

Samsung Advanced Institute of Technology Al Lab Montreal

Introduction - Networks for networks

Mellon

University



Introduction – Neuron symmetries

Neurons in an MLP can be reordered while maintaining exactly the same function [4]. Reordering neurons here means changing the preceding and following weights attached to the neuron accordingly.



Weight symmetries in a 2-layer MLP subnetwork.

Related works

- Overlook the inherent permutation symmetry
- Rely on intricate weight-sharing patterns to achieve equivariance
- Ignore the network architecture itself, limited to a single architecture

Neural networks are the new data! 🐸

More than 600,000 models on April 2024!





Keywords deep weight space graph neural networks transformers permutation equivariance neural functionals

(neuron i in layer l)

$$\mathbf{x}_i^{(l)} = \sigma \left(\mathbf{b}_i^{(l)} + \right)$$

Neural network as **neural graph**: Node *i* feature: $V_i^{(l)} \leftarrow \mathbf{b}$ Edge $j \to i$ feature: $\boldsymbol{E}_{ij}^{(l)} \leftarrow \mathbf{W}_{ij}^{(l)}$

- Different numbers of layers
- ✓ Residual connections





Graph Neural Networks for Learning Equivariant Representations of Neural Networks

Miltiadis Kofinas^{1*} Boris Knyazev² Yan Zhang² Yunlu Chen³

Gertjan J. Burghouts⁴

* Joint first and last authors

¹University of Amsterdam

²Samsung - SAIT AI Lab, Montreal

Networks for networks, a paradigm shift?		
Modern Paradigm		
Fit signal with an INR and save its parameters Process with parameter space networks		
Model zoos & <mark>NNs are the new data!</mark> Generative parameter space networks on the		

Neural networks as neural graphs



Neural graphs accommodate heterogeneous architectures: \checkmark Architectures with varying computational graphs

 \checkmark Different number of hidden dimensions

 \checkmark Different non-linearities

Node & edge features

Graph networks for neural networks

We adapt existing GNNs and transformers for neural graphs.

NG-GNN: We extend PNA [2] with an MLP that updates the edge features given the incident nodes' features and the previous layer's edge features.

NG-T: We extend Relational Transformer [3] with multiplicative interactions between node and edge features to algorithmically align it with the forward-pass of a neural network.

Positional embeddings

We impose order in the input & output nodes by adding **positional embeddings**. Hidden neurons in each layer are orderless, so they share the same positional embedding.



Probe features

Learn a set of sample inputs that we pass through the input neural network and monitor the intermediate activations and the output. We include these features as additional node features.



Heterogeneous architectures

- Non-linearities: Learned embeddings as added node features
- Residual connections: Additional edges with weight 1
- Normalization layers: Linear layers with diagonal weights
- Self-attention: 3-dimensional edge features (Q, K, V)



Efstratios Gavves¹ Cees G. M. Snoek¹ David W. Zhang^{1*}

³Carnegie Mellon University 4TNO







Ablation study. Importance of positional embeddings on MNIST INR classification.

Method	Accuracy in %
NG-GNN (Ours)	$91.4{\pm}0.6$
NG-GNN w/o positional embeddings	$83.9{\pm}0.3$
NG-T (Ours)	92.4\pm0.3
NG-T w/o positional embeddings	77.9 \pm 0.7

Predicting CNN generalization

Tasks: Predict the generalization performance of CNN classifiers based on their parame-

New dataset: We introduce CNN Wild Park, a dataset of heterogeneous CNNs that vary in the number of layers, kernel sizes, activation functions, and residual connections.

We measure performance using Kendall's au. Higher is better.

Method	CIFAR10-GS [9]	CIFAR10 Wild Par
NFN _{HNP} [10]	0.934 ± 0.001	
StatNN [9]	$0.915 {\pm} 0.002$	$0.719 {\pm} 0.010$
NG-GNN (Ours)	0.930 ± 0.001	0.804 ± 0.009
NG-T (Ours)	$0.935{\scriptstyle\pm0.000}$	$0.817 {\pm} 0.007$

Ablation study. Importance of non-linearity embeddings on predicting CNN generalization on CNN Wild Park.

Method	Kendall's $ au$ (†)
StatNN [9]	0.719 ± 0.010
NG-GNN (Ours)	0.804±0.009
NG-GNN w/o activation embedding	0.778±0.018
NG-T (Ours)	0.817 ± 0.007
NG-T w/o activation embedding	0.728 ± 0.010

Scan me!





[9]



Experiments



MNIST INR dilation

Learning to optimize

An exciting new application for neural graphs!

Task: Train a neural network (optimizer) that can optimize the weights of other neural networks (optimizee).



Figure credit: [1]

We train optimizers on Fashion MNIST, evaluate on Fashion MNIST & CIFAR10. We report the test image classification accuracy (%) after 1,000 steps.

Optimizer	FashionMNIST (validation task)	CIFAR-10 (test task)
Adam [5]	$80.97{\pm}0.66$	54.76±2.82
FF [7] I STM [6]	85.08 ± 0.14 85.69 ± 0.23	57.55 ± 1.06 59 10+0.66
NFN [10]	83.78 ± 0.58	$57.95{\pm}0.64$
NG-GNN (Ours)	$85.91{\pm}0.37$	$64.37{\scriptstyle\pm0.34}$
NG-T (Ours)	$86.52{\scriptstyle \pm 0.19}$	60.79 ± 0.51





Training (left) and testing (right) curves on CIFAR-10 for the baseline optimizers (Adam, FF, LSTM, NFN) and the optimizers trained with our NG-GNN and NG-T.

References

- Marcin Andrychowicz et al. "Learning to learn by gradient descent by gradient descent". In: NIPS. 2016.
- Gabriele Corso et al. "Principal neighbourhood aggregation for graph nets". In: NeurIPS. 2020. [2]
- [3] Cameron Diao and Ricky Loynd. "Relational Attention: Generalizing Transformers for Graph-Structured Tasks". In: ICLR. 2023.
- Robert Hecht-Nielsen. "On the algebraic structure of feedforward network weight spaces". In: Advanced Neural Computers. [4] Elsevier, 1990.
- Diederik P Kingma and Jimmy Ba. "Adam: A Method for Stochastic Optimization". In: ICLR. 2015.
- Luke Metz et al. "Tasks, stability, architecture, and compute: Training more effective learned optimizers, and using them to train themselves". In: arXiv preprint arXiv:2009.11243 (2020).
- Luke Metz et al. "Understanding and correcting pathologies in the training of learned optimizers". In: ICML. 2019.
- Aviv Navon et al. "Equivariant architectures for learning in deep weight spaces". In: ICML. 2023.
- Thomas Unterthiner et al. "Predicting Neural Network Accuracy from Weights". In: arXiv preprint arXiv:2002.11448 (2020). [10] Allan Zhou et al. "Permutation Equivariant Neural Functionals". In: *NeurIPS*. 2023.