

# Hyper-Representations as Generative Models: Sampling Unseen Neural Network Weights

Konstantin Schürholt<sup>1</sup>, Boris Knyazev<sup>2</sup>, Xavier Giró-i-Nieto<sup>3</sup>, Damian Borth<sup>1</sup>

<sup>1</sup> AI:ML Lab, School of Computer Science, University of St. Gallen

<sup>2</sup> Samsung - SAIT AI Lab, Montreal

<sup>3</sup> Institut de Robòtica i Informàtica Industrial, Universitat Politècnica de Catalunya



# Introduction

Learning from populations of Neural Network models is an emerging topic.

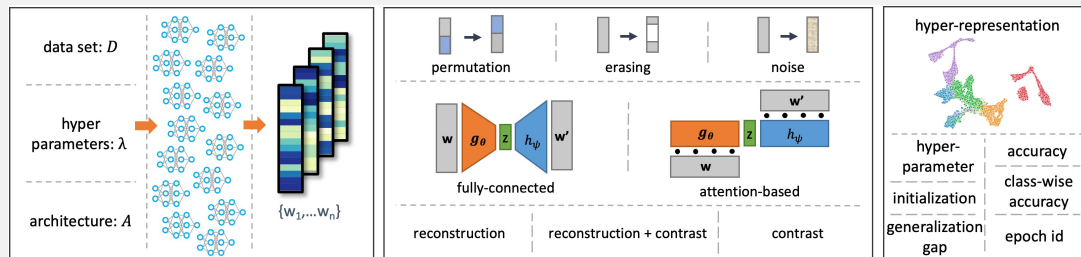
## Discriminative: predict model properties

- Predict: accuracy, generalization gap, hyperparameters
- Features: weights [Unterthiner et al., 2020; Martin et al., 2021], activations [Jiang et al., 2019], graph-metrics [Corneanu et al., 2020]

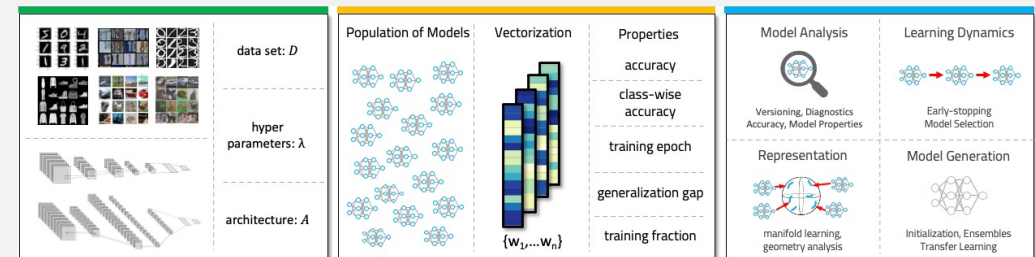
## Generative: generate new models

- HyperNetworks [Ha et al., 2016; Deutsch, 2018; Zhang et al., 2020; Knyazev et al., 2021; Zhmoginov et al., 2022; Ratzlaff and Fuxin, 2019.]
- Transfer Learning, Knowledge Distillation [Shu et al., 2021; Liu et al., 2019.]

### Hyper-Representations: SSL representations of NN weights [Schürholt et al., NeurIPS 2021]



### Model Zoos: Dataset of Diverse NN populations [Schürholt et al., NeurIPS 2022]



## This work: Generative Hyper-Representations

### Goal:

- Better initializations for fine-tuning and transfer-learning
- Knowledge distillation from populations
- Generate diverse ensembles

# Approach

## Zoo Generation Details

- Small CNNs: 3 conv, 2 FC layers
- ~2500 parameters
- 1000 models, trained for 25 epochs
- Initialized with different random seeds

## Hyper-Representation Details

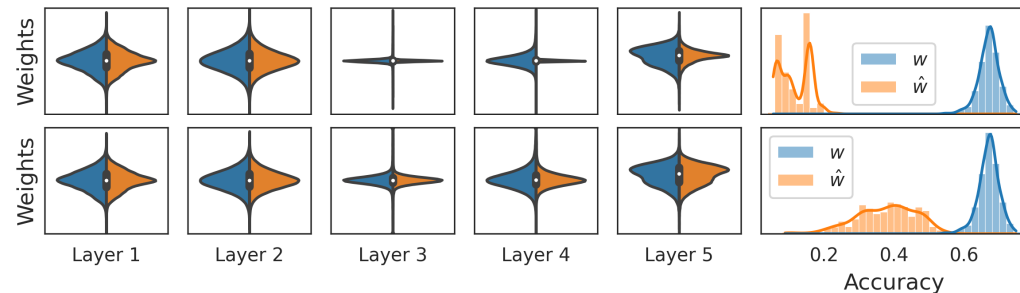
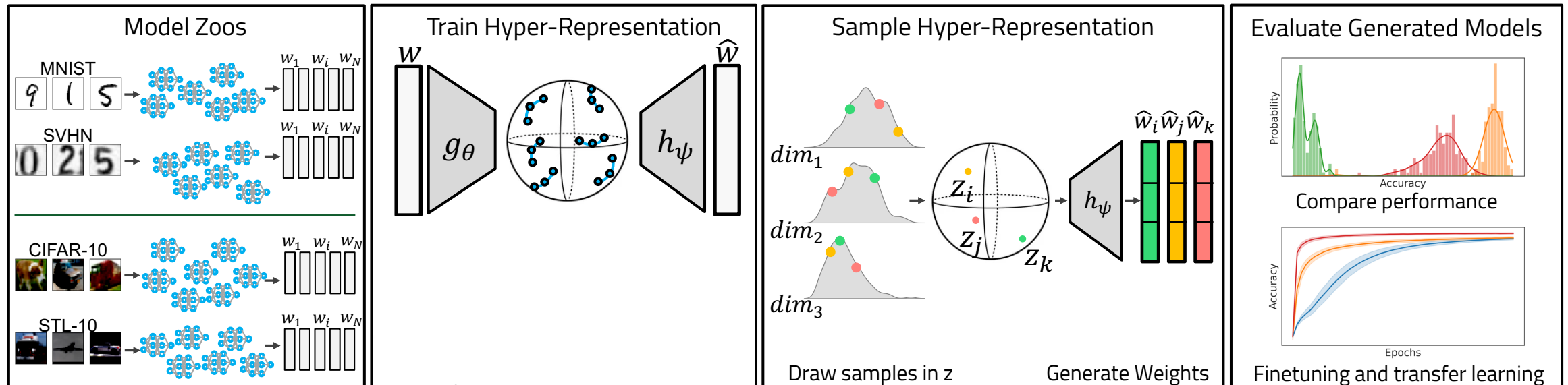
- Encoder-Decoder Transformer
- Trained with Reconstruction and Contrast

## Sampling Details

- Properties like accuracy are embedded in latent
- Problem: space is relatively high-dimensional
- We propose 3 methods to sample good models

## Evaluation Details

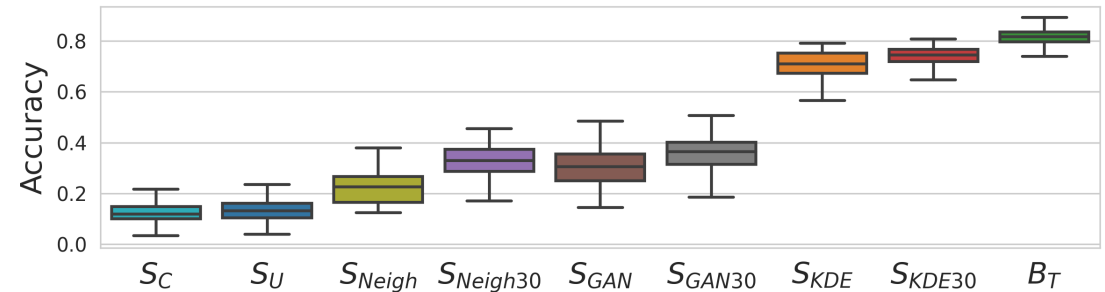
- Use sampled models as initialization:
  - finetuning in-distribution
  - transfer learning
  - generating diverse ensembles



$$\mathcal{L}_{MSE} = \frac{1}{MN} \sum_{i=1}^M \sum_{l=1}^L \left\| \frac{\hat{\mathbf{w}}_i^{(l)} - \mu_l}{\sigma_l} - \frac{\mathbf{w}_i^{(l)} - \mu_l}{\sigma_l} \right\|_2^2 = \frac{1}{MN} \sum_{i=1}^M \sum_{l=1}^L \frac{\|\hat{\mathbf{w}}_i^{(l)} - \mathbf{w}_i^{(l)}\|_2^2}{\sigma_l^2}$$

# Sampling initializations

Sampling methods are targeted:  
distinguish high / low accuracy



Sampled populations are better than (or comparable to) baselines:

- As initialization
- In finetuning (often after 1 ep better than 25 ep trained from scratch)

Method	Ep.	MNIST	SVHN	CIFAR-10	STL-10
$B_T$	0		$\approx 10\%$ (random guessing)		
$B_{KDE30}$	0	$63.2 \pm 7.2$	$10.1 \pm 3.2$	$15.5 \pm 3.4$	$12.7 \pm 3.4$
$S_{KDE30}$	0	<b><math>68.6 \pm 6.7</math></b>	<b><math>51.5 \pm 5.9</math></b>	<b><math>26.9 \pm 4.9</math></b>	<b><math>19.7 \pm 2.1</math></b>
$B_T$	1	$20.6 \pm 1.6$	$19.4 \pm 0.6$	$27.5 \pm 2.1$	$15.4 \pm 1.8$
$B_{KDE30}$	1	$83.2 \pm 1.2$	$67.4 \pm 2.0$	$39.7 \pm 0.6$	<b><math>26.4 \pm 1.6</math></b>
$S_{KDE30}$	1	<b><math>83.7 \pm 1.3</math></b>	<b><math>69.9 \pm 1.6</math></b>	<b><math>44.0 \pm 0.5</math></b>	<b><math>25.9 \pm 1.6</math></b>
$B_T$	25	$83.3 \pm 2.6$	$66.7 \pm 8.5$	$46.1 \pm 1.3$	$35.0 \pm 1.3$
$B_{KDE30}$	25	<b><math>93.2 \pm 0.6</math></b>	<b><math>75.4 \pm 0.9</math></b>	$48.1 \pm 0.6$	<b><math>38.4 \pm 0.9</math></b>
$S_{KDE30}$	25	$93.0 \pm 0.7$	$74.2 \pm 1.4$	<b><math>48.6 \pm 0.5</math></b>	$38.1 \pm 1.1$
$B_T$	50	$91.1 \pm 2.6$	$70.7 \pm 8.8$	$48.7 \pm 1.4$	$39.0 \pm 1.0$

# Sampling for New Tasks and Architectures

Sampled populations outperform or match baselines for transfer-learning

Method	SVHN to MNIST		
	Ep. 0	Ep. 1	Ep. 50
$B_T$	$10.0 \pm 0.6$	$20.6 \pm 1.6$	$91.1 \pm 1.0$
$B_F$	<b><math>33.4 \pm 5.4</math></b>	$84.4 \pm 7.4$	$95.0 \pm 0.8$
$S_{\text{KDE30}}$	$31.8 \pm 5.6$	<b><math>86.9 \pm 1.4</math></b>	<b><math>95.5 \pm 0.4</math></b>

Sampled weights generalize to changed architectures and outperform random initialization

Initialization	Epoch 1	Epoch 5	Epoch 50
3-conv (r. i.) + res-skip (r. i.)	$18.9 \pm 1.6$	$31.4 \pm 17$	$50.6 \pm 28$
3-conv (gen.) + res-skip (r. i.)	<b><math>34.5 \pm 14</math></b>	<b><math>60.5 \pm 21</math></b>	<b><math>68.0 \pm 21</math></b>
4-conv (r. i.)	$19.2 \pm 1.0$	$19.2 \pm 0.9$	$55.2 \pm 11$
4-conv (gen.)	<b><math>44.0 \pm 4.5</math></b>	<b><math>57.8 \pm 3.5</math></b>	<b><math>67.6 \pm 1.9</math></b>
4-conv + id.-skip (r. i.)	$18.9 \pm 1.0$	$19.6 \pm 1.7$	$56.4 \pm 7.9$
4-conv + id.-skip (gen.)	<b><math>48.0 \pm 4.0</math></b>	<b><math>59.9 \pm 2.5</math></b>	<b><math>66.4 \pm 1.7</math></b>

# Acknowledgements

Find our work at [hsg.ai/neurips22](https://hsg.ai/neurips22)

*Funding:*

- Google Research Scholar Award (Damian Borth)
- HSG Basic Research Fund
- MCIN/ AEI /10.13039/501100011033

# References

- [1] Thomas Unterthiner, Daniel Keysers, Sylvain Gelly, Olivier Bousquet, and Ilya Tolstikhin. Predicting Neural Network Accuracy from Weights. arXiv:2002.11448, February 2020.
- [2] Charles H Martin, Tongsu Serena Peng, and Michael W Mahoney. Predicting trends in the quality of state-of-the-art neural networks without access to training or testing data. *Nature Communications*, 12(1):1–13, 2021.
- [3] Yiding Jiang, Dilip Krishnan, Hossein Mobahi, Samy Bengio. Predicting the Generalization Gap in Deep Networks with Margin Distributions. ICLR 2019.
- [4] Ciprian Corneanu, Meysam Madadi, Sergio Escalera, Aleix Martinez. Computing the Testing Error Without a Testing Set. CVPR 2020.
- [5] David Ha, Andrew Dai, and Quoc V. Le. HyperNetworks. In arXiv:1609.09106 [Cs], 2016.
- [6] Lior Deutsch. Generating Neural Networks with Neural Networks. arXiv:1801.01952 [cs, stat], April 2018.
- [7] Chris Zhang, Mengye Ren, and Raquel Urtasun. Graph HyperNetworks for Neural Architecture Search. arXiv:1810.05749 [cs, stat], December 2020.
- [8] Boris Knyazev, Michal Drozdal, Graham W. Taylor, and Adriana Romero-Soriano. Parameter Prediction for Unseen Deep Architectures. In Conference on Neural Information Processing Systems (NeurIPS), 2021.
- [9] Andrey Zhmoginov, Mark Sandler, and Max Vladymyrov. HyperTransformer: Model Generation for Supervised and Semi-Supervised Few-Shot Learning. arXiv:2201.04182 [cs], January 2022.
- [10] Neale Ratzlaff and Li Fuxin. HyperGAN: A Generative Model for Diverse, Performant Neural Networks. In Proceedings of the 36th International Conference on Machine Learning, pages 5361–5369. PMLR, May 2019.
- [11] Yang Shu, Zhi Kou, Zhangjie Cao, Jianmin Wang, and Mingsheng Long. Zoo-tuning: Adaptive transfer from a zoo of models. In International Conference on Machine Learning, pages 457 9626–9637. PMLR, 2021.
- [12] Iou-Jen Liu, Jian Peng, and Alexander G. Schwing. Knowledge Flow: Improve Upon Your Teachers. In International Conference on Learning Representations (ICLR), April 2019.
- [13] Konstantin Schürholt, Dimche Kostadinov, and Damian Borth. Self-Supervised Representation Learning on Neural Network Weights for Model Characteristic Prediction. In Conference on Neural Information Processing Systems (NeurIPS), volume 35, 2021.