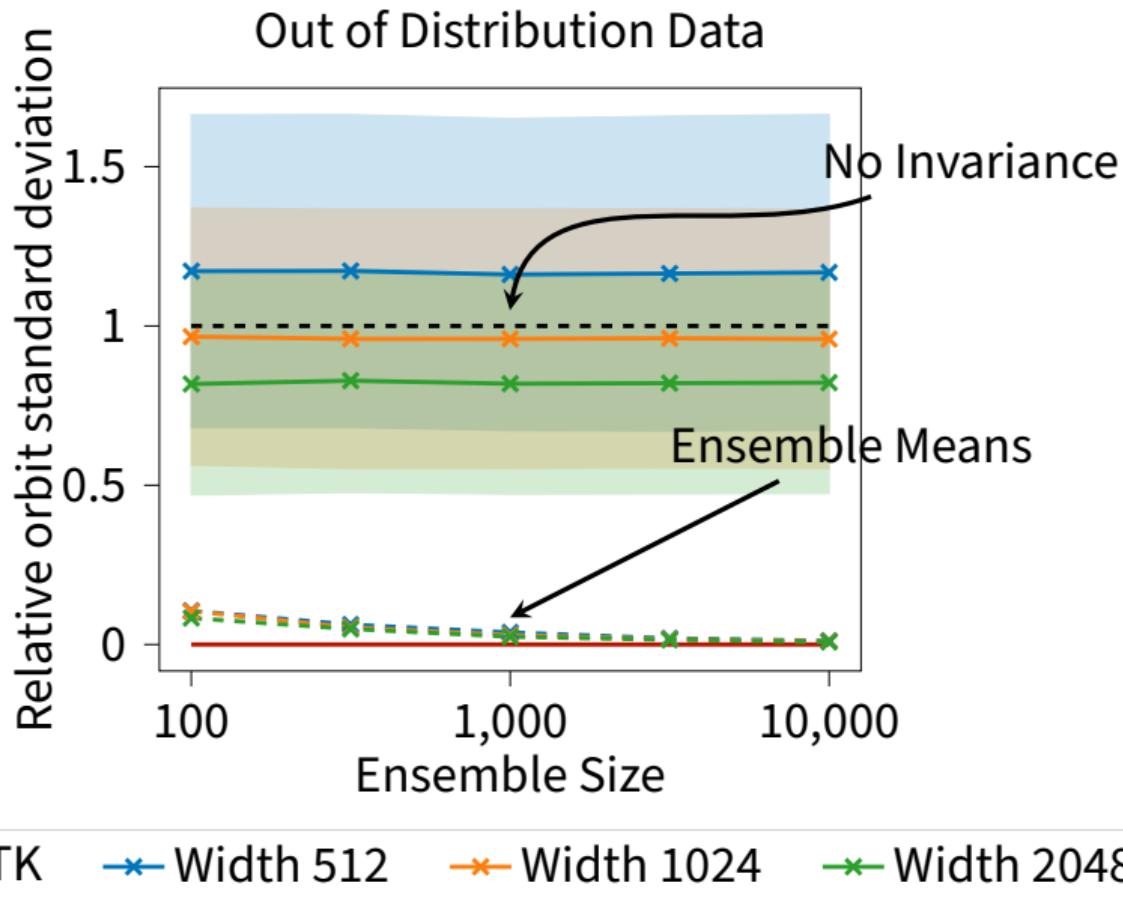
Out of Distribution Data



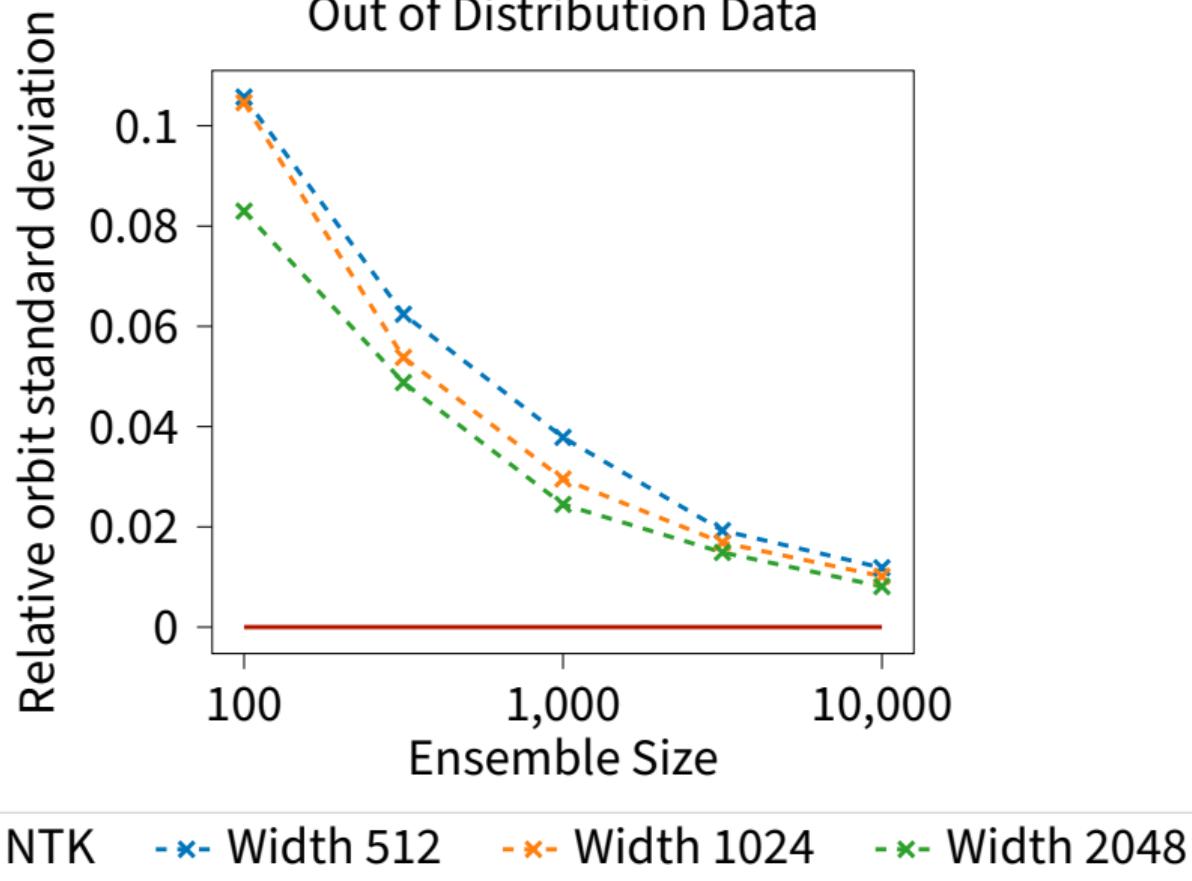
Out of Distribution Data



— NTK —*— Width 512 —*— Width 1024 —*— Width 2048

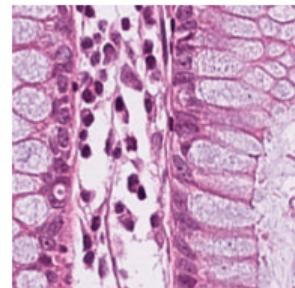


Out of Distribution Data



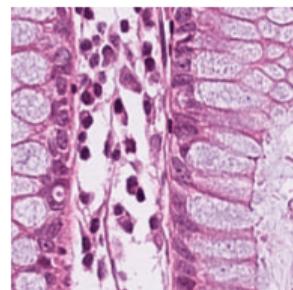
Histological slices

[Kather et al. 2018]



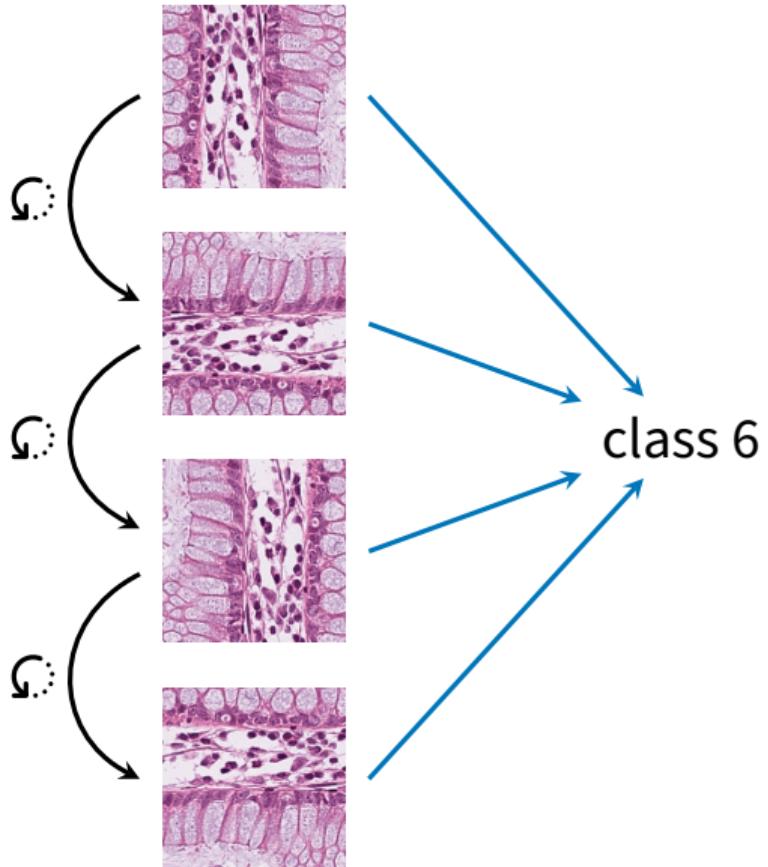
Histological slices

[Kather et al. 2018]

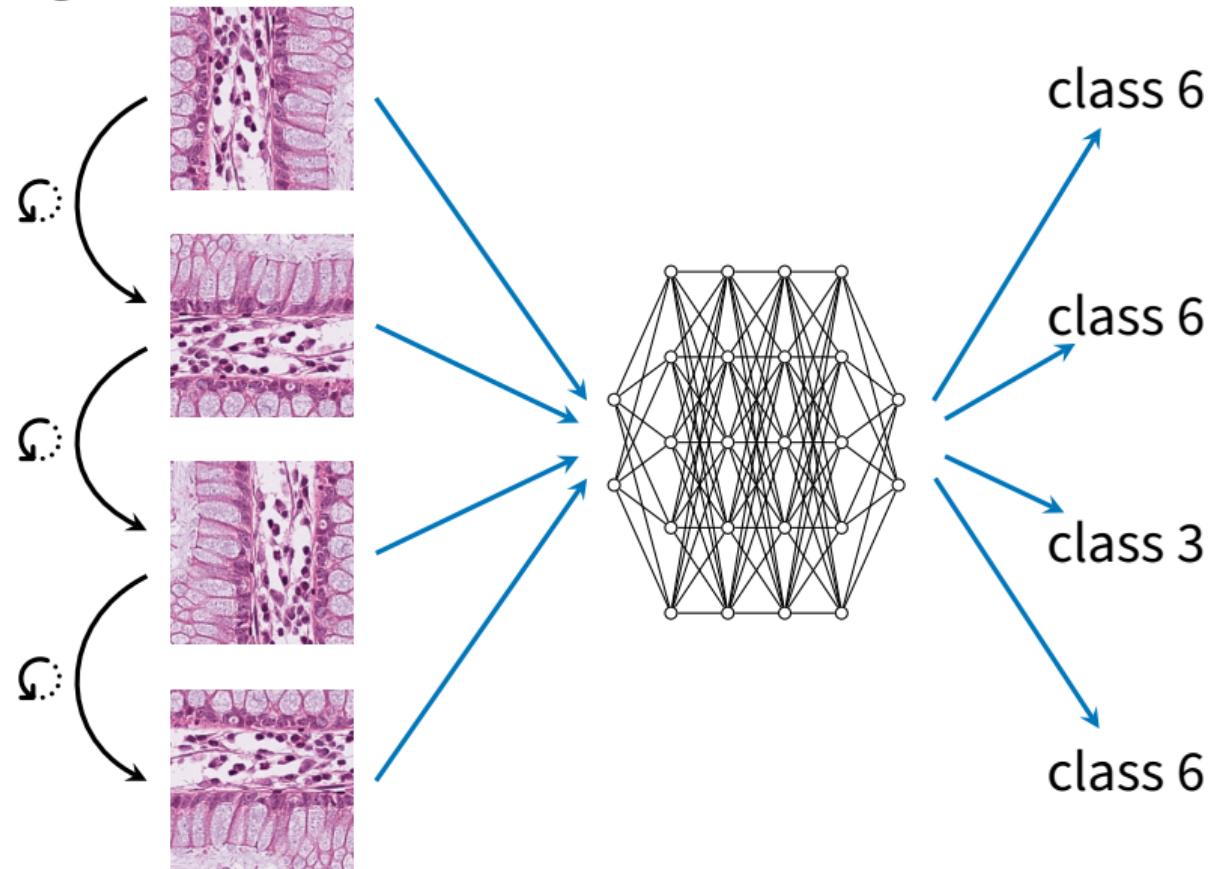


class 6

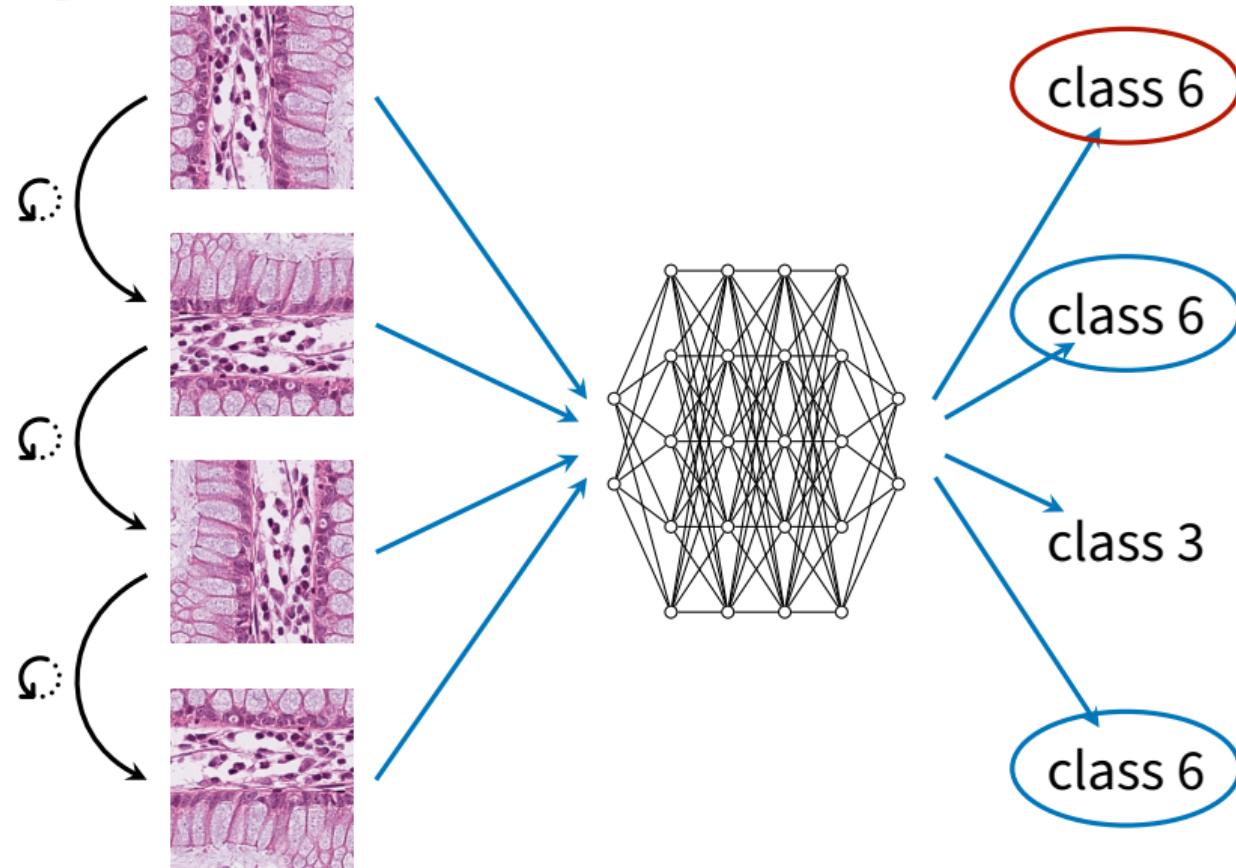
Histological slices



Histological slices

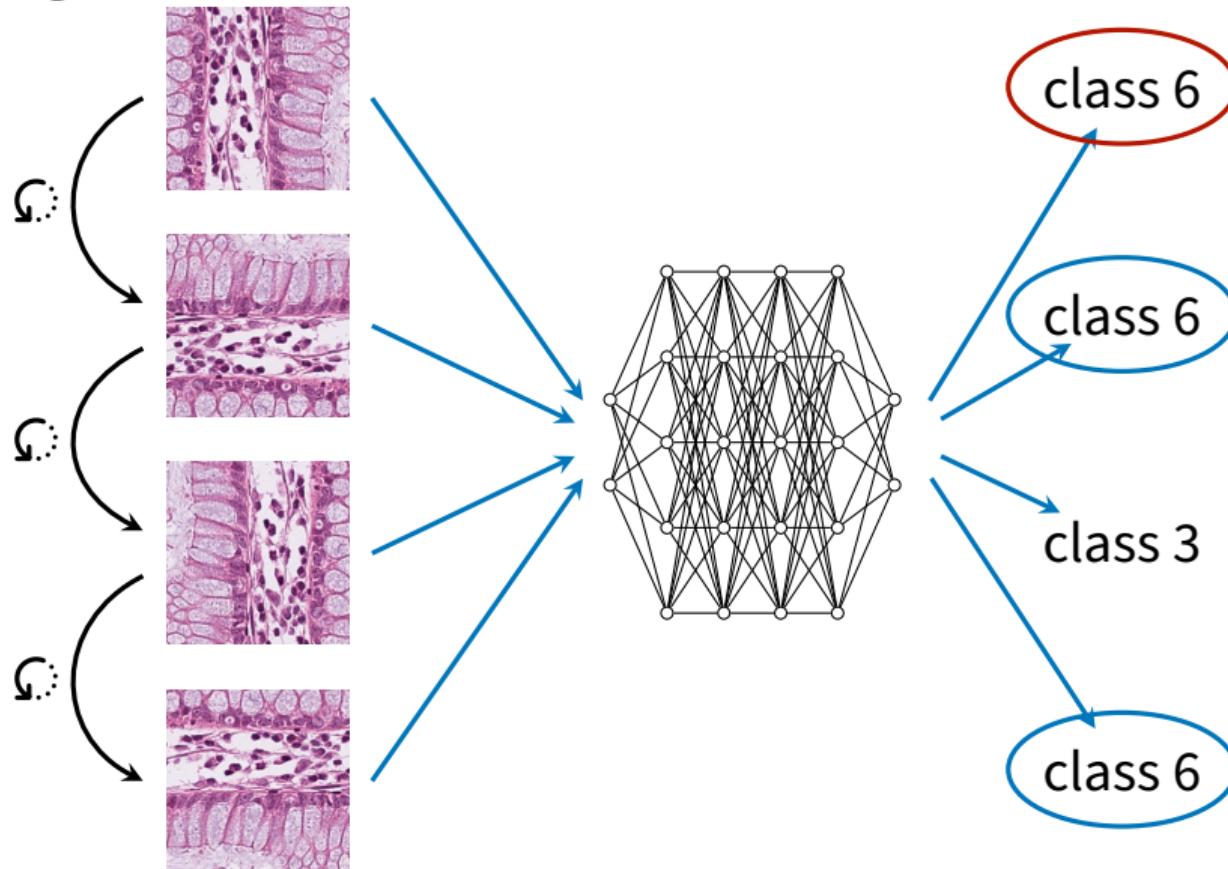


Histological slices

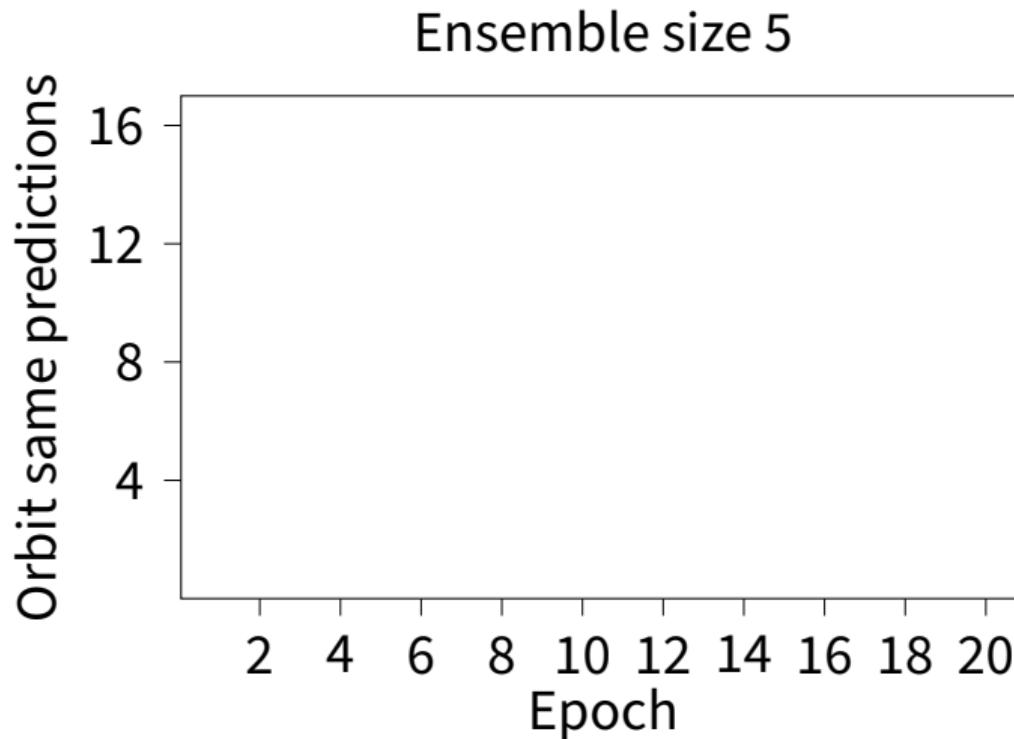


Histological slices

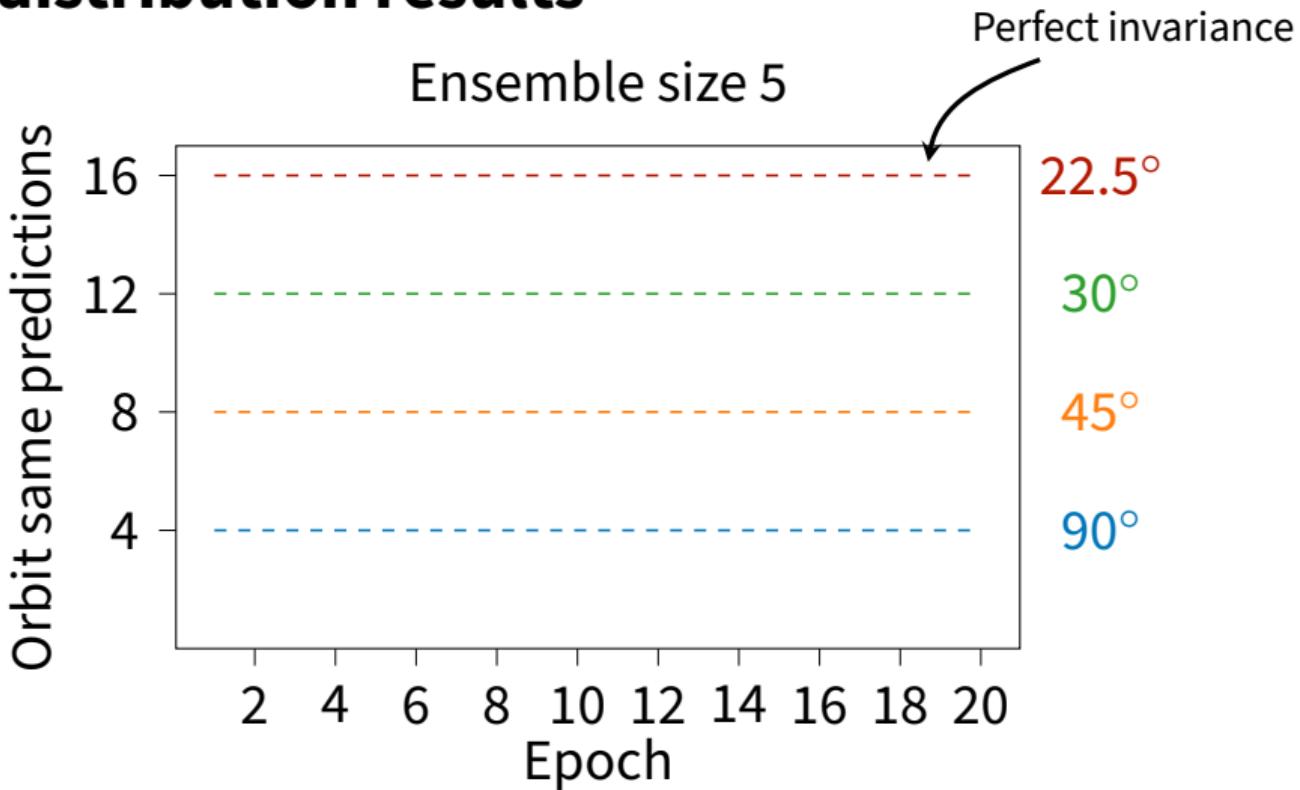
Orbit Same Predictions = 3



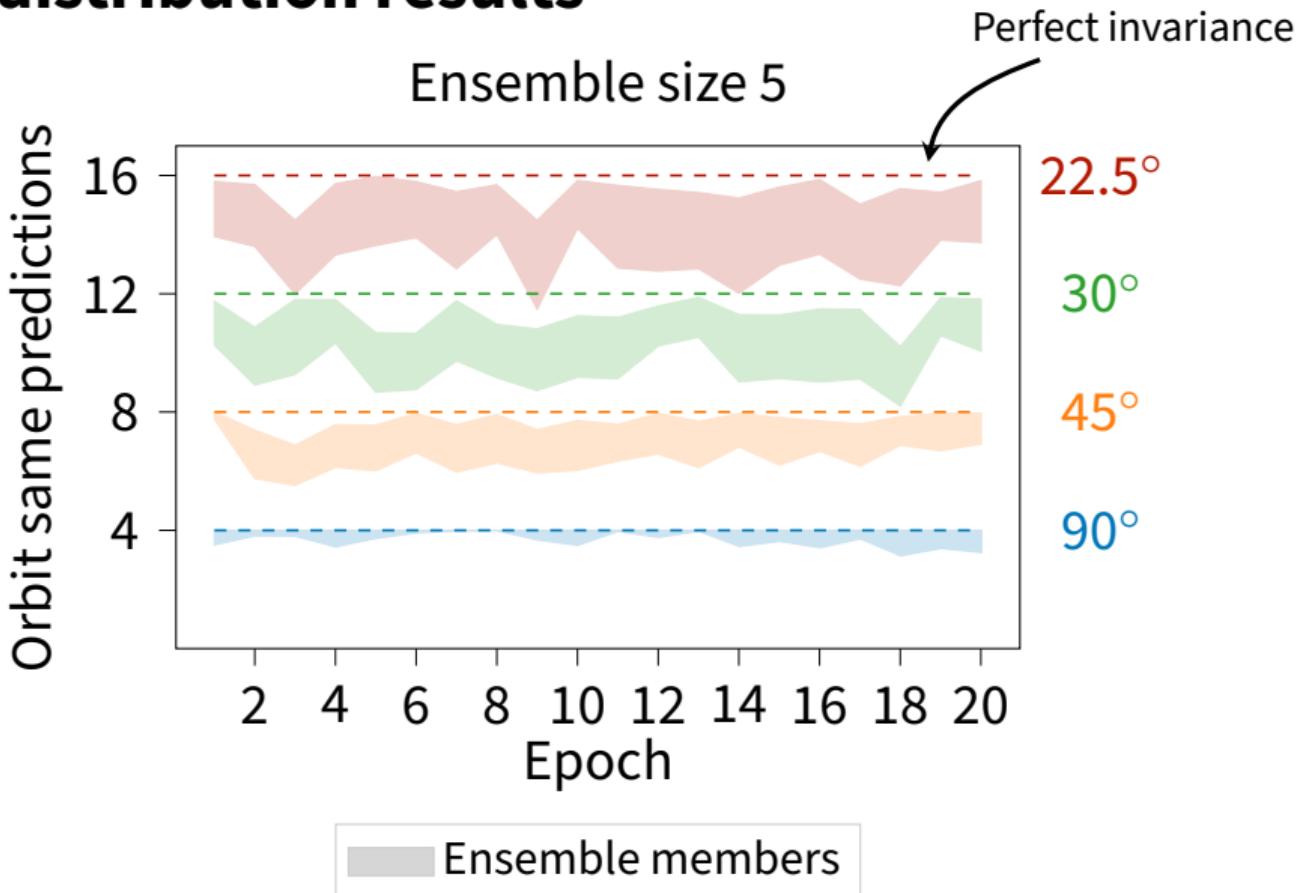
Out of distribution results



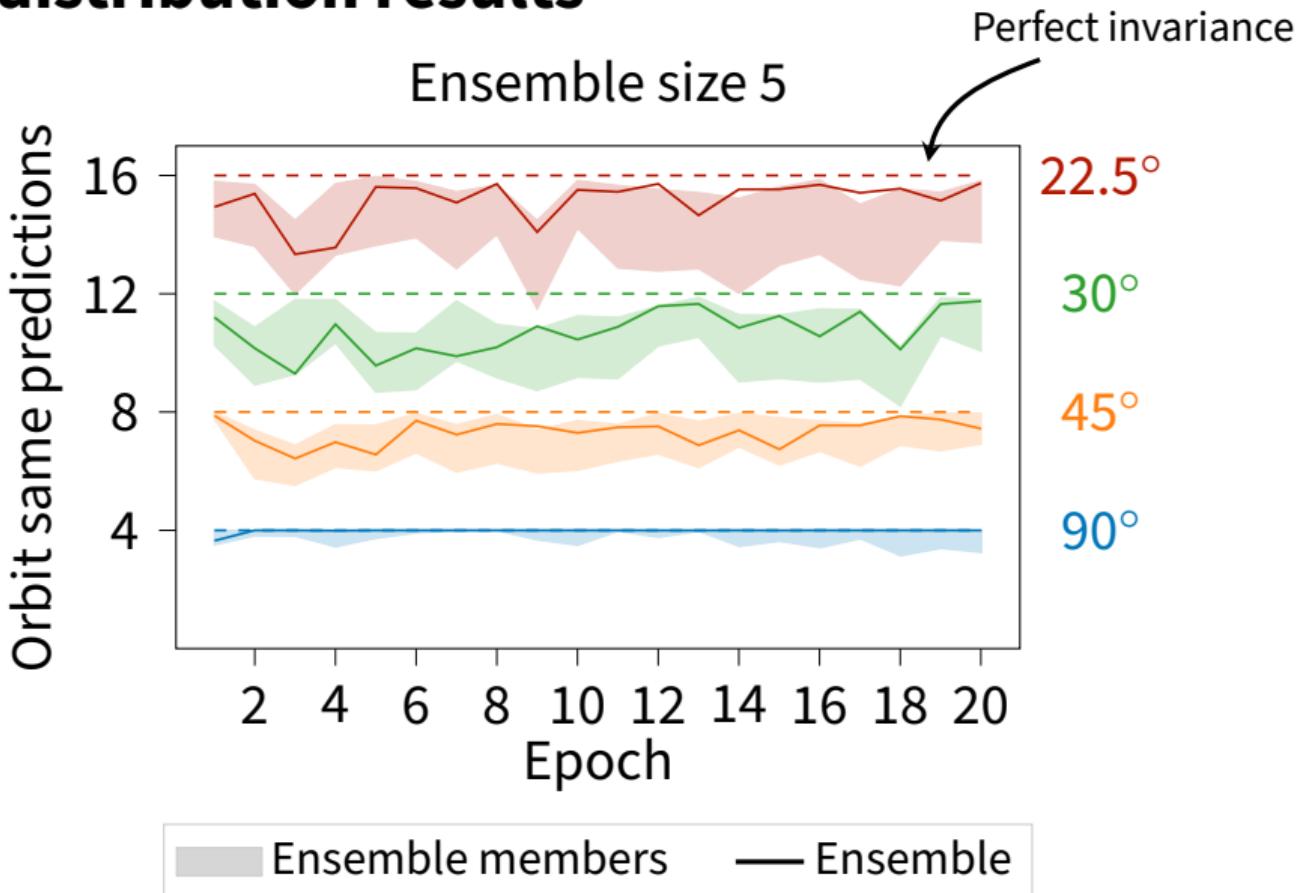
Out of distribution results



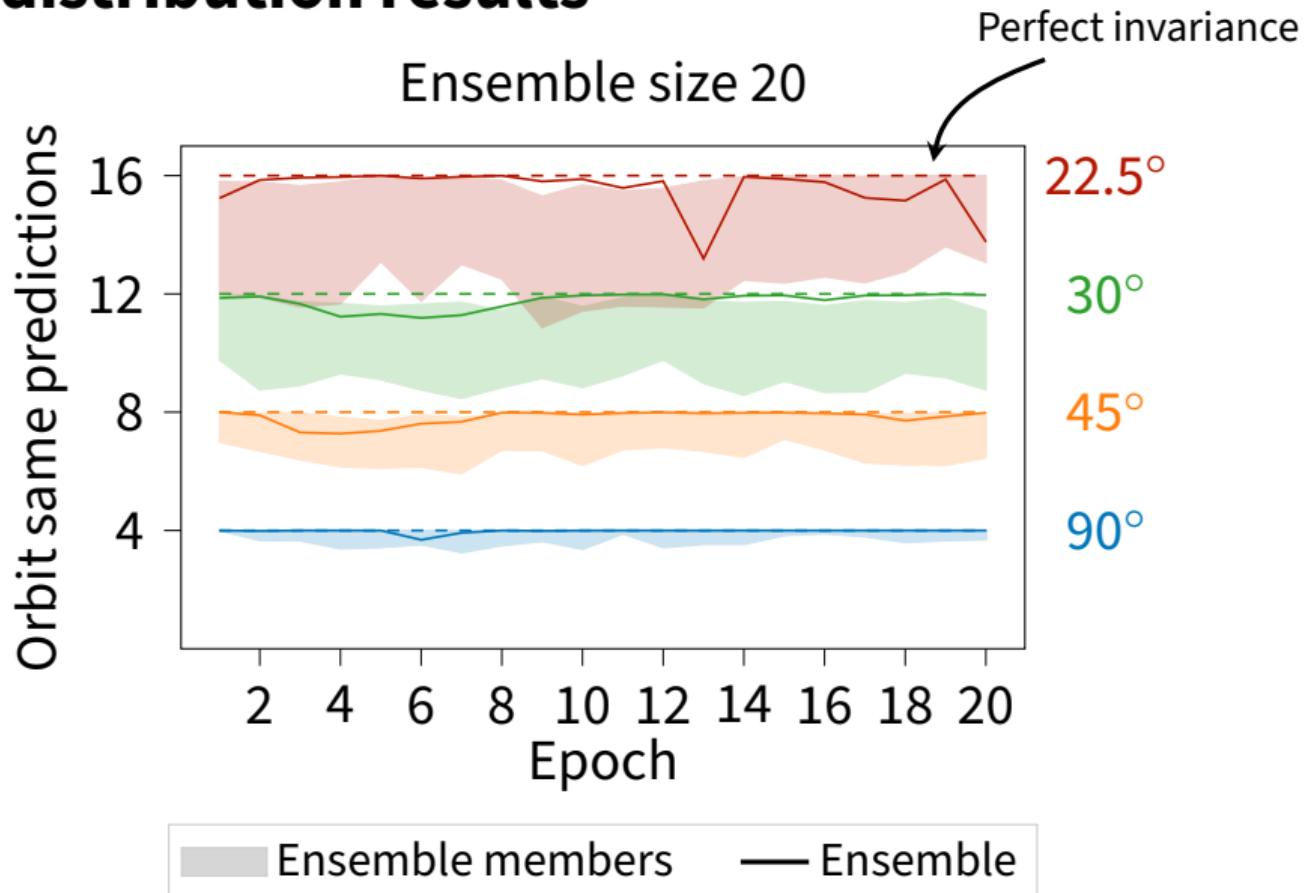
Out of distribution results



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

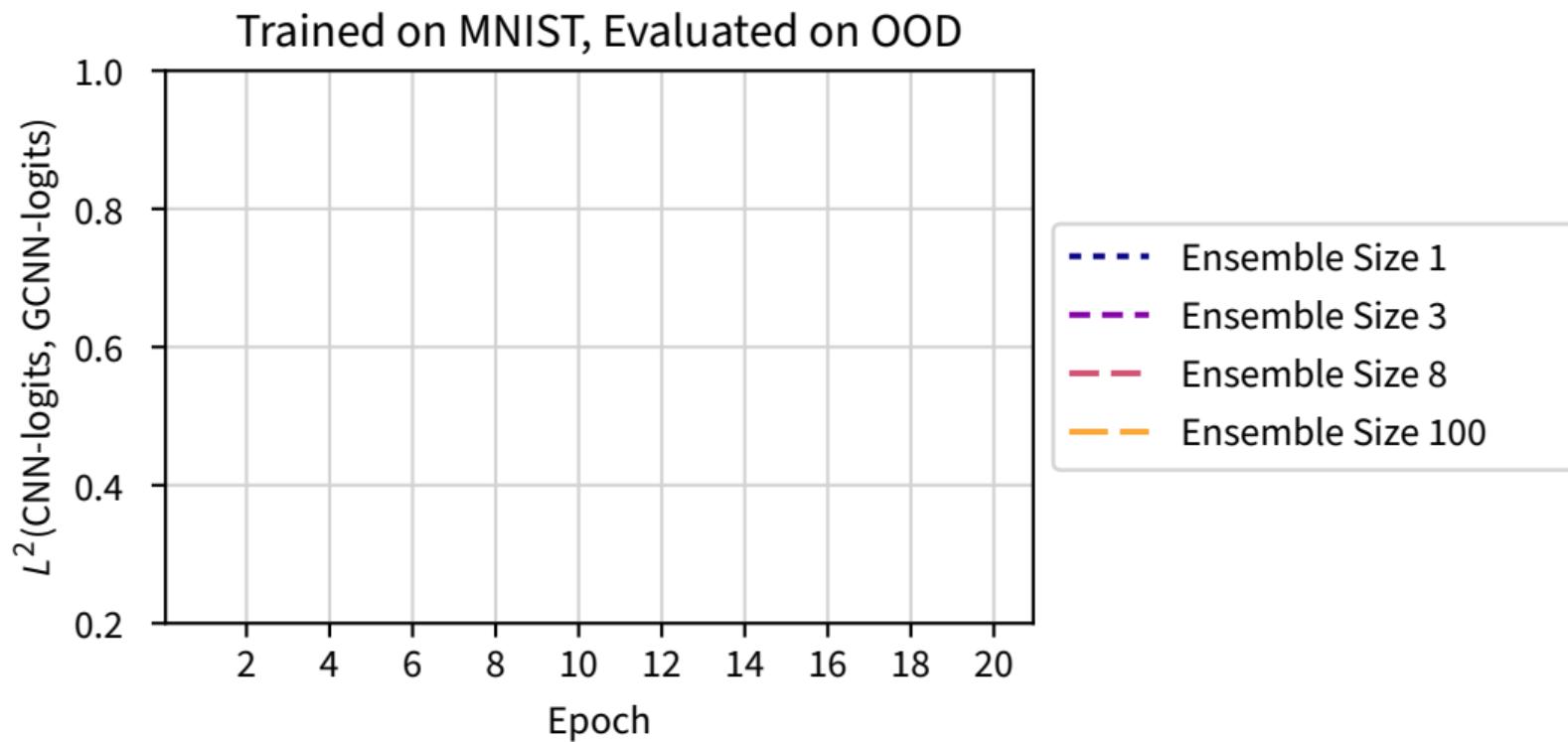
Orbit same predictions out of distribution:

	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

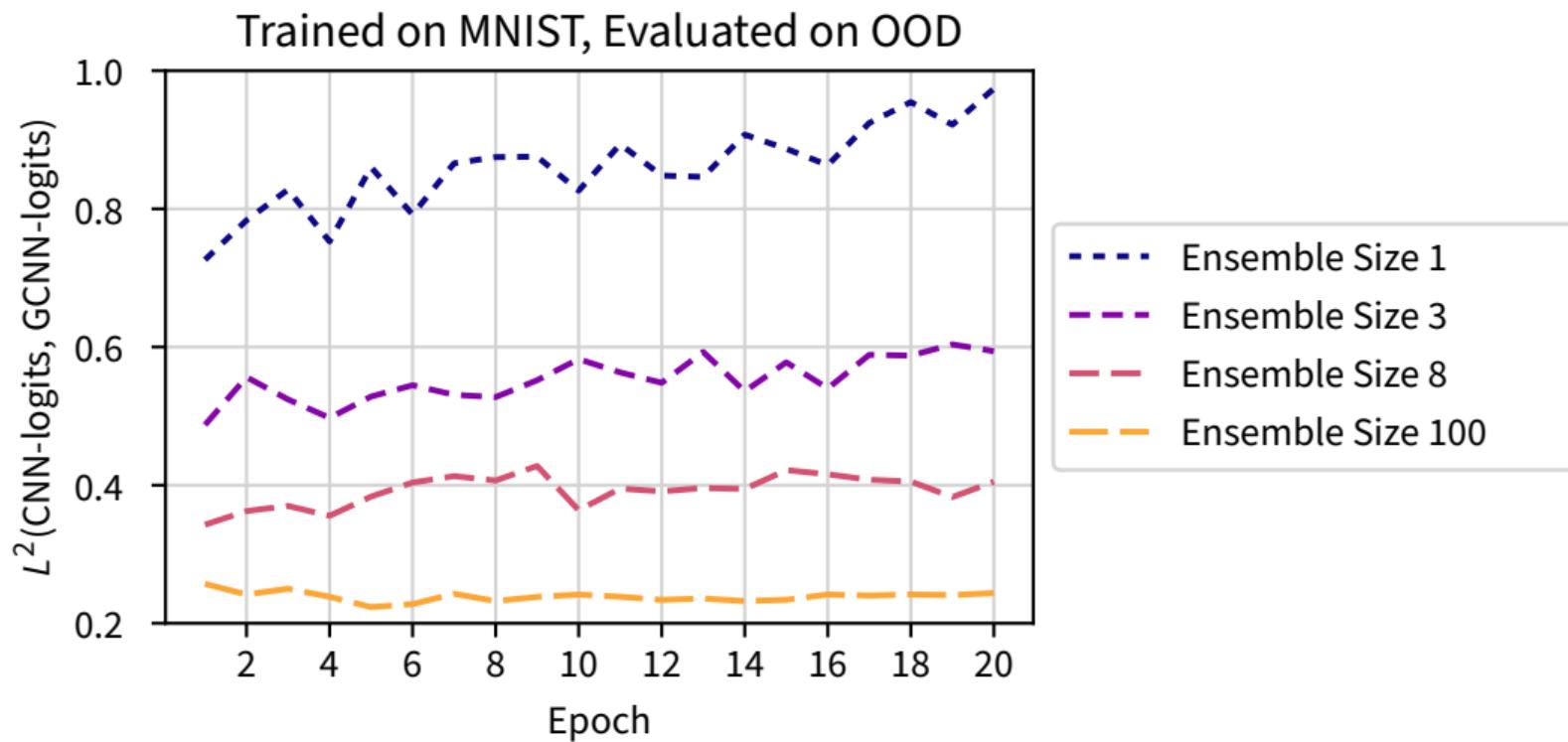
¹[Weiler et al. 2019], ²[Kaba et al. 2022]

Convergence of augmented CNNs to GCNNs

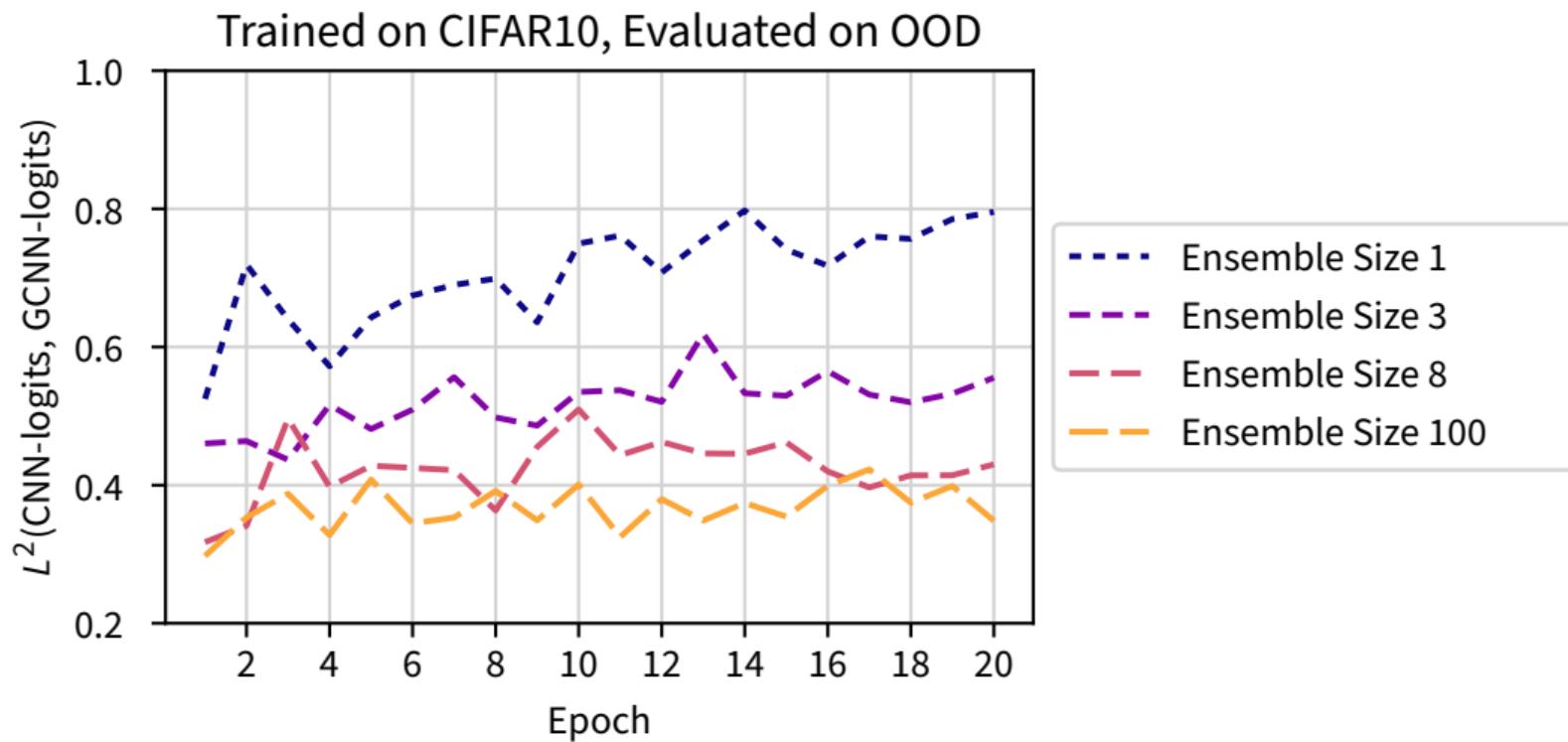
Convergence of augmented CNNs to GCNNs



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.

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

ICML 2025



Thank you

Backup

Empirical NTK

- Consider continuous gradient descent

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

Empirical NTK

- Consider continuous gradient descent

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

- Then, the network evolves according to

$$\frac{d\mathcal{N}_\theta(x)}{dt} = \sum_\mu \frac{\partial \mathcal{N}_\theta(x)}{\partial \theta_\mu} \frac{d\theta_\mu}{dt}$$

Empirical NTK

- Consider continuous gradient descent

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

- Then, the network evolves according to

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

Empirical NTK

- Consider continuous gradient descent

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

- Then, the network evolves according to

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

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

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

Infinite width limit

Consider an MLP in NTK parametrization

$$z^{(\ell)} = \frac{1}{\sqrt{n_{\ell-1}}} W^{(\ell)} \sigma(z^{(\ell-1)}(x)), \quad W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, \quad W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1)$$

Infinite width limit

Consider an MLP in NTK parametrization

$$z^{(\ell)} = \frac{1}{\sqrt{n_{\ell-1}}} W^{(\ell)} \sigma(z^{(\ell-1)}(x)), \quad W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, \quad W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1)$$

At infinite width

$$z_i^{(\ell)}(x) = \sqrt{n_{\ell-1}} \frac{1}{n_{\ell-1}} \sum_{j=1}^{n_\ell} W_{ij}^{(\ell)} \sigma(z_j^{(\ell-1)}(x))$$

Infinite width limit

Consider an MLP in NTK parametrization

$$z^{(\ell)} = \frac{1}{\sqrt{n_{\ell-1}}} W^{(\ell)} \sigma(z^{(\ell-1)}(x)), \quad W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, \quad W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1)$$

At infinite width

$$z_i^{(\ell)}(x) = \sqrt{n_{\ell-1}} \underbrace{\frac{1}{n_{\ell-1}} \sum_{j=1}^{n_\ell} W_{ij}^{(\ell)}}_{\text{mean}} \underbrace{\sigma(z_j^{(\ell-1)}(x))}_{\text{i.i.d.}}$$

Infinite width limit

Consider an MLP in NTK parametrization

$$z^{(\ell)} = \frac{1}{\sqrt{n_{\ell-1}}} W^{(\ell)} \sigma(z^{(\ell-1)}(x)), \quad W^{(\ell)} \in \mathbb{R}^{n_\ell \times n_{\ell-1}}, \quad W_{ij}^{(\ell)} \sim \mathcal{N}(0, 1)$$

At infinite width

$$z_i^{(\ell)}(x) = \sqrt{n_{\ell-1}} \underbrace{\frac{1}{n_{\ell-1}} \sum_{j=1}^{n_\ell} W_{ij}^{(\ell)} \sigma(z_j^{(\ell-1)}(x))}_{\text{mean}} \sim \underbrace{\mathcal{N}(0, \text{Cov}(z_i^{(\ell)}, z_j^{(\ell)}))}_{\text{NNGP}}$$

NTK at initialization

[Jacot et al. 2020]

When taking the layer widths to infinity sequentially, the empirical NTK $\Theta_{ij}^{\theta}(x, x')$ at initialization converges in probability to a deterministic kernel $\Theta(x, x')\delta_{ij}$

NTK at initialization

[Jacot et al. 2020]

When taking the layer widths to infinity sequentially, the empirical NTK $\Theta_{ij}^{\theta}(x, x')$ at initialization converges in probability to a deterministic kernel $\Theta(x, x')\delta_{ij}$

- The deterministic kernel is given in terms of a recursion over layers

NTK at initialization

[Jacot et al. 2020]

When taking the layer widths to infinity sequentially, the empirical NTK $\Theta_{ij}^{\theta}(x, x')$ at initialization converges in probability to a deterministic kernel $\Theta(x, x')\delta_{ij}$

- The deterministic kernel is given in terms of a recursion over layers
- For most common architectures, this recursion can be performed explicitly,
e.g. using neural-tangents Python package

[Novak et al. 2020]

Frozen NTK

[Jacot et al. 2020]

For a nonlinearity which is Lipschitz, twice differentiable and has bounded second derivative,

$$\Theta_{ij}^{\theta_t}(x, x') \rightarrow \Theta(x, x') \delta_{ij}$$

uniformly in t as the layer widths go to infinity sequentially.

Frozen NTK

[Jacot et al. 2020]

For a nonlinearity which is Lipschitz, twice differentiable and has bounded second derivative,

$$\Theta_{ij}^{\theta_t}(x, x') \rightarrow \Theta(x, x') \delta_{ij}$$

uniformly in t as the layer widths go to infinity sequentially.

- Intuitively, this happens because the weight updates vanish in the limit $n \rightarrow \infty$

Frozen NTK

[Jacot et al. 2020]

For a nonlinearity which is Lipschitz, twice differentiable and has bounded second derivative,

$$\Theta_{ij}^{\theta_t}(x, x') \rightarrow \Theta(x, x')\delta_{ij}$$

uniformly in t as the layer widths go to infinity sequentially.

- Intuitively, this happens because the weight updates vanish in the limit $n \rightarrow \infty$
- However, the network still learns

Kernel transformation

The neural tangent kernel Θ as well as the NNGP kernel K transform according to

$$\begin{aligned}\Theta(\rho(g)x, \rho(g)x') &= \rho_K(g)\Theta(x, x')\rho_K^\top(g), \\ K(\rho(g)x, \rho(g)x') &= \rho_K(g)K(x, x')\rho_K^\top(g),\end{aligned}$$

for all $g \in G$ and $x, x' \in X$.

Kernel transformation

The neural tangent kernel Θ as well as the NNGP kernel K transform according to

$$\begin{aligned}\Theta(\rho(g)x, \rho(g)x') &= \rho_K(g)\Theta(x, x')\rho_K^\top(g), \\ K(\rho(g)x, \rho(g)x') &= \rho_K(g)K(x, x')\rho_K^\top(g),\end{aligned}$$

for all $g \in G$ and $x, x' \in X$.

Hence, for MLPs,

$$\Theta(\rho(g)x, \rho(g)x') = \Theta(x, x') \quad \Rightarrow \quad \Theta(\rho(g)x, x') = \Theta(x, \rho^{-1}(g)x')$$

Permutation shift

- On the training data, group transformations permute the samples

$$\rho(g)x_i = x_{\pi_g(i)}, \quad \pi_g \in S_N$$

Permutation shift

- On the training data, group transformations permute the samples

$$\rho(g)x_i = x_{\pi_g(i)}, \quad \pi_g \in S_N$$

- Therefore, for a permutation of training samples associate to g

$$\begin{aligned}\Pi(g)\Theta(X, X) &= \Theta(\rho(g)X, X) \\ &= \Theta(X, \rho^{-1}(g)X) \\ &= \Theta(X, X)(\Pi^{-1}(g))^\top \\ &= \Theta(X, X)\Pi(g)\end{aligned}$$

NTKs for GCNNs

For GCNN-layers, define the NNGP and NTK via

$$K_{g,g'}^{(\ell)}(f, f') = \mathbb{E} \left[[\mathcal{N}^{(\ell)}(f)](g) \left([\mathcal{N}^{(\ell)}(f')] (g') \right)^T \right]$$

NTKs for GCNNs

For GCNN-layers, define the NNGP and NTK via

$$K_{g,g'}^{(\ell)}(f, f') = \mathbb{E} \left[[\mathcal{N}^{(\ell)}(f)](g) \left([\mathcal{N}^{(\ell)}(f')])(g') \right)^T \right]$$

$$\Theta_{g,g'}^{(\ell)}(f, f') = \mathbb{E} \left[\sum_{\ell'=1}^{\ell} \frac{\partial [\mathcal{N}^{(\ell)}(f)](g)}{\partial \theta^{(\ell')}} \left(\frac{\partial [\mathcal{N}^{(\ell)}(f')](g')}{\partial \theta^{(\ell')}} \right)^T \right]$$

NTKs for GCNNs

$$[\mathcal{N}^{(\ell)}(f)](g) = \int_G dg \kappa(g^{-1}h) [\mathcal{N}^{(\ell-1)}(f)](h)$$

The layer-recursion for a GCNN-layer is given by

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

NTKs for GCNNs

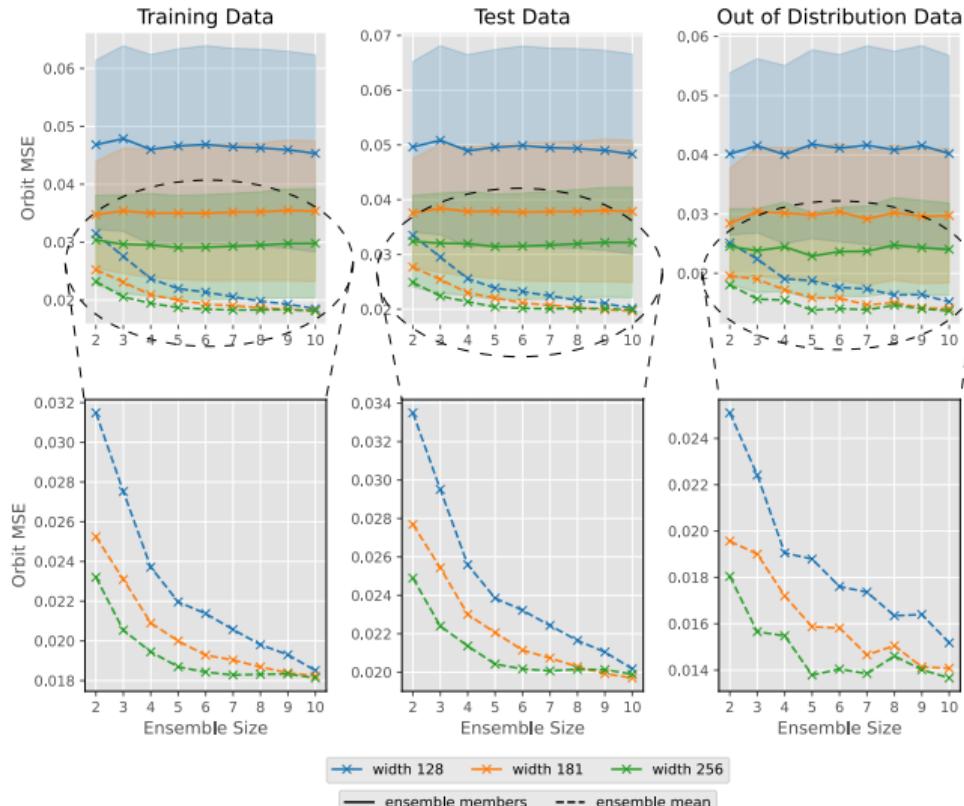
$$[\mathcal{N}^{(\ell)}(f)](g) = \int_G dg \kappa(g^{-1}h) [\mathcal{N}^{(\ell-1)}(f)](h)$$

The layer-recursion for a GCNN-layer is given by

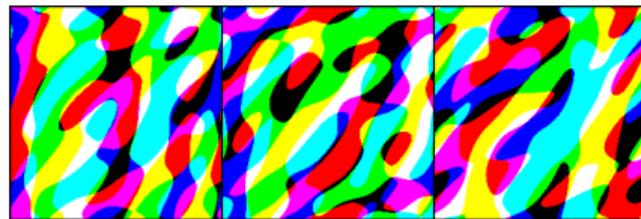
$$K_{g,g'}^{(\ell+1)}(f, f') = \frac{1}{|S_K|} \int_{S_K} dh K_{gh,g'h}^{(\ell)}(f, f')$$

$$\Theta_{g,g'}^{(\ell+1)}(f, f') = K_{g,g'}^{(\ell+1)}(f, f') + \frac{1}{|S_K|} \int_{S_K} dh \Theta_{gh,g'h}^{(\ell)}(f, f')$$

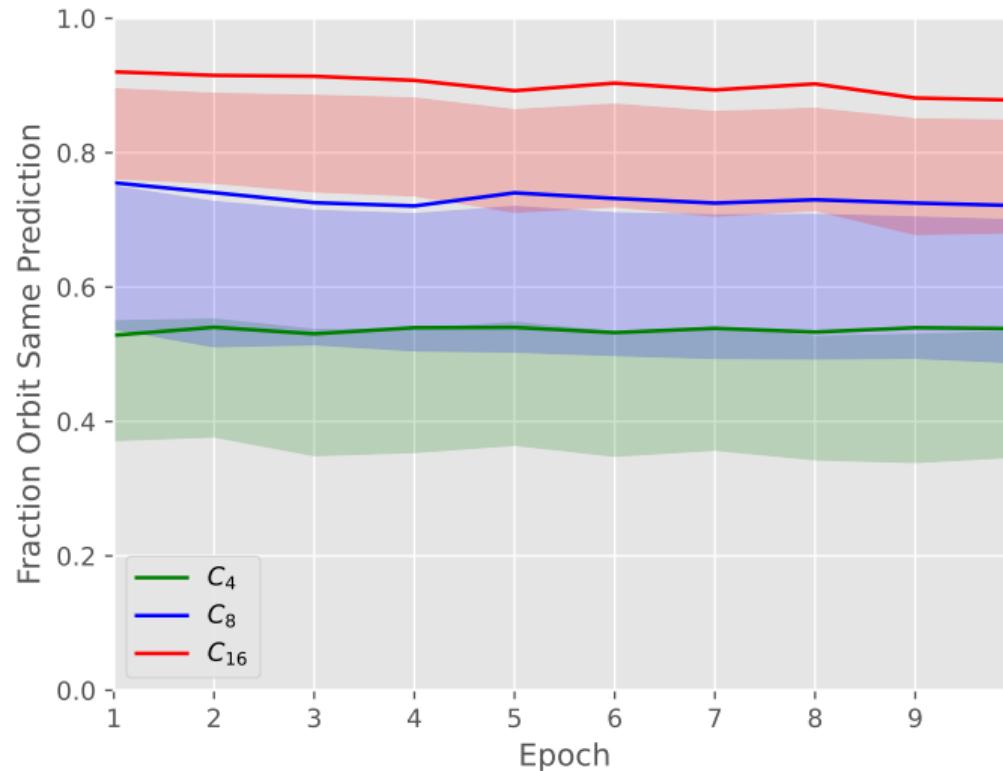
Emergent equivariance of cross products



Histological Data – OOD samples

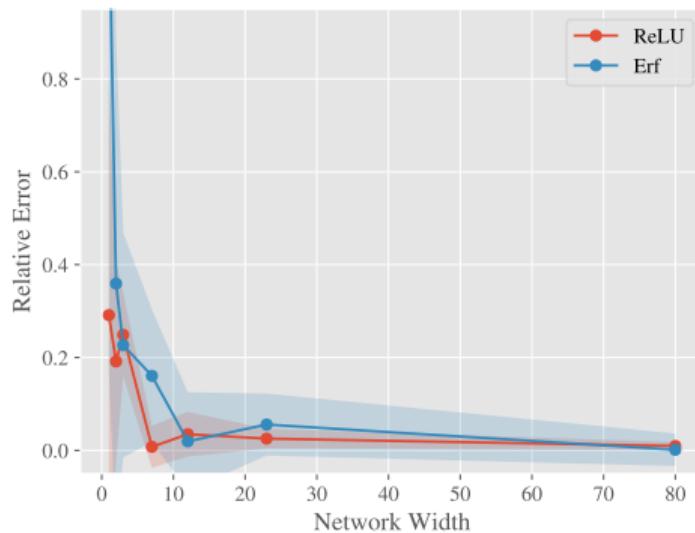


Emergent continuous symmetry on FashionMNIST

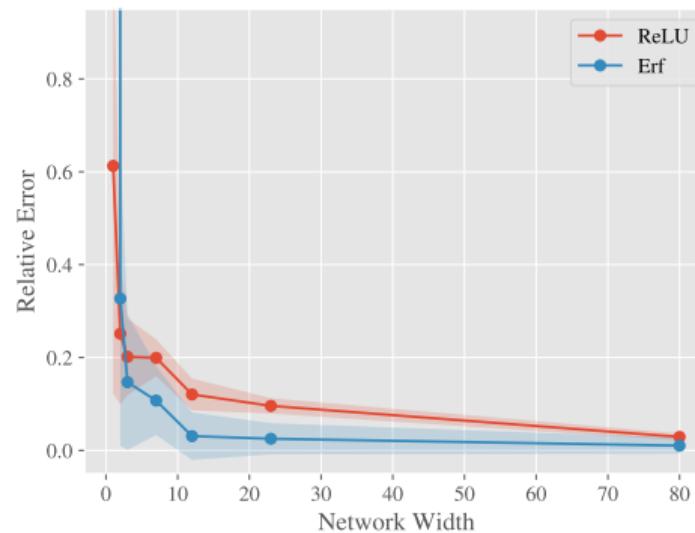


Kernel convergence

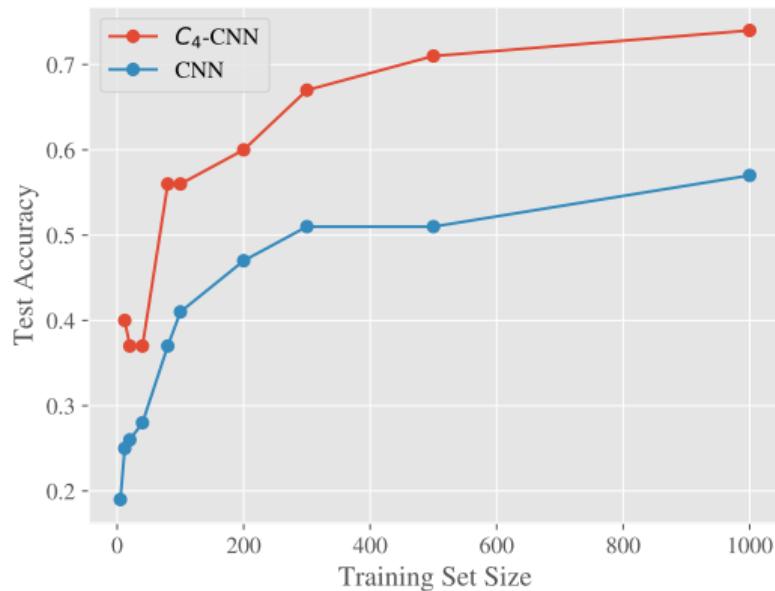
NNGP



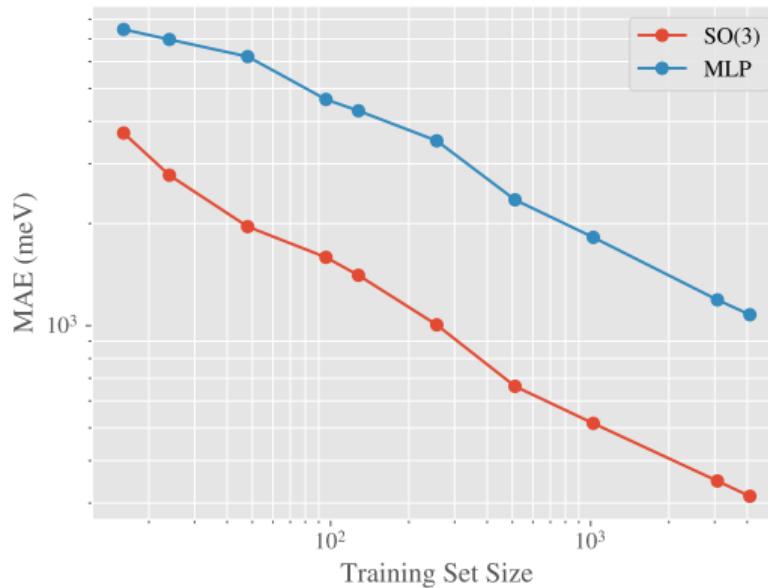
NTK



Equivariant NTKs for medical image classification

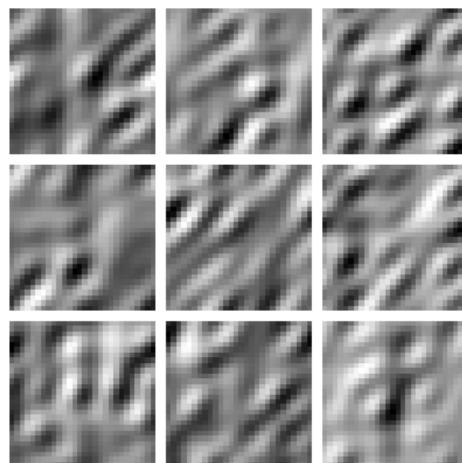


Equivariant NTKs for molecular property regression



OOD samples for CNN to GCNN convergence

MNIST



CIFAR10

